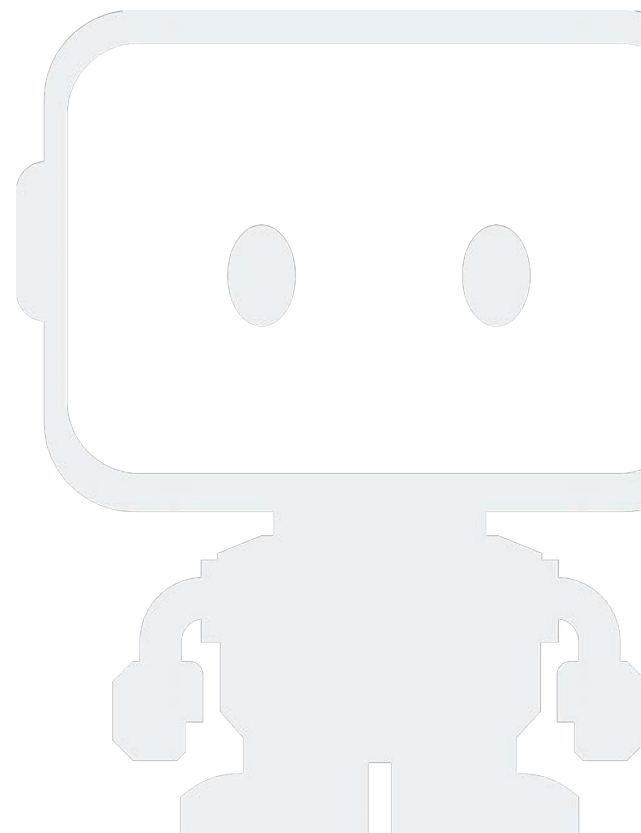


DataRobot

Paxata

Data Prep Documentation

Version: 2021.2 ; Updated: April 25, 2023



Paxata Data Prep

DataRobot Data Prep lets you gather, explore, and prepare data from multiple sources for machine learning. You can save and share your data and the steps you used to prepare it.

When you think about DataRobot Data Prep, think:

- Libraries where you save the datasets you prep
- Projects where you perform your data prep
- Data that you import, clean, and combine

These sections describe how to work with Data Prep to clean and prepare your data for machine learning:

TOPIC	DESCRIBES...
Getting started with Data Prep	Complete the Data Prep quickstart and tour the basics of the Data Prep Library and Projects .
Data Prep data connections	Configure Data Prep connectors.
Working with datasets	Import datasets and perform other dataset operations such as exporting , profiling , and updating .
Working with Project tools	Use Data Prep Project tools to clean and shape your data.
Working with column data	Update the data in columns using a wide range of column operations.
Connecting to data sources	Configure your data source connections so that you can import data from and export data to external systems.
Automation and operationalization	Employ workflow automations to reduce the number of repetitive tasks taken to produce AnswerSets.
Advanced topics	Create ClicktoPrep links , use interactive mode , and understand Data Prep infrastructure and application security .

Get started with Data Prep

These sections describe all you need to know to get started with a Data Prep project:

TOPIC	DESCRIBES...
Data Prep quickstart	Import your data and set up a project. Then, prep, export, and share your data.
Tour the basics of Data Prep	Tour the major components of the Data Prep application and learn how to get around.
Data Prep library	Add new datasets, manage existing datasets, and publish AnswerSets.
Data Prep projects	Explore and prepare your data in a project.

Prep your data

To prep your data using DataRobot Data Prep, you start by importing your data. You can import a local dataset or you can connect to an external data source. This quickstart walks you through importing a local dataset.

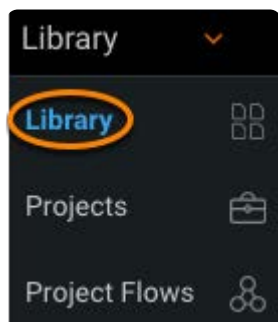
To complete the quickstart, you first log in to DataRobot Data Prep. Once you log in, complete these steps:

1. [Add data](#) to your library.
2. [Start a project](#).
3. [Prep your data](#) in a project.
4. [Publish your data](#) as an *AnswerSet*—a snapshot of your prepped data.
5. [Export your prepped data](#).

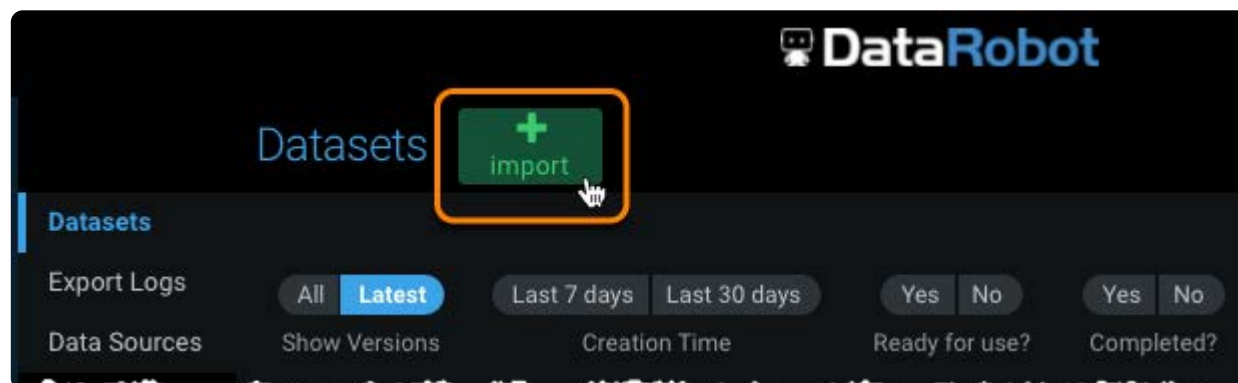
Add data to your Data Prep library

In this quickstart, you will import a local dataset into your library. You can also import the data directly into a project and you can import data from external data sources. To learn about these other options for importing, see [Work with datasets](#).

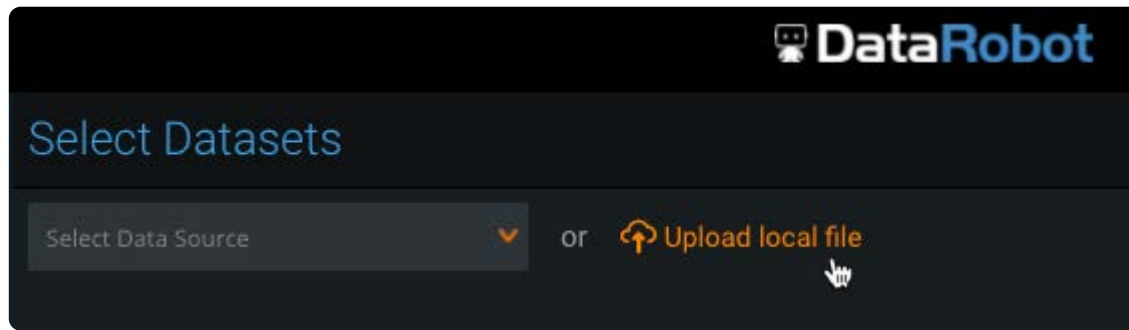
1. In DataRobot Data Prep, select **Library** on the top left.



2. On the top of the **Library** page, click **+ import**.



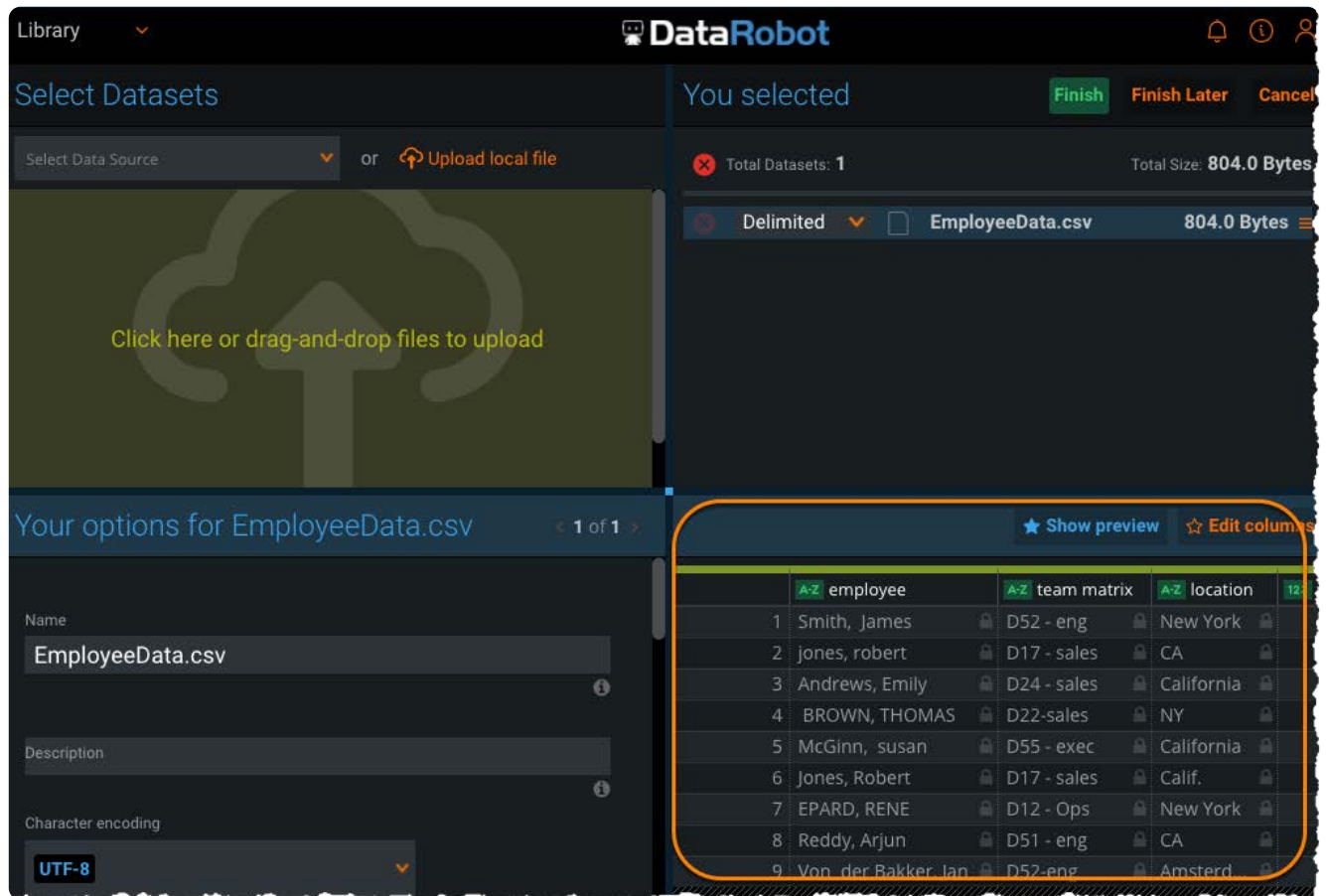
3. Click **Upload local file**.



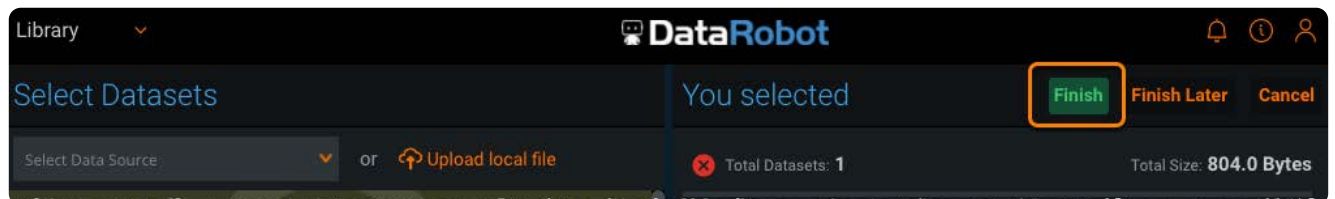
4. Browse for the file or drag the file to the drag-and-drop area.



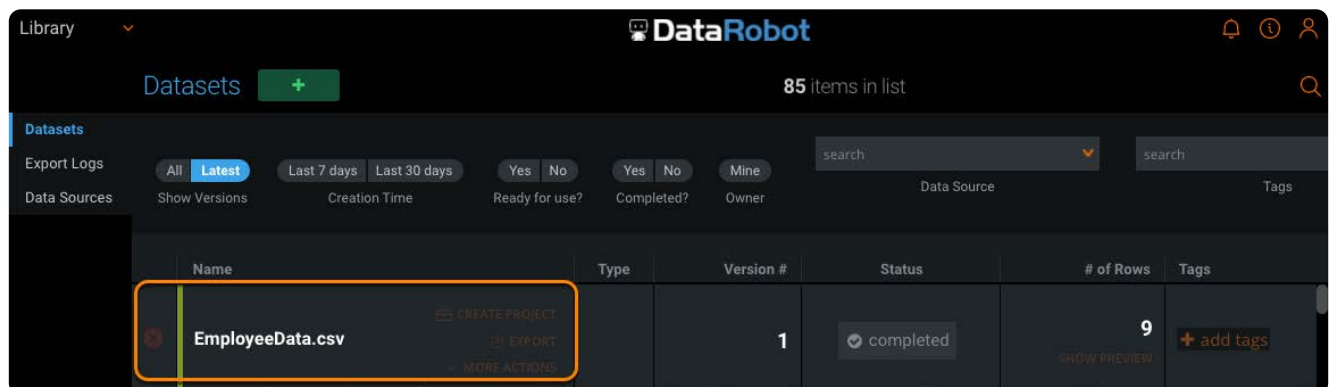
5. Check the preview of the dataset in the lower right.



6. If your data looks correct, click **Finish** on the top right.



Data Prep imports your dataset into the library and you can begin prepping it.



Start a Data Prep project

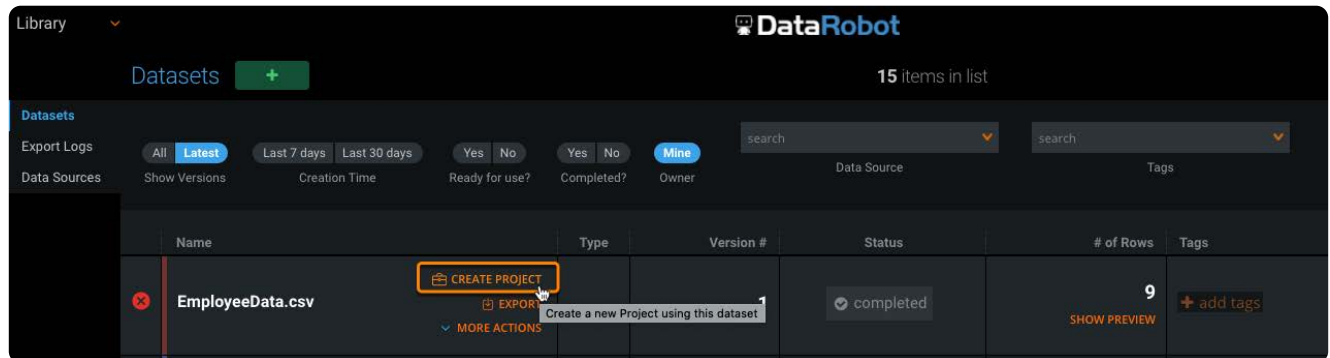
You can start a new project from:

- The **Library** page, where you select the dataset you want to use as the starting point for your project.

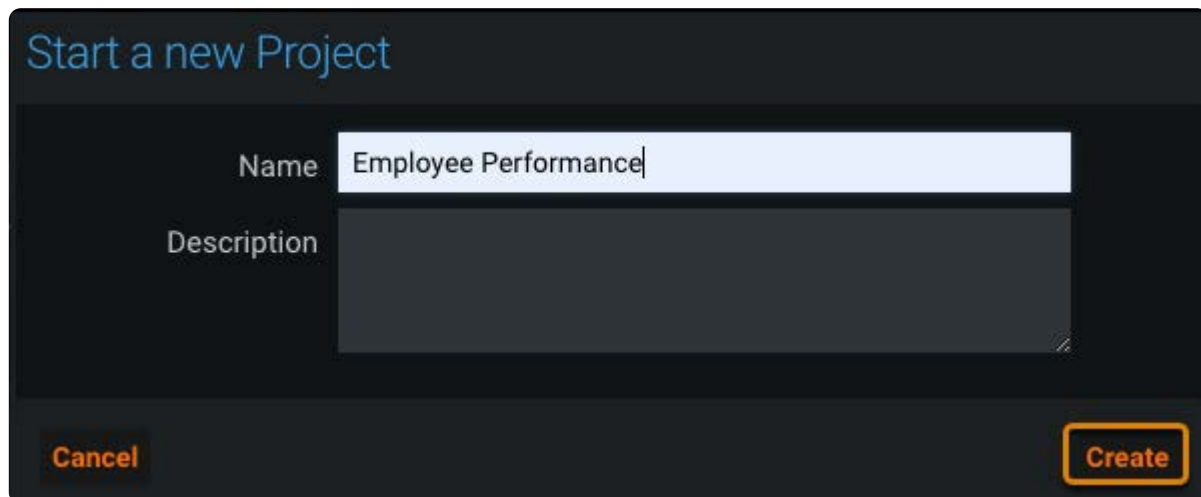
- The **Projects** page, where you start with an empty project and then add your data to it.

Start a new project from the library

1. Select **Library** on the top left.
2. Locate the dataset you uploaded and click **Create Project**.



3. In the **Start a new Project** dialog, enter the project **Name** and an optional **Description**.

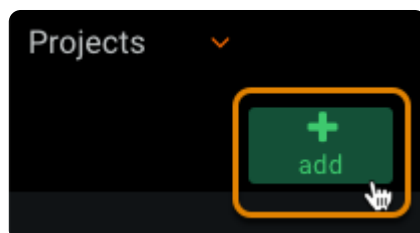


4. Click **Create**.

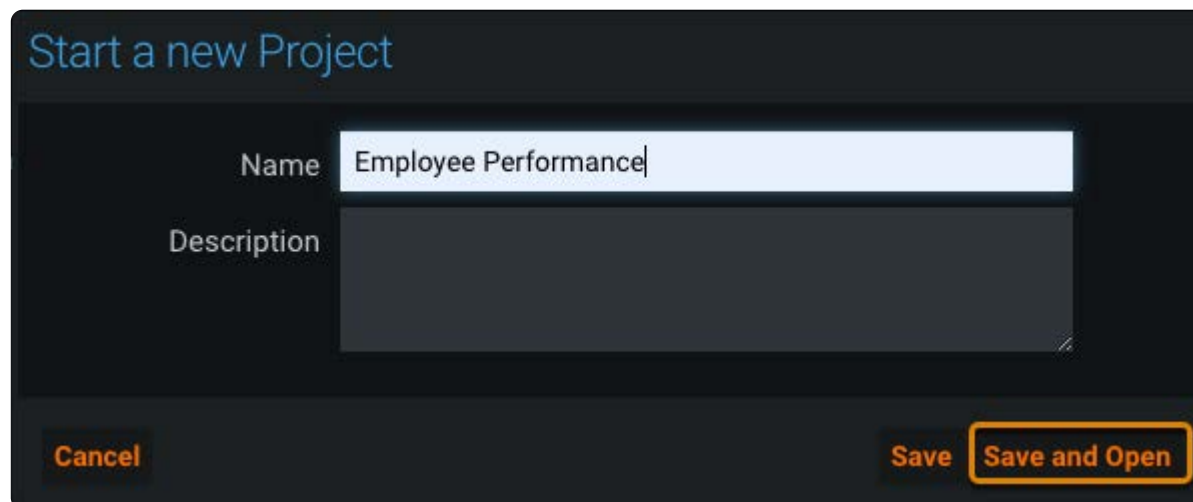
Start a new project from the Projects page

Rather than starting a project from the library, you can instead start from the **Projects** page:

1. Select **Projects** on the top left.
2. On the top of the **Library** page, click **+ add**.



3. Enter the project **Name** and an optional **Description**.



4. Click **Save and Open**.

Prep your data

Once you have started your project, you can begin prepping your data on the project preparation page.

1. Use the project **Tools** bar on the left to clean and transform your data:

	Sources	race	gender	age	weight	readmitted	admission_type_id	discharge_disposition_id
1		Caucasian	Female	[50-60]	?	false	Elective	Discharged to home
2		Caucasian	Female	[20-30]	[50-75]	false	Urgent	Discharged to home
3		Caucasian	Male	[80-90]	?	true	Not Available	Discharged/transferred to home with home health service
4		AfricanAmerican	Female	[50-60]	?	false	Emergency	Discharged to home
5		AfricanAmerican	Female	[50-60]	?	false	Emergency	Discharged to home
6		Caucasian	Male	[70-80]	?	false	Elective	Discharged to home
7		Caucasian	Female	[60-70]	?	false	Elective	Expired
8		Caucasian	Female	[50-60]	?	false	Emergency	Discharged to home
9		Caucasian	Male	[50-60]	?	false		Discharged to home
10		Caucasian	Male	[60-70]	?	true	Elective	Discharged to home
11		Caucasian	Male	[70-80]	?	false	Urgent	Discharged/transferred to a long term care hospital.
12		AfricanAmerican	Male	[80-90]	?	false	Emergency	Discharged/transferred to SNF
13		Caucasian	Female	[80-90]	?	false	Urgent	Discharged to home
14		Hispanic	Male	[60-70]	?	true	Emergency	Discharged/transferred to SNF
15		AfricanAmerican	Male	[70-80]	?	false	Emergency	Discharged/transferred to SNF
16		Caucasian	Male	[70-80]	?	false	Emergency	Discharged to home
17		AfricanAmerican	Male	[30-40]	?	false	Emergency	Discharged/transferred to another type of inpatient care institution
18		Caucasian	Male	[60-70]	?	true	Emergency	Discharged to home
19		AfricanAmerican	Male	[70-80]	?	true	Emergency	Discharged/transferred to another type of inpatient care institution
20		Caucasian	Male	[70-80]	?	true		Not Mapped
21		Caucasian	Female	[40-50]	?	true	Elective	Discharged to home
22		Caucasian	Female	[60-70]	?	false	Elective	Discharged to home
23		AfricanAmerican	Female	[70-80]	?	true	Emergency	Discharged to home
24		Caucasian	Female	[80-90]	?	false	Urgent	Discharged to home

See [Work with project tools](#) for detailed instructions.

2. Select the menu at the top of each column to apply column operations:

A-Z Admission Source		123 Days in Hospital	A-Z Payer Code
NA	FILTER values	7	MC
Physician Referra	SORT by ascending ▲	4	?
NA	by descending ▼	3	?
Physician Referra	CHANGE into ... »	14	MC
Transfer from an	COLUMN split	2	MC
Emergency Room	find + replace	6	?
Transfer from a S	duplicate	6	MC
Emergency Room	rename...	3	?
Emergency Room	fill... »	10	?
Emergency Room	WHITESPACE trim leading and trailing	6	?
Emergency Room	collapse consecutive	8	?
NA	OTHER cluster + edit...	1	?
Physician Referra		5	BC
Emergency Room		5	SP
Emergency Room		2	MC

See [Work with column data](#) for detailed instructions.

3. To view, rearrange, and mute your data prep steps, you can use the Steps tool:



See [Work with steps](#) for detailed instructions.

Publish an AnswerSet

When you're ready to save and share the data you prepped, you can publish it to the library as an *AnswerSet*. An AnswerSet is like a dataset but it is the published result of your data prep. Once published, you can reuse the AnswerSet in other projects or export the AnswerSet to share with other applications.

To publish an AnswerSet for a project:

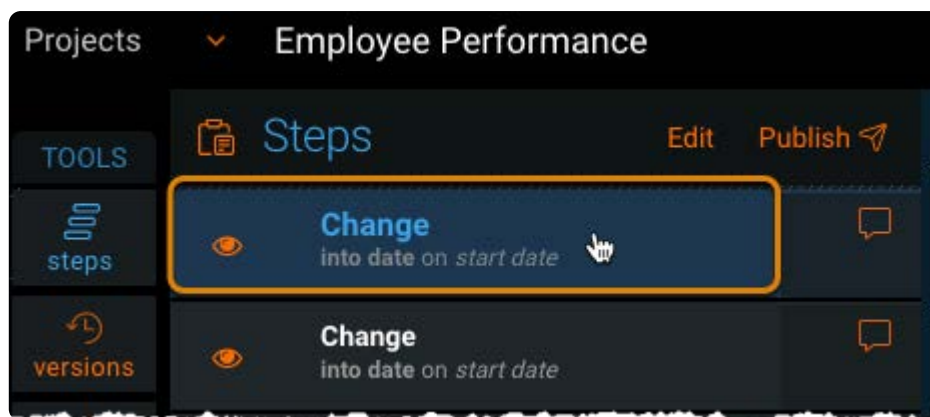
1. Click **steps** in the **Tools** bar.

The **Steps** pane opens.

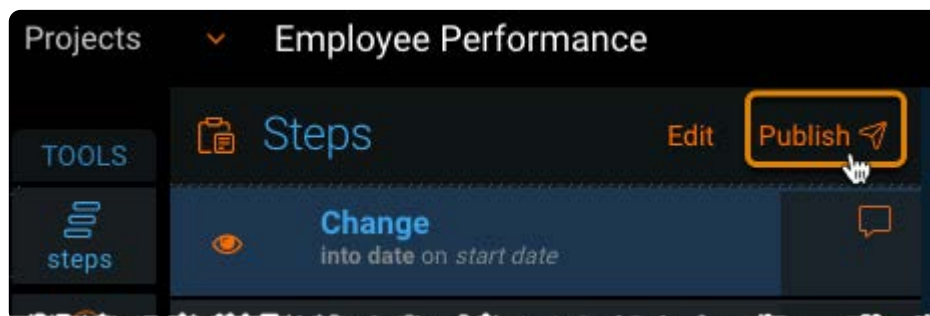


2. Click the step you want to publish an AnswerSet from.

Data Prep defaults to the last step in the project, which is the step at the top.

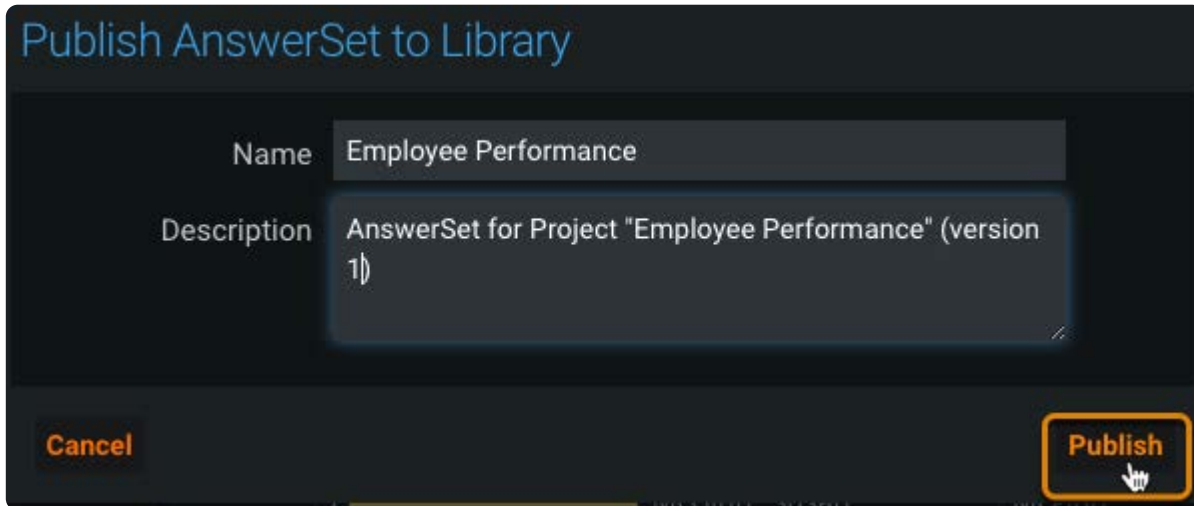


3. At the top of the **Steps** pane, click **Publish**.



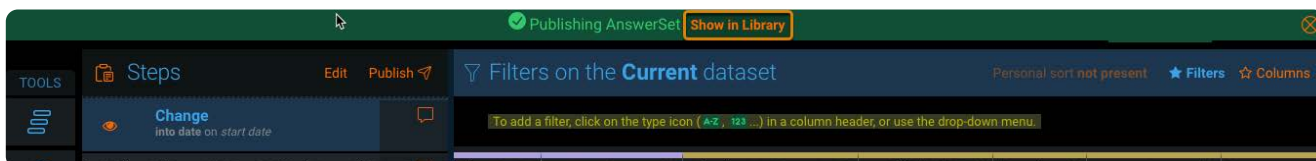
The **Publish AnswerSet to Library** window appears.

4. Enter a name for the AnswerSet in the **Name** field and an optional **Description**, then click **Publish**.

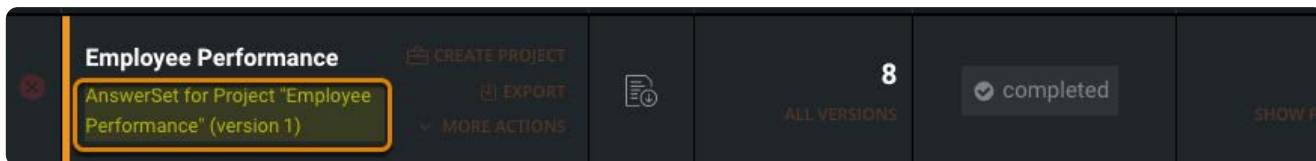


Data Prep publishes the AnswerSet to the library. The "Publishing AnswerSet" message appears.

5. Click **Show in Library** to view the AnswerSet in the library.



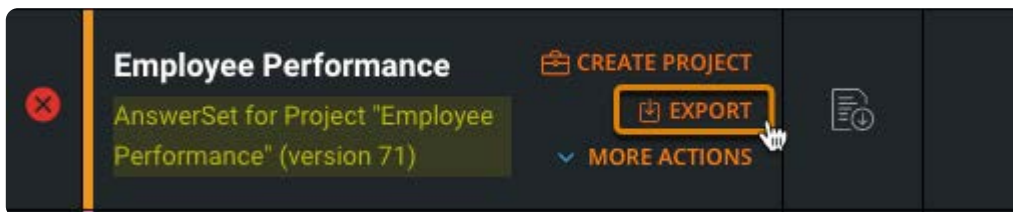
The AnswerSet includes the steps up to and including the step you selected.



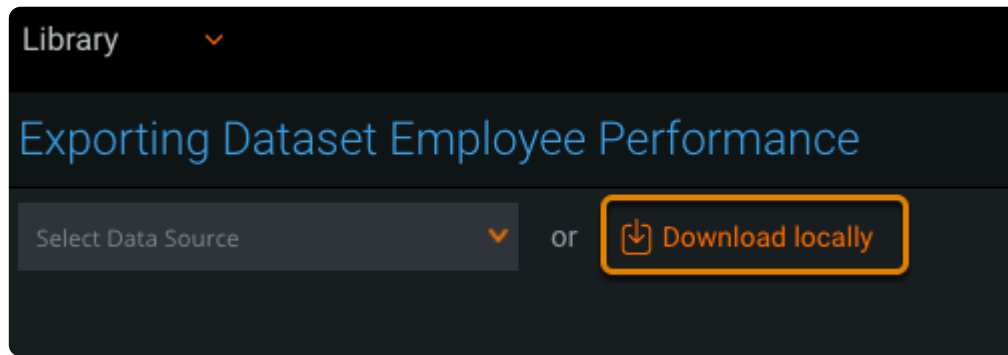
Export your prepped data

You can export datasets and AnswerSets locally or to a connected data source. These steps show how to download a local copy of a previously published AnswerSet.

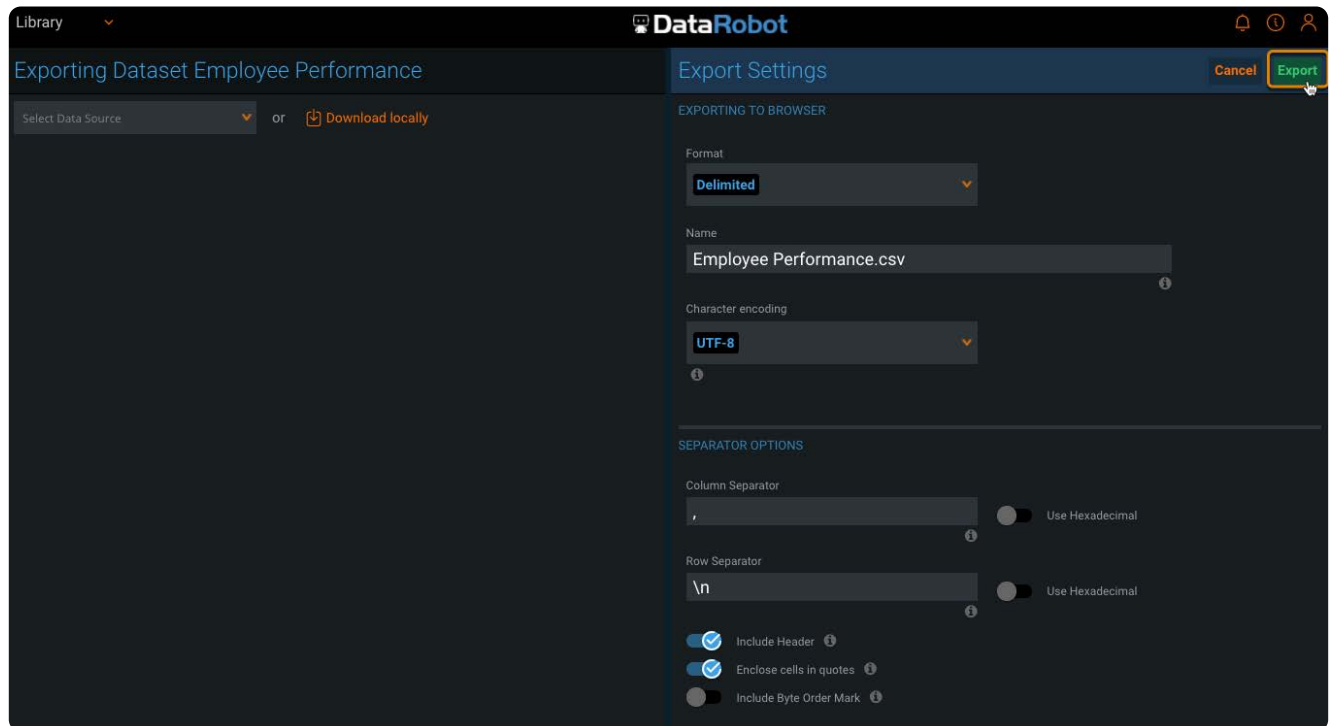
1. On the **Library** page, hover your mouse over the AnswerSet you want to export and click **Export**.



2. In the **Exporting** page, click **Download locally**.



3. In the **Export Settings** page, click **Export**.



The AnswerSet is downloaded to your computer as a CSV file. The **Export Logs** page appears.

Library ▼ DataRobot 🔔 🕒 👤

Exports 259 items in list 🔍

Datasets

Export Logs Last 7 days Last 30 days Yes No Mine

Finished Time Completed? Owner

Data Sources

Name	Dataset Name	Export Status	Export Destination	Created
Employee Performance.csv	Employee Performance	Complete	Local Download	May 21, 2021 1:22 pm by karen.germond@datarobot.com
Utah Housing Price Listings	Utah Housing Price Listings	Complete	AI_Catalog_FR	May 21, 2021 1:06 pm by felice.rando@datarobot.com
test AMAT	1.2 Demo:Lending Club_2 (Predict flow)	Complete	DataRobot	May 20, 2021 11:21 am by matthew.cohen@datarobot.com
Utah Housing Price Listings (Final Dataset)	Utah Housing Price Listings (Final Dataset)	Complete	AI Catalog AK Demo	May 19, 2021 12:58 pm by andrea.kropp@datarobot.com
LendingClub_raw.json	LendingClub_raw.csv	Complete	Local Download	May 18, 2021 9:37 pm by aman.sharma@datarobot.com
Hospital Admissions_raw.csv	Hospital Admissions_raw	Complete	Local Download	May 17, 2021 11:03 am by benjamin.miller@datarobot.com
LendingClub_raw.csv	LendingClub_raw.csv	Complete	Local Download	May 14, 2021 10:53 am by oleg.zarakhani@datarobot.com
US_Zip_Code_Validation_Ranges	US_Zip_Code_Validation_Ranges.csv	Complete	DataRobot AI Catalog	May 12, 2021 4:13 am

See [Export datasets](#) for details.

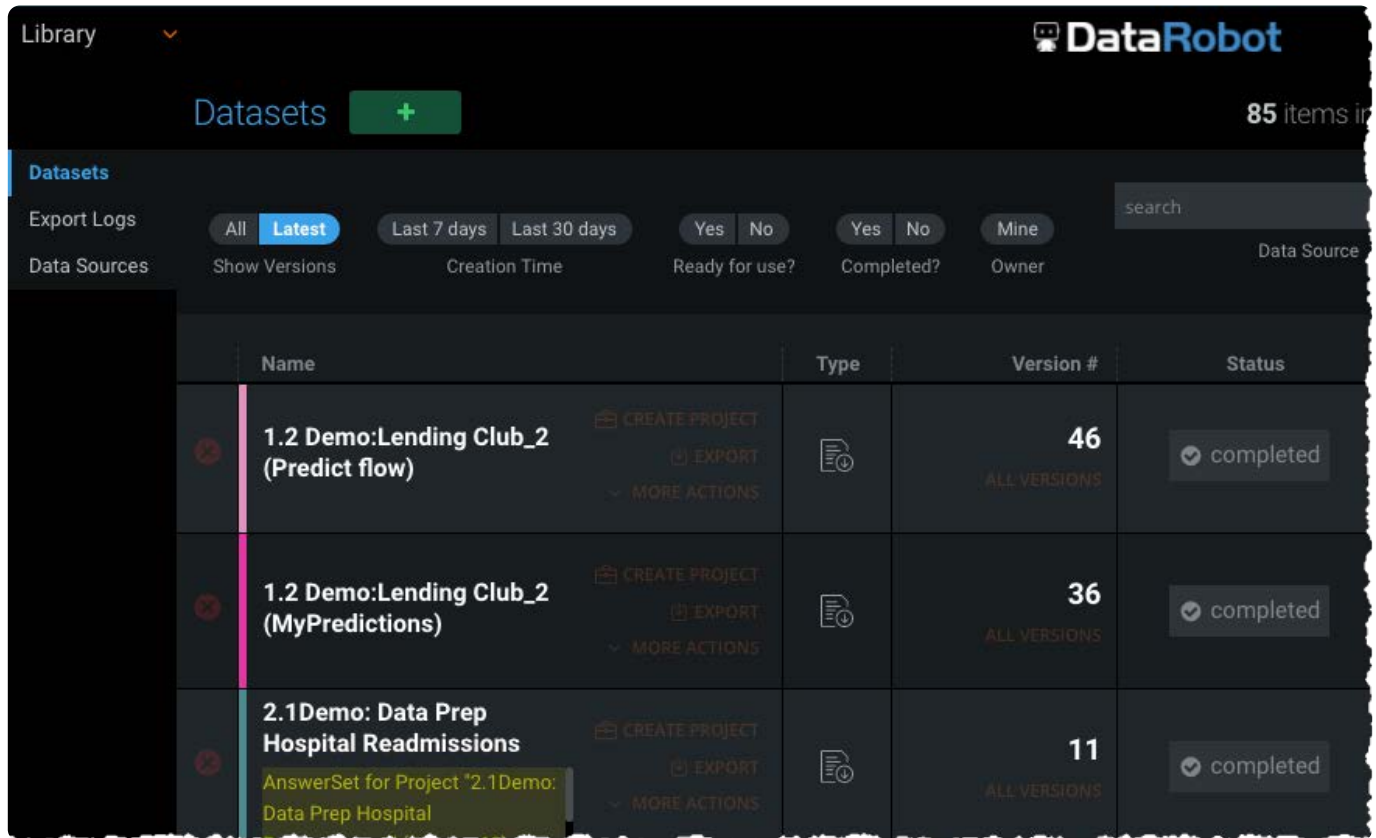
Tour the basics of Data Prep

In this topic, you'll take a look at the main components of the Data Prep application.

Library

The Data Prep library is where you:

- Add and manage datasets.
- Publish AnswerSets—your prepped datasets.
- Set up datasets for automation.
- Add new versions.
- Create profiles for your datasets.
- View warnings or errors that occur when you import datasets.

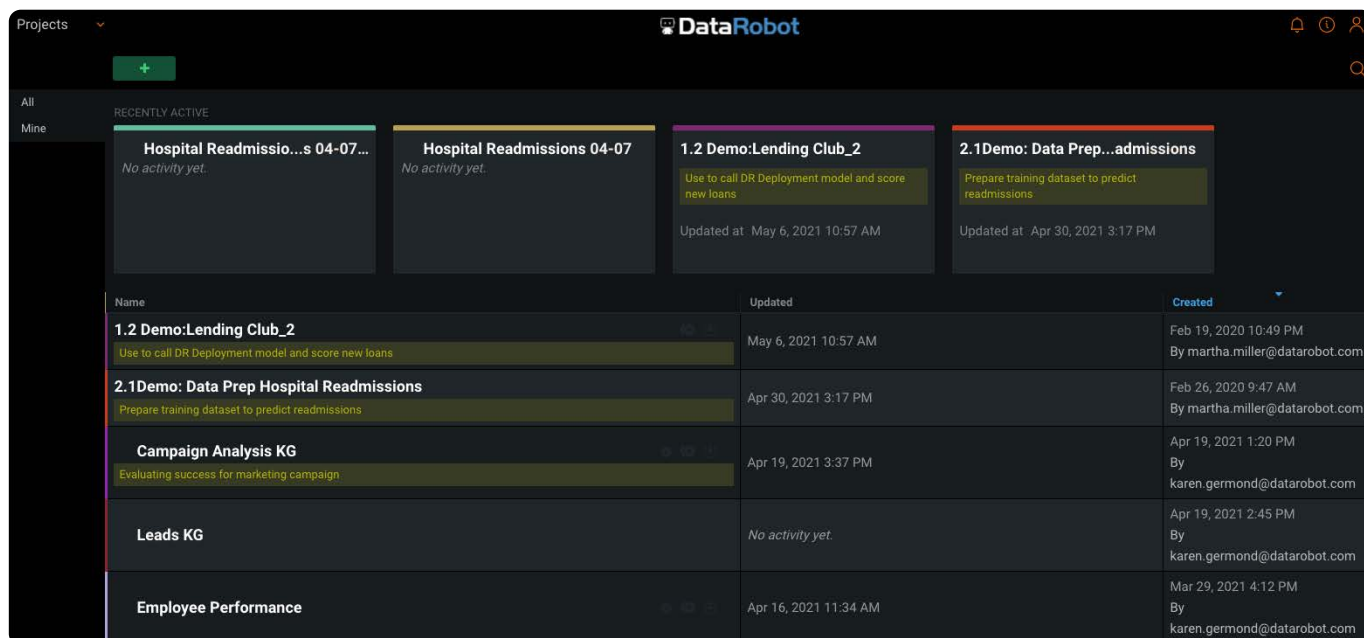


You import a dataset into the library, then you can begin prepping your data in a *Project*. When you finish prepping your data, you can publish it back to the library as an *AnswerSet*—a published dataset.

For a deeper look, see [Data Prep library](#).

Projects

The **Projects** page lists all projects that you have permission to view.



The screenshot shows the DataRobot Projects page. At the top, there's a 'Projects' header with a dropdown arrow and the DataRobot logo. Below the header, there's a sidebar with 'All' and 'Mine' options. The main content area is titled 'RECENTLY ACTIVE' and displays four project cards. Each card shows the project name, a description, and the last updated time. Below the cards is a table listing all projects with columns for Name, Updated, and Created.

Name	Updated	Created
1.2 Demo:Lending Club_2 Use to call DR Deployment model and score new loans	May 6, 2021 10:57 AM	Feb 19, 2020 10:49 PM By martha.miller@datarobot.com
2.1Demo: Data Prep Hospital Readmissions Prepare training dataset to predict readmissions	Apr 30, 2021 3:17 PM	Feb 26, 2020 9:47 AM By martha.miller@datarobot.com
Campaign Analysis KG Evaluating success for marketing campaign	Apr 19, 2021 3:37 PM	Apr 19, 2021 1:20 PM By karen.germond@datarobot.com
Leads KG	No activity yet.	Apr 19, 2021 2:45 PM By karen.germond@datarobot.com
Employee Performance	Apr 16, 2021 11:34 AM	Mar 29, 2021 4:12 PM By karen.germond@datarobot.com

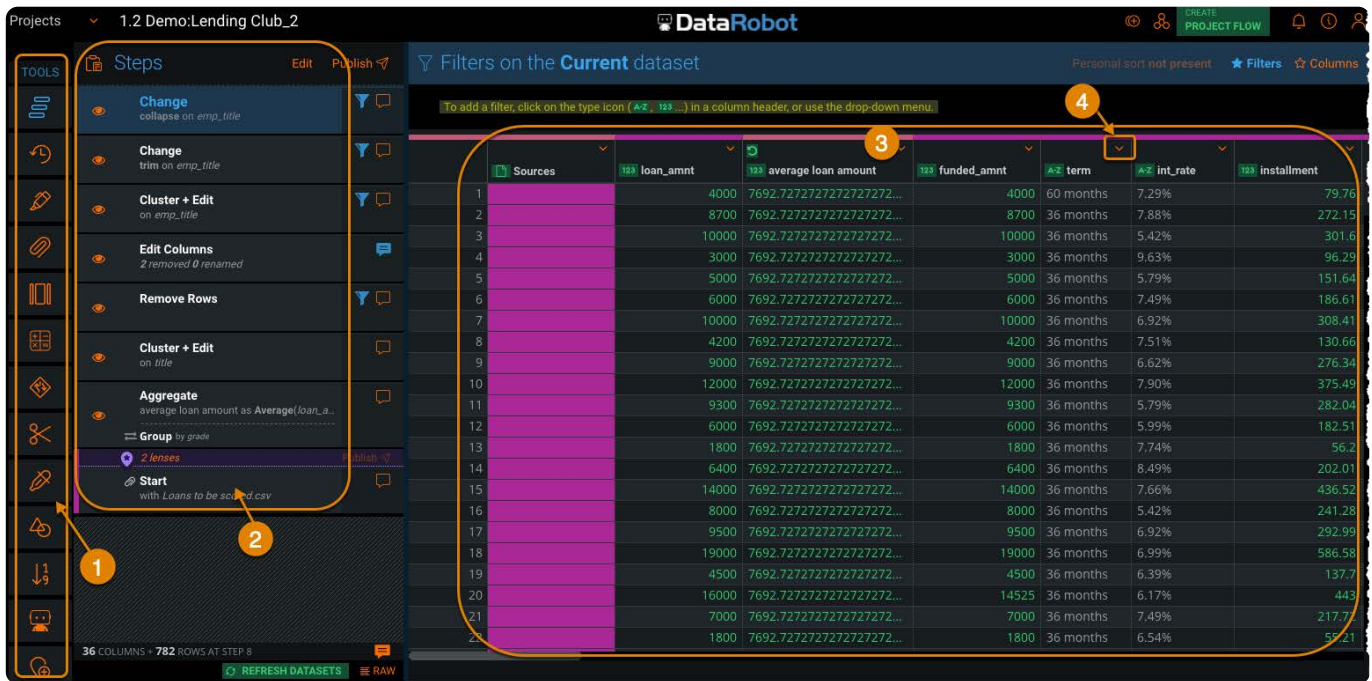
You add a dataset to a project, then you can explore the dataset and clean, transform, or combine it with other data.

You can publish changes to your library as an AnswerSet, which you can export or use within another project.

For a deeper look, see the [Data Prep projects page](#).

Project preparation

You open a project by clicking it on the **Projects** page or by starting a new project from the library. Once you open a project, you can begin to prep your data.



ELEMENT	DESCRIPTION
---------	-------------

- 1
Tools bar

On the left, you can access the project **Tools** bar that you will use to prep your data.
- 2
steps tool

The **Steps** tool saves each operation you perform so that you can replay, mute, and rearrange the steps.
- 3
Display Preview pane

Your data displays in the **Data Preview** pane.
- 4
column operations menu

Above each column, you can access the column operations to update the columns.

For a deeper look, see [Data Blue Project Preparation page](#).

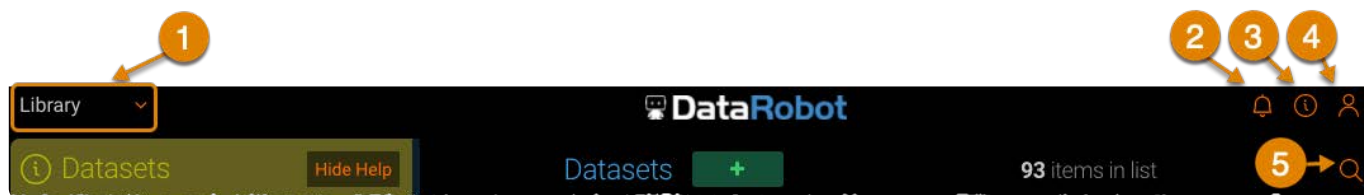
Data

You can import data from local files on your computer or from connected data sources. Your Data Prep system administrator must configure the data sources before you can import from them. Some examples of connected data sources are:

- Cloud storage like Amazon S3
- The Hadoop Distributed File System (HDFS)
- Relational databases like MySQL
- Secure File Transfer Protocol (SFTP)

Data Prep navigation

The Data Prep header provides navigation, help, and account management functions:



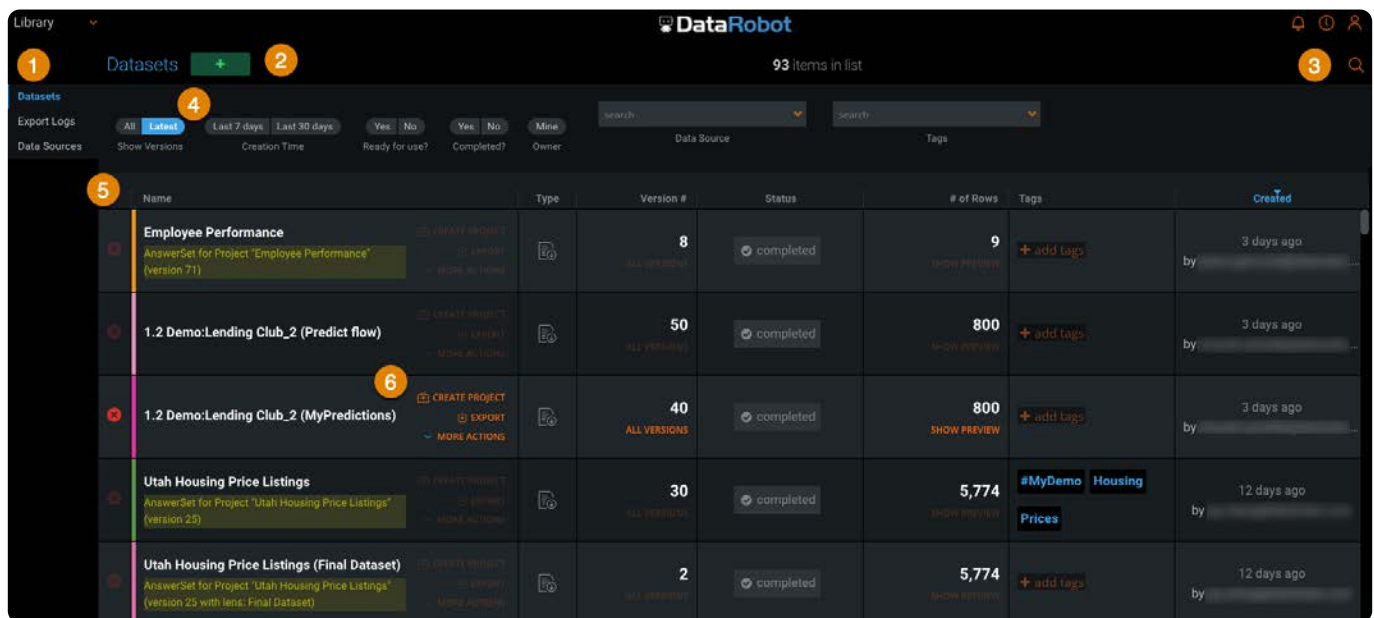
ELEMENT	DESCRIPTION
1	<p>Navigation menu</p> <p>Navigate between the Data Prep pages:</p> <ul style="list-style-type: none">• Library: Access your imported and published data.• Projects: Prepare your data.• Admin: Make connections to data sources and control users' permissions.• Project Flows: Automate data prep processes. <p>Note: The pages available to each user are based on the user's permissions.</p>
2	<p>Notification icon</p> <p>Indicates when Data Prep generates a warning or error. If highlighted, mouse over the icon to view the message.</p>
3	<p>Help</p> <p>Get Data Prep help.</p>
4	<p>User menu</p> <p>Access account-specific options like updating your password or logging out. You can also generate Tokens that are used to manage application access and authorization. Your Data Prep System Administrator will let you know when you need to generate tokens.</p>
5	<p>Search</p> <p>Page-specific search. For example, on the Library page, you can search for datasets and on the Projects page, you can search for projects.</p>

Data Prep library

On the **Library** page, you can add new datasets and manage existing datasets, including Data Prep AnswerSets that you publish from your projects. In the library, you can also export datasets, set them up for automation, add new versions, create profiles for your datasets, and view any warnings or errors that occurred when a dataset was imported.

The following tables describe the library layout, as well as the actions you take in the library to work with your datasets.

Library layout



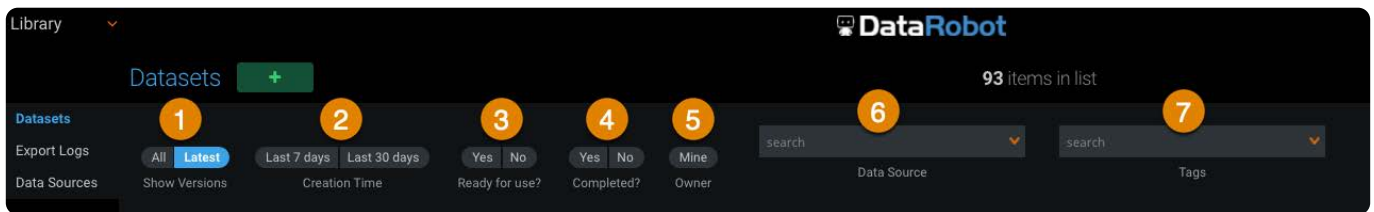
The following table describes the sections of the **Library** page.

ACTION	DESCRIPTION
<p>1 Library tabs</p>	<p>Select a tab:</p> <ul style="list-style-type: none"> • Datasets: Manage datasets. • Export Logs: View the logs generated during exports. • Data Sources: Add new data sources. Import and export from data sources that have been added to the library.

ACTION	DESCRIPTION
<p>2 Add a new dataset</p>	<p>To add a new dataset to your library, click Datasets + on the top left. On the Select Datasets page, import one or more datasets by selecting data sources or local datasets. If there are any errors during the import, a red warning icon displays adjacent to the dataset's listing on the page. For more information on any errors, mouse over the dataset's name and click Edit Details.</p>
<p>3 Search for a dataset on the page</p>	<p>To search for a dataset in your library, click the magnifying glass icon on the top right. In the field that displays, begin typing the name of the dataset you want to locate. Potential matches display as you continue to type.</p>
<p>4 Filter the datasets displayed on the page</p>	<p>Filter the list of datasets by categories such as version, creation time, and owner.</p>
<p>5 Sort columns</p>	<p>The Library page lists the filtered datasets. The library columns provide attributes for each dataset, including the type, version, status, number of rows, tags, and data created and by whom.</p>
<p>6 Dataset actions</p>	<p>You can perform actions on the datasets such as creating projects and exporting datasets. To delete a dataset, click the red X icon to the left of the dataset.</p>

Library filters

Use the filters at the top of the page to filter the list of datasets displayed on the page.



The following table describes your options for filtering datasets.

ACTION	DESCRIPTION
<p>1 Show Versions</p>	<p>Toggle to display all versions of every dataset and AnswerSet or only the latest version of each.</p>
<p>2 Creation Time</p>	<p>Select the last seven or 30 days.</p>

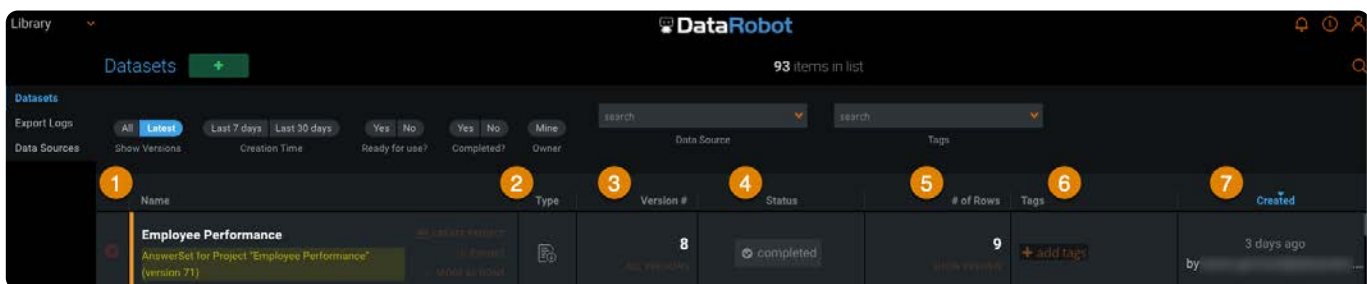
ACTION	DESCRIPTION
3 Ready for use?	This filter only displays when interactive mode is enabled for your Data Prep projects. It allows you to quickly see which datasets have finished loading their interactive portions and are ready for use in a project.
4 Completed?	Displays all datasets that have successfully finished importing into the library.
5 Owner	Displays the datasets and AnswerSets that you have imported and created from Data Prep projects.
6 Data Source	Click this field to display all of the data sources used to import datasets. You can select more than one data source by continuing to click and select.
7 Tags	Tags are descriptive words that enable you to organize your datasets. Click the Tags field to display all of the tags currently assigned to datasets in the library. To locate a dataset by a specific tag, type the tag name and click Enter . If you add multiple tags, only the datasets containing <i>all</i> of the search tags are returned as matches. To add a new tag for a dataset, hover over the dataset and click inside the add tags field for that dataset.

Library columns

You can create a new column display order for the **Library** page by clicking a column's header and dragging it to a new location. To sort the library list by a particular column, click that column name. To sort on multiple columns, hold **Shift** and click additional columns.

Note

When you reorder items on the **Library** page, the changes are temporary and not retained when you leave the **Library** page or refresh your browser.



The following table describes the columns on the **Library** page.

ACTION	DESCRIPTION
--------	-------------

1

Name

Displays the name of the dataset when it was imported into the library. To change the name, hover over **More Actions** for the dataset and click **edit details**. The **General** page for the dataset displays. You can change the dataset name and update metadata fields.

2

Type

Allows you to quickly identify which datasets in your library are AnswerSets that were created from Data Prep projects. If the interactive mode feature is enabled for your Data Prep projects, AnswerSets are represented with the partial icon to indicate that you were working in interactive mode when it was created.

3

Version #

Displays the number of versions for each dataset or AnswerSet. If there is more than one version, click **All Versions** for the dataset to view all versions. To return to viewing all datasets, click **All Datasets** on the top left.

Keep in mind, if you filter tags while drilled into the **All Versions** page, your search only applies to the **All Versions** page and not the entire **Library** page.

Version numbers do not necessarily correspond to the actual number of those datasets in the library.

Following are conditions under which a version number will not match the exact number of those datasets in the library:

- When an import is canceled before it completes, a version number is automatically generated and subsequent imports will simply be incremental version number additions.
- When a particular version of a dataset is deleted, the version numbers for the remaining datasets are not decremented.

ACTION	DESCRIPTION
--------	-------------

4

Status

Describes a dataset's load status as it's being imported into the library. In most cases, the status quickly progresses to "completed." However, for larger datasets, you will see interim states that indicate that your dataset is continuing to successfully import. The interim states you may see also depend on:

- **Whether the row count of the dataset can be predetermined prior to import:** In most cases, Data Prep knows the number of rows in a dataset before the import process even begins. However, there are cases where the count cannot be predetermined—for example, imports from Salesforce and queries on JDBC data sources.
- **Whether interactive mode is enabled for your projects:** When interactive mode is enabled, you'll notice the status icon has two concentric circles. The inner circle represents the interactive portion of your dataset. When the interactive portion is ready to be used in a project, the inner circle becomes a green check mark. The outer circle will then begin to fill green as the remainder of the dataset continues to load into the library. If any errors occur while importing the interactive portion or the remainder, a red warning icon displays in the respective concentric circle to indicate which part of the dataset failed to import into the library. See [Loading states](#) for examples of loading states. See [Failure states](#) for examples of the failure states.

You may see a "Pending" state in this column if you did not finish selecting the parsing options for the dataset. In this case, you will also see a **Click to Finish** button in the **Created** column. Click the button to open the **Import** page and finish the import.

5

of Rows

Displays the number of rows in a dataset. You can preview rows from a dataset by moving your mouse over the dataset and clicking the **show preview** link that displays in this column. When a dataset is currently in the import process and the row count is predetermined, the number displayed in this column continues to increase until the import is finished. If the dataset fails to import successfully, the number of rows that *successfully imported* are listed in this column. In this case, **show preview** displays a preview of those rows.

6

Tags






Tags are labels that you can add to your datasets to help organize your data. To add tags to a dataset, click in the Tags column for that dataset, type a tag name and click the **Add** link that displays or press the **Enter** key.

7




Created

Displays the user who imported the dataset and when it was imported. You may see a **Click to Finish** link in the column. This indicates the import was never initiated because the parse options were not finalized. Click this link to return to the **Import** page and finish the import process for the dataset.

Loading states

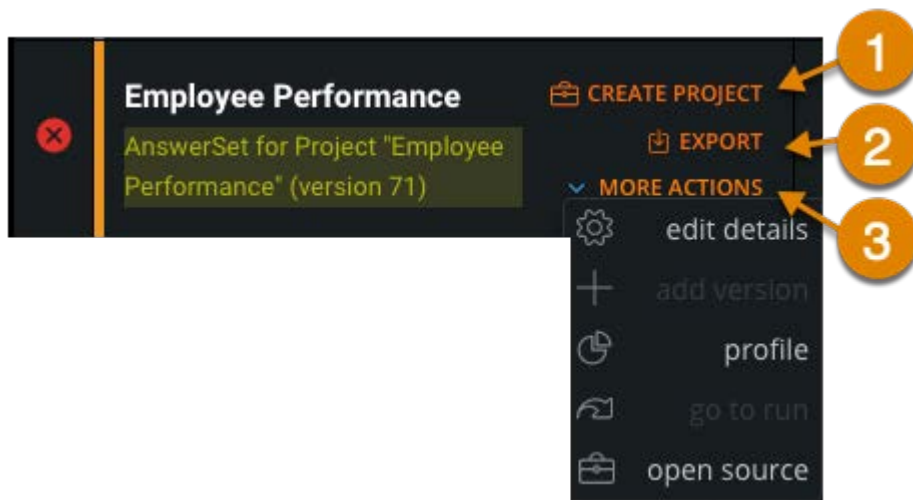
ICON	DESCRIPTION
	Icon displayed when interactive mode is <i>not</i> enabled and row count can be determined.
	Icon displayed when interactive mode is <i>not</i> enabled and row count cannot be predetermined.
	Icon displayed when interactive mode is enabled and row count cannot be predetermined.
	Icon displayed when interactive mode is enabled and row count cannot be predetermined.
	Icon displayed when loading is complete.

Failure states

ICON	DESCRIPTION
	Interactive mode not enabled: Dataset failed to import.
	Interactive mode: Interactive portion did not successfully import.
	Interactive mode: Interactive portion successfully completed but remainder of dataset failed to successfully import.

Actions you can take for a dataset

Three links that appear when you hover over a dataset provide you with the options you can take for that dataset.



The following table describes the actions you can perform on a dataset on the **Library** page.

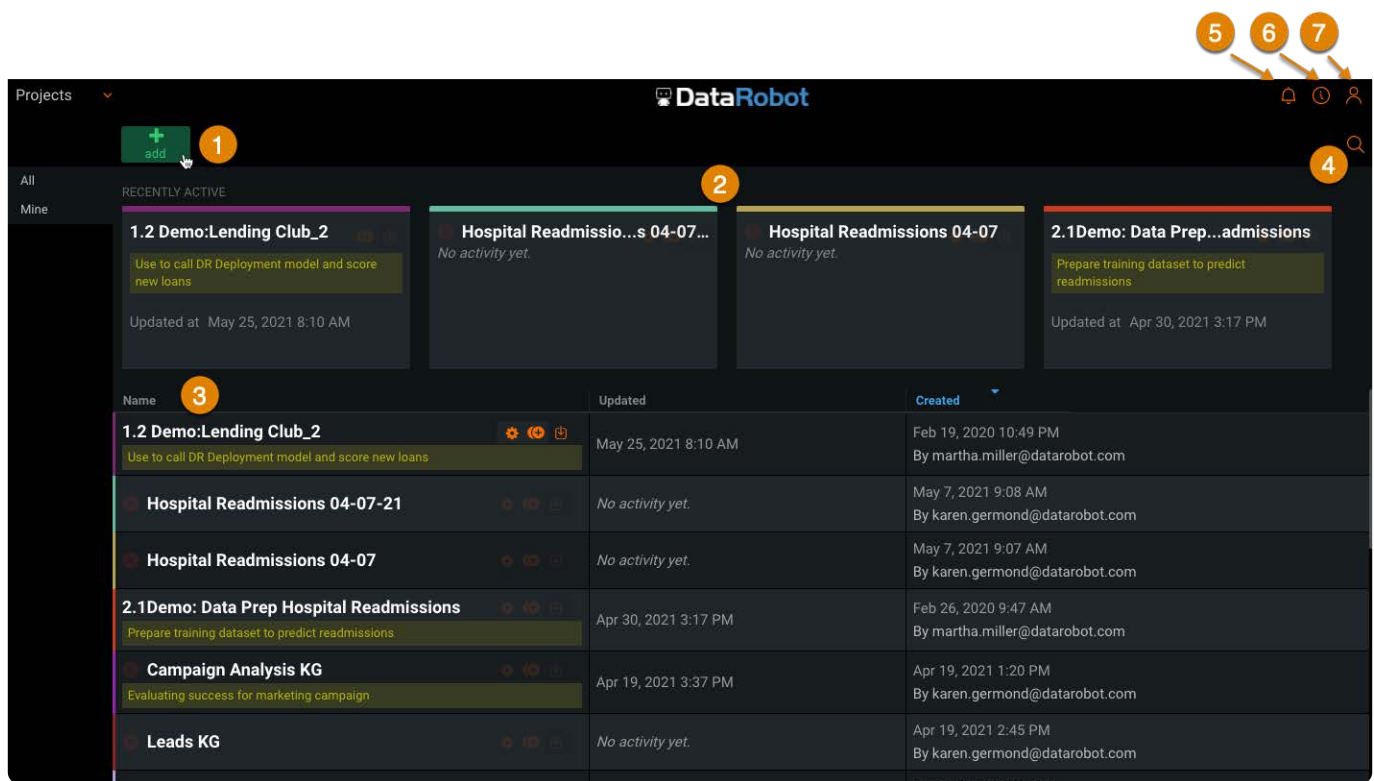
ACTION	DESCRIPTION
1 Create Project	Create a new project using the dataset as your base dataset.
2 Export	Export or download a dataset locally.
3 More Actions	Provides additional options, depending on the features that are enabled for your Data Prep application: <ul style="list-style-type: none">• edit details: Opens the dataset's General page. This is where you can update the dataset's name and metadata. This is also the page where you view warnings or errors that may have occurred during import. Datasets with warnings or errors are easy to locate in the list—they are flagged with a warning icon adjacent to the dataset name, the row color for the dataset is red, and the Status icon indicates a failure state.• add version: Add a new version of the existing dataset without overwriting the current version.• automate: Automate the dataset (if automation feature is enabled).• profile: Profile the dataset (if profiling feature is enabled).• open source: For any AnswerSet, open a project at the precise Step from which that AnswerSet was created.

Data Prep projects

In Data Prep, you explore and prepare your data in a *Project*. This topic walks you through the top-level **Projects** page, as well as the **Project Preparation** page for an individual project.

Projects page

The **Projects** page is the home page where you land after logging into Data Prep.



Note that if you want to open your [Library page](#) where your datasets are saved, click the menu in the top left and select **Library**.

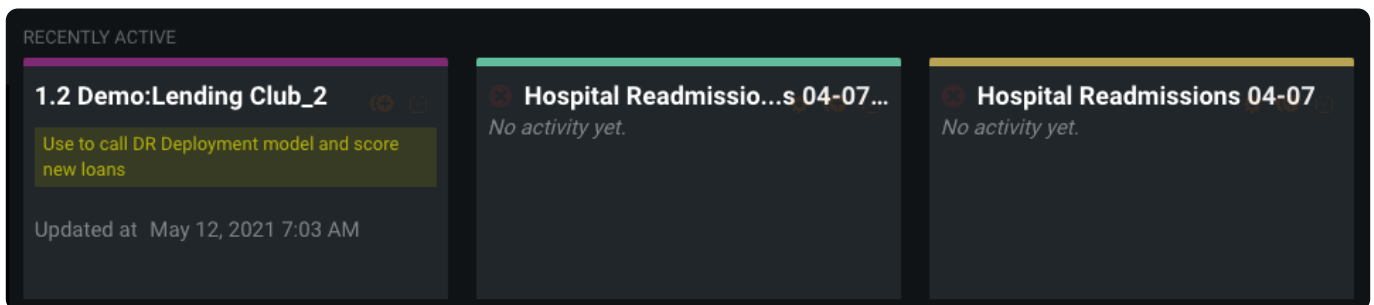
The following table describes the sections of the **Projects** page.

ACTION	DESCRIPTION
1 + Add Project	Create a new project. Provide a unique name and, optionally, a description. Once you create the project, you can import data .
2 Recently Active	Recently active projects are listed at the top of the page. Click a project's name to open it.

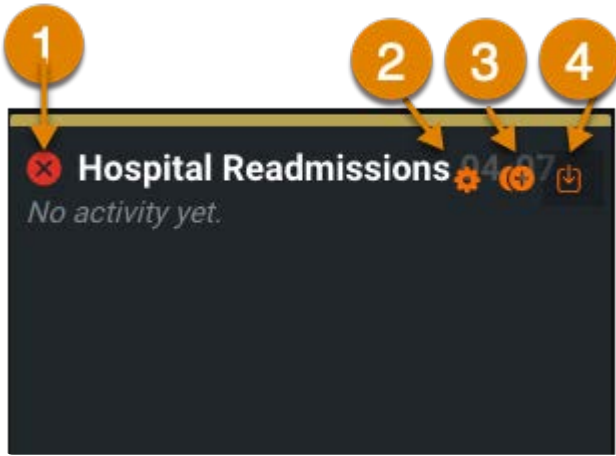
ACTION	DESCRIPTION
<p>3 All Projects inventory</p>	<p>All projects are listed here. Note they are ordered, by default, based on the most recently updated projects. You can order the list based on project names by clicking the Name column header. You can order the list based on the date the project was updated or created by clicking the column headers for Updated or Created.</p>
<p>4 Search</p>	<p>Click the magnifying glass icon to search for projects by name.</p>
<p>5 Notification icon</p>	<p>Indicates when the application has generated a message or error. When the icon is highlighted, mouse over it to view any messages.</p>
<p>6 Help</p>	<ul style="list-style-type: none"> • Show/Hide: Opens and closes the application's in-line help menu. • What's New: Opens a summarized list of all new features in the current application release. • Getting Started: Opens the new user guide. • Help Shelf: Opens the home page for all help documentation. • Feedback: Opens email to contact DataRobot's Customer Success team.
<p>7 User menu</p>	<ul style="list-style-type: none"> • My account: Displays your account information and provides the option to change your password. • Tokens: Generates tokens that are used to manage application access and authorization. Your DataRobot System Administrator will let you know when you need to generate tokens. • About: Displays details of the current application version numbers. • Logout: Logs you out of the application.

Recent projects

Your recent projects appear as tiles under **Recently Active**:



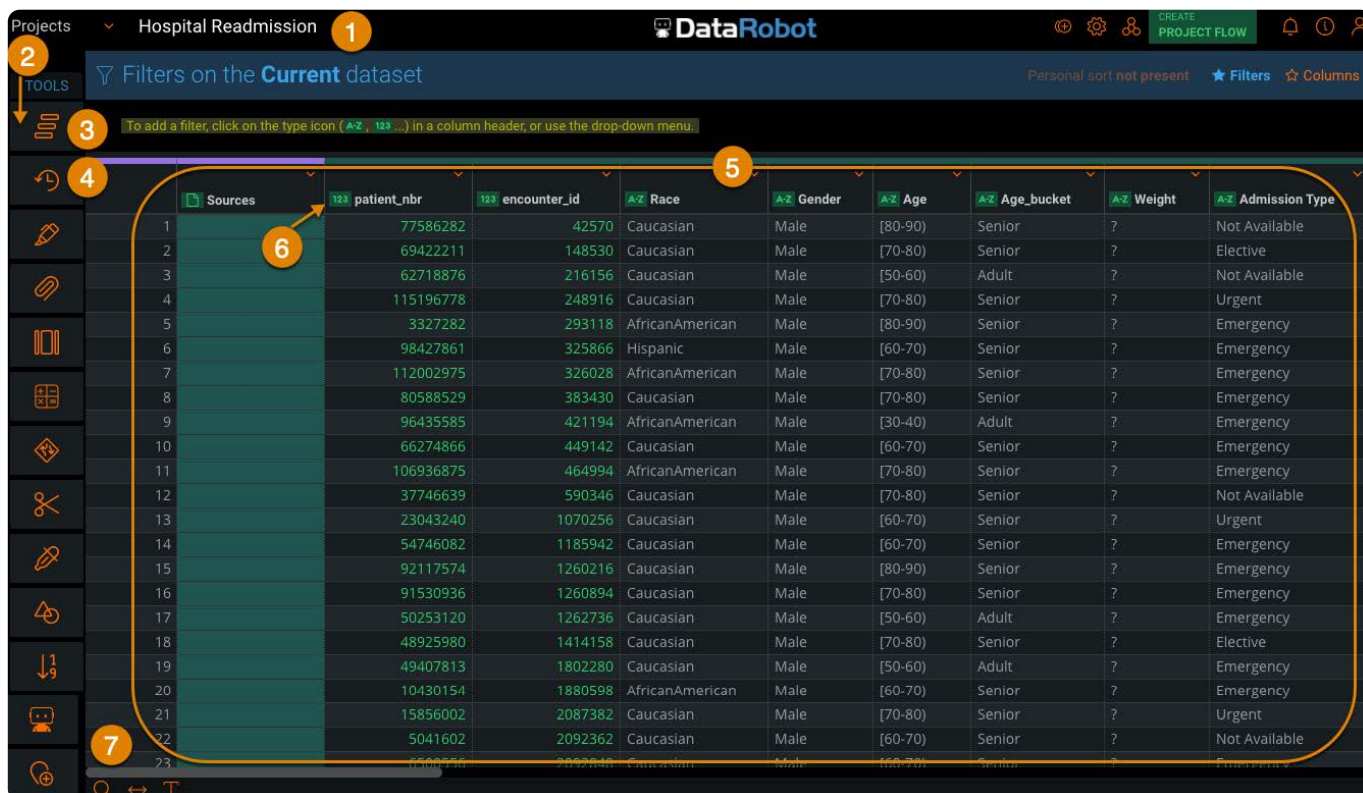
Each project tile has the following controls:



ACTION	DESCRIPTION
1 Delete Project	Click the red X icon on the upper left to delete a project. Click Okay in the confirmation box. You can only delete the projects you created.
2 Edit name or description	Edit a project's name and description.
3 Copy latest version of this Project	Duplicate the project.
4 Download Project Report	Download a Project Report file (in Word format). The Report file contains the project name and metadata, the column data, and the transformation steps used to prep the data.

Project Preparation page

You open a project by clicking on the dataset or its **Recently Active** tile on the **Projects** page, or by creating a project on the **Library** page. The **Project Preparation** page displays and you can prep your dataset.



The following table describes the sections of the **Project Preparation** page.

ACTION	DESCRIPTION
--------	-------------

- 1 Project name** The name you gave your project will display here.
- 2 Tools bar** Use the [project tools](#) to clean, shape, combine, and ultimately prep your data.
- 3 Steps tool** Every action you perform while prepping your data is logged as a *step*. The [Steps tool](#) allows you to:

 - View your steps in order.
 - Mute a step.
 - Edit what happens during a step.
 - Rearrange the order of your data preparation steps.
 - Delete steps.
- 4 Versions tool** Any time you save your project, a new version is created. The [Versions tool](#) lets you access previous versions of your project.
- 5 Data Preview** View changes to your data as you prep it.

ACTION	DESCRIPTION
6 Column menu	Click to open the column menu. Use column operations to clean and standardize your data.
7 Grid tools and status updates	Grid tools allow to you locate specific columns in your dataset, specify column widths and adjust how cell text displays. Status updates display when transformations that affect the Data Preview grid or filters are in progress. Note that the number of tasks displayed in the update messages may dynamically change as an operation progresses toward completion.

Connect to data sources for Data Prep

These pages contain the configuration instructions for each connector that Data Prep supports. Before diving in, you can [read about setting up connectors in general](#).

Select a connector below to learn how to configure it:

CONNECTOR	DESCRIPTION
Amazon Athena	Connect to AWS Athena as an import source.
Amazon S3	Import and export data against Amazon S3 object storage.
Amazon Redshift	Connect to Amazon Redshift as an import and export source.
Amazon DynamoDB	Connect to an Amazon DynamoDB for Library imports.
Cloudera CDH5 HDFS	Connect to Cloudera CDH 5.16 Hadoop File System (HDFS) for import and export.
Cloudera CDH6 HDFS	Connect to an Hadoop File System (HDFS) cluster for imports and exports.
Cloudera CDH5 Hive	Connect to Cloudera CDH 5.16 Hive for import and export.
Cloudera CDH6 Hive	Connect to a Hive database for imports and exports.
Cloudera CDH6 Impala	Connect to an Impala database for imports and exports.
Databricks	Connect to Databricks for Library imports and exports.
DataRobot	Connect to DataRobot for Library imports and exports.
Google Analytics	Connect to Google Analytics for browsing and importing available data.
Google BigQuery	Connect to BigQuery for importing and exporting available data.
Google Cloud Storage	Connect to Google Cloud Storage (GCS) for browsing and importing objects.
Google Cloud SQL	Connect to Cloud SQL using the Data Prep JDBC Connectors.
Google Drive	Connect to Google Drive for browsing and importing available data.

CONNECTOR	DESCRIPTION
Google Sheets	The Google Sheets Connector has been deprecated. Use the Google Drive Connector , which has support for importing and exporting Google Sheets.
Hortonworks HDP2 HDFS	Connect to a Hortonworks HDP 2.6.5 Hadoop File System (HDFS) for import and export.
Hortonworks HDP2 Hive	Connect to a Hortonworks HDP 2.6.5 Hive for import and export.
Hubspot	Connect to HubSpot for browsing and importing available data.
IBM DB2	Connect to IBM DB2 using the Data Prep JDBC Connector.
IBM Netezza	Connect to IBM Netezza using the Data Prep JDBC Connector.
JDBC	Import and export data leveraging Java Database Connectivity (JDBC) Drivers. Typically this connector is leveraged for import/export against relational databases, but many applications offer JDBC drivers.
Jira	Connect to Jira for browsing and importing available data.
Marketo	Connect to Marketo as an import source.
MicroStrategy	Connect to a MicroStrategy server for Library imports and exports.
MongoDB	Connect to MongoDB for browsing and importing available data.
MS Azure Data Lake Store (ADLS)	Connect to Azure Data Lake Storage (ADLS) for Library imports and exports.
MS Azure Data Lake Store Gen2 (ADLS Gen2)	Connect to Azure Data Lake Storage Gen2 for import and export.
MS Azure SQL	Connect to Azure SQL using the Data Prep JDBC Connector.
MS Azure Synapse Analytics	Connect to an Azure Synapse Analytics for Library imports and exports.
MS Dynamics 365	Connect to a Microsoft Dynamics 365 resource for Entity Set imports.
MS Sharepoint	Connect to a SharePoint site for Library imports and exports of files and SharePoint Lists.

CONNECTOR	DESCRIPTION
MS SQL Server	Connect to Microsoft SQL Server using the Data Prep JDBC Connector.
MS Windows Azure Blob Storage (WASB)	Connect to Azure Blob Storage Account for Library imports and exports.
MySQL	Connect to MySQL using the Data Prep JDBC Connector.
Netsuite	Connect to NetSuite for browsing and importing available data.
Network Share (SMB / Samba)	Connect to a network share using Server Message Block (SMB) protocol for imports and exports.
Oracle	Connect to Oracle using the Data Prep JDBC Connector.
Oracle Marketing Cloud (Eloqua)	Connect to Oracle Marketing Cloud for library imports.
PostgreSQL	Connect to PostgreSQL using the Data Prep JDBC Connector.
PowerBI	Connecting to PowerBI is configured in PowerBI Desktop. Learn more about how to connect.
REST API	Connect to a REST API to import a REST Resource.
Salesforce Lightning	Connect to a Salesforce org as an import source.
Salesforce Marketing Cloud	Connect to Salesforce Marketing Cloud for browsing and importing available data.
SAP HANA	Connect to SAP HANA using the Data Prep JDBC Connector.
Spark SQL	Connect to Spark SQL for browsing, importing, and exporting available data.
SFTP	Connect to an SSH File Transfer Protocol (SFTP) Server for Library imports and exports.
Snowflake	Connect to Snowflake using the Data Prep JDBC Connector.
Tableau .Hyper	Connect to Tableau as an export destination.
Tableau .tde (deprecated)	Connect to a Tableau Server and Tableau Online as an export destination.

CONNECTOR	DESCRIPTION
Teradata	Connect to Teradata using the Data Prep JDBC Connector.
Thoughtspot	Export your Data Prep AnswerSets to ThoughtSpot.
Vertica (HP)	Connect to Vertica using the Data Prep JDBC Connector.
Zendesk	Connect to Zendesk for browsing and importing available data.

Data Prep Connector setup for Data Prep

What are Data Prep Connectors?

Every Data Prep story starts and ends with Connectors. Being able to do data preparation is only valuable if you can get the data you need to prep and then can send that data where you need it after it's been prepped. Data Prep Connectors are the tools for getting data into and out of Data Prep.

Benefits of Data Prep Connectors

Straightforward data access for business users

Accessing data on disparate systems isn't very complicated for coders—most databases, file stores, and web services have well-developed, code-friendly interfaces that adhere to industry standards.

Data integration is hard for non-coding users

DataRobot has tackled this problem and has opened up as many data sources as possible to non-coding users of DataRobot Data Prep. Our goal is that a business analyst (non-coding user) can access any data in the organization they are authorized to use.

Browsing vs. Querying

One core aspect of enabling non-coding users is the browsing interface. Where other data prep or ETL solutions rely on SQL queries, every data source in Data Prep can be browsed and data can be imported with clicks.

Control and Governance

The business environment is significantly more fluid than IT infrastructure typically accommodates, but still, certain people should only have access to certain information and should only be able to send that information to certain places. The Connector framework allows large and complex organizations to ensure users can access only the information granted to them and can be configured simply for smaller organizations where speed and self-service are a priority.

Setup of Data Prep Connectors

Three Layers of Configuration

When setting up a Connector, there are three hierarchical levels of configuration, from highest to lowest: “Connector,” “Data Source,” and “Session.” If a field is filled out at a higher level, it won’t need to be filled out again downstream. Some fields may be alterable at a later stage, but that varies greatly across the Connectors.

Connector configuration

This level is typically created and managed by an Admin or IT and it exists to:

- Make a given Connector available to specific groups of users.
- Allow an administrator to enter information that users won't know and/or that will be the same across all users/data sources that rely on the Connector Config.
- It also allows an Admin to keep sensitive information secure from users who shouldn't have access, e.g. an SSH Key.

Data Source configuration

This level is typically created and managed by either individual users or admins, depending on how access to the source system data is being managed and it exists to:

- Contain all persistent configuration not already captured at the Connector Config level.
- Typically, this includes everything except for user credentials supplied at runtime for a shared Data Source Config.

Session configuration

This level is almost exclusively managed by individual users or ignored if not required and it exists to:

- Capture information at runtime of import/export.
- Typically, this is limited to user credentials.

Sharing controls

- Connector & Data Source Configs can be shared with groups within your tenant.
- These sharing controls also allow you to specify if members of the specified groups can Read, Update, or Delete the configuration and whether the users may perform imports and/or exports with the configuration.

Example Setups

The following are a few examples of business situations and how the Connector Framework can be set up to accommodate the needs of each team.

Example 1:

Business Situation

IT-managed SFTP Server authenticated by "SSH Key with Passphrase" where the key and passphrase are held by IT and several teams will need access to different directories.

Setup

- Connector Config
- IT will create one Connector Config and fill out SFTP Host & Port, SSH Key & Passphrase.
- Sharing: None
- Data Source Config
- Create a new Data Source for each team, specify the appropriate Root Directory.
- Sharing: Share each fully-configured Data Source as Read-only with the corresponding team and allow imports & exports if appropriate.
- Session Config
- N/A

Benefits of this approach

- If the credentials change, they only need to be managed in one place.
- IT can manage credentials and keep them private from users.
- Each team has the access they need without having to manage access control on the data source itself.

Example 2

Business Situation

Admin managed Salesforce Org where each user should access only the information they have permissions for in Salesforce and each user will need to run automation jobs within Data Prep.

Setup

- Connector Config
- The Salesforce Admin will create one Connector Config and fill out all relevant information except for User & Password.
- Sharing: Share this Config with each relevant group as Read-only.
- Data Source Config
- Each user should create their own Data Source config and fill out just their credentials so their setup persists and can be used in automation jobs.
- Sharing: None
- Session Config

- N/A

Benefits of this approach

- Admin level setup is completed by the admin and each user must only enter their username and password, the information they should have readily available.
- Each user's authorization is managed in Salesforce.

Amazon Athena Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to AWS Athena as an import source. The fields you are required to set up on the data source depend on how the connector was configured by your administrator.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple AWS Athena instances. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Amazon Athena Configuration

- **Athena Region:** The hosting region for AWS.
- **Access Key:** AWS account access key.
- **Secret Key:** AWS account secret key.

Query Results Storage Configuration

- **S3 Bucket Name:** The name of the S3 bucket in which Athena will store query results.
- **S3 Object Prefix:** Prefix under which Athena will store query results within the specified S3 bucket. See [How do I use folders in an S3 Bucket](#) for more information on prefixes.

- **Encryption Type:** AWS server-side encryption type.

About query results

When using Athena, each query result is stored in the configured S3 bucket. This is how Athena is designed to function and is expected behavior. When using Athena to import to Data Prep, your query results will be cleaned up by default when the connection closes. The Athena Connector is designed to perform this clean-up task so that you only have one instance of the query result, not two. Should you want the query results from your import to Data Prep to remain available in S3, simply run the query in Athena standalone and import the resulting file to Data Prep from S3.

Web Proxy Configuration

If you connect to AWS Athena through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to AWS Athena should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for Data Source.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server.

Leave **username** and **password** blank for an unauthenticated proxy connection.

Data Import Information

Via Browsing

Browsing is supported for this Connector and uses Athena queries to generate the browseable hierarchy. Please see the note below about Athena's cost structure.

Via SQL Query

Access the [SQL reference](#) for details.

Best Practices

When using Athena, you are charged for each query that you run. The amount that you are charged is based on the amount of data scanned by the query. For more information, see [Amazon Athena Pricing](#).

Amazon S3 Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector enables the ability to import and export data against Amazon S3 object storage. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple S3 buckets. Using a descriptive name can be a big help to users in identifying the appropriate data source. If you are a Data Prep SaaS customer, please inform Data Prep DevOps how you would like this set.

Amazon S3 Client Configuration

- **Bucket name:** An S3 bucket represents a collection of objects stored in Amazon S3. The connector requires the following permissions: `s3:ListBucket`, `s3:GetObject`, and (for export only) `s3:PutObject`. *In addition*, if there is a `SourceIP` condition block specified in your bucket policy, then you must include the IP addresses for your Main Core Server and Automation Core Server (if you have one).

See [AWS S3 Bucket Permission/Policy Details](#) at the bottom of this article for more details.

- **Prefix:** Limits results to only those keys that begin with the specified prefix.
- **Encryption type:** Server-side encryption type to be used. See [AWS Encryption Types](#) for more information.
- **Bucket region:** This option allows users to specify the region in which their S3 bucket is hosted or to choose that the connector should automatically determine the region.

Amazon S3 Authentication

These options specify how to authenticate with S3.

- **AWS Credentials:** The Access Key ID and Secret Key associated with the user's AWS Access Key. This is the default setting. See [AWS Security Credentials](#) for more details.
- **Instance Profile (IAM Role):** enables all users in this tenant to access AWS without needing to individually authenticate. See [Using Instance Profile \(IAM Role\) to Grant Access to AWS Resources on Amazon EC2](#) for more details.

Note

This connector will automatically retrieve credentials from the EC2 server instance.

- **IAM Cross Account:** enables access to S3 by assuming a role in another AWS account that has access to the configured S3 bucket. See [Cross Account Access](#) for more details.

Note

For the Instance Profile (IAM Role) and IAM Cross Account options, Data Prep must be installed on your Amazon EC2 hosts.

Web Proxy

If you connect to Amazon S3 through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to the Amazon S3 REST Endpoint should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for Data Source.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server. *Leave username & password blank for an unauthenticated proxy connection.

Miscellaneous Settings

Socket Timeout Seconds: The number of seconds to wait for a response from Amazon S3 on an established connection. The default value is 5 minutes. To handle the export of large files, increase the value.

Data Import and Export Information

Via Browsing

The Connector will present a browsable directory hierarchy starting at the location defined in the Prefix field.

The Connector also supports Wildcard and Glob importing which enables users to import multiple S3 data files into Data Prep as a single Dataset.

Via SQL Query

As S3 is a file store, SQL Queries are not supported for this data source. If you would like to directly query AWS S3 data, please reach out to your Customer Success contact regarding Data Prep's AWS Athena Connector.

AWS S3 Bucket Permission/Policy Details

This section reviews the permissions that must be assigned in your S3 bucket policy and what you are required to do if you have a SourceIP condition block specified in your bucket policy.

Required Permissions

The AWS S3 connector requires specific permissions in your S3 bucket policy to ensure that you can successfully import data from S3, publish to S3, and automate importing from an S3 source. In summary:

- The connector requires the *s3:ListBucket* permission on the bucket for browsing.
- For importing the bucket contents, Data Prep requires the permissions *s3:GetObject*
- For exporting to the bucket, Data Prep requires the permission *s3:PutObject*

Sample bucket policy example

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "Stmt1965292834357",
      "Effect": "Allow",
      "Principal": {
        "AWS": [
          "arn:aws:iam::123456781234:user/pax01",
          "arn:aws:iam::432143214321:user/pax02",
          "arn:aws:iam::121212343434:user/pax03"
        ]
      },
      "Action": "s3:ListBucket",
      "Resource": "arn:aws:s3::paxhh-session1"
    },
    {
      "Sid": "Stmt1965293102818",
      "Effect": "Allow",
      "Principal": {
        "AWS": [
          "arn:aws:iam::123456781234:user/pax01",
          "arn:aws:iam::432143214321:user/pax02",
          "arn:aws:iam::121212343434:user/pax03"
        ]
      },
      "Action": [
        "s3:DeleteObject",
        "s3:DeleteObjectVersion",
        "s3:GetObject",
        "s3:GetObjectAcl",
        "s3:GetObjectTorrent",
        "s3:GetObjectVersion",
        "s3:GetObjectVersionAcl",
        "s3:GetObjectVersionTorrent",
        "s3:PutObject",
        "s3:PutObjectAcl",
        "s3:PutObjectVersionAcl"
      ],
      "Resource": "arn:aws:s3::paxhh-session1/*"
    }
  ]
}
```

Minimum policy permissions

The minimum policy permissions for reading from an S3 bucket are:

```
{
  "Version": "2012-10-17",
  "Statement": [
    { "Effect": "Allow", "Action": "s3:ListBucket", "Resource":
      "arn:aws:s3::mybucketname"
    },
    { "Effect": "Allow", "Action": [ "s3:ListBucket", "s3:GetObject" ],
      "Resource": "arn:aws:s3::mybucketname/*" }
  ]
}
```

The minimum policy permissions for writing to an S3 bucket are:

```
{
  "Version": "2012-10-17",
  "Statement": [
    { "Effect": "Allow", "Action": "s3:ListBucket", "Resource":
      "arn:aws:s3:::mybucketname"
    },
    { "Effect": "Allow", "Action": [ "s3:ListBucket", "s3:GetObject",
      "s3:PutObject" ], "Resource": "arn:aws:s3:::mybucketname/*" }
  ]
}
```

For a detailed explanation of S3 buckets, refer to [Working with Amazon S3 Buckets](#).

SourceIP condition block

If there is a **SourceIP condition block** specified in your bucket policy, then you must include the IP addresses of your Data Prep cloud servers or Data Prep Core Server (depending on your Data Prep deployment) in the **SourceIP Condition block**. In addition, if you have a dedicated Data Prep server for automation, you must also include the automation server IP addresses in the **SourceIP Condition block**.

Please consult with Data Prep's Customer Success team to obtain the list of IP addresses for Data Prep cloud servers.

For details on the condition block element and examples, see [Specifying Conditions in a Policy](#) and [Identity and Access Management \(IAM\) Policy Elements Reference](#).

Amazon Redshift Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Amazon Redshift as an import and export source. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Redshift data warehouses. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Database URL

- **JDBC URL:** the JDBC connection string; the database name can be included in the URL.
- Example JDBC URL: `jdbc:redshift://examplecluster.abc123xyz789_._us-west-2.redshift.amazonaws.com:5439/dev?ssl=true`
- For more information on connection string options, please see this [AWS Documentation](#).

Visibility Settings

You can control the database, schemas, and tables that are shown to users when they browse a data source during import. For databases, schemas, and tables, you can choose to:

- "Show only" which returns only the databases, schemas or tables that you specify here.
- "Hide" which hides the databases, schemas, and tables that you specify here.
- "Show all" which is the default setting to display everything in the data source.

When you select the "Show only" or "Hide" options, a field is provided for specifying the databases, schemas or tables on which you want the option enforced.

Note

These settings are not enforced when users query against the data source; query results still return a complete list of matches. For example, if you choose to "hide" a specific database, users can still execute queries that pull data from tables within that database. However, that database will not be displayed to users when they browse the data source.

Import Configuration

- **Pre-import SQL:** a SQL statement to execute before beginning import, after the table's schema is determined.
- **Post-import SQL:** a SQL statement to execute after import completes.

Export Configuration

- **Pre-export SQL:** a SQL statement to execute before beginning export, after the table is created if auto-create is enabled.
- **Post-export SQL:** a SQL statement to execute after export completes.

Redshift Credentials

The user authentication can be done through a shared account or an individual account. Depending on your choice, the following fields are required:

- **Individual account:**
 - **User:** the username for an individual account used to authenticate to the database.
 - **Password:** the password for an individual account used to authenticate to the database.
- **Shared Account:**
 - **User:** the username for a shared account used to authenticate to the database.
 - **Password:** the password for a shared account used to authenticate to the database.
 - **Role:** if roles have been implemented for this database, authenticated users with this user role can perform queries after authentication.

Amazon S3 Client Configuration

- **Export using S3?:** This option specifies whether the connector will export to Redshift by uploading data into Amazon S3 and then copying it into Redshift, or by inserting data directly into Redshift.
- **Export using S3:** The connector will upload data into Amazon S3 and copy it into Redshift. This is the recommended approach for larger datasets as it will allow for a more performant export.
- **Bucket name:** An S3 Bucket represents a collection of objects stored in Amazon S3.
- **Prefix:** Limits results to only those keys that begin with the specified prefix.
- **Socket Timeout Seconds:** The number of seconds to wait for a response from an established S3 connection. The default value is 5 minutes and may need to be increased to handle the export of large files.
- **Export using SQL insert statements:** The connector will insert data directly into Redshift. This option will result in slower exports. If you plan to only perform imports from Redshift, you can select this option in order to not have to enter your S3 account details.

Note

The connector requires the s3:ListBucket permission on the bucket. Bucket contents require permissions s3:ListBucket, s3:GetObject, and (for export only) s3:PutObject. In addition, if there is a SourceIP condition block specified in your bucket policy, then you must include the IP addresses for your Data Prep server and any servers that are used to run automation jobs.

See [Amazon S3 Connector Setup](#) for more details.

Amazon S3 Authentication Settings

- **AWS Credentials:** The Access Key ID and Secret Key are associated with the user's AWS Access Key.
- **Instance Profile (IAM Role):** No additional fields are required.

See [AWS Security Credentials](#) for more details.

Web Proxy

If you connect to Amazon Redshift through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if connection to the Amazon Redshift REST Endpoint should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for Data Source.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server. *Leave username & password blank for an unauthenticated proxy connection.

Data Import & Export Information

Via Browsing

The browsing experience will vary and is based on the database, schemas, and tables visibility settings selected above along with the user credentials supplied.

Via SQL Query

As noted in the visibility section, the limits to what a user can import via query are only limited to their authorization as determined by the credentials supplied for the connection.

Amazon DynamoDB Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to an Amazon DynamoDB for Library imports. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple DynamoDB accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Web Proxy

If you connect to DynamoDB through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to DynamoDB should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for Data Source.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server.

Note: Leave username & password blank for an unauthenticated proxy connection.

AWS Configuration

- **AWS Region:** Select the region to be used when sending requests to the DynamoDB API.
- **AWS Authentication Type:** These options specify how to authenticate with AWS.
 - **AWS Credentials:** Requires each user to enter the Access Key ID and Secret Key associated with the user's AWS Access Key. This is the default setting. See AWS Security Credentials for more details.
 - **Instance Profile (IAM Role):** Enables all users in this tenant to access AWS without needing to individually authenticate. This authentication method is only available to customers deployed on AWS VPCs with EC2 servers configured to allow this kind of authentication. For more information on this approach, please see Using Instance Profile (IAM Role) to Grant Access to AWS Resources on Amazon EC2.

Note

This connector will automatically retrieve credentials from the EC2 server instance.

DynamoDB Table Settings

Sample Items: Specifies the number of records (the equivalent of "rows" in a traditional relational database) to use when determining the schema of the table for import.

Note: DynamoDB is not a relational database, but Data Prep transforms all data into a tabular format upon import. In order to do this, Data Prep looks at the first n documents of a DynamoDB table (DynamoDB is a document database) and determines what attributes are present so that it can treat those attributes as columns. The "Sample Items" value determines the number of documents to look at.

Data Import Information

Via Browsing

DynamoDB tables in the specified AWS region are available for import.

Via SQL Query

Not Supported.

FAQ / Troubleshooting / Common Issues

Certain permissions are required in order to import data from DynamoDB. They are:

- dynamodb:ListTables
- dynamodb:Scan

Cloudera CDH5 HDFS Connector for Data Prep

User Persona: Data Prep Admin, Data Source Admin, or IT/DevOps

Availability information

This Connector is not available to Data Prep SaaS customers.

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Cloudera CDH 5.16 Hadoop File System (HDFS) for import and export. The following fields are used to define the connection parameters.

Note

Configuring this Connector requires file system access on the Data Prep Server and a `core-site.xml` with the Hadoop cluster configuration. Please reach out to your Customer Success representative for assistance with this step.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple HDFS clusters. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Simple Configuration (only for Simple authentication)

- **Username:** The application web server will connect to your HDFS cluster as the username you provide here.

Configuration

- **Data Store Root Directory:** The 'parent directory' on your cluster where the Connector will read from and write to for import and export operations. This also supports import and export for sub-directories of the root.

Kerberos Configuration

The following parameters are required for Kerberos authentication.

- **Principal:** Kerberos Principal.
- **Realm:** Kerberos Realm.
- **KDC Hostname:** Kerberos Key Distribution Center Hostname.
- **Kerberos Configuration File:** Fully-qualified path of Kerberos configuration file on webserver.
- **Keytab File:** Fully-qualified path of Kerberos Keytab File on webserver.
- **Use Application User:** Check this box to read/write as the logged-in application user, or uncheck to use proxy user.
- **Proxy User:** The proxy used to authenticate with the cluster. Enter `${user.name}` as the proxy user. `${user.name}` works similar to selecting Use Application User but allows for more flexibility. For example:
 - To add a domain to the user's credentials, enter `\domain_name\${user.name}` in the Proxy User field. Data Prep will pass the username and the domain.
 - Example: `\Accounts\${user.name}` results in `AccountsJoe` (assuming Joe is the username).
 - To apply a text modifier to the username, add `.modifier` to the key `${user.name}`. The acceptable modifiers are: `toLower`, `toUpper`, `toLowerCase`, `toUpperCase`, and `trim`.
 - For example `${user.name.toLowerCase}` converts Joe into joe (assuming Joe is the username).

Data Import Information

Via Browsing

Supported

Via SQL Query

Not supported

Cloudera CDH6 HDFS Connector for Data Prep

User Persona: Data Prep Admin, Data Source Admin, or IT/DevOps

Availability information

This Connector is not available to Data Prep SaaS customers.

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configuring Data Prep

This connector allows you to connect to an HDFS cluster for imports and exports. The fields you are required to set up here depend on the authentication method you select—Simple or Kerberos. The type of authentication you select will apply to all data sources that you create based on a connector configuration.

Note

Configuring this Connector requires file system access on the Data Prep Server and a `core-site.xml` with the Hadoop cluster configuration. Please reach out to your Customer Success representative for assistance with this step.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple HDFS clusters. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Hadoop Cluster

- **Authentication Method:** Choose between Simple or Kerberos. The type of authentication you select will apply to all Data Sources that you create based on a connector configuration. See Simple or Kerberos Configuration section below for more details depending on your selection.
- **Cluster Core Site XML Path:** Fully qualified path of core-site.xml on webserver. Example: /path/to/core-site.xml
- **Cluster HDFS Site XML Path:** Fully qualified path of hdfs-site.xml on webserver. Example: /path/to/hdfs-site.xml
- **Native Hadoop Library Path:** Fully qualified path of native Hadoop libraries on webserver. Example: /path/to/libraries

Simple Configuration (only for Simple authentication)

- **Username:** The application web server will connect to your HDFS cluster as the username you provide here.

Kerberos Configuration

The following parameters are required for Kerberos and Hybrid authentication.

- **Principal:** Kerberos Principal.
- **Realm:** Kerberos Realm.
- **KDC Hostname:** Kerberos Key Distribution Center Hostname.
- **Kerberos Configuration File:** Fully-qualified path of Kerberos configuration file on webserver.
- **Keytab File:** Fully-qualified path of Kerberos Keytab File on webserver.
- **Use Application User:** Check this box to read/write as the logged-in application user, or uncheck to use proxy user.
- **Proxy User:** The proxy used to authenticate with the cluster. `${user.name}` can be entered as the proxy user. `${user.name}` works similar to selecting Use Application User but allows for more flexibility. For example:
 - To add a domain to the user's credentials, enter `\domain_name\${user.name}` in the Proxy User field. Data Prep will pass the username and the domain.
 - Example: `\Accounts\${user.name}` results in `AccountsJoe` (assuming Joe is the username).
 - To apply a text modifier to the username, add `.modifier` to the key `${user.name}`. The acceptable modifiers are: `toLower`, `toUpper`, `toLowerCase`, `toUpperCase`, and `trim`.
 - For example `${user.name.toLowerCase}` converts Joe into joe (assuming Joe is the username).

Configuration

- **Data Store Root Directory:** The 'parent directory' on your cluster where the Data Library will read from and write to for import and export operations. This also supports import and export for sub-directories of the root.
- **Map INT96 to Datetime:** Check to convert INT96 type fields to Datetime values on import.

Credentials

- **Hive User:** The username used to access Hive for Simple and Hybrid authentication.

- **Hive Password:** The password used to access Hive for Simple and Hybrid authentication.

Hive Options

- **Pre-Import SQL:** SQL to be executed before import process. This SQL may execute multiple times (for preview and import) and could be multiple SQL statements, newline-delimited.
- **Post-Import SQL:** SQL to be executed after import process. This SQL may execute multiple times (for preview and import) and could be multiple SQL statements, newline-delimited.

Note

As the Pre- and Post-Import SQL may be executed multiple times throughout the import process, please take care when specifying these values in the Connector/Datasource Configuration as they will be executed for every import performed with this configuration.*

- **Pre-Export SQL:** SQL to be executed before export process. This SQL will execute once and could be multiple SQL statements, newline-delimited.
- **Post-Export SQL:** SQL to be executed after export process. This SQL will execute once and could be multiple SQL statements, newline-delimited.

Data Import Information

Via Browsing

- Browse:
 - Delimited datasets: comma, tab...
 - XML
 - JSON
 - Excel: Xls and XLSX
 - Avro
 - Parquet
 - Fixed format
 - Browse to a file and select it for import
 - Supported data formats:
- Wildcard:
 - Globbing is supported

Via SQL Query

Using SQL Select queries

Export

Supported using one of the stream-based formats listed under [Import Via Browser](#).

Cloudera CDH5 Hive Connector for Data Prep

User Persona: Data Prep Admin, Data Source Admin, or IT/DevOps

Availability information

This Connector is not available to Data Prep SaaS customers.

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to a Cloudera CDH 5.16 Hive for import and export. The following fields are used to define the connection parameters. The fields you are required to set up here depend on the authentication method you select – Simple, Kerberos, or Hybrid. The type of authentication you select will apply to all Data Sources that you create based on a connector configuration.

Notes

Configuring this Connector requires file system access on the Data Prep Server and a `core-site.xml` with the Hadoop cluster configuration. Please reach out to your Customer Success representative for assistance with this step.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Hive databases. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Hadoop Cluster

- **HDFS User:** The username on the HDFS cluster used to write files for export to Hive.

Kerberos Configuration

The following parameters are required for Kerberos and Hybrid authentication.

- **Principal:** Kerberos Principal.
- **Realm:** Kerberos Realm.
- **KDC Hostname:** Kerberos Key Distribution Center Hostname.
- **Kerberos Configuration File:** Fully-qualified path of Kerberos configuration file on webserver.
- **Keytab File:** Fully-qualified path of Kerberos Keytab File on webserver.
- **Use Application User:** Check this box to read/write as the logged-in application user, or uncheck to use proxy user.
- **Proxy User:** The proxy used to authenticate with the cluster. `${user.name}` can be entered as the proxy user. `${user.name}` works similar to selecting Use Application User but allows for more flexibility. For example:
 - To add a domain to the user's credentials, enter `\domain_name\${user.name}` in the Proxy User field. Data Prep will pass the username and the domain.
 - Example: `\Accounts\${user.name}` results in `AccountsJoe` (assuming Joe is the username).
 - To apply a text modifier to the username, add `.modifier` to the key `${user.name}`. The acceptable modifiers are: `toLower`, `toUpper`, `toLowerCase`, `toUpperCase`, and `trim`.
 - For example `${user.name.toLowerCase}` converts Joe into joe (assuming Joe is the username).

Hive Configuration

- **JDBC URL:** The URL used to access Hive for import and registration of external tables. If Kerberos authentication is used, the following string must be added to the URL: `";auth=kerberos;hive.server2.proxy.user=${user.name}"`
 - If a proxy user is used, then the string `${user.name}` must be replaced with the proxy username
- **Hive File Location:** The location on the HDFS cluster used to store Hive files for external tables.

Credentials

- **Hive User:** The username used to access Hive for Simple and Hybrid authentication.
- **Hive Password:** The password used to access Hive for Simple and Hybrid authentication.

Hive Options

- **Pre-Import SQL:** SQL to be executed before import process. This SQL may execute multiple times (for preview and import) and could be multiple SQL statements, newline-delimited.
- **Post-Import SQL:** SQL to be executed after import process. This SQL may execute multiple times (for preview and import) and could be multiple SQL statements, newline-delimited.

Note

As the Pre- and Post-Import SQL may be executed multiple times throughout the import process, please take care when specifying these values in the Connector/Datasource Configuration as they will be executed for every import performed with this configuration.*

- **Pre-Export SQL:** SQL to be executed before export process. This SQL will execute once and could be multiple SQL statements, newline-delimited.
- **Post-Export SQL:** SQL to be executed after export process. This SQL will execute once and could be multiple SQL statements, newline-delimited.

Data Import Information

Via Browsing

Not Supported

Via SQL Query

Using SQL Select queries

Cloudera CDH6 Hive Connector for Data Prep

User Persona: Data Prep Admin, Data Source Admin, or IT/DevOps

Availability information

This Connector is not available to Data Prep SaaS customers.

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configuring Data Prep

This connector allows you to connect to a Hive database for imports and exports. The fields you are required to set up here depend on the authentication method you select—Simple, Kerberos, or Hybrid. The type of authentication you select will apply to all Data Sources that you create based on a connector configuration.

Note

Configuring this Connector requires file system access on the Data Prep Server and a `core-site.xml` with the Hadoop cluster configuration. Please reach out to your Customer Success representative for assistance with this step.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Hive databases. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Hadoop Cluster

- **Authentication Method:** Choose between Simple, Kerberos, or Hybrid. The type of authentication you select will apply to all Data Sources that you create based on a connector configuration.
- **Cluster Core Site XML Path:** Fully qualified path of core-site.xml on webserver. Example: /path/to/core-site.xml
- **Cluster HDFS Site XML Path:** Fully qualified path of hdfs-site.xml on webserver. Example: /path/to/hdfs-site.xml
- **Native Hadoop Library Path:** Fully qualified path of native Hadoop libraries on webserver. Example: /path/to/libraries
- **HDFS User:** The username on the HDFS cluster used to write files for export to Hive.

Hive Configuration

- **JDBC URL:** The URL used to access Hive for import and registration of external tables. If Kerberos authentication is used, the following string must be added to the URL: ";auth=kerberos;hive.server2.proxy.user=\${user.name}"
- If a proxy user is used, then the string \${user.name} must be replaced with the proxy username
- **Hive File Location:** The location on the Hadoop cluster used to store Hive files for external tables.

Kerberos Configuration

The following parameters are required for Kerberos and Hybrid authentication.

- **Principal:** Kerberos Principal.
- **Realm:** Kerberos Realm.
- **KDC Hostname:** Kerberos Key Distribution Center Hostname.
- **Kerberos Configuration File:** Fully-qualified path of Kerberos configuration file on webserver.
- **Keytab File:** Fully-qualified path of Kerberos Keytab File on webserver.
- **Use Application User:** Check this box to read/write as the logged-in application user, or uncheck to use proxy user.
- **Proxy User:** The proxy used to authenticate with the cluster. \${user.name} can be entered as the proxy user. \${user.name} works similar to selecting Use Application User but allows for more flexibility. For example:
 - To add a domain to the user's credentials, enter \domain_name\\${user.name} in the Proxy User field. Data Prep will pass the username and the domain.
 - Example: \Accounts\\${user.name} results in AccountsJoe (assuming Joe is the username).
 - To apply a text modifier to the username, add .modifier to the key \${user.name}. The acceptable modifiers are: toLower, toUpper, toLowerCase, toUpperCase, and trim.
 - For example \${user.name.toLowerCase} converts Joe into joe (assuming Joe is the username).

Credentials

- **Hive User:** The username used to access Hive for Simple and Hybrid authentication.
- **Hive Password:** The password used to access Hive for Simple and Hybrid authentication.

Visibility Settings

You can control the schemas and tables that are shown to users when they browse a data source during import. For schemas and tables you can choose to:

- **"Show only"** which returns only the schemas or tables that you specify here.
- **"Hide"** which hides the schemas and tables that you specify here.
- **"Show all"** which is the default setting to display everything in the data source.

When you select the "Show only" or "Hide" options, a field is provided for specifying the schemas or tables on which you want the option enforced.

Note

These settings are not enforced when users query against the data source; query results still return a complete list of matches. For example, if you choose to "hide" a specific schema, users can still execute queries that pull data from tables within that schema. However, that schema will not be displayed to users when they browse the data source.

Import Configuration

- **Query Prefetch Size:** Number of rows per batch.
- **Max Column Size:** The maximum size in Unicode characters allowed for any value for both import and export. Values larger than this will be replaced by null.
- **PRE-IMPORT SQL:** SQL to be executed before import process. This SQL may execute multiple times (for preview and import) and could be multiple SQL statements, newline-delimited.
- **POST-IMPORT SQL:** SQL to be executed after import process. This SQL may execute multiple times (for preview and import) and could be multiple SQL statements, newline-delimited.

Note

As the Pre- and Post-Import SQL may be executed multiple times throughout the import process, please take care when specifying these values in the Connector/Datasource Configuration as they will be executed for every import performed with this configuration.*

Export Configuration

- **PRE-EXPORT SQL:** SQL to be executed before export process. This SQL will execute once and could be multiple SQL statements, newline-delimited.
- **POST-EXPORT SQL:** SQL to be executed after export process. This SQL will execute once and could be multiple SQL statements, newline-delimited.

Data Import and Export Information

Via Browsing

Browse to a table and "Select" the table for import.

Via SQL Query

Using SQL Select queries

Cloudera CDH6 Impala Connector for Data Prep

User Persona: Data Prep Admin, Data Source Admin, or IT/DevOps

Availability information

This Connector is not available to Data Prep SaaS customers.

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configuring Data Prep

This connector allows you to connect to an Impala database for imports and exports. The fields you are required to set up here depend on the authentication method you select—Simple, Kerberos, or Hybrid. The type of authentication you select will apply to all Data Sources that you create based on a connector configuration.

Note

Configuring this Connector requires file system access on the Data Prep Server and a `core-site.xml` with the Hadoop cluster configuration. Please reach out to your Customer Success representative for assistance with this step.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Impala databases. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Hadoop Cluster

- **Authentication Method:** Choose between Simple, Kerberos, or Hybrid. The type of authentication you select will apply to all Data Sources that you create based on a connector configuration.
- **Cluster Core Site XML Path:** Fully qualified path of core-site.xml on webserver. Example: /path/to/core-site.xml
- **Cluster HDFS Site XML Path:** Fully qualified path of hdfs-site.xml on webserver. Example: /path/to/hdfs-site.xml
- **Native Hadoop Library Path:** Fully qualified path of native Hadoop libraries on webserver. Example: /path/to/libraries
- **HDFS User:** The username on the HDFS cluster used to write files for export to Impala.

Impala Configuration

- **JDBC URL:** The URL used to access Impala for import and registration of external tables. If Kerberos authentication is used, the following string must be added to the URL: ";auth=kerberos;impala.server2.proxy.user=\${user.name}"
- If a proxy user is used, then the string \${user.name} must be replaced with the proxy username
- **Impala File Location:** The location on the Hadoop cluster used to store Impala files for external tables.

Kerberos Configuration

The following parameters are required for Kerberos and Hybrid authentication.

- **Principal:** Kerberos Principal.
- **Realm:** Kerberos Realm.
- **KDC Hostname:** Kerberos Key Distribution Center Hostname.
- **Kerberos Configuration File:** Fully-qualified path of Kerberos configuration file on webserver.
- **Keytab File:** Fully-qualified path of Kerberos Keytab File on webserver.
- **Use Application User:** Check this box to read/write as the logged-in application user, or uncheck to use proxy user.
- **Proxy User:** The proxy used to authenticate with the cluster. \${user.name} can be entered as the proxy user. \${user.name} works similar to selecting Use Application User but allows for more flexibility. For example:
 - To add a domain to the user's credentials, enter \domain_name\\${user.name} in the Proxy User field. Data Prep will pass the username and the domain.
 - Example: \Accounts\\${user.name} results in AccountsJoe (assuming Joe is the username).
 - To apply a text modifier to the username, add .modifier to the key \${user.name}. The acceptable modifiers are: toLower, toUpper, toLowerCase, toUpperCase, and trim.
 - For example \${user.name.toLowerCase} converts Joe into joe (assuming Joe is the username).

Credentials

- **Impala User:** The username used to access Impala for Simple and Hybrid authentication.
- **Impala Password:** The password used to access Impala for Simple and Hybrid authentication.

Visibility Settings

You can control the schemas and tables that are shown to users when they browse a data source during import. For schemas and tables you can choose to:

- **"Show only"** which returns only the schemas or tables that you specify here.
- **"Hide"** which hides the schemas and tables that you specify here.
- **"Show all"** which is the default setting to display everything in the data source.

When you select the "Show only" or "Hide" options, a field is provided for specifying the schemas or tables on which you want the option enforced.

Note

These settings are not enforced when users query against the data source; query results still return a complete list of matches. For example, if you choose to "hide" a specific schema, users can still execute queries that pull data from tables within that schema. However, that schema will not be displayed to users when they browse the data source.

Import Configuration

- **Query Prefetch Size:** Number of rows per batch.
- **Max Column Size:** The maximum size in Unicode characters allowed for any value for both import and export. Values larger than this will be replaced by null.
- **PRE-IMPORT SQL:** SQL to be executed before import process. This SQL may execute multiple times (for preview and import) and could be multiple SQL statements, newline-delimited.
- **POST-IMPORT SQL:** SQL to be executed after import process. This SQL may execute multiple times (for preview and import) and could be multiple SQL statements, newline-delimited.

Note

As the Pre- and Post-Import SQL may be executed multiple times throughout the import process, please take care when specifying these values in the Connector/Datasource Configuration as they will be executed for every import performed with this configuration.

Export Configuration

- **PRE-EXPORT SQL:** SQL to be executed before export process. This SQL will execute once and could be multiple SQL statements, newline-delimited.
- **POST-EXPORT SQL:** SQL to be executed after export process. This SQL will execute once and could be multiple SQL statements, newline-delimited.

Data Import and Export Information

Via Browsing

Browse to a table and "Select" the table for import.

Via SQL Query

Using SQL Select queries

Databricks Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Databricks for Library imports and exports. The following fields are used to define the connection parameters. It has been certified against Databricks on Azure and AWS.

This connector enables Import via browse, query and export operation.

All actions are performed over JDBC connection, except the data loading directly into Databricks storage (i.e., ADLS Gen2 or S3 bucket, depending on Databricks service provider) at the time of export.

Note

The Databricks connector for Data Prep has an import size limit of 8GB. Importing files larger than 8GB may result in an error.

General

Name: Name of the data source as it will appear to users in the UI.

Description: Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Databricks accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Databricks Server Configuration

- **Databricks Service Provider:** Set this property based on which type of Databricks service you want to connect. We support Databricks on Azure and AWS.
- Databricks on Azure
- Databricks on AWS
- **Databricks Server Settings Type:** Set this property based on how you want to configure the datasource to connect to Databricks.
- Basic
- Advanced
- **Databricks Server:** The hostname of the server hosting the Databricks service.
- **Databricks Port:** The port of the Databricks server.
- **Use SSL:** Set this property to the value specified in the 'hive.server2.use.SSL' property of your Hive configuration file (hive-site.xml).
- **Transport Mode:** Set this property to the value specified in the 'hive.server2.transport.mode' property of your Hive configuration file (hive-site.xml).
- **HTTP Path:** In HTTP Transport Mode, set this property to specify the path component of the URL endpoint. This property should be set to the value specified in the 'hive.server2.thrift.http.path' property of the Hive configuration file (hive-site.xml).
- **Timeout:** Seconds to wait until an operation times out. If set to zero, operations do not timeout.
- **JDBC Url:** In Advance settings, set the entire JDBC Url to connect to Databricks. Please refer to the CData JDBC driver documentation documentation for more details.

Databricks Server Authentication Configuration

- **User:** The username used to authenticate with Databricks server. Usually the username is 'token'.
- **Password:** The Personal Access Token used to authenticate with Databricks. Personal Access Token can be obtained by navigating to the User Settings page of your Databricks instance and selecting the Access Tokens tab.

Databricks Log Settings

- **Verbosity:** The verbosity level that determines the amount of detail included in the log file. This is very useful to debug an issue in production.
- **Logfile:** A path for driver log file within Pax server. All the directories in the specified path should pre-exist.

Databricks Server Export Storage Layer Configuration

Azure

- **ADLS Gen2 Data Store Root Directory:** The apparent root path accessible by this connector. Use '/' to store the Databricks data within root folder of ADLS Gen2 file system.

- **ADLS Gen2 Storage Account Name:** The Subdomain Name of your unique Azure URL. This Storage Account must be associated and accessible by Databricks cluster. ADLS Gen2 Storage account names must be between 3 and 24 characters in length and may contain numbers and lowercase letters only. Your ADLS Gen2 storage account name must be unique within Azure. No two storage accounts can have the same name.
- **ADLS Gen2 File System Name:** The name of the ADLS Gen2 file system where you want to store the Databricks data within the storage account. This is sometimes called the 'container' name.
- **Authentication Type:** The type of authentication you want to connect to ADLS Gen2 storage, either "Storage Account Access Key" or "Active Directory Username/Password."
- **ADLS Gen2 Storage Account Access Key:** Enter the Storage Account Access Key in the field. This is sometimes referred to as a "Shared Key."
- **Active Directory Username/Password:** Enter the Azure Directory username and password associated with your account.

Note

You must grant access for Data Prep to read and write data within your Microsoft account, otherwise you will get an error while attempting to connect. To grant access, click **Test Data Source** in the Connector set-up pane and follow the **Grant Access** link. This brings you to your Microsoft account where you can log in and grant access. Then, come back to Data Prep to continue.

AWS

- **S3 Bucket Name:** An S3 Bucket name, where you want to store the Databricks data in Amazon S3. This S3 bucket must be associated and accessible by Databricks cluster.
- **S3 Object Prefix:** The apparent root path accessible by this connector. Use "/" to store the Databricks data within root folder of S3 Bucket.
- **Authentication type:** The authentication method for accessing S3 bucket.
- **AWS Credentials:** Requires each user to enter the Access Key ID and Secret Key associated with the user's AWS Access Key. This is the default setting.
- **Instance Profile (IAM Role):** Enables all users in this tenant to access AWS without needing to individually authenticate.
- **IAM Cross Account:** Enables access to S3 by assuming a role in another AWS account that has access to the configured S3 bucket.

Important

For the Instance Profile (IAM Role) and IAM Cross Account options, Data Prep must be installed on your Amazon EC2 hosts.

- **Encryption type:**
 - None
 - SSE-S3
 - SSE-KMS
- **Bucket Region Locator:** S3 AWS bucket region locator strategy.

• **Socket Timeout Seconds:** The number of seconds to wait for a response from Amazon S3 on an established connection. The default value is 5 minutes. To handle export of large files, increase the value.

• **Browse:**

- View list of available databases and tables.

• **Import:**

- Browse:
 - Browse to a table (partitioned and non-partitioned) and click the name for import.
- Query:
 - Using a legal SQL Select Query

• **Export:**

- Browse to a database and export the table.

Configuration Layout

Databricks on Azure with ADLS Gen2 Storage

The screenshot shows the configuration interface for Databricks on Azure with ADLS Gen2 Storage. The interface is dark-themed and organized into several sections:

- GENERAL:** Includes fields for Name and Description.
- DATABRICKS SERVER CONFIGURATION:** Includes dropdowns for Databricks on Azure (Basic), Databricks Service Provider, Databricks Server Settings Type, Databricks Server (127.0.0.1), Databricks Port (10000), Use SSL, Transport Mode, and Timeout (3600).
- DATABRICKS SERVER AUTHENTICATION CONFIGURATION:** Includes fields for token (User), Personal Access Token, and User.
- DATABRICKS LOG SETTINGS:** Includes fields for Verbosity and Logfile (/tmp/paxata/databricks-driver.log).
- DATABRICKS SERVER EXPORT STORAGE LAYER CONFIGURATION:** Includes fields for ADLS Gen2 Data Store Root Directory, ADLS Gen2 Storage Account Name, ADLS Gen2 File System Name, and Authentication Type.

Databricks on AWS with S3 bucket Storage

Databricks

Connector type

GENERAL

Databricks on Azure

Name Description

DATABRICKS SERVER CONFIGURATION

Databricks on Azure Basic

Databricks Service Provider Databricks Server Settings Type

127.0.0.1 10000

Databricks Server Databricks Port

Use SSL Transport Mode

3600

Timeout

DATABRICKS SERVER AUTHENTICATION CONFIGURATION

token Set your personal access token obtained from Databricks

User Personal Access Token

DATABRICKS LOG SETTINGS

Verbosity /tmp/paxata/databricks-driver.log

Logfile

DATABRICKS SERVER EXPORT STORAGE LAYER CONFIGURATION

/

ADLS Gen2 Data Store Root Directory ADLS Gen2 Storage Account Name

ADLS Gen2 File System Name Authentication Type

Import via Browse

Databricks

Connector type

GENERAL

Databricks on Azure

Name Description

DATABRICKS SERVER CONFIGURATION

Databricks on Azure **Basic**

Databricks Service Provider Databricks Server Settings Type

127.0.0.1 10000

Databricks Server Databricks Port

Use SSL Transport Mode

3600

Timeout

DATABRICKS SERVER AUTHENTICATION CONFIGURATION

token Set your personal access token obtained from Databricks

User Personal Access Token

DATABRICKS LOG SETTINGS

Verbosity Logfile

/tmp/paxata/databricks-driver.log

DATABRICKS SERVER EXPORT STORAGE LAYER CONFIGURATION

ADLS Gen2 Data Store Root Directory ADLS Gen2 Storage Account Name

ADLS Gen2 File System Name Authentication Type

Export via Browse

Databricks

Connector type

GENERAL

Databricks on Azure

Name Description

DATABRICKS SERVER CONFIGURATION

Databricks on Azure **Basic**

Databricks Service Provider Databricks Server Settings Type

127.0.0.1 10000

Databricks Server Databricks Port

Use SSL Transport Mode

3600

Timeout

DATABRICKS SERVER AUTHENTICATION CONFIGURATION

token Set your personal access token obtained from Databricks

User Personal Access Token

DATABRICKS LOG SETTINGS

Verbosity Logfile

DATABRICKS SERVER EXPORT STORAGE LAYER CONFIGURATION

ADLS Gen2 Data Store Root Directory ADLS Gen2 Storage Account Name

ADLS Gen2 File System Name Authentication Type

Configure Databricks cluster

In addition to configuring your Databricks connector in DataRobot, you must also add Spark configurations to your Databricks cluster:

1. Navigate to the **Configuration** tab of your DataBricks cluster and expand **Advanced Options**.
2. In the **Spark** tab, add and save the following configuration settings:

```
spark.sql.legacy.parquet.datetimeRebaseModeInRead LEGACY
```

```
spark.driver.maxResultSize 12g
```

The screenshot shows the Databricks Configuration page. At the top, there are navigation tabs: Configuration (highlighted), Notebooks (0), Libraries, Event Log, Spark UI, Driver Logs, Metrics, Apps, and Spark Cluster UI - Master. Below the tabs, the 'Worker Type' section shows 'Standard_DS3_v2' with '14 GB Memory, 4 Cores', 'Min Workers' set to 2, and 'Max Workers' set to 8. There is a checkbox for 'Spot instances' which is unchecked. The 'Driver Type' section also shows 'Standard_DS3_v2' with '14 GB Memory, 4 Cores'. Below this, there is a 'DBU / hour: 2.25 - 6.75' label and a 'Standard_DS3_v2' button. The 'Advanced Options' section is expanded, showing 'Azure Data Lake Storage Credential Passthrough' with a note 'Available on Azure Databricks Premium Learn more' and an unchecked checkbox 'Enable credential passthrough for user-level data access'. Below this, there are tabs for 'Spark' (highlighted), Tags, Logging, Init Scripts, JDBC/ODBC, and Permissions. The 'Spark Config' section contains several configuration properties: 'spark.sql.legacy.parquet.datetimeRebaseModeInRead LEGACY', 'spark.hadoop.fs.azure.createRemoteFileSystemDuringInitialization true', 'spark.hadoop.hive.server2.idle.session.timeout 300000', 'spark.hadoop.fs.azure.account.key.apacadlsgen2storage.dfs.core.windows.net 6exz/bLu2XzYuBTHv22/1zo6eqjsU7rvEMNuwfGLLwS+ae/nyWZmpO9hVqD8oMArHurMhFwKn3FMwDo9JBaWlQ==', and 'spark.driver.maxResultSize 12g'. The 'Environment Variables' section shows 'No environment variables'.

Databricks Connector Known Issues and Limitations

The following list of features might not work in some production environments. These issues will be fixed in an upcoming release.

- Authentication with Azure Databricks instance and ADLS GEN2 storage using Active Directory credentials.
- Authentication with AWS Databricks instance and Amazon S3 service storage using Cross account bucket ARN.
- Authentication with AWS Databricks instance and Amazon S3 service storage using IAM Role enabled.
- Import tables from AWS Databricks instance whose data is encrypted with SSE-KMS in an unencrypted S3 bucket.
- Import tables from AWS Databricks instance whose data is encrypted with SSE-S3 and SSE-KMS in an encrypted S3 bucket.
- Import tables from AWS Databricks instance after authentication with Cross account and IAM role.
- Export to AWS Databricks instance after authentication with Cross account and IAM role.

DataRobot Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to DataRobot for Library imports and exports. This connection also allows you to [create and access a DataRobot project](#) directly from Data Prep.

The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple DataRobot accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

DataRobot Configuration

- **Server URL:** The server URL for DataRobot. For example: <https://app.datarobot.com>

Authentication Configuration

- **Authentication Type:** Select the authentication type to use:
 - **API Key**
 - **Key:** The DataRobot API Key

User Credentials

- **Email:** Email or username for authenticating with DataRobot.
- **Password:** Password for authenticating with DataRobot.
- **Note:** Multi-factor authentication is not supported and will result in an error.

Data Import & Export Information

Via Browsing

- To import from the **AI Catalog**, **select the AI Catalog** option to view all available datasets and select the desired dataset to see a preview and adjust the import settings.

Note

If during import you receive an error stating you do not have permissions to download datasets from the AI Catalog, you need to adjust your settings in DataRobot. Login to DataRobot and click the **user icon** at the top-right → **Settings** → **Optional Products** → check **Enable AI Catalog Downloads** → **Save**. Then come back to DataRobot Data Prep to continue.

- To export from the AI Catalog, select the **AI Catalog** option, then click **Select**. Name the dataset and click **Export**.
- To export a dataset directly to DataRobot and create a project in one step, see [Create a DataRobot project](#).

Via SQL Query

Not supported

FAQ/Troubleshooting/Common Issues

What if I export the model I've generated in DataRobot and want to run that code where my data lives?

Data Prep has over 50 other Connectors and can likely still send the prepped data to the appropriate location. If Data Prep does not support Connectivity to the service/storage location you require, please reach out to your Customer Success Representative.

Why can't I import my dataset?

- Issue 1: Data Prep has designed the integration with DataRobot's **AI Catalog** to only support the importing of "Snapshotted" datasets. The data contained in "Not snapshotted" datasets are not actually stored in DataRobot and are retrieved upon usage. In Data Prep's case, that would mean DataRobot would fully import a dataset from the data source and only then would Data Prep begin importing that dataset. For "Not snapshotted" datasets, it's much more efficient to pull the data directly from the data source into Data Prep. To determine if your dataset is a "Snapshot", go to the **AI Catalog**, select the dataset in question and look at the "Status" in the right-hand panel.

- Issue 2: If you receive an error stating you do not have permissions to download datasets from the **AI Catalog**, you need to adjust your settings in DataRobot. Login to DataRobot and click the **user icon** at the top-right > **Settings** > **Optional Products** > check **Enable AI Catalog Downloads** > **Save**. Then come back to DataRobot Data Prep to continue.
- Issue 3: If you receive an error stating "Mapping for not found, expected one of []", you need to adjust your settings in DataRobot. Login to DataRobot and click the **user icon** at the top-right > **Settings** > **CSV export** > uncheck **Include BOM** > **Save**. Then come back to DataRobot Data Prep to continue.

When I export a new version of my dataset, does it appear as such in the **AI Catalog**?

Yes, versions of datasets with the same name will appear under the *Version History* tab of the **AI Catalog**, rather than as a new dataset.

Requirements for data exports to the AI Catalog

Datasets exported to DataRobot must meet the following criteria:

- At least 100 rows
- At least 2 columns
- Have valid column names

Google Analytics Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Google Analytics for browsing and importing available data. The following parameters are used to configure the connection.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Google Analytics accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Google Analytics Configuration

- **OAuth Verifier Key:** The verifier key used to authenticate with Google Analytics. To obtain the verifier key, click 'Test Data Source' and follow the link to grant access to Google Analytics. After allowing access, you will be redirected to a page that displays an access code. Copy the code into this field.
- **Profile:** The Google Analytics profile or view to connect to. This can be set to either the Id or website URL for the profile.

Web Proxy

Select the web proxy options for connecting to Google Analytics.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to Google Analytics should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy Host:** The hostname or IP address of the proxy server.
- **Proxy Port:** The port of the proxy server.
- **Proxy Username and Proxy Password:** User credentials for an authenticated proxy connection. Leave these blank for an unauthenticated proxy connection.

Data Import Information

Via Browsing

- Browse a list of pre-defined datasets in Google Analytics and select a dataset, click Select for import.
- Pre-defined datasets:
 - [Accounts](#): Lists all Accounts to which the user has access.
 - [AdWords](#): Retrieves AdWords data.
 - [Ecommerce](#): Retrieves Ecommerce data.
 - [Events](#): Retrieves Event data.
 - [GoalCompletions](#): Retrieves Goal Completion data.
 - [Profiles](#): Lists all Profiles to which the user has access.
 - [Segments](#): Lists all Segments to which the user has access.
 - [SiteContent](#): Retrieves internal Site Content data.
 - [SiteSearch](#): Retrieves internal Site Search data.
 - [SiteSpeed](#): Retrieves internal Site Speed data.
 - [Traffic](#): Retrieves all Traffic data.
 - [WebProperties](#): Lists Web Properties to which the user has access.
- Export is not supported.

Via SQL Query

- Using a [SQL Select](#) query.
- The default behavior of the CData JDBC Driver is to retrieve data for the last 7 days. In order to customize the time window you can set the StartDate and EndDate values in your query directly. The supported inputs for StartDate and EndDate in the Google Analytics API are 'today', 'yesterday', 'NdaysAgo' (where N is some number), and an exact date.

For example:

```
SELECT * FROM Traffic WHERE StartDate='2020-01-01' AND EndDate='5daysAgo'
```

See http://cdn.cdata.com/help/DAE/jdbc/pg_table-sitecontent.htm for more details.

Google BigQuery Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to [BigQuery](#) for importing and exporting available data. The fields you are required to set up here depending on how the connector was configured by your administrator.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple BigQuery accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

BigQuery Configuration

- **OAuth Verifier Key:** The verifier key used to authenticate with BigQuery. To obtain the verifier key, click "Test Data Source" and follow the link to grant access to BigQuery. After allowing access, you will be redirected to a page that displays an access code. Copy the code into this field.
- **Profile:** The ID of the [GCP Project](#) to which you will connect.
- **Automatically Create Table (optional):** If enabled, Data Prep will drop the table whose name matches the name of the exported dataset, if one already exists, and recreate the table using the exported dataset. If disabled, Data Prep will expect that the table is already created and will try to export it.

Google Cloud Storage Configuration for Export

These fields are necessary to perform export to BigQuery. If you intend to only import, you can leave these blank.

Note

They must either both be provided or both left blank.

- **Google Cloud Storage Bucket Name:** Google Cloud Storage bucket name to be used as a staging area for export.
- **Google Cloud Storage JSON Web Token:** Content of JSON Web Token (JWT) to be used to connect to Google Cloud Storage.

Web Proxy

If you connect to BigQuery through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to BigQuery should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy Host:** The hostname or IP address of the proxy server.
- **Proxy Port:** The port of the proxy server.
- **Proxy Username and Proxy Password:** User credentials for an authenticated proxy connection. Leave these blank for an unauthenticated proxy connection.

Data Import Information

Via Browsing

- View datasets and tables within the project specified in your configuration. The project will appear as the top-level directory in the browsing view.
- Browse to a table within a dataset and "Select" the table for import.

Via SQL Query

- Using a [SQL Select Query](#).

Usage

Each table name in a query must be single-quoted, with any dot separation occurring outside the single-quotes.

Valid syntax `SELECT * FROM `my-project`.`paxata`.`test``

Invalid syntax `SELECT * FROM `my-project.paxata.test``

Google Cloud Storage Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Google Cloud Storage (GCS) for browsing and importing objects. The following fields are used to create a connection to the data source.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple GCS accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Google Cloud Storage Configuration

- **Bucket Name:** A Google Cloud Storage bucket represents a collection of objects in Google Cloud Storage.
- **Object Prefix:** Prefix is a folder/sub-folder in the bucket. Select the prefix you want to use in the bucket. Default value to view all objects is "/".
- **JSON Web Token:** JSON Web Token for Google Cloud Storage is required for authenticating the account. Provide the JWT file content for establishing a secured connection with Google Cloud Storage. For more details on the JWT, see the Google documentation for [Using OAuth 2.0 for Server to Server Applications](#).

Web Proxy Configuration

- If you connect to Google Cloud Storage through a proxy server, these fields define the proxy details.
 - **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to Google Cloud Storage should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
 - **Proxy host:** The host name or IP address of the web proxy server.
 - **Proxy port:** The port on the proxy server for Data Source.
 - **Proxy username:** The username for the proxy server.
 - **Proxy password:** The password for the proxy server. *Leave username & password blank for an unauthenticated proxy connection.

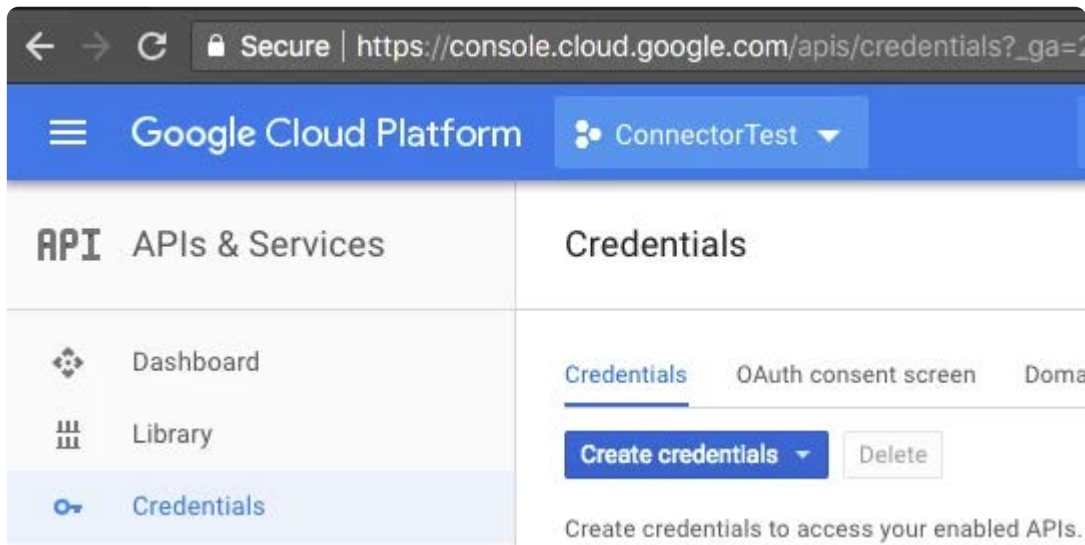
How to Authenticate with Google

The Data Prep Google Cloud Storage Connector leverages Service Account authentication.

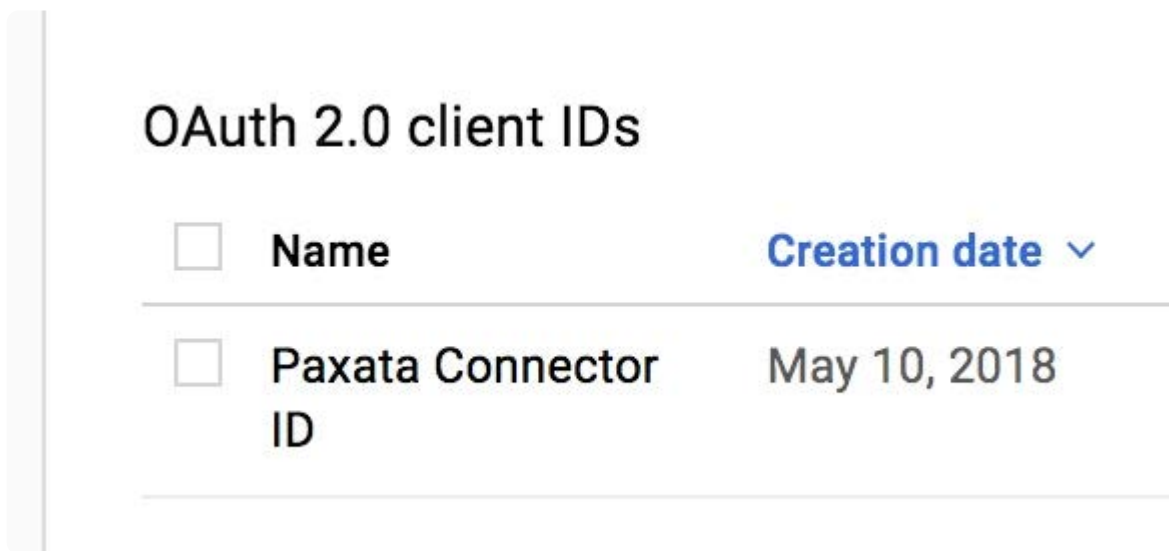
In order to access Google Cloud Storage using Data Prep you must:

1. Create a Google Service Account for the Cloud Storage service:
 - a. Open the list of credentials in the Google Cloud Platform Console: <https://console.cloud.google.com/apis/credentials>.
 - b. Click **Create credentials**.
 - c. Select **Service account key**.
 - d. In the **Create service account key** window, click the drop-down box below **Service account**, then click **New service account**.
 - e. Enter a name for the service account in **Name**.
 - f. Choose a [Cloud Storage Role](#) that grants the service account the desired level of access.
 - g. Use the default **Service account ID** or generate a different one.
 - h. Select the **Key type: JSON**.
 - i. Click **Create**.

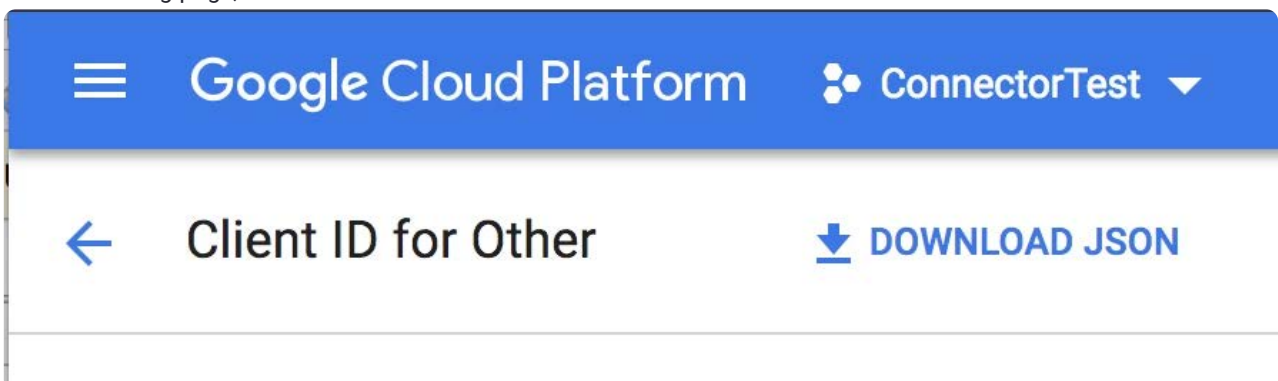
A **Service account created** window is displayed and the private key for the **Key type** you selected is downloaded automatically. Remember the downloaded credential location. 10. Click **Close**. 2. Download the JSON credential for an existing Service Account for the Cloud Storage service: 1. Log in to the Google Console using the end-user account: <https://console.cloud.google.com/apis/credentials>.
 - j. Ensure that the correct Project is selected in the dropdown list.



2. Scroll down to the "OAUTH 2.0 client IDs" section.



- a. Click the Name of the existing ID that you plan to use in the Connector.
- b. On the resulting page, click the "DOWNLOAD JSON" link.



- c. Remember the downloaded credential location.

For additional reference, please see: <https://cloud.google.com/storage/docs/authentication#generating-a-private-key>

Data Import Information

Via Browsing

Browse directories and files within the configured Bucket/Prefix.

Google Cloud SQL Connector for Data Prep

User Persona: Data Prep User or Data Prep Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configuring Data Prep

The Data Prep JDBC connectors ([MySQL](#), [PostgreSQL](#), and [SQL Server](#)) can be utilized to connect to Cloud SQL. See the [JDBC Connector Documentation](#) along with the [Cloud SQL documentation](#) for details that will enable you to configure the connection properly.

Google Drive Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Google Drive for browsing and importing available data. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Google Drive accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Google Drive Configuration

- **OAuth Verifier Key:** The verifier key used to authenticate with Google Drive. To obtain the verifier key, click "Test Data Source" and follow the link to grant access to Google Drive. After allowing access, you will be redirected to a page that displays an access code. Copy the code into this field.

Web Proxy

If you connect to Google Drive through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to Google Drive should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.

- **Proxy Host:** The hostname or IP address of the proxy server.
- **Proxy Port:** The port of the proxy server.
- **Proxy Username and Proxy Password:** User credentials for an authenticated proxy connection. Leave these blank for an unauthenticated proxy connection.

Import from Google Sheets

- When browsing to import a sheet, each Google Sheet is listed (identified with the prefix "GSheet") as a directory item (i.e. a folder), wherein you will find a separate data file for each sheet to import. For example, if you had one Google Sheet that had 5 separate sheets within, you use this browsing interface to select each of the sheets individually for import rather than selecting the workbook as a whole.
- It is important to note the following criteria:
 - The name of the Google Sheet to be imported cannot have the % character.
 - Be aware of [size limits for Google Drive](#) files.

[Connect to data sources for Data Prep](#) > Google Sheets Connector

Google Sheets Connector for Data Prep

The Google Sheets Connector has been deprecated. Please use the [Google Drive Connector](#), which has support for importing and exporting Google Sheets.

Hortonworks HDP2 HDFS Connector for Data Prep

User Persona: Data Prep Admin, Data Source Admin, or IT/DevOps

Availability information

This Connector is not available to Data Prep SaaS customers.

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configuring Data Prep

This connector allows you to connect to a Hortonworks HDP 2.6.5 Hadoop File System (HDFS) for import and export. The following fields are used to define the connection parameters.

Note

Configuring this Connector requires file system access on the Data Prep Server and a `core-site.xml` with the Hadoop cluster configuration. Contact your Customer Success representative for assistance with this step.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple HDFS clusters. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Simple Configuration (only for Simple authentication)

- **Username:** The application web server will connect to your HDFS cluster as the username you provide here.

Configuration

- **Data Store Root Directory:** The 'parent directory' on your cluster where the Connector will read from and write to for import and export operations. This also supports import and export for sub-directories of the root.

Kerberos Configuration

The following parameters are required for Kerberos authentication.

- **Principal:** Kerberos Principal.
- **Realm:** Kerberos Realm.
- **KDC Hostname:** Kerberos Key Distribution Center Hostname.
- **Kerberos Configuration File:** Fully-qualified path of Kerberos configuration file on webserver.
- **Keytab File:** Fully-qualified path of Kerberos Keytab File on webserver.

The Proxy User and Use Application User options allow you to specify the account to impersonate. See [this documentation](#) for more information about impersonation in HDFS. You have three options here: use a specific proxy user, a proxy user with modifiers, or the individual application user.

- **Proxy User:** Here you can either specify the user account that will be impersonated for all connections or check the Use Application User box to impersonate the user account of the individual Data Prep user who runs the connector. Note that the Proxy User field is not enabled if Use Application User is checked. Entering `${user.name}` as the proxy user works similarly to selecting Use Application User but allows for more flexibility because you can add modifiers or additional text. For example:
 - To add a domain to the user's credentials, enter `\domain_name\${user.name}` in the Proxy User field. Data Prep will pass the username and the domain.
 - Example: `\Accounts\${user.name}` results in `Accounts\Joe` (assuming Joe is the username).
 - To apply a text modifier to the username, add `.modifier` to the key `${user.name}`. The acceptable modifiers are: `toLower`, `toUpper`, `toLowerCase`, `toUpperCase`, and `trim`.
 - Example: `${user.name.toLowerCase}` converts Joe into joe (assuming Joe is the username).

Data Import Information

Via Browsing

Supported

Via SQL Query

Not supported

Hortonworks HDP2 Hive Connector for Data Prep

User Persona: Data Prep Admin, Data Source Admin, or IT/DevOps

Availability information

This Connector is not available to Data Prep SaaS customers.

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to a Hortonworks HDP 2.6.5 Hive for import and export. The following fields are used to define the connection parameters.

Note

Configuring this Connector requires file system access on the Data Prep Server and a `core-site.xml` with the Hadoop cluster configuration. Contact your Customer Success representative for assistance with this step.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Hive directories. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Hadoop Cluster

- **HDFS User:** The username on the HDFS cluster used to write files for export to Hive.

Kerberos Configuration

The following parameters are required for Kerberos Authentication.

- **Principal:** Kerberos Principal.
- **Realm:** Kerberos Realm.
- **KDC Hostname:** Kerberos Key Distribution Center Hostname.
- **Kerberos Configuration File:** Fully-qualified path of Kerberos configuration file on webserver.
- **Keytab File:** Fully-qualified path of Kerberos Keytab File on webserver.

The Proxy User and Use Application User options allow you to specify the account to impersonate. See [this documentation](#) for more information about impersonation in HDFS. You have three options here: use a specific proxy user, a proxy user with modifiers, or the individual application user.

- **Proxy User:** Here you can either specify the user account that will be impersonated for all connections or check the Use Application User box to impersonate the user account of the individual Data Prep user who runs the connector. Note that the Proxy User field is not enabled if Use Application User is checked. Entering `${user.name}` as the proxy user works similarly to selecting Use Application User but allows for more flexibility because you can add modifiers or additional text. For example:
 - To add a domain to the user's credentials, enter `\domain_name\${user.name}` in the Proxy User field. Data Prep will pass the username and the domain.
 - Example: `\Accounts\${user.name}` results in `Accounts\Joe` (assuming Joe is the username).
 - To apply a text modifier to the username, add `.modifier` to the key `${user.name}`. The acceptable modifiers are: `toLower`, `toUpper`, `toLowerCase`, `toUpperCase`, and `trim`.
 - Example: `${user.name.toLowerCase}` converts Joe into joe (assuming Joe is the username).

Hive Configuration

When you export data using the Hive connector, a file is written into HDFS and then an external table is created in Hive through the Hive JDBC driver. The Proxy User field specifies the user account to impersonate when writing a file into HDFS, but in order to do an impersonation in Hive, you must also specify the user in the JDBC URL.

- **JDBC URL:** The URL used to access Hive for import and registration of external tables. If Kerberos authentication is used, the following string must be added to the URL `";auth=kerberos;hive.server2.proxy.user=${user.name}`.
- If a proxy user is used, then the string `${user.name}` must be replaced with the proxy username.
- **Hive File Location:** The location within HDFS used to store Hive files for external tables.

Credentials

- **Hive User:** The username used to access Hive for Simple and Hybrid authentication.

- **Hive Password:** The password used to access Hive for Simple authentication.

Hive Options

- **Pre-Import SQL:** SQL to be executed before import process. This SQL may execute multiple times (for preview and import) and could be multiple SQL statements, newline-delimited.
- **Post-Import SQL:** SQL to be executed after import process. This SQL may execute multiple times (for preview and import) and could be multiple SQL statements, newline-delimited.

Note

As the Pre- and Post-Import SQL may be executed multiple times throughout the import process, take care when specifying these values in the Connector/Datasource Configuration as they will be executed for every import performed with this configuration.*

- **Pre-Export SQL:** SQL to be executed before export process. This SQL will execute once and could be multiple SQL statements, newline-delimited.
- **Post-Export SQL:** SQL to be executed after export process. This SQL will execute once and could be multiple SQL statements, newline-delimited.

Data Import Information

Via Browsing

Not Supported

Via SQL Query

Using SQL Select queries

HubSpot Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to HubSpot for browsing and importing available data. The following parameters are used to configure the connection.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI. form Data Prep DevOps how you would like this set.

Tip

You can connect Data Prep to multiple Hubspot accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Web Proxy

If you connect to HubSpot through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if connection to HubSpot should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for HubSpot.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server.*Leave username & password blank for an unauthenticated proxy connection.

HubSpot Configuration

- **OAuth Verifier Key:** The verifier key used to authenticate with HubSpot. To obtain the verifier key, click "Test Data Source" and follow the link to grant access to HubSpot. After allowing access, you will be redirected to `http://localhost:33333` (`http://localhost:33333/`, but a web page will not be displayed. Copy the value of the "code" URL parameter from the URL. This is the verifier key used to authenticate the connector. Paste it into the verifier key field.

Data Import Information

Via Browsing

View list of predefined datasets and "Select" the dataset for import.

Via SQL Query

Using a legal SQL Select Query.

IBM DB2 Connector for Data Prep

User Persona: Data Prep Admin or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

The ability to connect to IBM DB2 is part of the Data Prep JDBC Connector. This topic provides details specific to setting up a connection to IBM DB2. See also the [JDBC Connector Documentation](#) to configure the connection.

Example JDBC URI:

```
jdbc:db2://db2._yourdb2instancedomain.com:yourDB2instanceport/Data  
Prep:ConnectionRetryCount=3;LoginTimeout=10;com.ibm.db2.jcc.DB2BaseDataSource.keepAliveTimeOut=20;
```

Technical Specs

Driver Specs

- **IBM DB2 Database driver name & version:**
 - Driver classname: com.ibm.db2.jcc.DB2Driver
 - Version: 11.5
- **Supported IBM DB2 Database versions:**
 - Supports JDBC 3 and JDBC 4 standard - All Software Editions

Driver Documentation

- **General Driver Documentation:** <https://www.ibm.com/support/pages/download-initial-version-115-clients-and-drivers>

IBM Netezza Connector for Data Prep

User Persona: Data Prep Admin or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

The ability to connect to IBM Netezza is part of the Data Prep JDBC Connector. This topic provides details specific to setting up a connection to IBM Netezza. See also the [JDBC Connector Documentation](#) to configure the connection.

Example JDBC URI:

```
jdbc:netezza://YourNetezzaIPAddress:YourNetezzaPortnumber/MYDB
```

Technical Specs

Driver Specs

- **IBM Netezza Database driver name & version:**
 - Version: 7.2.1.0
- **Supported IBM Netezza Database versions:**
 - 7.0.x, 7.1.x, 7.2.x

Driver Documentation

- **General Driver Documentation:** https://www.ibm.com/support/knowledgecenter/SSULQD_7.2.1/com.ibm.nz.datacon.doc/c_datacon_introduction.html
- **Version Compatibility** https://www.ibm.com/support/knowledgecenter/SSULQD_7.2.1/com.ibm.nz.datacon.doc/c_datacon_release_compatibility_matrix.html

JDBC Connector for Data Prep

User Persona: Data Prep Admin or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This Connector enables the ability to import and export data leveraging JDBC Drivers. Typically this connector is leveraged for import/export against relational databases, but many applications offer JDBC drivers. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple JDBC sources. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Database URI

- **JDBC URI:** The JDBC connection string as defined by the driver being used. For more details on connection string options, please see the documentation for the driver you intend to use. JDBC connection strings typically take the form:

```
jdbc:///;/;
```

Database Visibility

You can control the database, schemas, and tables that are shown to users when they browse a data source during import. For databases, schemas, and tables, you can choose to:

- "Show only" which returns only the databases, schemas or tables that you specify here.
- "Hide" which hides the databases, schemas, and tables that you specify here.
- "Show all" which is the default setting to display everything in the data source.

When you select the "Show only" or "Hide" options, a field is provided for specifying the databases, schemas or tables on which you want the option enforced.

Note

These settings are not enforced when users query against the data source; query results still return a complete list of matches. For example, if you choose to "hide" a specific database, users can still execute queries that pull data from tables within that database. However, that database will not be displayed to users when they browse the data source.

Import Configuration

- **Query Prefetch Size:** The number of rows per batch during import.

Note

A larger batch size will increase throughput for large imports, but setting this to too large a value can cause the Connector process to run out of memory. The default value for this field is 10,000 rows per batch. The amount of memory granted to a Connector process varies greatly depending on the size of your Data Prep installation and your Data Prep Administrator should be consulted before setting this field to anything greater than the default.

- **Max Column Size:** The maximum length, in Unicode characters, for any column. Values larger than this will be replaced with 'null'.
- **Pre-Import SQL:** SQL to be executed before import process. This SQL may execute multiple times (for preview and import) and could be multiple SQL statements, newline-delimited.
- **Post-Import SQL:** SQL to be executed after import process. This SQL may execute multiple times (for preview and import) and could be multiple SQL statements, newline-delimited.
- **Perform Count Query:** This selector allows users to prevent the Connector from executing count queries upon import, which can be very slow on some database tables. To disable counting the number of rows while importing, set this to "False".

Note

As the Pre- and Post-Import SQL may be executed multiple times throughout the import process, please take care when specifying these values in the Connector/Datasource Configuration as they will be executed for every import performed with this configuration.

Export Configuration

- **Export Batch Size:** The number of rows per batch during export.

Note

A larger batch size will increase throughput for large exports, but setting this to too large a value can cause the Connector process to run out of memory. The default value for this field is 10,000 rows per batch. The amount of memory granted to a Connector process varies greatly depending on the size of your Data Prep installation and your Data Prep Administrator should be consulted before setting this field to anything greater than the default.

- **Max VARCHAR Size:** The maximum width for a VARCHAR column.

Note

This connector will attempt to export a column using the CLOB type if the maximum VARCHAR width is exceeded and the database supports to CLOB type.

- **Automatically Create Table:** ENABLED | DISABLED
- **ENABLED:** Data Prep will automatically create a new table when exporting a data set. If a table exists, Data Prep will drop the existing table before creating a new table with the same name.
- **DISABLED:** Data Prep will not automatically create a new table when exporting a data set. The Connector assumes that a table exists that matches the exported dataset's name & format. The exported data will be appended to the existing table.
- **Pre-Export SQL:** SQL to be executed before export process. This SQL will execute once and could be multiple SQL statements, newline-delimited.
- **Post-Export SQL:** SQL to be executed after export process. This SQL will execute once and could be multiple SQL statements, newline-delimited.

Credentials

- **User:** Username used to access the data source.
- **Password:** Password used to access the data source.
- **Role:** Some applications allow you to specify a role when connecting. Enter your role value here if required, or this field may be left blank.

Data Import & Export Information

Import Via Browsing

- Browse Databases, Schemas and/or tables based on your configuration settings. "Select" the table for import.

Import Via SQL Query

- Requires a SQL Select Query for your database.
- Example: `SELECT * FROM "SAMPLE_DATA"."TPCH_SF1"."CUSTOMER"`

Export Information

- Browse Databases and/or Schemas based on your configuration settings
- There may be errors if you attempt to export directly under a catalog and not a schema when the database supports catalogs, schemas, and tables.
- (Optional) Edit the name used for the table name.

Tier 2 Support for JDBC

Non-SaaS customers of Data Prep are able to supply and install their own drivers in for use with the JDBC Connector. This functionality is not currently available for SaaS customers.

Data Prep has a test kit that can help assess whether a given JDBC driver will be a good fit for use in Data Prep. The Data Prep Customer Success team will run this test kit on your behalf.

Note

Success with this tool does not certify the supplied JDBC driver for official support. Drivers that pass tests using this tool are not guaranteed to work with current or future versions of the Data Prep application. The JDBC Test Kit is not a comprehensive JDBC test suite tool.

What does “Tier 2” mean?

If the JDBC Connector Test Kit passes all tests, Data Prep supports usage of the combination of:

- A specific Data Prep JDBC connector version and a specific Data Prep Core Server version. These numbers generally match the Test Kit version, but this is not always the case.
- A specific database/application version.
- A specific JDBC driver version.
- The datatypes tested.
- Only the successfully tested features - There are currently three: Import via query, Import via browse, and Export.

What else does this mean:

- Data Prep does NOT certify the database/application for JDBC or other connector connectivity.
- Data Prep does NOT explicitly or formally test against the database/application.
- Data Prep does NOT test new versions nor Service Packs for potential regressions of Tier 2 data sources.

- Data Prep does NOT make Tier 2 JDBC sources available in Data Prep Cloud. Only the certified (Tier 1) JDBC sources are available in Data Prep Cloud.

How do I set it up?

If you have a JDBC driver you would like to use with Data Prep, please follow these steps for using it under the JDBC Connector as a Tier 2 data source:

1. Start by reaching out to your Customer Success representative and provide the driver you would like to use, they can test the driver for you using the test kit.
2. If the test is successful, the Customer Success representative will help you install the driver in the correct directory on the Data Prep Core Server and add the driver to the registry.
3. From there, you can refer to your driver documentation and the details above to configure the JDBC Connector to work with your new driver.

Jira Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Jira for browsing and importing available data. The following parameters are used to configure the connection.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Jira accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Jira Configuration

- **Jira URL:** Your Jira URL in the format `https://your-site-name.atlassian.net`.
- **Username:** The email address of the user connecting to Jira.
- **Authentication Type:** Type of authentication you want to use, either password or API token.
- **Password:** Password for connecting to Jira.
- **API Token:** API token for connecting to Jira. This is used for Cloud Jira only. For assistance generating an API Token, please see this [Jira Documentation](#).
- **Timeout:** The number of seconds to wait until a timeout error cancels a running operation.

Web Proxy

If you connect to Jira through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to Jira should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for Data Source.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server. *Leave username & password blank for an unauthenticated proxy connection.

Data Import Information

Via Browsing

- Browse to a dataset and click the dataset name for import. Please notes that some datasets may not be accessible for all users. If an account doesn't have access to one of the pre-defined datasets, they will encounter an error when they try to preview or import it.

Via SQL Query

- Using a legal SQL Select Query.

Marketo Connector for Data Prep

User Persona: Data Prep User or Marketo Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Marketo as an import source. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Marketo accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source. If you are a Data Prep SaaS customer, inform Data Prep DevOps how you would like this set.

Marketo Configuration

- **Marketo REST Endpoint:** Your Marketo REST API Endpoint.
- **OAuth Client ID:** The Client ID of your service in Marketo.
- **OAuth Client Secret:** The Client Secret of your service in Marketo.
- **Timeout:** The number of seconds to wait until a timeout error cancels a running operation. The default is 60 seconds, for most situations this is adequate.

Web Proxy

If you connect to Marketo through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to the Marketo REST Endpoint should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for Data Source.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server.

Tip

Leave username and password blank for an unauthenticated proxy connection.

Configure Marketo

Marketo Connector leverages the Marketo REST API and OAuth. This step must be completed by a Marketo Admin.

In order to establish connectivity, the Marketo Administrator must:

- Create a "Custom Service" in Marketo.
- Obtain OAuth credentials for the "Custom Service" to identify the client (Data Prep Connector) that is accessing Marketo.

Create a "Custom Service" in Marketo

A custom service is required to connect to Marketo using the REST API. The below steps will guide you in creating a custom service.

1. Navigate to the **admin** area of your Marketo application.
2. Click **Users & Roles** in the Security section.
3. Select the **Roles** tab and click **New Role** to create a new Role.
4. Enter a **Role Name** and select the permissions for the Role. The **Access API** permissions are specific to the REST API.
5. Now that an API Role is created, select the Users tab and click **Invite New User**.
6. Enter the new user information and select the role that was just created with API access. The API Only option can be selected to denote the user as an API Only user.
7. Now that a new user has been created, a new service will need to be created. Click the **LaunchPoint** option (**Admin > Integration > LaunchPoint**).
8. Click **New Service**.
9. Select the Custom service type and enter a display name and description.
10. Select the user you created.

Obtain the OAuthClientId and OAuthClientSecret Values

To obtain the OAuthClientId and OAuthClientSecret, navigate to the LaunchPoint option on the Admin area. Click the View Details link for the desired service. A window containing the authentication credentials is displayed.

Obtain the REST Endpoint URL

The REST Endpoint can be found on your Marketo Admin area on the Integration -> Web Services option in the REST API section. Note the Identity Endpoint will not be needed.

Data Import Information

Via Browsing & Via SQL Query

See the following table for the list of visible Objects and example SQL queries.

MARKETO CONNECTOR DATA OBJECTS:

Some objects in Marketo may or may not be present based on your configuration and some may run into the 500 column import maximum (this is a default setting, please reach out to your Data Prep Customer Success contact if this is inadequate). There are also several SQL Query examples in the table below.

ELEMENT	DESCRIPTION
Activities	Custom Activities for your Marketo organization. <ul style="list-style-type: none">• Each custom activity contained within your Marketo organization will be returned as it's own object.• Each table name will be prefixed with 'Activity' followed by the name of your custom activity.
ActivityBulkExports	Returns a list of activity export jobs that were created in the past 7 days.
Campaigns	Campaigns for a Marketo organization.
Channels	Channels for a Marketo organization.
Companies	Companies for a Marketo organization. This object is only available for Marketo subscriptions which do not have a native CRM sync enabled.
CustomObjects	Custom objects for a Marketo organization.
Emails	Emails for a Marketo organization.

ELEMENT	DESCRIPTION
LeadBulkExports	List of lead export jobs that were created in the past 7 days.
LeadPartitions	Lead Partitions for a Marketo organization.
Leads	<p>Leads for a Marketo organization.</p> <p>Based on your organization's Marketo configuration, the Lead object may contain > 500 columns. You may receive a truncation error when attempting to import > 500 columns using the default Data Prep Library configuration.</p>
Lists	Lists for a Marketo organization.
NamedAccounts	<p>Named Accounts for a Marketo organization.</p> <ul style="list-style-type: none"> • This Marketo object can only be imported using SQL using the "Create Query" option. • Queries <i>must</i> include a filter that leverages the "=" operator. • Example: <code>SELECT * FROM NamedAccounts WHERE State='CA'</code>
Opportunities	<p>Opportunities for a Marketo organization.</p> <ul style="list-style-type: none"> • This table is only available for Marketo subscriptions which do not have a native CRM sync enabled. • Queries MUST include a filter that leverages the "=" operator. • Example: <code>SELECT * FROM Opportunities WHERE State='CA'</code>
OpportunityRoles	<p>Opportunity Roles for a Marketo organization.</p> <ul style="list-style-type: none"> • This table is only available for Marketo subscriptions which do not have a native CRM sync enabled. • Queries MUST include a filter that leverages the "=" operator. • Example: <code>SELECT * FROM OpportunityRoles WHERE ExternalOpportunityId='Opportunity1'AND LeadId='1'AND Role='MyRole'</code>
Programs	<p>Programs for a Marketo organization.</p> <ul style="list-style-type: none"> • Tag and Cost columns are not returned when browsing all Programs (such as performing a <code>SELECT * query</code>). These columns are only returned when filtering by a specific Program Id or Name. • Tag and Cost Columns are not returned in this case: <code>SELECT * FROM Programs.</code> • Tag and Cost Columns are returned in this case: <code>SELECT * FROM Programs WHERE Id='1001'.</code>

ELEMENT	DESCRIPTION
SalesPersons	<p>Sales Persons for a Marketo organization.</p> <ul style="list-style-type: none">• This table is only available if the SalesPerson API is enabled.• A filter must be specified when retrieving companies. Valid filters are any searchable columns that include Id, ExternalSalesPersonId, or Email.• Example: <code>SELECT * FROM SalesPersons WHERE ExternalSalesPersonId='sales@company.com'</code>
Tags	Tags for a Marketo organization

MicroStrategy Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to a MicroStrategy server for Library imports and exports.

The following is information on the parameters used to create the connector.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple MicroStrategy accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

MicroStrategy Configuration

In the MicroStrategy Configuration section, provide the information used to locate and connect to the MicroStrategy server.

- **Server Hostname:** You can use either the fully qualified hostname, including the domain name, or the IP address of the MicroStrategy Host.
- **Server Port:** The port number of the MicroStrategy server.
- **Use SSL:** Select whether to use Secure Sockets Layer for connecting to the MicroStrategy server.

Credentials

- **Username and Password:** The username and password used for authenticating to the MicroStrategy server.
- **Authentication Mode:** The authentication mode used to connect to MicroStrategy Server.

Export Configuration

- **Export Batch Size:** The batch size to be used while exporting the MicroStrategy server.

Web Proxy

If you connect to MicroStrategy through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to MicroStrategy should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy Host:** The hostname or IP address of the proxy server.
- **Proxy Port:** The port of the proxy server.
- **Proxy Username and Proxy Password:** User credentials for an authenticated proxy connection. Leave these blank for an unauthenticated proxy connection.

Import/Export Information

- MicroStrategy does not support boolean data-types. Hence, while exporting from Data Prep, all boolean values will be exported as strings ("true" or "false").
- **Exporting reports:**
 - MicroStrategy reports can be imported, but the MicroStrategy Connector does not support the ability to update a report upon Export.
- **Exporting to cubes:**
 - The MicroStrategy Connector does not support exporting data to a cube with more than one table.
 - The MicroStrategy Connector utilizes REST API, thus it does not support updating cubes upon export that were made using the MicroStrategy Web UI (or by any means other than the REST API).
 - Exporting to an existing cube will not be possible if the schema (in the Data Prep AnswerSet) does not match the schema of the cube.
 - When a column is added in the AnswerSet, new or updated rows will be appended to the cube but the new column will not be added to the cube.
 - To overwrite a cube that has dependencies, ensure that MicroStrategy permissions are set up to allow overwriting of dependent cubes.

MongoDB Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to MongoDB for browsing and importing available data. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple MongoDB deployments. Using a descriptive name can be a big help to users in identifying the appropriate data source.

MongoDB Configuration

- **Server:** The hostname or IP address of the server hosting the MongoDB instance. If connecting to a replica set, use the hostname or IP address of one of the servers.
- **Port:** The port for connecting to MongoDB. The default port is 27017.
- **User:** The MongoDB user.
- **Password:** The MongoDB user's password.
- **Use SSL:** This field sets whether SSL is enabled.
- **Timeout:** The number of seconds to wait until an operation times out. If set to 0, operations never time out.

Database Type Configuration

- **Database Type:** Connect to a standalone MongoDB instance or to a replica set
- **Database Name:** MongoDB database name. The database name is required when connecting to a standalone instance. When connecting to a replica set, if this field is left blank, all available databases will be displayed in the import UI
- **If connecting to a replica set:**
 - **Replica Set:** A comma-separated list of secondary servers (server:port). This allows you to specify multiple servers in addition to the one configured in Server and Port
 - **Read Preference:** Strategy for reading from a replica set. See [Read Preferences](#) for more details.

Data Import Information

Via Browsing

Browsing will allow you to view a list of databases and collections within MongoDB. If connecting to a MongoDB replica set, all databases can be browsed if no database was specified in the configuration. If a database was specified, only the collections within that database can be browsed.

Via SQL Query

This Connector is built on top of JDBC, using a driver provided by CData, and therefore allows data to be imported using SQL SELECT queries. Note that queries use SQL as documented in [CData's driver documentation](#), not the JavaScript-based DSL used in the MongoDB shell.

Best Practices

- Data Prep uses MongoDB as its metadata store. Using this Connector to connect to Data Prep's own MongoDB replica set is not an intended or recommended use case for this Connector. We cannot guarantee the performance and correct functioning of the connector and of Data Prep itself when it is used in this way.
- If your intent is to read Data Prep metadata, please periodically create a backup of the data in Data Prep MongoDB metadata store and restore the backup in a separate instance of MongoDB, then point the Connector at that instance.

MS Azure Data Lake Storage (ADLS) Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Azure Data Lake Storage (ADLS) for imports and exports. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Azure Data Lake Storage accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Azure Data Lake Storage Configuration

- **ADL URI:** The URI for the ADL site.
- **Root Directory:** Specifies the top-level of the directory structure from which import/export of data is enabled.
- **Application ID:** The application ID for the ADL site.
- **OAuth 2.0 Token Endpoint:** The OAuth 2.0 Token Endpoint for the ADL site.
- **Application Access Key Value:** The Application Access Key Value for the ADL site. See Issue 1 in the FAQ/Troubleshooting/Common Issues section for more information.

FAQ/Troubleshooting/Common Issues

Can we have both ADLS Gen1 and ADLS Gen2 Connectors in the same Data Prep account?

Yes. The two Connectors can coexist and will not interfere with each other.

Issue: When you Test Connection, it fails and reports a "base64" issue.

How to fix it: In March 2020, Azure changed the format of the **Application Access Key Value**. The new format does not work for authentication, so you'll need to use the Azure command line to set the Base64 encoded version of the **Application Access Key Value**.

In the Azure Portal:

1. Create a new Application Access Service Account.
2. Copy the generated Access Key value.
3. Base64 encode the password.
4. Mac example: `echo -n " | openssl base64`
5. Windows: Use a tool like Base64 Encoder.
6. Create an Azure Command to reset the password to a Base64 encoded version of the password:
7. `az ad sp credential reset --name --credential-description "" --append --years 2 -p "" -o=jsonc`
8. Open a Command prompt in the Azure Portal and paste the command from Step 4.
9. Ensure that this service account has appropriate ACLs for the storage. You will receive ACL errors if permissions are not appropriate.
10. Configure **Application Access Key Value** using the new Base64 encoded password.

MS Azure Data Lake Storage Gen2 (ADLS Gen2) Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Azure Data Lake Storage Gen2 for import and export. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Azure Data Lake Storage Gen2 accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Azure Data Lake Storage Gen2 Configuration

- **Data Store Root Directory:** The apparent root path accessible by this connector. Use "/" to access all files in the file system.
- **Azure Storage Account Name:** The Subdomain Name of your unique Azure URL. Storage account names must be between 3 and 24 characters in length and may contain numbers and lowercase letters only. Your storage account name must be unique within Azure. No two storage accounts can have the same name.
- **File System Name:** The name of the file system within the storage account. This is sometimes called the "container" name.

Azure Data Lake Storage Gen2 Authentication Settings

From the drop-down, select the preferred authentication method for ADLS Gen2 storage and fill out the required fields.

- **Storage Account Access Key:** Enter the Storage Account Access Key in the field. This is sometimes referred to as a “Shared Key”.
- **Active Directory Username/Password:** Enter the Azure Directory username and password associated with your account.

Note

You must grant access for Data Prep to read and write data within your Microsoft account, otherwise, you will get an error while attempting to connect. To grant access, click on the 'Test Data Source' button in the Data Source set-up panel and follow the 'Grant Access' link. This will bring you to your Microsoft account where you can log in and grant access. Then, come back to Data Prep to continue.

Data Import Information

Via Browsing

The connector will present a browsable directory hierarchy starting at the location defined in the Data Store Root Directory field.

Via SQL Query

Not Supported

FAQ/Troubleshooting/Common Issues

Can we have both ADLS Gen1 and ADLS Gen2 connectors in the same Data Prep account?

Yes. The two connectors can coexist and will not interfere with each other.

MS Azure SQL Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

The Data Prep JDBC Connector allows you to connect to the Microsoft Azure SQL Database. This topic provides details for setting up a connection to Azure SQL Database. See also the [JDBC Connector Documentation](#) for guidelines on configuring the connection.

Example JDBC URI

```
jdbc:sqlserver://
```

```
serverName.database.windows.net:serverPortNumber;encrypt=true;trustServerCertificate=false;hostNameInCertificate=*.database.windows.net;loginTimeout=30;data
```

Example JDBC URI for multi-factor authentication

Following is an example of a JDBC URI for configurations that use Active Directory multi-factor authentication:

```
jdbc:sqlserver://
```

```
<serverName>.database.windows.net:<serverPortNumber>;encrypt=true;trustServerCertificate=false;hostNameInCertificate=*.database.windows.net;loginTimeout=30;
```

See [Set up Active Directory multi-factor authentication](#) for more information.

Technical Specs

Driver Specs

Microsoft Azure SQL Database driver name and version:

- Driver Classname: `com.microsoft.sqlserver.jdbc.PxMSSQLDriver`
- Version: 9.2.1.jre8

Supported Microsoft Azure SQL Database versions:

- Azure SQL Database
- Azure SQL Managed Instance (Extended Private Preview)

Set up Active Directory multi-factor authentication

The MS Azure SQL Connector supports Active Directory (AD) multi-factor authentication (MFA), allowing users to authenticate using two or more verification factors interactively.

To support MFA interactively, add the following connection properties to the JDBC URL:

```
authentication=ActiveDirectoryInterActive  
clientid=<azure_registered_app's_client_id>
```

Register the multi-factor authentication app

To set up multi-factor authentication, you need to register an app that allows the Microsoft identity platform to provide authentication and authorization services for the Data Prep (Paxata) MFA application and its users.

Tip

To set up the app in the Microsoft identify platform, follow the guidelines below and also refer to the [Quickstart: Register an application with the Microsoft identity platform](#).

The application client ID from the MFA app is required in the JDBC URL so that once you complete the MFA process, the Microsoft identity platform will redirect you to the Data Prep (Paxata) instance's endpoint (specified in the [Redirect URI format](#) section).

During app registration on the Microsoft identify platform, you must do the following:

- Ensure that the app is authorized to call the Azure SQL Database API when an AD user grants permission. The following shows an example of the settings in the Microsoft identity platform:

Search (Cmd+/) Refresh Got feedback?

Overview
Quickstart
Integration assistant

Manage

Branding
Authentication
Certificates & secrets
Token configuration
API permissions
Expose an API
App roles
Owners
Roles and administrators | Preview
Manifest

Support + Troubleshooting
Troubleshooting
New support request

Configured permissions

Applications are authorized to call APIs when they are granted permissions by users/admins as part of the consent process. The list of configured permissions should include all the permissions the application needs. [Learn more about permissions and consent](#)

+ Add a permission ✓ Grant admin consent for Paxata

API / Permissions name	Type	Description	Admin consent requ...	Status
▼ Azure SQL Database (1)				
user_impersonation	Delegated	Access Azure SQL DB and Data Warehouse	No	...
▼ Microsoft Graph (1)				
User.Read	Delegated	Sign in and read user profile	No	✓ Granted for Paxata ...

To view and manage permissions and user consent, try [Enterprise applications](#).

user_impersonation	Delegated	Access Azure SQL DB and Data Warehouse	No	...
▼ Microsoft Graph (1)				
User.Read	Delegated	Sign in and read user profile	No	✓ Granted for Paxata ...

To view and manage permissions and user consent, try [Enterprise applications](#).

- On the **Platform configurations** page, under **Mobile and desktop applications**, set up the redirect URI. After completing MFA with the Azure SQL database, the Microsoft identity platform uses the URI to redirect and send security tokens from your client to the Data Prep (Paxata) application.

Overview
Quickstart
Integration assistant

Manage

Branding
Authentication
Certificates & secrets
Token configuration
API permissions
Expose an API
App roles
Owners
Roles and administrators | Preview
Manifest

Support + Troubleshooting
Troubleshooting
New support request

Platform configurations

Depending on the platform or device this application is targeting, additional configuration may be required such as redirect URIs, specific authentication settings, or fields specific to the platform.

+ Add a platform

Mobile and desktop applications Quickstart Docs ?

Redirect URIs

The URIs we will accept as destinations when returning authentication responses (tokens) after successfully authenticating users. Also referred to as reply URLs. [Learn more about Redirect URIs and their restrictions](#)

- https://login.microsoftonline.com/common/oauth2/nativeclient
- https://login.live.com/oauth20_desktop.srf (LiveSDK)
- msal12732c27-4306-44fb-862e-fb0d480dd2bd//auth (MSAL only)
- https://pax-installation-ad-mfa-rajeev-kumar-dev-eks.paxata.ninja/interactive-msal-token
- https://pax-installation-subhabrata-rajeev-kumar-dev-eks.paxata.ninja/interactive-msal-token
- http://localhost:8080/interactive-msal-token

Add URI

Note

You need a separate app registration per instance to manage your organization's account. If your organization has multi-tenant access to Data Prep, you can have one MFA app for all of the tenants. You just need to add the redirect URIs for each tenant on the **Platform configurations** page.

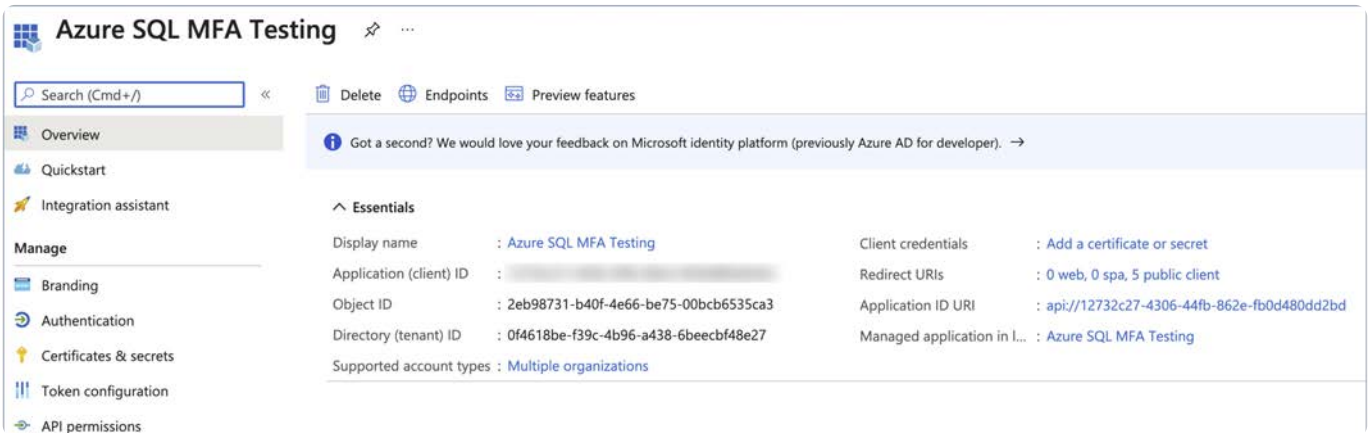
Redirect URI format

http(s)://<paxata_instance_host_name>/interactive-msal-token

Example Redirect URI

https://datarobot.paxata.com/interactive-msal-token

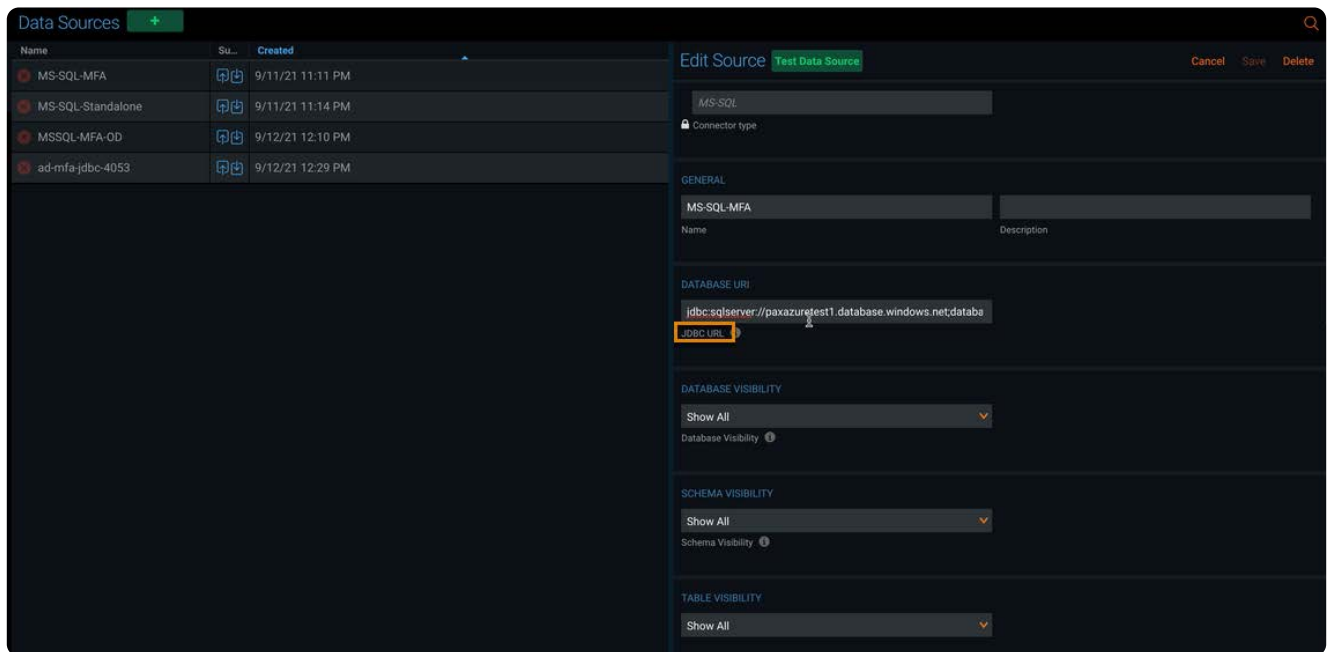
Once you have set up the app, the configured permissions look similar to the following:



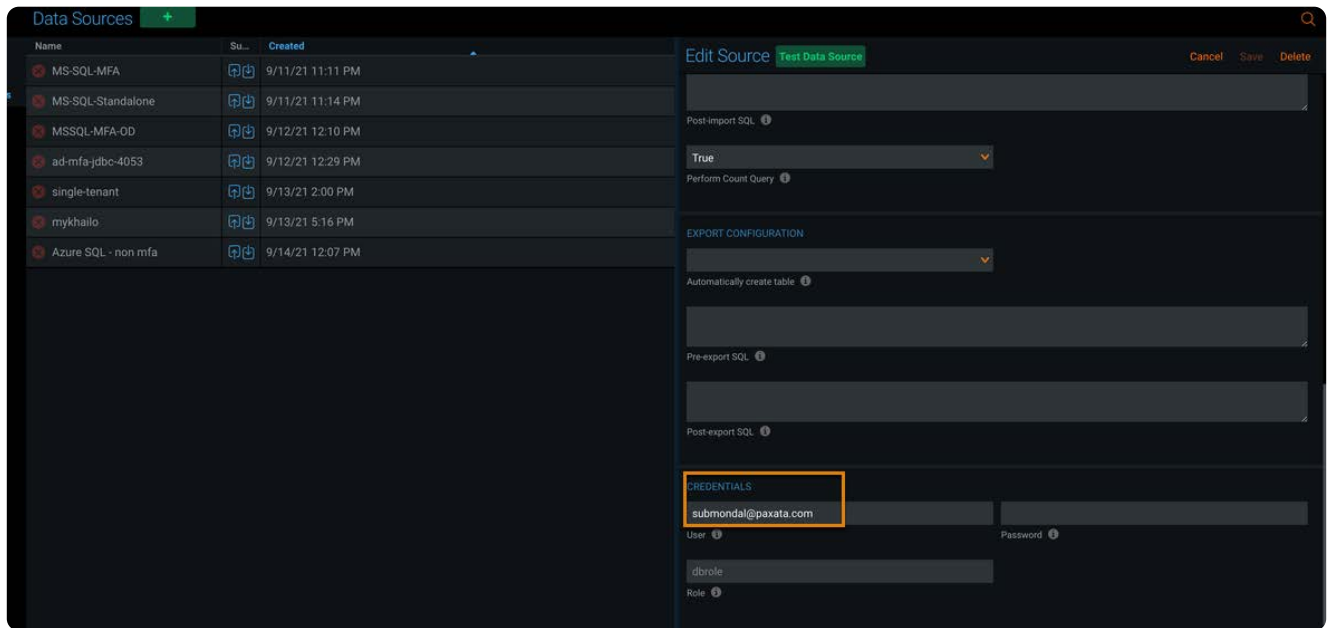
Multi-factor authentication process

1. Configure the JDBC URL on the Data Prep (Paxata) **Data Sources** page with these settings:

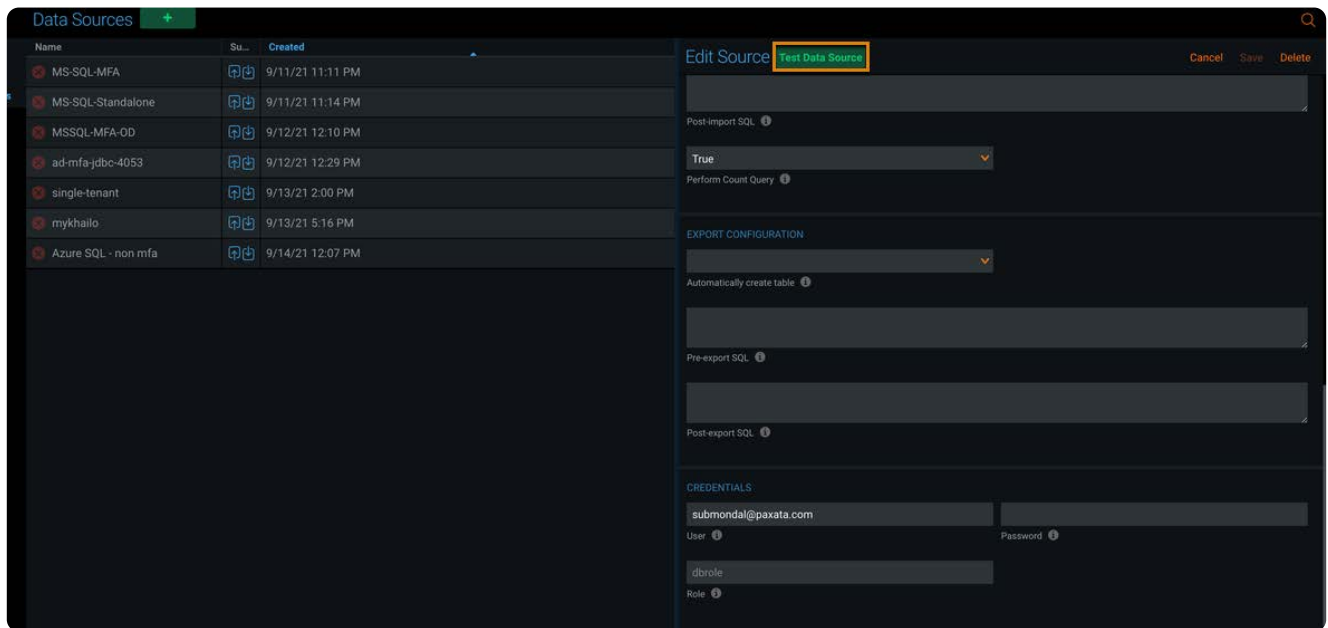
```
authentication=ActiveDirectoryInterActive  
clientid=<azure_registered_app's_client_id>
```



2. For **Credentials**, enter only the Active Directory user name. The password is not needed.



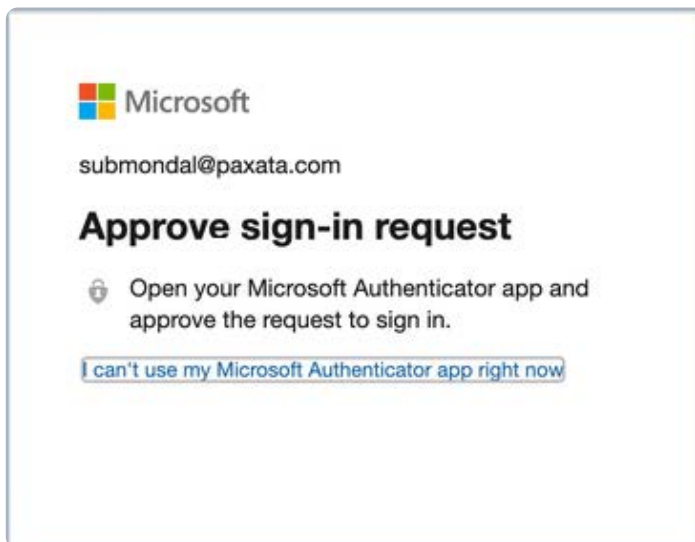
3. Click **Test Data Source**.



4. A new tab opens in the browser with a login screen set to the Active Directory user name and password. Click **Sign in**.



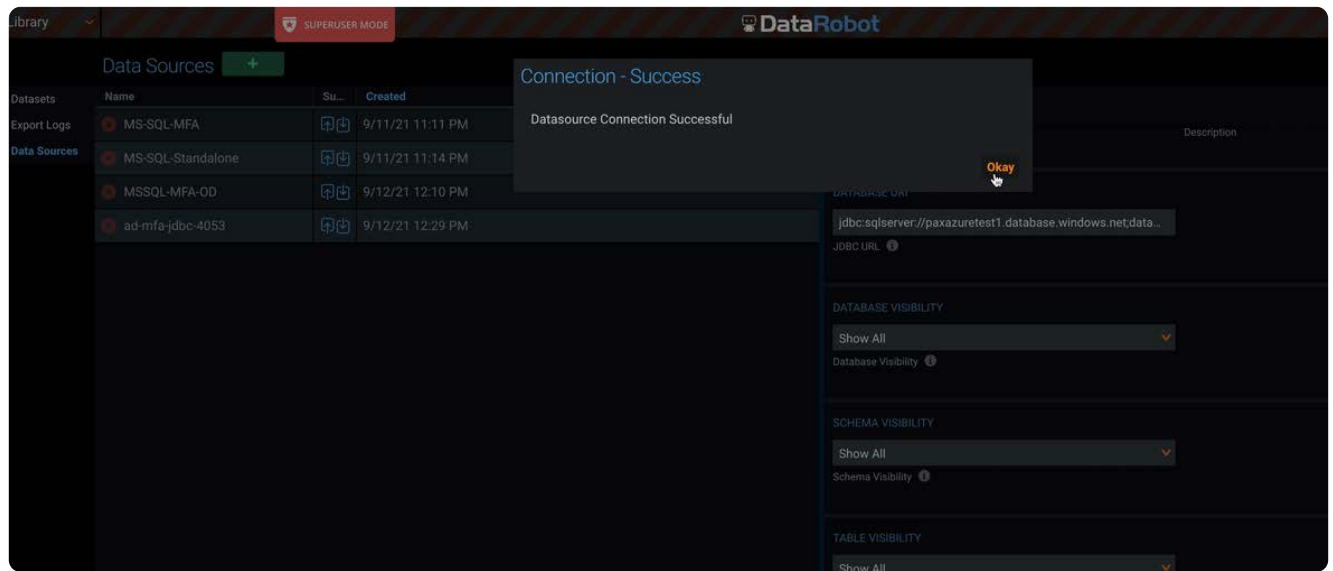
5. Once your user name and password are authenticated, perform the second verification factor and any additional required verification factors (the verification factors are based on how your MFA account was set up).



6. Once you verify all factors, the Microsoft identity platform redirects you to the Data Prep application.



7. Close the window and return to the **Data Source** configuration page.



Note

If the MFA authentication process has not completed within two minutes, you will receive a message that either the wrong client ID was specified in the JDBC URL, the authentication token has not been received yet, or you do not have access to the requested resource.

Driver Documentation

- General Driver Documentation: <https://docs.microsoft.com/en-us/sql/connect/jdbc/using-the-jdbc-driver>
- Driver Compatibility: <https://docs.microsoft.com/en-us/sql/connect/jdbc/microsoft-jdbc-driver-for-sql-server-support-matrix>

Additional Documentation

- General information on connecting to Azure SQL Database: <https://docs.microsoft.com/en-us/azure/sql-database/sql-database-connect-query-java>

MS Azure Synapse Analytics Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to an Azure Synapse Analytics for Library imports and exports. For export, the connector uploads data into Azure Data Lake service and then exposes the data as a table in the SQL Data Warehouse.

The following is information on the parameters used to create the connector.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Azure Synapse Analytics accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source. If you are a Data Prep SaaS customer, inform Data Prep DevOps of how you would like this set.

Database URL

The Database URL is where you provide the Java Database Connectivity (JDBC) connection string. This string tells Data Prep where the database is for import and export. You can include the schema name in the URL.

Visibility Settings

You can control the database, schemas, and tables that are shown to users when they browse a data source during import. For databases, schemas, and tables, you can choose to:

- **"Show only"** which returns only the databases, schemas or tables that you specify here.
- **"Hide"** which hides the databases, schemas, and tables that you specify here.
- **"Show all"** which is the default setting to display everything in the data source.

When you select the "Show only" or "Hide" options, a field is provided for specifying the databases, schemas or tables on which you want the option enforced.

Note

These settings are not enforced when users query against the data source; query results still return a complete list of matches. For example, if you choose to "hide" a specific database, the query results will still include that database in the returned results. However, that database will not be displayed to users when they browse the data source.

Import Configuration

The import configurations let you specify how data is imported into Data Prep.

- **Query Fetch Size:** Number of rows per batch to use when fetching data on import.
- **Pre-import SQL:** A SQL statement to execute before beginning import, after the table's schema is determined.
- **Post-import SQL:** A SQL statement to execute after import completes.

Export Configuration

The export configurations let you specify how data is exported out of Data Prep.

- **Pre-export SQL:** a SQL statement to execute before beginning export, after the table is created if auto-create is enabled.
- **Post-export SQL:** a SQL statement to execute after export completes.
- **External Data Source Name:** Name of the [external data source in SQL Data Warehouse](#) for accessing data in Azure Data Lake.

Note

The maximum VARCHAR size is 8000 characters. Values exceeding 8000 characters are exported as an empty string.

Credentials

The credentials settings allow you to specify a single user account to authenticate with when connecting to the data source.

Azure Data Lake Configuration

The Azure Data Lake Configuration is where you provide the settings Data Prep needs to connect to your Azure Data Lake.

- **ADL URI:** The URI for the ADL site.
- **Root Directory:** Specifies the top-level of the directory structure from which import/export of data is enabled.
- **Application ID:** The application ID for the ADL site.
- **OAUTH 2.0 Token Endpoint:** The OAUTH 2.0 TOKEN ENDPOINT for the ADL site.
- **Application Access Key Value:** The Application Access Key Value for the ADL site.

Data Import Information

Via Browsing

Browse to a table and "Select" the table for import.

- Supported data formats:
 - Delimited datasets: comma, tab...
 - XML
 - JSON
 - Excel: Xls and XLSX
 - Avro
 - Parquet
 - Fixed format
 - Export
 - Only one export format is available, JDBC
 - (Optional) Edit the name used for the Snowflake table name

Via SQL Query

Using a legal SQL Select Query

MS Dynamics 365 Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to a Microsoft Dynamics 365 resource for Entity Set imports. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Microsoft Dynamics 365 resources. Using a descriptive name can be a big help to users in identifying the appropriate data source. If you are a Data Prep SaaS customer, inform Data Prep DevOps of how you would like this set.

Microsoft Dynamics 365 Settings

In this section, provide the information used to identify and connect to the Microsoft Dynamics 365 resource. These fields are required.

- **Tenant Domain Name/ID:** This is your Microsoft Azure Active Directory Tenant ID or Domain Name.
- **Resource URL:** This is your Microsoft Dynamics 365 resource URL.

Proxy Settings

If you connect to MS Dynamics 365 through a proxy server (optional), these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if connection to the MS Dynamics 365 REST Endpoint should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for Data Source.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server.*Leave username & password blank for an unauthenticated proxy connection..

Authentication Settings

In this section, provide the information used to authenticate and authorize the Microsoft Dynamics 365 resource. These fields are required.

- **Client ID:** Application ID of the registered application (Web App/API or Native) in Microsoft Azure Active Directory. For information on how to register your application please see [this documentation](#).
- **Authentication Type:** Select Client Credential if your application is registered in Azure AD as Web App/API application. Otherwise, select User Credential if your application is registered in Azure AD as Native application.
 - **Client Secret:** If Authentication Type is selected as Client Credential, this needs to be provided for authentication.
 - **Username and Password:** If Authentication Type is selected as User Credential, these need to be provided for authentication.

Data Import Information

Via Browsing

View list of CRM objects available for import.

- Select object name from the list to enable import.
- The Dynamics Connector will automatically handle paginated results from the Dynamics 365 Web API.

Via SQL Query

Not supported

MS SharePoint Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to a SharePoint site for Library imports and exports of files and SharePoint Lists. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Sharepoint sites. Using a descriptive name can be a big help to users in identifying the appropriate data source.

SharePoint Configuration

- **SharePoint site URL:** The URL for the SharePoint site.
- To create a connector for a specific SharePoint site, include “/sites/” in the URL. Example: `https://acme.sharepoint.com/sites/sales-department`
- For a SharePoint on a company's internal network only, you may need to add the port to the URL. Example: `http://sp.your-organization.com:8080/sites/department`
- **SharePoint Edition:** The edition of SharePoint being used, either online or on-premise.
 - For SharePoint Online, you will need to select **Simple Authentication** or **Single Sign-On** as the authentication mode.
 - For Single Sign-On, you will also need to provide the **SSO Domain** of the user being authenticated.
- **Username:** The username for an individual or shared account used to authenticate with SharePoint.

- For SharePoint Online, this is typically in the form of an email address. For a SharePoint on a company's internal network only, you may need to provide the domain with the username. For example: Accounts/JDoe
- **Password:** The password used to authenticate with SharePoint.

Web Proxy

If you connect to SharePoint through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to SharePoint should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The hostname or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server. *Leave username & password blank for an unauthenticated proxy connection.

Data Import Information

Browse

- The Connector will present a browsable directory hierarchy of files and Lists.
- You can expect the hierarchy to contain datasets similar to those presented in the Site's "Site Contents" page.

Query

- Not Supported.

Import

- File Import: Supported.
- Sharepoint List Import: Supported.

Export

- File Export: Supported.
- List Export: Supported.

FAQ/Troubleshooting/Common Issues

Issue #1

- **Problem:** After import, your data displays HTML tags in your Data Prep columns, ("

" "

", "", etc)

. Any field set as "Enhanced Rich Text" will import this way because SharePoint is storing HTML in the field. SharePoint renders this HTML as formatted text in the List view, but when Data Prep receives the raw text from the database, it has the HTML tags and is not formatted text.

- **Solution:** In SharePoint, go to the list you're importing. Go to the options page for the column in question. Here you can specify the type of text to allow. Set this to "Plain Text".

Issue #2

- **Problem:** After import, a column is missing from your Data Prep columns.
- **Solution:** In SharePoint Lists, the "Title" column is a required, default column with embedded configuration parameters. If the "Title" column is renamed, then this results in a mismatch between the View Name and Database Name for the column. This mismatch prevents Data Prep from importing the column. Either of the following options will allow Data Prep to import the column:
 - (a) In SharePoint, reset the column name back to "Title" and import it again.
 - (b) If you need a different column name, before importing your SharePoint, copy the data from the "Title" column into a new column with the name you want and hide the "Title" column.

MS SQL Server Connector for Data Prep

User Persona: Data Prep Admin or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

The ability to connect to Microsoft SQL Server is part of the Data Prep JDBC Connector. This topic provides details specific to setting up a connection to Microsoft SQL Server. See also the [JDBC Connector Documentation](#) to configure the connection.

Example JDBC URI:

```
jdbc:sqlserver://cs-mssql1-db1.yourMsSqlHost:yourMssqlPort
```

Technical Specs

Driver Specs

- Microsoft SQL Database driver name & version:
 - Driver Classname: com.microsoft.sqlserver.jdbc.SQLServerDriver
 - Version: 7.4
- Supported Microsoft SQL Database versions:
 - Microsoft SQL Server 2019
 - Microsoft SQL Server 2017
 - Microsoft SQL Server 2016
 - Microsoft SQL Server 2014
 - Microsoft SQL Server 2012

Driver Documentation

- General Driver Documentation: <https://docs.microsoft.com/en-us/sql/connect/jdbc/using-the-jdbc-driver>

Driver Compatibility: <https://docs.microsoft.com/en-us/sql/connect/jdbc/microsoft-jdbc-driver-for-sql-server-support-matrix>

MS Windows Azure Blob Storage (WASB) Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Azure Blob Storage Account for Library imports and exports. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Azure Blob accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Azure Blob Storage Configuration

- **Data Store Root Directory:** The root path on the Data Store accessible by this connector. Use "/" to access all files in a Container.
- **Azure Storage Account Name:** A storage account name may contain lowercase letters and numbers.
- **Blob Service Container Name:** A container organizes a set of blobs, similar to a folder in a file system. All blobs reside within a container.
- **Map INT96 to Datetime:** For INT96 fields, convert to Datetime values during import. Specifically, this allows Data Prep to read Parquet files written by Impala.

Azure Blob Storage Authentication Settings

- **Authentication Type:** Two authentication methods are supported: Shared Key and Shared Access Signatures. See [Azure Storage Services Authentication Types](#) for more details.
- **Shared Key:** Authenticate with your account access key.
- **Shared Access Signatures (SAS):** Authenticate with shared access signatures (SAS) token.

Data Import Information

Via Browsing

Supported

Via SQL Query

Not Supported

MySQL Connector for Data Prep

User Persona: Data Prep Admin or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

The ability to connect to MySQL is part of the Data Prep JDBC Connector. This topic provides details specific to setting up a connection to MySQL. See also the [JDBC Connector Documentation](#) to configure the connection.

Example JDBC URI:

```
jdbc:mysql://yourMysqlSubdomain.yourMysqlDomain.com:yourMysqlPort/yourDatabaseName
```

Technical Specs

Driver Specs

- **MySQL Database driver name & version:**
 - Driver classname: com.mysql.jdbc.Driver
 - Version: mysql-connector-java-5.1.39-bin
- **Supported IBM DB2 Database versions:**
 - For supported versions, see: <https://dev.mysql.com/doc/connector-j/5.1/en/connector-j-versions.html> [dev.mysql.com]

Driver Documentation

- **General Driver Documentation:** <https://dev.mysql.com/doc/connector-j/5.1/en/> [ibm.com]

NetSuite Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to NetSuite for browsing and importing available data. The following parameters are used to configure the connection.

General

- **Name:** Name of the data source as it will appear to users in the UI
- **Description:** Description of the data source as it will appear to users in the UI

Tip

You can connect Data Prep to multiple Netsuite accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

NetSuite Configuration

- **NetSuite Account ID:** Your NetSuite account ID. You can find this by logging in to NetSuite, clicking the Setup tab → Integration → Web Services Preferences.
- **Username:** The user of the NetSuite account used to authenticate.
- **Password:** The password of the NetSuite user.
- **Role ID:** The Internal ID of the role that will be used to connect to NetSuite. Leave this empty to use the user's default role. To find the Role ID, follow the steps above for finding the Account ID, then select the username from the drop-down list, and finally select the corresponding Role associated with the user account.
- **Timeout:** The number of seconds to wait for an operation, the default value is 300 seconds.

Web Proxy Configuration

If you connect to Data Source through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to NetSuite should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The hostname or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for Data Source.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server.

Note

Leave username and password blank for an unauthenticated proxy connection.

Data Import Information

Via Browsing

- Browse to an Object and "Select" the table for import.

Via SQL Query

- Using a legal SQL Select Query.

Network Share SMB Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to a network share using Server Message Block (SMB) protocol for imports and exports. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple SMB shares. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Configuration

- **Share Hostname:** is the server hostname.
- **Share Port:** is the server port number. The default SMB port is 445.
- **Share Name:** is the name of the share. The name of the share that the user intends to use for import or export. The Share Name is not a path and must not contain ` ` characters. Space characters are allowed.

Credentials

User authentication can be done through a shared account or an individual account. If credentials are not configured with the Data Source, the user will be prompted for credentials.

- **Username:** The username used to authenticate on the share.
- **Password:** The password used to authenticate on the share.
- **User Domain:** The user's domain for connecting to the share.
 - Also known as the SMB Workgroup.
 - In an AD/LDAP-managed account, this is the AD domain to which the account belongs.
 - Domain account structure: \
 - For non AD/LDAP-managed account, Domain can be blank.

Timeouts

- **Read/Write Timeout:** The timeout (in seconds) for reading/writing from/to the share.

Data Import Information

Via Browsing

Browse directories and files within the configured Share.

Via SQL Query

As SMB is a file store, SQL Queries are not supported for this data source.

FAQ/Troubleshooting/Common Issues

Data Prep tests this Connector against Windows Shares (read-only and read/write) using LDAP and Linux Shares using local server accounts. There are many ways to host and configure an SMB service and this sometimes can cause issues establishing a connection. When there are issues, please ask your administrator for a copy of the configuration file for the server (with sensitive values removed) and a copy of the log output for the server as well.

Oracle Connector for Data Prep

User Persona: Data Prep User or Data Prep Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

The ability to connect to Oracle is part of the Data Prep JDBC Connector. This topic provides details specific to setting up a connection to Oracle Database. See also the [JDBC Connector Documentation](#) to configure the connection.

Example JDBC URI:

```
jdbc:oracle:thin:yourOracleHost:yourOraclePort:XE
```

Technical Specs

Driver Specs

- **Oracle Database driver version:** 19.3
- **Supported Oracle Database versions:** 12.2.0.1, 18c, 19c

Driver Documentation

- **General Driver Documentation:** <https://www.oracle.com/database/technologies/appdev/jdbc.html>
- **Driver FAQ:** <https://www.oracle.com/database/technologies/faq-jdbc.html>

FAQ/Troubleshooting/Common Issues

Data Types Handling Troubleshooting

One rare item that has the potential to cause issues during export to an Oracle Database relates to a differing in how datatypes are handled between Data Prep and Oracle Database. Data Prep is designed to handle columns with mixed data types, where

Oracle is not. Data Prep determines a column data type based on the predominant data type present in the column. This becomes an issue when Data Prep identifies the datatype in a column is boolean and exports the column as a boolean even though it still contains mixed datatypes. This should be addressed in Data Prep before exporting to Oracle.

Oracle Marketing Cloud (Eloqua) Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Oracle Marketing Cloud for library imports. The following parameters are used to configure the connection.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Oracle Marketing Cloud accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Oracle Marketing Cloud Configuration

- **Company:** The company for which the Oracle Marketing Cloud account is registered
- **User:** The user of the Oracle Marketing Cloud account.
- **Password:** The password of the Oracle Marketing Cloud user.
- **Timeout:** The number of seconds to wait for an operation. The default value is 300 seconds. To allow for a longer wait time, increase this value. If set to 0, operations never time out.

Web Proxy Configuration

If you connect to Oracle Marketing Cloud through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to the Oracle Marketing Cloud account should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for .
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server.*Leave username & password blank for an unauthenticated proxy connection

Data Import Information

Via Browsing

See the [Oracle Marketing Objects table](#).

Via SQL Query

- Using a validSQL Select Query.
- Query Limitations: Oracle Marketing Cloud has limited support on multiple conditions for querying.
- SINGLE WHERE CLAUSE: In many tables the AND condition between filters is not supported. If AND is included the server might not return any results, nor will it throw an error.
- DATETIME NO EQUALS: The "=" operator cannot be used with datetime columns in the WHERE clause. However, the "<", ">", ">=", and "<=" operators can be used to filter on datetime columns.
- UPDATEDAT ONLY GREATER THAN: Additionally, the UpdatedAt column only supports filtering with the '>' operator.
- Objects that exhibit these limitations are identified in the [Oracle Marketing Objects table](#).

Oracle Marketing Objects

NAME	DESCRIPTION	LIMITATIONS
Account	Browse and query Accounts.	
AccountGroup	Browse and query Account group.	
Campaign	Browse and query Campaigns.	
Contact	Browse and query Contacts.	

NAME	DESCRIPTION	LIMITATIONS
ContactEmailSubscription	Browse and query all email group subscription status for the specific contact.	
ContactSegment	Browse and query Contact Segments.	
ContentSection	Browse and query Content Sections.	
Custom	Browse and query custom.	
Email	Browse and query Emails.	SINGLE WHERE CLAUSE DATETIME NO EQUALS UPDATEDAT ONLY GREATER THAN
EmailFooter	Browse and query Email Footers.	
EmailGroup	Browse and query Email Groups.	SINGLE WHERE CLAUSE DATETIME NO EQUALS UPDATEDAT ONLY GREATER THAN
EmailHeader	Browse and query Email Headers.	SINGLE WHERE CLAUSE DATETIME NO EQUALS UPDATEDAT ONLY GREATER THAN
Event	Browse and query Event.	
ExternalActivity	Browse and query External Activities.	
ExternalAsset	Browse and query External Assets.	SINGLE WHERE CLAUSE DATETIME NO EQUALS UPDATEDAT ONLY GREATER THAN
Folder	Browse and query Folders.	
Form	Browse and query Forms.	
Hyperlink	Browse and query Hyperlinks.	

NAME	DESCRIPTION	LIMITATIONS
LandingPage	Browse and query Landing Pages.	SINGLE WHERE CLAUSE DATETIME NO EQUALS UPDATEDAT ONLY GREATER THAN
Microsite	Browse and query Microsites.	SINGLE WHERE CLAUSE DATETIME NO EQUALS UPDATEDAT ONLY GREATER THAN
OptionList	Browse and query Option Lists.	SINGLE WHERE CLAUSE DATETIME NO EQUALS UPDATEDAT ONLY GREATER THAN
AccountField	Browse and query Account Fields.	
AccountView	Browse and query Account Views.	
Activity_Bounceback	Browse and query Bounceback Activities.	
Activity_CampaignMembership	Browse and query Campaign Membership Activities.	
Activity_EmailClickThrough	Browse and query Email Click Through Activities.	
Activity_EmailOpen	Browse and query Email Open Activities.	
Activity_EmailSend	Browse and query Email Send Activities.	
Activity_EmailSubscribe	Browse and query Email Subscribe Activities.	
Activity_EmailUnsubscribe	Browse and query Email Unsubscribe Activities.	
Activity_FormSubmit	Browse and query Form Submit Activities.	
Activity_PageView	Browse and query Page View Activities.	
Activity_WebVisit	Browse and query Web Visit Activities.	
CampaignElement	Browse and query Campaign Elements.	

NAME	DESCRIPTION	LIMITATIONS
CampaignField	Browse and query Campaign Fields.	
CampaignFolder	Browse and query Campaign Folders.	
ContactField	Browse and query Contact Fields.	
ContactFilter	Browse and query Contact Filters.	
ContactFilterFolder	Browse and query Contact Filter Folders.	
ContactList	Browse and query Contact Lists.	SINGLE WHERE CLAUSE DATETIME NO EQUALS UPDATEDAT ONLY GREATER THAN
ContactListFolder	Browse and query Contact List Folders.	
ContactScoringModelFolder	Browse and query Contact Scoring Model Folders.	
ContactSegmentData	Browse and query Eloqua Contact Segment Data.	
ContactSegmentFolder	Browse and query Contact Segment Folders.	
ContactView	Browse and query Contact Views.	
ContentSectionFolder	Browse and query Content Section Folders.	
Dependencies	Browse and query Dependencies.	
DynamicContent	Browse and query Dynamic Contents.	
DynamicContentFolder	Browse and query Dynamic Content Folders.	
EmailDeployment	Browse and query email deployments.	
EmailFolder	Browse and query Email Folders.	SINGLE WHERE CLAUSE DATETIME NO EQUALS UPDATEDAT ONLY GREATER THAN
EmailFooterFolder	Browse and query Email Footer Folders.	

NAME	DESCRIPTION	LIMITATIONS
EmailHeaderFolder	Browse and query Email Header Folders.	
ExternalType	Browse and query External Types.	
FieldMerge	Browse and query Field Merges.	
FieldMergeFolder	Browse and query Field Merge Folders.	
FormElement	Browse and query Form Elements.	
FormFolder	Browse and query Form Folders.	
FormProcessingStep	Browse and query Form Processing Steps.	
HyperlinkFolder	Browse and query Hyperlink Folders.	
Image	Browse and query Image.	
ImageFolder	Browse and query Image Folders.	
ImportedFile	Browse and query Imported Files.	
ImportedFileFolder	Browse and query Imported File Folders.	
LandingPageFolder	Browse and query Landing Page Folders.	
LeadScoringModel	Browse and query Query Lead Scoring Models.	
PageTag	Browse and query Page Tags.	
PageTagGroup	Browse and query Page Tag Groups.	
Program	Browse and query Programs.	
Style	Browse and query Styles.	
Template	Browse and query Templates.	
TemplateCategory	Browse and query Template Categories.	

NAME	DESCRIPTION	LIMITATIONS
TrackedUrl	Browse and query Tracked Urls.	
User	Browse and query Users.	SINGLE WHERE CLAUSE DATETIME NO EQUALS UPDATEDAT ONLY GREATER THAN
Visitor	Browse and query Visitors.	SINGLE WHERE CLAUSE DATETIME NO EQUALS UPDATEDAT ONLY GREATER THAN
VisitorProfileField	Browse and query all of the visitor profile fields.	

Technical Specs

Bulk API

This Connector will automatically negotiate usage of the Bulk API when possible.

PostgreSQL Connector for Data Prep

User Persona: Data Prep User or Data Prep Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

The ability to connect to PostgreSQL is part of the Data Prep JDBC Connector. This article provides details specific to setting up a connection to PostgreSQL. See also the [JDBC Connector Documentation](#) to configure the connection properly.

Example JDBC URI:

```
jdbc:postgresql://yourPostgresHost:yourPostgresPort/yourDatabaseName
```

Technical Specs

Driver Specs

- **PostgreSQL Database driver name & version:**
- Driver class name: org.postgresql.Driver
- Version: 42.2.8
- ****Supported PostgreSQL Database versions:****
 - 8.2 or newer

Driver Documentation

- **General Driver Documentation:** <https://jdbc.postgresql.org/documentation/head/index.html>

PowerBI Connector for Data Prep

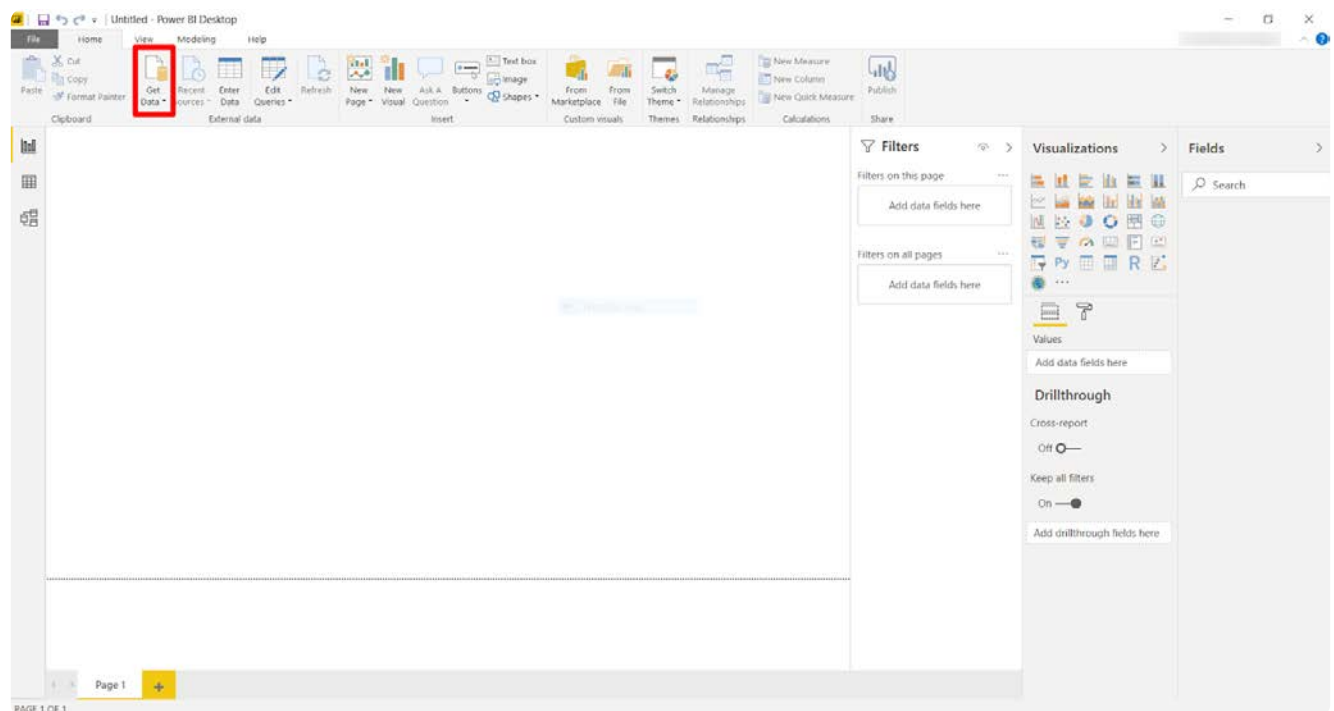
User Persona: Data Prep User

Note

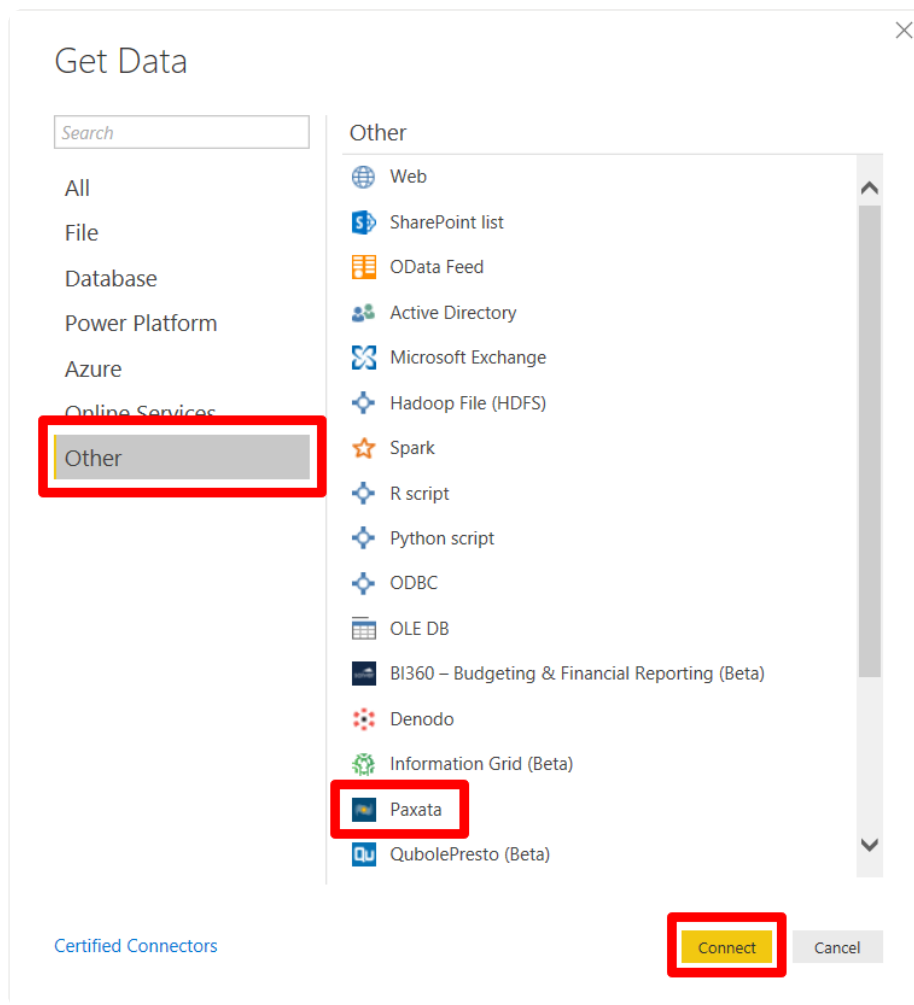
This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

To Configure

1. Open PowerBI Desktop
2. Click **Get Data**.



3. Click **Other > Data Prep > Connect**.



4. If you receive a warning about "Connecting to a third-party service", click **Continue**.

5. Enter the Data Prep URL you log in at and click **OK**.

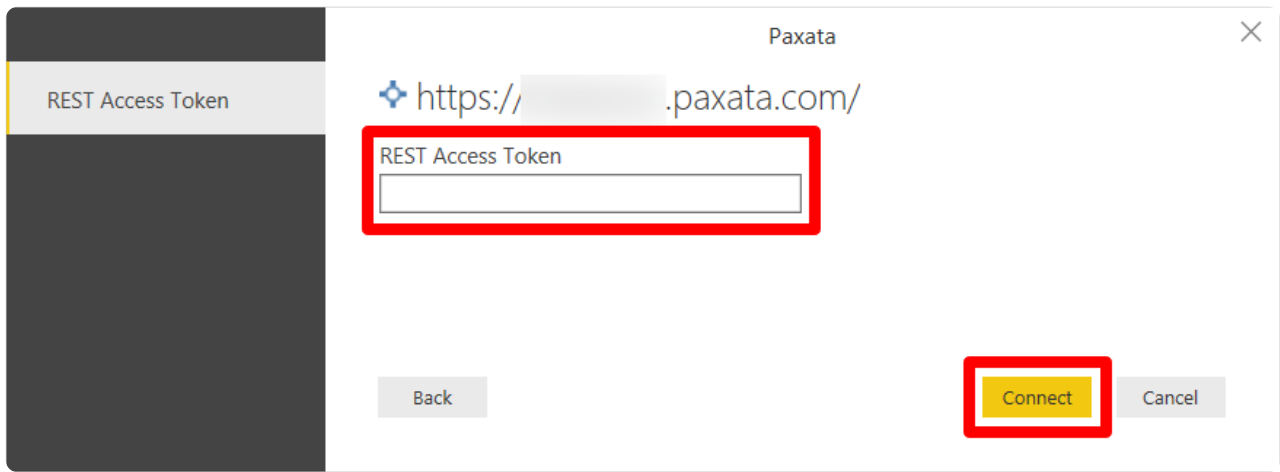
a. This may be "https://yourdomain.Data Prep.com/domain/paxauth" if you're on a Kubernetes-based deployment, or "https://yourdomain.Data Prep.com/" if you're on a non-Kubernetes based deployment.

b. Don't forget to add "https://" at the beginning

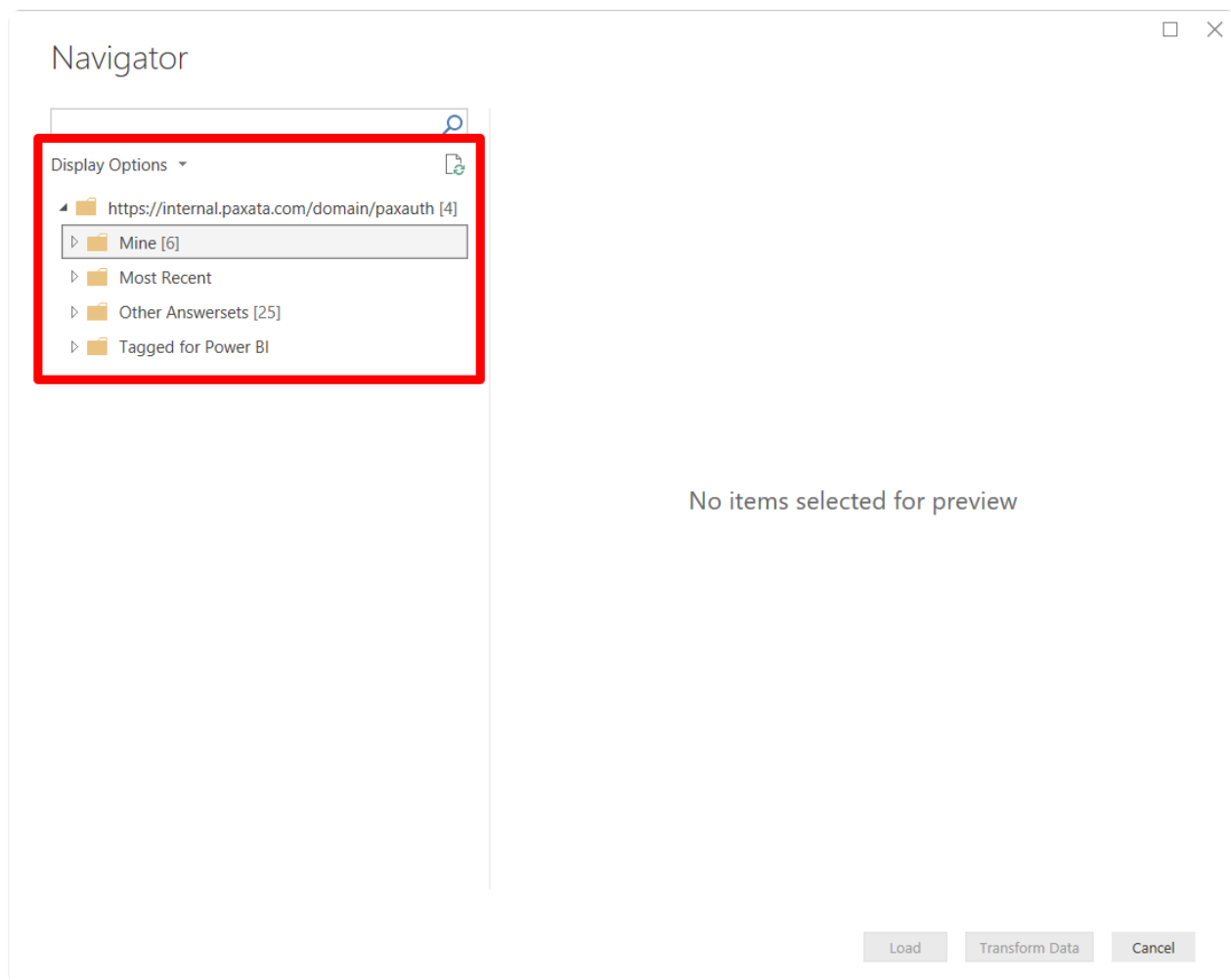
6. Enter your Data Prep REST API Token and click **Connect**.

a. If you do not have a REST API token, go to the **User** menu and click **Tokens**, where you can generate a token to manage access and authorization.

b. Save this token, as it cannot be retrieved, only deleted and regenerated, if lost.



7. It may take a minute or so for the authentication to process. Upon completion, you'll land on the screen where you can browse the Data Prep Data Library for Data Files and AnswerSets to import.



Best Practices

- For ad hoc reporting, using this Connector is the recommended path.
- For automated reporting/updating published dashboards, there are several approaches worth considering as automating a process on a Desktop is not possible. Here's one approach that we've seen be successful in several situations:

- Automate your data preparation in Data Prep and export to a storage location that can be directly queried by a published PowerBI Dashboard, such as ADLS Gen2.
- Using a separate automation on the PowerBI side, import the Data Prep AnswerSet into PowerBI from the shared location.

If you have further questions, contact your Data Prep Customer Success Representative.

REST API Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to a REST API to import a REST Resource. The following is information on the parameters used to create the connector.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can use the REST API Connector to connect Data Prep to multiple sources and potentially multiple instances of the same source. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Web Proxy

If you connect to your REST API source through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if connection to the REST Endpoint should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for Data Source.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server. *Leave username & password blank for an unauthenticated proxy connection.

REST API Configuration

In this section, provide the information used to locate the REST API resource.

For examples of how to set this up, see [REST API Authentication Configuration](#).

- **Base URL:** Base URL of the REST API. The base URL must include the protocol (http/https), hostname (port number is optional) and context path.
 - Example: `http(s)://api.domain.com(:port)/rest/v1`
- **Resources:** Multiple REST resources to be imported. Each line should contain a single REST resource configuration in `name:path?query` format.
 - The `name` is the user-visible name for the resource to be imported and is required for a REST resource configuration. This name will be presented in the Browse user-interface, for example: Account Details.
 - The `path` is the path to the resource and is required for a REST resource configuration. This path should start with a slash (/) and optionally has multiple segments separated by a slash (/), for example: `/resource/sub-category`
 - The `query` is an optional filtering criterion to use while retrieving the resource and is optional for a REST resource configuration. The query syntax must be key=value pairs delimited by '&', for example: `criteria=active&order=desc` or `jql=status=done`.

REST API Authentication Configuration

In this section, provide the information used to authenticate to the REST API service endpoint.

- **Authentication Type:** Select one of the options based on your requirement.
- **No Auth:** if the REST API doesn't require any authentication.
- **Basic Auth:** if the REST API allows authentication with Username and Password.
- **Bearer Token:** if the REST API allows authentication through Bearer Token. In the case of Bearer Token, each web service may provide access to or the generation of tokens differently and the web service's documentation should explain how to find it.
- **Username and Password:** If Authentication Type is selected as Basic Authentication, these fields are provided for authentication. Some web services only require one field or the other, so while most will require both fields, the configuration page allows them both to be blank. This may cause an error while authenticating to the data source, but will not cause form validation errors when saving the Data Source.
- **Bearer Token:** If Authentication Type is selected as Bearer Token, this needs to be provided for authentication. The user must know how to obtain this token as every system will handle this differently. Obtaining this token may also require Administrator help.

REST API Test Connection & Operation Configuration

- **Test Connection & Operation Method:** The HTTP method used in a request to determine if the Data Prep connector can connect to the REST API service and what method will be used when the Connector requests a resource. Selecting "Automatic" will try HEAD, GET and POST to test the connection and is the best option if you're unsure which method to select.

- The selected method is also used for actual import, if `GET` or `HEAD` succeeds in the test, then `GET` will be used for import, if `POST` is successful, then `POST` will be used for the actual import.
- **Connection Timeout:** Timeout (in milliseconds) for connecting to REST API.

Data Import Information

Via Browsing

Will present the Resources in the import workflow as the importable data set using the Resource Name as defined in the Resource list.

Via SQL Query

Not Supported.

Technical Specs

Pagination

- This Connector supports RFC 5988 pagination of REST datasets: <https://tools.ietf.org/html/rfc5988>
- For paginated REST responses, each paginated response contains HTTP Headers that identify the URL for the next page of results.
- When a paginated dataset is requested, the REST Connector will automatically identify that the dataset is paginated and follow data links.
 - Automatically extract the HTTP link for the next page of data.
 - Return the results from the current page of results.
 - Execute a call to obtain the next page of results.
- During Import using the Data Prep UI, we present only 1 page of data values in order to allow for rapid presentation of the Preview as well as to reduce hits against rate-limited APIs.
- During Import, the Connector will:
 - Automatically extract the HTTP link for the next page of data
 - Return the results from the current page of results
 - Execute a call to obtain the next page of results.

Performance

The performance of the REST Connector is very dependent upon the implementation of the REST API that it leverages.

- Best performance is found for REST APIs that support returning an entire dataset per REST API invocation. This is typical of APIs that leverage chunked transfer encoding. In this scenario, the REST Connector executes a single API call to obtain a full dataset.

- REST API's that leverage pagination reduce performance by requiring additional REST API calls.
 - Pagination style: RFC 5899
 - Each response contains N records and an HTTP Header containing a URL that points to the next batch.
- Review your REST API documentation to identify the maximum page size that you can configure in order to reduce the number of API calls.
- Example: GitHub REST API
- APIs can have rate limitations. When importing large datasets that are paginated, it is not uncommon to run into limitations on the number of REST calls made within a window of time. For example, GitHub allows for 5000 requests per hour and Google Drive allows for 1000 requests per 100 seconds.

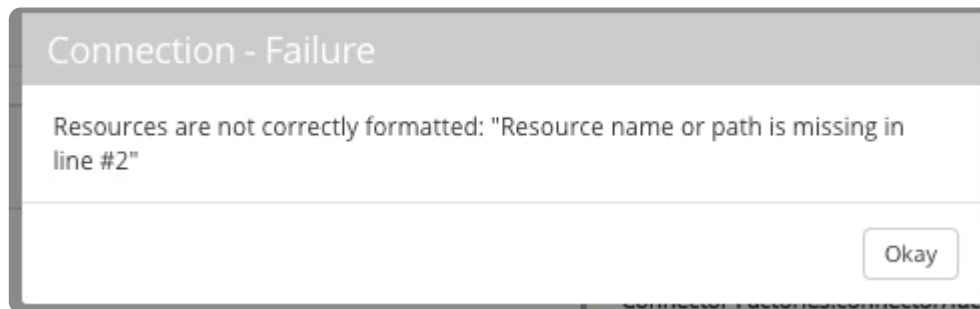
FAQ/Troubleshooting/Common Issues

Is OAuth authentication supported?

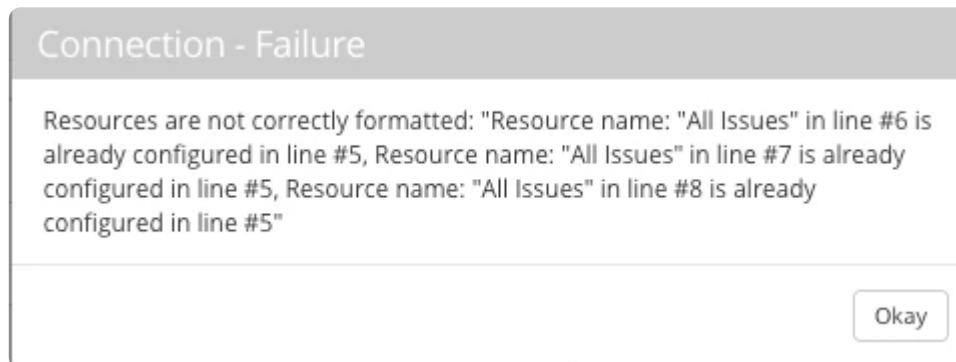
Not at this time. Currently, only username/password and token authentication methods are supported. Many data sources only allow for OAuth authentication and those sources would be unsupported at this time. Please contact your Data Prep Client Success if you find that you are unable to connect to a DataSource for this reason.

What do the "Test Connection" messages mean?

- Test Connection verifies that each entry in the Resource List matches the expected format.
- Failure of an entry to match the expected format results in an error indicating the identified format issue and the entry number.



- Failure to use unique Dataset Name for each Resource entry results in a format validation failure.



- After testing confirms proper formatting of resources, only the first entry on the Resource List is used to verify the connectivity is configured correctly.

Example Configurations

The following are real-world examples of how the REST API Connector has been used. Please feel free to use any of these in your account, but please be advised that companies may change their APIs at any time, sometimes without notice, and these are not fully supported data sources. This means that Data Prep will be unable to help you troubleshoot any issues you may have with these configurations and they may become out of date.

Simple Learning Example

There are many simple, unauthenticated REST API resources on the web that were created with the expressed purposes of learning, testing, and prototyping. One of those sources is JSONPlaceholder . This example may be over simplified, but is intended to demonstrate the building blocks of how to connect to a RESTful web service using the REST API Connector.

There are no rate limits posted and no pagination as the datasets are small.

CONFIGURATION

- Base URL: <https://jsonplaceholder.typicode.com>
- Resources:
- Posts:/posts
- Comments:/comments
- Albums:/albums
- Authentication: No Authentication
- REST API Test Configuration: Automatic

After clicking "Test Data Source" to confirm your setup is working and clicking "Save", you can now use this Data Source to import data into Data Prep.

GitHub Example

GitHub is a cloud-based software source code repository that provides a rich, but rate-limited REST API. [GitHub REST API Documentation](#). (Log in to GitHub before clicking this link.)

RATE LIMITS

- GitHub API rate limit reference: https://developer.github.com/v3/rate_limit/ (Log in to GitHub before clicking this link.)
- GitHub allows for 5000 requests per hour, specific limits vary per service.
- Rate limits vary for unauthenticated and authenticated users.

PAGINATION

- GitHub supports RFC 5899 pagination of REST datasets.
- For paginated REST responses, the user receives one page of data (30 entries for /search API) per API call.
- Users can override the number of results per page up to 100 using the "per_page" API parameter.

- Setting per-page result count to the maximum allowable setting will enable higher throughput of data import by reducing the number of REST API calls.

CONFIGURATION

- Base URL: <https://api.github.com> Log in to GitHub before clicking this link.
- Resources:
 - Mozilla Repos: retrieve a list of software repositories that match a search for "mozilla".
 - Note: This query will exhaust a user's quota for the /search API when using default 30 records per call.
 - Expected result count > 6600
 - Mozilla Repos: /search/repositories?q=mozilla
 - Mozilla Repos Page 33+: example of performing a search set starting at page 33 of results.
 - Mozilla Repos Page 33: search/repositories?q=mozilla&page=33
 - Square Repos: retrieve a list of the repositories that belong to the Square organization.
 - Square Repos: /orgs/square/repos
 - Organizations: retrieve the paginated list of all GitHub organizations using 100 records per request. *WARNING: This will run for a long time to pull 2+ million entries.*
 - Organizations: /organizations?per_page=100

Jira Example

Jira is a Project and Issue tracking software typically used by software development teams. [Jira REST API Documentation](#).

RATE LIMIT, PAGINATION, AND SETUP

- Rate limits will vary by subscription level.
- Jira REST API is limited to return at most 100 results per page of data.
- Jira does not support RFC 5899 pagination.
- Jira Cloud instances may require users to create a JIRA REST API token.
- Create a token: <https://confluence.atlassian.com/cloud/api-tokens-938839638.html>
- Authentication type: Basic Authentication.
- Username in username field.
- API Token in password field.

CONFIGURATION

- Base URL:
- On-Premise: [https://\(hostname\):\(port\)/rest/api/](https://(hostname):(port)/rest/api/)
- Jira Cloud: [https://\(your-domain\).atlassian.net/rest/api/](https://(your-domain).atlassian.net/rest/api/)
- Resources
- All Project List:
 - All Projects: /project

Example JQL resource: run a Jira JQL query to retrieve 200 To-Do task Items

- JQL query had to be URL encoded before pasting into configuration.
- Connector To Do Tasks:/search?
jql=Project%3D_yourProject_%20and%20statusCategory%3D%22To%20Do%22&maxResults=200
- Authentication: Basic (username/password)

PAGINATION FOR JIRA

Jira does not support RFC 5899 [pagination](#). In order to support pagination for JIRA:

- Define Datasource Resource entries that specify pages of data:
- Use "maxResults=100" to maximize the number of entries per REST call
- Use "startAt=N" to specify the starting point. N starts at 0
- Example: 4 pages of 100 search results
 - All Issues 0:/search?jql=&startAt=0&maxResults=100 All Issues 1:/search?jql=&startAt=100&maxResults=100 All Issues 2:/search?jql=&startAt=200&maxResults=100 All Issues 3:/search?jql=&startAt=300&maxResults=100
- Use Data Prep's Wildcard feature to select all pages of data
- Wildcard pattern = "All Issues*"

The screenshot shows the 'Select Datasets' interface. On the left, a search for '/All Issues*' is performed, resulting in four datasets: 'All Issues 1', 'All Issues 0', 'All Issues 3', and 'All Issues 2'. A 'COMBINE ALL RESULTS' button is visible. On the right, the 'You selected' panel shows the selected dataset: 'JSON /All Issues* (4 datasets)'. Below this, a preview of the JSON data is shown, with columns for 'id' and 'type'. The data is as follows:

	id	type
1	55613	htt
2	55607	httj
3	55606	httj
4	55604	httj
5	55603	httj

- After JSON parsing, which flattened the JSON and duplicated some rows to account for subtasks we obtained 557 rows of data X 298 columns

US Census Data Example

The Census data web site is not a REST API, but we can use our REST API Connector to retrieve data over HTTP

- Base URL: <https://www2.census.gov/>
- Resource examples for American Community Survey 2002

ACS_2002_Midwest:/acs2002/2007_prod_release1/BaseTablesSubjectTables/Region/
MidwestRegionBaseTables02000US2.csv

ACS_2002_US_OH_Franklin:/acs2002/2007_prod_release1/BaseTablesSubjectTables/States/Ohio/StateCounty/
FranklinCountyOhio/BaseTables05000US39049.csv

ACS_2002_Base_California:/acs2002/2007_prod_release1/BaseTablesSubjectTables/States/California/
CaliforniaBaseTables04000US06.csv

- Authentication: None
- Use Data Prep's Wildcard feature to select all pages of "ACS 2002 Base" data Wildcard Pattern:
- "ACS 2002 Base A*": Import Alabama, Alaska, Arkansas and Arizona files as 1 dataset.
- "ACS 2002 Base*": Import All matching 'ACS 2002 Base' files as 1 dataset.

Salesforce Lightning Connector for Data Prep

User Persona: Data Prep User or Salesforce Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to a Salesforce org as an import source. The following fields are used to define the connection parameters.

Note

Using Salesforce Lightning Edition is *not* a prerequisite for using this Connector. The Connector will work with any Salesforce org that supports the REST API v40.0 or later.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Salesforce orgs (Sandbox, Dev, UAT, etc). Using a descriptive name can be a big help to users in identifying the appropriate data source. If you are a Data Prep SaaS customer, inform Data Prep DevOps of how you would like this set.

Web Proxy

If you connect to Salesforce through a proxy server, these fields define the proxy details.

- **Web Proxy:** If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.

- **Proxy port:** The port on the proxy server for Data Source.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server. *Leave username & password blank for an unauthenticated proxy connection.

User Credentials

The user authentication can be done through a Shared Account or an Individual Account. If you choose to authenticate with an individual account, the user will be prompted to enter credentials to access this Data Source. If you choose to authenticate with a shared account, the following fields are required.

- **Salesforce URL:** The URL for Salesforce. By default, <https://login.salesforce.com>.
- If you're connecting to a Salesforce Sandbox, set the URL to: <https://test.salesforce.com>.
- **Session Security:** The session security options to use when connecting to Salesforce. Select 'API Security Token' if Salesforce is not configured to use [Trusted IP Ranges](#). If 'Trusted IP Ranges' is selected, an API Security Token value will not be required.
- **User:** The username for a shared account used to authenticate with Salesforce.
- **Password:** The password for a shared account used to authenticate with Salesforce.
- **API Security Token:** A user-specific Salesforce API security token that allows API-based access. This token is typically emailed to the Salesforce user each time the password is changed. It is not required if 'Trusted IP Ranges' is selected under Session Security.
- **OAuth App Fields** To set up Data Prep as an OAuth App in your Salesforce org, please follow the steps below in the Configuring Salesforce section below. Once you've completed these steps, you will be able to find the values required.
- **Consumer Key:** A value used by the Data Prep Salesforce Lightning connector to identify itself to Salesforce. In OAuth 2.0 the value is referred to as the client_id.
- **Consumer Secret:** A secret used by the Data Prep Salesforce Lightning connector to establish ownership of the Consumer Key. In OAuth 2.0 the value is referred to as client_secret.

Configuring Salesforce:

Salesforce Lightning Connector leverages the Salesforce REST API and OAuth. This step must be completed by a Salesforce Admin.

In order to establish connectivity:

- A Salesforce Administrator must create a "Connected App" in Salesforce.
- A Salesforce Administrator must obtain OAuth credentials for the "Connected App" to identify the client (Data Prep Connector) that is accessing Salesforce.
- Each individual user must authenticate.

Create a "Connected App" in Salesforce

Please use the Salesforce documentation linked immediately below for background on Connected Apps and navigation inside Salesforce. Use the guide below for how to configure the specified fields.

Salesforce Instructions: https://help.salesforce.com/articleView?id=connected_app_create.htm&type=5

- **Enable OAuth Settings:** Select this option.
- **Callback URL:** Salesforce requires a callback URL to be defined, but it is not required for Data Prep.
- For simplicity, please use your Data Prep URL.
- **Selected OAuth Scopes:**
 - Under "**Available OAuth Scopes**", select "Access and manage your data (api)" and click the "Add" button.
 - **Require Secret for Web Server Flow:** Select this option as Data Prep will securely store your OAuth Consumer Secret.
 - Be sure to hit "Save".

OAuth credentials for the "Connected App"

Upon Saving:

- You will be presented with a page that contains the OAuth credentials for the "Connected App".
- **Consumer Key:** Copy and save this.
- **Consumer Secret:** Click the "Click to reveal" link.
- ***Trusted IP Range for OAuth Web server flow***
 - If the customer wishes to further restrict access to a list of IP Addresses, Data Prep server IP addresses would be added here.

Data Import Information

Via Browsing

- All Objects are listed in alphabetical order.
- Custom Objects will be listed alphabetically as well and will include "__c" at the end.
- Objects with a Namespace will also be listed alphabetically under their Namespace.
- Folder of Reports at the top.

Via SQL Query

Using Queries with the Salesforce Connector relies on Salesforce's own query language, known as SOQL, for Salesforce Object Query Language. For information on this, please see: https://developer.salesforce.com/docs/atlas.en-us.soql_sosl.meta/soql_sosl/sforce_api_calls_soql.htm.

Here several example Queries:

1. Query which has escape characters in it

```
SELECT Id FROM Account WHERE Name LIKE 'Ter%'
```

2. Query with functions in it select count(Id) from Account

```
SELECT Name, MAX(Amount), MIN(Amount) FROM Opportunity GROUP BY Name
```

3. Query with Alias

```
SELECT a.Id, c.Id, c.name FROM Contact c, c.Account a WHERE a.name = 'MyriadPubs'
```

4. null in where clause

```
SELECT AccountId FROM Event WHERE ActivityDate != null
```

5. SubQuery in Where clause

```
SELECT Id, Name FROM Account WHERE Id IN ( SELECT AccountId FROM Opportunity WHERE StageName = 'Closed Lost' )
```

6. with GroupBy clause

```
SELECT LeadSource, COUNT(Name) FROM Lead GROUP BY LeadSource
```

7. Field as alias SELECT Name n FROM Opportunity

8. Query with TYPEOF SELECT TYPEOF What WHEN Account THEN Phone, NumberOfEmployees WHEN Opportunity THEN Amount, CloseDate ELSE Name, Email END FROM Event

9. Relationship query (to parent)

```
SELECT Contact.FirstName, Contact.Account.Name FROM Contact
```

10. Relationship query (to child)

```
SELECT Account.Name, (SELECT Contact.LastName FROM Account.Contacts) FROM Account
```

11. Relationship query with custom objects

```
SELECT Opportunity__c, Id, Opportunity__r.Name, Opportunity__r.Owner.Manager.Email, Opportunity__r.Owner.Email FROM Opportunity_Change__c
```

12. SubQuery in Select clause

```
SELECT Amount, Id, Name, ( SELECT Quantity, ListPrice, PricebookEntry.UnitPrice, PricebookEntry.Name FROM OpportunityLineItems ) FROM Opportunity
```

Best Practices

- Exporting to Salesforce:
- Most organizations have a defined process for bulk updating information within Salesforce, for this reason, the Data Prep Salesforce Connector only supports import.
- To export data to Salesforce, download a CSV file of your data locally and upload to Salesforce following your organization's guidelines.
- For more information on how to bulk load data into Salesforce, please see: https://help.salesforce.com/articleView?id=data_import_wizard.htm&type=5

Technical Specs

- Utilizes Salesforce REST API v40.0

FAQ/Troubleshooting/Common Issues

Salesforce Custom Reports Row Limit

Salesforce API limits Salesforce Reports to a result set of 2000 rows.

- https://help.salesforce.com/articleView?id=rd_reports_limits.htm&type=5T

Column Data Types Handling

When importing data from Salesforce, Data Prep will handle column data types differently when browsing vs. querying:

- Browse and Import:
 - Column data types are identified using Salesforce object metadata and mapped to Data Prep's internal types.
 - SOQL query import:
 - Salesforce SOQL results do not return the specific datatypes for columns and Data Prep does not parse queries or query results to determine column data types. As a result, all query results are interpreted as Text.

Example

When importing Salesforce Opportunity SObject row with "CloseDate" column containing "2013-11-13", CloseDate would be imported as:

- Browse and Import: DateTime
- "2013-11-13T00:00:00.000-08:00"
- SOQL Query: Text
- "2013-11-13"

Salesforce Marketing Cloud Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Salesforce Marketing Cloud for browsing and importing available data. The following parameters are used to configure the connection

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Salesforce Marketing Cloud accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source. If you are a Data Prep SaaS customer, inform Data Prep DevOps how you would like this set.

Salesforce Marketing Cloud Configuration

- **User:** The Salesforce Marketing Cloud user.
- **Password:** The password of the Salesforce Marketing Cloud user.
- **Timeout:** The number of seconds to wait for an operation. There is no default value.

Web Proxy Configuration

If you connect to Salesforce Marketing Cloud through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if connection to the Salesforce Marketing Cloud should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for Sale.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server.

Note

Leave username & password blank for an unauthenticated proxy connection.

Data Import & Export Information

Via Browsing

Browse to a table and "Select" the table for import.

Via SQL Query

Using a legal SQL Select Query

- For ts-date-time values, only > and < are supported in WHERE clauses.
- For non-ts-date-time values, =, !=, <>, >=, <=, IN are supported in WHERE clauses.

Salesforce Objects

NAME	TYPE	DESCRIPTION
Account	Table	A Marketing Cloud account.
AccountUser	Table	An individual user within an account. This table does not support deletes.
BusinessUnit	Table	A unit within a larger Enterprise or Enterprise 2.0 account. This table supports queries and updates only.
ContentArea	Table	A ContentArea represents a defined section of reusable content.

NAME	TYPE	DESCRIPTION
DataExtension	Table	Represents a data extension within an account.
Email	Table	Represents an email in a Marketing Cloud account.
EmailSendDefinition	Table	Record that contains the message information, sender profile, delivery profile, and audience information.
FileTrigger	Table	Reserved for future use. This table does not support deletes.
FilterDefinition	Table	Defines an audience based on specified rules in a filter. This table does not support inserts.
ImportDefinition	Table	Defines a reusable pattern of import options. This table does not support inserts.
List	Table	A marketing list of subscribers.
Portfolio	Table	Indicates a file within the Portfolio of a Marketing Cloud account.
ProgramManifestTemplate	Table	Reserved for future use. This table does not support deletes or inserts.
QueryDefinition	Table	Represents a SQL query activity accessed and performed by the SOAP API. This table does not support updates or inserts.
ReplyMailManagementConfiguration	Table	Details configuration settings for the reply mail management in an account. This table does not support deletes.
Send	Table	Used to send email and retrieve aggregate data. This table does not support deletes or updates.
SendClassification	Table	Represents a send classification in a Marketing Cloud account.
SenderProfile	Table	The send profile used in conjunction with an email send definition.

NAME	TYPE	DESCRIPTION
SMSTriggeredSend	Table	Indicates a single instance of an SMS triggered send. This table does not support deletes or updates.
Subscriber	Table	A person subscribed to receive email or SMS communication.
SuppressionListDefinition	Table	A suppression list that can be associated with different contexts.
TriggeredSendDefinition	Table	To create or update a TriggeredSendDefinition where the list ID is the All Subs List ID, you need the Email - Subscribers - All Subscribers - View and SendEmailToList permissions.
Automation	View	Defines an automation that exists within Automation Studio for an account.
BounceEvent	View	Contains SMTP and other information pertaining to the specific event of an email message bounce.
ClickEvent	View	Contains time and date information, as well as a URL ID and a URL, regarding a click on a link contained in a message.
DataExtensionField	View	Represents a field within a data extension.
DataExtensionTemplate	View	Represents a data extension template within an account.
DataFolder	View	Represents a folder in a Marketing Cloud account
DoubleOptInMOKeyword	View	The DoubleOptInMOKeyword object defines an MO keyword, allowing a mobile user to subscribe to SMS messages using a double opt-in workflow.
FileTriggerTypeLastPull	View	Reserved for future use.
ForwardedEmailEvent	View	Indicates a subscriber used the Forward To A Friend feature to send an email to another person.
ForwardedEmailOptInEvent	View	Specifies an opt-in event related to a Forward To A Friend event.

NAME	TYPE	DESCRIPTION
HelpMOKeyword	View	Defines actions associated with the HELP SMS keyword for an account.
ImportResultsSummary	View	A retrieve-only object that contains status and aggregate information on an individual import started from an ImportDefinition.
LinkSend	View	Provides information about a link in a send.
ListSend	View	Specifies retrieve-only properties associated with the list(s) for a completed send.
ListSubscriber	View	Retrieves subscribers for a list or lists for a subscriber.
MessagingVendorKind	View	Contains the vendor details for an SMS (short message service) or voice messaging vendor. Deprecated.
NotSentEvent	View	Contains information on when email message failed to be sent.
OpenEvent	View	Contains information about the opening of a message send by a subscriber.
PrivateIP	View	The PrivateIP object contains information on private IP address to be used as part of messages sends.
Publication	View	Reserved for future use.
PublicationSubscriber	View	Describes subscriber on a publication list.
PublicKeyManagement	View	Reserved for future use.
ResultItem	View	Contains results of asynchronous API call.
ResultMessage	View	Message containing results of async call.
Role	View	Defines roles and permissions assigned to a user in an account.
SendEmailMOKeyword	View	Defines the action that sends a triggered email message to the email addresses defined in an MO message.

NAME	TYPE	DESCRIPTION
SendSMSMOKeyword	View	Defines actions to take when the specified MO keyword is received.
SendSummary	View	A retrieve only object that contains summary information about a specific send event.
SentEvent	View	Contains tracking data related to a send, including information on individual subscribers.
SMSMTEvent	View	Contains information on a specific SMS message sent to a subscriber.
SMSSharedKeyword	View	Contains information used to request a keyword for use with SMS messages in a Marketing Cloud account.
SMSTriggeredSendDefinition	View	Defines the send definition for an SMS message.
SubscriberList	View	Use to retrieve lists for a specific subscriber.
SubscriberSendResult	View	Reserved for future use.
SuppressionListContext	View	Defines a context that a SuppressionListDefinition can be associated with.
SurveyEvent	View	Contains information on when a survey response took place.
Template	View	Represents an email template in a Marketing Cloud account.
TimeZone	View	Represents a specific time zone in the application.
TriggeredSendSummary	View	Summary of results for a specific triggered send.
UnsubEvent	View	Contains information regarding a specific unsubscription action taken by a subscriber.
UnsubscribeFromSMSPublicationMOKeyword	View	Defines keyword used by a subscriber to unsubscribe from an SMS publication list.

SAP HANA Connector for Data Prep

User Persona: Data Prep Admin or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

The ability to connect to SAP HANA is part of the Data Prep JDBC Connector. This topic provides details specific to setting up a connection to SAP HANA. See also the [JDBC Connector Documentation](#) to configure the connection.

Example JDBC URI:

```
jdbc:sap://_yourSAPHanaHostOrIP_:_yourSAPHannaPort_/?currentschema=_yourRootSchema_
```

Technical Specs

Driver Specs

- **SAP HANA Database driver version:**
 - Version: 2.4.63
- **Supported SAP HANA Database versions:**
 - SAP HANA 1.0 and SAP HANA 2.0 databases

Driver Documentation

- **General Driver Documentation:** <https://help.sap.com/viewer/0eec0d68141541d1b07893a39944924e/2.0.04/en-US>

Spark SQL Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to Spark SQL for browsing, importing, and exporting available data. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Spark SQL instances. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Spark SQL Server Configuration

- **Spark SQL Server:** The hostname or IP address of the server hosting the Spark SQL database.
- **Spark SQL Port:** The port for the Spark SQL database.
- **Use SSL:** Set this property to the value specified in the 'hive.server2.use.SSL' property of your Hive configuration file (hive-site.xml).
- **Transport Mode:** Set this property to the value specified in the 'hive.server2.transport.mode' property of your Hive configuration file (hive-site.xml).
- **HTTP Path:** This property is used to specify the path component of the URL endpoint when using HTTP Transport Mode. This property should be set to the value specified in the 'hive.server2.thrift.http.path' property of the Hive configuration file (hive-site.xml).

Spark SQL Server Authentication Configuration

- **User:** The username used to authenticate with Spark SQL. For Databricks, set to 'token'.
- **Password:** The password used to authenticate with Spark SQL. For Databricks, set to your personal access token (value can be obtained by navigating to the User Settings page of your Databricks instance and selecting the Access Tokens tab).

Data Import Information

Via Browsing

- View to a table and "Select" the table for import.

Via SQL Query

- Supports importing using a legal SQL Select Query.

SFTP Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This Connector allows you to connect to an SSH File Transfer Protocol (SFTP) Server for Library imports and exports. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple SFTP servers. Using a descriptive name can be a big help to users in identifying the appropriate data source. If you are a Data Prep SaaS customer, inform Data Prep DevOps how you would like this set.

SFTP Host

If the SFTP Host section appears in the Add Source/Edit Source form, provide the information used to locate and connect to the SFTP host.

- **SFTP Hostname:** You can use either the fully qualified hostname, including the domain name, or the IP address of the SFTP server.
- **SFTP Port:** The socket port for the SFTP server. The protocol specifies port 22 as default.
- **Automatic Host Key Verification:** Automatically accept the host key from the SFTP server.
- *Selected:* This option enables SFTP Connector to automatically trust connections to your SFTP server. This is equivalent to setting `StrictHostKeyChecking=no` in SSH.

- *Deselected (default setting)*: This option disables automatic trust configuration to SFTP HOSTNAME. This is selected as the default option as it represents the higher security configuration.
- **Keep Alive**: Enable/Disable session activity to prevent a timeout.
 - *Selected (default setting)*: Enables periodic background communication between SFTP Connector and SFTP server to keep the connection from being closed by the server during browse, import, and export.
 - *Deselected*: The duration of the connection is managed by the SFTP server configuration. Idle connections may be terminated by the server. In this configuration, it is best to avoid lapsed inactivity when browsing to import/export data.
- **Data Compression**: Enable data compression during transfer.
 - *Selected (default setting)*: Enables ZLIB compression of data during transfer between the SFTP server and Data Prep, resulting in an increase in transfer speed for most datasets. In the event that ZLIB compression cannot be negotiated between Data Prep and server, the connection will fall back to uncompressed transfer automatically.
 - *Deselected*: Disables ZLIB compression of data during transfer between SFTP server and Data Prep
- **Socket Timeout Seconds**: The number of seconds to wait for SFTP command execution (list directory, create directory, logout...). The default value is 30 seconds. To allow for longer wait, increase this value.
- This option will most likely be used when the SFTP server directories contain very large lists of data files.

Configuration

- **Root Directory**: Defines the top-level directory to be presented in Data Prep's browse interfaces for import and export. Users can see files and directories within this directory in the browsing interface.

Authentication

The SFTP connector can authenticate using password authentication or SSH keys (with or without a Passphrase). Here are the options:

- **User Credentials**: This is a username and password combination.
 - *USERNAME*: The username for authenticating with the SFTP server
 - *PASSWORD*: The password associated with the provided username
- **SSH Key Without Passphrase**: This option only requires that you paste in the SSH Key
 - *USERNAME*: The username for authenticating with the SFTP server
 - *SSH PRIVATE KEY*: The contents of the SSH private key associated with the username
- **SSH Key With Passphrase**: Paste in the SSH Key and enter the Passphrase
 - *USERNAME*: The username for authenticating with the SFTP server
 - *SSH PRIVATE KEY*: The contents of the SSH private key associated with the username
 - *PASSPHRASE*: The encryption passphrase for your Private Key

Data Import & Export Information

Via Browsing

- The Connector will present a browsable directory hierarchy starting at the location defined in the ROOT DIRECTORY field.
- The Connector also supports Wildcard & Glob importing, this enables users to import multiple SFTP data files into Data Prep as a single Dataset.

Via SQL Query

- As SFTP is a file store, SQL Queries are not supported for this data source.

Technical Specs

- We test this Connector against a standard, non-configured Linux implementation of OpenSSH

FAQ/Troubleshooting/Common Issues

Note that SFTP is as much a protocol as it is a type of storage. If you have an “SFTP Server”, what you really have is a storage location that interfaces with the web using the SSH File Transfer Protocol. This is an important distinction as anything (web services, SFTP service providers, etc) can expose data to the web using this protocol. These services might be using different implementations of SFTP or they may do things behind the scenes that a traditional SFTP Server would not. All this is to say that SFTP servers may have custom behavior that presents challenges either in connecting or importing data.

Here’s one example of where this type of variance from standard SFTP caused some challenges: A customer was using the SFTP Connector to pull data from one of their vendors. The vendor was using a service that exposed data via SFTP, but would then delete each datafile after being read. When Data Prep provides a preview of data upon import, this is done by querying the data source for a small chunk of the data present. This caused the system to delete the file before it could be fully imported.

Snowflake Data Warehouse Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

The Snowflake connector allows you to use a JDBC-based connection for Library imports and exports. The following fields are used to create a connection.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Snowflake accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source. If you are a Data Prep SaaS customer, inform Data Prep DevOps how you would like this set.

DataBase URI

- **JDBC URI:** The JDBC connection string. The schema name can be included in the URI. For more information, see: <https://docs.snowflake.net/manuals/user-guide/jdbc-configure.html>

Database, Schema, and Table Visibility

You can control the database, schemas, and tables that are shown to users when they browse a data source during import. For each, you can choose to:

- "Show only" which returns only the databases, schemas or tables that you specify here.
- "Hide" which hides the databases, schemas, and tables that you specify here.
- "Show all" which is the default setting to display everything in the data source.

When you select the "Show only" or "Hide" options, a field is provided for specifying a comma-separated list of the databases, schemas or tables on which you want the option enforced.

Note

These settings are not enforced when users query against the data source using SQL; query results still return a complete list of matches. For example, if you choose to "hide" a specific database, users can still execute queries that pull data from tables within that database. However, that database will not be displayed to users when they browse the data source.

Import Configuration

- **Query Prefetch Size for queries:** The batch size used (in rows) when prefetching data on import.
- **Max Column Size:** The maximum size in Unicode characters allowed in a cell when importing or exporting data.
- **Pre-import SQL:** A SQL statement to execute before beginning import, after the table's schema is determined. This SQL will also be executed before data is previewed in the Data Prep UI.
- **Post-import SQL:** A SQL statement to execute after import completes. This SQL will also be executed before data is previewed in the Data Prep UI.

Export Configuration

- **Export method:** Select the method for exporting data from Snowflake. Both of these methods are specific to how Snowflake imports the data. For more details on these options, please refer to the Snowflake documentation linked behind the listed option. There are two options:
- **Internal stage:** Write data to a file on a Snowflake internal stage before loading the data into a table. This method is recommended for larger datasets since it will be faster than Direct SQL. **Stage Types:**
 - **Temporary:** The stage created will be dropped at the end of the session in which it was created. The stage is managed by Snowflake, so no further configuration is required.
 - **Permanent:** Specifies the name of a Stage that has already been created in Snowflake.
 - **Stage Name:** specifies the name for an existing named Internal Stage in Snowflake. See the Snowflake [Identifier Syntax](#).
- **Direct SQL:** Use SQL insert statements to export data. For larger datasets, this approach is slower than using an internal stage.
 - **Export Batch Size:** The batch size used when exporting data if the Direct SQL export method has been selected.

- **Max VARCHAR Size:** The maximum size allowed for a VARCHAR column allowed on export. Values greater than this size will be replaced with null when data is loaded into the Snowflake table.
- **Automatically create table:** Create a new table automatically on export. If enabled, Data Prep will drop the table whose name matches the name of the exported dataset, if one already exists, and recreate the table using the exported dataset. If not enabled, Data Prep will not create a new table, but will instead load the exported data into the table whose name matches the name of the exported dataset.
- **Pre-export SQL:** A SQL statement to execute before beginning export, after the table is created if auto-create is enabled.
- **Post-export SQL:** A SQL statement to execute after export completes.

Note

Data Prep exports using the "TIMESTAMP_LTZ(9)" type. If a table was created using a different timestamp, exporting Data Prep data to a column with a mismatching timestamp type will result in an error. The error reads: "There was an error while performing the export. Reason: SQL compilation error: Expression type does not match column data type, expecting TIMESTAMP_#### but got TIMESTAMP_LTZ(9) for column _Column_Name."

To correct this, do one of the following:

- Allow Data Prep to create the table and then perform the export, or
- Create your table with TIMESTAMP_LTZ(9) and then perform the export.

Credentials

The user authentication can be done through a Shared Account or an Individual Account. If you choose to authenticate with an Individual Account, the user will be prompted to enter a username and password to access this Data Source. If you choose to authenticate with a Shared Account, the following fields are required.

User: The username for a shared account used to authenticate with the database.

Password: The password for a shared account used to authenticate with the database.

Role: The role to set for the session using 'USE ROLE'. The specified role should be an existing role that has already been assigned to the user. Roles can also be specified using the JDBC URI, but the value specified in the Role field will supersede a role specified in the URI. Information on Snowflake Roles can be found here: <https://docs.snowflake.net/manuals/sql-reference/sql/use-role.html>

Data Import Information

Via Browsing

The browsing experience will vary and is based on the database, schemas, and tables visibility settings selected above along with the user credentials supplied.

Via SQL Query

As noted in the database, schemas, and tables visibility section, the limits to what a user can import via query are only limited to their authorization as determined by the credentials supplied for the connection.

Queries can be performed using legal SQL Select statement as defined here: <https://docs.snowflake.net/manuals/sql-reference/sql/select.html>

Example:

```
SELECT * FROM "SNOWFLAKE_SAMPLE_DATA"."TPCH_SF1"."CUSTOMER"
```

Tableau Hyper Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configuring Data Prep

This connector allows you to connect to Tableau as an export destination. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Tableau accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source. If you are a Data Prep SaaS customer, inform Data Prep DevOps how you would like this set.

Tableau Configuration

- **Tableau Server URL:** An http or https URL where the Tableau Server is running. Include port if necessary. For Tableau Online, the URL must contain the pod name, such as 10az, 10ay, or us-east-1.
- **Site:** Tableau site in which Data Prep will publish a data source.
- **Project:** Tableau project in which Data Prep will publish a data source.
- **Overwrite Tableau Data Source:** Overwrite a data source that has the same name, if one already exists.

Tableau Credentials

- **Username:** Username or email address for authenticating with Tableau.
- **Password:** Password for authenticating with Tableau.

Web Proxy

If you connect to Tableau through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to the Tableau REST Endpoint should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for Data Source.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server. Leave username & password blank for an unauthenticated proxy connection.

Data Export Information

Via Browsing

Files are published into the specified Site and Project defined in the Tableau Connector/Datasource configuration.

No directories will be visible in the browsing UI. Click the 'Select' button in the export view to export into the specified Site and Project.

Note that AnswerSets are converted to Tableau Hyper Extract (.hyper) files and then published into the specified Site and Project you define here.

Via SQL Query

Not supported.

Technical Specs

Deployment Scenarios

The Tableau Hyper connector writes .hyper files using native libraries from the Tableau SDK. There are two supported deployment scenarios that allow you to use these native libraries:

1. Use the native libraries that are bundled in the connector. This is the recommended approach.
2. Install the Tableau SDK on the host system and remove the native libraries that are extracted from the connector zip.

If you need additional information or help installing the SDK for the Tableau Hyper Connector, please reach out to your Customer Success representative.

FAQ/Troubleshooting/Common Issues

Why are there two Tableau Connectors and what's the difference?

Tableau ".tde" is an older extract type and was implemented with version 1.0 of Tableau's extract API. The Tableau ".hyper" Connector is newer and is based on version 2.0 of Tableau's extract API. We built the ".hyper" Connector as a net-new Connector so both Connectors could coexist in an account enabling customers who have many versions of Tableau running to choose the best Connector for each instance of Tableau.

Why can't my peer certificate be authenticated?

SSL communication between servers requires an identification method and something that verifies the identity assertion. *Keystores* (for example, jks) contain certificates and/or private keys that hold the identity assertion used when communicating with other servers. *Truststores* (for example, cacerts) are used when communicating with other servers to confirm it's safe to establish a connection.

If your Data Prep truststores do not have the necessary CA certificates to validate those that DataRobot receives from server connections, you may receive the following error:

`PEER CERTIFICATE CANNOT BE AUTHENTICATED WITH GIVEN CA CERTIFICATES` To resolve this error, update the Data Prep truststore for HTTPS connections over SSL. Use the tabs below to update the Data Prep truststore for your connection type:

Connections such as Tableau will use web-based rather than Java-based connections. Therefore, use the Linux-based CA-certificates truststore.

To update this truststore, do the following:

1. If dynamic CA isn't enabled yet, enable it:

```
sudo update-ca-trust force-enable .
```

2. Add your root CA and any intermediate certs to `/etc/pki/ca-trust/source/anchors/`.

3. Add the certs to the truststore: `sudo update-ca-trust extract`.

Connections such as JDBC will be Java-based and use the JKS truststore. By default, Data Prep points to its own truststore in `/usr/local/paxata/server/config/px.properties`:

```
px.ssl.truststore=/usr/local/paxata/server/paxata.jks  
px.ssl.truststore.password=OBF:ioevnaaksreu923rh
```

Note

Verify this truststore before proceeding.

To use the cacerts truststore:

1. Add a root CA and intermediate certificates to the cacerts truststore.

```
sudo keytool -importcert -alias <app>_rootca -file  
<app>_rootca.crt -keystore /usr/lib/jvm/java-X.X.X-  
openjdk/jre/lib/security/cacerts  
Password: changeit
```

2. Create the OBF obfuscated password for the cacerts truststore. This assumes the password for cacerts is the default `changeit`.

```
java -cp <path to jetty-util>/jetty-util-X.X.X.vXXXXX.jar  
org.eclipse.jetty.util.security.Password changeit  
OBF:f0wrhz9rshc8srh32  
MDF:irewqnsd9re320dfk
```

3. Edit `/usr/local/paxata/server/config/px.properties` to use the cacerts truststore.

```
px.ssl.truststore=/usr/lib/jvm/java-X.X.X-openjdk/jre/  
lib/security/cacerts
```

```
px.ssl.truststore.password=OBF:f0wrhz9rshc8srh32
```

4. Restart the Data Prep (Paxata) core server:

```
sudo systemctl restart paxata-server .
```

Tableau TDE Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to connect to a Tableau Server & Tableau Online as an export destination. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Snowflake accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source.

Tableau Server

- **Tableau Server URL:** An http/https URL where the Tableau Server is running. Include port if necessary.
- **Site:** Tableau site in which Data Prep will publish a data source.
- **Project:** Tableau project in which Data Prep will publish a data source.

Tableau Credentials

- **Username:** This is your Tableau username.
- **Password:** This is your Tableau password.

Data Export Information

Via Browsing

- Note that the browsing interface will not show any directory items from Tableau. Files are published into the specified Site and Project defined in the Tableau Connector/Datasource configuration, not to a destination selected at export.
- Note that AnswerSets are converted to Tableau Data Extract (.tde) files and then published into the specified Site and Project you defined in the Connector/Datasource configuration.

Via SQL Query

- Not supported.

Technical Specs

- The Tableau SDK (native libraries) must be installed on the Data Prep server host in order to enable communication to the Tableau Server and to export TDE files.
- [Tableau Connector SDK](#)

FAQ/Troubleshooting/Common Issues

- It's not uncommon for there to be Tableau permissions issues when exporting from Data Prep to Tableau.

Teradata Connector for Data Prep

User Persona: Data Prep Admin or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

The ability to connect to Teradata is part of the Data Prep JDBC Connector. This topic provides details specific to setting up a connection to Teradata. See also the [JDBC Connector Documentation](#) to configure the connection.

Example JDBC URI:

```
jdbc:teradata://yourTeradataHost/TMODE=ANSI,COP=OFF
```

Technical Specs

Driver Specs

- **Teradata Database driver name & version:**
 - Driver classname: com.teradata.jdbc.TeraDriver
 - Version: 16.20.00.12
- **Supported Teradata Database versions:**
 - 16.20, 16.10, 15.10, 15.0, 14.10

Driver Documentation

- **General Driver Documentation:** <https://teradata-docs.s3.amazonaws.com/doc/connectivity/jdbc/reference/current/frameset.html>
- **Driver FAQ:** <https://teradata-docs.s3.amazonaws.com/doc/connectivity/jdbc/reference/current/faq.html>

ThoughtSpot Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

This connector allows you to export your Data Prep AnswerSets to ThoughtSpot. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple ThoughtSpot instances. Using a descriptive name can be a big help to users in identifying the appropriate data source.

ThoughtSpot Server Configuration

- **Hostname:** ThoughtSpot server hostname or IP address.
- **Port:** ThoughtSpot server port. The standard port is 22.

ThoughtSpot User Configuration

- **User:** The user on the ThoughtSpot server. (This is the SSH username, not the username for logging in on ThoughtSpot web app. See the [Thoughtspot documentation](#) for more information on login credentials.)
- **Password:** The password for the user on the ThoughtSpot server.

Export Configuration

- **Behavior for Exporting into Existing Table:** Select the action that Data Prep should take upon export if a table with the same name already exists in the selected database and schema.
 - **Append New Data to Existing Table:** Data Prep will add the new rows of data to the existing table.
 - **Drop Existing Table and Create New Table:** Data Prep will drop all existing rows and columns, and create a new table with the new data.
 - **Empty Target Table and Add New Data:** Data Prep will keep the existing table but remove the row values and fill in the new data.

Note

If there is no existing table that matches the name specified in export, then all of the options are equivalent: all of them will create a new table and load all of the rows from the dataset.

Data Export Information

Via Browsing

View list of databases and schemas within ThoughtSpot.

Via SQL Query

Not Supported.

Best Practices

During the set up of the Connector Config and Data Source Config, you have the option to select the behavior for the export. These settings are configured once and used over and over again. With this, if you want one export to truncate the table and one to append, you should create one Connector Config where you leave this setting blank, then create two Data Source Configs, one where you specify the append behavior, the other where you specify truncate. Be sure to give the Data Sources clear names.

Vertica Connector for Data Prep

User Persona: Data Prep Admin or Data Source Admin

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configure Data Prep

The ability to connect to HP Vertica is part of the Data Prep JDBC Connector. This topic provides details specific to setting up a connection to Vertica. See also the [JDBC Connector Documentation](#) to configure the connection.

Example JDBC URI:

```
jdbc:vertica://yourVerticaHost:_yourVerticaPort/yourDatabaseName_
```

Technical Specs

Driver Specs

- **HP Vertica driver version:** 9.2.1
- **Supported HP Vertica versions:** 9.2.x

Driver Documentation

- **General Driver Documentation:** <https://www.vertica.com/docs/9.2.x/HTML/Content/Authoring/ConnectingToVertica/ClientJDBC/ProgrammingJDBCClientApplications.htm>
- **SQL Reference:** <https://www.vertica.com/docs/9.2.x/HTML/Content/Authoring/SQLReferenceManual/SQLReferenceManual.htm>

Zendesk Connector for Data Prep

User Persona: Data Prep User, Data Prep Admin, Data Source Admin, or IT/DevOps

Note

This document covers all configuration fields available during connector setup. Some fields may have already been filled out by your Administrator at an earlier step of configuration and may not be visible to you. For more information on Data Prep's connector framework, see [Data Prep Connector setup](#). Also, your Admin may have named this connector something else in the list of Data Sources.

Configuring Data Prep

This connector allows you to connect to Zendesk for browsing and importing available data. The following fields are used to define the connection parameters.

General

- **Name:** Name of the data source as it will appear to users in the UI.
- **Description:** Description of the data source as it will appear to users in the UI.

Tip

You can connect Data Prep to multiple Zendesk accounts. Using a descriptive name can be a big help to users in identifying the appropriate data source. If you are a Data Prep SaaS customer, inform Data Prep DevOps how you would like this set.

Web Proxy

If you connect to Zendesk through a proxy server, these fields define the proxy details.

- **Web Proxy:** 'None' if no proxy is required or 'Proxied' if the connection to the Zendesk REST Endpoint should be made via a proxy server. If a web proxy server is required, the following fields are required to enable a proxied connection.
- **Proxy host:** The host name or IP address of the web proxy server.
- **Proxy port:** The port on the proxy server for Data Source.
- **Proxy username:** The username for the proxy server.
- **Proxy password:** The password for the proxy server. *Leave username & password blank for an unauthenticated proxy connection.

Zendesk Configuration

- **Zendesk URL:** Your Zendesk URL in the format `https://your-domain.zendesk.com`.
- **Username:** The email address of the user connecting to Zendesk.
- **Authentication Type:** Type of authentication you want to use, either Password or API token.
 - **Password:** Password for connecting to Zendesk. To be provided when Password is selected as the authentication type.
 - **API Token:** API token for connecting to Zendesk. To be provided when API Token is selected as the authentication type.
- **Timeout:** The number of seconds to wait until a timeout error cancels a running operation.

Data Import Information

Via Browsing

Browse to a table and "Select" the table for import.

Via SQL Query

Using a legal SQL Select Query

Zendesk Objects

NAME	DESCRIPTION
AccountSettings	Account Settings in Zendesk.
ActivityStream	ActivityStreams in Zendesk.
AppLocations	AppLocations in Zendesk.
Attachments	Show ticket attachments in Zendesk.
Automations	Automations in Zendesk.
Brands	Brands in Zendesk.
Collaborators	Collaborators in Zendesk.
CustomAgentRoles	CustomAgentRoles in Zendesk.
GroupMemberships	Group Memberships in Zendesk.

NAME	DESCRIPTION
Groups	Groups in Zendesk.
Holidays	Schedules in Zendesk.
JobStatus	A status record is created when somebody kicks off a job such as updating multiple tickets. * You can access the job status data for an hour after a particular job is created, after which the data is no longer available.
Locales	Locales in Zendesk.
Macros	Macros in Zendesk.
MonitoredTwitterHandles	Monitored Twitter handles in Zendesk.
OrganizationFields	Organization Fields in Zendesk.
OrganizationMemberships	Organization Memberships in Zendesk.
OrganizationSubscriptions	Organization Subscriptions in Zendesk.
Organizations	Organizations in Zendesk.
Requests	Requests in Zendesk.
SatisfactionRatings	Requests in Zendesk.
Schedules	Schedules in Zendesk.
Sessions	Sessions in Zendesk.
SharingAgreements	SharingAgreements in Zendesk.
SlaPolicies	SlaPolicies in Zendesk.
SupportAddresses	Support Addresses in Zendesk.
SuspendedTickets	Suspended Tickets in Zendesk.
TicketAudits	TicketAudits in Zendesk.

NAME	DESCRIPTION
TicketComments	TicketComments belonging to a specified ticket in Zendesk.
TicketFields	TicketFields in Zendesk.
TicketForms	Ticket Forms in Zendesk.
TicketMetricEvents	TicketMetricEvents in Zendesk.
TicketMetrics	TicketMetrics in Zendesk.
Tickets	Tickets.
Triggers	Triggers in Zendesk.
UserFields	User Fields.
UserIdentities	User Identities.
UserRelatedInformation	UserRelatedInformation in Zendesk.
Users	Users in Zendesk.
Views	Views in Zendesk.

Work with datasets

These sections describe how to manage your datasets in Data Prep:

TOPIC	DESCRIBES...
Import datasets	How to import data from a data source or from a local file. Learn about what settings you can adjust before importing.
Export datasets	How to export datasets and AnswerSets . Learn about what settings you can adjust before exporting.
Profile datasets	How to generate a profile of a dataset, including information about the quality of the data.
Update datasets with new data	How to update the data in an existing dataset .
Update project datasets	Learn how to refresh and replace datasets using the Steps tool.

Import datasets

Importing data into Data Prep is the first step to getting your data ready for machine learning. During the import process you can:

- Select multiple datasets from a variety of data sources.
- Combine datasets together into one dataset.
- Choose which columns in a dataset to import.
- Select extensionless files.
- Import datasets from zipped or compressed folder.
- Change the format used to analyze and structure your data.

Using the Import page

After you select a dataset for import, the page divides into four quadrants called panes.

The screenshot shows the DataRobot interface for importing a dataset. The page is divided into four panes:

- Pane 1: Select Datasets** - Shows options to select a data source or upload a local file. A large cloud icon with an upward arrow and the text "Click here or drag-and-drop files to upload" is visible.
- Pane 2: You selected** - Shows the selected dataset "bank.csv" with a total size of 366.7 KB. Buttons for "Finish", "Finish Later", and "Cancel" are present.
- Pane 3: Your options for bank.csv** - Shows configuration options for the dataset, including Name (bank.csv), Description, Character encoding (UTF-8), and Rows to process for schema (1000).
- Pane 4: Data preview** - Shows a table with 10 rows of data. The columns are: age, job, marital, education, default, and balance. The data is as follows:

	age	job	marital	education	default	balance
1	30	unemployed	married	primary	no	178
2	33	services	married	secondary	no	478
3	35	managem...	single	tertiary	no	135
4	30	managem...	married	tertiary	no	147
5	59	blue-collar	married	secondary	no	
6	35	managem...	single	tertiary	no	74
7	36	self-emplo...	married	tertiary	no	30
8	39	technician	married	secondary	no	14
9	41	entrepren...	married	tertiary	no	22
10	43	services	married	primary	no	-8

Following is an overview of each pane of the **Import** page.

#	ELEMENT	DESCRIPTION
1	Select Datasets pane	Select the datasets you want to import from this pane. You can: <ul style="list-style-type: none">• Select multiple datasets from local files and connected data sources.• Search and query connected data sources for datasets.• Combine multiple datasets into one glob for import.
2	You selected pane	After you select a dataset, your dataset is listed in this pane. You can: <ul style="list-style-type: none">• See a list of datasets that have been selected for import.• Select a dataset to preview or update the import options.• Quickly identify datasets with potential import errors.• Change the format used to analyze and structure your data for import.• Import the same dataset multiple times with different import options.
3	Your options pane	In many cases, your data will import easily into Data Prep. In some cases, you may need or want to adjust the import options. This is where you make those adjustments.
4	Preview pane	Here is where you preview of your data. As you select datasets from the You selected pane, update the format, or change import options, the Preview pane displays how the selected dataset will look once imported. From here, you can also choose which columns to import.

Snapshot of the import process

Following is a quick snapshot of how to import your datasets into Data Prep.

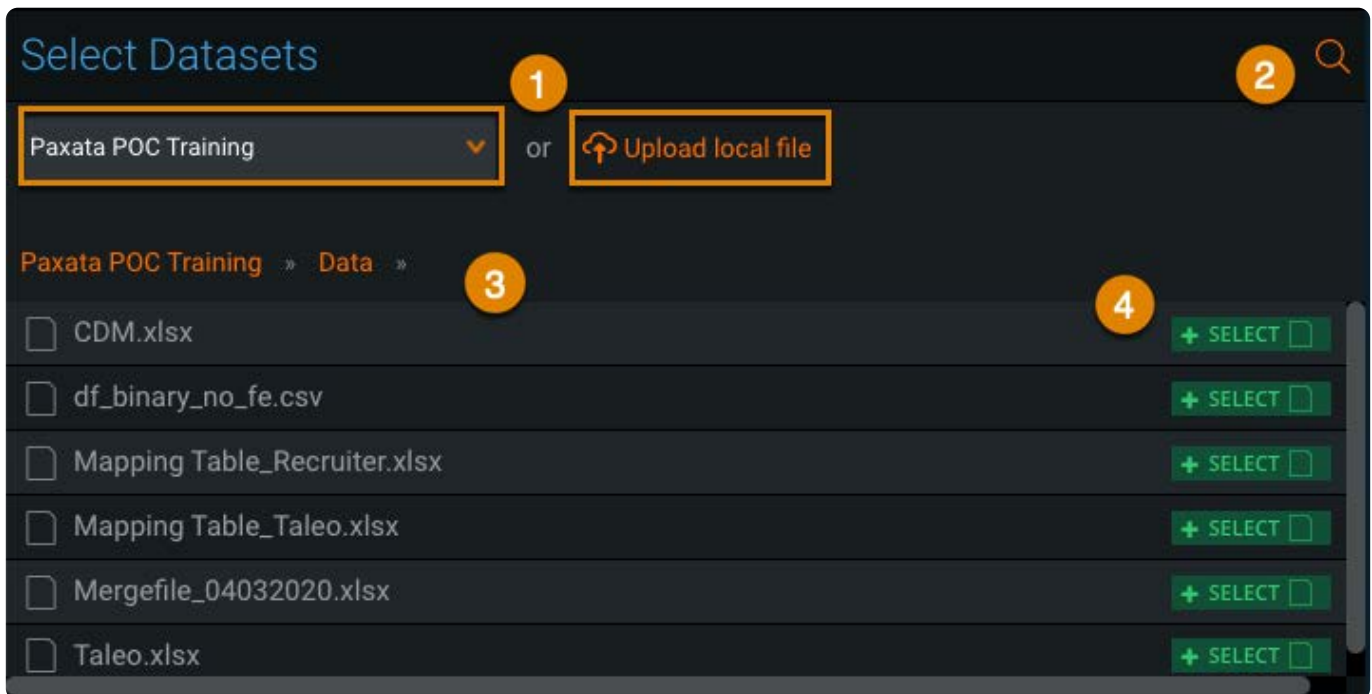
1. On the **Library** page, click **Import**.
2. On the **Import** page, you can [select datasets](#), [search for datasets](#), or [combine datasets](#).
3. Check the preview of the dataset. Does your data look correct?
 - If your data is correct, you can continue adding datasets until you have selected all your datasets for import.
 - If your data is incorrect, try [adjusting the import settings](#).
4. Click **Finish**.

Your data is imported as a dataset and ready to be prepped in a project.

Select datasets

You can import datasets from a local file or a connected data source. This section provides more detail on how to select one or more datasets for import.

Select Datasets pane



Following is an overview of the elements of the **Select Datasets** pane.

#	ELEMENT	DESCRIPTION
1	Data Source options	<p>Maybe you need to import a dataset from Amazon S3, JDBC, or some other data source. Maybe you just want to import a spreadsheet you saved on your computer. Either way, this is where you do it.</p> <p>The Data Source list lets you select a configured data source. Your administrator must connect to the data source before you can access it.</p> <p>Click Upload local file to select a dataset from your computer.</p>
2	Search	<p>For times when you want to find a specific dataset or a group of similar datasets, enter search criteria. The Search field accepts wildcard characters. This will help you find specific and similarly named datasets. See Search for datasets for details.</p>
3	Datasets Lists pane	<p>The contents of the selected data source are listed here. In this example, the data source has six items: one comma-separated value (CSV) file and five excel files.</p>
4	Select	<p>When you see a dataset you want to import, click Select. The dataset is listed in the You selected pane and will be imported when you click Finish.</p>

Select datasets from local files

To select a dataset from a file on your computer or shared network drive:

1. Click the **Upload File** pane and select a dataset or drag a file to the pane.

The dataset is added to the list in the **You selected** pane. Data Prep displays the **Your options** pane for the dataset and a preview of the dataset.

2. To add more datasets, click any additional dataset you want to include in the import.

The additional datasets are added to the list in the **You selected** pane.

Select datasets from a data source

To select a dataset from a connected data source:

1. Click **Select Data Source** and choose the data source you want to use.
2. Locate the dataset you want to import.

To locate your dataset using search, see [Search for datasets](#).

3. To select a dataset, click **Select**.

The dataset is added to the list in the **You selected** pane. Data Prep displays the **Your options** pane for the dataset and a preview of the dataset.

4. To add more datasets from the currently selected data source, click any additional dataset you want to include in the import.
5. To add more datasets from a different data source, repeat steps 1 - 3 for each data source.

The additional datasets are added to the datasets list in the **You selected** pane.

Search for datasets

You can search for datasets by typing the name of the dataset or entering a query string. The search is case sensitive and only the results that exactly match your search criteria are returned. You can use wildcard characters to locate a dataset when you aren't sure what the exact name is or to locate similarly named datasets.

Search for a dataset

To search for a dataset:

1. Select a data source and click the **Search** icon on the top right of the **Select Datasets** pane.

The **Search** icon appears only when you select a data source, not when you upload a local file.

2. In the **Wildcard Search** field, type your search criteria.

Datasets that match your search criteria exactly are returned.

See [Wildcard characters](#) for help setting search criteria.

Query a database

To query a database:

1. Click **Select Data Source** and choose the data source you wish to use.
2. Click **Create Query** on the bottom right of the **Select Datasets** pane.
3. In the **Query String** field, type your search criteria.

To search with wildcard characters, see [Wildcard characters](#).

Datasets that match your search criteria exactly are returned.

Wildcard characters

Following are the wildcard characters you can use to search for datasets.

CHARACTER	MATCHES
*	Any number of characters, including none
?	A single character
[0-9] or [a-z]	A character in the range given in the bracket
[123] or [abc]	A character listed inside the bracket

Example searches using wildcards

Here are some example searches and the results:

SEARCH EXAMPLE	RETURNS
*	All the datasets
*.csv	Datasets with a '.csv' file extension.
a?b.csv	Datasets that that are named aac.csv, abc.csv, ..., azc.csv..
a*z.csv	Datasets that begin with a lowercase 'a' and end with 'z.csv' regardless of what characters or how many characters are between
a[0-9].csv	Datasets that are named a0.csv, a1.csv, a2.csv, ..., a9.csv

SEARCH	RETURNS
EXAMPLE	
a[a-z].csv	Datasets that are named aa.csv, ab.csv, ..., az.csv
a[abc].csv	Datasets that are named aa.csv, ab.csv, ac.csv

Combine datasets

Data Prep can combine multiple datasets into one *glob* to be imported. A glob is the result of appending multiple datasets into one dataset during import. This section provides more information on how to glob multiple datasets together prior to import.

Guidelines for combining datasets

Following are some guidelines that help make globbing datasets a success.

- Datasets can only be globbed from the same data source.
- Datasets can only be globbed through a wildcard search.
- Datasets being globbed together should have the same structure (number of columns and type of data).

Data sources that support globbing

For a list of data sources and file formats that are supported for globbing, review the Platform Support matrix in the latest [Data Prep Release Notes](#).

Create a glob

To combine multiple datasets into one glob:

1. Click **Select Data Source** and select the data source.
2. Use [search](#) to locate the datasets you want to combine.
3. Click **Combine All Results**.

The datasets are combined into one glob. The glob is added to the datasets list in the **You selected** pane. The name of the glob defaults to the search criteria. Data Prep displays the **Your options** pane for the glob and a preview of the glob.

Preview a dataset before import

To change the dataset in the preview, from the **You selected** pane, click the dataset you want to preview.

The **Preview** pane displays the selected dataset.

By default, Data Prep displays a preview of the last selected dataset.

Add a dataset again

During import, there might be times when you need to apply different import options to the same dataset. This is especially true when you need to import more than one Excel worksheet from the same Excel file.

To add a dataset with different import options:

1. From the **You selected** pane, click the **More** button (three vertical dots) of the dataset you want to add again.
2. Click **Add Again**.

The dataset is added to the list in the **You selected** pane.

3. [Adjust the import settings](#) as needed.

Adjust import settings

Once a dataset is selected, Data Prep analyzes your data to determine the right settings for the best results. But data isn't a one-size-fits all kind of thing. Sometimes, you need to tweak the settings to get them just right. This section provides information on how to adjust some of the more universal settings of a dataset prior to import. For specific information about a setting, hover your cursor over the help tip (question mark) button.

Following are a few of the frequent and more basic settings you can adjust:

ACTION	STEPS
Add a tag.	In the Your options pane, type or select the tag from the Tags list.
Add a column to show the source file lineage.	<p>In the Your options pane, toggle the Add column to show source file button.</p> <p>The new Source File column is added to the end of the dataset showing the path of the source file for every imported row.</p>
Change the format of the dataset.	In the You selected pane, select the format you want to apply to the dataset from the Format menu. See Supported Formats for more information.
Change the name of a dataset.	In the Your options pane, type the new name in the Name field. The dataset name is updated in the You selected pane.
Exclude columns from import.	<p>In the Preview pane:</p> <ol style="list-style-type: none">1. Click Edit columns.2. Deselect the columns you don't want to import.3. Click Show preview. <p>The deselected columns are removed from the Preview.</p>

ACTION	STEPS
Import additional worksheets from the same Excel file.	For each additional worksheet: <ol style="list-style-type: none"> <li data-bbox="597 254 1154 285">1. In the You selected pane, add the Excel file again. <li data-bbox="597 306 1438 373">2. In the Your options pane, select the worksheet to import from the Worksheet menu.
Rearrange the columns.	In the Preview pane: <ol style="list-style-type: none"> <li data-bbox="597 474 829 506">1. Click Edit columns. <li data-bbox="597 527 1458 558">2. Click the up arrow or down arrow until the column is in the position you want. <li data-bbox="597 579 837 611">3. Click Show preview.
Rename a column.	In the Preview pane: <ol style="list-style-type: none"> <li data-bbox="597 705 829 737">1. Click Edit columns. <li data-bbox="597 758 1125 789">2. Click Edit (pencil icon) and type the new name. <li data-bbox="597 810 837 842">3. Click Show preview.

Supported Formats

For file-based connectors, the common formats are listed in the following table. Data Prep's Intelligent Ingest identifies the format of the file by looking into the contents of the file instead of relying on the file extension. Even if your file does not have an extension or has an incorrect extension, Data Prep correctly identifies the format.

COMMON FORMAT	IMPORT SUPPORT FOR WILDCARDS AND GLOBBING
Delimited files (CSV, TSV, etc.)	Yes
Fixed-width column data	Yes
JSON	Yes
XML	Yes
Apache Avro	Yes
Microsoft Excel (XLS, XLSX)	No. See Wildcard characters and Guidelines for combining datasets .
SAS BDAT	Yes

Data Prep supports the import of compressed files in one of the following formats: Deflate, LZ4, Snappy, ZIP, Gzip, or Bzip. In general, the decompressed file must be a common format as listed in the previous table.

Additionally, connectors that support Parquet files also support compressed versions of Parquet files.

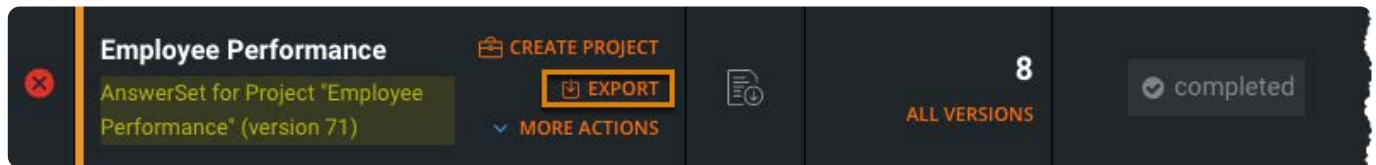
Note

When importing a ZIP file that contains multiple files, the largest file in the compressed set is automatically identified and selected for import to the Library.

Export datasets

You can export datasets and AnswerSets out of Data Prep to a file on your computer or to a connected data source. Exporting your data gives you the flexibility to share your data with other people or other systems.

Export a dataset



To export a dataset or AnswerSet:

1. On the **Library** page, hover over the dataset you want to export and click **Export**.
2. On the **Exporting** page, select a **Data Source** or click **Download locally**.
3. In the **Export Settings** pane, [adjust the settings](#) as needed.
4. Click **Export**.

The dataset is exported to the location you specified.

Adjust the export settings

The export settings allow you to define the structure of your dataset for export. The settings available for a dataset depend on the format you've chosen for the export.

The following are a few of the frequent and more basic settings you can adjust:

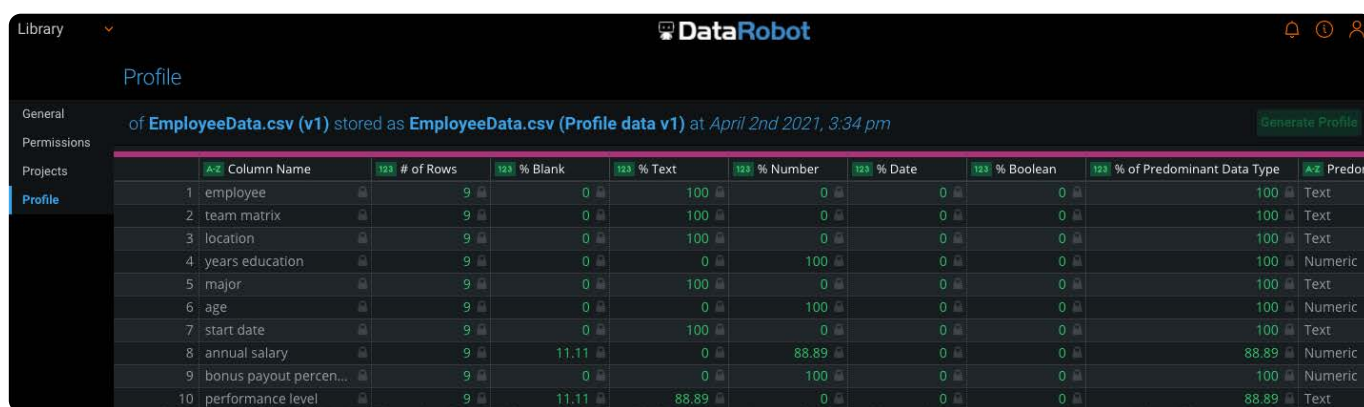
ACTION	STEPS
Adjust the format of the dataset.	Select the format you want to apply to the dataset from the Format menu.
Update the name of the dataset.	Type the new name in the Name field.
Update the character set used.	Select a different character encoding from the Character encoding menu.
Note: Change from an English character set to Japanese character set.	Note: The byte order mark (BOM), also known as the unicode signature BOM, particularly applies to extended character sets like UTF-16 and UTF-32.

Profile datasets

Note

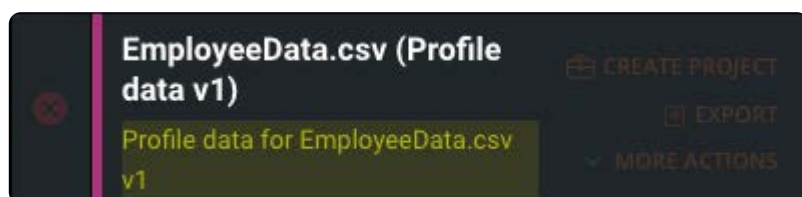
Your Data Prep Administrator must enable this feature in your application.

When you profile a dataset, you generate statistics about the data in that dataset. The results are displayed on the dataset's **Profile** page:



	A-Z Column Name	123 # of Rows	123 % Blank	123 % Text	123 % Number	123 % Date	123 % Boolean	123 % of Predominant Data Type	A-Z Predom
1	employee	9	0	100	0	0	0	100	Text
2	team matrix	9	0	100	0	0	0	100	Text
3	location	9	0	100	0	0	0	100	Text
4	years education	9	0	0	100	0	0	100	Numeric
5	major	9	0	100	0	0	0	100	Text
6	age	9	0	0	100	0	0	100	Numeric
7	start date	9	0	100	0	0	0	100	Text
8	annual salary	9	11.11	0	88.89	0	0	88.89	Numeric
9	bonus payout percen...	9	0	0	100	0	0	100	Numeric
10	performance level	9	11.11	88.89	0	0	0	88.89	Text

Your profile is also automatically saved in the library with a name to indicate that it's a profile type AnswerSet:



How can I use profiles of my data?

Getting data is like getting a package, only you don't know what's in the package. With a package, you have a packing slip so you don't have to dig in and tear through everything to know what is in it. With Data Prep, you can create a profile of your data so that you can quickly understand it.

Data profiles are essential for determining the quality of data in a dataset before you begin working with that data. For example, you can quickly determine if there are mixed types in the data, nulls, non-printable characters, and patterns that don't belong.

1	A	\$1.00	#	Z
2	3	\$5.00	@	Y
3	C	\$0.75	%	
4	D	\$200.00	&	W
5	E	\$0.50	*	V

Based on the data profile, you can address quality issues by bringing the data into a Data Prep project.

As you continue to update versions of a dataset in the library—through either manual or automated import—you can continue to profile each subsequent version. In this way, you can monitor the data quality, version over version, and you can remediate as necessary.

What is the meaning of each column that I see in a profile's AnswerSet?

When you profile a dataset, the result is an AnswerSet that has a row to represent each column in your dataset. Each column in the profile AnswerSet provides statistics about the columns in your dataset.

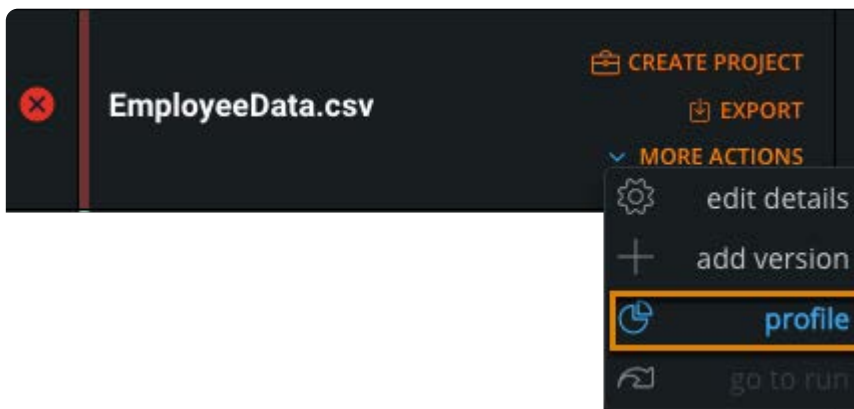
Following are the statistics included in the data profile for each column:

COLUMN NAME	DEFINITION
# of Rows	Total number of rows in the dataset
% Blank	Percentage of blanks in the column
% Text	Percentage of text values in the column
% Number	Percentage of numeric values in the column
% Date	Percentage of date values in the column

COLUMN NAME	DEFINITION
% Boolean	Percentage of boolean values in the column
% of Predominant Data Type	Percentage of values in the column that contain the most dominant data type
Predominant Data Type	Most dominant data type in the column
# of Unique Values	Number of unique values in the column
# of Phonetically Unique Values (Metaphone)	Number of unique values in the column after clustering like values using the metaphone (sounds like) algorithm—for example "Good Samaritan" and "Good Samaritan" are clustered to count as the same value.
% Possible Phonetic Duplicates (Metaphone)	Number of phonetically unique values (metaphone) / the number of unique values. This ratio indicates the possibility of duplicates in the column. A higher number indicates a higher probability of values that are potential duplicates, and indicates that you may need to do a cluster-and-edit operation on the column to identify duplicate values.
Top 5	Top five most common values in the column
Min String Length	Shortest string length in the column
Max String Length	Longest string length in the column
Avg String Length	Average string length for strings in the column
# of NAs or NONEs or NULLs	Number of times the column contains "na" or "none" or null
% All Upper Case	Percentage of cells containing all upper case characters
% All Lower Case	Percentage of cells containing all lower case characters
% with Non Standard ASCII Chars	Percentage of cells containing non printable characters, for example the control character
% with HTML Tags	Percentage of cells containing HTML tags

COLUMN NAME	DEFINITION
Avg # of Consecutive Spaces	Average number of consecutive spaces that are found in the column
% Negative Numbers	Percentage of cells containing negative
% Zeros	Percentage of cells containing the value zero

Create a profile for a dataset



To create a data profile:

1. On the **Library** page, hover over the dataset for which you want to create a data profile.
2. Click **More Actions** and select **profile**.
3. On the **Profile** page, click **Generate Profile** on the top right.

The profile appears in the **Profile** pane. In addition, the profile is automatically saved as an AnswerSet in the library.

Note

The library preview of the AnswerSet is limited to the first 100 rows of the profile.

Update datasets with new data

Data is always changing. Even if you just imported data into Data Prep, there's a chance the data is outdated. Updating the data in a dataset allows you to import a dataset as a new version of an existing dataset. After you update a dataset, you can use the new version in an existing project.

When updating a dataset, you can update to a completely different dataset with new values, structure, and format. Or you can update to a dataset where:

- Only the values have changed—the structure and format are the same.
- The format or structure have changed, for example, columns are added or removed.

Update a dataset with new data

To update a dataset with new data:

1. On the **Library** page, hover over the dataset you want to update and click **More Actions**.
2. Select **add version**.
3. Locate and select the dataset you want to import from the **Select Data Source** list or click **Upload a local file**.

Note: If a SQL statement was used during the initial import, the SQL statement is retained and can be used again to update the data in the dataset.

The dataset is added to the list in the **You selected** pane. Data Prep displays the **Your options** pane for the dataset and a preview of the dataset.

4. Check the preview of the dataset and [adjust the import settings](#) as necessary.
5. Click **Finish**.

Your data is imported as a new version and is ready to be prepped in a project.

Update project datasets

When you add a dataset into your project, either as the base dataset, through a look up, or an append, you are identifying a specific version of the dataset to use in your project. If newer versions of the datasets become available in the library, your project doesn't automatically use the newer versions; the work you've done in your project and the subsequent results may depend on the specific dataset versions you initially selected.

Many times this works well. Other times, you may want to update the project datasets with newer versions.

There are two methods for updating a project's datasets:

- Refresh a project dataset to the newest version of an existing dataset.
- Replace a project dataset with another dataset.

A dataset refresh updates the project data to use the most current version of a dataset.

For example, if you start your project with a dataset that is version 1 and over time newer versions of the dataset are imported into the library (either through manual import or automation,) you have the option to refresh the dataset in your project to use the newest version.

Refresh the datasets

To refresh a dataset to the latest version:

1. In the project, click **steps** in the **Tools** bar.
2. At the bottom of the **Steps** tool, click **Refresh Datasets**.

The **Refresh datasets** pane appears. All datasets that can be refreshed are selected by default.

3. Select the datasets to refresh; you can select **All** or select individual datasets.
4. Click **Save**.

The project data is updated to the most current versions of the selected datasets.

When can a dataset be refreshed?

A dataset can be refreshed when:

- There is a [newer version of the dataset](#) in the library.
- The [Interactive Mode](#) feature is enabled and the interactive portion size for datasets has changed.

When a dataset can be refreshed in a project, you are provided visual cues:

- If the **Refresh Datasets** button is green, a newer version of one or more datasets used by your project is detected. If the button is gray, there are no newer versions of your project's datasets.
- If the **Use Latest** button located on the **Refresh datasets** pane is green, a newer version of the dataset is available.
- The **file details** link opens a **Version Information** pane that allows you to quickly determine the number of new rows and columns in the dataset's latest version. If your project is in Interactive Mode and the dataset contains more rows than the interactive portion, you will also see a column for Interactive that lists the number of rows you can bring into the project. This number is important because it allows you to quickly determine if the Interactive portion has been increased or decreased, and then you can determine if you do want to refresh the dataset.

Note

All Data Prep projects have a maximum project row limit that is set by the Data Prep System Administrator. If you are close to reaching that limit, and your Administrator cannot increase it, you can selectively choose which datasets to update with latest versions so that you can continue bringing newer data into your project without exceeding the project row limit.

If you deselect a dataset, the **Use Latest** button turns dark gray. This indicates that there is a newer version of the dataset and that you have chosen not to update the dataset.

When there are no new versions for the dataset, the **Use Latest** button is light gray.

Replace a dataset

Unlike refreshing the data, replacing a dataset lets you decide which dataset or which specific version of a dataset to use in your project. For example, if you started a project with version 1 of a dataset and five additional versions were imported, replacing a dataset lets you pick the exact version to use, which may not be the latest version. Replacing the dataset also gives you the power to entirely change which dataset is being used in your project.

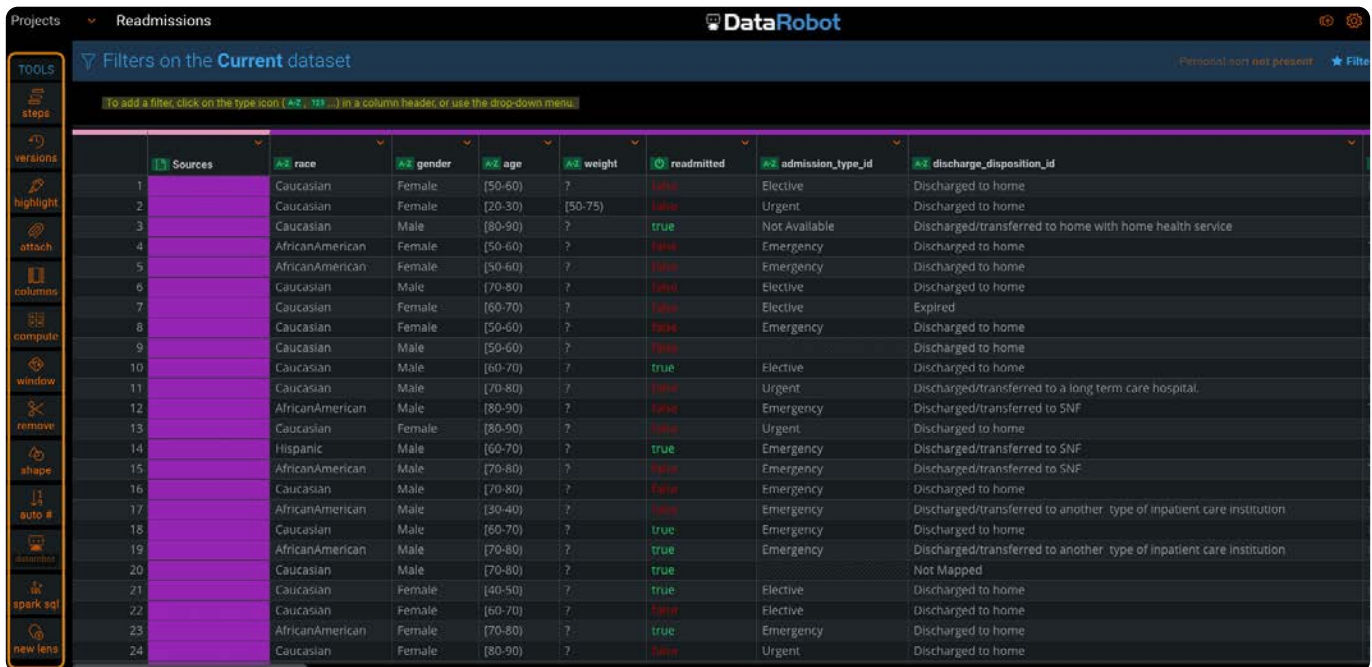
To replace a dataset used in a project:

1. In the project, click **steps** in the **Tools** bar.
2. In the **Steps** tool, click the step with the dataset you want to update and click **Edit** at the top.
The project returns to the state it was in when the selected step was created.
3. Above the **Data preview** pane, click the name of the dataset you want to update.
4. On the **Select Datasets** page, select the dataset you want to use.
 - To select a previous version of a dataset, on the dataset, click **All Versions**. On the version you want to use, click **Select**.
 - To select a different dataset in your library, on the dataset, click **Select**.
5. Click **Save**.

The project data is updated.

Work with project tools

Data Prep provides many tools to help you clean, shape, combine, and ultimately prep your data. Access these tools in the project **Tools** bar:



The screenshot shows the DataRobot interface for a project named 'Readmissions'. The main area displays a table with 24 rows and 8 columns. The columns are: Sources, race, gender, age, weight, readmitted, admission_type_id, and discharge_disposition_id. The 'readmitted' column contains boolean values (true/false). The 'discharge_disposition_id' column contains various discharge codes. A 'TOOLS' sidebar is visible on the left with icons for steps, versions, highlight, attach, columns, compute, window, remove, shape, auto #, spark sql, and new lens. The top bar shows 'Filters on the Current dataset' and a 'Personal icon not present' notification.

	Sources	race	gender	age	weight	readmitted	admission_type_id	discharge_disposition_id
1		Caucasian	Female	[50-60]	?	false	Elective	Discharged to home
2		Caucasian	Female	[20-30]	[50-75]	false	Urgent	Discharged to home
3		Caucasian	Male	[80-90]	?	true	Not Available	Discharged/transferred to home with home health service
4		AfricanAmerican	Female	[50-60]	?	false	Emergency	Discharged to home
5		AfricanAmerican	Female	[50-60]	?	false	Emergency	Discharged to home
6		Caucasian	Male	[70-80]	?	false	Elective	Discharged to home
7		Caucasian	Female	[60-70]	?	false	Elective	Expired
8		Caucasian	Female	[50-60]	?	false	Emergency	Discharged to home
9		Caucasian	Male	[50-60]	?	false		Discharged to home
10		Caucasian	Male	[60-70]	?	true	Elective	Discharged to home
11		Caucasian	Male	[70-80]	?	false	Urgent	Discharged/transferred to a long term care hospital.
12		AfricanAmerican	Male	[80-90]	?	false	Emergency	Discharged/transferred to SNF
13		Caucasian	Female	[80-90]	?	false	Urgent	Discharged to home
14		Hispanic	Male	[60-70]	?	true	Emergency	Discharged/transferred to SNF
15		AfricanAmerican	Male	[70-80]	?	false	Emergency	Discharged/transferred to SNF
16		Caucasian	Male	[70-80]	?	false	Emergency	Discharged to home
17		AfricanAmerican	Male	[30-40]	?	false	Emergency	Discharged/transferred to another type of inpatient care institution
18		Caucasian	Male	[60-70]	?	true	Emergency	Discharged to home
19		AfricanAmerican	Male	[70-80]	?	true	Emergency	Discharged/transferred to another type of inpatient care institution
20		Caucasian	Male	[70-80]	?	true		Not Mapped
21		Caucasian	Female	[40-50]	?	true	Elective	Discharged to home
22		Caucasian	Female	[60-70]	?	false	Elective	Discharged to home
23		AfricanAmerican	Female	[70-80]	?	true	Emergency	Discharged to home
24		Caucasian	Female	[80-90]	?	false	Urgent	Discharged to home

These pages describe the project tools:

TOPIC	DESCRIBES...
Edit steps	View, edit, add, rearrange, and delete steps within a project.
Manage project versions	Review the version histories of a project.
Highlight data	Highlight patterns, spaces, and numeric ranges, then use column operations to update values.
Attach datasets	Perform lookup, join, and append operations on your dataset.
Update columns	Edit the column names, their order, and their availability in a project.
Compute columns	Compute columns using functions to create a new column.
Group with windows	Group rows to perform aggregations, shifting, and ranking.

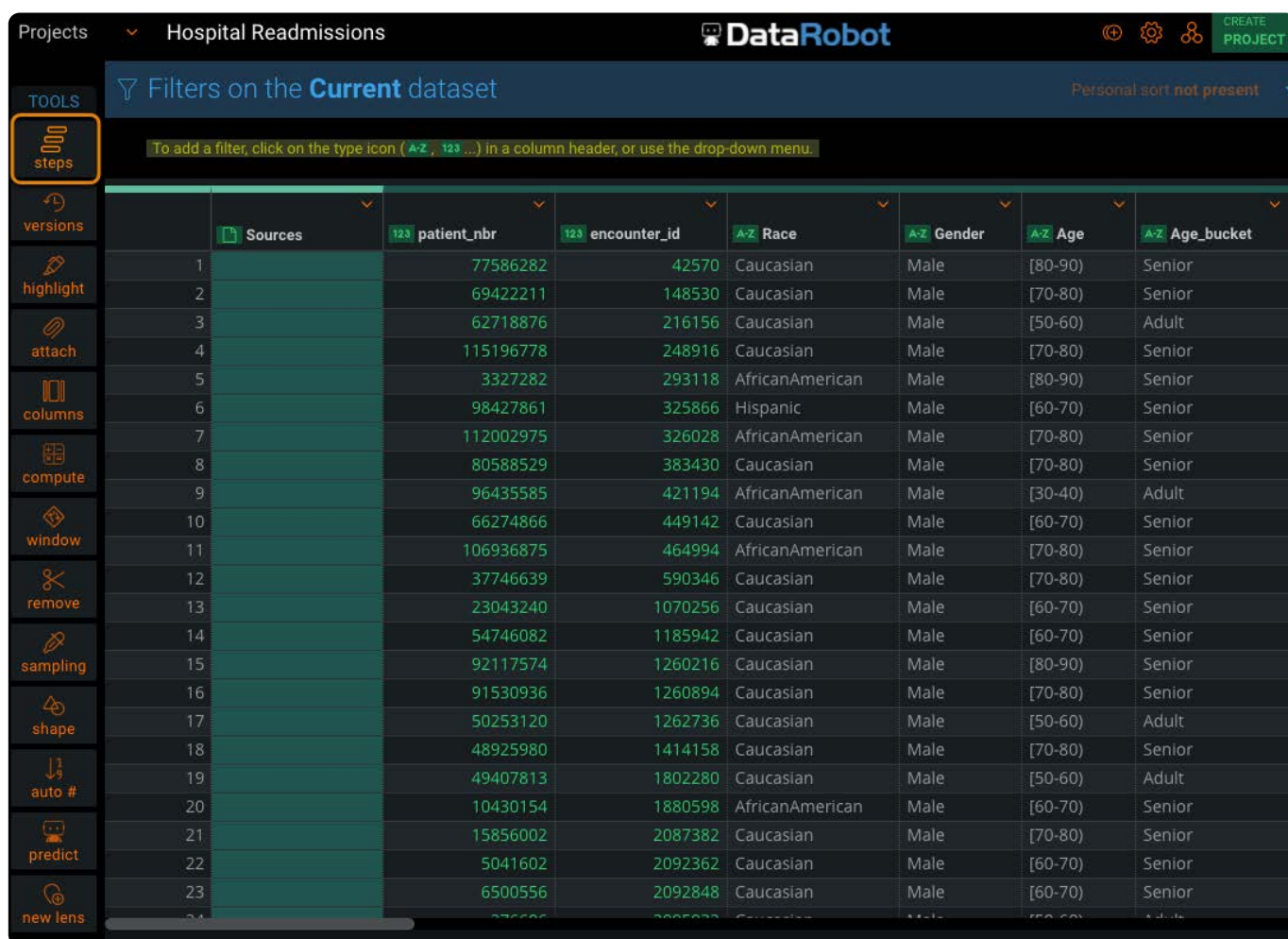
TOPIC	DESCRIBES...
Remove rows	Remove rows that don't meet your needs.
Sample data	Sample project data.
Shape data	Use data shaping tools to deduplicate, group, transpose, and pivot data.
Autonumber rows	Generate a new column that autonumbers the rows of your dataset.
Make predictions	Learn how to score your data using the predict tool.
Transform your data with Spark SQL	Perform data transformations using Spark SQL.
Create a DataRobot project from a DR lens	Use the Data Prep Build tool to create a DataRobot machine learning project from a Data Prep DR lens.
Use lenses to select steps for publishing	Identify steps to be published to AnswerSets using lenses.

Work with steps

The Data Prep **Steps** tool lets you view, edit, add, rearrange, and delete steps within a project. From the **Steps** tool, you can export a snapshot of your data, called an *AnswerSet*. You can export the results of all of your data prep steps or you can select a particular step and export the AnswerSet resulting from the data prep activities up to the selected step.

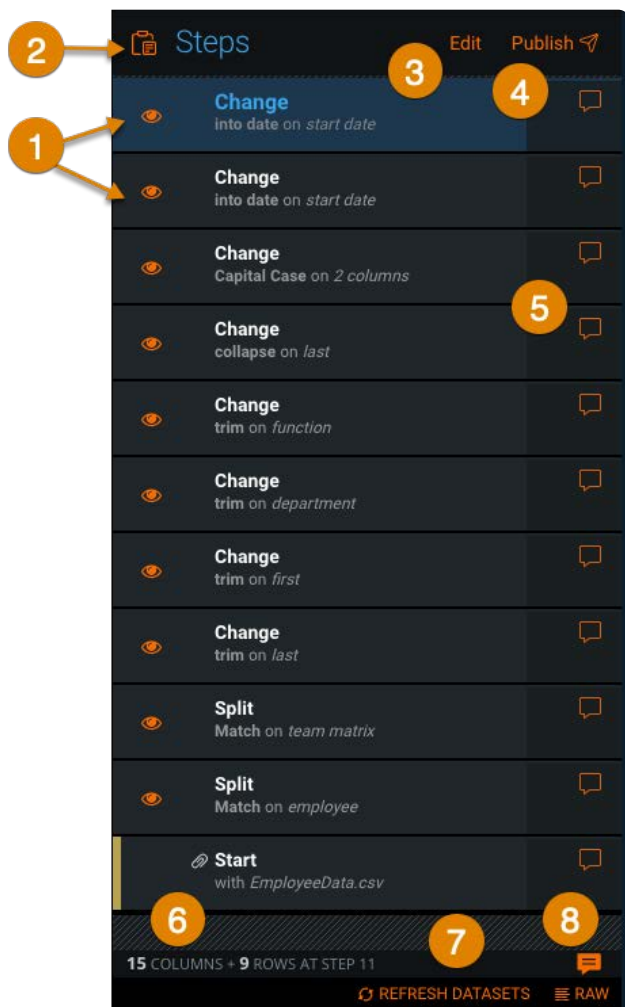
Work with the Steps tool

To access the **Steps** tool, click **steps** in the project **Tools** bar:



Following is an overview of the elements in the **Steps** pane.

The steps display sequentially. The first step is at the bottom and the most recent step is at the top.



ELEMENT	DESCRIPTION
---------	-------------

- | | |
|---|--|
| 1 | Steps history Displays the steps taken so far in the project, from the Start step at the bottom up to the most recent step at the top of the pane. |
| 2 | Copy Steps Copies the project steps to a clipboard or a file for reuse in a different project. See Reuse project steps for details. |
| 3 | Edit Opens the Steps tool in editor mode so that you can make changes to any step in the project. |
| 4 | Publish Publishes your project as an AnswerSet in the Library. An AnswerSet is the published result of your data prep. |
| 5 | Annotations Provides an annotation for any step in your project. See what you can do with the steps tool for details. |
| 6 | Refresh datasets Refreshes datasets in the project so that they use the most current versions. See Update project datasets for details. |

ELEMENT	DESCRIPTION
7	Project stats Displays stats about your project—the number of columns and rows at any step in your project. Click on any step to view the stats for that step.
8	Debugging Displays a raw script of your project in JSON format. Use this is only for debugging purposes if there is ever a problem with your project.

What you can do with the Steps tool

Following are descriptions of the actions you can perform with the **Steps** tool.

ACTION	DESCRIPTION
View	Click a step to display the data as it appeared at that particular step. Click the eye icon next to a step to mute (hide) the step from the project without actually deleting it. Click the eye again to show that step again.
Edit	Use the Steps tool to edit steps already committed to the project at any time. When you select a step, the data adjusts to display how it appeared at that particular step in the process. Click Edit in the upper right of the Steps tool to edit the currently selected step. (Be sure to click Save in the upper right following any edits, otherwise your changes will not be kept.) When in edit mode, the Steps tool displays the number of active Data Filtergrams for each individual step of the project, also.
Add a Step	Anytime you save an action within your project, it is added to the steps at the top of the Steps tool in reverse chronological order (most recent actions are at the top). It's also possible to add a new step at any historical point in the project. Whenever you click a step in the Steps tool, your data displays as it historically appeared at the completion of that particular step. If you take a new action on the data while this historical view is shown, that new action (step) will be added into the Steps tool directly following the step you initially selected.
Rearrange Steps	To rearrange steps, you click, drag, and drop steps within the Steps tool. As you rearrange steps, the data automatically refreshes to reflect the new changes that result from the rearrangement. If the rearrangement of steps causes errors with any of the steps, an alert displays on the Steps button, as well as on any individual steps that have errors.
Delete Steps	You can click the eye icon to mutes (hides) steps within the project without completely deleting the step. If you click Edit in the upper right corner of the Steps tool, you will notice that the eye icon disappears, and when you hover over an individual step, an X icon displays instead. Click this X icon to delete the step from the project.

ACTION	DESCRIPTION
Add Annotations	You can provide an annotation for any step in your project. When you click the annotation button that appears on the right side of a step, a free-text field opens where you can enter up to 1,000 characters for the annotation.
Replayability	You can replay the same steps against new data. Once you have built your project with the steps necessary to produce your AnswerSet, edit the first step of the project and replace the Base dataset with a new Base dataset. This will result in all of the steps in the project automatically running on the new Base dataset.

Reuse project steps using copy and paste

With the **Steps** tool, you can copy steps from a Data Prep project and use them elsewhere in the same project or copy them into another Data Prep project. When you copy the steps, you can copy them to your computer's clipboard for a one-click paste. Optionally, you can copy the steps to a file, which you can save for later use or for sharing with other Data Prep users.

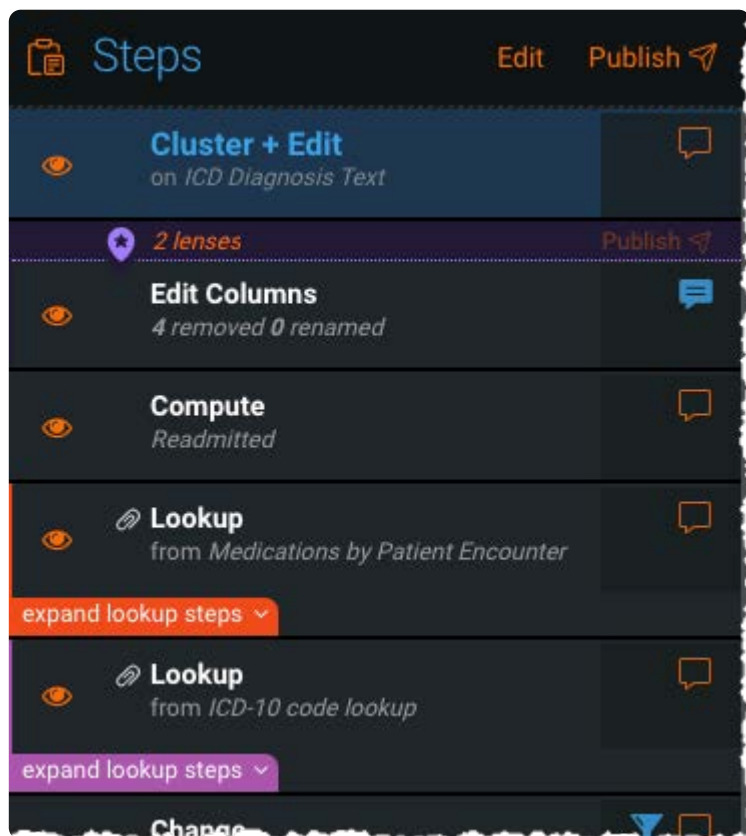
Note

- You can never copy the base dataset (the first step) in a project.
- If you expand a Lookup or Append step, you'll notice that you can select only the transformation steps that were applied during the import of the Lookup or Append. Those transformation steps are treated as individual steps when you paste them into another project. They cannot be nested under a Lookup or Append in the destination project.

Copy steps

To copy steps from a project:

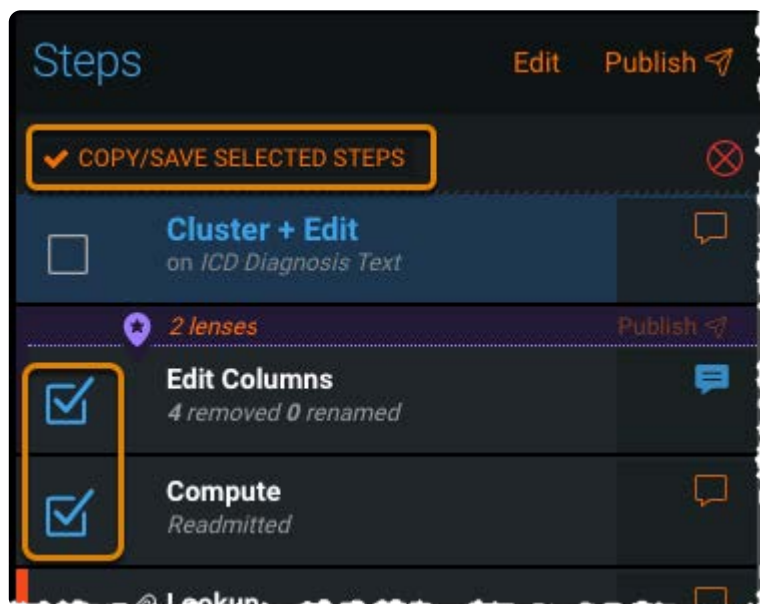
1. From the **Steps** pane, click the **Copy Steps** (clipboard) icon.



2. Click **Select Steps**.



3. Click the checkboxes for each step you want to copy and click **Copy/Save Selected Steps**.



Data Prep copies the steps to the clipboard. You can paste them into the current project, another project, or you can download them to a .txt file for later reuse.

These Steps have been copied to the clipboard.

You can now paste these here or into another Project (📄 → 📄 PASTE) or download them to use at a later time (📄 → 📁 LOAD).

Download as:

Comments:

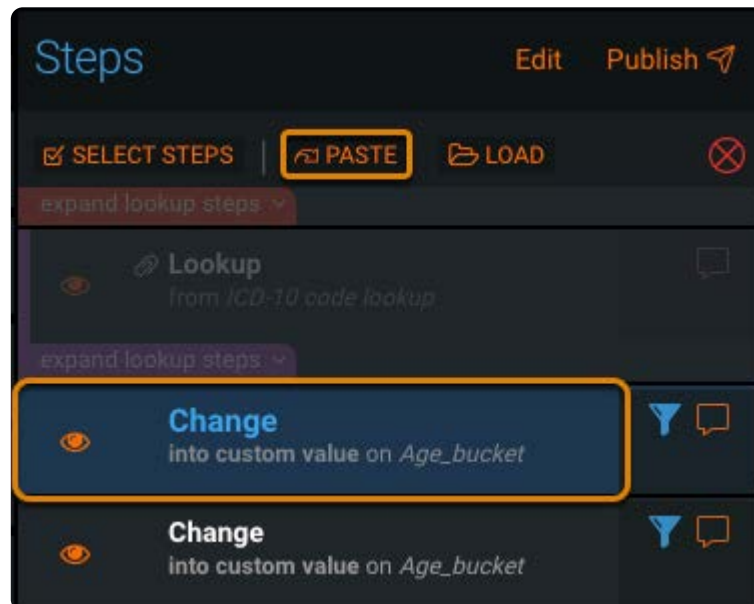
Download **OK**

Paste copied steps

1. From the **Steps** pane, click the **Copy Steps** (clipboard) icon.



2. Click the step in your project after which you want to paste the step (the circled **Change** step in this example) and click **Paste**.



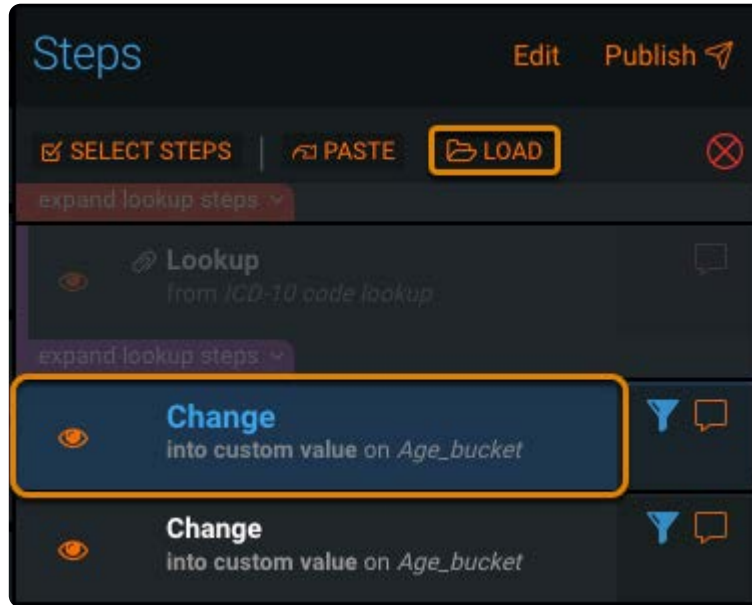
Data Prep pastes the copied steps after the selected step.

Load and paste copied steps

1. From the **Steps** pane, click the **Copy Steps** (clipboard) icon.



2. Click the step in your project after which you want to paste the step (the circled **Change** step in this example) and click **Load**.



3. Navigate to the saved steps file and click **Open**.

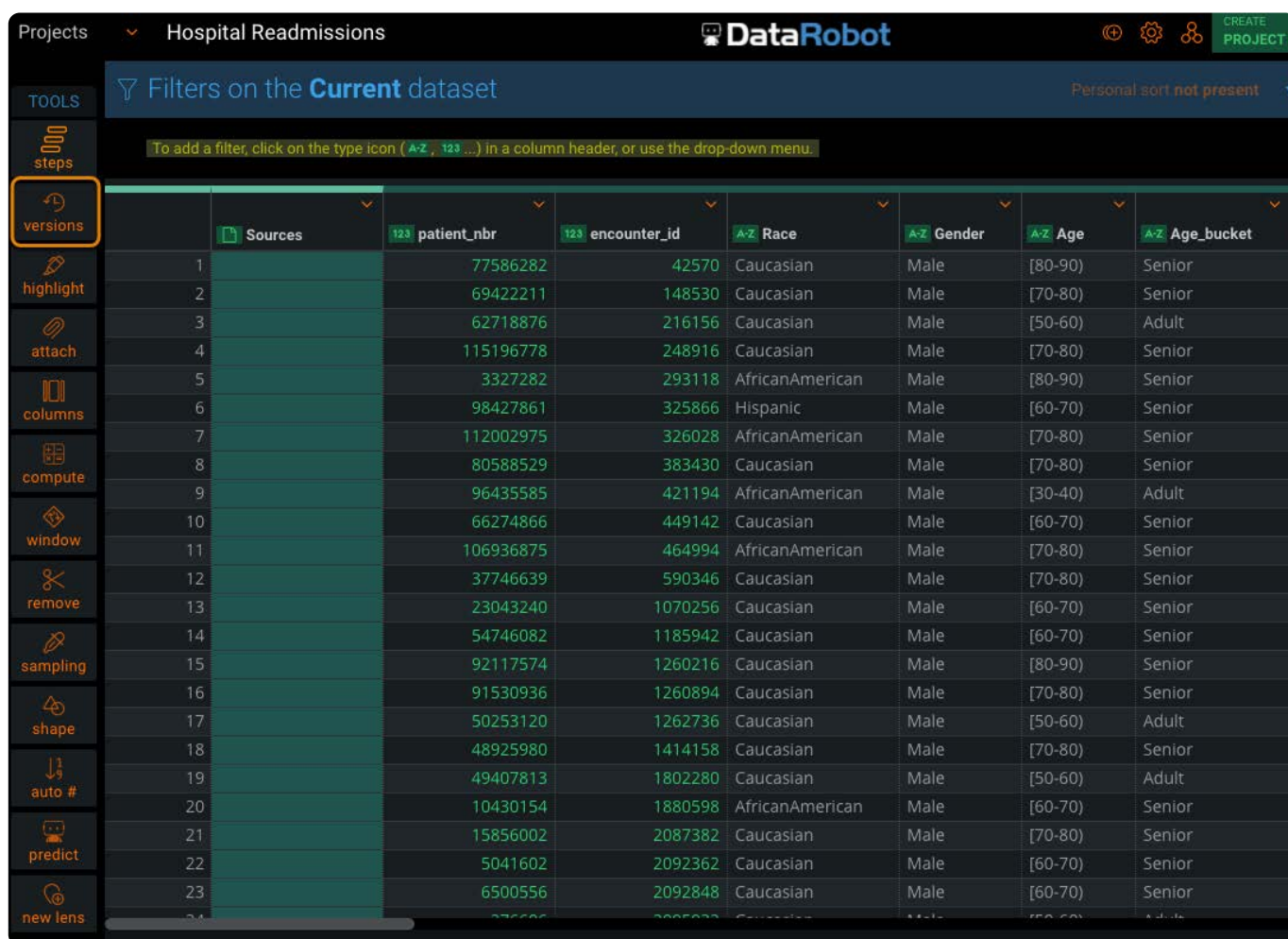
Data Prep pastes the copied steps after the selected step.

Manage project versions

Every time you take an action in your Data Prep project—for example, adding a step, removing a step, rearranging steps—a new version of your project is created.

Work with the Versions tool

To access the **Versions** tool, click **versions** in the project **Tools** bar:

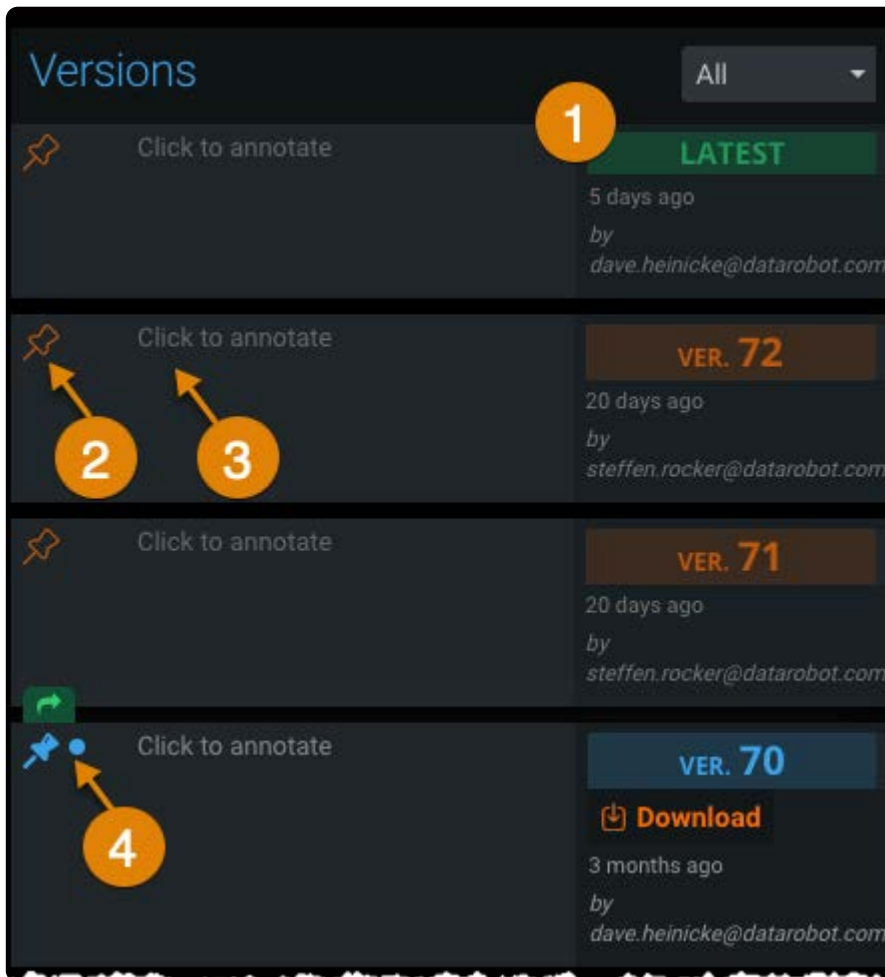


The screenshot shows the DataRobot interface for a project named "Hospital Readmissions". The "Tools" bar on the left has the "versions" tool highlighted. The main area displays a data table with the following columns: Sources, patient_nbr, encounter_id, Race, Gender, Age, and Age_bucket. The table contains 24 rows of data.

	Sources	patient_nbr	encounter_id	Race	Gender	Age	Age_bucket
1		77586282	42570	Caucasian	Male	[80-90]	Senior
2		69422211	148530	Caucasian	Male	[70-80]	Senior
3		62718876	216156	Caucasian	Male	[50-60]	Adult
4		115196778	248916	Caucasian	Male	[70-80]	Senior
5		3327282	293118	AfricanAmerican	Male	[80-90]	Senior
6		98427861	325866	Hispanic	Male	[60-70]	Senior
7		112002975	326028	AfricanAmerican	Male	[70-80]	Senior
8		80588529	383430	Caucasian	Male	[70-80]	Senior
9		96435585	421194	AfricanAmerican	Male	[30-40]	Adult
10		66274866	449142	Caucasian	Male	[60-70]	Senior
11		106936875	464994	AfricanAmerican	Male	[70-80]	Senior
12		37746639	590346	Caucasian	Male	[70-80]	Senior
13		23043240	1070256	Caucasian	Male	[60-70]	Senior
14		54746082	1185942	Caucasian	Male	[60-70]	Senior
15		92117574	1260216	Caucasian	Male	[80-90]	Senior
16		91530936	1260894	Caucasian	Male	[70-80]	Senior
17		50253120	1262736	Caucasian	Male	[50-60]	Adult
18		48925980	1414158	Caucasian	Male	[70-80]	Senior
19		49407813	1802280	Caucasian	Male	[50-60]	Adult
20		10430154	1880598	AfricanAmerican	Male	[60-70]	Senior
21		15856002	2087382	Caucasian	Male	[70-80]	Senior
22		5041602	2092362	Caucasian	Male	[60-70]	Senior
23		6500556	2092848	Caucasian	Male	[60-70]	Senior
24		376666	3005000	Caucasian	Male	[50-60]	Adult

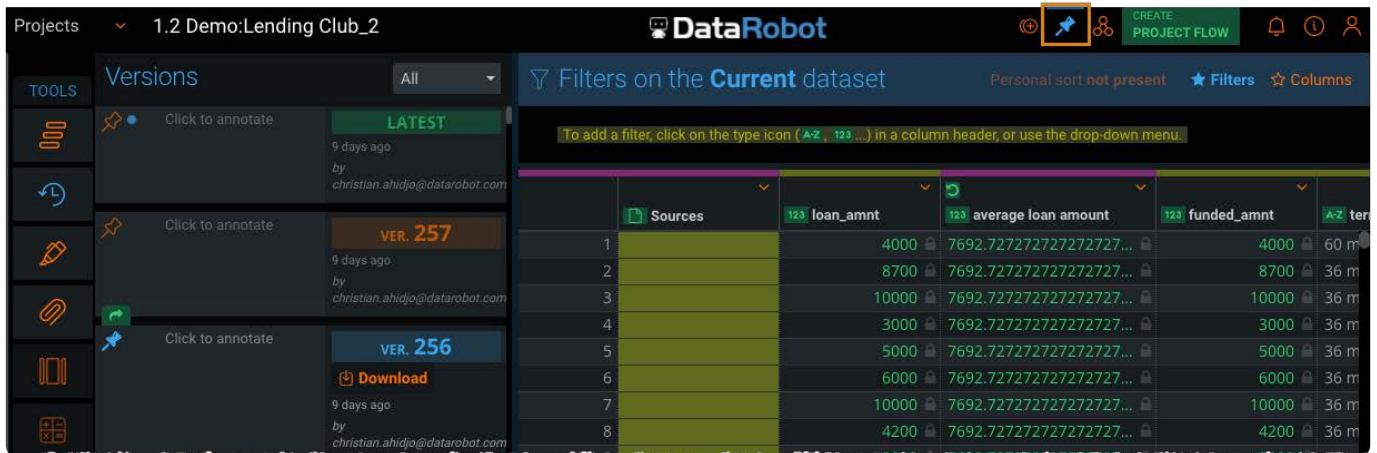
All versions are listed in the **Versions** pane. Click any version in the pane to view your project at that point in its history.

Following is an overview of the elements in the **Versions** pane:

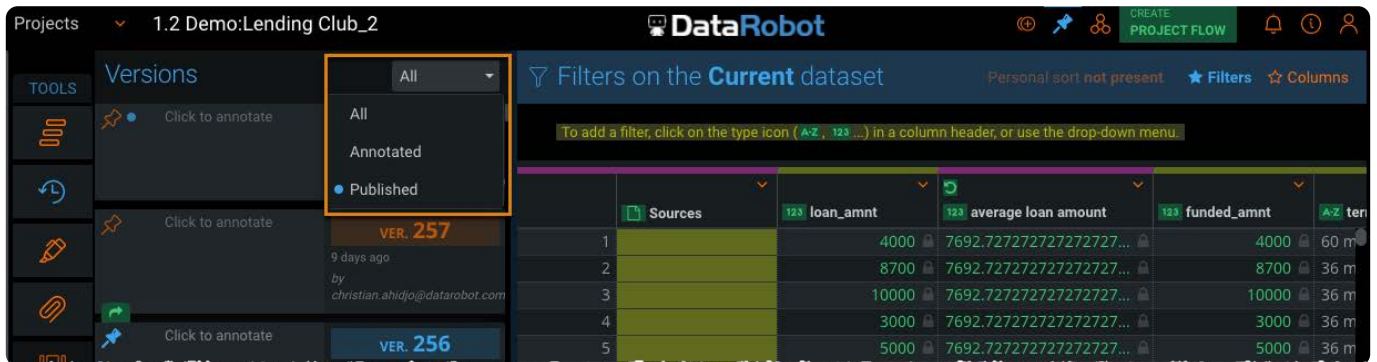


ELEMENT	DESCRIPTION
1	<p>Latest version</p> <p>The Latest version of your project is listed at the top of the pane.</p>
2	<p>pin icon</p> <p>Click the pin icon to view a previous version. You can also promote any earlier version of your project to be the latest version. To promote an earlier version:</p> <ol style="list-style-type: none"> 1. Click the orange pin adjacent to the version's name. 2. Select CLICK TO USE THIS VERSION. <p>The version is promoted to the most recent in the Versions pane.</p>
3	<p>Click to annotate</p> <p>You may find it helpful to annotate a project version to reflect the steps you took to render that particular version, for example: "Append Customer ID data." To annotate a version, select the Click to annotate text for the version and provide an annotation.</p>
4	<p>AnswerSet</p> <p>A blue dot appears adjacent to any version that has been published as an AnswerSet to the data library.</p>

Any time you are viewing a previous version of your project, a pin icon displays at the top of the project Preparation page to remind you that you're not viewing the most recent version. If you click the pin icon, you are immediately returned to the most recent version of your project.



By default, all versions of your project are shown in the pane. However, you can limit the types of versions that are displayed. From the drop-down menu in the top of the **Versions** pane, choose to view only **Annotated** Versions that have not been published or view only your **Published** versions.

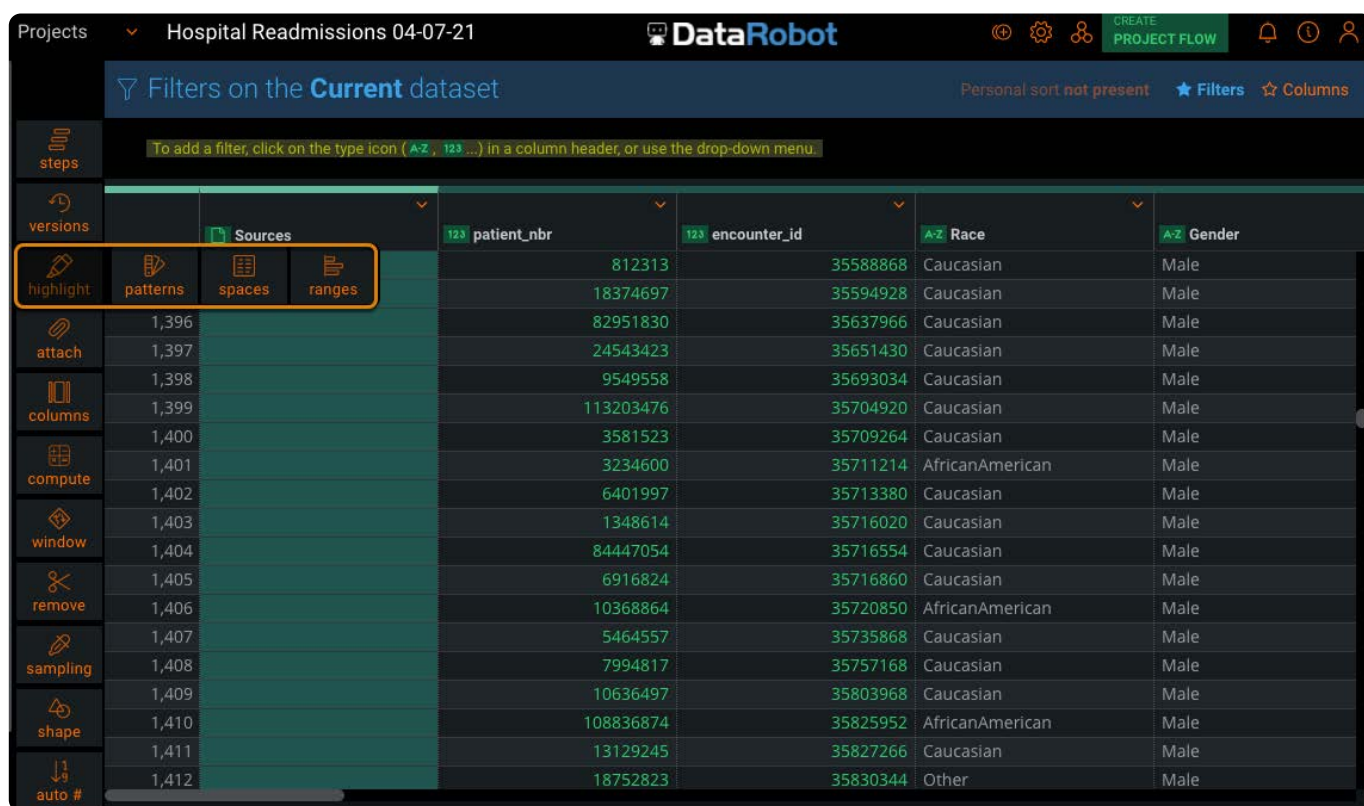


Highlight data

The Data Prep **Highlight** tools provide visual queues that help you to better understand your data. You can highlight patterns, spaces, and numeric ranges, then use column operations to update values.

Work with the Highlight tools

To work with highlighting, hover over the **Highlight** tools in the project **Tools** bar and select the [patterns](#), [spaces](#), or [ranges](#) tool.



The screenshot shows the DataRobot interface for a project named "Hospital Readmissions 04-07-21". The current dataset is displayed with columns: patient_nbr, encounter_id, Race, and Gender. The 'highlight' tool is selected in the 'Tools' bar, and the 'patterns', 'spaces', and 'ranges' sub-tools are also visible. The 'patient_nbr' and 'encounter_id' columns are highlighted in green, indicating that the 'highlight' tool is active. The 'Race' and 'Gender' columns are sorted by 'A-Z'.

patient_nbr	encounter_id	Race	Gender
812313	35588868	Caucasian	Male
18374697	35594928	Caucasian	Male
82951830	35637966	Caucasian	Male
24543423	35651430	Caucasian	Male
9549558	35693034	Caucasian	Male
113203476	35704920	Caucasian	Male
3581523	35709264	Caucasian	Male
3234600	35711214	AfricanAmerican	Male
6401997	35713380	Caucasian	Male
1348614	35716020	Caucasian	Male
84447054	35716554	Caucasian	Male
6916824	35716860	Caucasian	Male
10368864	35720850	AfricanAmerican	Male
5464557	35735868	Caucasian	Male
7994817	35757168	Caucasian	Male
10636497	35803968	Caucasian	Male
108836874	35825952	AfricanAmerican	Male
13129245	35827266	Caucasian	Male
18752823	35830344	Other	Male

Note

You can enable all Highlight tools simultaneously.

Highlight patterns

The **Highlight patterns** tool detects cells in your dataset that may be similar to one another and provides a color code to flag those cells.

For example, cells with nearby matching data values share the same color. The **Highlight patterns** tool is especially useful when used in conjunction with the **Filter** pane and **Columns** tool.

Highlight spaces

The **Highlight spaces** tool automatically highlights space characters in your dataset. Space characters appear as gray boxes. Space characters include the following: space, tab, carriage return, new line character, vertical tab character, form feed character.

Highlight ranges

The **Highlight ranges** tool provides automatic coloration of numeric columns to indicate where a value falls within the range of all values in the column.

The range displays similar in appearance to a bar graph:

- **All positive values:** Suppose the range of values is from 50 to 100. For a cell with value 75, a shaded bar extends halfway across the cell because 75 is halfway between 50 and 100. For cell value 100, the cell is entirely shaded.
- **All negative values:** All negative values display with the same basic appearance as for all positive values. For range -100 to -50, cell value -75 has its cell halfway shaded and cell value -50 has its cell entirely shaded.
- **Both positive and negative values:** If the column includes negative and positive values, for negative values, the shaded bar extends to the left from the point that represents zero. For positive values, the shaded bar extends to the right from the point that represents zero. If zero is not the midway point of the actual range of values, zero is proportionally positioned to the left or right of center. For example, if your range is from -10 to 1000, the point that represents zero is near the left edge of the cell.

Note

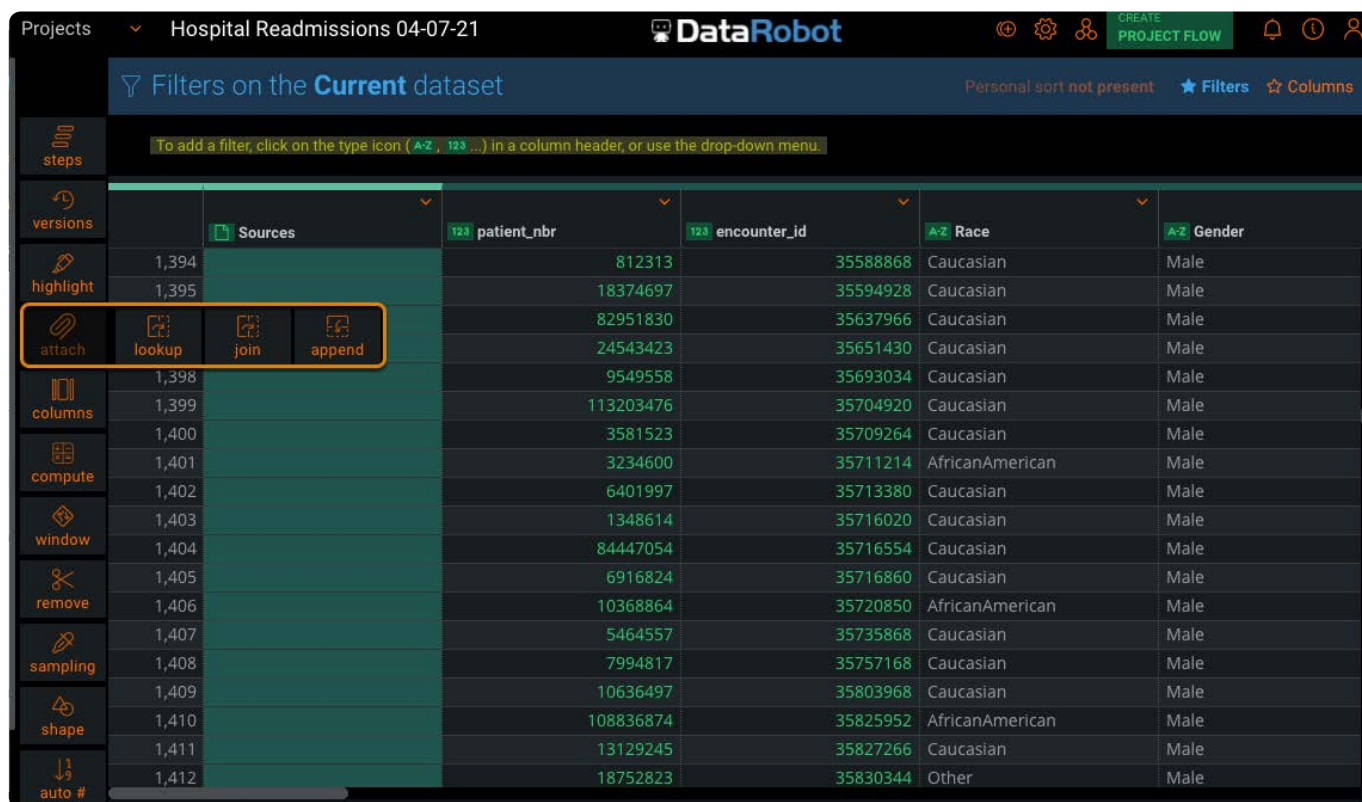
The highlighting does not affect your project's published AnswerSet. See [Export datasets](#) to learn about AnswerSets.

Attach datasets

The Data Prep **Attach** tools lets you perform lookup, join, and append operations on your dataset.

Work with the Attach tools

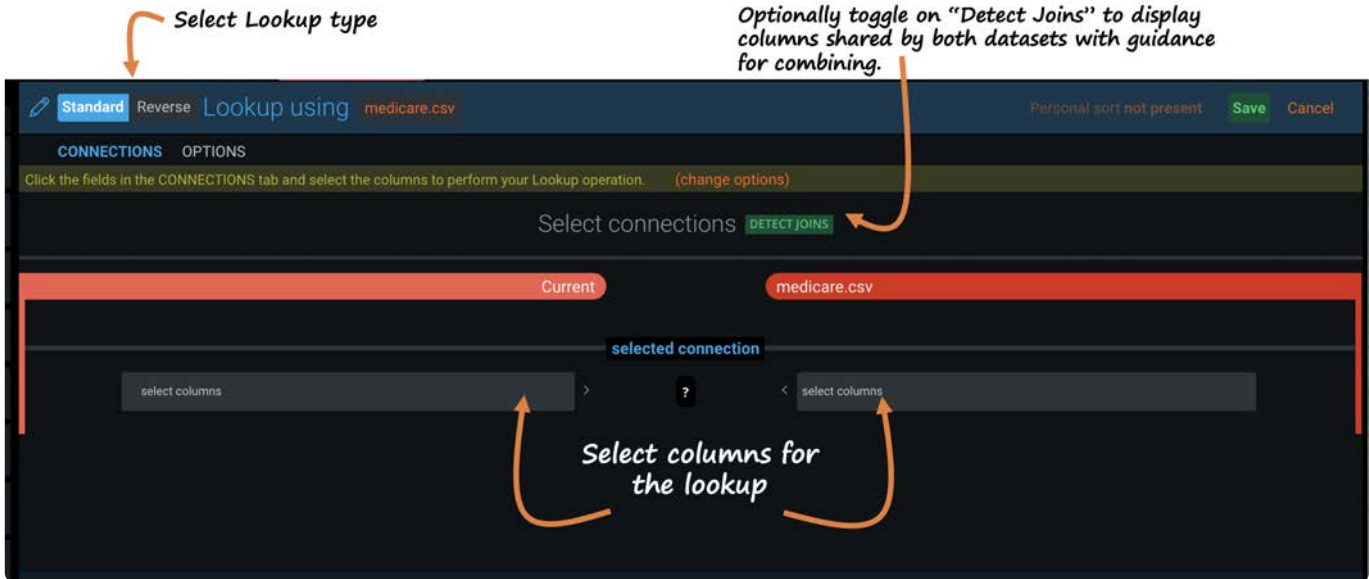
To attach datasets, hover over the **Attach** tool in the project **Tools** bar and select the [lookup](#), [join](#), or [append](#) tools.



Work with the Lookup tool

The **Lookup** tool provides a lookup type operation, which is comparable to an MS Excel VLOOKUP. A reverse lookup is also supported.

After you select the **Lookup** tool, select a **Lookup Source** dataset from your data library. You then select the columns for the lookup operation. The **Current** column refers to your base dataset, and when you click in each column, the available columns from each dataset display.



Note

If you use click to select the green **Detect Joins** option, the **Connections** menu displays the columns that your two datasets share. Additionally, a percentage score provides guidance for how best to combine the datasets. See [Lookup scoring calculations](#) to learn how scores are calculated for **Detect Joins** operations.

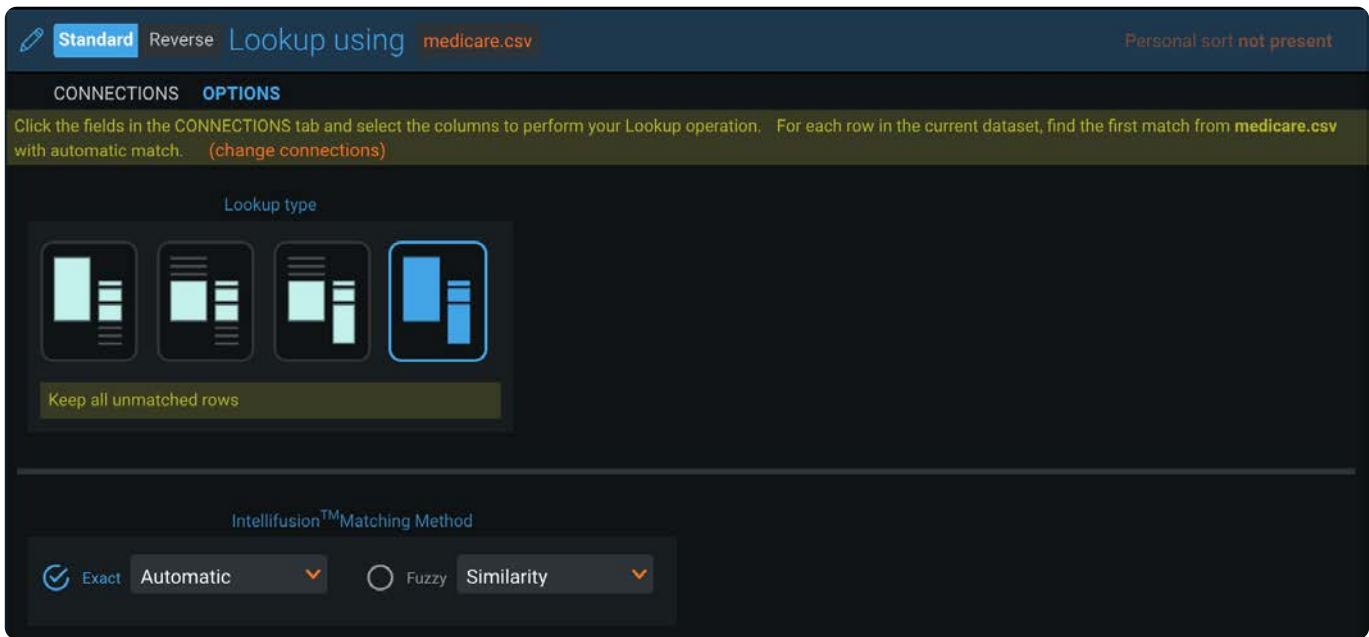
After selecting the columns for the lookup, a preview of the data as it will be combined is displayed on the grid:

1		100010	195	1	25.6	100010--	100010	100010	5/12/13
2					92	100001--	100001	100001	5/12/13
3					123	100005--	100005	100005	5/12/13
4					21	100003--	100003	100003	5/12/13
5					90.9	100007--			
6						100006--	100006	100006	5/12/13
7					21	100005--	100005	100005	5/12/13
8						100006--			
9		100009	176	1	36.5	100009--	100009	100009	5/12/13
10		100007	176	2	4.8	100007--			
11		100000	170	1	2.8	100000--	100000	100000	5/12/13
12		100000	170	1	91.9	100000--	100000	100000	5/12/13
13		100008	167	5	1.8	100008--			
14		100005	166	3	123	100005--	100005	100005	5/12/13
15		100010	165	3	50	100010--	100010	100010	5/12/13
16		100001	164	3		100001--	100001	100001	5/12/13
17		100006	159	5	4.5	100006--			

Both datasets are represented in Sources column to provide a visual indication of where blanks occur.

Next, click the **Options** tab to define:

- The **Lookup Type**, which defines what to do with any unmatched rows.
- The **Matching Method** algorithm to use for the lookup operation. Note: the Fuzzy Matching Method is only available for a Standard type lookup.



When you are satisfied with how your lookup is previewed on the grid, click the green **Save** button to complete the lookup operation.

Lookup scoring calculations

When Data Prep calculates lookup scores for the **Detect Joins** option, it considers two factors:

- **Selectivity:** The extent to which each column in the lookup connection is populated with unique values
- **Overlap:** The percentage of rows that match in the lookup connection

These two factors produce a percentage score to reflect the quality of the connection. The higher the percentage score, the better the connection. However, scores of less than 100% do not necessarily indicate a problem with the proposed combine operation. In fact, there are legitimate use cases where the percentage may be less than 1%.

Rather than being an absolute indicator of “right” or “wrong,” the percentage should serve as a sanity check that should match some level of expectation you have regarding the amount of data the two datasets have in common.

Following are examples of data that produce low and high scores:

Low Scores

BASE (DRIVING) DATASET	LOOKUP (SOURCE) DATASET	EXPLANATION
a,b,c	a,a,b,b,c,c	The lookup has duplicate values.
a,b,c.	. c,d,e,f	The lookup has only one overlap value—'c'.

High Scores

BASE (DRIVING) DATASET	LOOKUP (SOURCE) DATASET	EXPLANATION
a,b,c	a,b,c,d,e,f,g	All lookup values are unique.
a,a,b,b,c,c	a,b,c	All lookup values overlap.

Work with the Join tool

The **Join** tool supports the following join types:

- **Standard Join** that combines all matches across both datasets—comparable to a SQL join.
- **Range Join** that combines datasets based on a marker column that is matched against two columns representing a range in the joining dataset.
- **Cross Product Join** that combines all rows from both datasets. Important: the Cross-Product join significantly increases the number of rows that will be added to your project because every row from the joining dataset is brought into the project for each row in the Base dataset.

Note

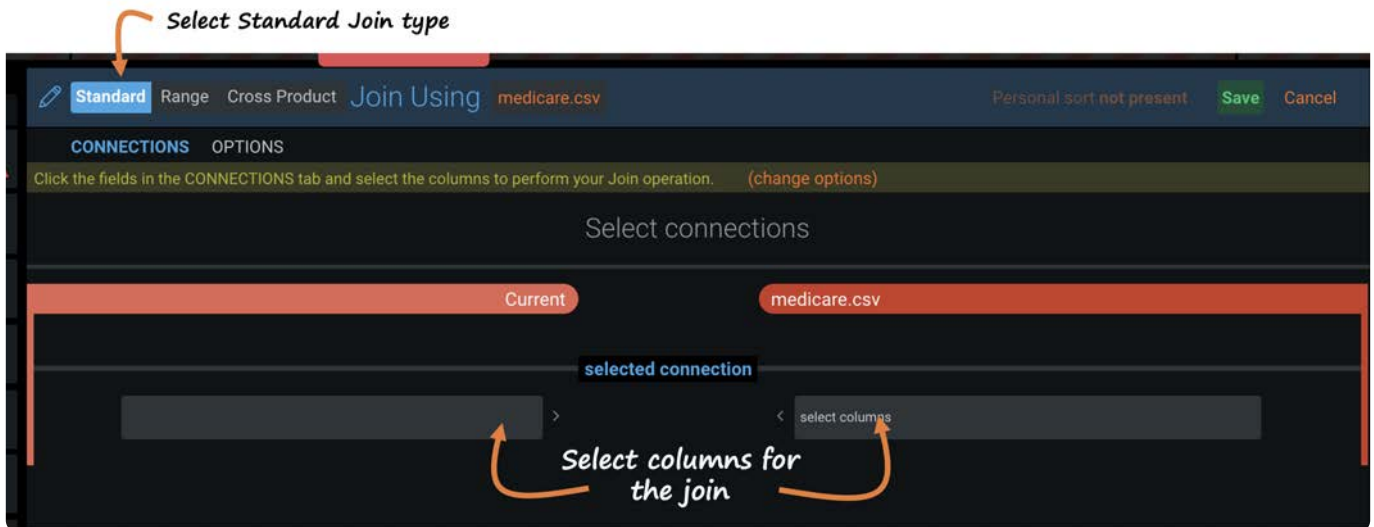
The **Join** tool must be enabled by your System Administrator in order for it to be available in the UI.

Warning

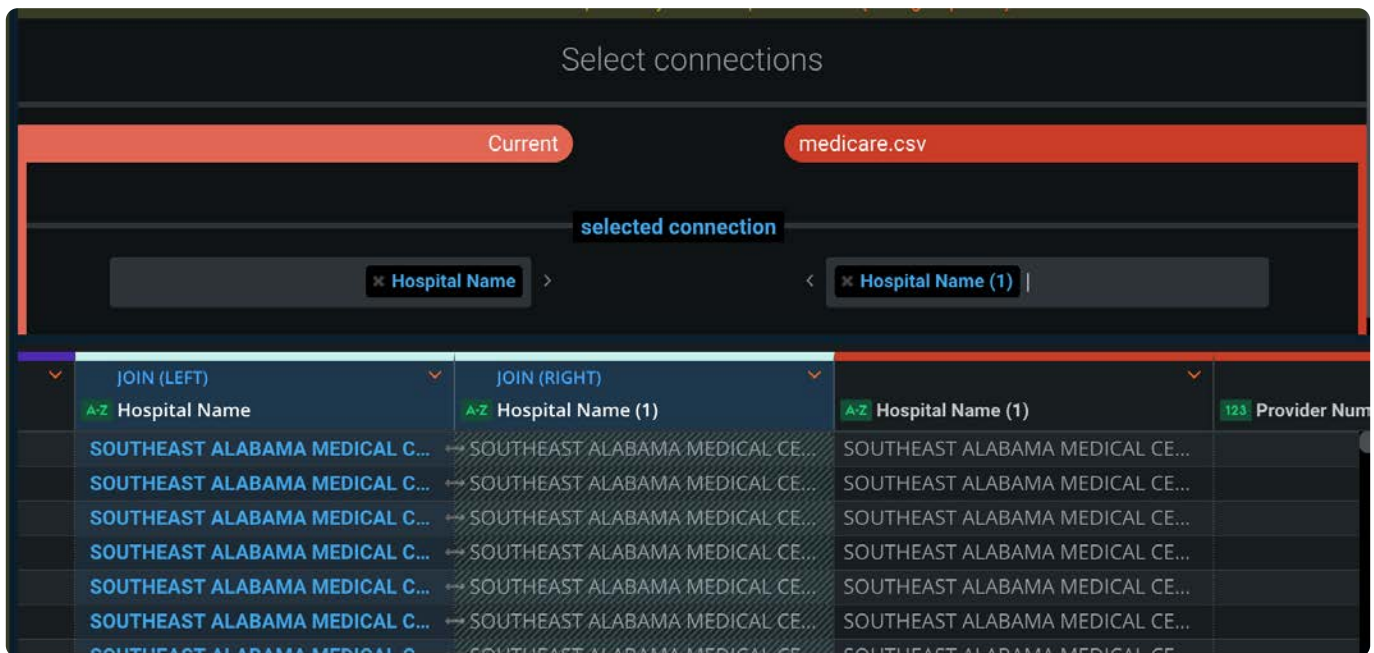
[Sampling a dataset](#) using a "sampling seed" immediately after performing a join may produce non-deterministic results.

Standard Join

After you select the **Join** tool, select the **Standard** type, then select the columns for the join operation. The **Current** column refers to your base dataset. When you click in each column, the available columns from each dataset display. Begin by selecting the columns on which to create the join:

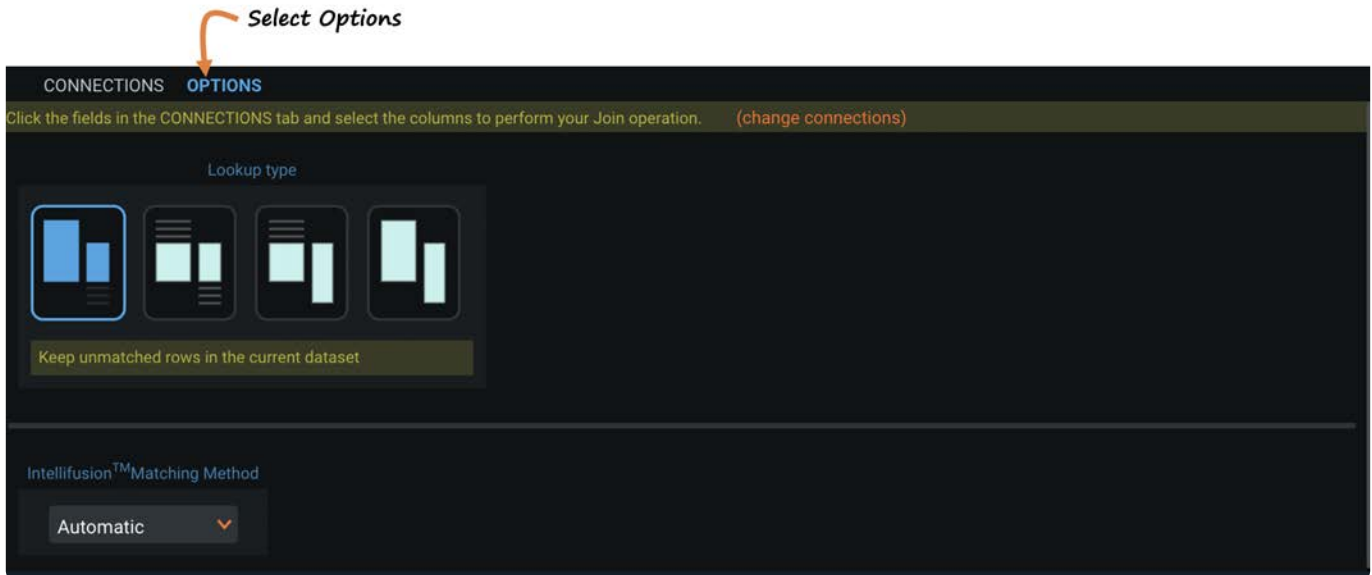


After selecting the columns for the join, a preview of the data as it will be combined is displayed on the grid:



Next, click the **Options** tab to define:

- The **Lookup Type**, which defines what to do with any unmatched rows.
- The **Matching Method** algorithm to use for the join operation.



When you are satisfied with how your lookup is previewed on the grid, click the green **Save** button to complete the lookup operation

Range Join

The Range Join feature allows you to join datasets together based on a "marker" column in your base dataset that is matched against two separate "range" columns in the adjoining dataset.

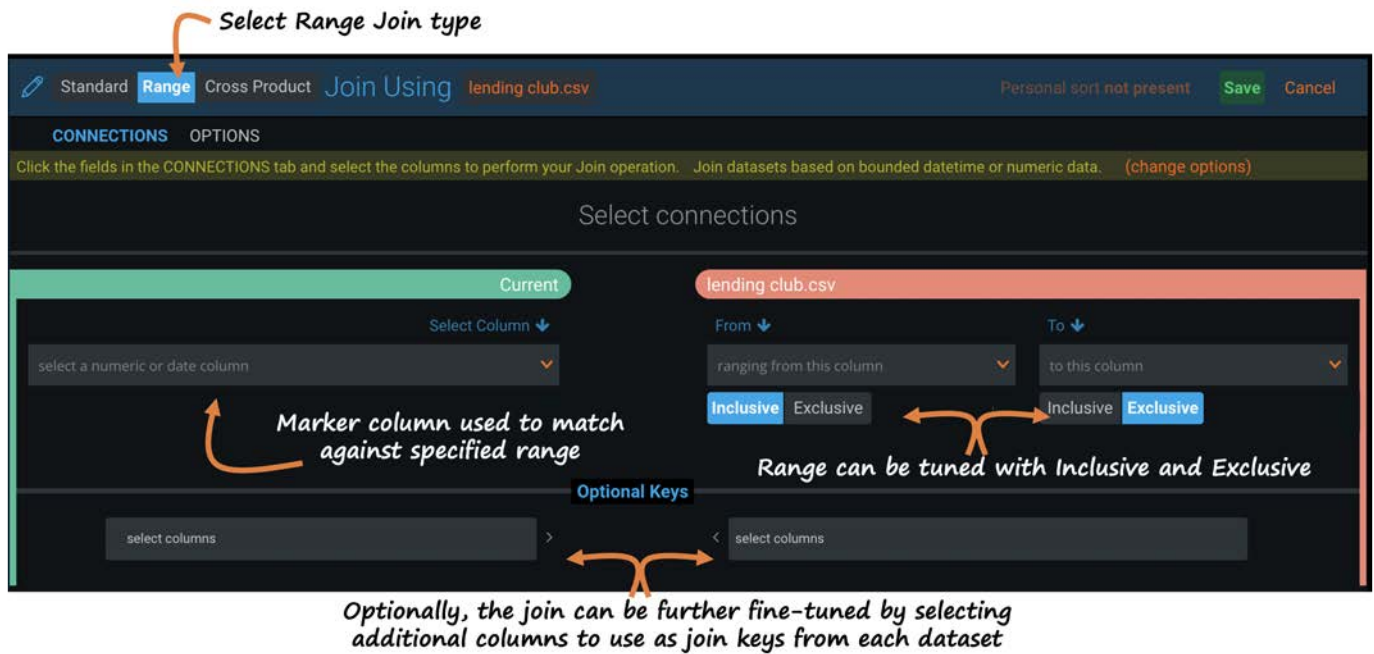
Note

The Range Join feature is only for column types numeric and date. You can always use the **Change into** operation in the column menu to convert columns into these types if the data in those columns can be validly converted.

Note

The Range Join feature must be enabled by your System Administrator in order for it to display as a join option in the UI.

After you select the **Join** tool, select the **Range** option, then select the columns for the join operation. The **Current** column refers to the marker column in your base dataset. It is used to match against the lower and upper range of values in the adjoining dataset:



There are two ways to utilize the Range Join feature:

- **Keyless:** A keyless range join groups together entries that fall within a certain range of the marker value.
- **Keyed:** *Optional Keys* are extra columns you can select from the base and adjoining dataset to more carefully increase the specificity of the join criteria.

After selecting the columns for the join, a preview of the data as it will be combined is displayed on the grid.

Example

Our base dataset is a lending club dataset that contains typical information about individual loans. To that base, we are joining the `loan_riskfactors.csv` dataset that contains averaged risk factor of loans based on grade and date. The goal of our project is to determine the risk factor of each loan in lending club base dataset with the industry averages provided by the risk factor dataset.

Here we are using the column "loan_date" as our marker column in the lending club base dataset and matched it against the "StartDate" and "EndDate" columns from the `loan_riskfactors.csv` dataset. Notice the grid displays a preview of the resulting join:

Standard **Range** Cross Product Join Dataset loan_riskfactors.csv Personal sort not present

CONNECTIONS OPTIONS

Click the fields in the CONNECTIONS tab and select the columns to perform your Join operation. Join datasets based on bounded datetime or numeric data. (change options)

Select connections

Current loan_riskfactors.csv

Select Column ↓ loan_date

From ↓ StartDate

To ↓ EndDate

Inclusive Exclusive Inclusive Exclusive

Optional Keys

select columns > select columns <

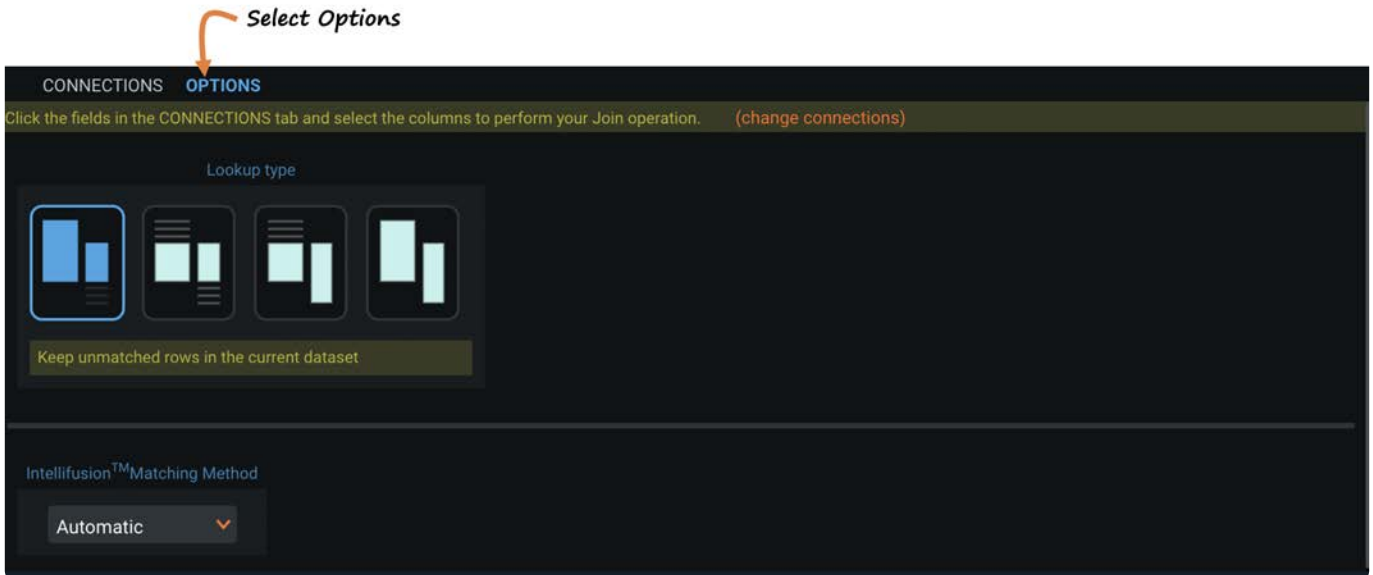
zip_code	addr_state	loan_date	Start Date	End Date	GRADE	Start Date	End Date
3xx	TX	2002-10-01 00:00:00.000Z	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z	A	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z
3xx	TX	2002-10-01 00:00:00.000Z	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z	B	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z
3xx	TX	2002-10-01 00:00:00.000Z	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z	C	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z
3xx	TX	2002-10-01 00:00:00.000Z	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z	D	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z
5xx	NJ	2000-03-01 00:00:00.000Z	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z	A	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z
5xx	NJ	2000-03-01 00:00:00.000Z	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z	B	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z
5xx	NJ	2000-03-01 00:00:00.000Z	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z	C	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z
5xx	NJ	2000-03-01 00:00:00.000Z	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z	D	2000-01-01 00:00:00.000Z	2009-12-31 00:00:00.000Z
3xx	OH	1988-10-01 00:00:00.000Z	1980-01-01 00:00:00.000Z	1989-12-31 00:00:00.000Z	A	1980-01-01 00:00:00.000Z	1989-12-31 00:00:00.000Z
3xx	OH	1988-10-01 00:00:00.000Z	1980-01-01 00:00:00.000Z	1989-12-31 00:00:00.000Z	B	1980-01-01 00:00:00.000Z	1989-12-31 00:00:00.000Z
3xx	OH	1988-10-01 00:00:00.000Z	1950-01-01 00:00:00.000Z	1989-12-31 00:00:00.000Z	C	1950-01-01 00:00:00.000Z	1989-12-31 00:00:00.000Z
3xx	OH	1988-10-01 00:00:00.000Z	1950-01-01 00:00:00.000Z	1999-12-31 00:00:00.000Z	D	1950-01-01 00:00:00.000Z	1999-12-31 00:00:00.000Z
3xx	CT	2000-01-01 00:00:00.000Z					
5xx	CO	1988-10-01 00:00:00.000Z	1980-01-01 00:00:00.000Z	1989-12-31 00:00:00.000Z	A	1980-01-01 00:00:00.000Z	1989-12-31 00:00:00.000Z
5xx	CO	1988-10-01 00:00:00.000Z	1980-01-01 00:00:00.000Z	1989-12-31 00:00:00.000Z	B	1980-01-01 00:00:00.000Z	1989-12-31 00:00:00.000Z
5xx	CO	1988-10-01 00:00:00.000Z	1950-01-01 00:00:00.000Z	1989-12-31 00:00:00.000Z	C	1950-01-01 00:00:00.000Z	1989-12-31 00:00:00.000Z
5xx	CO	1988-10-01 00:00:00.000Z	1950-01-01 00:00:00.000Z	1999-12-31 00:00:00.000Z	D	1950-01-01 00:00:00.000Z	1999-12-31 00:00:00.000Z
2xx	CA	1992-07-01 00:00:00.000Z	1990-01-01 00:00:00.000Z	1999-12-31 00:00:00.000Z	A	1990-01-01 00:00:00.000Z	1999-12-31 00:00:00.000Z
2xx	CA	1992-07-01 00:00:00.000Z	1990-01-01 00:00:00.000Z	1999-12-31 00:00:00.000Z	B	1990-01-01 00:00:00.000Z	1999-12-31 00:00:00.000Z
2xx	CA	1992-07-01 00:00:00.000Z	1990-01-01 00:00:00.000Z	1999-12-31 00:00:00.000Z	C	1990-01-01 00:00:00.000Z	1999-12-31 00:00:00.000Z

Tip

If the resulting dataset contains too many rows after the join, remember that you can use the **Optional Keys** to create a keyed join that increases the specificity of the join criteria, thereby reducing the number of resulting rows.

Next, click the **Options** tab to define:

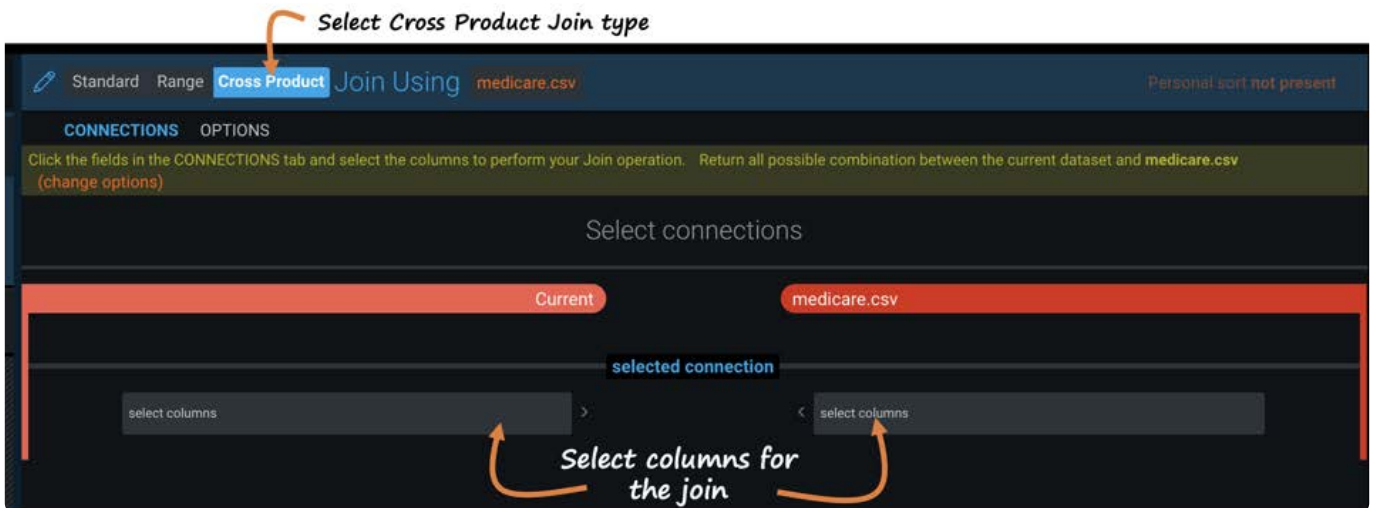
- The **Lookup Type**, which defines what to do with any unmatched rows.
- The **Matching Method** algorithm to use for the join operation.



When you are satisfied with how your lookup is previewed on the grid, click the green **Save** button to complete the lookup operation.

Cross Product Join

After you select the **Join** tool, select the **Cross Product** type, then select the columns for the join operation. The **Current** column refers to your base dataset, and when you click in each column, the available columns from each dataset display. Begin by selecting the columns on which to create the join:



After selecting the columns for the join, a preview of the data as it will be combined is displayed on the grid:

Standard Range **Cross Product** Join Using medicareSHORT.xlsx Personal sort not present

CONNECTIONS OPTIONS

Click the fields in the CONNECTIONS tab and select the columns to perform your Join operation. Return all possible combination between the current dataset and medicareSHORT.xlsx (change options)

Select connections

Current medicareSHORT.xlsx

selected connection

Provider ID Provider Number

A-Z Hospital Ownership	A-Z Emergency Services	LOOKUP (RIGHT) 123 Provider ID	LOOKUP (RIGHT) 123 Provider Number	A-Z Hospital Name (1)
Government - Hospital District or Authority	Yes			SOUTHEAST ALABAMA MEDICAL CE..
Government - Hospital District or Authority	Yes			SOUTHEAST ALABAMA MEDICAL CE..
Government - Hospital District or Authority	Yes			SOUTHEAST ALABAMA MEDICAL CE..
Government - Hospital District or Authority	Yes			SOUTHEAST ALABAMA MEDICAL CE..
Government - Hospital District or Authority	Yes			SOUTHEAST ALABAMA MEDICAL CE..

Note

Unlike the other join types, there are no lookup or matching options for a Cross Product join because all rows will be matched in this operation.

When you are satisfied with how your lookup is previewed on the grid, click the green **Save** button to complete the lookup operation.

Note

The Cross-Product join significantly increases the number of rows that will be added to your project because every row from the joining dataset is brought into the project for each row in the Base dataset. If the number of rows exceeds your project's row limit, then the Join operation will fail with an error message. In this case, reduce the number of rows in your datasets before combining them or speak with your System Administrator regarding the project row limit.

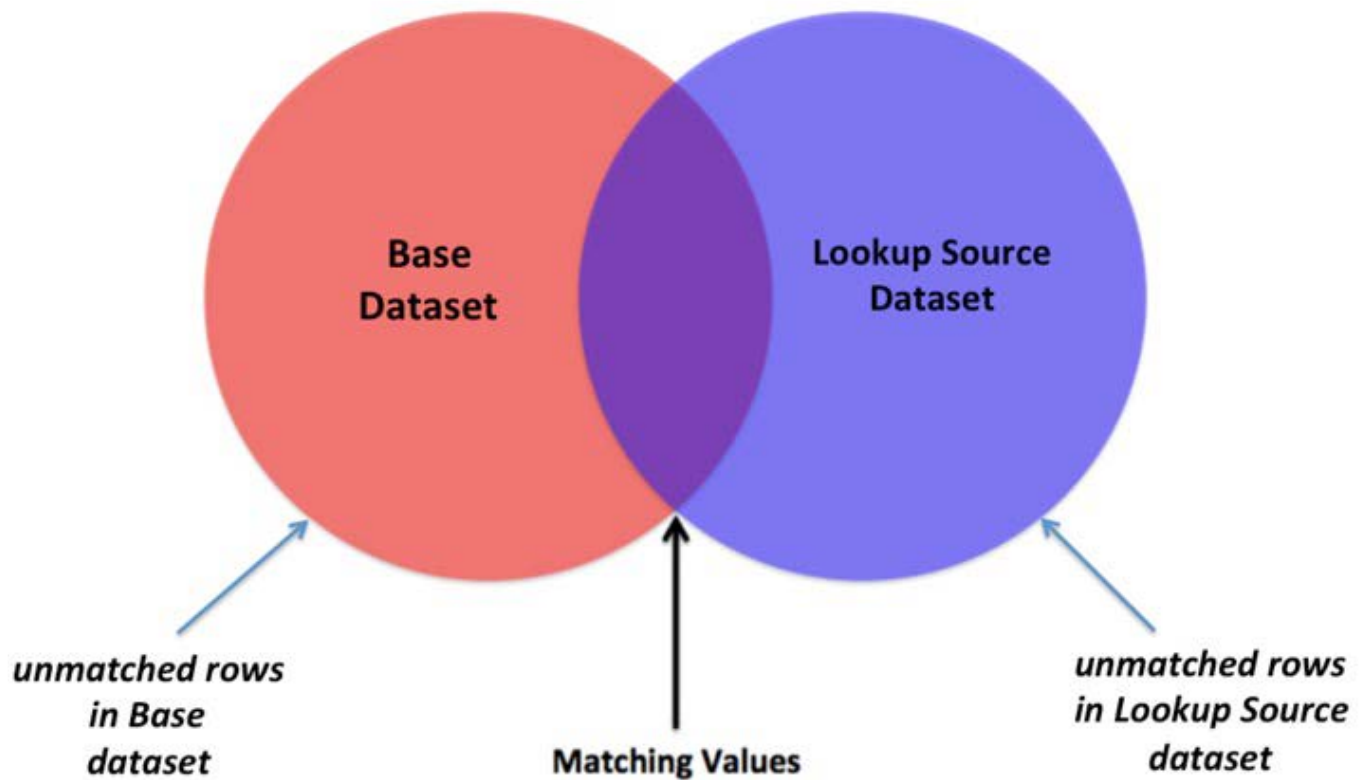
Lookup Type options for Lookups and Joins

The **Lookup Type** options specify what to do with unmatched rows. There are four options: Left Lookup, Inner Lookup, Right Lookup, and Outer Lookup.

Note

When you mouse over the icons for each **Lookup Type** option, its associated type displays.

One of the easiest ways to understand the different Lookup Types is by using a Venn diagram:



Left Lookup

A Left Lookup returns a combined dataset containing all rows with matching values in the specified columns, as well as all rows from the left (Base) dataset that had no corresponding matching value in the right (Lookup) dataset.

AUSTIN	
BOSTON	
WASHINGTON D.C.	
NEW YORK CITY	NEW YORK CITY
LOS ANGELES	LOS ANGELES
MIAMI	MIAMI
DALLAS	DALLAS
SAN FRANCISCO	SAN FRANCISCO

Inner Lookup

An Inner Lookup returns a combined dataset containing only those rows which contain matching values within the specified columns:

NEW YORK CITY	NEW YORK CITY
LOS ANGELES	LOS ANGELES
MIAMI	MIAMI
DALLAS	DALLAS
SAN FRANCISCO	SAN FRANCISCO

Right Lookup

A Right Lookup is similar to a Left Lookup in that it returns a combined dataset containing all rows with matching values in the specified columns; however, rather than return the rows with non-matching values from the left (Base) dataset, it returns all non-matching rows from the right (Lookup) dataset instead.

NEW YORK CITY	NEW YORK CITY
LOS ANGELES	LOS ANGELES
MIAMI	MIAMI
DALLAS	DALLAS
SAN FRANCISCO	SAN FRANCISCO
	SEATTLE
	CHICAGO
	HOUSTON

Outer Lookup

An Outer Lookup returns a combined dataset containing all rows with matching values in the specified columns, as well as all rows from both datasets that contained no matching values:

AUSTIN	
BOSTON	
WASHINGTON D.C.	
NEW YORK CITY	NEW YORK CITY
LOS ANGELES	LOS ANGELES
MIAMI	MIAMI
DALLAS	DALLAS
SAN FRANCISCO	SAN FRANCISCO
	SEATTLE
	CHICAGO
	HOUSTON

Matching methods for Lookups and Joins

Select an Intellifusion™ matching method algorithm to customize how to match text values in columns that have a connection: **Exact** or **Fuzzy**. Each option is explained below.

Exact option

There are three choices for the Exact option: [Automatic Match](#), [Exact Match](#), and [Custom Match](#).

Automatic Match

Automatic Match (the default option) ignores case, word order, and punctuation around words for matching text values. Generally speaking, space characters and punctuation define word boundaries, but there are important exceptions. This matching method ignores punctuation around words, so *.ave.* matches *AVE*. However, Automatic Match detects special kinds of words that have punctuation inside the word. Any punctuation inside a word does not define a word boundary, and the punctuation must match exactly in both datasets. **Automatic Match** detects the following words with inner punctuation:

- **Numbers with periods:** Currency and floating point numbers such as 12.34. Currency symbols and punctuation around the number are ignored. Comparison is text comparison not numeric, so 3.0 and 3 do not match. Tip: To reduce issues, use the "Parse cell text into numbers" option during import or convert a column to numbers using the column drop-down menu.
- **Email addresses:** Email addresses are one word including internal periods.

- **Acronyms with periods:** Acronyms such as U.S.A. include inner punctuation and counts as one word. However, U.S.A. does not match USA because punctuation inside a word is not ignored when comparing words.

Example of Automatic Match

The following table shows the behavior of Automatic Match. The first two columns show example values. The third column shows if those values match with Automatic Match. For examples where Exact Match and Automatic Match return different results, the answers appears in bold.

Val1	Val2	Automatic Match
Mary T	Mary T	yes
Jan16	Jan16	yes
5,6,7	5,6,7	yes
Mary T	Mary M	no
U.S.	US	no
Mary T	Mary	no
Mary T	Marry T	no
1,2,3	1 2 3	no
1.234	123.4	no
5 6 7	6 7 5	yes
Doe, J	J DOE	yes
A, B	A B	yes
A, B, C	C B A	yes
Main St	MAIN ST.	yes

Exact Match

Use the **Exact Match** option if every character in the two values must exactly match.

Custom Match

Use the **Custom Match** option for making precise adjustments. This option allows you to selectively choose how to handle: word order, case, whitespace, and specific punctuation values. This matching method ignores or keeps word order and case depending on your selection. These options are commonly used for data containing names. Any combination of Ignore and Keep for word order and case can be selected. Click the **Ignore** and **Keep** buttons to see an example of how word order and case is handled based on your selections.

For Whitespace, specify whether to ignore, keep or split on when whitespace is encountered. The **Split On** button splits your data into separate terms on whitespace. Whitespace includes characters from the space bar, tab key and carriage returns. The Whitespace options are commonly used to increase matches for data containing address information. Click the **Ignore**, **Keep** and **Split On** buttons to see an example of how a whitespace is handled based on your selection.

Punctuation options allow you to ignore, keep or split on for specific punctuation. The default punctuation values are: comma and hyphen. Click the white **Other** button to add to the default values. A new field appears for adding more punctuation. To delete the defaults or any punctuation you have added, move your mouse over a punctuation field and click the orange "X" that appears in the top of that field. If you delete all of the punctuation fields, the default is to keep all punctuation.

Settings that you specify for word order, case, whitespace, and punctuation apply only to the current Lookup data prep step. These settings do not apply to other Lookup steps.

Fuzzy option

The **Fuzzy** option is only available for a **Standard Lookup** and must be enabled by your System Administrator in order for it to display as an option in the UI.

The **Fuzzy** option employs an edit distance algorithm to predict the probability of likely matches between two join keys that you select.

Example of a Fuzzy option

In this example, you have a base dataset with company names and you want to enrich that data with addresses for each company.

You have a second dataset, with address information, but the company names do not match closely enough to your base dataset.

In this case, you can use the **Fuzzy** option to quickly create a join for the company names without having to do any cleanup ahead of time on the "company name" column in either dataset:

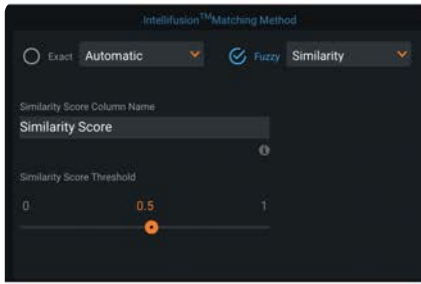
Your base dataset has company names:

	Sources	A-Z Company
1		Accel partners
2		Bussiness Corp of Americas
3		Camden consultants
4		Everlast Incorporated
5		Madrid Group
6		Next wave
7		PPL INC
8		Outland group of America
9		Kimshi LTD
10		Young, Phillips & Marks

Your enrichment dataset with addresses has company names listed this way:

	A-Z Company Name	A-Z Address	A-Z City	A-Z State
1	Accel partners	123 New Avenue	West Hollyw...	CA
2	Bussiness Corp of America	622 Tripp Road	Woodside	CA
3	Camden consultant	42 Short North R...	Clearwater	FL
4	Everlast INC	345 West Hocking	New Orleans	LA
5	Madrid group LLC	11 Marsh Road	Middleton	VA
6	Nextwave	589 Covington Ro...	Newberry	MA
7	PPL incorporated	16 High Street	Paris	TX
8	Outland group of America	888 Scott Road	Burlington	VA
9	Kimshi Limited	16 W Town Street	Sun Valley	ID
10	Young, Phillips and Marks	555 Middle Way	New York	NY

You tune the **Fuzzy** options as follows:



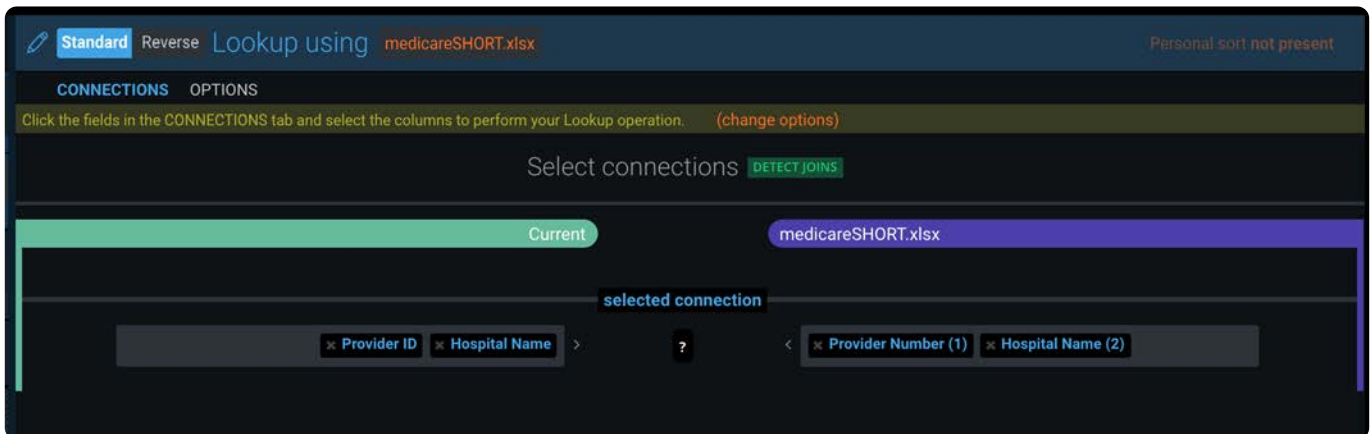
The Similarity Score column is automatically created for the Fuzzy lookup option and provides a measure of how close the join keys match—with the value of “1” being a perfect match. You can change the column’s name by entering a new name in the field.

The threshold slider allows you to change the threshold probability for matches. The data grid updates as you move the slider so that you can see how the Similarity Score is affected.

	Sources	Company	LOOKUP (LEFT)	LOOKUP (RIGHT)	Similarity Score	Company Name	Address	City	State
1		Accel partners	Accel partners	Accel partners	1	Accel partners	123 New Avenue	West Hollywood	CA
2		Bussiness Corp of Americas	Bussiness Corp of Americas	Bussiness Corp of America	0.9995068907737732	Bussiness Corp of America	622 Tripp Road	Woodside	CA
3		Camden consultants	Camden consultants	Camden consultant	0.9989712238311768	Camden consultant	42 Short North Road	Clearwater	FL
4		Everlast Incorporated	Everlast Incorporated	Everlast INC	0.9387755392981934	Everlast INC	345 West Hocking	New Orleans	LA
5		Madrid Group	Madrid Group	Madrid group LLC	0.9791666865348817	Madrid group LLC	11 Marsh Road	Middleton	VA
6		Next wave	Next wave	Nextwave	0.977777791023254	Nextwave	589 Covington Road	Newberry	MA
7		PPL INC	PPL INC	PPL incorporated	0.89453125	PPL incorporated	16 High Street	Paris	TX
8		Outland group of America	Outland group of America	Outland group of America	1	Outland group of America	888 Scott Road	Burlington	VA
9		Kimshi LTD	Kimshi LTD	Kimshi Limited	0.9591936929321289	Kimshi Limited	16 W Town Street	Sun Valley	ID
10		Young, Phillips & Marks	Young, Phillips & Marks	Young, Phillips and Marks	0.974928081035614	Young, Phillips and Marks	555 Middle Way	New York	NY

After you find the best threshold match, click **Save** to join the data.

Note that you can simultaneously create multiple join keys for the fuzzy option by adding more column names in the Connections panel:



However, both of your datasets must contain the exact same number of columns. Otherwise, you will receive an error and cannot proceed.

After the data is joined, you can apply additional options:

- Use the **Column tool** to remove the Similarity Score column and one of the Company name columns.
- Use a **Filtergram** on the remaining Company column to ensure the name formats are all correct. If you find names that you want to edit, use the **Find and Replace** option to correctly update those names.

Examples of Lookups and Joins

Example: Reverse Lookup

A Reverse Lookup operation is useful when you want to achieve a one-to-many relationship between your current dataset and an enrichment dataset. This Reverse operation uses your current base dataset as the Lookup table—meaning only the first match from the base dataset is brought into the selected enrichment dataset.

In this example, you have two datasets:

- An **Orders** dataset with unique **Order IDs**

	123 ORDER_ID	123 ORDER_DATE	123 CUSTOMER_ID	123 NUM_ITEMS	123 TOTAL_AMT
1	100000	5/12/13	10876	8	90.9
2	100001	5/12/13	10017	4	59.4
3	100002	5/12/13	18388	1	4.8
4	100003	5/12/13	11137	2	21
5	100004	5/12/13	13079	1	1.8
6	100005	5/12/13	14005	3	36.5
7	100006	5/12/13	11037	2	50
8	100007	5/12/13	12986	4	123
9	100008	5/12/13	12571	2	25.6
10	100009	5/12/13	15153	1	4.5
11	100010	5/12/13	12196	3	16.5

- An **Order Details** dataset that has all of the ordering details associated with an **Order ID**

	123 ORDER_ID	123 LINEITEM_ID	123 PRODUCT_ID	123 QUANTITY	123 LIST_PRICE	123 DISCOUNT	123 EXT_PRICE	123 SHIP_DATE
1	100000	1	170	1	5	30	3.5	2013-05-13T00:00:00.000Z
2	100000	2	126	1	6	0	6	2013-05-13T00:00:00.000Z
3	100000	3	140	3	5	20	12	2013-05-14T00:00:00.000Z
4	100000	4	182	1	8	20	6.4	2013-05-14T00:00:00.000Z
5	100000	5	170	1	5	0	5	2013-05-13T00:00:00.000Z
6	100000	6	182	2	8	0	16	2013-05-16T00:00:00.000Z
7	100000	7	155	1	2	0	2	2013-05-15T00:00:00.000Z
8	100000	8	193	4	10	0	40	2013-05-13T00:00:00.000Z
9	100001	1	164	3	10	0	30	2013-05-13T00:00:00.000Z
10	100001	2	143	3	6	0	18	2013-05-17T00:00:00.000Z
11	100001	3	138	4	1	40	2.4	2013-05-14T00:00:00.000Z
12	100001	4	193	1	10	10	9	2013-05-17T00:00:00.000Z
13	100002	1	126	1	6	20	4.8	2013-05-13T00:00:00.000Z
14	100003	1	121	3	1	0	3	2013-05-13T00:00:00.000Z
15	100003	2	188	1	20	10	18	2013-05-13T00:00:00.000Z
16	100004	1	139	1	2	10	1.8	2013-05-14T00:00:00.000Z
17	100005	1	128	1	5	10	4.5	2013-05-13T00:00:00.000Z
18	100005	2	152	5	1	0	5	2013-05-15T00:00:00.000Z
19	100005	3	166	3	10	10	27	2013-05-15T00:00:00.000Z

You create a new project with the **Orders** dataset as your Base dataset and you want to enrich it with data from **Order Details** based on the "Order ID" column in both datasets. In this case, you select the **Lookup Reverse** option, and then choose "Order_ID" as the column for the matching the datasets:

Individual Orders by Order ID

ORDER_ID	ORDER_DATE	CUSTOMER_ID	NUM_ITEMS	TOTAL_AMT
100000	5/12/13	10876	8	90.9
100001	5/12/13	10017	4	59.4
100002	5/12/13	18388	1	4.8
100003	5/12/13	11137	2	21
100004	5/12/13	13079	1	1.8
100005	5/12/13	14005	3	36.5



Order Details by Order ID

ORDER_ID	QUANTITY	PRODUCT_ID	SHIP_DATE
100000	1	170	2013-05-13 10:00:00...
100000	1	126	2013-05-13 10:00:00...
100000	3	140	2013-05-14 10:00:00...
100000	1	182	2013-05-14 10:00:00...
100000	1	170	2013-05-13 10:00:00...
100000	2	182	2013-05-16 10:00:00...
100000	1	155	2013-05-15 10:00:00...
100000	4	193	2013-05-13 10:00:00...
100001	3	164	2013-05-13 10:00:00...
100001	3	143	2013-05-17 10:00:00...
100001	4	138	2013-05-14 10:00:00...
100001	1	193	2013-05-17 10:00:00...
100002	1	126	2013-05-13 10:00:00...
100003	3	121	2013-05-13 10:00:00...
100003	1	188	2013-05-13 10:00:00...
100004	1	139	2013-05-14 10:00:00...
100005	1	128	2013-05-13 10:00:00...
100005	5	152	2013-05-15 10:00:00...
100005	3	166	2013-05-15 10:00:00...

Note

To keep any unmatched (blank) rows containing the "Order_ID", choose an Outer or Left Lookup type. Alternatively, you can discard those rows by selecting an Inner Lookup type. Refer to Lookup Type, which defines what to do with any unmatched rows.

Important

A **Reverse Lookup** affects the number of rows that will be added to your base dataset and, possibly, the sort order.

Example: Join datasets

This selection combines all matches across both datasets. The Join operation is useful when you want to achieve a many-to-many relationship between two datasets. Note: this selection affects the number of rows that will be added to your project.

In this example, you have two datasets:

- **Transactions** with a set of purchase transaction IDs, each row indicating a book that was purchased. Note that there may be multiple transactions IDs with the same book because that book may have sold multiple times.

	Sources	123 Transaction ID	A-Z Book
1		1001	The Elements of Style
2		1002	All the President's Men
3		1003	The C Programming Language
4		1004	The Communist Manifesto
5		1005	Harry Potter and the Sorcerer's Stone
6		1006	The Elements of Style
7		1007	Capitalism and Freedom
8		1008	Harry Potter and the Sorcerer's Stone
9		1009	The Elements of Style
10		1010	Capitalism and Freedom
11		1011	Harry Potter and the Sorcerer's Stone
12		1012	The C Programming Language
13		1013	Harry Potter and the Sorcerer's Stone
14		1014	Capitalism and Freedom

- **Books** with a set of books and their authors. Note that each book may have multiple authors, such that one book may be listed in multiple rows, once for each co-author.

	Sources	A-Z Book	A-Z Author
1		The Elements of Style	Strunk
2		The Elements of Style	White
3		The C Programming Language	Kernighan
4		The C Programming Language	Ritchie
5		The Communist Manifesto	Marx
6		The Communist Manifesto	Engels
7		Capitalism and Freedom	Friedman
8		All the President's Men	Woodward
9		All the President's Men	Bernstein
10		Harry Potter and the Sorcerer's Stone	Rowling

Your objective is to determine the number of transactions per author. And because there are multiple authors for many of the books—instead of a single, unique value for author—the "Join" is the option you want to select for combining these datasets.

Create a new project with the **Transactions** dataset. Then do a standard Join operation with the **Books** dataset, and choose "Book" as the column for joining the datasets:

	Sources	BOOK-TRANSACTIONS	BOOK-TRANSACTIONS	BOOKS	BOOKS
		Transaction ID	Book	Book (1)	Author
1		1003	The C Programming Language	The C Programming Language	Kernighan
2		1003	The C Programming Language	The C Programming Language	Ritchie
3		1012	The C Programming Language	The C Programming Language	Kernighan
4		1012	The C Programming Language	The C Programming Language	Ritchie
5		1005	Harry Potter and the Sorcerer's Stone	Harry Potter and the Sorcerer's Stone	Rowling
6		1008	Harry Potter and the Sorcerer's Stone	Harry Potter and the Sorcerer's Stone	Rowling
7		1011	Harry Potter and the Sorcerer's Stone	Harry Potter and the Sorcerer's Stone	Rowling
8		1013	Harry Potter and the Sorcerer's Stone	Harry Potter and the Sorcerer's Stone	Rowling
9		1004	The Communist Manifesto	The Communist Manifesto	Marx
10		1004	The Communist Manifesto	The Communist Manifesto	Engels
11		1002	All the President's Men	All the President's Men	Woodward
12		1002	All the President's Men	All the President's Men	Bernstein
13		1001	The Elements of Style	The Elements of Style	Strunk
14		1001	The Elements of Style	The Elements of Style	White
15		1006	The Elements of Style	The Elements of Style	Strunk
16		1006	The Elements of Style	The Elements of Style	White
17		1009	The Elements of Style	The Elements of Style	Strunk
18		1009	The Elements of Style	The Elements of Style	White
19		1007	Capitalism and Freedom	Capitalism and Freedom	Friedman
20		1010	Capitalism and Freedom	Capitalism and Freedom	Friedman
21		1014	Capitalism and Freedom	Capitalism and Freedom	Friedman

Note

Performing a Join can substantially increase the number of rows in the resulting dataset. If the number of rows exceeds your project's row limit, then the Join operation will fail with an error message. In this case, reduce the number of rows in your datasets before joining them or speak with your System Administrator regarding the project row limit.

Example: Cross Product Join to return all combinations

This selection combines all rows from both of the datasets in the operation.

In this example, you have three datasets:

- **Customer Master** with all of your customers IDs.

	A-Z	CUSTOMER_ID	
1		Customer_Alfa	🔒
2		Customer_Kilo	🔒
3		Customer_Oscar	🔒
4		Customer_Echo	🔒
5		Customer_Delta	🔒
6		Customer_Foxtrot	🔒
7		Customer_Tango	🔒
8		Customer_Bravo	🔒
9		Customer_Hotel	🔒

- **Products Master** with all of your products and associated IDs.

	123	PRODUCT_ID	A-Z	Product	
1		120	🔒	ADDE chair	🔒
2		126	🔒	AINA fabric	🔒
3		140	🔒	FADO table la...	🔒
4		150	🔒	FEJS wall clock	🔒
5		155	🔒	HENSVIK cabi...	🔒
6		160	🔒	JANINGE armc...	🔒
7		170	🔒	KLYSA wall clo...	🔒
8		180	🔒	MAJJE throw	🔒
9		182	🔒	MUSTIG glass	🔒

- **Customer Orders** with the Customer IDs, and the associated Products and Quantity that were purchased.

	A-Z CUSTOMER_ID	123 prod ID	123 qty
1	Customer_Alfa	120	1
2	Customer_Kilo	120	1
3	Customer_Oscar	140	3
4	Customer_Echo	140	1
5	Customer_Delta	150	1
6	Customer_Foxtrot	150	2
7	Customer_Tango	160	1
8	Customer_Bravo	170	4
9	Customer_Hotel	180	3

Your objective is to create a master project that includes all customers and all products, and then determine all of the products that each customer has *not* purchased.

Create a new project with the **Customer Master** as your base dataset. Then do a **Cross-Product** join operation with the **Products Master** dataset. Choose "Customer_ID" and "Product_ID" as the columns for the cross-product operation.

You now have a master dataset that has all of your customers and all of your products:

	A-Z CUSTOMER_ID	123 prod ID	123 qty
1	Customer_Alfa	120	1
2	Customer_Kilo	120	1
3	Customer_Oscar	140	3
4	Customer_Echo	140	1
5	Customer_Delta	150	1
6	Customer_Foxtrot	150	2
7	Customer_Tango	160	1
8	Customer_Bravo	170	4
9	Customer_Hotel	180	3

Next, do a **Standard Lookup** operation with the **Customer Orders** dataset. Choose "Customer_ID" and "Product_ID" as the columns for the lookup, and be sure to keep all unmatched rows. Refer to **Lookup Type**, which defines what to do with any unmatched rows.

The resulting dataset enables you to easily identify all of the products that each customer has not yet purchased:

	Sources	CUSTOMER_MASTER A-Z CUSTOMER_ID	PRODUCTS_MASTER 123 PRODUCT_ID	PRODUCTS_MASTER A-Z Product
1		Customer_Alfa	120	ADDE chair
2		Customer_Alfa	126	AINA fabric
3		Customer_Alfa	140	FADO table lamp
4		Customer_Alfa	150	FEJS wall clock
5		Customer_Alfa	155	HENSVIK cabinet
6		Customer_Alfa	160	JANINGE armchair
7		Customer_Alfa	170	KLYSA wall clock
8		Customer_Alfa	180	MAJJE throw
9		Customer_Alfa	182	MUSTIG glass
10		Customer_Kilo	120	ADDE chair
11		Customer_Kilo	126	AINA fabric
12		Customer_Kilo	140	FADO table lamp
13		Customer_Kilo	150	FEJS wall clock
14		Customer_Kilo	155	HENSVIK cabinet
15		Customer_Kilo	160	JANINGE armchair
16		Customer_Kilo	170	KLYSA wall clock
17		Customer_Kilo	180	MAJJE throw
18		Customer_Kilo	182	MUSTIG glass

Each time you update the **Customer Orders** dataset for this project, you automatically see the data on the grid updated to reflect the new purchases. To capture all of the purchase information over time, you can create a [project lens](#) to publish an AnswerSet for each result.

Note

The **Cross-Product** operation significantly increases the number of rows that will be added to your project because every row from the Lookup source is brought into the project for each row in the Base dataset. If the number of rows exceeds your project's row limit, then the Join operation will fail with an error message. In this case, reduce the number of rows in your datasets before combining them or speak with your System Administrator regarding the project row limit.

Work with the Append tool

Using the Append tool, you can select an additional dataset that contains rows that you want to add to the end of your base dataset. You can customize column matching between the two datasets. If all columns from the appended dataset are matched to columns in the current dataset, the set of columns in your original dataset are unchanged in the result dataset. If you leave columns from the appended dataset unmatched, those columns become new columns in the result dataset.

Update columns

As you prepare your data in Data Prep, you will find times when you need to make changes to your columns. You can use the **Columns** tool to edit the column names, their order, and their availability in a project.

Note

The **Columns** tool lets you work with whole rows but if you need to work with column *data*, use the column menus, as well as the **Filter** and **Column** panes. See [Work with column data](#) for details.

Work with the Columns tool

The **Columns** tool serves multiple purposes. It displays the columns currently in your project and each column's source and type—String, Number, DateTime. The **Columns** tool also gives you the ability to:

- Rename columns.
- Reorder columns.
- Remove columns.

To access the **Columns** tool, click **columns** in the **Tools** bar:

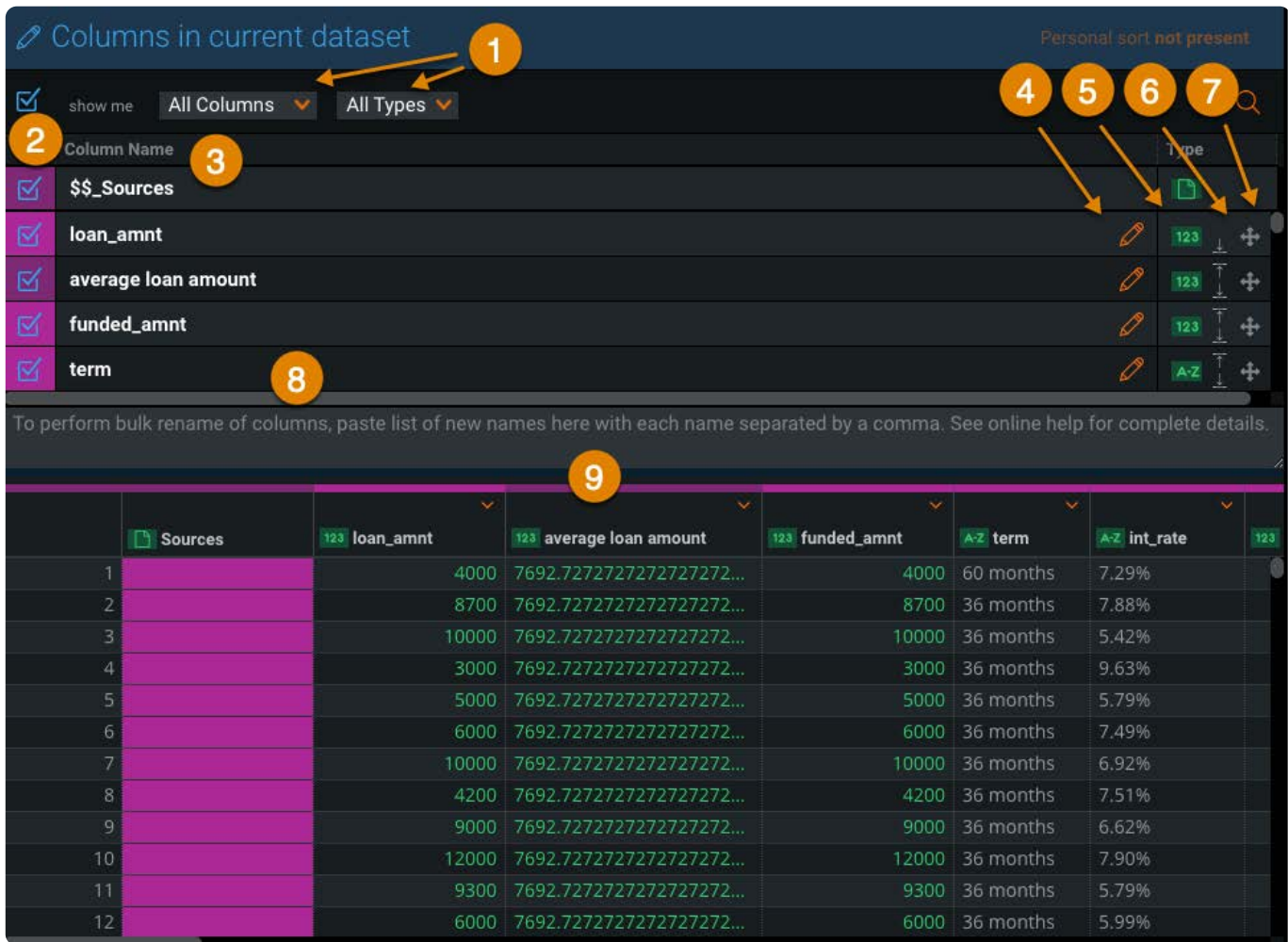
Projects Hospital Readmissions DataRobot

Filters on the Current dataset Personal-sort not present

To add a filter, click on the type icon (A-Z, 123 ...) in a column header, or use the drop-down menu.

	Sources	123 patient_nbr	123 encounter_id	A-Z Race	A-Z Gender	A-Z Age	A-Z Age_bucket
1		77586282	42570	Caucasian	Male	[80-90)	Senior
2		69422211	148530	Caucasian	Male	[70-80)	Senior
3		62718876	216156	Caucasian	Male	[50-60)	Adult
4		115196778	248916	Caucasian	Male	[70-80)	Senior
5		3327282	293118	AfricanAmerican	Male	[80-90)	Senior
6		98427861	325866	Hispanic	Male	[60-70)	Senior
7		112002975	326028	AfricanAmerican	Male	[70-80)	Senior
8		80588529	383430	Caucasian	Male	[70-80)	Senior
9		96435585	421194	AfricanAmerican	Male	[30-40)	Adult
10		66274866	449142	Caucasian	Male	[60-70)	Senior
11		106936875	464994	AfricanAmerican	Male	[70-80)	Senior
12		37746639	590346	Caucasian	Male	[70-80)	Senior
13		23043240	1070256	Caucasian	Male	[60-70)	Senior
14		54746082	1185942	Caucasian	Male	[60-70)	Senior
15		92117574	1260216	Caucasian	Male	[80-90)	Senior
16		91530936	1260894	Caucasian	Male	[70-80)	Senior
17		50253120	1262736	Caucasian	Male	[50-60)	Adult
18		48925980	1414158	Caucasian	Male	[70-80)	Senior
19		49407813	1802280	Caucasian	Male	[50-60)	Adult
20		10430154	1880598	AfricanAmerican	Male	[60-70)	Senior
21		15856002	2087382	Caucasian	Male	[70-80)	Senior
22		5041602	2092362	Caucasian	Male	[60-70)	Senior
23		6500556	2092848	Caucasian	Male	[60-70)	Senior

The Columns in current dataset pane displays:



Following is an overview of the elements you will work with when you edit the columns in your project.

ELEMENT	DESCRIPTION
1	<p>Column filters</p> <p>Filter the Columns in current dataset list by:</p> <ul style="list-style-type: none"> • Selected columns • Renamed columns • Data type
2	<p>Column selectors</p> <p>Clear the selector for each column you want to remove. Check the selector to retain the column. Hiding a column removes it from the AnswerSet when you publish it.</p>
3	<p>Columns in current dataset list</p> <p>View your columns and the type of data they contain. The columns are listed in the order they appear in your data.</p>
4	<p>edit column name</p> <p>Update the name of a column.</p>
5	<p>Type</p> <p>View the data type of the column.</p>

	ELEMENT	DESCRIPTION
6	Move to top / Move to bottom	Move the column to the beginning or end of your dataset.
7	Move	Drag the column to a new location.
8	bulk rename	Rename all columns using a single, comma-separated string.
9	Data Preview pane	View the data in your project. You will see your data change as you prep it.

Rename columns

Rename an individual column

To change the name of a column:

1. From the **Tools** bar, click **columns**.

The **Columns in current dataset** pane appears.

2. Click the name of the column you want to rename or click the pencil icon.
3. Type the new name for the column and click **Enter**.

The **Old Column Name** section appears and displays the column's original name. The **Data Preview** pane displays the updated column name.

If you change your mind about the new name, click **Reset** to reset the column name to the original name.

4. Click **Save** on the top left.

Your change is saved as a step in your project. The column is updated in the **Data Preview** pane.

Rename columns in bulk

The bulk renaming feature allows you to rename all columns at once.

To rename a list of columns, follow this step:

- In the bulk rename field below the **Columns in current dataset** list, type the new column names separated by commas.


The column names in the list update accordingly. You can also paste in new column names from header files separated by commas to quickly rename all of the columns in your dataset.

Reorder columns

To change the location of a column:

1. From the **Tools** bar, click **columns**.

The **Columns in current dataset** pane appears.

2. In the **Type** section, position your pointer over the Move  icon for the column you want to move and drag the column to a new location.

The **Data Preview** pane displays the column in its new position.

You can instead use the **Move to top**  or **Move to bottom**  icons to move the column to the first or last position.

3. Click **Save**.

Your change is saved as a step in your project. The column is updated in the **Data Preview** pane.

Remove columns

The following procedure shows how to remove a column from a project.


Warning

After you remove a column, it's no longer available for use in the project. You won't be able to use the column for subsequent steps and errors will occur in any subsequent steps that rely on a column you removed. You can make a removed column available again by returning to the original step where it was removed and selecting it again to include in your data. You can hide a column rather than remove it using the [hide column operation](#).

To remove a column:

1. From the **Tools** bar, click **columns**.

The **Columns in current dataset** pane appears.

2. Clear the column selector  on the left for each column you want to remove.

The column is shaded in the **Columns in current dataset** list and is removed from the **Data Preview** pane below.

3. Click **Save**.

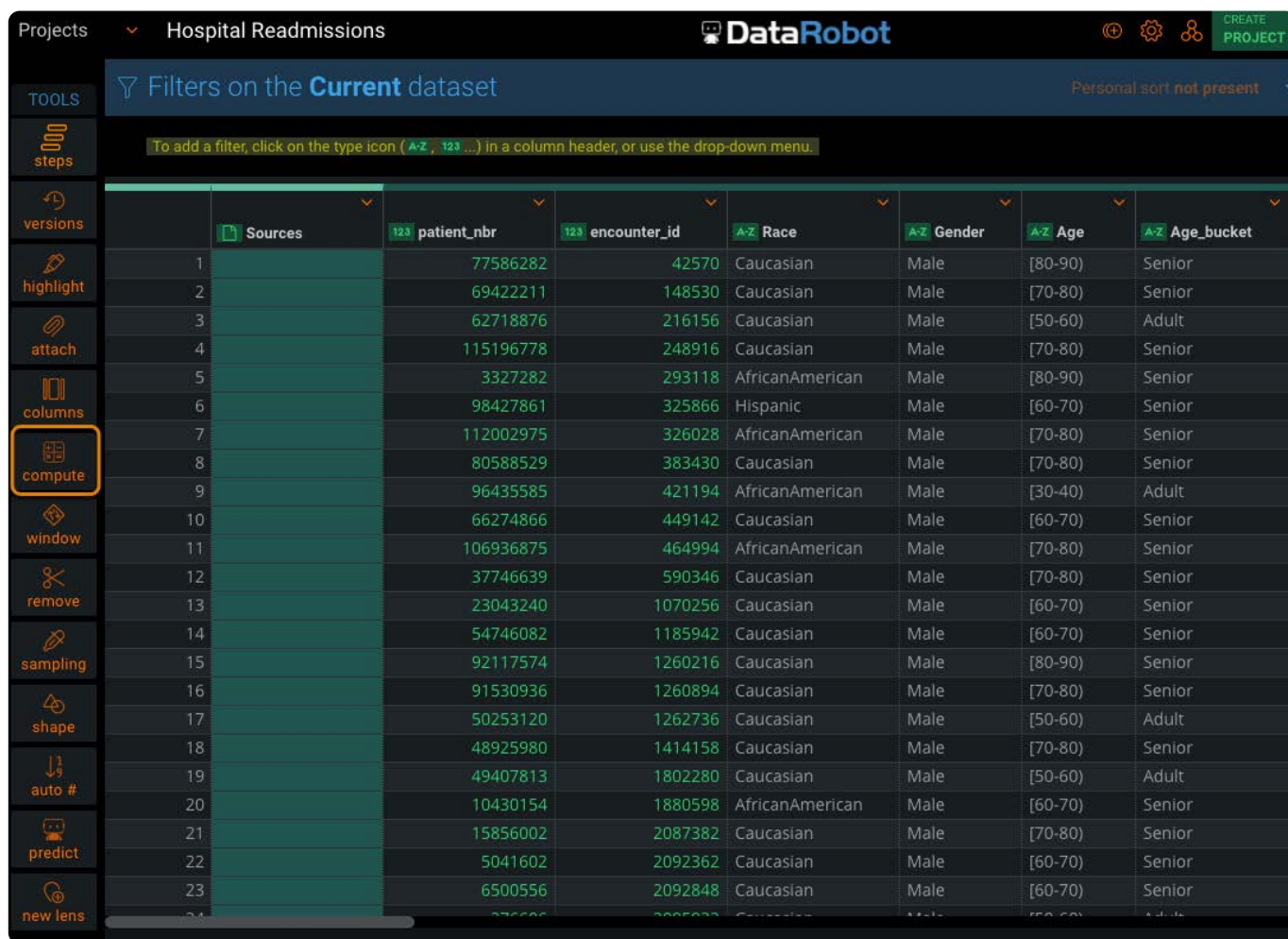
Your change is saved as a step in your project.

Compute columns

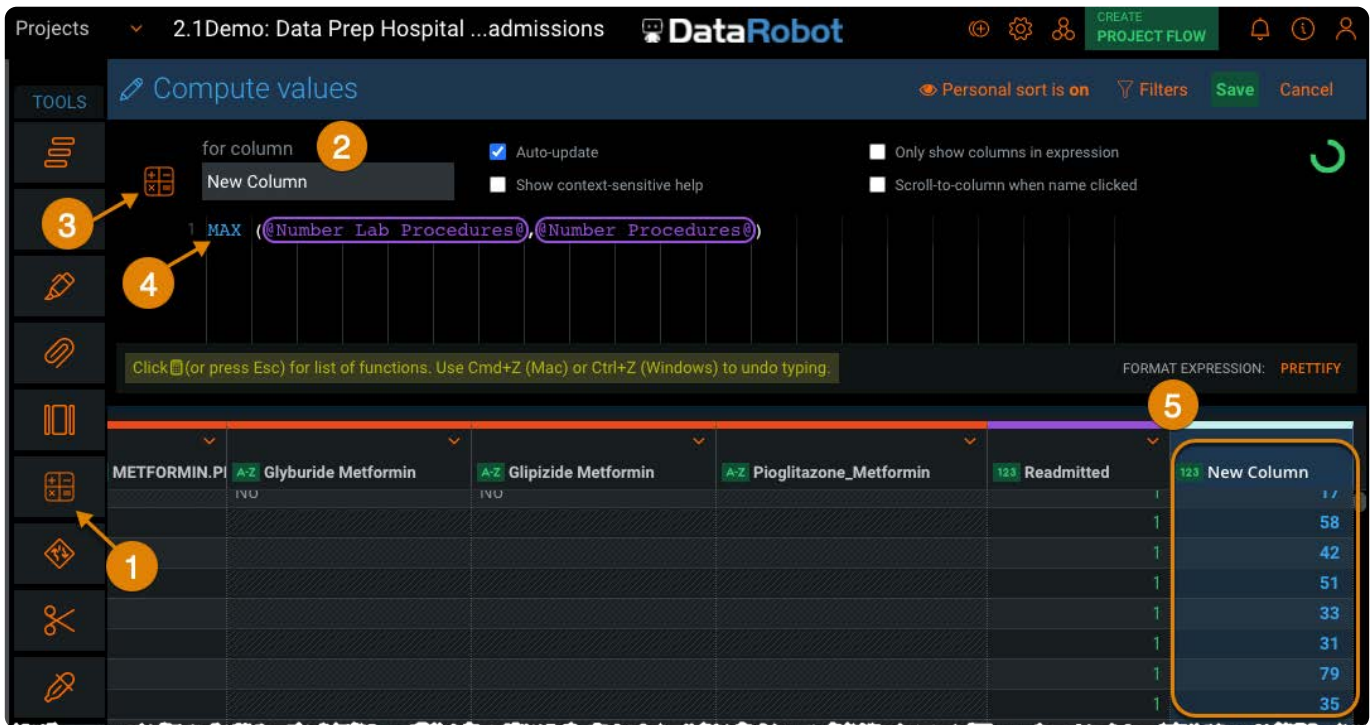
With Data Prep, you can add new columns by applying functions to existing columns in your dataset.

Work with the Compute tool

To access the Data Prep **Compute** tool, click **compute** in the project **Tools** bar:



Following is an overview of the elements in the **Compute values** pane:



ELEMENT	DESCRIPTION
---------	-------------

- | | | |
|---|----------------------------|---|
| 1 | Compute
tool | Click compute to access the Compute values pane. |
| 2 | for column
field | Enter the name of the new column. |
| 3 | functions
menu | Click to access the functions menu which contains a list of the supported functions. See How to add a computed column and Supported functions for details. |
| 4 | formula row | Use this row to build the formula to compute values for your new column. Enter simple formulas or build advanced formulas using columns and functions. This example shows the MAX function. Notice the column titles are enclosed between two '@' (@) signs. If there is an error in your formula, an error message appears below the formula row explaining the error.
Note: Data Prep is case sensitive. When typing column titles in the formula row, the case of your input must match the case of the column's name. |
| 5 | computed
column | The new column that will contain the result of your formula. |

ELEMENT

DESCRIPTION

6

timer

The timer provides a countdown to indicate how long before the Data Preview refreshes. The timer appears after you pause your typing, and resets again after either you continue entering the expression or the grid refreshes. This feature is especially useful when you're constructing a large computed column expression and you don't want the grid to refresh after every key stroke. If you prefer to disable the timer, deselect the **Auto-update** check box adjacent to the new column field. An **Update** button is then available in place of the timer and clicking it manually refreshes the grid.

How to add a computed column

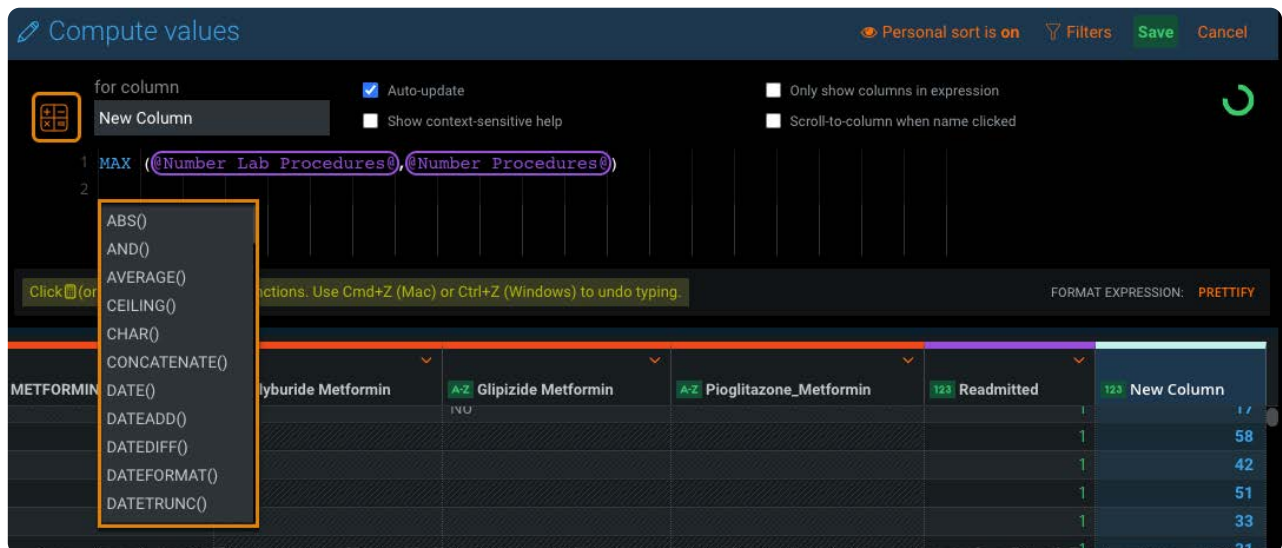
To add a computed column to your project:

1. Click **compute** in the **Tools** bar.

The **Compute values** pane displays.

2. Enter the name of the new column in the **for column** field.
3. Use one of the following methods to add the function to the formula:

- Click the functions icon at the top-left of the **Compute values** pane and select a function from the list:



- In the formula row, press **ESC** and select a function from the list.
- In the formula row, type the function.

See [Supported functions](#) for information on the available functions.


4. Use one of the following methods to add the column between the parentheses in the formula:

- From the Data Preview, click the column name.
- In the formula row, type **@**, then press **ESC**, and select the column from the list.
- In the formula row, type the name of the column between two **@** characters. For example, if the column's title is Hire Date, type **@Hire_Date@**.

5. Click **Save**.

The new column with the computed values is committed to your project.

Note

If there is an error in your formula, an error icon () displays in the Steps tool. In this case, you can create and save a lens but you cannot publish it to an AnswerSet.

Syntax rules

- Put spaces around operator symbols, for example, `1 + 1` not `1+1`
- Put double quotes around text, for example `"Hello"`.

If you need to literally type a quote sign, type a backslash before it. To type a backslash, type an additional backslash before it (because the backslash character needs to be escaped with a backslash). For example, to enter the text: `Go to "C:`

`\windows"` Type the following: `"Go to \"C:\\windows\""`

Supported functions

The following is a list of supported functions. Select a function to see its description, syntax, and an example.

Date/Time functions

- [DATE](#)
- [DATEADD](#)
- [DATEDIFF](#)
- [DATEFORMAT](#)
- [DATETRUNC](#)
- [DATEVALUE](#)
- [DAY](#)
- [DAYOFWEEK](#)
- [DAYOFYEAR](#)
- [ENDOFMONTH](#)
- [FROMUNIXTIME](#)
- [HOUR](#)
- [MAXDATE](#)
- [MIDNIGHT](#)
- [MINDATE](#)
- [MINUTE](#)
- [MONTH](#)
- [NETWORKDAYS](#)

- NOW
- QUARTER
- SECOND
- SETTIMEZONE
- TODAY
- WEEKOFYEAR
- WORKDAY
- YEAR

Informational functions

- FIRSTNONBLANK
- ISBLANK
- ISDATE
- ISNULL
- ISNUMBER
- ISTEXT

Logical functions

- AND
- IF
- IFERROR
- NOT
- OR

Mathematical functions

- Multiply, divide a column by any number
- Add to, subtract from a column by any number
- ABS
- CEILING
- EXP
- FACTORIAL
- FLOOR
- INT
- LN
- LOG

- LOG10
- MOD
- POWER
- ROUND
- ROUNDDOWN
- ROUNDPERC
- ROUNDUP
- SIGN
- SQRT
- SUM

Statistical functions

- AVERAGE
- MAX
- MEDIAN
- MIN
- MODE
- STDEV
- STDEVP
- VAR
- VARP

Text functions

- CHAR
- CONCATENATE
- FIND
- HASHVALUE
- LEFT
- LEN
- LOWER
- MID
- PADLEFT
- PADRIGHT
- REGEXP
- REPEAT
- REPLACE

- [REVERSE](#)
- [RIGHT](#)
- [SEARCH](#)
- [STR](#)
- [SUBSTITUTE](#)
- [TRIM](#)
- [TRIMLEFT](#)
- [TRIMRIGHT](#)
- [UPPER](#)
- [VALUE](#)

Comparison operators

- equal to
- greater than
- greater than or equal to
- less than
- less than or equal to
- not equal to

See [Comparison operators](#) for details.

Custom functions

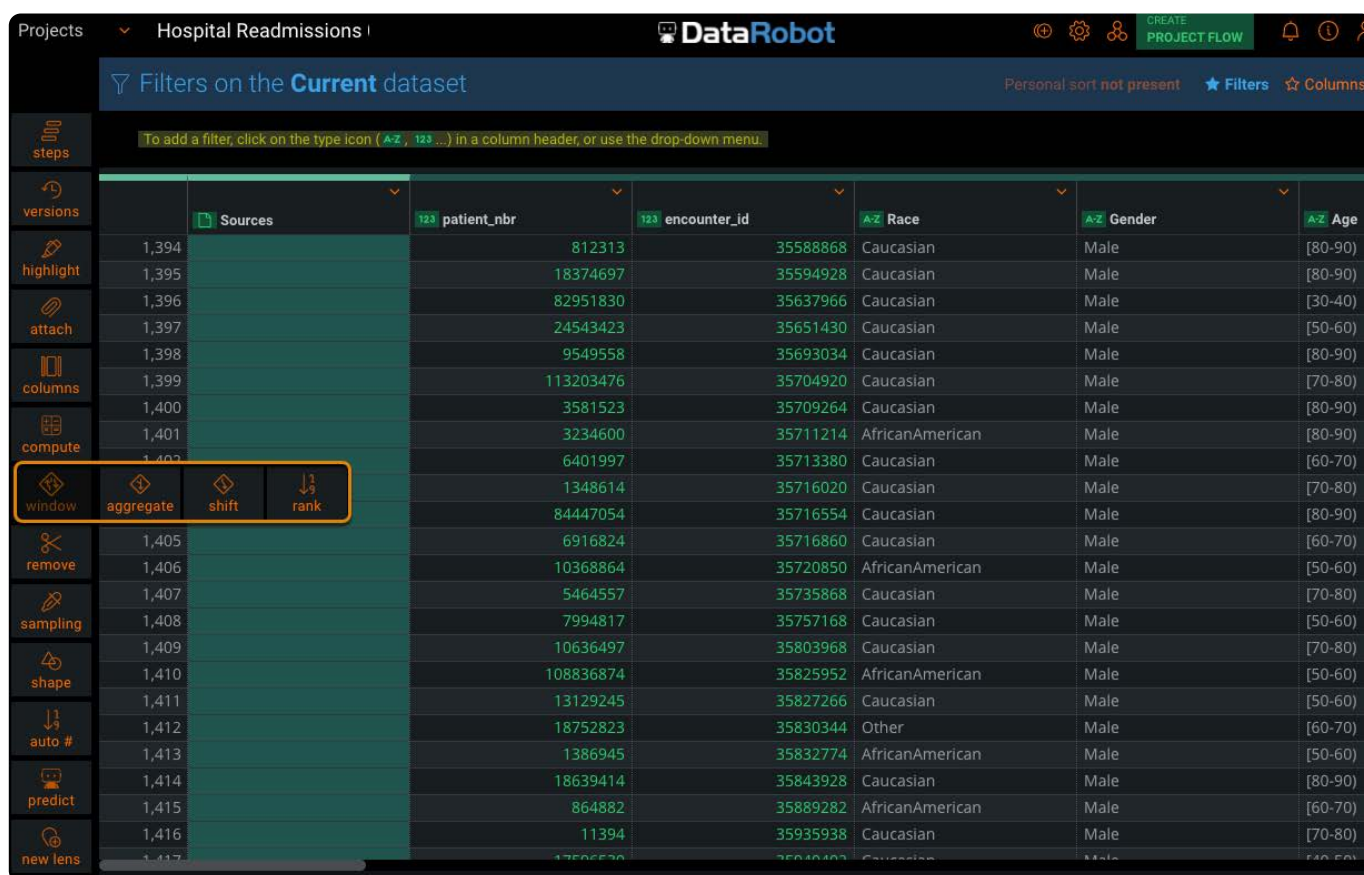
If your organization has developed and installed custom functions, you can add new columns by applying a custom function to existing columns in your dataset. See [Custom computed column functions](#) for details.

Group with windows

Data Prep Window functions are a set of tools that enable you to group sets of rows called a *window* for the purpose of performing a specific function.

Work with windows

To work with windows, hover over the **Window** tool in the project **Tools** bar and select the [aggregate](#), [shift](#), or [rank](#) tools.



Note

If you do not see the **Window** tool, contact your System Administrator who can enable this option for you.

Work with the Aggregate tool

The **Aggregate** tool lets you group sets of rows for the purpose of computing a specific function. In Data Prep, instead of writing formulas, calculations are created through point-and-click Steps.

In Data Prep, a window is defined as the set of rows participating in a computation. Windows can be identified as:

- Fixed window: A grouping of rows based on a shared similar value in one or more columns.
- Sliding window: A grouping of rows relative to the current row (i.e. rolling or moving functions).

Note

Aggregation in **Group By** is different than the Window Aggregate function. One main distinction is that in GroupBy, the number of rows is reduced to only leave the aggregate row value, while Window functions compute the aggregates on each row of the dataset.

Define your window aggregates

The screenshot shows the 'Compute' configuration window. At the top, it says 'Compute Average on territory as New Column'. Below this are four rows of configuration options: 'Window Grouped by' (group by column(s)), 'Window Sorted by' (sort column(s)), 'Starting from' (top of window), and 'Ending at' (bottom of window). Annotations with arrows point to these fields: 'Average' is labeled 'Select the function type from the drop-down.', 'territory' is labeled 'Select the existing column that you want to perform the function on.', 'New Column' is labeled 'Name the new column that will show the results of your computations.', and the 'Starting from' and 'Ending at' fields are grouped together and labeled 'Set the Window Boundaries Make fixed or sliding windows by defining where they should begin and end relative to each row.' On the left, a separate annotation 'Form your Windows Group by: select which column(s) contains the values by which you want rows aggregated. Sort by: optionally, add column(s) to change the order of your row values.' points to the 'Window Grouped by' and 'Window Sorted by' fields.

1. Create your output column.

Notice that a preview of this new column appears next to the column you select to do the function on.

Compute **Average** on **territory** as **New Column**

Window Grouped by: group by column(s)

Window Sorted by: sort column(s)

Starting from: top of window

Ending at: bottom of window

The new column appears next to the column being computed and updates as changes are made.

	A-Z territory	123 New Column	A-Z Region	123 Regional Sales (\$M)
1	NA	→	Central US	500
2	NA	→	Western US	200
3	EMEA	→	Africa	150
4	EMEA	→	UAE	340
5	EMEA	→	UK	210
6	APAC	→	Australia	400
7	APAC	→	China	600
8	APAC	→	Japan	378

- **Function type:** Drop-down includes Average, Sum, Count, First, Last, Min, Max, Median. Average is the default.
- **Pick the column** to do the computation on. The column type must be numeric for Average, Sum, Min, Max, and Median. Count, First, and Last can be any column type.
- **Name** the output column.

2. Group and Sort your rows.

- **Group by:** From the drop-down, select which column(s) contains the values by which you want to group your rows by. Select at least one, but you may include multiple columns here to make more/smaller windows. Aggregated rows in the column(s) will show up in ascending order.
- **Sort by** (optional): From the drop-down, select the column(s) to define the order by which your rows within the aggregate should be stacked. Click the blue arrow on the selected column to toggle between sorting the values in ascending or descending order. This is especially important for the First and Last functions, as the selection of the First/Last is done after the sort is applied.

3. Set the boundaries for the function.

- **Top of window:** For each row in the aggregate, the function will begin with the value in the first row of the window to which it belongs.
- **Bottom of window:** For each row in the aggregate, the function will end with the value in the last row of the window to which it belongs.
- **Current row with offset of __ rows:** Selecting this as the top, bottom, or both boundaries will create a rolling window that changes with relation to the row being computed (the current row). The computation will begin/end with the specified number of rows above or below the row being calculated. For instance, if you select the current row with offset of 0 as starting and ending, the computation for each row will just be based on that single row value. If you offset by -1, it will compute the function, including the row itself and 1 row before. See the examples below.


• Possible Options:

- Top of window to bottom of window (Fixed window)
- Top of window to current row with offset (Sliding window)
- Current row with offset to bottom of window (Sliding window)
- Current row with offset to current row with offset (Sliding window)

Examples

Retail Company Sales

The simple dataset below shows regional sales for a retail company.

		A-Z territory	A-Z Region	123 Regional Sales (\$M)
1		NA	Central US	500
2		NA	Western US	200
3		EMEA	Africa	150
4		EMEA	UAE	340
5		EMEA	UK	210
6		APAC	Australia	400
7		APAC	China	600
8		APAC	Japan	378

USE CASE 1: FIXED WINDOW

The company wants to compare the Average Regional Sales between territories to compare their performances.

With the settings in the image below, the rows have been grouped by *Territory* and sorted in ascending order alphabetically by their *Region*.

Setting the boundaries as the top of the window through the bottom window is indicative of a fixed window. This means the average is computed for each row using all of the values in the *Regional Sales* column within the aggregate. The new preview column, *Average Regional Sales*, shows up next to the *Regional Sales* column which is being used for the computation.

Compute Average on Regional Sales (\$M) as Average Regional Sal...

Window Grouped by territory ASC

Window Sorted by Region ASC

Starting from top of window

Ending at bottom of window

	territory	Region	Regional Sales (\$M)	Average Regional Sales
1	APAC	Australia	400	459.33333333333333...
2	APAC	China	600	459.33333333333333...
3	APAC	Japan	378	459.33333333333333...
4	EMEA	Africa	150	233.33333333333333...
5	EMEA	UAE	340	233.33333333333333...
6	EMEA	UK	210	233.33333333333333...
7	NA	Central US	500	350
8	NA	Western US	200	350

You can now see that all rows in a given window have the same calculated average value.

Thus, the company can now compare the *Average Regional Sales per Territory* and find that, on average, the APAC territory is performing the best and EMEA is performing the worst.

Company Salaries

The simple dataset below shows employee salaries in a company.

	▼	▼	▼	▼
	📄	A-Z Department	123 EmployeeNumber	123 Salary
1		sales	11	5200
2		sales	7	4200
3		sales	4	4500
4		sales	5	6000
5		sales	10	5500
6		hr	2	4000
7		hr	6	3500
8		research	1	4800
9		research	3	5000
10		research	8	4800

USE CASE 2: SLIDING WINDOW

The company wants to analyze salaries within the Sales department compared to the average salary of the two people previously hired in that department.

To do this, group the employees based on their *Department*, and stack them in the order they were hired using their *Employee Number* as the sorting order. Instead of finding the average for the entire Sales department, create a sliding (or rolling) window that shows the average salary for each Sales employee and the Sales employee hired immediately prior. Set the window to start one row before each employee (with offset of -1) and ending with that row (offset of 0).

Compute Average on Salary as Average Salary

Window Grouped by Department ASC

Window Sorted by EmployeeNumber ASC

Starting from current row with offset of -1 rows

Ending at current row with offset of 0 rows

	Department	EmployeeNumber	Salary	Average Salary
1	hr	2	4000	4000
2	hr	6	3500	3750
3	research	1	4800	4800
4	research	3	5000	4900
5	research	8	4800	4900
6	sales	4	4500	4500
7	sales	5	6000	5250
8	sales	7	4200	5100
9	sales	10	5500	4850
10	sales	11	5200	5350

Now, look at the Sales department to see that the window aggregate changes for each computation. Each average is based on the previous and current row.

	Department	EmployeeNumber	Salary	Average Salary
1	hr	2	4000	4000
2	hr	6	3500	3750
3	research	1	4800	4800
4	research	3	5000	4900
5	research	8	4800	4900
6	sales	4	4500	4500
7	sales	5	6000	5250
8	sales	7	4200	5100
9	sales	10	5500	4850
10	sales	11	5200	5350

Employee 4 was the first hire in the Sales Department, so the window only includes this row.

The rest of the aggregates are comprised of the row itself and one above.

Using this information, the Company determines employee number 7 is the only employee in the Sales department who earns a salary lower, on average, than the previously hired Sales team members.

USE CASE 3: SLIDING WINDOW

The company wants to analyze how the average salary in a department changes over time.

To do this, create a sliding window that starts at the top of the window and grows with each row that follows. This is similar to creating a RunningAverage function in Excel, but without needing to write a formula including \$ to set your absolute cells.

Like the previous example, group the rows together by *Department* and sort the *Employee Numbers* in ascending order. But unlike the previous example, the top boundary remains as a fixed starting point while the bottom of the window changes with relation to the row being calculated (i.e current row). Thus, set the window to start from the top row and end at the current row with an offset of 0. See the windows created below for the Sales Department:

	Department	EmployeeNumber	Salary	Average Salary
1	hr	2	4000	4000
2	hr	6	3500	3750
3	research	1	4800	4800
4	research	3	5000	4900
5	research	8	4800	4866.666666666...
6	sales	4	4500	4500
7	sales	5	6000	5250
8	sales	7	4200	4900
9	sales	10	5500	5050
10	sales	11	5200	5080

Using the output column of this sliding window function, you can see that over time the average salary in the Sales Department has generally increased. Also, an interesting trend emerges from this calculation which shows that the variance between the averages is getting smaller as time goes on, meaning salaries are becoming more consistent.

USE CASE 4: SLIDING WINDOW

The company wants a running total of salaries within each department for the purpose of analyzing which department's total salary exceeds \$15,000.

Grouping by *Department* and sorting *Employee Number* in descending order will stack your rows with the most recent hires at the top of your aggregate. Set a sliding window by starting the window with the current row being calculated (offset by 0) and

end it at the bottom of the window. This is similar to creating a Running Total in Excel, but without needing to write a formula including \$ to set your absolute cells.

Compute Sum on Salary as Total Payroll

Window Grouped by Department asc

Window Sorted by EmployeeNumber DESC

Starting from current row with offset of 0 rows

Ending at bottom of window

	Department	EmployeeNumber	Salary	Total Payroll
1	hr	6	3500	7500
2	hr	2	4000	4000
3	research	8	4800	14600
4	research	3	5000	9800
5	research	1	4800	4800
6	sales	11	5200	25400
7	sales	10	5500	20200
8	sales	7	4200	14700
9	sales	5	6000	10500
10	sales	4	4500	4500

For this running total, it's important to have employees sorted by the order they were hired. To have the most recent hire at the top of the aggregate, sort by *Employee Number* in descending order.

As the window slides to exclude employees hired most recently, it becomes apparent that the total pay was below \$15,000 in the green window, which was before Employee 10 was hired.

As illustrated in the image above, the running total for the HR and Research departments are both under the \$15,000 amount being analyzed. However, the current total payroll for the Sales Department is \$25,400. Upon further examination, you can locate when the total salary threshold was crossed. The green window has salaries for Employees 4, 5, and 7 and is below the \$15,000 threshold. However, when Employee 10 is hired, the threshold is crossed as indicated by the blue window.

Work with the Shift tool

The **Shift** tool creates a new column next to the original column being referenced with the row values shifted up or down by the set number of cells. You may be familiar with this action as Lag or Lead in SQL. Upon selecting the **Shift** tool, identify the column to shift and give it a new name. Then, define the direction and amount of rows to offset by. Optionally, you may use the Group and Sort fields to order the row values as needed.

For example, if a company wants to compare monthly sales to the previous month, it would be useful to see the sales figures for consecutive months next to each other for comparison. Use the Shift Down feature to create a new "Previous Month" column based on the original Sales column, with an offset of 1 row. Now, with the Sales values side-by-side the company can easily create a Computed Column in order to quantify the difference between the sales figures month to month.

The new column preview, named "Previous Month", is created next to the original column being selected to shift.

A-Z Month	123 Sales	123 Previous Month
January	1000	
February	1200	1000
March	900	1200
April	875	900
May	1250	875
June	1125	1250
July	1000	1125
August	1050	1000
September	900	1050
October	1550	900
November	1900	1550
December	2000	1900

Each row value is shifted down by 1 cell. Note that the first row is blank because no previous value exists.

Work with the Rank tool

The **Rank** tool includes the ability to add a column to your dataset to denote Rank, Dense Rank, or Row Number.

- **Rank:** Two or more rows share the same rank if they have tied values. The next row will have the next nonconsecutive value. For example, two matching cells would both be given the rank of 1, and the next cell will be ranked 3. If you use the Sort by option, the ranking is based on the sort order. If you use the Group by option, the ranking is specified within a group and resets to 1 for each top entry of the next grouping of rows.
- **Dense Rank:** Similar to Rank, each tied row value will share the same rank. However, the next unique value will simply continue with the next consecutive integer, without regard to the number of values that shared the rank before it. For example, two matching cells would both be given the rank of 1, and the next unique cell will be ranked 2. The same Sort and Group by rules apply as the rank operation.
- **Row Number:** This function assigns a sequential ID to each row, starting with 1, based on the sort order of the window (see Image 1 below). Use the Group or Sort fields to change the order of the rows, and the number values will start with 1 for each partition created (see Image 2 below).

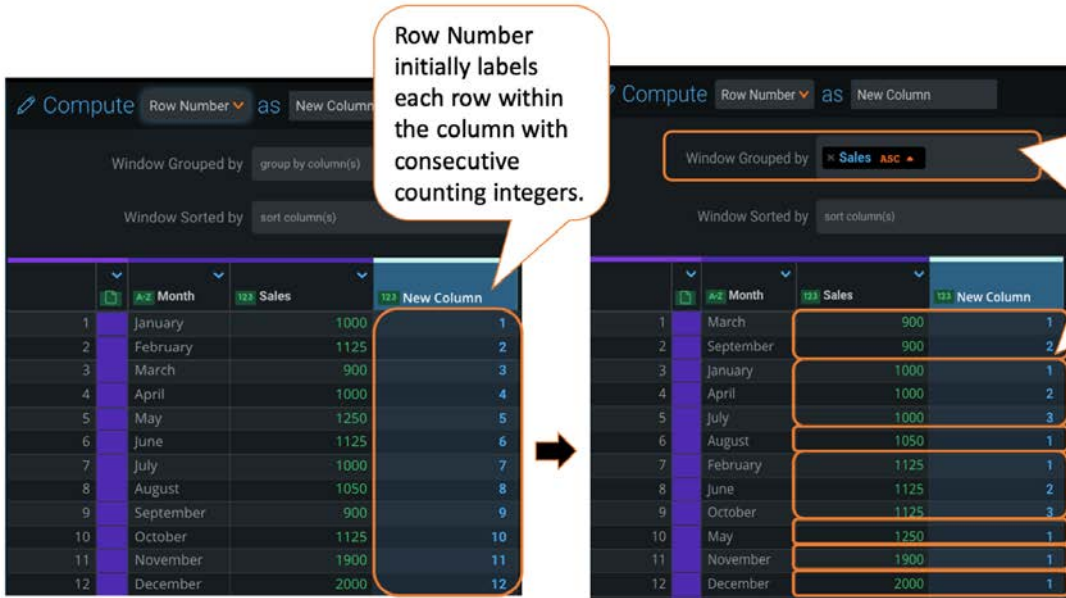


Image 1

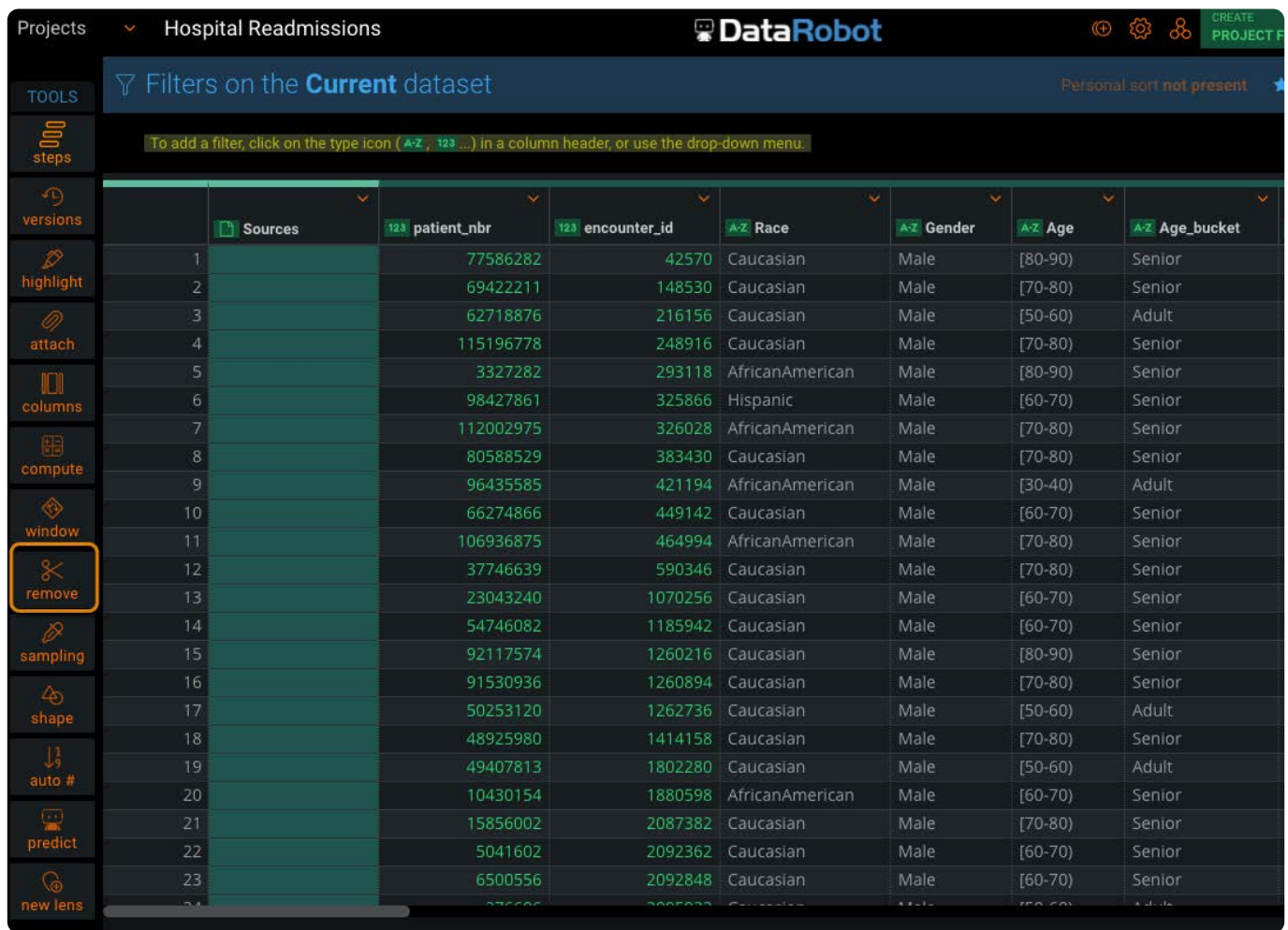
Image 2

Remove rows

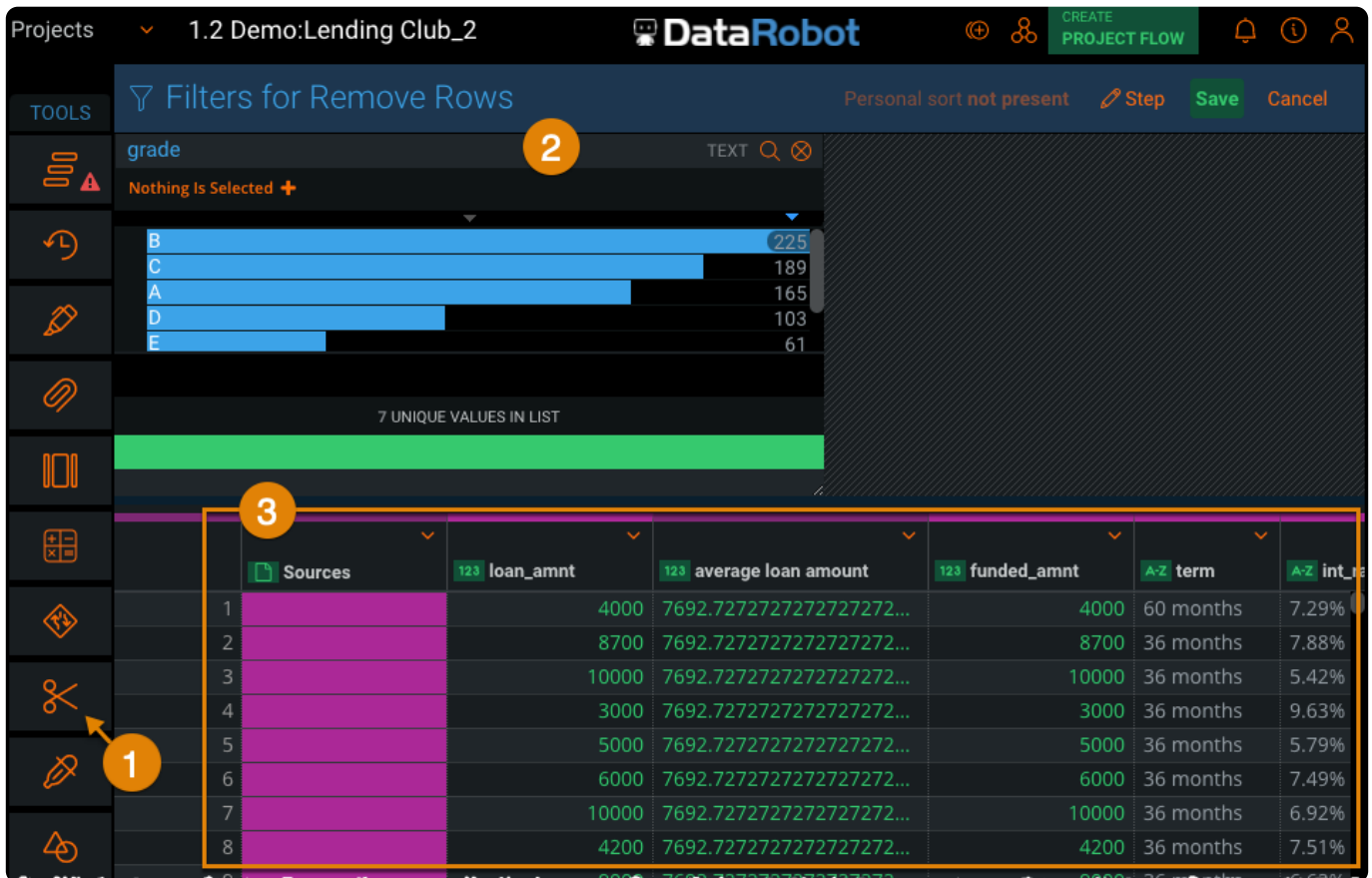
As you prepare your data in Data Prep, you will find times when you want to retain a specific subset of data. The best way to accomplish this is to remove the rows of data that don't meet your needs.

Work with the remove tool

To access the **remove** tool, click **remove** in the project **Tools** bar:



Following is an overview of the elements you work with when you remove rows from your project:



ELEMENT	DESCRIPTION
---------	-------------

- 1 Remove tool** Click **remove** to access the **Remove rows** pane and remove rows from your project.
- 2 Filter for Remove Rows pane** Filters your data down to a subset that you want to isolate for the remove operation. To access this pane, click **Filters** on the top right.
- 3 Data Preview pane** Displays the data in your project and how it changes as you prep it.

Remove rows

To remove rows from your data:

- Click the **Filters** link on the top right to add a Filtergram and isolate the rows you want to remove.
The **Data Preview** displays the records that match your filter criteria. See [Data Filtergrams](#) for details on working with Filtergrams.
- Click **remove** in the **Tools** bar.
The **Filters on the Current Dataset** pane becomes the **Filters for Remove Rows** pane.
- Click **Save**.

The filtered rows are removed. The **Data Preview** is blank because the filter used is still active and the data has been removed.

4. To display your updated dataset, do one of the following:

- On the Filtergram, click **x clear**.
- Close the Filtergram.

The Filtergram is cleared. The **Data Preview** displays the updated data.

Tip

Publish the data from the removed rows to a separate AnswerSet for later reference. See [Capture data from removed rows](#) for details.

Note

If you update or add to your dataset after removing rows, the remove rows Step will be applied to the new data in the following ways:

- If you selected histogram ranges or individual values to remove rows, rows from the new data will only be removed if they satisfy that exact condition.
- If you selected rows with a string search or dynamic percentile range, then the tool will recompute based on the new data.

Capture data from removed rows

To add a lens to publish the rows you removed to an [AnswerSet](#):

1. From **Tools**, click **steps**.

The **Steps** pane appears.

2. Click the Step prior to the **Remove Rows** Step.

The **Data Preview** displays the records that match the criteria of your filter.

3. Add a [Filtergram](#) to isolate the rows you removed.

The **Data Preview** displays the records that match the criteria of your filter.

4. From **Tools**, click [lens > new lens](#).

5. Click **Save**.

The lens is added to the project.

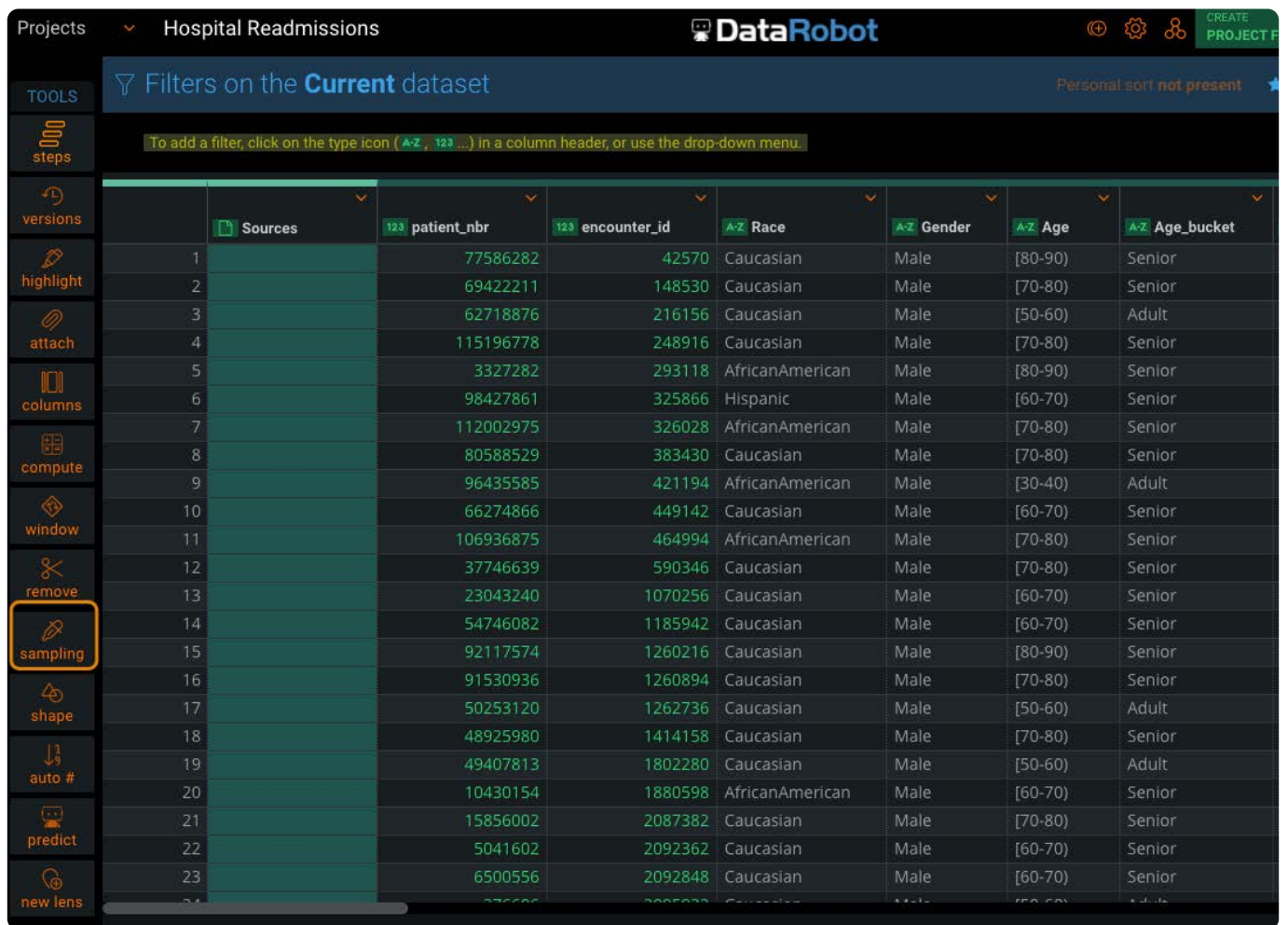
6. To publish the lens, click **Publish** from the lens.

Sample datasets

You may find it useful to work with a sample of a dataset before bringing all the data into your Data Prep project. For large datasets, this can make initial exploration and discovery easier. The **Sampling** tool also gives you the flexibility to filter down to a specific set of rows in your data, and then sample on the remainder.

Work with the Sampling tool

To access the **Sampling** tool, click **sampling** in the project **Tools** bar:



You may want to sample a very large dataset for initial discovery before bringing all of the data into your project. The **Sampling** tool also gives you the flexibility to filter down to a specific set of rows in your data, and then sample on the remainder.


Note

If you choose to sample your data, you are only shown the patterns, lookup combinations, and aggregations for that sample. When your exploration is complete, you can easily remove the sampling operation by either muting or deleting it in the **Steps** pane.

Sampling methods

Sampling can be based on a percentage of your dataset or a specific number of rows in the dataset.

- **Percentage-based sampling:** Perform a random and repeatable sample across your dataset based on the percentage you specify. You can also choose to specify a column in your dataset that is used for generating the sample. In this case, only the data in the column is used for determining the sample.
- **Row-based sampling:** Perform a random and repeatable sample across your dataset based on the number of rows you specify. The number of rows you specify is divided by the total number of rows in your dataset. A subset sample of your data is returned. If you are performing row-based sampling as a data prep step in your project, the number of rows you specify is divided by the total number of rows in your dataset from the previous step.

For both types of sampling, you can save the "sampling seed" number to ensure that you can repeat your sampled subset of data. You can also click the green reseed  icon to produce a different subset sample of your data. For an optimal sample, your dataset should exceed 100k rows.

Sample using percentage

To create a sample based on percentage of your dataset:

1. From the **Tools** bar, click **columns**.

Projects Hospital Readmissions DataRobot

Filters on the Current dataset


To add a filter, click on the type icon (A-Z, 123...) in a column header, or use the drop-down menu.


	Sources	patient_nbr	encounter_id	Race	Gender	Age	Age_bucket
1		77586282	42570	Caucasian	Male	[80-90]	Senior
2		69422211	148530	Caucasian	Male	[70-80]	Senior
3		62718876	216156	Caucasian	Male	[50-60]	Adult
4		115196778	248916	Caucasian	Male	[70-80]	Senior
5		3327282	293118	AfricanAmerican	Male	[80-90]	Senior
6		98427861	325866	Hispanic	Male	[60-70]	Senior
7		112002975	326028	AfricanAmerican	Male	[70-80]	Senior
8		80588529	383430	Caucasian	Male	[70-80]	Senior
9		96435585	421194	AfricanAmerican	Male	[30-40]	Adult
10		66274866	449142	Caucasian	Male	[60-70]	Senior
11		106936875	464994	AfricanAmerican	Male	[70-80]	Senior
12		37746639	590346	Caucasian	Male	[70-80]	Senior
13		23043240	1070256	Caucasian	Male	[60-70]	Senior
14		54746082	1185942	Caucasian	Male	[60-70]	Senior
15		92117574	1260216	Caucasian	Male	[80-90]	Senior
16		91530936	1260894	Caucasian	Male	[70-80]	Senior
17		50253120	1262736	Caucasian	Male	[50-60]	Adult
18		48925980	1414158	Caucasian	Male	[70-80]	Senior
19		49407813	1802280	Caucasian	Male	[50-60]	Adult
20		10430154	1880598	AfricanAmerican	Male	[60-70]	Senior
21		15856002	2087382	Caucasian	Male	[70-80]	Senior
22		5041602	2092362	Caucasian	Male	[60-70]	Senior
23		6500556	2092848	Caucasian	Male	[60-70]	Senior

The **Sample using** pane appears.

Sample using **Percentage** Rows on Select column (optional) Personal sort is off Filters Save Cancel

BY PERCENTAGE SAMPLING SEED

% 

2. Click **Percentage** if it is not already selected.
3. Optionally select a column.
The sampling percentage is based on the selected column
4. In the **By Percentage** field, enter the percentage of the dataset that you want included in the sample.
5. Optionally click the green **reseed**  icon.
6. Click **Save**.

Sample using rows

To create a sample based on percentage of your dataset:

1. From the **Tools** bar, click **columns**.

Projects Hospital Readmissions DataRobot

Filters on the Current dataset

To add a filter, click on the type icon (A-Z, 123...) in a column header, or use the drop-down menu.

	Sources	patient_nbr	encounter_id	Race	Gender	Age	Age_bucket
1		77586282	42570	Caucasian	Male	[80-90)	Senior
2		69422211	148530	Caucasian	Male	[70-80)	Senior
3		62718876	216156	Caucasian	Male	[50-60)	Adult
4		115196778	248916	Caucasian	Male	[70-80)	Senior
5		3327282	293118	AfricanAmerican	Male	[80-90)	Senior
6		98427861	325866	Hispanic	Male	[60-70)	Senior
7		112002975	326028	AfricanAmerican	Male	[70-80)	Senior
8		80588529	383430	Caucasian	Male	[70-80)	Senior
9		96435585	421194	AfricanAmerican	Male	[30-40)	Adult
10		66274866	449142	Caucasian	Male	[60-70)	Senior
11		106936875	464994	AfricanAmerican	Male	[70-80)	Senior
12		37746639	590346	Caucasian	Male	[70-80)	Senior
13		23043240	1070256	Caucasian	Male	[60-70)	Senior
14		54746082	1185942	Caucasian	Male	[60-70)	Senior
15		92117574	1260216	Caucasian	Male	[80-90)	Senior
16		91530936	1260894	Caucasian	Male	[70-80)	Senior
17		50253120	1262736	Caucasian	Male	[50-60)	Adult
18		48925980	1414158	Caucasian	Male	[70-80)	Senior
19		49407813	1802280	Caucasian	Male	[50-60)	Adult
20		10430154	1880598	AfricanAmerican	Male	[60-70)	Senior
21		15856002	2087382	Caucasian	Male	[70-80)	Senior
22		5041602	2092362	Caucasian	Male	[60-70)	Senior
23		6500556	2092848	Caucasian	Male	[60-70)	Senior


The **Sample using** pane appears.

Sample using Percentage Rows

Personal sort not present Filters Save Cancel

BY ROWS SAMPLING SEED

1000 rows 1623103215367

2. Click **Percentage** if it is not already selected.
3. Optionally select a column.
The sampling percentage is based on the selected column
4. In the **By Percentage** field, enter the percentage of the dataset that you want included in the sample.
5. Optionally click the green **reseed**  icon.
6. Click **Save**.

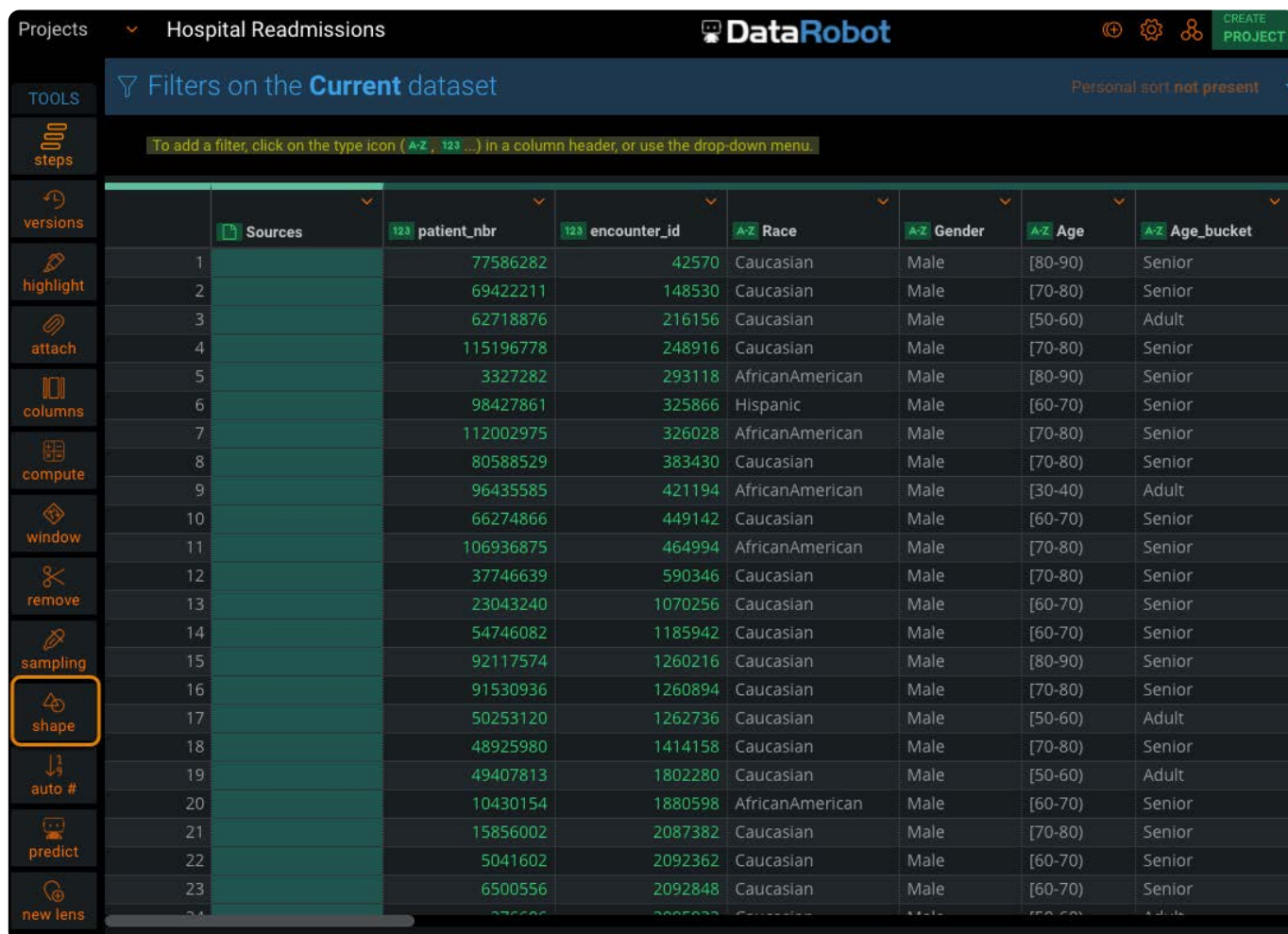
Shape data

Data Prep provides a **Shape** tool that lets you:

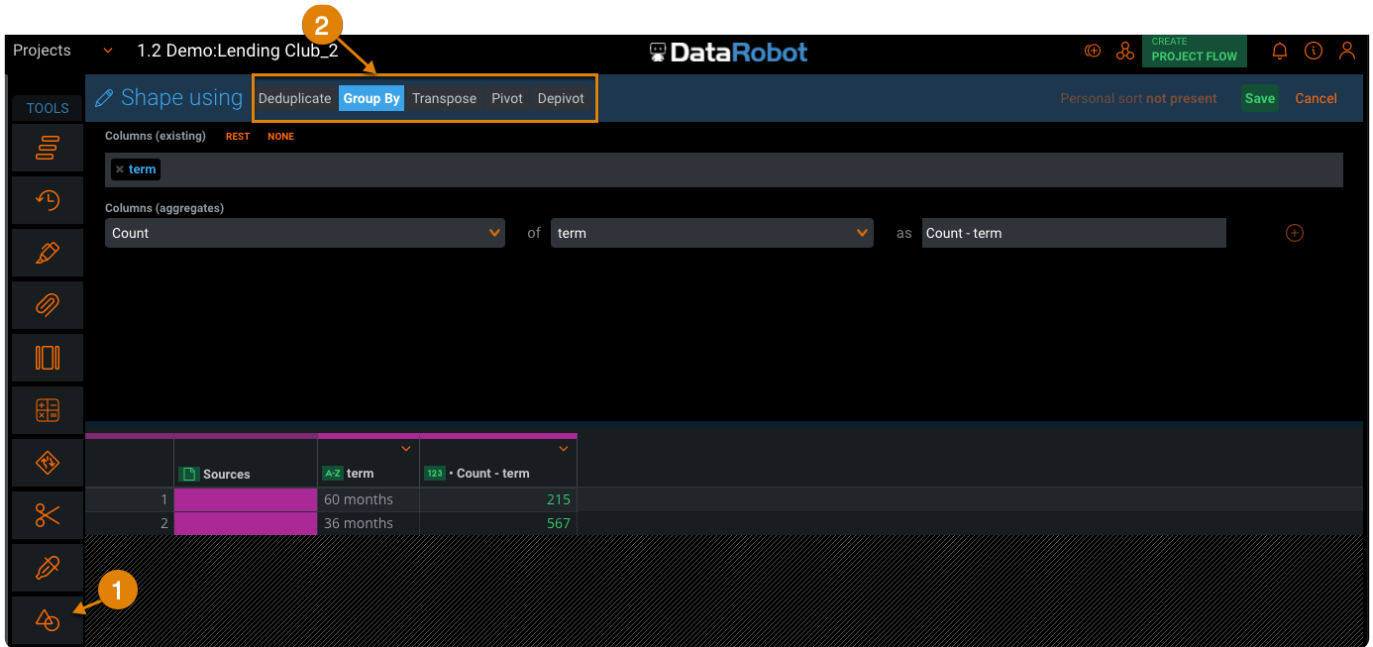
- Deduplicate
- Group By
- Transpose
- Pivot
- Depivot

Work with the Shape tool

To access the **shape** tool, click **shape** in the project **Tools** bar:



Following is an overview of the elements in the **Shape using** pane:

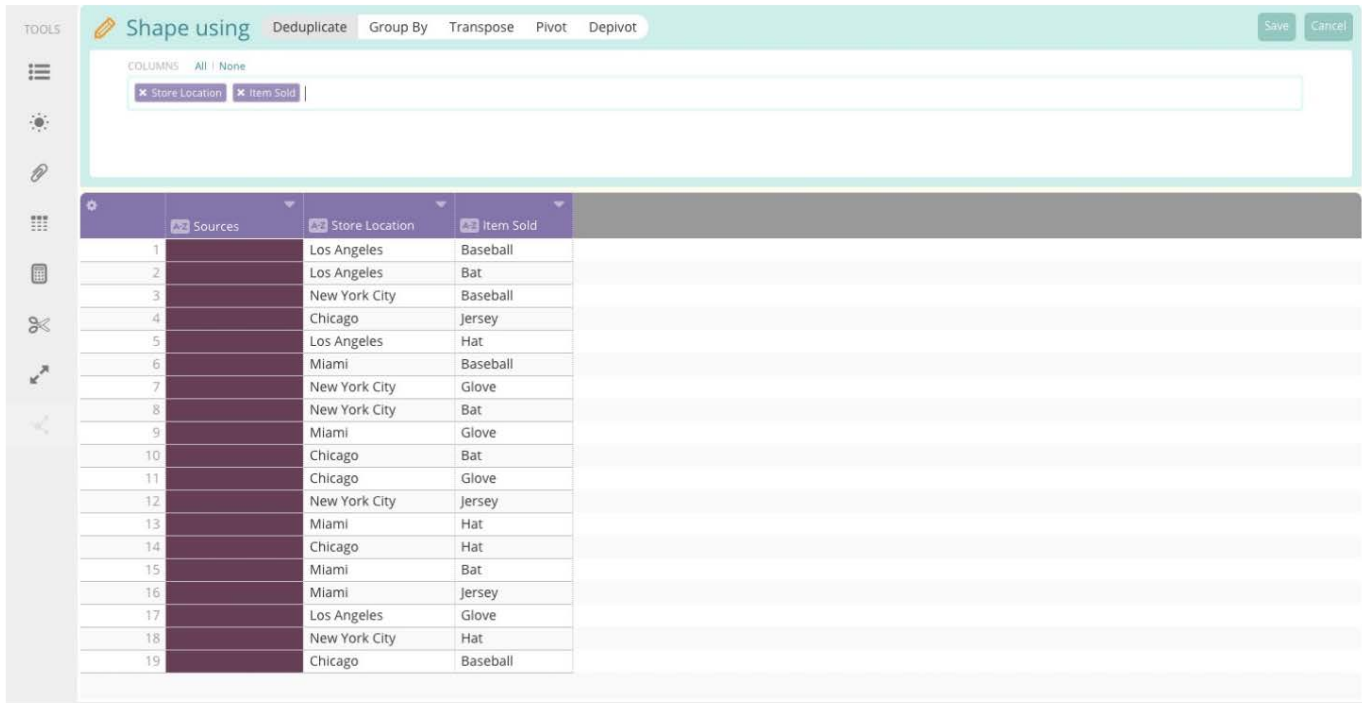


ELEMENT	DESCRIPTION
---------	-------------

- 1** **Shape tool** Click **shape** to access the **Shape using** pane.
- 2** **Shape using pane** Select a shaping tool in the **Shape using** pane:
 - [Deduplicate](#)
 - [Group By](#)
 - [Transpose](#)
 - [Pivot](#)
 - [Depivot](#)

Deduplicate

The **Deduplicate** function searches your data for rows that match each other exactly, and condenses them into a single row, removing duplicate entries from your data.



In the **Deduplicate** pane, you have the option to add or remove columns from the **Columns** field. Any columns added to this field are included in the deduplication process. As you add columns, they appear in the **Data Preview**.

Note

Only those columns added to the **Columns** field are included in the deduplication process and resulting output; any columns not selected are not considered during the deduplication process and are removed from the data upon its completion. Once you are satisfied with the parameters of the deduplication function, click **Save** to finalize the process and commit it to your project.

You can also use the **Fuzzy** option with the **Deduplicate** function. When enabled, you see this option as an alternate selection to **Exact** match.



Fuzzy deduplication finds matching rows using a fuzzy algorithm. This means similar values that are not exact matches are grouped together and deduplicated, for example:

These similar values would be grouped together with the Fuzzy algorithm...

	Sources	Fname	Lname	Address	City	State
1		Jennifer	jones	123 street	Santa clara	CA
2		Jennifer	Jones	123 street	santa clara	CA
3		Jenni	smith	123 street	santa clara	ca
4		jenni	jones	123 street	santa clara	CA

...providing this Deduplicated result!

Shape using **Deduplicate** Group By Transpose Pivot Depivot with Exact **Fuzzy** match

COLUMNS All | None

× Fname × Lname × Address × City × State

	Sources	Fname	Lname	Address	City	State
1		Jennifer	jones	123 street	santa clara	CA

Similarly, fuzzy deduplication groups similar items even if they have blank values:

These similar values would be grouped together with the Fuzzy algorithm...

	Sources	Fname	Lname	Address	City	State
1		Jennifer	jones	123 street	Santa clara	CA
2		Jennifer	Jones	123 street	santa clara	
3		Jenni	Jones	123 street		
4		jenni	jones	123 street	santa clara	CA

...providing this Deduplicated result!

Shape using **Deduplicate** Group By Transpose Pivot Depivot with Exact **Fuzzy** match

COLUMNS All | None

× Fname × Lname × Address × City × State

	Sources	Fname	Lname	Address	City	State
1		Jennifer	jones	123 street	santa clara	CA

Group By

The **Group By** function lets you use a number of different aggregate functions on existing columns within the dataset. When **Group By** is selected, a pane appear above the **Data Preview** where you can specify which columns to include in the Group By process, which columns you want to perform an aggregate function on, which aggregate function to use, as well as specify the name of the new aggregate column being created.

The **Data Preview** displays your selections, highlighted in blue, so that you can see how it affects the data. Remember that only those columns included in the **Columns (aggregates)** field remain in your data following the **Group By** function. Those columns included are used when identifying duplicate rows for grouping.

	Sources	Store Location	Item Sold	Item Type	Sum - Transaction Amount
1		Los Angeles	Baseball	Equipment	474.81
2		Los Angeles	Bat	Equipment	6899.54
3		New York City	Baseball	Equipment	462.84
4		Chicago	Jersey	Uniform	749.85
5		Los Angeles	Hat	Uniform	1024.59
6		Miami	Baseball	Equipment	1001.49
7		New York City	Glove	Equipment	559.93
8		New York City	Bat	Equipment	2849.81
9		Miami	Glove	Equipment	2879.64
10		Chicago	Bat	Equipment	1799.88
11		Chicago	Glove	Equipment	1839.77
12		New York City	Jersey	Uniform	3649.27
13		Miami	Hat	Uniform	299.88
14		Chicago	Hat	Uniform	1574.37
15		Miami	Bat	Equipment	1649.89
16		Miami	Jersey	Uniform	1449.71
17		Los Angeles	Glove	Equipment	3759.53
18		New York City	Hat	Uniform	224.91

For a list of available aggregate functions, see [Group by aggregate functions](#).

These operations are called aggregate because they find matching rows in the dataset and then combine them into one row. A matching row is defined as one that, excluding the reference column, share the same values in a column-by-column examination. The reference column is excluded from the column-by-column examination because its value are submitted to the aggregate function in order to produce the reference column value in the single-row result.

Transpose

The **Transpose** functions lets you switch your rows and columns (imagine rotating your data 90 degrees).

The transpose function allows you to select any one column to create new column headers with. The values of the selected column become the new column headers, while all other column headers shift and become row headers (unless the user removes them during the transpose process). In this new transposition, if there is more than one value that matches a particular set of headings, Data Prep displays the last available value from the original dataset.

Here we have a simple sales dataset:

TOOLS **Filters on the Current dataset**

To add a filter, click on the type icon (A-Z, 123...) in a column header, or use the drop-down menu.

	Sources	Store Location	Region	Item Sold	Item Type	Quantity	Transaction Amount	Transaction No.	Year
1		Los Angeles	West	Baseball	Equipment	50	199.5	1	2014
2		Los Angeles	West	Bat	Equipment	15	2249.85	2	2014
3		New York City	East	Baseball	Equipment	25	99.75	3	2014
4		Chicago	Central	Jersey	Uniform	10	499.9	4	2014
5		Los Angeles	West	Hat	Uniform	10	249.9	5	2014
6		Miami	East	Baseball	Equipment	18	71.82	6	2014
7		New York City	East	Glove	Equipment	7	559.93	7	2013
8		New York City	East	Bat	Equipment	19	2849.81	8	2014
9		Los Angeles	West	Baseball	Equipment	37	147.63	9	2014
10		Miami	East	Glove	Equipment	22	1759.78	10	2014
11		Chicago	Central	Bat	Equipment	12	1799.88	11	2014
12		Los Angeles	West	Hat	Uniform	4	99.96	12	2014
13		Chicago	Central	Glove	Equipment	10	799.9	13	2014
14		New York City	East	Jersey	Uniform	19	949.81	14	2014
15		Los Angeles	West	Bat	Equipment	10	1499.9	15	2013
16		Miami	East	Hat	Uniform	12	299.88	16	2013
17		Miami	East	Baseball	Equipment	64	255.36	17	2014
18		Miami	East	Glove	Equipment	3	239.97	18	2014
19		Miami	East	Baseball	Equipment	40	159.6	19	2013
20		New York City	East	Jersey	Uniform	12	599.88	20	2014
21		Chicago	Central	Glove	Equipment	7	559.93	21	2014
22		Los Angeles	West	Bat	Equipment	18	2699.82	22	2014
23		Chicago	Central	Hat	Uniform	21	624.78	23	2014

After clicking the **shape** tool, select **Transpose** at the top of the **Shape using** pane:

TOOLS **Shape using** Deduplicate Group By Transpose Pivot Depivot Save Cancel

COLUMN LABEL:

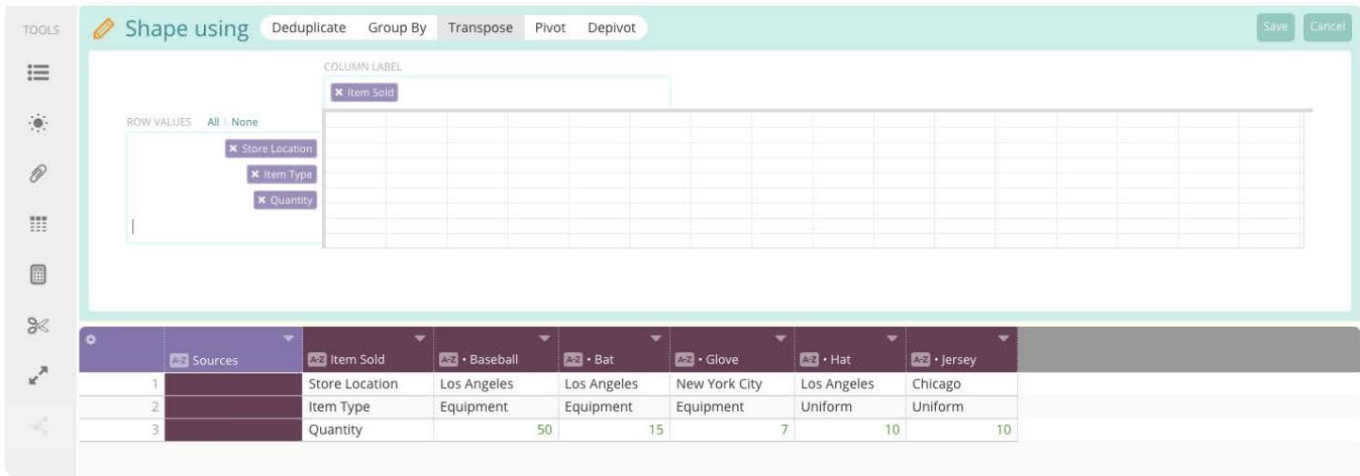
ROW VALUES: All | None

Item Sold
 Item Type

	Sources	Store Location	Chicago	Los Angeles	Miami	New York City
1		Item Sold	Jersey	Baseball	Baseball	Baseball
2		Item Type	Uniform	Equipment	Equipment	Equipment

You can choose one column for the column header (the values of this chosen column become the new column headers) as well as any columns you would like to include as rows in the newly transposed data. Pay attention to the grid below the **Shape using** pane as it displays a preview of how the data will appear based on the options you have chosen for the transpose process.

Be aware, depending on which column you choose for your transpose, you may not end up with a complete picture of your data. Notice that in the example, there is more than one instance where the "Item Sold" value is "Baseball". If you were to transpose with the "Item Sold" as your new columns, Data Prep will only display the last value in the dataset, as shown here:

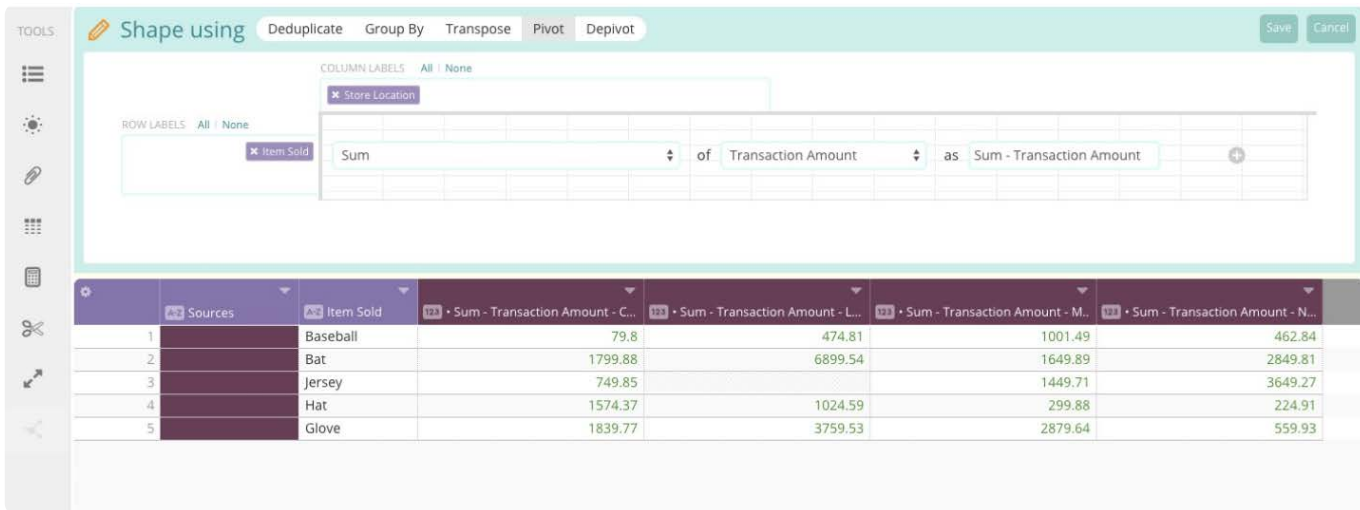


Once you are satisfied with the options you have chosen for the transpose process, click **Save** in the upper right to finalize and commit the process to your project.

Pivot

The **Pivot** function is similar to the **Transpose** function in that it allows you to rearrange your column headers into a row header. The **Pivot** function also allows you to perform aggregate functions on a selected column to be displayed as the body of data within your pivot table. Unlike the **Transpose** function, with **Pivot**, you can choose multiple headings as column headers.

When you open the **shape** tool and select **Pivot**, the Pivot options display. You can choose which columns to use as your column headers, as well as which columns to include as rows in the pivot table to be created. Finally, you can choose which aggregate function to apply to which column to create the body of the pivot table. If you click the **+** button, you can add more than one aggregation to your pivot table.



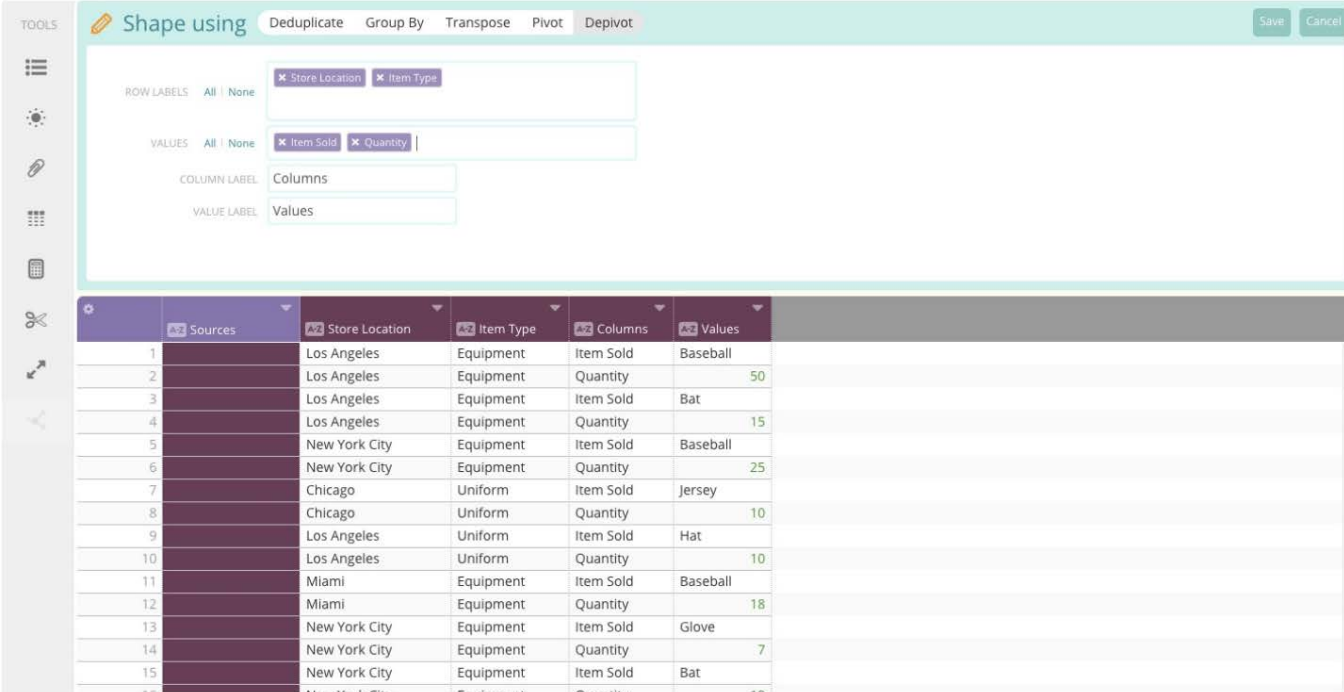
The selected aggregate function is applied to the chosen column, and then displays the total based on the intersection of your chosen column and row headers. For a list of the available aggregate functions, see [Group by aggregate functions](#).

Pay attention to the grid below the **Shape using** pane as it displays a preview of how the data will appear based on the options you have chosen for the pivot process.

Once you are satisfied with the options you have chosen for the transpose process, click **Save** in the upper right to finalize and commit the process to your project.

Depivot

The **Depivot** function is a powerful function that lets you take your data and stack columns into two columns—one column containing the column headers that the original data came from and a second column containing the value from that column. For users who have never used a function like **Depivot**, it can be difficult to understand exactly what is happening during the depivot process.



The screenshot shows the 'Depivot' tool interface. At the top, there are tabs for 'Deduplicate', 'Group By', 'Transpose', 'Pivot', and 'Depivot'. The 'Depivot' tab is active. Below the tabs, there are four configuration fields: 'ROW LABELS' (set to 'All | None'), 'VALUES' (set to 'All | None'), 'COLUMN LABEL' (set to 'Columns'), and 'VALUE LABEL' (set to 'Values'). Below these fields is a preview table with 15 rows and 5 columns. The columns are labeled 'Sources', 'Store Location', 'Item Type', 'Columns', and 'Values'. The 'Columns' and 'Values' columns are highlighted in a darker shade, indicating they are the selected columns for the depivot process.

	Sources	Store Location	Item Type	Columns	Values
1		Los Angeles	Equipment	Item Sold	Baseball
2		Los Angeles	Equipment	Quantity	50
3		Los Angeles	Equipment	Item Sold	Bat
4		Los Angeles	Equipment	Quantity	15
5		New York City	Equipment	Item Sold	Baseball
6		New York City	Equipment	Quantity	25
7		Chicago	Uniform	Item Sold	Jersey
8		Chicago	Uniform	Quantity	10
9		Los Angeles	Uniform	Item Sold	Hat
10		Los Angeles	Uniform	Quantity	10
11		Miami	Equipment	Item Sold	Baseball
12		Miami	Equipment	Quantity	18
13		New York City	Equipment	Item Sold	Glove
14		New York City	Equipment	Quantity	7
15		New York City	Equipment	Item Sold	Bat

In the Depivot window, you are given multiple options to configure.

- **Row Labels:** Select which columns you want to remain static in your data. These columns will not be "stacked" during the depivot process.
- **Values:** Choose which columns from your data will be included (or "stacked") during the depivot process.
- **Column Label:** Designate a column name for the new column containing the column labels from the columns selected in the **Values** field.
- **Value Label:** Designate a column name for the new column containing the values from the columns selected in the **Values** field.

In the previous example, the preview displays how the **Depivot** function takes the selected columns from the **Values** field and stacks them into the two new columns. To visualize exactly what is happening, see the following table:

LOCATION	COLUMN	VALUE
Los Angeles	Q1 SALES	\$1,250
Los Angeles	Q2 SALES	\$2,700
Los Angeles	Q3 SALES	\$2,465
Palo Alto	Q1 SALES	\$975
Palo Alto	Q2 SALES	\$1,680
Palo Alto	Q3 SALES	\$1,800
Seattle	Q1 SALES	\$1,100
Seattle	Q2 SALES	\$2,355
Seattle	Q3 SALES	\$2,890

The columns have been color-coded to see where specific values are coming from and moving to during the depivot process.

In this example, the **LOCATION** column has been selected as the only value for the **Row Label**. Due to the stacking of the column values during the depivot process, there are now duplicate rows for each location. This is to allow each value from the columns selected in the **Values** field to have its own row. You can see that by using the **Depivot** function, all of the sales numbers are in a single column. In this format, [Group By functions](#) can be useful and can provide insight into your data that may otherwise go unnoticed.

Group by aggregate functions

Data Prep aggregate functions let you combine sets of rows and compute a specific function on the referenced column.

These operations are called aggregate because they find matching rows in the dataset and then combine them into one row. A matching row is defined as one that, excluding the reference column, shares the same values in a column-by-column examination. The reference column is excluded from the column-by-column examination because its value are submitted to the aggregate function in order to produce the reference column value in the single-row result.

The following sections describe the aggregate functions available when you use **Group By** shaping operations:

FUNCTION	PURPOSE
array	Condenses duplicate rows into a single row of data and assembles the reference column data into a single, comma-separated string.
average	Condenses duplicate rows into a single row of data and displays the average of the numbers in the reference column.
count	Condenses duplicate rows into a single row of data and displays the number of duplicate rows in the reference column.
count (numbers only)	Condenses duplicate rows into a single row of data and displays the number of duplicate rows but only counts numeric values.
count distinct	Condenses duplicate rows into a single row of data and displays the number of unique values in the reference column.
first	Condenses duplicate rows into a single row of data and displays the first value that appeared for the duplicate rows.
last	Condenses duplicate rows into a single row of data and displays the last value that appeared for the duplicate rows.
max	Condenses duplicate rows into a single row of data and displays the largest of the numbers in the reference column.
min	Condenses duplicate rows into a single row of data and displays the smallest of the numbers in the reference column.

FUNCTION	PURPOSE
median	Condenses duplicate rows into a single row of data and displays the median of the numbers in the reference column.
stdev	Condenses duplicate rows into a single row of data and displays the standard deviation of the numeric values in the reference column.
stdevp	Condenses duplicate rows into a single row of data and displays the standard deviation for a population contained within the reference column.
sum	Condenses duplicate rows into a single row of data and displays the sum of the numbers in the reference column.
var	Condenses duplicate rows into a single row of data and displays the variance of the numeric values in the reference column.
varp	Condenses duplicate rows into a single row of data and displays the variance of a population contained within the reference column.

array

Of the available aggregate functions available, **array** is unique in that it operates on both text and numeric values. Rather than perform a mathematical operation on collapsed rows, all values in the reference column (the column to which **array** is applied) are temporarily stored. As the unique single row is created, **array** assembles the reference column data from the set into a single, comma-separated string within the column.

Excluding the column to which the **array** is being applied, all rows are examined column-by-column in order to find identical rows. The header of the reference column receives a name change to become *Array of*.

Examples

The small dataset below will be used to show how **array** operates.

COLUMN A	COLUMN B	COLUMN C
1	two	5
1	two	6
two	two	7
1	two	4

Example 1

Applying the **array** function to *Column C* reduces the row count from four to two. The value in the column *Array of Column C* shows the sum of of the *Column C* values in the duplicate rows that were collapsed during the operation.

COLUMN A	COLUMN B	COLUMN C
1	two	5.0, 6.0, 4.0
two	two	7.0

Example 2

Applying the **array** function to *Column A* results in no row count reduction since examination of *Column B* and *Column C* values in each row reveals that each is already unique. The value in the column *Array of Column A* therefore displays each row with its original value—however, numbers have been converted into text.

COLUMN A	COLUMN B	COLUMN C
1.0	two	5
two	two	7
1.0	two	6
1.0	two	4

average

The **average** aggregate function finds an average of the numbers in the reference column (the column to which it is applied) as all duplicate rows are collapsed into unique single rows. Excluding the column to which **average** is being applied, all rows are examined column-by-column in order to find identical rows.

A mathematical average (also called an “arithmetic mean” or simply “mean”) is calculated by adding all of the numeric values in a set and then dividing the resulting sum by the number of items the set contained. Note that attempts to apply **average** to a text value in the reference column will result in a 0 for that row. The header of the reference column receives a name change to become *Average of* .

Related aggregate functions are [median](#) and [mode](#).

Examples

The small dataset below will be used to show how **average** operates. All numbers are numeric values (not text) in order to illustrate behavior of this aggregate function.

COLUMN A	COLUMN B	COLUMN C
1	two	5
1	two	6
two	two	7
1	two	4

Example 1

Applying the **average** function to *Column C* reduces the row count from four to two. The value in the column *Average of Column C* shows the average of of the *Column C* values in the duplicate rows that were collapsed during the operation: $(5 + 6 + 4) \div 3 = 5$ while $7 \div 1 = 7$.

COLUMN A	COLUMN B	COLUMN C
1	two	5
two	two	7

Example 2

Applying the average function to *Column A* results in no row count reduction since examination of *Column B* and *Column C* values in each row reveals that each is already unique. The value in the column *Average of Column A* therefore displays each row with its original numeric value (in this case, a number 1) since the **average** function found no duplicate rows that could be involved in a mathematical operation. The 0 replaces the text "two" because the **average** operation cannot be applied to text values.

COLUMN A	COLUMN B	COLUMN C
1	two	5
1	two	6
1	two	4
0	two	7

count

The **count** aggregate function returns the number of duplicate rows in a dataset. Excluding the column to which it is being applied (the reference column), all rows are examined column-by-column in order to find duplicate rows. Those rows that contain duplicate data are collapsed into a single unique row. The reference column receives a name change to become *Count of* and the number that appears in the column indicates the number of duplicate rows that have been collapsed.

Examples

The small dataset below will be used to show how **count** operates.

COLUMN A	COLUMN B	COLUMN C
one	two	5
one	two	6
two	two	7
one	two	4

Example 1

Applying the **count** function to *Column C* reduces the row count from four to two. The value in the column *Count of Column C* shows the count of (number of times) duplicate rows that were collapsed during the operation.

COLUMN A	COLUMN B	COLUMN C
one	two	3
two	two	1

Example 2

Applying the **count** function to *Column A* results in no row count reduction since examination of *Column B* and *Column C* values in each row reveals that each is already unique. The value in the column *Count of Column A* therefore shows a value of 1 for each of the four rows.

COLUMN A	COLUMN B	COLUMN C
1	two	5
1	two	6

COLUMN A	COLUMN B	COLUMN C
1	two	7
1	two	4

count (numbers only)

The **count (numbers only)** aggregate function operates exactly the same as the [count](#) function. However, **count (numbers only)** only counts numeric values and ignores text values during the counting process.

count distinct

The **count distinct** aggregate function returns the number of unique values in the column being counted, as opposed to the [count](#) function which counts all values.

first

The **first** aggregate function searches the data for duplicate rows (based on the columns chosen by the user) and condenses them into a single row of data. The **first** function will then display the first value within the data that appeared for the duplicate rows. Any values following the first value will be lost during this process.

Example

The small dataset below will be used to show how **first** operates.

COLUMN A	COLUMN B	COLUMN C
one	two	5
one	two	6
two	two	7
one	two	4

Applying the **first** function to *Column C* reduces the row count from four to two. The value in the column *First of Column C* shows the first value from *Column C* of the duplicate rows that were collapsed during the operation.

COLUMN A	COLUMN B	COLUMN C
one	two	5
two	two	7

last

The **last** aggregate function searches the data for duplicate rows (based on the columns chosen by the user) and condenses them into a single row of data. The **last** function will then display the last value within the data that appeared for the duplicate rows. Any values prior to the last value will be lost during this process.

Example

The small dataset below will be used to show how **last** operates.

COLUMN A	COLUMN B	COLUMN C
one	two	5
one	two	6
two	two	7
one	two	4

Applying the Last function to *Column C* reduces the row count from four to two. The value in the column *Last of Column C* shows the last value from *Column C* of the duplicate rows that were collapsed during the operation.

COLUMN A	COLUMN B	COLUMN C
one	two	4
two	two	7

max

The **max** aggregate function returns the largest of the numbers in the reference column (the column to which it is applied) as all duplicate rows are collapsed into unique single rows. Excluding the column to which **max** is being applied, all rows are examined column-by-column in order to find identical rows.

The counterpart to this function is **min** (minimum).

Example

The small dataset below will be used to show how **max** operates. All numbers are numeric values (not text) in order to illustrate behavior of this aggregate function.

COLUMN A	COLUMN B	COLUMN C
1	two	5
1	two	6
two	two	7
1	two	4

Applying the **max** function to *Column C* reduces the row count from four to two. The value in the column *Max of Column C* shows the maximum of of the *Column C* values in the duplicate rows that were collapsed during the operation.

In the dataset that is returned by the **max** function (shown below) number 6 in the first row resulted from the set of numbers {4, 5, 6}. Each of these numbers exists in this set because each was a member of an identical row when *Column A* and *Column B* were examined. (*Column C* was excluded from this examination because it is the reference column.) In this set of three numbers, 6 is the largest—therefore it became the value shown in the reference column.

The number 7 results from a set with a single number {7} because there were no duplicate rows that could contribute other numbers. Since 7 is both the minimum and maximum of the single-row set, the function returns a 7 for that row.

min

The **min** aggregate function returns the smallest of the numbers in the reference column (the column to which it is applied) as all duplicate rows are collapsed into unique single rows. Excluding the column to which the **min** is being applied, all rows are examined column-by-column in order to find identical rows.

The counterpart to this function is **max** (maximum).

Example

The small dataset below will be used to show how **min** operates. All numbers are numeric values (not text) in order to illustrate behavior of this aggregate function.

COLUMN A	COLUMN B	COLUMN C
1	two	5

COLUMN A	COLUMN B	COLUMN C
1	two	6
two	two	7
1	two	4

Applying the **min** function to *Column C* reduces the row count from four to two. The value in the column *Min of Column C* shows the minimum of of the *Column C* values in the duplicate rows that were collapsed during the operation.

In the dataset that is returned by the **min** function (shown below) number 4 in the first row resulted from the set of numbers {4, 5, 6}. Each of these numbers exists in this set because each was a member of an identical row when *Column A* and *Column B* were examined. (*Column C* was excluded from this examination because it is the reference column.) In this set of three numbers, 4 is the smallest—therefore it became the value shown in the reference column.

The number 7 results from a set with a single number {7} because there were no duplicate rows that could contribute other numbers. Since 7 is both the minimum and maximum of the single-row set, the function returns a 7 for that row.

COLUMN A	COLUMN B	COLUMN C
1	two	4
two	two	7

median

The **median** aggregate function finds the median of the numbers in the reference column (the column to which median is applied) as all duplicate rows are collapsed into unique single rows. Excluding the column to which **median** is being applied, all rows are examined column-by-column in order to find identical rows.

A median value is one in the middle of a range of numbers ordered from lowest to highest value. This means that half of the numbers are to the “right” of the value returned and half of the numbers are to the “left” of the value returned. Where there is an even set of numbers (i.e., no single number in the middle of the collection), the function calculates an average on the pair of numbers in the middle of the range (i.e., the two numbers on either side of the mid-point).

Note that attempts to apply **median** to a text value in the reference column will result in a error for that row. The header of the reference column receives a name change to become *Median of* .

Related aggregate functions are [average](#) and [mode](#).

Example

The small dataset below will be used to show how **median** operates. All numbers are numeric values (not text) in order to illustrate behavior of this aggregate function.

COLUMN A	COLUMN B	COLUMN C
one	two	5
one	two	6
two	two	7
one	two	4

Applying the **median** function to *Column C* reduces the row count from four to two. The value in the column *Median of Column C* shows the median of the *Column C* values in the duplicate rows that were collapsed during the operation.

In the dataset that is returned by the **median** function (shown below) number 5 in the first row resulted from the ordered set of numbers {4, 5, 6}. Each of these numbers exists in this set because each was a member of an identical row when *Column A* and *Column B* were examined. (*Column C* was excluded from this examination because it is the reference column.) In this set of three numbers, 5 is the middle value with one number in the set on either side of it.

The number 7 results from a set with a single number {7} because there were no duplicate rows that could contribute other numbers. Since 7 is in the middle of the set (there are zero numbers on either side of it) the function returns a 7 for that row.

COLUMN A	COLUMN B	COLUMN C
one	two	5
two	two	7

stdev

The **stdev** (standard deviation) aggregate function computes the standard deviation—the amount of variation from the average—that exists within a sample set of data. This aggregate function calculates the standard deviation of the numeric values in the reference column (the column to which **stdev** is applied) among those rows that are identical.

All duplicate rows are discovered by conducting a column-by-column examination (excluding the column to which **stdev** is being applied) before they are collapsed into unique single rows. For each resulting single row, the value in the reference column from the contributing duplicate row becomes part of the standard deviation calculation. The header of the reference column receives a name change to become *Stdev of*.

If there are text values in the reference column, they are ignored within the **stdev** calculation. Note also that the aggregate **stdev** function requires at least two values. In other words, there must be at least two identical rows available for aggregation for each unique single row that is returned. Rows that occur only once will contribute only one value in the reference column for the calculation—resulting in an error.

The standard deviation for data is the square root of its variance. If the set under analysis represents all data points (referred to as a “population”), use of `stdevp` is recommended for a more accurate result. A related function that deals with statistical variance is `varp`.

Example

The dataset below will be used to show how `stdev` operates. All numbers are numeric values (not text) in order to illustrate behavior of this aggregate function.

COLUMN A	COLUMN B	COLUMN C
one	two	0.2
one	two	0.1
one	two	1.1
one	two	0.2
one	two	0.6
one	one	0.2
one	one	0.27
one	two	0.2
one	two	0.4

As shown in the table below, applying the `stdev` function to *Column C* reduces the row count from nine to two. The value in the column *Stdev of Column C* shows the standard deviation of the *Column C* sample data values in the duplicate rows that were collapsed during the operation.

COLUMN A	COLUMN B	COLUMN C
one	two	0.3511884584284246
one	one	0.049497474683058325

stdevp

The `stdevp` (standard deviation for a population) aggregate function computes the standard deviation—the amount of variation from the average—that exists within an entire set of data (a population). This aggregate function calculates the

standard deviation for a population using the numeric values in the reference column (the column to which **stdevp** is applied) among those rows that are identical.

All duplicate rows are discovered by conducting a column-by-column examination (excluding the column to which **stdevp** is being applied) before they are collapsed into unique single rows. For each resulting single row, the value in the reference column from the contributing duplicate row becomes part of the standard deviation calculation for the population. The header of the reference column receives a name change to become *StdevP of* .

If there are text values in the reference column, they are ignored within the **stdevp** calculation. Note also that the aggregate **stdevp** function requires at least two values. In other words, there must be at least two identical rows available for aggregation for each unique single row that is returned. Rows that occur only once will contribute only one value in the reference column for the calculation—resulting in an error.

If the set under analysis represents only a sample of data, use of [stdev](#) is recommended for a more accurate result. Other related functions that deal with statistical variance are [var](#) and [varp](#).

Example

The dataset below will be used to show how **stdevp** operates. All numbers are numeric values (not text) in order to illustrate behavior of this aggregate function.

COLUMN A	COLUMN B	COLUMN C
one	two	0.2
one	two	0.1
one	two	1.1
one	two	0.2
one	two	0.6
one	one	0.2
one	one	0.27
one	two	0.2
one	two	0.4

As shown in the table below, applying the **stdevp** function to *Column C* reduces the row count from nine to two. The value in the column *StdevP of Column C* shows the standard deviation of the population contained in *Column C* among the duplicate rows that were collapsed during the operation.

COLUMN A	COLUMN B	COLUMN C
one	two	0.32513733362117264
one	one	0.034999999999999996

sum

The **sum** aggregate function executes an addition operation on the numbers in the reference column (the column to which it is applied) as all duplicate rows are collapsed into unique single rows. Excluding the column to which **sum** is being applied, all rows are examined column-by-column in order to find identical rows. Note that attempts to apply **sum** to a text value in the reference column will result in a 0 for that row. The header of the reference column receives a name change to become *Sum of* .

Examples

The small dataset below will be used to show how **sum** operates. All numbers are numeric values (not text) in order to illustrate behavior of this aggregate function.

COLUMN A	COLUMN B	COLUMN C
1	two	5
1	two	6
two	two	7
1	two	4

Example 1

Applying the **sum** function to *Column C* reduces the row count from four to two. The value in the column *Sum of Column C* shows the sum of the *Column C* values in the duplicate rows that were collapsed during the operation: $5 + 6 + 4 = 15$ while $7 + 0 = 7$.

COLUMN A	COLUMN B	COLUMN C
1	two	15
two	two	7

Example 2

Applying the **sum** function to *Column A* results in no row count reduction since examination of *Column B* and *Column C* values in each row reveals that each is already unique. The value in the column *Sum of Column A* therefore displays each row with its original numeric value (in this case, a number 1) since the **sum** function found no duplicate rows that could be involved in an adding operation. The 0 replaces the text “two” because the **sum** operation cannot be applied to text values.

COLUMN A	COLUMN B	COLUMN C
1	two	5
1	two	6
1	two	4
0	two	7

var

The **var** (variance) aggregate function estimates dispersion—how much the values are spread out—within a sample set of data. This aggregate function calculates the variance of the numeric values in the reference column (the column to which it is applied) among those rows that are identical.

All duplicate rows are discovered by conducting a column-by-column examination (excluding the column to which **var** is being applied) before they are collapsed into unique single rows. For each resulting single row, the value in the reference column from the contributing duplicate row becomes part of the variance calculation. The header of the reference column receives a name change to become *Var of*.

If there are text values in the reference column, they will be ignored within the **var** calculation. Note also that the aggregate **var** function requires at least two values. In other words, there must be at least two identical rows available for aggregation for each unique single row that is returned. Rows that occur only once will contribute only one value in the reference column for the calculation—resulting in an error.

If the set under analysis represents all data points (referred to as a “population”) use of **varp** is recommended for a more accurate result. Related functions that deal with statistical variance are **stdev** and **stdevp**.

Example

The dataset below will be used to show how **var** operates. All numbers are numeric values (not text) in order to illustrate behavior of this aggregate function.

COLUMN A	COLUMN B	COLUMN C
one	two	0.2
one	two	0.1

COLUMN A	COLUMN B	COLUMN C
one	two	1.1
one	two	0.2
one	two	0.6
one	one	0.2
one	one	0.27
one	two	0.2
one	two	0.4

As shown in the table below, applying the **var** function to *Column C* reduces the row count from nine to two. The value in the column *Var of Column C* shows the variance of the *Column C* sample data values in the duplicate rows that were collapsed during the operation.

COLUMN A	COLUMN B	COLUMN C
one	two	0.12333333333333334
one	one	0.00245

varp

The **varp** (variance for a population) aggregate function computes dispersion—how much the values are spread out—for an entire set of data (a population). This aggregate function calculates the variance for a population contained within the reference column (the column to which it is applied) among those rows that are identical.

All duplicate rows are discovered by conducting a column-by-column examination (excluding the column to which **varp** is being applied) before they are collapsed into unique single rows. For each resulting single row, the value in the reference column from the contributing duplicate row becomes part of the variance calculation for the population. The header of the reference column receives a name change to become *VarP of*.

If there are text values in the reference column, they will be ignored within the **varp** calculation. Note also that the aggregate **varp** function requires at least two values. In other words, there must be at least two identical rows available for aggregation for each unique single row that is returned. Rows that occur only once will contribute only one value in the reference column for the calculation—resulting in an error.

If the set under analysis represents only a sample set of data, use of **var** is recommended for a more accurate result. Related functions that deal with statistical variance are [stdev](#) and [stdevp](#).

Example

The dataset below will be used to show how **varp** operates. All numbers are numeric values (not text) in order to illustrate behavior of this aggregate function.

COLUMN A	COLUMN B	COLUMN C
one	two	0.2
one	two	0.1
one	two	1.1
one	two	0.2
one	two	0.6
one	one	0.2
one	one	0.27
one	two	0.2
one	two	0.4

As shown in the table below, applying the **varp** function to *Column C* reduces the row count from nine to two. The value in the column *VarP of Column C* shows the variance of the population contained in *Column C* among the duplicate rows that were collapsed during the operation.

COLUMN A	COLUMN B	COLUMN C
one	two	0.10571428571428572
one	one	0.001225

Autonumber rows

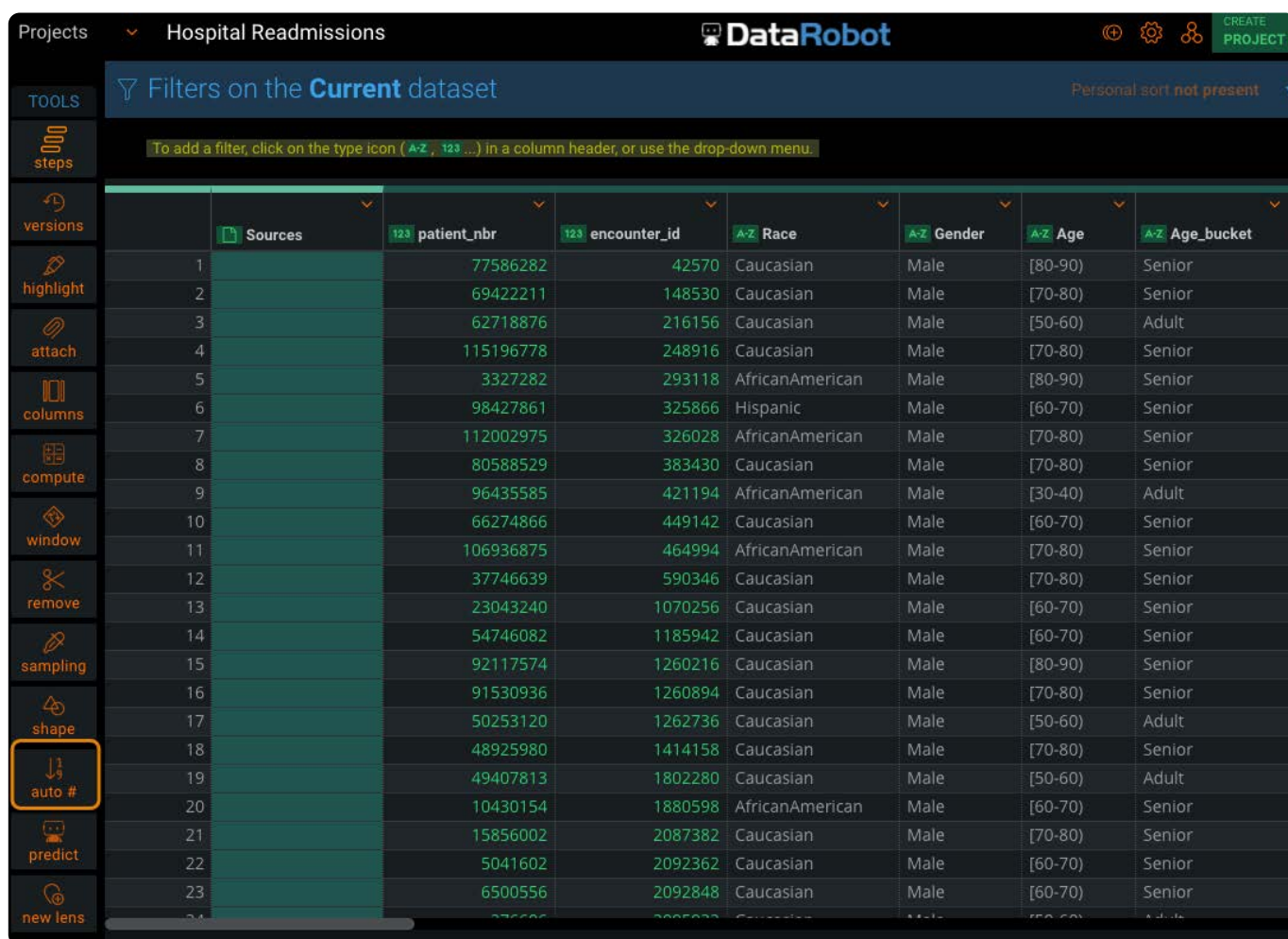
The Data Prep **Auto #** tool lets you assign each row a number. This tool is helpful if you need to give each row a unique identifier. The **Auto #** tool creates a new column that autonumbers the rows.

The **Auto #** tool gives each row a unique identifier so it is useful when you need to:

- Track your dataset's original order.
- Assign row identifiers to your dataset.

Work with the Auto # tool

To access the **Auto #** tool, click **auto #** in the project **Tools** bar:

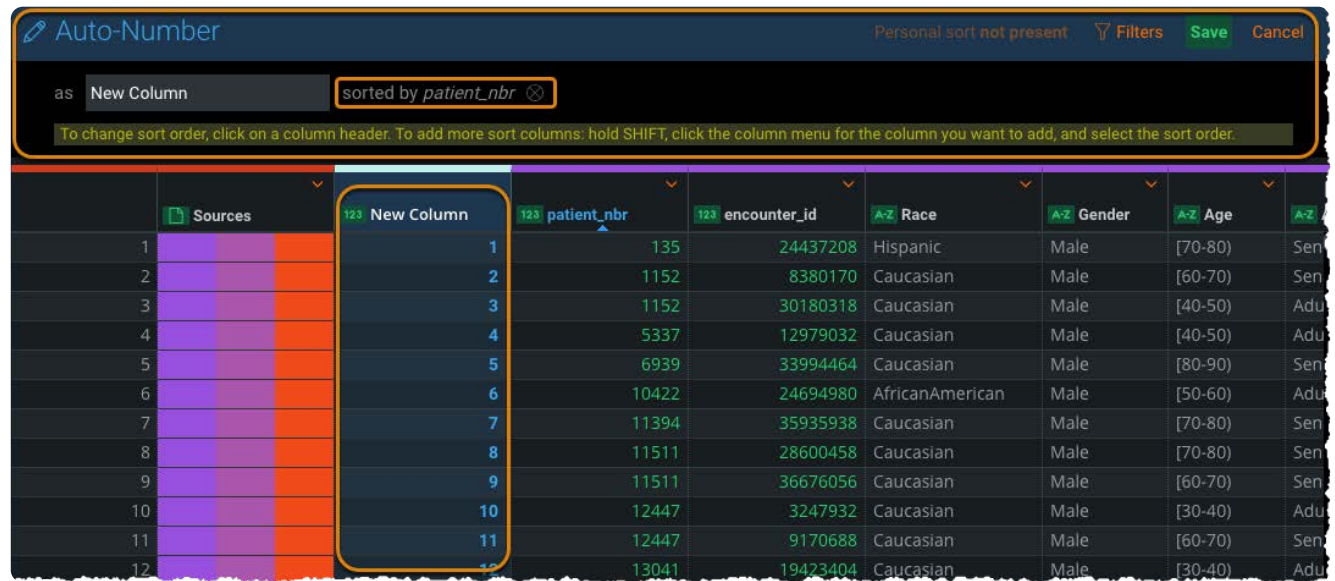


Add an autonumbered column

To add an autonumbered column:

1. From the **Tools** bar, click **auto #**.

The **Auto-Number** pane appears.



2. Enter a name for the autonumbered column you are adding.

The default name is "New Column(#)" where # starts as 1 and increments for each new autonumbered column.

3. To bind your autonumbered column to the sort order of an existing column in your dataset, click the header of the column you want to select.

The column is listed in the **sorted by** field.

4. To add more sort columns, hold SHIFT and, from the column you want to add, hover over the column menu and click **Sort by ascending** or **Sort by descending**.

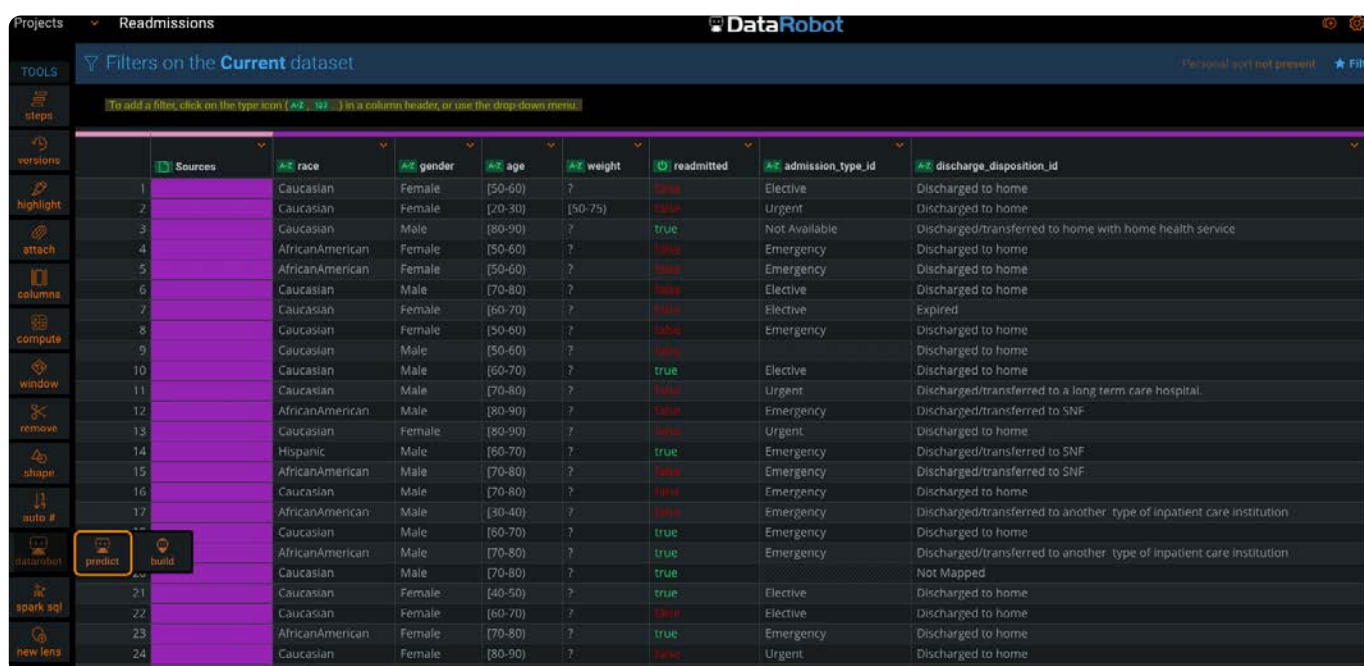
5. To remove columns you have selected for sorting, click the X icon adjacent to the column name in the **Auto-Number** pane.

Make predictions in Data Prep

When you have data that needs to be scored against a deployed Machine Learning (ML) model in DataRobot, the Data Prep **Predict** tool is how you generate the score.

Work with the Predict tool

To access the **Predict** tool, click the DataRobot icon in the **Tools** bar and select **predict**:



To generate the score, provide your DataRobot API token which is used to retrieve a list of your DataRobot deployments.

Projects Diabetes Readmission Paxata

TOOLS Predict values as New Column

DEPLOYMENTS OPTIONS

User token

Provide API token

Sources	patient_nbr	encounter_id	race	gender	age
1	86047875	64410	Caucasian	Male	[80-90]
2	63555939	15738	Caucasian	Male	[60-70]
3	77586282	42570	Hispanic	Male	[60-70]
4	108662661	84222	Caucasian	Male	[60-70]
5	107389323	89682	AfricanAmerican	Male	[70-80]
6	69422211	148530	Caucasian	Male	[70-80]
7	98427861	325866	Caucasian	Male	[60-70]
8	96435585	421194	Caucasian	Male	[50-60]
9	37746639	590346	Caucasian	Male	[60-70]
10	113848434	604188	Caucasian	Male	[70-80]
11	93232917	630342	Caucasian	Male	[70-80]
12	24370299	685086	Caucasian	Male	[70-80]
13	60679647	927786	Caucasian	Male	[50-60]
14	54746082	1185942	Asian	Male	[80-90]
15	92117574	1260216	AfricanAmerican	Male	[70-80]

predict tool

Note

To obtain your token, navigate to **User Settings > Developer Tools > API Keys**.

Next, select the deployment. Your data is scored against the model in this deployment. If the model used for scoring is a Time Series model, you must indicate this by checking the **Time Series Model** checkbox. Then, in the **Options** tab, specify the **Forecast Point** and, optionally, the **Series Id**. See [Options](#) for details.

Note

Deployments for custom models are not currently supported.

By default, the new column for the prediction score is created as "Target" in the dataset. To change this name, click the **Options** tab and provide a different name in the **Prediction Column** field.

After you select the deployment, the prediction runs. The new column is created and provides the prediction score. In addition, the "Target Prediction Value" column is also generated to provide the associated prediction value for each score. For multiclass predictions, the prediction values are returned per classification. For example, if classifying images into "apple", "orange" or "pear", then three additional columns are returned—one value for each corresponding score.

Examples of use case prediction values

- Predict the probability that a hospital patient may be readmitted after discharge. The prediction column will contain a binary value of 1 or 0 to indicate if the patient is likely to be readmitted or not readmitted.
- Classify a set of images into one of three fruits: oranges, pears, or apples. The prediction column will contain one of three values: orange, pear, apple.
- Forecast sales based on forecast dates. The prediction column in this case will contain the sales dollar amount.

For binary and time series prediction deployments, the **Options** tab provides additional options. See [Options](#) for details.

Options

For Times Series predictions, you must also provide the forecast point, which is the point you are making a prediction from—a relative time “if it was now...” DataRobot trains models using all potential forecast points in the training data. In production, it is typically the most recent time.

Important

The format of this date must be ISO 2014-08-12T00:00:00Z.

Optionally, if your dataset has multiseried data, for example a dataset that contains multiple time series to forecast the sales for multiple stores, then you can specify a column as a Series Id to group the data and return the predictions separately for each group.

For binary predictions, the **Options** tab provides prediction explanations that help you to understand why a prediction was returned, for example, “Why did this patient score a 1 for possibility of readmission?” or “Why was this image identified as an apple?”

When **Explanations** is enabled, five new columns are generated per explanation in the project:

- **feature**: The name of the feature contributing to the prediction.
- **feature value**: The value the feature took on for the row.
- **strength**: The amount this feature's value affected the prediction.
- **qualitative**: A human-readable description of how strongly the feature affected the prediction. For example: ++++; -; +
- **label**: Describes what output was driven by this prediction explanation. For regression projects, it is the name of the target feature. For classification projects, it is the class whose probability, if increased, would correspond to a positive strength of this prediction explanation.

Additionally, Low and High Threshold values can be set so that explanations are only generated for scores outside of the threshold.

See [Prediction Explanations](#) for complete details on values returned for predictions.

Transform your data with Spark SQL

Data Prep provides a tool to let you transform your data using Spark SQL. Spark SQL provides a [library of functions](#) for you to prep, clean, and transform your data.

Note

Your Data Prep Administrator must enable this feature in your application.

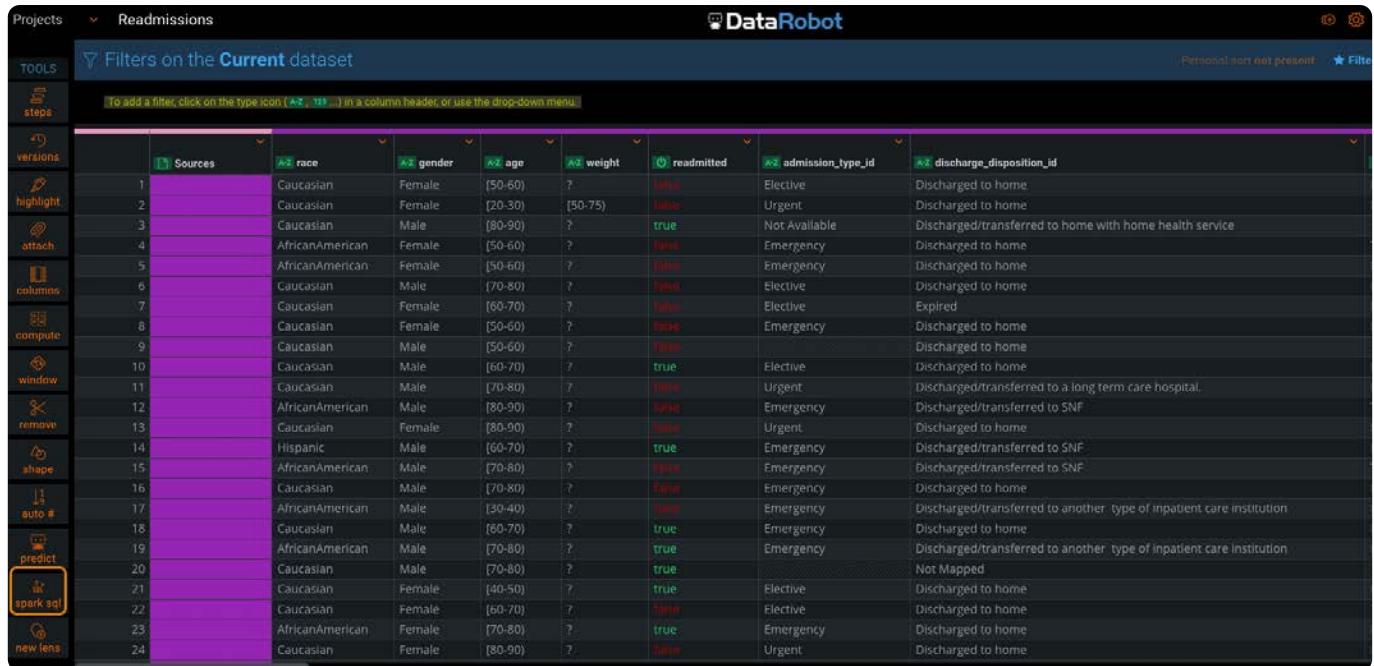
The following sections describe how to use the Spark SQL tool in Data Prep. For a discussion of the SQL statements supported, see [Data Prep Spark SQL Guidelines](#).

Tip

You can also select and transform your data using the AI Catalog. See [Prepare data in AI Catalog with Spark SQL](#).

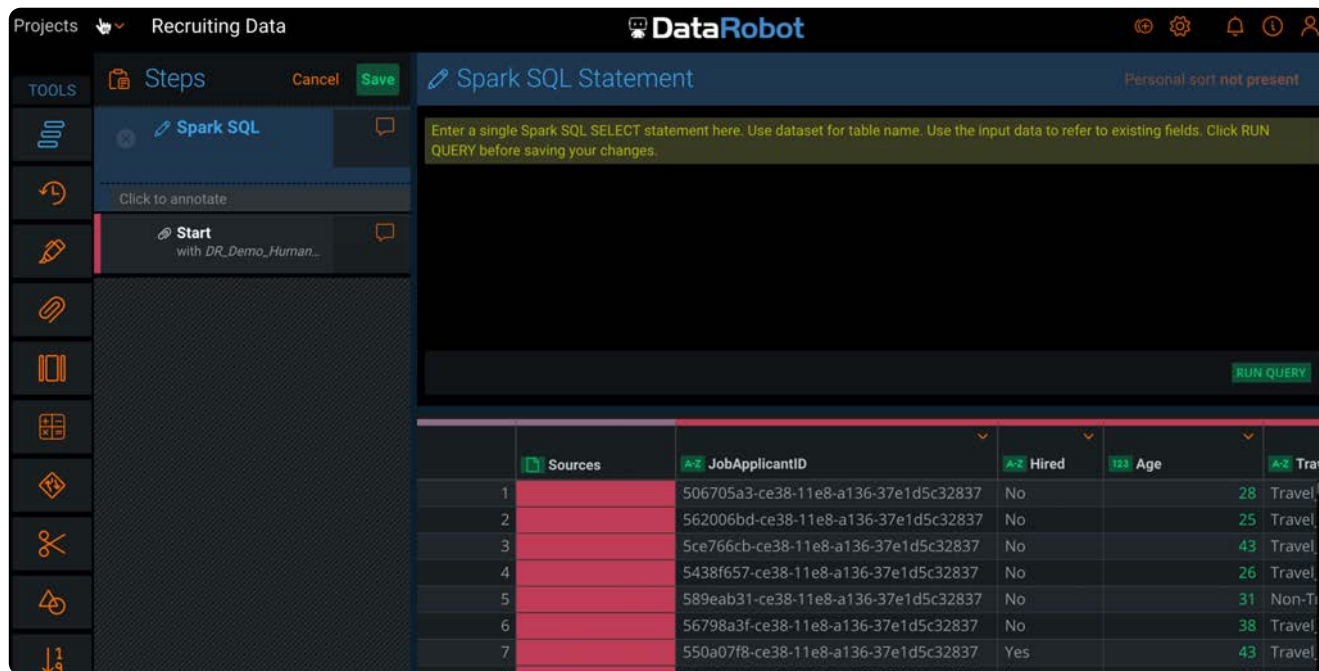
Work with Spark SQL

To access the **Spark SQL** tool, click **spark sql** in the project **Tools** bar:



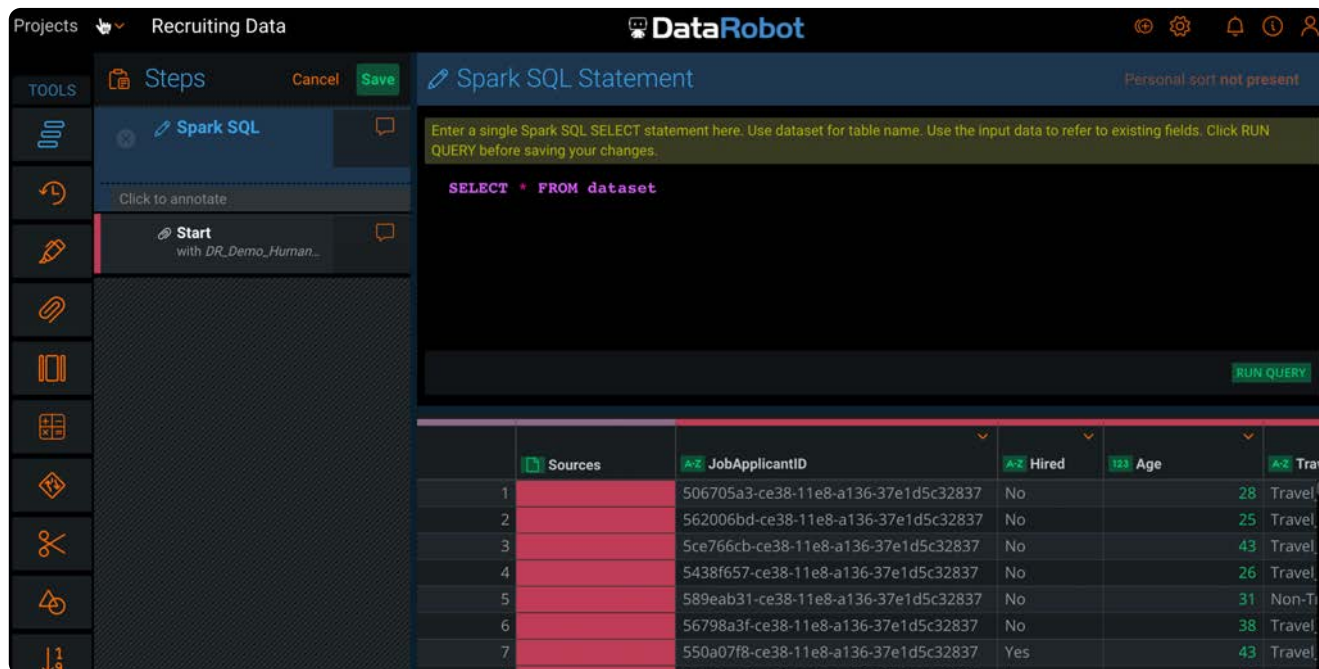
Add a Spark SQL statement

1. From the **Tools** bar, click **spark sql**. The **Spark SQL Statement** pane appears.



	Sources	JobApplicantID	Hired	Age	Travel
1		506705a3-ce38-11e8-a136-37e1d5c32837	No	28	Travel
2		562006bd-ce38-11e8-a136-37e1d5c32837	No	25	Travel
3		5ce766cb-ce38-11e8-a136-37e1d5c32837	No	43	Travel
4		5438f657-ce38-11e8-a136-37e1d5c32837	No	26	Travel
5		589eab31-ce38-11e8-a136-37e1d5c32837	No	31	Non-Ti
6		56798a3f-ce38-11e8-a136-37e1d5c32837	No	38	Travel
7		550a07f8-ce38-11e8-a136-37e1d5c32837	Yes	43	Travel

2. Enter a Spark SQL statement. See the [Data Prep Spark SQL guidelines](#) for usage details.



	Sources	JobApplicantID	Hired	Age	Travel
1		506705a3-ce38-11e8-a136-37e1d5c32837	No	28	Travel
2		562006bd-ce38-11e8-a136-37e1d5c32837	No	25	Travel
3		5ce766cb-ce38-11e8-a136-37e1d5c32837	No	43	Travel
4		5438f657-ce38-11e8-a136-37e1d5c32837	No	26	Travel
5		589eab31-ce38-11e8-a136-37e1d5c32837	No	31	Non-Ti
6		56798a3f-ce38-11e8-a136-37e1d5c32837	No	38	Travel
7		550a07f8-ce38-11e8-a136-37e1d5c32837	Yes	43	Travel


3. Click **Run Query** on the lower right of the **Spark SQL Statement** pane to validate your query. If the query is successful, the results display below. View the results to ensure that the query is functioning as expected.

If the query is unsuccessful, an error message displays below the query:

The screenshot shows the DataRobot interface with a 'Spark SQL' step in the 'Steps' tool. The step has an error icon (a red triangle with an exclamation mark) and a message: 'SQL statement is not valid'. The error details are: 'mismatched input 'ELECT' expecting '(''. The SQL statement entered is 'ELECT * FROM dataset'. Below the error, a table of data is displayed with the following columns: Sources, JobApplicantID, Hired, Age, and Travel. The table contains 7 rows of data.

	Sources	JobApplicantID	Hired	Age	Travel
1		506705a3-ce38-11e8-a136-37e1d5c32837	No	28	Travel
2		562006bd-ce38-11e8-a136-37e1d5c32837	No	25	Travel
3		5ce766cb-ce38-11e8-a136-37e1d5c32837	No	43	Travel
4		5438f657-ce38-11e8-a136-37e1d5c32837	No	26	Travel
5		589eab31-ce38-11e8-a136-37e1d5c32837	No	31	Non-T
6		56798a3f-ce38-11e8-a136-37e1d5c32837	No	38	Travel
7		550a07f8-ce38-11e8-a136-37e1d5c32837	Yes	43	Travel

4. Click **Save** to save the query. You can save queries with errors and return later to resolve the errors.

If the SQL query contains an error, the Spark SQL step in the Steps tool displays an error icon () . Click the icon to view the error message.

Note

After saving the Spark SQL step in the Steps tool, you might need to make changes to a previous step or add a new step before the Spark SQL step. In this case, click the Spark SQL step to edit it, click **Run Query**, and save the query again.

Data Prep Spark SQL Guidelines

Structured Query Language (SQL) is a declarative language designed for managing data stored in relational databases. Spark SQL is a component of Spark that allows queries to be written against registered DataFrames (data organized into named columns) in much the same way as SQL is used to query databases. Data Prep supports a [library of Spark SQL functions](#) for use with the Spark SQL tool.

Use case

The example dataset featured in this section contains a sample of job application data with the target feature, `Hired`.

JobApplicantID	Hired	Age	TravelPreference	HiringDepartment	DistanceFrom	EducationLevel	EducationField	Gender	Role	Internships	Over18	StandardHours	Summary
506705a3-ce38-11e8-	No	28	Travel_Frequently	Research & Develo	4	3	Engineering	Male	Research Scientist	9	Y	80	Hi there. I am ;
562006bd-ce38-11e8-	No	25	Travel_Rarely	Research & Develo	3	4	Medical	Male	Research Scientist	0	Y	80	Hi there. I am ;
5ce766cb-ce38-11e8-	No	43	Travel_Rarely	Human Resources	4	3	Life Sciences	Male	Human Resources	4	Y	80	Hi there. I am ;
5438f657-ce38-11e8-	No	26	Travel_Rarely	Sales	5	3	Other	Male	Sales Representativ	0	Y	80	Hi there. I am ;
589eab31-ce38-11e8-	No	31	Non-Travel	Sales	16	3	Life Sciences	Male	Sales Executive	1	Y	80	Hi there. I am ;
56798a3f-ce38-11e8-	No	38	Travel_Rarely	Research & Develo	25	2	Life Sciences	Male	Research Director	3	Y	80	I am a Life Scie
550a07f8-ce38-11e8-	Yes	43	Travel_Rarely	Research & Develo	8	4	Other	Male	Research Scientist	9	Y	80	I am a Other e
52447ab2-ce38-11e8-	Yes	25	Travel_Rarely	Research & Develo	9	3	Medical	Male	Research Scientist	1	Y	80	Recently I com
5b2194be-ce38-11e8-	No	34	Travel_Rarely	Research & Develo	1	4	Life Sciences	Male	Research Scientist	2	Y	80	Recently I com
599c1ce1-ce38-11e8-	No	27	Travel_Rarely	Research & Develo	1	2	Medical	Male	Laboratory Technicia	1	Y	80	Recently I com
585d8dd3-ce38-11e8-	No	32	Travel_Frequently	Research & Develo	2	2	Life Sciences	Male	Laboratory Technicia	0	Y	80	Hi there. I am ;
5cdb0cb8-ce38-11e8-	Yes	45	Travel_Rarely	Research & Develo	1	4	Engineering	Male	Healthcare Represe	1	Y	80	I am a Enginee
4fd933fa-ce38-11e8-	No	37	Travel_Rarely	Research & Develo	11	3	Medical	Male	Laboratory Technicia	4	Y	80	Recently I com
536044e5-ce38-11e8-	No	27	Non-Travel	Research & Develo	9	3	Medical	Male	Research Scientist	6	Y	80	Hi there. I am ;
50003dd0-ce38-11e8-	No	42	Travel_Rarely	Research & Develo	28	3	Life Sciences	Male	Research Director	3	Y	80	Hi there. I am ;
5b15fdd2-ce38-11e8-	Yes	34	Travel_Rarely	Research & Develo	29	3	Medical	Male	Laboratory Technicia	4	Y	80	Recently I com
5628df08-ce38-11e8-	No	30	Travel_Rarely	Research & Develo	6	3	Engineering	Male	Laboratory Technicia	0	Y	80	Hi there. I am ;
5ad72a0a-ce38-11e8-	No	39	Travel_Frequently	Sales	1	3	Marketing	Male	Sales Executive	0	Y	80	Recently I com
50844fc-ce38-11e8-	No	47	Travel_Frequently	Sales	27	2	Life Sciences	Male	Sales Executive	4	Y	80	Hi there. I am ;
5b384d75-ce38-11e8-	No	33	Travel_Rarely	Research & Develo	25	3	Engineering	Male	Manufacturing Direc	3	Y	80	Recently I com
515b5104-ce38-11e8-	No	42	Travel_Rarely	Research & Develo	24	3	Medical	Male	Manufacturing Direc	1	Y	80	Hi there. I am ;
50485ef4-ce38-11e8-	No	41	Travel_Rarely	Sales	9	3	Marketing	Male	Sales Executive	8	Y	80	Hi there. I am ;
5bf05cb9-ce38-11e8-	No	35	Travel_Frequently	Research & Develo	25	4	Life Sciences	Male	Research Scientist	1	Y	80	Hi there. I am ;
5088e020-ce38-11e8-	No	29	Non-Travel	Research & Develo	1	4	Medical	Male	Manufacturing Direc	0	Y	80	Hi there. I am ;
5387c3d4-ce38-11e8-	Yes	42	Travel_Frequently	Research & Develo	9	2	Medical	Male	Laboratory Technicia	1	Y	80	I am a Medical
50f59a74-ce38-11e8-	No	36	Travel_Rarely	Research & Develo	25	2	Life Sciences	Male	Research Director	3	Y	80	Hi there. I am ;
519bf949-ce38-11e8-	No	32	Travel_Rarely	Research & Develo	1	1	Life Sciences	Male	Research Scientist	1	Y	80	Recently I com
579be332-ce38-11e8-	No	27	Travel_Rarely	Human Resources	17	3	Other	Male	Human Resources	1	Y	80	Hi there. I am ;
55090e1-ce38-11e8-	No	50	Travel_Rarely	Research & Develo	2	3	Medical	Male	Research Director	5	Y	80	Hi there. I am ;
50f93f6-ce38-11e8-	No	33	Travel_Rarely	Research & Develo	4	4	Medical	Male	Laboratory Technicia	9	Y	80	I am a Medic

Query guidelines

Once a dataset is loaded into Data Prep, it is registered as a DataFrame and can now be queried using the `dataset` alias in a Spark SQL statement.

To use Spark SQL to shape your data, [enable the Spark SQL tool](#) and enter SQL queries in the **Spark SQL Statement** pane.

Data Prep allows only `SELECT` queries. Build the queries using the column names, for example:

```
SELECT
  EducationLevel,
  Hired
FROM dataset
```

See [Prohibited keywords and functions](#) for other restrictions.

Sample queries

The following are sample Spark SQL queries used in Data Prep, along with descriptions of the results.

Note

In Data Prep, SQL statements are case sensitive, adhering to the common SQL convention of using all caps for keywords and lowercase for variable names.

Example 1

```
SELECT * FROM dataset
```

The updated dataset includes all columns and all rows from the dataset.

Example 2

```
SELECT * FROM dataset
```

The result is: `ParseException`

Example 3

```
SELECT * from doesNotExist
```

The result is: `NoSuchTableException`

Example 4

```
SELECT Hired FROM dataset
```

The updated dataset includes only the `Hired` column (containing whether or not the application was successful) and all rows.

Example 5

```
SELECT  
  EducationLevel,  
  Hired  
FROM dataset  
WHERE EducationLevel = 5
```

The updated dataset includes only the `Hired` and `EducationLevel` columns and only those rows where the `EducationLevel` is 5.

Example 6

```
SELECT  
  EducationLevel,
```

```
CASE WHEN Hired = 'No' THEN 0 ELSE 1 END
AS HiredNum
FROM dataset
WHERE EducationLevel = 5
```

The updated dataset includes `EducationLevel` and a numeric version of the `Hired` column for rows where the `EducationLevel` is 5.

Example 7

```
SELECT
  EducationLevel,
  avg(CASE WHEN Hired = 'No' THEN 0 ELSE 1 END) AS acceptance_rate
FROM dataset
GROUP BY EducationLevel
ORDER BY EducationLevel
```

The updated dataset includes `EducationLevel` and the average acceptance rate in `EducationLevel` groups, ordered by `EducationLevel`.

Example 8

```
SELECT
  EducationLevel,
  avg(CASE WHEN Hired = 'No' THEN 0 ELSE 1 END) AS acceptance_rate,
  std(CASE WHEN Hired = 'No' THEN 0 ELSE 1 END) AS acceptance_rate_std
FROM dataset
GROUP BY EducationLevel
ORDER BY EducationLevel
```

The updated dataset includes `EducationLevel` and the average and standard deviation of the acceptance rate in `EducationLevel` groups, ordered by `EducationLevel`.

Example 9

```
SELECT
  EducationLevel,
  length(Summary) AS length_summary
FROM dataset
```

Here we surface the `EducationLevel` and the length of the cover letter that went along with the application.

Example 10

```
SELECT
  EducationLevel,
  avg(length(Summary)) AS avg_length_summary,
  std(length(Summary)) AS std_length_summary,
  std(CASE WHEN Hired = 'No' THEN 0 ELSE 1 END) AS acceptance_rate_std,
  avg(CASE WHEN Hired = 'No' THEN 0 ELSE 1 END) AS acceptance_rate
```

```
FROM dataset  
GROUP BY EducationLevel  
ORDER BY EducationLevel
```

This example puts it all together— the updated dataset includes `EducationLevel`, the average and standard deviation of the acceptance rate, and the average and standard deviation of the summary length in `EducationLevel` groups, ordered by `EducationLevel`.

Prohibited keywords and functions

Some SparkSQL 2.4.0 functions present a security risk. Data Prep prohibits the use of commands and methods that raise a potential security risk.

Click the tabs below to view lists of prohibited keywords and functions.

Keywords

Functions

ALTER
(DATABASE|
SCHEMA)
ALTER [TABLE |
VIEW]
ALTER VIEW
CREATE
(DATABASE|
SCHEMA)
CREATE
FUNCTION
CREATE TABLE
CREATE VIEW
DROP
DATABASE
DROP
FUNCTION
DROP [TABLE |
VIEW]
MSCK REPAIR
TABLE or ALTER
TABLE
RECOVER
PARTITIONS
TRUNCATE
TABLE
USE

LOAD DATA
EXPLAIN

ADD FILE
ADD JAR
ANALYZE TABLE
CACHE TABLE
CLEAR CACHE
DESCRIBE
DATABASE
DESCRIBE
FUNCTION
DESCRIBE
TABLE
LIST FILE
LIST JAR

REFRESH
REFRESH TABLE
RESET
SET
SHOW
COLUMNS
SHOW CREATE
TABLE
SHOW
(DATABASES)
SCHEMAS)
SHOW
FUNCTIONS
SHOW
PARTITIONS
SHOW TABLES
or SHOW TABLE
EXTENDED
SHOW
TBLPROPERTIES
UNCACHE
TABLE
DESCRIBE
CREATE TEMP
VIEW USING
CREATE TABLE

Considerations

INSERT [INTO |

OVERWRITE] The Data Prep Spark SQL tool supports the [commands](#) and [functions](#) available in [Spark SQL 2.4.0](#).

CREATE TABLE queries are allowed. All other SQL operations that might change the current dataset are prohibited (e.g., INSERT, UPDATE, DELETE).

OVERWRITE Built-in functions that could potentially hurt the system are also prohibited (e.g., `reflect`, `java_method`). See [Prohibited keywords and functions](#).

Each SQL statement is executed under its own Spark session, so there is no possibility for data leaking between SQL executions.

- For data type handling:
 - Infinity is not supported and is treated as null.
 - If a column has mixed data types, the value is converted to the specified type if convertible, otherwise, it is set to null.
 - When converting a SQL result into a Data Prep dataset, all data types from the SQL columns are converted to Data Prep data types. If the type is not supported, it is converted to String type.

Work with project tools > Create a DataRobot project using a DR lens

Create a DataRobot project using a DR lens

After you finish prepping your data, you can create a DataRobot project directly from Data Prep and start modeling in DataRobot. To do this, use the Data Prep **Build** tool.

Note

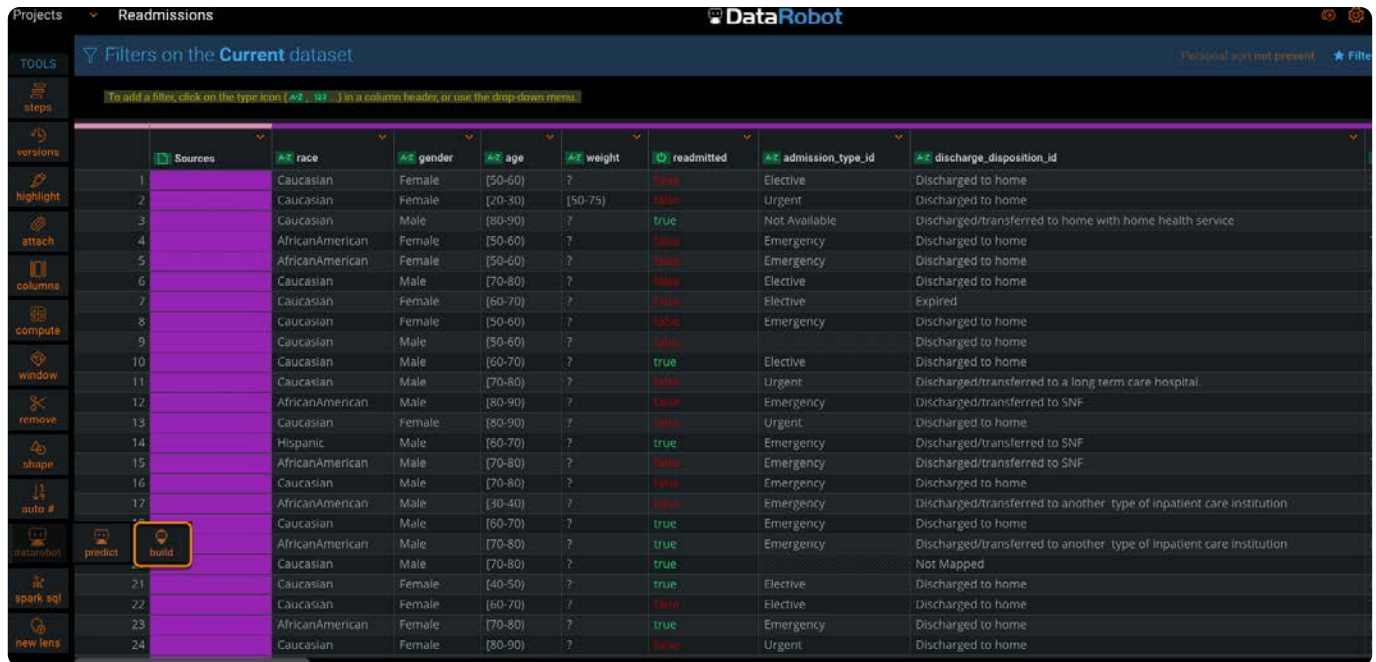
To create a DataRobot project from the Data Prep application, the DR lens feature must be enabled. If you do not see the DataRobot **build** tool in the **Tools** bar, contact your Data Prep System Administrator.

Work with the Build tool

Note

A [DataRobot connector](#) must be configured in order for you to use the **Build** tool to create a DataRobot project. The DataRobot connector must be updated to the latest version.

To access the **Build** tool, click the DataRobot icon in the **Tools** bar and select **build**:



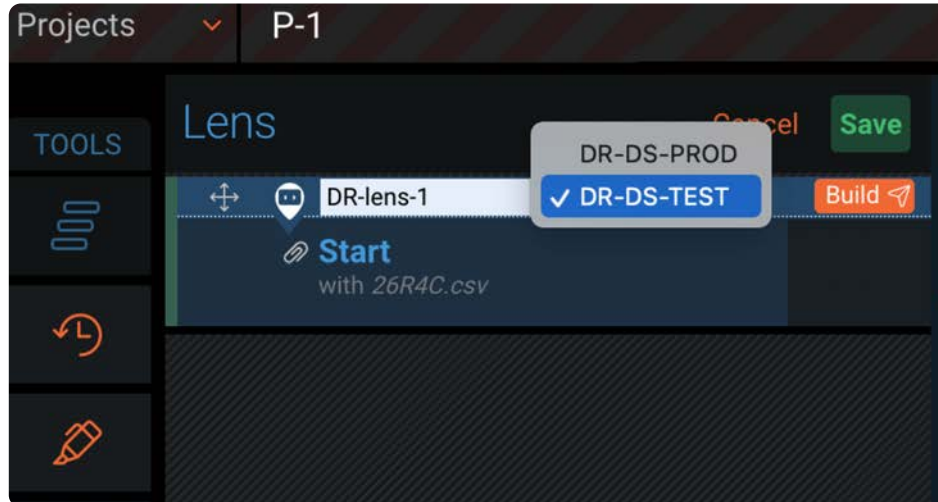
To create a DataRobot project, you first create a *DR lens*. A DR lens is similar to a standard *lens*—they both generate an [AnswerSet](#) based on the steps in the lens, but the DR lens creates a DataRobot project based on the AnswerSet.

Note

You cannot create an [Automatic Project Flow \(APF\)](#) from a DR lens. Your project must have a standard lens to create an APF.

Create a DataRobot project from a DR lens

1. Click the DataRobot icon in the **Tools** bar and select **build**.
2. Enter the name of the DR lens and select a DataRobot connector.



Note

If a DataRobot connector has not been configured, the connector dropdown list is empty. Contact your administrator to request that a DataRobot connector be configured.

3. Click **Build**. A notification displays at the top of the window indicating that the DataRobot project is being created.
4. Once the DataRobot project has been created, a success message displays at the top of the window with a "Click here" link. Click the link to access the machine learning project that has been created in DataRobot. The DataRobot project is named after the DR lens name you specified.

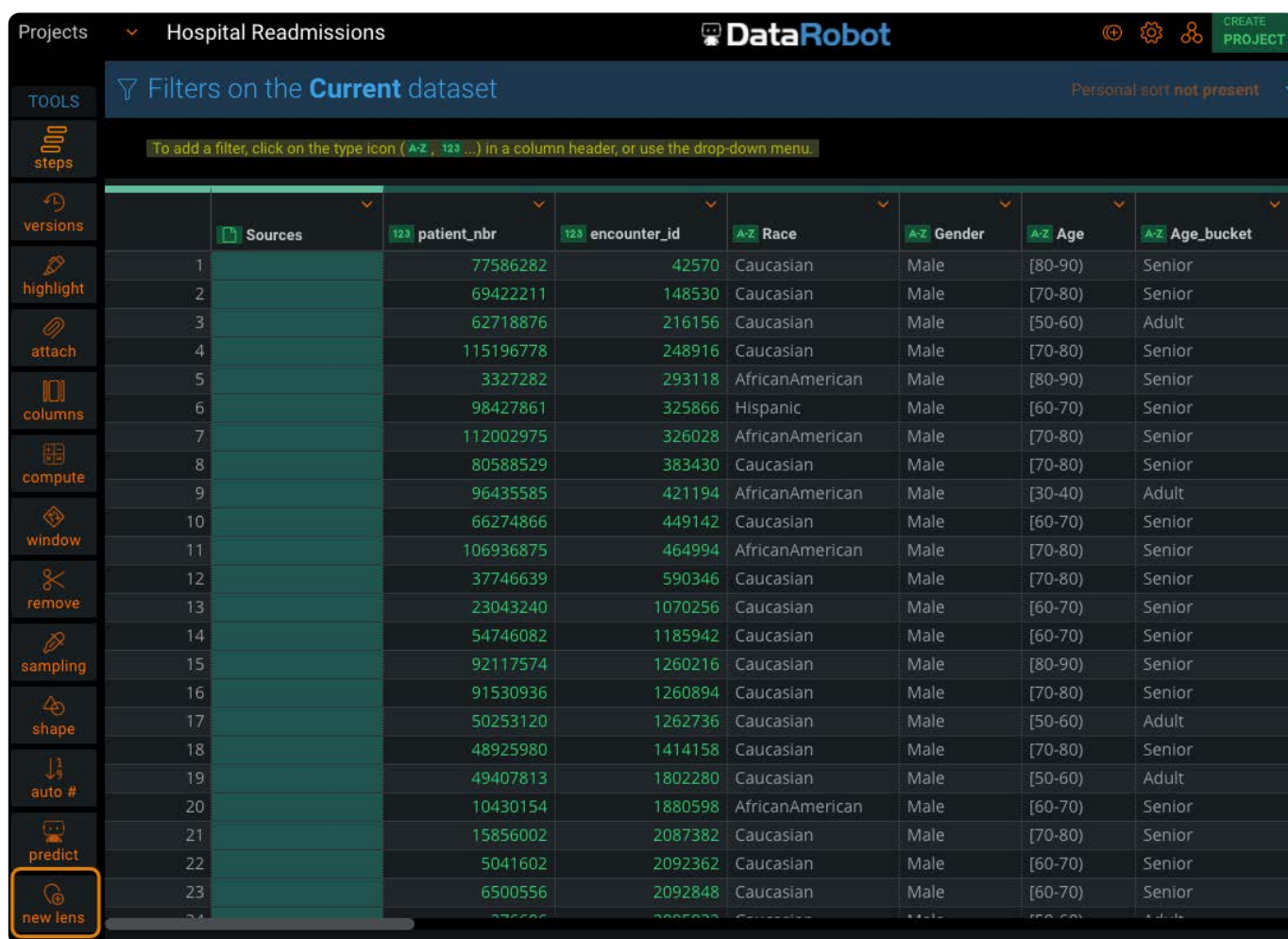
If the DataRobot project cannot be created, an error message displays. Review the details in the export log by selecting **Library > Export Logs**.
5. Click **Save** to save the DR lens step in the Steps tool.

Use lenses for publishing

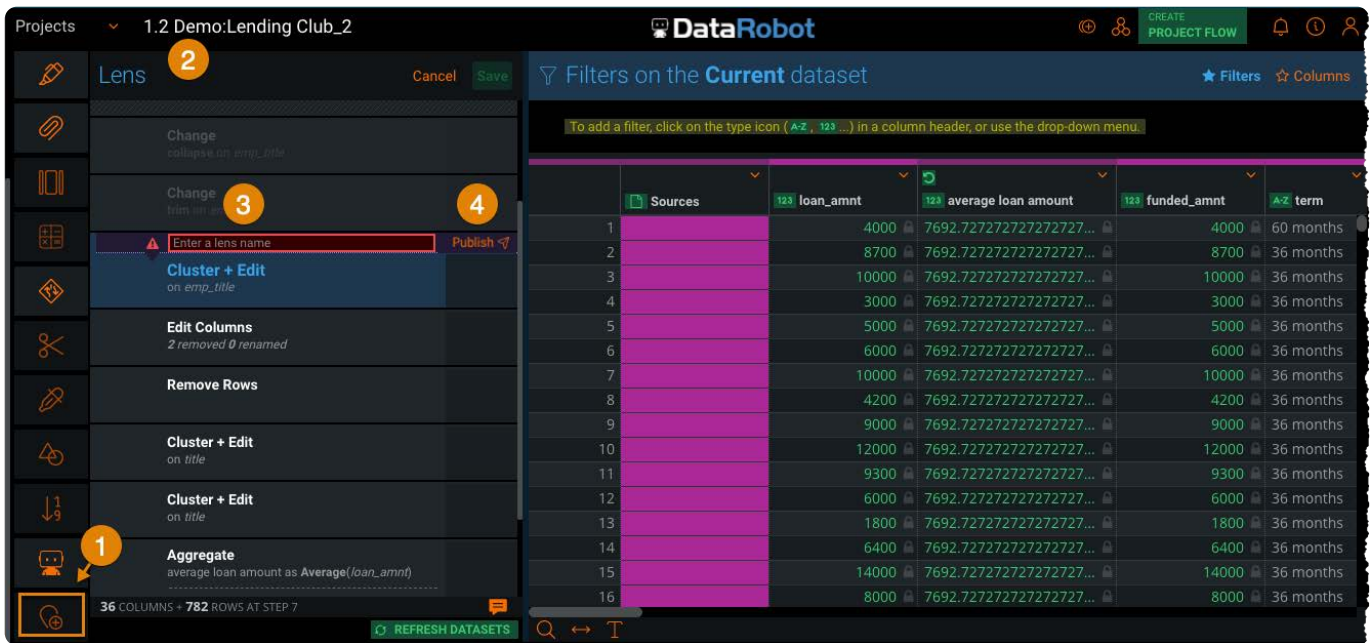
You use lenses to create publishing points from steps in your Data Prep project. When you publish from a lens, the resulting AnswerSet is a snapshot of your dataset at that particular step in the project. By default, the AnswerSet is saved to your data library.

Work with the New lens tool

To access the **New lens** tool, click **new lens** in the project **Tools** bar:



The following is an overview of the elements you work with when removing rows from your project:



ELEMENT	DESCRIPTION
---------	-------------

- 1
New lens tool
Click **step** in the **Tools** bar and select a step. Next, click **new lens** to access the **Lens** pane.
- 2
Lens pane
Lets you create a lens and publish it as an AnswerSet.
- 3
Enter a lens name
Enter the name of the lens and click **Save**.
- 4
Publish
After you set your lens, click **Publish** to save the state of your data to an **AnswerSet**.


Add a lens

To add a lens:

1. Click **steps** in the **Tools** bar and click the step where you want to add the lens.
2. Click **new lens** in the Tools bar.
3. In the **Lens** pane, enter a unique lens name and click **Save**.
4. Optionally, save to an Answerset by clicking **Publish**.

Tips for using lenses

- You can add a lens to any step or sub-step in your project, for example, to the Import step of an Append.
- You can drag an existing lens to any step or add an existing lens multiple times.

- All lenses persist as part of the project steps and are public to anyone who shares your project.
- A lens name must be unique because it is used to name the resulting AnswerSet.
- If there is an error in a formula created using the [Compute tool](#), an error icon () displays in the Steps tool. In this case, you can create and save a lens but you cannot publish it to an [AnswerSet](#).

The lenses you create are retained in project versions, and you can publish AnswerSets from lenses in older versions of your project.

Lenses are also essential for project automation because they define the publishing points to use for automated jobs. When you set up a project for automation, you select lenses and configure a corresponding schedule to automatically publish AnswerSets to your data library. Therefore, in order to automate a project, you require at least one lens in the project. For more help on automating a project, see [Automation and operationalization](#).

Following are examples of when to use lenses.

Isolate rows in your dataset

You can use a lens to isolate rows in your dataset that need further investigation. To do so, you add a lens on a step and you filter rows that you want to isolate from your current dataset. Name the lens and click **Publish**. The resulting AnswerSet is published to the data library and includes only the isolated rows you can investigate later. Now you can create a new step to remove those rows from your current dataset.

View before and after aggregation

To view your data before and after aggregation, you can add a lens to publish the current dataset prior to shaping your data. Name the lens and click **Publish**. The resulting AnswerSet of pre-aggregated data is published to the data library. Create a Shape step, then add a lens to publish the resulting dataset. You now have two AnswerSets that reflect your data before and after the aggregation.

Schedule a project for automation.

To schedule a project for automation, add a lens for every step in your project where you want to create a publishing point. Name each lens with a unique name to describe the output generated from that publishing point. Set up automation to use the lenses for publishing AnswerSets to the data library based on the schedule you configure. See [Automation and operationalization](#) for details.


Computed column functions

You can add columns to your Data Prep project by applying functions to existing columns. The following pages provide syntax and examples for the computed column functions used with the [Compute tool](#).

TOPIC	DESCRIBES ...
Date/Time functions	Apply date/time functions to existing columns to create new columns.
Informational functions	Query column data to check for particular values, for example, you can find the first nonblank column value or check for blanks and null values.
Logical functions	Evaluate logical functions on column values, such as AND, OR, and NOT.
Mathematical functions	Evaluate math functions on column values and store the new column values as numeric data types.
Statistical functions	Evaluate statistical functions on column values, such as AVERAGE, MAX, MIN, MEDIAN, MODE, and STDEV.
Text functions	Operate on text columns using functions like CHAR for converting from ASCII values to a character, CONCATENATE for combining strings, and FIND for finding a string inside of a string.
Comparison operators	Test logical operations using functions like =, >, >=, and <>.
Custom column functions	Apply column functions developed by your organization (on-premise installations only).

Date/time computed column functions

This section provides syntax and examples for the date/time computed column functions you can use with the Data Prep [Compute tool](#).

To use a date/time function, the value must be stored as a datetime data type. Datetime data types can be identified by the datetime data type icon  in the header row. If the value is not stored as a datetime data type, convert the value to a datetime data type using the `DATEVALUE` function. See the `DATEVALUE ()` section of this article.

DATE

Takes three separate arguments and combines them to form a date in a new DateTime column.

Syntax

```
DATE(YEAR, MONTH, DATE)
```

- YEAR is four-digit value
- MONTH is two-digit value
- DATE is two-digit value

Example

```
DATE(@year@, @month@, @day@)
```

year	month	day	New Column
1999	5	8	1999-5-08T00:00:00.000Z
1999	6	8	1999-6-08T00:00:00.000Z
1999	7	8	1999-7-08T00:00:00.000Z
1999	8	8	1999-8-08T00:00:00.000Z

Notes on use

Leading zeros for MONTH and DATE are not supported, for example:

```
DATE(1999,05,08) should be expressed as DATE(1999,5,8)
```

DATEADD

Calculates the date that is so many days, weeks, months from a given date.

Syntax

DATEADD(DATETIME, INCREMENT, INTERVAL)

- `DATETIME` is the date you want to start with.
- `INCREMENT` is the number you provide to be added to the `DATETIME`.
- `INTERVAL` is the interval (minutes, days, years, etc.) to add by. The following is a list of the recognized values for `INTERVAL`:
 - Years
 - Months
 - Weeks
 - Days
 - Hours
 - Minutes
 - Seconds
 - Millis

Example

DATEADD(@Date Received@, 6, "months")

Date Received	New Column
2015-08-24T06:36:33.000Z	2016-02-24T06:36:33.000Z
2011-09-08T07:38:59.000Z	2012-03-08T07:38:59.000Z
2012-09-03T07:13:18.000Z	2013-03-03T07:13:18.000Z

Notes on use

The `DATETIME` you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object. The `INCREMENT` provided must be an integer. Millis accepts a maximum of +/- 2147483647.

DATEDIFF

Calculates the days, weeks, months between two dates.

Syntax

DATEDIFF(DATETIME_1, DATETIME_2, INTERVAL)

- `DATETIME_1` is the date you want to start with.
- `DATETIME_2` is the date you want to end with.
- `INTERVAL` is the interval type (minutes, days, years, etc.) you want returned. The following is a list of the recognized values for the `INTERVAL` value:
 - Years
 - Months

- Weeks
- Days
- Hours
- Minutes
- Seconds
- Millis

Example

DATEDIFF(@Date Received@, @Date Shipped@, "months")

Date Received	Date Shipped	New Column
2015-08-24T06:36:33.000Z	2016-02-24T06:36:33.000Z	6
2011-09-08T07:38:59.000Z	2012-04-08T07:38:59.000Z	7
2012-09-03T07:13:18.000Z	2013-04-03T07:13:18.000Z	7

Notes on use

The `DATETIME` you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object. The `INCREMENT` provided must be an integer. Millis accepts a maximum of +/- 2147483647.

It is recommended you use the latest datetime value for the `DATETIME_2`. If you enter the earliest date as the `DATETIME_2` value, the `DATEDIFF` function will return a negative number.

`DATEDIFF` always rounds the result down to the nearest whole number. For example, if the difference between two dates is 3 years and 11 months, the `DATEDIFF` function returns the difference as 3 years.

DATEFORMAT

Converts a value stored as a datetime data type to a text sting in a given format.

Syntax

DATEFORMAT(DATETIME, FORMAT)

- `DATETIME` is the date you want to convert.
- `FORMAT` is format you want the `DATETIME` converted to.

Example

DATEFORMAT(@Date Received@, "dd-MMM-yyyy HH:mm")

Date Received	New Column
2015-08-24T06:36:33.000Z	August 24, 2015
2011-09-08T07:38:59.000Z	September 08, 2011
2012-09-03T07:13:18.000Z	September 03, 2012

Notes on use

The `DATETIME` you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

DATETRUNC

Removes the unwanted detail of a timestamp and rounds it to the interval you want. This provides the same output as the SQL `DATE_TRUNC()` function. Use case: you want to explore trends in your Community user signups and you need to aggregate signup event data by the time each event occurred. You're only interested in signups by year, month, or day but not the hour, minute, and millisecond. Use `DATETRUNC` to remove the portion of the timestamp that you don't need.

Syntax

`DATETRUNC(x)` where `x` can be any one of the following arguments.

- minutes
- month
- weeks
- days
- hours
- seconds

Example

`DATETRUNC(@DATE@, "months")`

DATE	New Column
2019-01-10T00:00:00.000Z	2019-01-01T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-01T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-01T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-01T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-01T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-01T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-01T00:00:00.000Z

DATEVALUE

Converts a datetime text string to a datetime object so you can use it for calculations.

Syntax

```
DATEVALUE(DATETIME, FORMAT, TIME_ZONE)
```

- `DATETIME` is the datetime as a text string.
- `FORMAT` is format of the `DATETIME`.
- `TIME_ZONE` is the time zone you want associated with the datetime object.

Example

```
DATEVALUE(@Date@, "yyyy-MMM-dd hh:mm a", "GMT-05:00")
```

Notes on use

Use the `DATEVALUE` function to convert a text column into a date column, or a date that you type into a date object. With the resulting data object, you can use Data Prep date functions, for example return the number of days or years between two dates (see next example below for Date Manipulation).

Date objects can store a date, time, or a combination of date and time.

To convert text to a Data Prep date object: Specify its format in Data Prep date format syntax. Repeated characters indicate the length of the field, such as `yyyy` means a 4-digit year.

For column `DateCol` with input text that specifies February 28, 2012 as:

```
2012/28/02
```

Convert to date object:

```
DATEVALUE(@DateCol@, "yyyy/dd/MM")
```

The date format must match your input data.

If February 28, 2012 looks like:

```
2012-15-02
```

Use date format:

```
"yyyy-dd-MM"
```

If February 28, 2012 looks like:

```
2-28-12
```

Use date format:

```
"dd-MM-yy"
```

If the time 1:29 pm looks like:

13:29

Use time format:

"HH:mm"

If the time 1:29 pm looks like:

01:29PM

Use time format:

"hh.mmaa"

Advanced example

If the input text is a date and time, separated by the letter T, then a time zone:

2012-02-28T09:29:00-05:00

For letter characters that literally appear in input text, surround the letter with single straight quotes. Use the following date format:

"yyyy-MM-dd'T'HH:mm:ssZZ"

Date manipulation

Use the `DATEDIFF` function to calculate differences in ts-date-time values between two Data Prep date objects. Calculate days between August 1, 1998 and a date column:

```
DATEDIFF(DATEVALUE("01-AUG-1998", "dd-MMM-yyyy"), @MyDate@, "days")
```

DAY

Extracts the day from a date.

Syntax

```
DAY(DATETIME)
```

`DATETIME` is the date you want to extract the day from.

Example

```
DAY(@Date@)
```

Date	123 New Column
2011-01-15T06:37:40.000Z	15
2011-01-23T07:09:58.000Z	23
2011-01-30T07:27:56.000Z	30

Notes on use

The `DATETIME` you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

Returned values range from 1 to 31.

DAYOFWEEK

Returns the day of the week from a date.

Syntax

```
DAYOFWEEK(DATETIME)
```

`DATETIME` is the date you want to evaluate.

Example

```
DAYOFWEEK(@Date@)
```

Date	123 New Column
2011-01-15T06:37:40.000Z	6
2011-01-23T07:09:58.000Z	7
2011-01-30T07:27:56.000Z	7

Notes on use

The `DATETIME` you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

Returned values range from 1 (Monday) to 7 (Sunday).

DAYOFYEAR

Returns the day of the year from a date.

Syntax

```
DAYOFYEAR(DATETIME)
```

DATETIME is the date you want to evaluate.

Example

DAYOFYEAR(@Date@)

Date	New Column
2011-01-15T06:37:40.000Z	15
2011-01-23T07:09:58.000Z	23
2011-01-30T07:27:56.000Z	30

Notes on use

The DATETIME you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

Returned values range 1 to 365 (366 on a leap year).

ENDOFMONTH

Returns the datetime for the last day of the month in a new DateTime column. This provides the same output as Excel's EOMONTH function.

Syntax

ENDOFMONTH(DATE_TIME)

DATE_TIME is a DateTime object.

Example

ENDOFMONTH(@Date@)

DATE	New Column
2019-01-10T00:00:00.000Z	2019-01-31T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-31T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-31T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-31T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-31T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-31T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-31T00:00:00.000Z

FROMUNIXTIME

Returns a ts-date-time object from a Unix timestamp. This provides the same output as the MySQL FROM_UNIXTIME() function.

Syntax

FROMUNIXTIME(MILLISECONDS)

MILLISECONDS is the int value represented as milliseconds.

Example

FROMUNIXTIME(@UNIX TIME STAMP@)

UNIX TIME STAMP	New Column
831877766	1970-01-10 15:04:37.766Z
834556166	1970-01-10 15:49:16.166Z
837148166	1970-01-10 16:32:28.166Z

HOUR

Extracts the hour from a time.

Syntax

HOUR(DATETIME)

DATETIME is the time you want to extract the hour from.

Example

HOUR(@Date@)

Date	New Column
2011-01-15T06:37:40.000Z	6
2011-01-23T07:09:58.000Z	7
2011-01-30T07:27:56.000Z	7

Notes on use

The DATETIME you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

Returned values range from 0 (12:00 am) to 23 (11:00 pm).

MAXDATE

Compares two or more dates and returns the latest date in the comparison.

Syntax

MAXDATE(DATETIME_1, [DATETIME_2, ...])

- DATETIME_1 is the first date.
- DATETIME_2, ... [optional] are the additional dates.

Example

MAXDATE(@Target Ship Date@ ,@Date Shipped@)

Target Ship Date	Date Shipped	New Column
2015-09-23T06:36:33.000Z	2015-09-29T06:36:33.000Z	2015-09-29T06:36...
2011-10-08T07:38:59.000Z	2011-10-16T07:38:59.000Z	2011-10-16T07:38...
2012-10-03T07:13:18.000Z	2012-09-16T07:13:18.000Z	2012-10-03T07:13...

Notes on use

The DATETIME you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

Here's how the MAXDATE function will respond to some common scenarios:

- If only one date is provided, the provided date is returned.
- The time zone of all the dates are temporarily converted to the same time zone to determine the latest date. The conversion is neither a permanent nor a visual transformation.
- Cells with text strings are ignored. Blank cells are ignored.
- Cells with errors are ignored.
- If no datetime objects are found, a blank cell is returned.

MIDNIGHT

Resets the given time to midnight (00:00).

Syntax

MIDNIGHT(DATETIME)

DATETIME is the time you want to reset.

Example

MIDNIGHT(@Date@)

Date	New Column
2011-01-15T06:37:40.000Z	2011-01-15T00:00:00.000Z
2011-01-23T07:09:58.000Z	2011-01-23T00:00:00.000Z
2011-01-30T07:27:56.000Z	2011-01-30T00:00:00.000Z

Notes on use

The `DATETIME` you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

The time zone isn't affected.

MINDATE

Compares two or more dates and returns the earliest date in the comparison.

Syntax

`MINDATE(DATETIME_1, [DATETIME_2, ...])`

- `DATETIME_1` is the first date.
- `DATETIME_2, ...` [optional] are the additional dates.

Example

`MINDATE(@Target Ship Date@ ,@Date Shipped@)`

Target Ship Date	Ship Date	New Column
2015-09-23T06:36:33.000Z	2015-09-29T06:36:33.000Z	2015-09-23T06:36...
2011-10-08T07:38:59.000Z	2011-10-16T07:38:59.000Z	2011-10-08T07:38...
2012-10-03T07:13:18.000Z	2012-09-16T07:13:18.000Z	2012-09-16T07:13...

Notes on use

The `DATETIME` you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

Here's how the `MINDATE` function will respond to some common scenarios:

- If only one date is provided, the provided date is returned.
- The time zone of all the dates are temporarily converted to the same time zone to determine the latest date. The conversion is neither a permanent nor a visual transformation.
- Cells with text strings are ignored. Blank cells are ignored.
- Cells with errors are ignored.
- If no datetime objects are found, a blank cell is returned.

MINUTE

Extracts the minute from a time.

Syntax

```
MINUTE(DATETIME)
```

DATETIME is the time you want to extract the minute from.

Example

```
MINUTE(@Date@)
```

Date	New Column
2011-01-15T06:37:40.000Z	37
2011-01-23T07:09:58.000Z	9
2011-01-30T07:27:56.000Z	27

Notes on use

The DATETIME you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

Returned values range from 0 to 59.

MONTH

Extracts the month from a date.

Syntax

```
MONTH(DATETIME)
```

DATETIME is the date you want to extract the month from.

Example

```
MONTH(@Date@)
```

Date	New Column
2011-03-20T06:03:57.000Z	3
2011-06-25T07:32:34.000Z	6
2012-08-06T08:23:39.000Z	8

Notes on use

The `DATETIME` you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

Returned values range from 1 (January) to 12 (December).

NETWORKDAYS

Returns the number of working days between two ts-date-time objects. This provides the same output as Excel's `NETWORKDAYS` function.

Syntax

```
NETWORKDAYS(DATE_TIME_START, DATE_TIME_END)
```

- `DATE_TIME_START` is a ts-date-time object for start date.
- `DATE_TIME` is a ts-date-time object for end date.

Example

```
NETWORKDAYS(@DATE@, DATE(2019,1,12))
```

DATE	New Column
2019-01-10T00:00:00.000Z	2
2019-01-10T00:00:00.000Z	2
2019-01-10T00:00:00.000Z	2
2019-01-10T00:00:00.000Z	2

NOW

Returns the current date and time.

Syntax

```
NOW(TIME_ZONE) returns the current date and time.
```

`TIME_ZONE` , optional, sets the time zone.

Example

```
NOW("GMT-03:00")
```

New Column
2018-01-22 16:47:49.330 -03:00
2018-01-22 16:47:49.330 -03:00
2018-01-22 16:47:49.330 -03:00

Notes on use

The `DATETIME` you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

If a time zone is not specified in the function, the returned datetime object will default to Greenwich Mean Time (GMT). See the [Date and Time Syntax](#) article for a list of time zones and their appropriate syntax.

QUARTER

Returns the quarter as an integer from a given ts-date-time object.

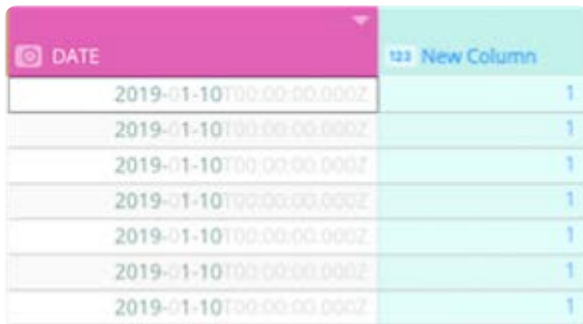
Syntax

```
QUARTER(DDATE_TIME)
```

`DATE_TIME` is a ts-date-time object.

Example

```
QUARTER(@DATE@)
```



The screenshot shows a table with two columns. The first column is labeled 'DATE' and contains seven rows of ISO 8601 timestamps: '2019-01-10T00:00:00.000Z', '2019-01-10T00:00:00.000Z', '2019-01-10T00:00:00.000Z', '2019-01-10T00:00:00.000Z', '2019-01-10T00:00:00.000Z', '2019-01-10T00:00:00.000Z', and '2019-01-10T00:00:00.000Z'. The second column is labeled 'New Column' and contains the value '1' for every row, indicating that all the provided timestamps fall within the first quarter of the year.

DATE	New Column
2019-01-10T00:00:00.000Z	1
2019-01-10T00:00:00.000Z	1
2019-01-10T00:00:00.000Z	1
2019-01-10T00:00:00.000Z	1
2019-01-10T00:00:00.000Z	1
2019-01-10T00:00:00.000Z	1
2019-01-10T00:00:00.000Z	1

SECOND

Extracts the seconds from a time.

Syntax

```
SECOND(DATETIME)
```

`DATETIME` is the time you want to extract the seconds from.

Example

```
SECOND(@Date@)
```

Date	123 New Column
2011-03-20T06:03:57.000Z	57
2011-06-25T07:32:34.000Z	34
2012-08-06T08:23:39.000Z	39

Notes on use

The `DATETIME` you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

Returned values range from 0 to 59.

SETTIMEZONE

Changes the time zone of a time to the time zone you specify.

Syntax

```
SETTIMEZONE(DATETIME, TIME_ZONE)
```

- `DATETIME` is the time you want to set the time zone of.
- `TIME_ZONE` is the time zone you want associated with the datetime object.

Example

```
SETTIMEZONE(@Date Received@, "GMT-3:00")
```

Date Received	New Column
2015-08-24T06:36:33.000Z	2015-08-24T06:36:33.000-03:00
2011-09-08T07:38:59.000Z	2011-09-08T07:38:59.000-03:00
2012-09-03T07:13:18.000Z	2012-09-03T07:13:18.000-03:00

Notes on use

The `DATETIME` you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

The transformation doesn't change the time, it simply assigns a new time zone to the existing time. See the Date and Time Syntax article for a list of time zones and their appropriate syntax.

TODAY

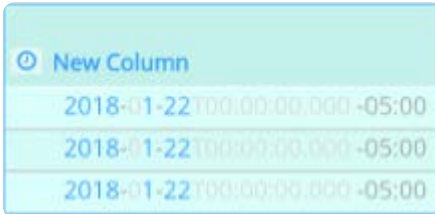
Returns the current date, doesn't include the time.

Syntax

TODAY()

Example

TODAY()



New Column
2018-01-22T00:00:00.000-05:00
2018-01-22T00:00:00.000-05:00
2018-01-22T00:00:00.000-05:00

WEEKOFYEAR

Returns the week number as an integer from a given ts-date-time object. This provides the same output as Excel's `WEEKNUM` function.

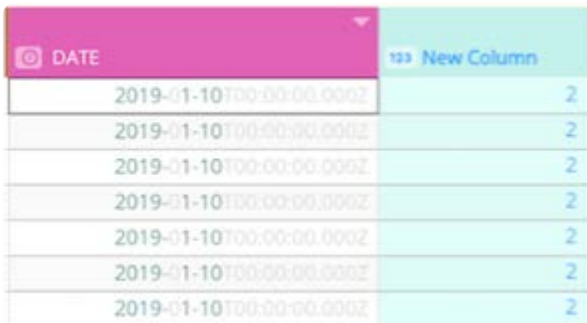
Syntax

WEEKOFYEAR(DATE_TIME)

DATE_TIME is a ts-date-time object.

Example

WEEKOFYEAR(<@DATE@>)



DATE	New Column
2019-01-10T00:00:00.000Z	2
2019-01-10T00:00:00.000Z	2
2019-01-10T00:00:00.000Z	2
2019-01-10T00:00:00.000Z	2
2019-01-10T00:00:00.000Z	2
2019-01-10T00:00:00.000Z	2
2019-01-10T00:00:00.000Z	2

WORKDAY

Returns a number that represents a date that is the indicated number of working days before or after a date (the starting date). Working days exclude weekends and any dates identified as holidays. This provides the same output as Excel's `WORKDAY` function. Use `WORKDAY` to exclude weekends or holidays when you calculate invoice due dates, expected delivery times, or the number of days of work performed.

Syntax

WORKDAY(STARTDATE, DAYS)

- STARTDATE is a date that represents the start date.
- DAYS is the number of nonweekend and nonholiday days before or after start date. A positive value for days yields a future date; a negative value yields a past date.

Example

WORKDAY(@DATE@,12)

DATE	New Column
2019-01-10T00:00:00.000Z	2019-01-28T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-28T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-28T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-28T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-28T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-28T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-28T00:00:00.000Z
2019-01-10T00:00:00.000Z	2019-01-28T00:00:00.000Z

YEAR

Extracts the year from a date.

Syntax

YEAR(DATETIME)

DATETIME is the date you want to extract the year from.

Example

Year()

Date	New Column
2012-07-20T07:15:24.000Z	2012
2016-03-25T07:06:44.000Z	2016
2014-01-30T07:16:25.000Z	2014

Notes on use

The DATETIME you provide must be a datetime object, a column that contains a datetime object, or a function that returns a datetime object.

Informational computed column functions

This section provides syntax and examples for the informational computed column functions you can use with the Data Prep [Compute tool](#). The informational functions let you query column values and create new columns based on the results.

FIRSTNONBLANK

Compares the values of two or more columns and returns the first non-blank value. This function provides the same output as Excel's FIRSTNONBLANK function.

Syntax

```
FIRSTNONBLANK(ARGUMENT_1, [ARGUMENT_2, ...])()
```

- ARGUMENT_1 is the first column.
- ARGUMENT_2, ... [optional] are the additional columns.

Example

```
FIRSTNONBLANK(@Current Employer@, @Previous Employer@, @School@)
```

Current Employer	Previous Employer	School	New Column
Banana Inc	Mermaidhut	Camden College	Banana Inc
CloudCo			CloudCo
		Greendale Community College	Greendale Community College
	BansheeElectronics	University of New York	BansheeElectronics
	PyramidIndustries		PyramidIndustries

Notes on use

If only one column is specified, the value of the provided column is returned.

If no non-blank values are found, the `FIRSTNONBLANK` function will return an empty cell (unless you include a final argument for what value to display in the output when no non-blank values are found).

ISBLANK

Checks for blank or null values within a specified column. If a blank or null value is found, the value `TRUE` is returned.

Syntax

```
ISBLANK(ARGUMENT)
```

ARGUMENT is the column to check.

Example

ISBLANK(@Column@)

123 Column	New Column
Rufus Daniel	false
	true
	true
1789	false
Bryant Carr	false

ISDATE

Syntax

ISDATE(ARGUMENT)

ARGUMENT is the column to check.

Example

ISDATE(@Column@)

A-Z Column	New Column
2013-01-02T00:01:...	true
	false
3/6/2014	false

Notes on use

The values must be datetime objects, not datetime text strings. Datasets imported from an Excel spreadsheet will automatically import dates as datetime objects. Dates from all other sources need to be converted to a datetime object using the `DATEVALUE` function. See the `DATEVALUE()` section of this article.

ISNULL

Checks for blanks or null values within a specified column. If a blank or null value is found, the value `TRUE` is returned.

Syntax

ISNULL(ARGUMENT)

ARGUMENT is the column to check.

Example

ISNULL(@Column@)

Column	New Column
Bryant Carr	false
	true
2013-01-02T00:01:00.000Z	false
Kelli Martinez	false
Rufus Daniel	false
	true

ISNUMBER

Checks for numeric values within a specified column. If a numeric value is found, the value TRUE is returned.

Syntax

ISNUMBER(ARGUMENT)

ARGUMENT is the column to check.

Example

ISNUMBER(@Column@)

Column	New Column
3/6/2014	false
2013-01-02T00:01:...	false
6	true
7	true

ISTEXT

Checks for text within a specified column. If a blank or null is found, the value TRUE is returned.

Syntax

ISTEXT(ARGUMENT)

ARGUMENT is the column to check.

Example

ISTEXT(@Column@)

A-Z Column	New Column
2013-01-02T00:01:...	false
Rufus Daniel	true
Kelli Martinez	true
	false
3/6/2014	true

Logical computed column functions

This section provides syntax and examples for the logical computed column functions you can use with the Data Prep [Compute tool](#). Use the logical functions to evaluate logical functions on column values. The new column contains `TRUE` or `FALSE` depending on the function results.

AND

Evaluates whether all arguments within an expression evaluate to `TRUE`. If the arguments do evaluate to `TRUE`, the value `TRUE` is returned.

Syntax

```
AND(ARGUMENT_1, [ARGUMENT_2, ...])
```

- `ARGUMENT_1` is the argument to evaluate.
- `ARGUMENT_2, ... [optional]` are the additional arguments.

Example

```
AND(@Column_A@, @Column_B@, @Column_C@)
```

Column_A	Column_B	Column_C	New Column
false	true	true	false
true	true	true	true
false	true	false	false
false	false	false	false

Notes on use

The `ARGUMENT`s you provide must be either a `TRUE` or `FALSE` value, a column that contains either value, or a function that returns either value.

The `AND` function is case insensitive, so it treats `True`, `TRUE`, and `true` the same way. Similarly, `False`, `FALSE`, and `false` are treated the same.

IF

Allows you to specify a different output depending on whether or not a given statement is true.

Syntax

IF(CONDITION, TRUE_VALUE, FALSE_VALUE)

- **CONDITION** is the expression you want to evaluate.
- **TRUE_VALUE** is the value the function returns if the **CONDITION** is true.
- **FALSE_VALUE** is the value that is returned if the **CONDITION** is not true.

Example

IF(@Current Employer@ = 0, "N/A", @Current Employer@)

Current Employer	New Column
0	N/A
Banana Computers Inc	Banana Computers Inc
0	N/A
Banana	Banana
CloudCo inc	CloudCo inc
Microstuff	Microstuff
0	N/A
0	N/A
Acme	Acme
0	N/A

Notes on use

The **IF** function is ideal in cases where a set of values need to be created based on information in one or more other columns.

The **CONDITION** must provide either a **TRUE** or **FALSE** value. Other functions can be incorporated as part of the **CONDITION**. Another **IF** function can be used as one or both of the values. This allows for very fine-grained control over the returned value. In most cases, the **CONDITION** will include an operator, see the Comparison operators () section of this article.

IFERROR

Allows you to specify a different output depending on whether or not a given statement is true.

Syntax

IFERROR(ARGUMENT, VALUE)

- **ARGUMENT** is the column you want to check.
- **VALUE** is the value to return if the column cell contains an error.

Example

IFERROR(@New Column@, "N/A")

123 New Column	123 New Column (1)
!	N/A
Banana Computers Inc	Banana Computers Inc
!	N/A
Banana	Banana
CloudCo inc	CloudCo inc
Microstuff	Microstuff
!	N/A

Notes on use

The `VALUES` you provide can be a text string or numeric value, a column that contains a text string or numeric value, or a function that returns a text string or numeric value.

For a cell where no error is found, the cell's original value is returned.

NOT

Reverses the result of an expression that results in a `TRUE` or `FALSE` value.

Syntax

`NOT(ARGUMENT)`

`ARGUMENT` is the `TRUE` or `FALSE` value you want to reverse.

Example

`NOT(@Column@)`

Column	New Column
false	true
true	false

Notes on use

The `ARGUMENT` you provide must be either a `TRUE` or `FALSE` value, a column that contains either value, or a function that returns either value.

The `NOT` function is case insensitive, so it treats `True`, `TRUE`, and `true` the same way. Similarly, `False`, `FALSE`, and `false` are treated the same.

OR

Determines if at least one value within an expression is `TRUE`. If one value is `TRUE`, the value `TRUE` is returned.

Syntax

OR(ARGUMENT_1, [ARGUMENT_2, ...])

- ARGUMENT_1 is the first argument you want to evaluate.
- ARGUMENT_2, ... [optional] are the additional columns.

Example

OR(@Column_A@, @Column_B@, @Column_C@)

Column_A	Column_B	Column_C	New Column
false	true	true	true
true	true	true	true
false	true	false	true
false	false	false	false


Notes on use

The ARGUMENT s you provide must be either a TRUE or FALSE value, a column that contains either value, or a function that returns either value.

The OR function is case insensitive, so it treats True, TRUE, and true the same way. Similarly, False, FALSE, and false are treated the same.

Mathematical computed column functions

This section provides syntax and examples for the mathematical computed column functions you can use with the Data Prep [Compute tool](#). For statistical functions, see [Statistical computed column functions](#).

To use a math function, you must store the value as a numeric data type. Numeric data types can be identified by the numeric data type  icon in the header row. If the value is not stored as a numeric data type, convert the value to numeric format using the [VALUE function](#).

In addition to the mathematical functions listed in this section, the following standard mathematical operations are supported:

- Multiply, divide a column by any number.
- Add to, subtract from a column by any number

Examples using column name "Revenue":

```
@Revenue@ * 100
```

```
@Revenue@ / 100
```

```
@Revenue@ + 100
```

```
@Revenue@ - 100
```

ABS

Returns the absolute value (ABS) of a real number.

In mathematical notation, absolute value is indicated with a bar on either side. For example, the absolute value of x is written as $|x|$.

Syntax

```
ABS(VALUE)
```

VALUE is the value you want to find the absolute value of.

Example

```
ABS(@Column@)
```

Column	New Column
6	6
14	14
-10	10
-11	11

Notes on use

The `VALUE` you provide must be a real number, a column that contains a real number, or a function that returns a real number.

`ABS` can be thought of as the distance from zero a given number has on a number line. For `ABS`, positive and negative do not matter. The number's distance from zero is the same, or absolute, regardless of whether the number is to the right of zero (positive) or to the left of zero (negative). In mathematical notation, absolute values is indicated with a bar on either side. For example, the absolute value of x would be written as $|x|$.

CEILING

Returns a given number rounded up to whole number.

Syntax

`CEILING(VALUE)`

`VALUE` is the value you want to round.

Example

`CEILING(@Column A value@)`

EXP

Returns the exponential for the specified value.

Syntax

`EXP(NUMBER)`

`NUMBER` is any real number.

Example

`EXP(@Column A value@)`

Column A value	New Column
6.588	726.3267627508812
9.43	12456.526731608414
11.345	84541.68455061226
14.796	2665760.6580085587
20	485165195.4097903

FACTORIAL

Returns the product of an integer and all the integers below it.

Syntax

FACTORIAL(NUMBER)

NUMBER is any real number.

Example

FACTORIAL(@Column A value@)

Column A value	New Column
6.588	1464.461395243295301844992
9.43	468217.34343705300007777749
11.345	36056362.17229731793933645445682699
14.796	640499653925.3259018035792303107524
20	2.43290200817664E+18

FLOOR

Returns a given number rounded down to whole number.

Syntax

FLOOR(VALUE)

VALUE is the value you want to round.

Example

FLOOR(@Column A value@)

Column A value	New Column
6.588	6
9.43	9
11.345	11
14.796	14
20	20

INT

Rounds a real number down to the next integer that is less than or equal to it.

Syntax

INT(VALUE)

VALUE is the real number you want to round down.

Example

INT(@Column@)

Column	New Column
6.1	6
14.11	14
-9.88	-10
-10.88	-11

Notes on use

The VALUE you provide must be a real number, a column that contains a real number, or a function that returns a real number.

LN

Returns the natural logarithm of a number. Natural logarithms are based on the constant e (2.71828182845904). This provides the same output as Excel's LN function.

Syntax

LN(NUMBER)

NUMBER is positive real number for which you want the natural logarithm.

Example

LN(@Column A value@)

Column A value	New Column
6.588	1.8852498123153938
9.43	2.2438960966453663
11.345	2.428777118231805
14.796	2.6943568739702077
20	2.995732273553991

LOG

Returns the logarithm of a number to the base you specify. This function provides the same output as Excel's LOG function.

Syntax

LOG(NUMBER, BASE)

- NUMBER is positive real number for which you want the natural logarithm.
- BASE is the base of the logarithm.

Example

LOG(@Column A value@,2)

Column A value	New Column
6.588	2.719840555064268
9.43	3.237257770900372
11.345	3.5039847038976126
14.796	3.8871353004619085
20	4.321928094887363

LOG10

Returns the base-10 logarithm of a number. This function provides the same output as Excel's LOG10 function.

Syntax

LOG10(NUMBER)

NUMBER is positive real number for which you want the natural logarithm.

Example

LOG10(@Column A value@)

Column A value	New Column
6.588	0.8187535904977168
9.43	0.9745116927373284
11.345	1.0548045002209547
14.796	1.1701443226433565
20	1.3010299956639813

MOD

Returns the remainder after number is divided by divisor. The result has the same sign as divisor. This provides the same output as Excel's MOD function.

Syntax

MOD(NUMBER, DIVISOR)

- NUMBER is any real number.
- DIVISOR is any real number.

Example

MOD(@Column A value@,3)

Column A value	New Column
6.588	0.588
9.43	0.43
11.345	2.345
14.796	2.796
20	2

POWER

A Math/Trigonometric function that computes and returns the result of a number raised to a power. This function provides the same output as Excel's POWER function.

Syntax

POWER(NUMBER, POWER)

- NUMBER is the base number that is any real number.
- POWER is the exponent, any real number, to which the base number is raised.

Example

POWER(@Column A value@,3)

123 Column A value	123 New Column
6.588	285.930689472
9.43	838.5618069999999
11,345	1460.2038886250002
14.796	3239.1642303359995
20	8000

ROUND

Rounds a number to the number of decimal places you specify.

Syntax

ROUND(VALUE, PLACES)

- VALUE is the real number you want to round.
- PLACES is the number of decimal places to round to.

Example

ROUND(@Column@, 2)

123 Column	123 New Column
6.3141592653589793238462643	6.31
14.34557519189487725623089073	14.35
-9.62300888156922481138448284	-9.62
-10.62300888156922481138448284	-10.62

Notes on use

The VALUE you provide must be a real number, a column that contains a real number, or a function that returns a real number.

The PLACES value must be a positive integer (not negative or a contain decimal). Values with less decimal places than what you specified won't be affected by the ROUND function.

ROUNDDOWN

Rounds a number down to the number of decimal places you specify.

Syntax

ROUNDDOWN(VALUE, PLACES)

- VALUE is the real number you want to round down.
- PLACES is the number of decimal places to round down to.

Example

`ROUNDDOWN(@Column@, 3)`

Column	New Column
6.3141592653589793238462643	6.314
14.34557519189487725623089073	14.345
-9.62300888156922481138448284	-9.623
-10.62300888156922481138448284	-10.623

Notes on use

`ROUNDDOWN` is similar to `ROUND` except that it always rounds down.

The `VALUE` you provide must be a real number, a column that contains a real number, or a function that returns a real number.

The `PLACES` value must be a positive integer (not negative or contain a decimal). Values with fewer decimal places than what you specified won't be affected by the `ROUNDDOWN` function. If zero decimal places are specified, `ROUNDDOWN` operates like the `INT` function.

ROUNDPERC

Rounds a real percentage value (between -100 and 100) to the number of decimal places you specify, ensuring that values near 0% and $\pm 100\%$ have at least the number of decimal places you specify and as many, but no more, decimal places to distinguish the value from true 0% or $\pm 100\%$ values.

Syntax

`ROUNDPERC(VALUE, PLACES)`

- `VALUE` is the real percentage value you want to round.
- `PLACES` is the minimum number of decimal places to round to.

Example

`ROUNDPERC(@Column@, 2)`

Column	New Column
6.3141592653589793238462643	6.31
14.34557519189487725623089073	14.35
-9.62300888156922481138448284	-9.62
-10.62300888156922481138448284	-10.62

Notes on use

The `VALUE` you provide must be a real percentage value between -100 and 100, a column that contains a real percentage value, or a function that returns a real percentage value.

The `PLACES` value must be an integer (no decimals). Values with fewer decimal places than what you specified won't be affected by the `ROUNDPERC` function.

ROUNDUP

Rounds a number up to the number of decimal places you specify.

Syntax

`ROUNDUP(VALUE, PLACES)`

- `VALUE` is the real number you want to round up.
- `PLACES` is the number of decimal places to round up to.

Example

`ROUNDUP(@Column@, 3)`

Column	New Column
6.3141592653589793238462643	6.315
14.34557519189487725623089073	14.346
-9.62300888156922481138448284	-9.624
-10.62300888156922481138448284	-10.624

Notes on use

`ROUNDUP` is similar to `ROUND` except that it always rounds up.

The `VALUE` you provide must be a real number, a column that contains a real number, or a function that returns a real number.

The `PLACES` value must be a positive integer (not negative or contain a decimal). Values with fewer decimal places than what you specified won't be affected by the `ROUNDUP` function.

SIGN

Determines the sign of a number. Returns 1 if the number is positive, zero (0) if the number is 0, and -1 if the number is negative. This provides the same output as Excel's `SIGN` function.

Syntax

`SIGN(NUMBER)`

`NUMBER` is any real number.

Example

SIGN(@Column A value@)

Column A value	New Column
-4.44	-1
9.43	1
-11.345	-1
14.796	1
-20	-1

SQRT

Returns a positive square root. This function provides the same output as Excel's SQRT function.

Syntax

SQRT(NUMBER)

NUMBER is any positive number for which you want to calculate the square root.

Example

SQRT(@Column A value@)

Column A value	New Column
6.588	2.5667099563448925
9.43	3.0708305065568173
11.345	3.368233958619858
14.796	3.8465569019579053
20	4.47213595499958

SUM

Adds the given numeric values together.

Syntax

SUM(VALUE_1, [VALUE_2, ...])

- VALUE_1 is the first value.
- VALUE_2, ... [optional] are the additional values.

Example


SUM(@Column_A@, @Column_B@, @Column_C@)

123 Column_A	123 Column_B	123 Column_C	123 New Column
6	14	36	56
6	17	36	59
7	11	24	42
6	13	24	43

Notes on use

The `VALUE` you provide must be a numeric value, a column that contains numeric values, or a function that returns a numeric value.

Statistical computed column functions

This section provides syntax and examples for the statistical computed column functions you can use with the Data Prep [Compute tool](#). To use a statistics function, you must store the value as a numeric data type. Numeric data types can be identified by the numeric data type  icon in the header row. If the value is not stored as a numeric data type, convert the value to numeric format using the [VALUE function](#).

AVERAGE

Calculates the value equal to the sum of a list of numbers divided by the number of items in the list.

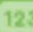

Syntax

```
AVERAGE(VALUE_1, [VALUE_2, ...])
```

- VALUE_1 is the first value.
- VALUE_2, ... [optional] are the additional values.

Example

```
AVERAGE(@Column_A@, @Column_B@, @Column_C@)
```

 Column_A	 Column_B	 Column_C	 New Column
5	15	10	10
7	14		10.5
7	21	5	11
5	10	11	8.66666666666666...

Notes on use

The [VALUE](#) you provide must be a numeric value, a column that contains numeric values, or a function that returns a numeric value.

MAX

Returns the greatest (maximum) value from a set of values.

Syntax

MAX(VALUE_1, [VALUE_2, ...])

- VALUE_1 is the first value.
- VALUE_2, ... [optional] are the additional values.

Example

MAX(@Column_A@, @Column_B@, @Column_C@)

123 Column_A	123 Column_B	123 Column_C	123 New Column
5	17	7	17
6	17		17
6	10	4	10
7	12	13	13

Notes on use

The VALUE you provide must be a numeric value, a column that contains numeric values, or a function that returns a numeric value.

MEDIAN

Returns the number that exists in the middle of a range of numbers ordered from lowest to highest value.

Syntax

MEDIAN(VALUE_1, [VALUE_2, ...])

- VALUE_1 is the first value.
- VALUE_2, ... [optional] are the additional values.

Example

MEDIAN(@column_A@, @column_B@, @column_C@, @column_C@, @column_E@)

123 Column A	123 Column B	123 Column C	123 Column D	123 Column E	123 New Column
3	4	5	2	10	4
10	4	1	0	0	1
4	2	0	7	0	2
6	1	0	4	4	4
9	3	0	4	3	3
4	4	7	4	6	4
2	3	8	9	10	8
2	3	10	3	1	3

Notes on use

The `VALUE` you provide must be a numeric value, a column that contains numeric values, or a function that returns a numeric value.

In a range with an even set of numbers, the median is the center number, half the numbers are to the right of the value returned and half of the numbers are to the left of the value returned. Where there is no single number in the middle of the range), MEDIAN calculates the average on the two numbers on either side of the mid-point.

Note

MEDIAN is different than AVERAGE. AVERAGE is an arithmetic mean—calculated by adding up a set of numbers and then dividing by the number of values in the set. MEDIAN takes the value at the center of the range. In those number ranges that exhibit a balance in the distribution of values within the collection, the MEDIAN and AVERAGE calculations may coincide; in skewed distributions, the values will be different.

MIN

Returns the smallest (minimum) value from a set of values.

Syntax

`MIN(VALUE_1, [VALUE_2, ...])`

- `VALUE_1` is the first value.
- `VALUE_2, ... [optional]` are the additional values.

Example

`MIN(@Column_A@, @Column_B@, @Column_C@)`

123 Column_A	123 Column_B	123 Column_C	123 New Column
5	15	10	5
7	14		7
7	21	5	5
5	10	11	5

Notes on use

The **VALUE** you provide must be a numeric value, a column that contains numeric values, or a function that returns a numeric value.

MODE

Returns the value that occurs most frequently in a set of numbers.

Syntax

MODE(VALUE_1, VALUE_2, [VALUE_3, ...])

- VALUE_1 is the first value.
- VALUE_2 is the second value.
- VALUE_3, ... [optional] are the additional values.

Example

MODE(@Column_A@, @Column_B@, @Column_C@)

123 Column_A	123 Column_B	123 Column_C	123 New Column
3	9	9	9
8	13	12	8
7	11	11	11
6	19	11	6

Notes on use

The **VALUE** you provide must be a numeric value, a column that contains numeric values, or a function that returns a numeric value.

If multiple numbers have an equal occurrence count greater than one, the value returned is the number (of those that are equal in frequency) that appears first in the set (reading from left to right). If no number appears more than once, the function returns an error.

The most common problem related to `MODE` is when the provided set of numbers has no duplicates. At least one number must appear two or more times for the function to successfully evaluate. If the minimum number of arguments (two) are used, then each argument must evaluate to the same number or an error will occur. As you might expect, larger sets of numbers with a more limited variation will reduce the chance of `MODE` returning an error.

STDEV

Estimates the standard deviation, how much variation from the average, that exists within a sample set of data.

Syntax

`STDEV(VALUE_1, [VALUE_2, ...])`

- `VALUE_1` is the first value.
- `VALUE_2, ...` [optional] are the additional values.

Example

`STDEV(@Column_A@, @Column_B@, @Column_C@)`

123 Column_A	123 Column_B	123 Column_C	123 New Column
0.66	0.28	0.12	0.2773685875028...
0.66	0.34	0.12	0.2715388247255...
0.77	0.22	0.08	0.3647373484212...
0.66	0.26	0.08	0.2968725877094...

Notes on use

The `VALUE` you provide must be a numeric value, a column that contains numeric values, or a function that returns a numeric value.

The standard deviation for data is the square root of its variance. If the set under analysis represents all data points (referred to as a population), use `STDEVP` instead.

STDEVP

Estimates the standard deviation, how much variation from the average, that exists within the entire set (population) of data.

Syntax

`STDEVP(VALUE_1, [VALUE_2, ...])`

- `VALUE_1` is the first value.
- `VALUE_2, ...` [optional] are the additional values.

Example

STDEVP(@Column_A@, @Column_B@, @Column_C@)

123 Column_A	123 Column_B	123 Column_C	123 New Column
0.66	0.28	0.12	0.2264705033528...
0.66	0.34	0.12	0.2217105219775...
0.77	0.22	0.08	0.2978067979225...
0.66	0.26	0.08	0.2423954528359...

Notes on use

The `VALUE` you provide must be a numeric value, a column that contains numeric values, or a function that returns a numeric value.

If the set under analysis represents only a sample of the data points, use `STDEV` instead.

VAR

Estimates how much dispersion exists (how much the values are spread out) within a samples set of data.

Syntax

VAR(VALUE_1, [VALUE_2, ...])

- `VALUE_1` is the first value.
- `VALUE_2, ...` [optional] are the additional values.

Example

VAR(@Column_A@, @Column_B@, @Column_C@)

123 Column_A	123 Column_B	123 Column_C	123 New Column
0.55	0.3	0.1	0.0508333333333...
0.55	0.34	0.07	0.0579
0.66	0.34	0	0.1089333333333...
0.66	0.2	0.04	0.1036

Notes on use

The `VALUE` you provide must be a numeric value, a column that contains numeric values, or a function that returns a numeric value.

If the set under analysis represents all data points (referred to as a population), use `VARP` instead.

VARP

Estimates how much dispersion exists (how much the values are spread out) within the entire set (population) of data.

Syntax

```
VARP(VALUE_1, [VALUE_2, ...])
```

- `VALUE_1` is the first value.
- `VALUE_2, ...` [optional] are the additional values.

Example

```
VARP(@Column_A@, @Column_B@, @Column_C@)
```

123 Column_A	123 Column_B	123 Column_C	123 New Column
0.55	0.3	0.1	0.033888888888888...
0.55	0.34	0.07	0.0386
0.66	0.34	0	0.072622222222222...
0.66	0.2	0.04	0.069066666666666...

Notes on use

The `VALUE` you provide must be a numeric value, a column that contains numeric values, or a function that returns a numeric value.

If the set under analysis represents only a sample of data points, use `VAR` instead.

Text computed column functions

This section provides syntax and examples for the text computed column functions you can use with the Data Prep [Compute tool](#).

CHAR

Returns the character for the specified ASCII value.

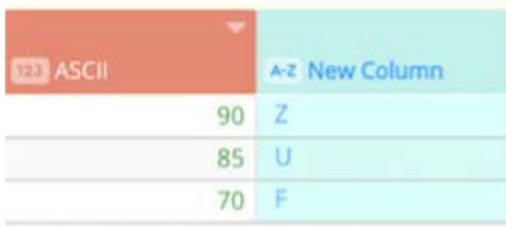
Syntax

```
CHAR(INT)
```

INT is the ASCII value for the character to return.

Example

```
CHAR(ASCII)
```



The screenshot shows a table with two columns: 'ASCII' and 'New Column'. The 'ASCII' column contains the values 90, 85, and 70. The 'New Column' column contains the corresponding characters 'Z', 'U', and 'F'.

ASCII	New Column
90	Z
85	U
70	F

CONCATENATE

To concatenate a series of text strings into a single text string, you can use the + operator. This example combines two columns with a hyphen between them:

```
@Last@ + "-" + @First@
```

Alternatively, you can use the `CONCATENATE` function:

```
CONCATENATE(@Last@ , "-", @First@ )
```

Syntax

```
CONCATENATE(String_1, [String_2, ...])
```

- `String_1` is the first value.
- `String_2, ...` [optional] are the additional strings.

Example

CONCATENATE(@Applicant Last@ , " ", @Applicant First@ , " of " , @City@)

Applicant First	Applicant Last	City	New Column
Maximo	Ehmann	Wolfdale	Ehmann, Maximo of Wolfdale
Velia	Goldman	Wolfdale	Goldman, Velia of Wolfdale
Nick	Walters	Glennville	Walters, Nick of Glennville
Rachel	Ray	Hull	Ray, Rachel of Hull

Notes on use

The `STRING` you provide can be a text string or numeric value, a column that contains a text string or numeric value, or a function that returns a text string or numeric value.

FIND

Determines if one word (or string of text) can be found in a second piece of text. If found, the `FIND` function returns the numeric position of the text string within the second text string. The characters in the second string are counted and the number indicates the character where the first piece of text begins its overlap in the second string of text.

This function allows for an optional third argument, a number. This number indicates the position (by number of characters) where you want the search to begin in the second string. If the third argument is omitted, then the second string is searched beginning at the first character.

If the first string is not found in the second string, then the function returns `0`.

Syntax

FIND(`STRING_1`, `STRING_2`, [`VALUE`])

- `STRING_1` is the string you want to find.
- `STRING_2` is the string you want to search in.
- `VALUE`, optional, is the numeric position in `STRING_2` you want to start your search.

Examples

FIND("Tech",@School@)

School	New Column
UC Sunnydale	0
Pacific Tech	9
Blue Mountain State	0
South Harmon Institute of Technology	27

FIND("the", "The quick sly fox jumped over the lazy brown dog laying next to the other dog.") returns a value of **31**.

FIND("dog", "The quick sly fox jumped over the lazy brown dog laying next to the other dog.") returns a value of **46**, which corresponds to the first occurrence of "dog" in the second string.

FIND("dog", "The quick sly fox jumped over the lazy brown dog laying next to the other dog.", 47) returns a value of **75** because the third argument value of **47** pushes the start of the search past character **46** (where the first "dog" occurs) and forces the function to find the second "dog" in the string.

Notes on use

The `STRING` you provide must be a text string, a column that contains a text string, or a function that returns a text string. Likewise, the `VALUE` must be a numeric value, a column that contains a numeric value, or a function that returns a numeric value.

If `STRING_1` occurs multiple times in `STRING_2`, `FIND` only indicates the position of the first match—not successive matches in the pair.

The `FIND` function is case sensitive, so it treats `True`, `TRUE`, and `true` separately.

This function matches pieces of text—not just words. Therefore, the text "jump" will be determined to be in the string "jumped" at position 1.

Text characters, not just words, can be used as search strings and discovered in the second string.

HASHVALUE

Transforms a text string to make fuzzy matching easier.

Syntax

HASHVALUE(`STRING`, `OPTION`, [`VALUE`])

- `STRING` is the string you want to transform.
- `OPTION` is the algorithm to use for the transformation. Available options are: `METAPHONE`
- `NGRAM FINGERPRINT`
- `VALUE`, used with `NGRAM`, specifies the number of ngrams to use.

Example

HASHVALUE(@Current Employer@, "metaphone")

A-Z Current Employer	A-Z New Column
Boogle	PKL
CloudCo inc	KLTKNK
Self	SLF
Banana Inc	PNNNK
Acme inc.	AKMNK

Notes on use

The `STRING` you provide must be a text string, a column that contains text strings, or a function that returns a text string. Both `OPTION` and `VALUE` are treated as strings and must be surrounded by quotation marks, i.e. "metaphone."

`HASHVALUE` uses algorithms to generate hashes based on provided string values. The algorithms used are also used by the `Cluster + Edit` column operation to find close matches between values within a column. See [Cluster + Edit](#) for more information on `METAPHONE`, `NGRAM`, and `FINGERPRINT`.

LEFT

Returns a given number of characters starting from the left-most (beginning) position of a text string.

Syntax

```
LEFT(String, Value)
```

- `STRING` is the string you want to search.
- `VALUE` is how many characters to return. The default is 1.

Example

```
LEFT(@School@,4)
```

A-Z School	A-Z New Column
UC Sunnydale	UC S
Pacific Tech	Paci
Blue Mountain State	Blue
South Harmon Institute of Technology	Sout

Notes on use

The `STRING` you provide must be a text string, a column that contains text strings, or a function that returns a text string.

LEN

Counts the number of characters in a text string.

Syntax

```
LEN(String)
```

`STRING` is the text string you want to evaluate.

Example

```
LEN(@School@)
```

A-Z School	123 New Column
Coolidge College	16
The University of Los Angeles	29
Camden College	14
South Harmon Institute of Technology	36

Notes on use

The `STRING` you provide must be a text string, a column that contains text strings, or a function that returns a text string.

LOWER

Converts text in column to all lowercase.

Syntax

```
LOWER(String, LOCALE)
```

- `STRING` is the string or column that you want to convert to lowercase.
- `LOCALE` (optional) is the locale, which may need to be specified in order to output required characters for the lowercase.

Refer to <https://www.oracle.com/java/technologies/javase/jdk8-jre8-suported-locales.html> for the supported locale values.

Example

```
LOWER(@Values@, "tr")
```

A-Z Values	A-Z New Column
IAŞLIK	iaşlık
IAŞLIK	iaşlık
IAŞLIK	iaşlık

MID

Returns a given number of characters from the middle of a text sting.

Syntax

```
MID(String, VALUE_1, VALUE_2)
```

- `STRING` is the text string you want to evaluate.
- `VALUE_1` is the start position.
- `VALUE_2` is the number of characters to return.

Example:

MID(@School@,4, 5)

A-Z School	A-Z New Column
Blue Mountain State	e Mou
Pennbrook University	nbroo
Hillman College	lman

Notes on use

The `STRING` you provide must be a text string, a column that contains text strings, or a function that returns a text string. The values provided must be a numeric value, a column that contains numeric values, or a function that returns a numeric value.

PADLEFT

Pads a string with a specified character, for the specified number of times. This provides same output as MySQL LPAD.

Syntax

`PADLEFT(STRING, NUMBER, VALUE)`

- `STRING` or column is the value to pad.
- `NUMBER` is the number of times to replace with the `VALUE`.
- `VALUE` is the literal replacement value.

Example

`PADLEFT(@set@, 10, "-")`

A-Z set4	A-Z New Column
test1	----test1
test2	----test2
test3	----test3

PADRIGHT

Pads a string with a specified character, for the specified number of times. This provides same output as MySQL LPAD and RPAD.

Syntax

`PADRIGHT(STRING, NUMBER, VALUE)`

- `STRING` or column is the value to pad.
- `NUMBER` is the number of times to replace with the `VALUE`.
- `VALUE` is the literal replacement value.

Example

```
PADRIGHT(@set@, 10, "-")
```

A-Z set4	A-Z New Column
test1	test1-----
test2	test2-----
test3	test3-----

REGEXP

Executes a search and replace on a text string using regular expressions. This function is based on Java Regex.

Tip

To determine if a string of text appears in another piece of text without conducting any replacements, see [FIND](#). The `FIND` function has the advantage of a slightly more usable syntax—however, the trade-off is that `FIND` is slightly less powerful in pattern matching.

Syntax

```
REGEXP(STRING_1, STRING_2, STRING_3)
```

- `STRING_1` is the text string you want to search.
- `STRING_2` is the text you are searching for.
- `STRING_3` is the text you want to replace `STRING_2`.

The three arguments are required. `STRING_1` must be a text string, a column that contains text strings, or a function that returns a text string. `STRING_2` and `STRING_3` are composed of character combinations that define the search and replace activity.

Note

In regular expressions, there are 12 characters with special meanings:

`\ ^ . | ? * + () []` and open curly brace.

If you want to search for these actual characters and not their special meanings, add a double backslash (not a single backslash) before it. For example, to search for asterisk characters with a regular expression, type `"\" not ""`. To search for a backslash character with a regular expression, type four backslash characters.

Examples

Convert a space character to an underscore

REGEXP(@School@," ","_")

A-Z School	A-Z New Column
Adams College	Adams_College
California University	California_University
Adams College	Adams_College
University of New York	University_of_New_York

Replace a text string with another textstring

REGEXP(@ProductID@ , "ABC", "DEF")

Convert a slash to a hyphen

REGEXP("The/quick/sly/fox." , "/" , "-") returns *The-quick-sly-fox.*

Convert a backslash (special character) to a hyphen

REGEXP(@ProductID@ , "\\\" , "-")

To convert an asterisk (special character) to a hyphen

REGEXP(@ProductID@ , "*" , "-")

Delete characters from Column 1 that are not numbers

REGEXP(_column1_ , "[^0-9]" , "")

Examples of Extract and Replace patterns

COMMAND	RETURNS
RegexpExtract("replace me", "e m")	"e m"
RegexpExtract("replace me", "e.?m")	"e m"
RegexpExtract("replace me", "r.*c")	"replac"
RegexpExtract("123123456789", "(123)+456(.*)")	"123123456789"
RegexpExtract("123123456789", "(123)+456(.*)", 0)	"123123456789"
RegexpExtract("123123456789", "(123)+456(.*)", 1)	"123"
RegexpExtract("123123456789", "(123)+456(.*)", 2)	"789"

COMMAND	RETURNS
<code>RegexExtract("456789", "(123)*456(.*)", 2)</code>	"789"
<code>RegexReplace("replace me", "e m", "---")</code>	"replac---e"
<code>RegexReplace("replace me", "e.?m", "---")</code>	"replac---e"
<code>RegexReplace("replace me", "r.*c", "--")</code>	"--e me"
<code>RegexReplace("123123456789", "(123)+456(.*)", "---")</code>	"---"
<code>RegexReplace("123123456789", "abc", "---")</code>	"123123456789"

Notes on use

For more guidance on Regex pattern matching, refer to

<https://docs.oracle.com/javase/8/docs/api/java/util/regex/Pattern.html>

REPEAT

Repeats a specified string N number of times.

Syntax

```
REPEAT(VALUE,REPEAT)
```

- VALUE is the string or column to locate and repeat.
- REPEAT is the number of times to repeat the VALUE .

Example

```
REPEAT(@set4@, 3)
```

A-Z set4	A-Z New Column
test1	test1test1test1
test2	test2test2test2
test3	test3test3test3
test4	test4test4test4

REPLACE

Replaces part of a text string, based on the number of characters you specify, with a different text string.

Syntax

```
REPLACE(VALUE, START NUM, NUM CHARS, NEW VALUE)
```

- **VALUE** is the text or the column in which you want to replace characters.
- **START NUM** is the start position of the character in the **VALUE** that you want to replace.
- **NUM CHARS** is the number of characters in the text that you want to replace with the new string.
- **NEW VALUE** is the replacement value. Note this is case sensitive.

Example

```
REPLACE(@timestamp@,10,5," ")
```



timestamp	New Column (1)
2019-01-23T05:24:56	2019-01-2 24:56
2019-01-23T05:24:56	2019-01-2 24:56
2019-01-23T05:24:56	2019-01-2 24:56
2019-01-23T05:24:56	2019-01-2 24:56
2019-01-23T05:24:56	2019-01-2 24:56

Notes on use

Use **REPLACE** when you want to replace any text that occurs in a specific location in a text string; use **SUBSTITUTE** when you want to replace specific text in a text string. For example: `REPLACE(@Hospital Name@, Search(@Hospital Name@,"ADVOCATE"), 8, "ALPHA")`

REVERSE

Reverses the specified string.

Syntax

```
REVERSE(STRING)
```

STRING is the column's value or string to reverse.

Example

```
REVERSE(@set4@)
```

RIGHT

Returns a give number of characters starting from the right-most (end) position of a text string.

Syntax

RIGHT(String, Value)

- `String` is the string to search.
- `Value` is how many characters to return. The default is 1.

Example

RIGHT(@School@,4)

A-Z School	A-Z New Column
Pacific Tech	Tech
Grand Lakes University	sity
Coolidge College	lege
South Central Louisiana State University	sity

Notes on use

The `String` you provide must be a text string, a column that contains text strings, or a function that returns a text string.

SEARCH

Searches for a specified string and returns the index of the string. If not found, returns value of -1.

Syntax

SEARCH(Value, String)

- `Value` is the text or the column in which you want to substitute characters
- `String` is the string to search.

Example

SEARCH(@Hospital Name@, "ADVENTIST")

Hospital Name	New Column (1)
ADVANCED SURGICAL HOSPITAL	-1
ADVENTIST BOLINGBROOK HOSPITAL	1
ADVENTIST GLENOAKS	1
ADVENTIST LA GRANGE MEMORIAL HOSPITAL	1
ADVENTIST MEDICAL CENTER	1
ADVENTIST MEDICAL CENTER	1
ADVENTIST MEDICAL CENTER - REEDLEY	1
ADVOCATE BROMENN MEDICAL CENTER	-1
ADVOCATE CHRIST HOSPITAL & MEDICAL CEN...	-1

Notes on use

SEARCH can be combined with REPLACE .

Example

```
REPLACE(@Hospital Name@, Search(@Hospital Name@,"ADVOCATE"), 8, "ALPHA")
```

STR

Converts the data in the argument into a text string.

Syntax

```
STR(VALUE)
```

VALUE is the value you want to convert to a text string.

Example

```
STR(@Date@)
```

Date	A-Z New Column
2016-03-19T00:00:00.000Z	2016-03-19T00:00:00.000Z
2012-06-30T00:00:00.000Z	2012-06-30T00:00:00.000Z
2013-12-28T00:00:00.000Z	2013-12-28T00:00:00.000Z

Notes on use

The Value you provide must be a numeric value, a column that contains a numeric value, or a function that returns a numeric value.

The STR function is useful for converting a numeric value into text or for ensuring that a column of mixed text and number values is treated entirely as a column of text so that other text functions can successfully be executed against it.

SUBSTITUTE

Substitutes new text for old text in a text string.

Syntax

```
SUBSTITUTE(VALUE, OLD TEXT, NEW TEXT)
```

- VALUE is the text or the column in which you want to substitute characters.
- OLD TEXT is the text you want to replace. Note this is case sensitive.
- NEW TEXT is the text you want to use to replace OLD TEXT. Note this is case-sensitive.

Example

SUBSTITUTE(@Hospital Name@,"CREIGHTON","Merton")

Hospital Name	New Column
ALBANY MEMORIAL HOSPITAL	ALBANY MEMORIAL HOSPITAL
ALBANY VA MEDICAL CENTER	ALBANY VA MEDICAL CENTER
ALBEMARLE HOSPITAL AUTHORITY	ALBEMARLE HOSPITAL AUTHORITY
ALBERT EINSTEIN MEDICAL CENTER	ALBERT EINSTEIN MEDICAL CENTER
ALEGENT CREIGHTON HEALTH BERGAN MERCY MEDICAL CTR	ALEGENT MERTON HEALTH BERGAN MERCY MEDICAL ...
ALEGENT CREIGHTON HEALTH CREIGHTON UNIVERSITY MED	ALEGENT MERTON HEALTH MERTON UNIVERSITY MED
ALEGENT CREIGHTON HEALTH IMMANUEL MEDICAL CENTER	ALEGENT MERTON HEALTH IMMANUEL MEDICAL CENT...
ALEGENT CREIGHTON HEALTH LAKESIDE HOSPITAL	ALEGENT MERTON HEALTH LAKESIDE HOSPITAL
ALEGENT CREIGHTON HEALTH MEMORIAL HOSPITAL, SCHUYL	ALEGENT MERTON HEALTH MEMORIAL HOSPITAL, SCH...

Notes on use

Use `SUBSTITUTE` when you want to replace specific text in a text string; use `REPLACE` when you want to replace any text that occurs in a specific location in a text string.

TRIM

Removes all leading and trailing spaces for the specified string.

Note

The `TRIM` function was designed to trim the 7-bit ASCII space character (value 32) from text. In the Unicode character set, there is an additional space character called the nonbreaking space character that has a decimal value of 160. This character is commonly used in Web pages as the HTML entity, ` `. By itself, the `TRIM` function does not remove this nonbreaking space character.

Syntax

`TRIM(String)`

`String` is the value you want to you want trimmed.

The column can be specified as the `String` value as in the following example.

Example

`TRIM(@Company@)`

Company	New Column
Apple corp	Apple corp
Apple corporation	Apple corporation
Apple computers	Apple computers

TRIMLEFT

Returns the string stripped of whitespace from the left end of the string.

Syntax

```
TRIMLEFT(String)
```

String is the column's values you want to you want trimmed.

Example

```
TRIMLEFT(@Company@)
```

A-Z Company	A-Z New Column
Apple corp	Apple corp
Apple corporation	Apple corporation
Apple computers	Apple computers

TRIMRIGHT

Returns the string stripped of whitespace from the right end of the string.

Syntax

```
TRIMRIGHT(String)
```

String is the column's values you want to you want trimmed.

Example

```
TRIMRIGHT(@Company@)
```

A-Z Company	A-Z New Column
Apple corp	Apple corp
Apple corporation	Apple corporation
Apple computers	Apple computers

UPPER

Converts text in column to all uppercase.

Syntax

```
UPPER(String, LOCALE)
```

- String is the string or column that you want to convert to uppercase.
- LOCALE (optional) is the locale, which may need to be specified in order to output required characters for the uppercase.

Refer to <https://www.oracle.com/java/technologies/javase/jdk8-jre8-suported-locales.html> for the supported locale values.

Example

UPPER(@Values@, "tr")

A-Z Values	A-Z New Column
iaşlık	IAŞLIK
iaşlık	IAŞLIK
iaşlık	IAŞLIK

VALUE

Converts numbers stored as a string value into a numeric value.

Syntax

VALUE(STRING)

STRING is the numbers, stored as a text string, you want to convert to a numeric value.

Example

VALUE(@COLUMN@)

Column A value	New Column
6.588	1464.461395243295301844992
9.43	-468217.34343705300007777749
11.345	36056362.17229731793933645445682699
14.796	640499653925.3259018035792303107524
20	2.43290200817664E+18

Notes on use

The STRING you provide must be a number stored as text string, a column that contains a number stored as text string, or a function that returns a number stored as text string.

If STRING contains characters other than numbers, the functions returns an error. A single period (decimal point) is allowed within the argument in order to create a real number.

The VALUE function is useful for converting a text value into numbers to ensure that a column of number values is treated as a column of numbers so that number based functions can successfully be executed against it.

Comparison operators

Use comparison operators to test logical conditions. They are most commonly used within the first argument of the [IF](#) function in order to generate a `TRUE` or `FALSE` value.

The following are the operators you can use in Data Prep:

OPERATOR	DEFINITION	EXAMPLE THAT RETURN TRUE
=	Equal to	1 + 2 = 3
>	Greater than	3 > 2
>=	Greater than or equal to	11 >= 10 11 >= 11
<	Less than	2 < 3
<=	Less than or equal to	10 <= 11 10 <= 10
<>	Not equal to	2 <> 3

Comparison operators with numeric values

Using comparison operators to conduct comparisons between numeric values is straightforward. Bear in mind, however, that the two values to be compared must both be of the same data type. The text value "3" is not the same as the numeric value 3.

To safeguard against mixing data types, use the [VALUE](#) function to convert numbers stored as text to numeric value. For example, "3" = 3 would evaluate to `FALSE`, but `VALUE("3") = 3` would evaluate to `TRUE`.

Comparison operators with text values

The most commonly used comparison operator with text is = (equals). It is used to determine if two text strings are the same. Note that like other string functions that perform matching (such as [FIND](#)), it is case sensitive. In other words, it treats "The" as a different string than "the". For the comparison to be true, the two pieces of text must match *exactly*—including capitalization. Use of <> (not equal to) follows the same pattern as use of the = (equals). It is also case sensitive when examining text strings.

It may be surprising to note that even comparisons that include < (less than) and > (greater than)—including <= (less than or equals to) and >= (greater than or equals to)—can be used on text values. Characters are represented by a numeric value and since no two characters are the same, no two characters share the same numeric value.

Predicting the behavior of text comparisons requires some additional information about how printable characters are encoded by computers.

Custom computed column functions

Note

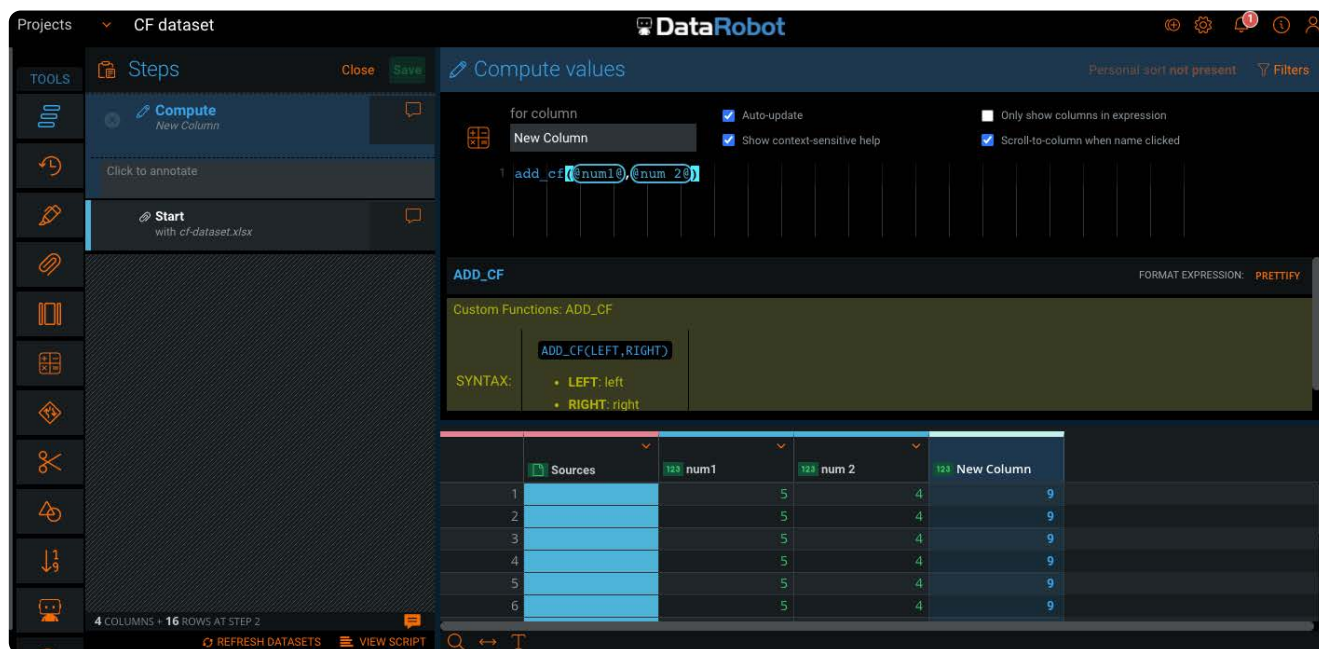
Data Prep custom functions are only available for on-premise and virtual private cloud installations. Your Data Prep Administrator must enable this feature in your application. Contact DataRobot Support for documentation on developing Data Prep custom functions.

If your organization has developed and installed custom functions, you can add new columns by applying a custom function to existing columns in your dataset. In the following section, you'll learn how to use the [Compute tool](#) with a custom function.

Use a custom function

To use custom functions with the Compute tool:

1. In Data Prep, click **compute** in the **Tools** bar.
2. In the **Compute values** pane, name the new column.
3. In the formula row beneath the column name, type the name of a custom function. The usage details for the function display beneath the formula row.
4. Select the columns to be used as variables in the function. Follow the usage guidelines to form your expression.



In this example, the custom function is named `add_cf`, and `num1` and `num2` are columns in the dataset:

```
add_cf(@num1@,@num2@)
```

If the formula contains errors, the Compute tool displays an error message as it does for built-in functions. The errors also display in the Compute step in the Steps tool until the errors are resolved.

5. Verify that the function is operating as intended by viewing the new column.

6. Click **Save** to save the new column.

Work with column data

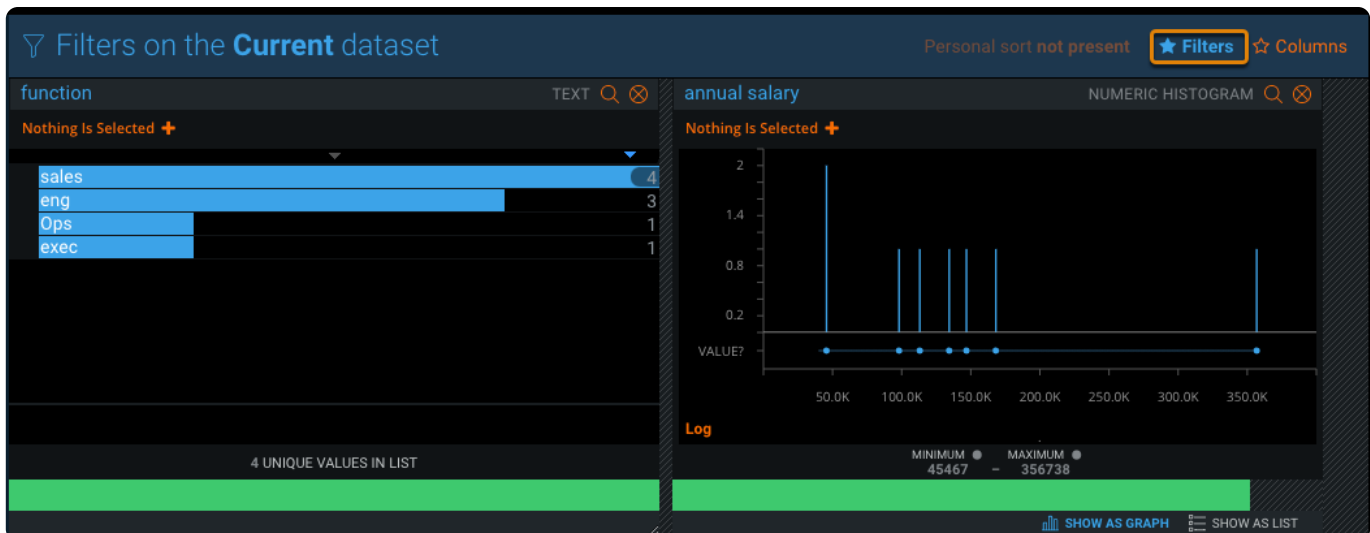
Data Prep provides many methods for working with columns. The topics in this section describe the operations you can perform on column data, for example, highlighting and changing values, finding and replacing values, and filtering your data. To do so, you work with the **Filters** pane, the **View Columns** pane, and the column operations available in the menu above each column.

Tip

This section covers working with column *data*. To manage whole columns, use the **columns** tool on the project **Tools** bar. The **columns** tool lets you update column names, change the order of columns, and remove them from the project. See [Update columns](#) for details.

Filters pane

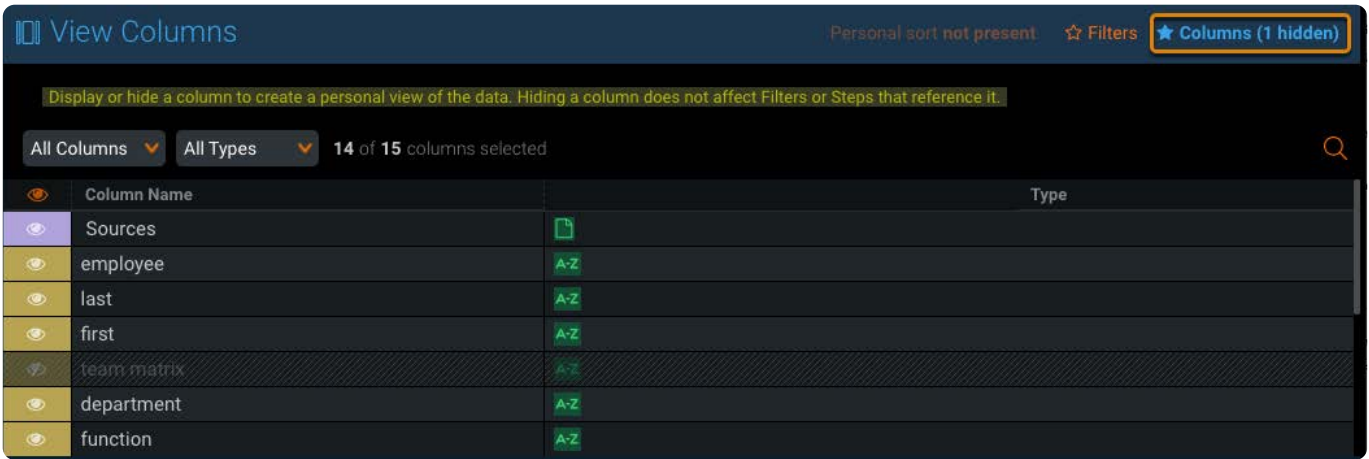
The **Filters** pane displays *Data Filtergrams* for selected columns. A Data Filtergram is a combination of a filter and a histogram. In this example, *function* is a text column and *annual salary* is a numeric column:



Click **Filters** on the top-right of the data preview pane to display the **Filters** pane. See [Data Filtergrams](#) to learn more.

View Columns pane

The **View Columns** pane displays the column types and lets you hide columns. In this example, the *team matrix* column is hidden:



Click **Columns** on the top-right of the data preview pane to display the **View Columns** pane. See [Hide columns](#) to learn how to use the **View Columns** pane to remove columns from the data preview pane. To remove columns permanently from the project, use the [columns tool](#).

Column operations

You access most of the tools for operating on column data in the menu above each column:



These pages describe the column operations:

TOPIC	DESCRIBES...
Transform data types	Learn how Data Prep identifies and transforms data types.
Change column values	Change values in columns, for example, by changing cases, changing data types, trimming spaces, and managing blanks.

TOPIC	DESCRIBES...
Find and replace in columns	Find and replace text within a single column or multiple columns.
Hide columns	Hide columns to create a personal view or to prepare to publish an AnswerSet. If you want to actually remove columns, see Update columns .
Filter data	Explore and filter your data using Data Filtergrams.
Detect and transform date formats	Learn how to work with date formatting in Data Prep.
Split columns	Split columns based on specified character strings, character count, or regular expressions.
Fill columns	Populate blank cells within a column based on the known values directly preceding or following the blanks.
Normalize using cluster and edit	Normalize data and identify inconsistencies and errors in a column.
View column lineage	Identify the project steps that resulted in the selected column.

Data type transformations

Data Prep's data transformation capabilities support heterogeneous data types—meaning that data types are automatically identified at the cell level when importing a dataset into the Data Prep library. The ability to support heterogeneous data types, and within the same Data Prep project column, is powerful because it enables you to bring all of the data into your project. The mixed data types with inherent data quality issues can then be easily homogenized and harmonized with Data Prep, as explained in the Best Practices section below. In other strongly-typed applications that do not support heterogeneous data, the source data must be homogenized, using a different tool, before your data prep work can even begin.

The purpose of this article is to explain when and how Data Prep determines data types for cells and columns, and the best practices for working with heterogeneous data in a Data Prep project.

How does Data Prep identify data types?

When you import a dataset into the Data Prep library, all data, in all cells, is automatically identified as one of the following data types during the import process:

- Numeric
- Boolean
- String or Text
- Date Time (under conditions described below)

Data Prep does this through an algorithm that follows these rules:

1. If the value is null, ignore the value.
2. If the value is, exactly, either "true" or "false", treat the value as a Boolean.
3. If the value can be programmatically read as a number, treat it as Numeric.
4. All other values default to String.

For example, you have a dataset with 10 columns and 1 million rows of data. This translates to a total of 10 million cells. In this case, Data Prep identifies the data type for each one of the 10 million cells following the algorithm rules above.

Then, the data type *for each column* is cast based on the predominant data type that resides within each cell of a column.

Date Time values

As a rule, due to the inherent complexities required to analyze and resolve for the multitude of date time formats, Data Prep does not identify date time values in flat files. However, there are some exceptions to this rule. Under the following conditions, cell data will be recognized as type date time when imported from:

- Database table (JDBC, Hive, etc) *and* the database provides a schema

- Parquet file *and* the Parquet format provides the schema
- Microsoft Excel file *and* the Excel format specifies a data type for each cell

Heterogeneous data

How does Data Prep determine the column type for heterogeneous data within the same column? Returning to the example above, in 1 million rows of data, there's a good possibility that, *within the same column*, there is data belonging to different data types—for example string and numeric values could be mixed in the same column. In this case, Data Prep has further logic to determine how to cast the column's data type. Let's use another very simple example to illustrate that logic.

Here is a column of data that has 15 rows. The first 9 rows are identified as type Numeric, and the remaining 6 values are identified as type String. Notice the column type has been cast as Numeric:

Mixed values	
1	1
2	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9
10	A
11	B
12	C
13	D
14	E
15	F

During the import process, the algorithm calculation is automatically performed to detect that 60% of values are Numeric and 40% Strings. The predominant data type in a column then determines how Data Prep will cast the column type—in this case, the column is cast as type Numeric.

What if there's a tie for data types in a column?

In the event there is a tie—meaning 50% of the column values are one type while 50% are another type—the calculation logic provides these additional rules to break the tie:

COLUMN DATA MIX	WINNER
50% Boolean and 50% Date	Boolean
50% Boolean and 50% Numeric	Boolean
50% Boolean and 50% String	Boolean
50% Date and 50% Numeric	Numeric
50% Date and 50% String	String
50% Numeric and 50% String	String

To summarize, in order of predominance, ties are broken are as follows:

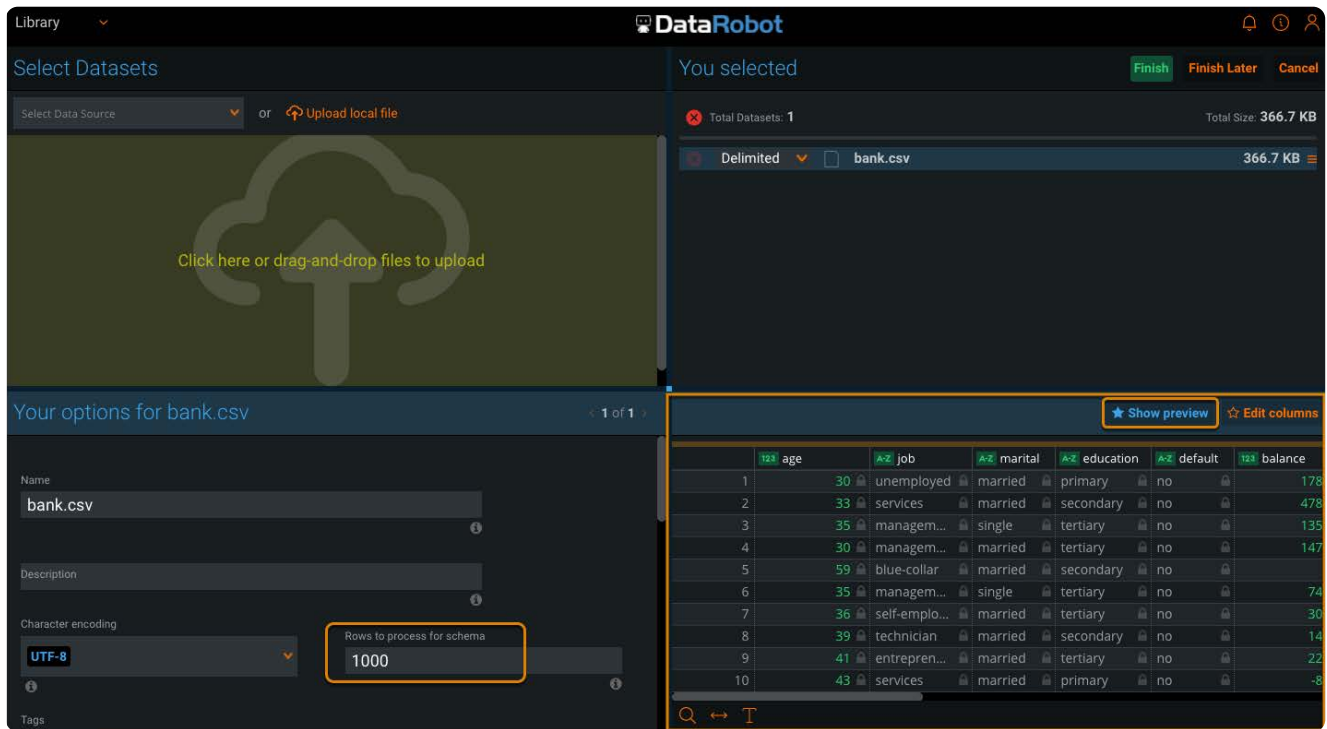
1. Boolean
2. String
3. Numeric
4. Date

Effects of adding data

What happens if the predominant data type in a column changes when you bring in new data to my Data Prep library or an existing project? Column type inferencing and subsequent casting only occurs during the import process to the Data Prep library.

There are two potential scenarios when a column type may not accurately reflect the predominant data type in a column:

- During import into the Data Prep library, the first 1,000 rows of data are used to inference the column type. As a general rule, Data Prep has found that 1,000 rows of data—a configurable value—are sufficient to accurately inference and cast column types for your datasets. These first 1,000 rows are informally known as the "preview" state and it's the state you see in the application while a dataset is loading—either for the first time, or as an updated version for an existing dataset:



There may be unusual cases in which the predominant data type for a column changes after the first 1,000 rows. In this case, the column will remain cast using those first 1,000 rows. You can configure the number of rows in the preview by changing the value of the **Rows to process for schema** field in the options pane. Although the number of rows is configurable during import, Data Prep best practices recommend that you use [Filtergrams](#) to identify and address data quality issues. See [Best Practices](#) for details.

- After a lookup or an append operation in an existing project, the predominant type for a column may change based on the data that comes into the column as a result of the operation. Because inferencing for column type occurs only during the import process, the column type, as it was originally cast, will remain despite the new predominant type. However, Data Prep best practices recommends that, as part of your standard data harmonization practices, you always use [Filtergrams](#) after blending your data from multiple sources to identify and address data quality issues.

Best practices

How do I use Data Prep to locate and remediate data typing issues in my data?

Data Prep was built from the very beginning to identify and address such data quality issues. Typically, as soon as a dataset is imported into the library or appended in a project, the next recommended step is to harmonize the data type so that data quality is enhanced. Data harmonization is one of the key aspects of data preparation and Data Prep provides you with visual indicators and tools like [Filtergrams](#) for your harmonization exercise.

Example

Following an append into an existing project, the predominant data type changes for this column from "numeric" to "string". The column was initially and correctly cast as numeric during the import process and subsequent use in this project. After the append operation, notice the column type remains string type, though the predominant type is now numeric. However, a visual indicator—numeric values are right-justified—allows you to quickly notice a typing disparity in the column:

	Sources	Account Name
1		AAA
2		AAA rentals
3		AAA daily rentals
4		AAA montly rentals
5		1234
6		5678
7		11223344
8		11223344
9		11223344
10		1234
11		1234
12		1234

When opening a [Filtergram](#) on this column, you can quickly determine the values in the column that are not "valid" for this type of column:

Account Name TEXT

NOTHING IS SELECTED +

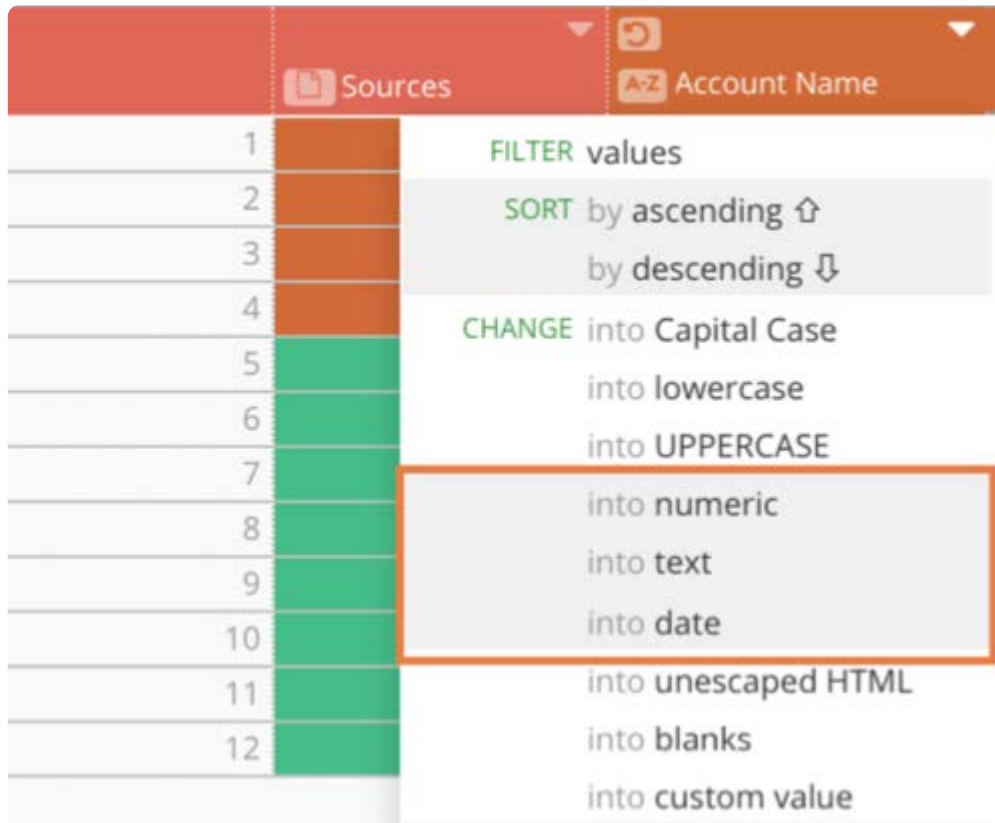
1234	4
11223344	3
AAA rentals	1
AAA montly rentals	1
AAA daily rentals	1
AAA	1
5678	1

7 UNIQUE VALUES IN LIST

VALID + INVALID

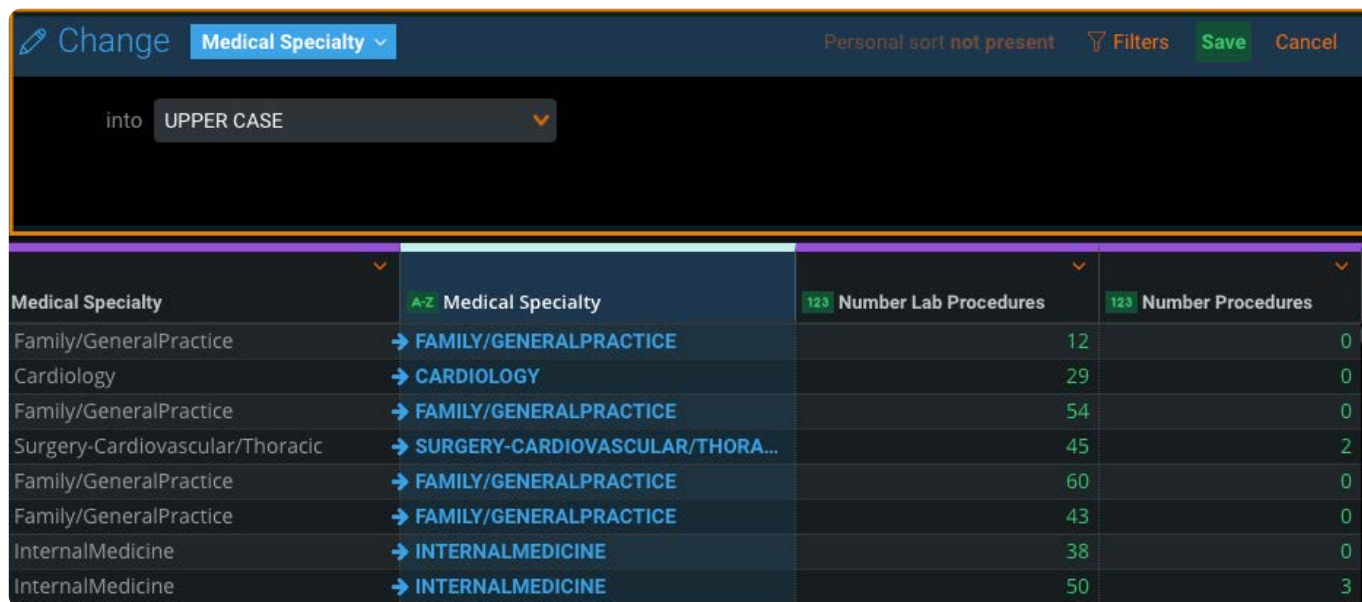
The red bar on the Filtergram indicates there are non-conforming data types. Click the **Invalid** link to view only the values that do not conform.

After you've filtered the display to identify the "invalid" data types—all types other than the predominant type—you can create a [lens](#) in your project to generate an AnswerSet that lists only those non-conforming values. You can then use the AnswerSet to assist in your remediation process for those values. If after reviewing the "invalid" types, you want to convert the column type to another data type, you can do this using the column menu:



Change column values

In Data Prep, you can change data values using the **Change into** operation in the column menu. This example shows a change operation on the "Medical Specialty" column that changes text to uppercase.




Use the **Change into** operation to select one or more columns and then change the data in those columns to:

- Capital case
- Lowercase
- Uppercase
- Numeric values
- Text
- Date
- Unescaped HTML
- Blanks
- Custom values
- Trim leading and trailing spaces from cells in the column
- Collapse consecutive, multiple spaces into a single space

Change values in a single column

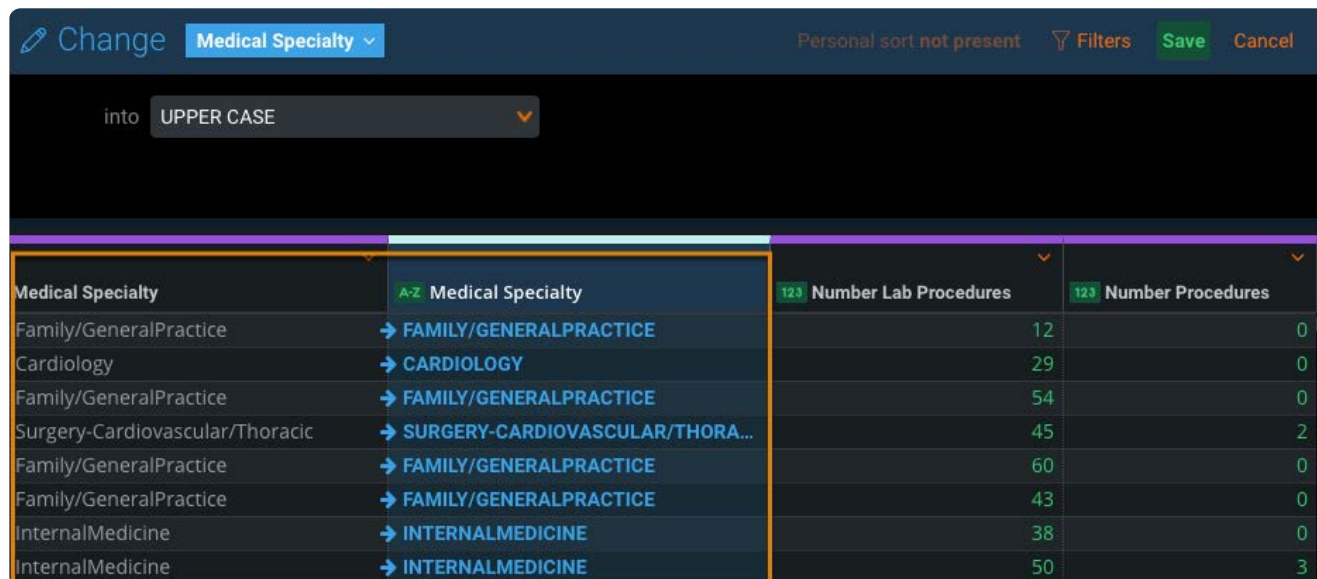
To find and replace text in a single column:

1. Locate the column where you want to change values.

2. Hover over the column menu icon , then hover over **Change into** and select the change you want to make.



Data Prep generates a copy of the original column that reflects the changes, for example:



Medical Specialty	A-Z Medical Specialty	123 Number Lab Procedures	123 Number Procedures
Family/GeneralPractice	→ FAMILY/GENERALPRACTICE	12	0
Cardiology	→ CARDIOLOGY	29	0
Family/GeneralPractice	→ FAMILY/GENERALPRACTICE	54	0
Surgery-Cardiovascular/Thoracic	→ SURGERY-CARDIOVASCULAR/THORA...	45	2
Family/GeneralPractice	→ FAMILY/GENERALPRACTICE	60	0
Family/GeneralPractice	→ FAMILY/GENERALPRACTICE	43	0
InternalMedicine	→ INTERNALMEDICINE	38	0
InternalMedicine	→ INTERNALMEDICINE	50	3


3. Click **Save** at the top to accept the changes.

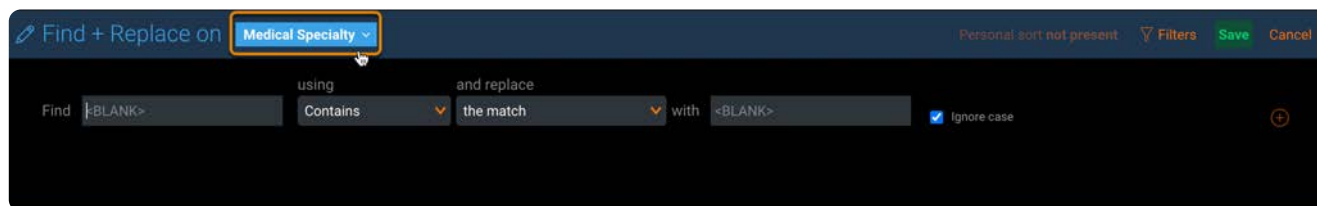
Change values in multiple columns

If you need to perform a change values across your entire dataset or a specific set of columns, you use the *advanced* **Change into** pane. Here are some examples where the advanced function is useful:

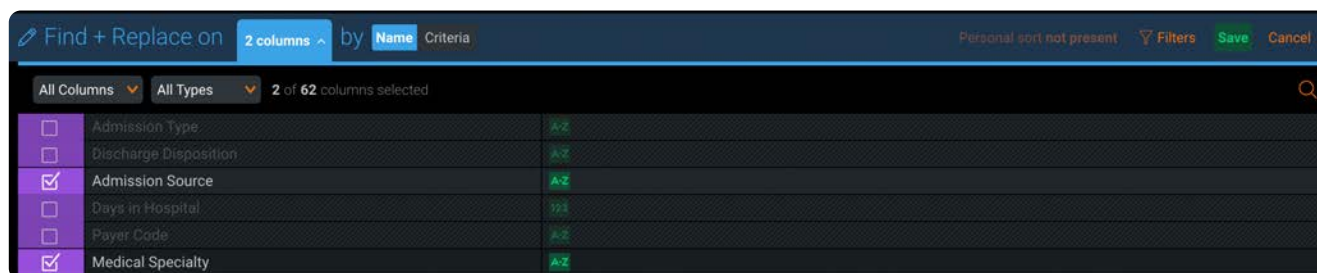
- The dataset has both "incorporated" and "Inc" everywhere. You want to standardize the entire dataset to have only the "Inc" value.
- The dataset has "incorporated" everywhere, and for the most part that's accurate. But you need to change the value to "Inc" for *some specific columns* in the dataset.
- You've pulled two datasets into your Project—one has "NA" and the other uses blanks to represent non-applicable values. You want to change all of the "NA" values into blanks.

To find and replace across multiple columns:

1. Hover over any column's menu icon  and click **find + replace**.
2. Click the column name that appears in the **Find + Replace** pane.



3. In the advanced **Find + Replace** pane that displays, click the check box next to each column that you want to include in the find and replace operation.



The rest of the steps for find and replace across multiple columns are the same as the steps for find and replace for a single column. See [Find and replace](#).

Changing values by Name or Criteria

In the advanced **Find + Replace** pane, you can select multiple columns by either *Name* or *Criteria*.

Change values by Name:

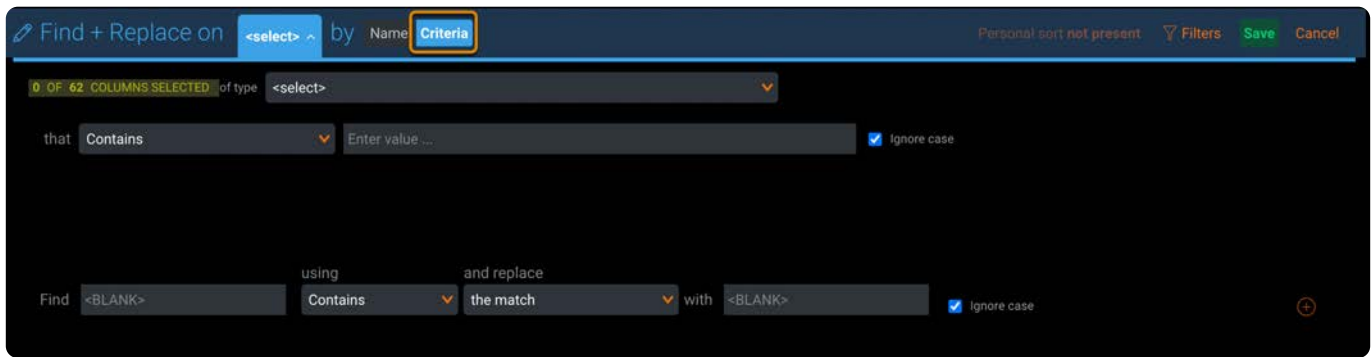
Finding and replacing by *Name* applies the replace operation *only* to the specific columns you select.

To select columns by *Name*:

- Click the check box adjacent to the column(s) that you want to select.
- Click the top-most check box to select all columns.
- Use the Columns and Types filters at the top of the panel to quickly filter down to the columns you want to select for the operation.
- Use the search function to locate a column by name.

Change values by Criteria:

Finding and replacing by *Criteria* applies the replace operation to *any column* that meets the criteria you specify.



For example, if you have String type columns in your dataset and you specify the replace operation for String type columns, then all existing columns of this type in your dataset—and any new String type columns that are introduced to the dataset prior to this Step—will be dynamically replaced.

To select columns based on criteria:

- Optionally specify the data type of the column—Boolean, DateTime, Number or String.
- Optionally specify the pattern for the column name—contains, starts with, equals or ends with.

Notice the header message updates to indicate the number of columns you have selected based on that criteria. You may later notice the number of selected columns increases or decreases if new data is brought into an earlier Step that introduces or removes columns that meet your criteria.

Note

If you switch between the **Name** and **Criteria** options before saving the replace operation, Data Prep retains your selections and provides a **Restore last selection** link that returns you to your initial selection method.

Example: Change into Numeric

This column operation converts all numbers stored as text strings into numeric values. By doing this, mathematical operations can be performed on values in this column when, as numbers stored as text, these actions would otherwise be considered invalid.

Numbers stored as strings appear left-aligned within a cell and in black text; numbers stored as numeric values are right-aligned and appear in green.

When this operation is applied to cells that cannot be converted to numeric values, it will have no effect. In a column with both text and numbers on different rows, only those rows that can be converted will be changed.

If a value that appears suitable for conversion is not successfully converted, it is likely that there are non-number characters somewhere in the cell.

The following are examples of characters that can inhibit the transformation:

- Leading or trailing spaces. These can be removed by using the column operation for **"White Space trim leading and trailing"** before you apply the **"Transform into numeric operation"**.

The "White Space trim leading and trailing" operation examines all rows for spaces at both the beginning and end of the text string. Where it finds them, they are removed—leaving only the value in the cell.

- Intermediate characters (such as commas or spaces.) Operations such as [Column split](#) or a [Compute columns](#) that uses [REGEX](#) may be required, first, in order to successfully create a column of numeric values.

Note

A single period (".") in a cell of numbers will be interpreted as a decimal point. These strings will be able to be converted into numeric values without requiring any other operations.

Find and replace in columns


The Data Prep find and replace operations allows you to locate and replace text within specified columns. This example shows a find and replace operation on the "Medical Specialty" column where ? values are replaced by blanks.

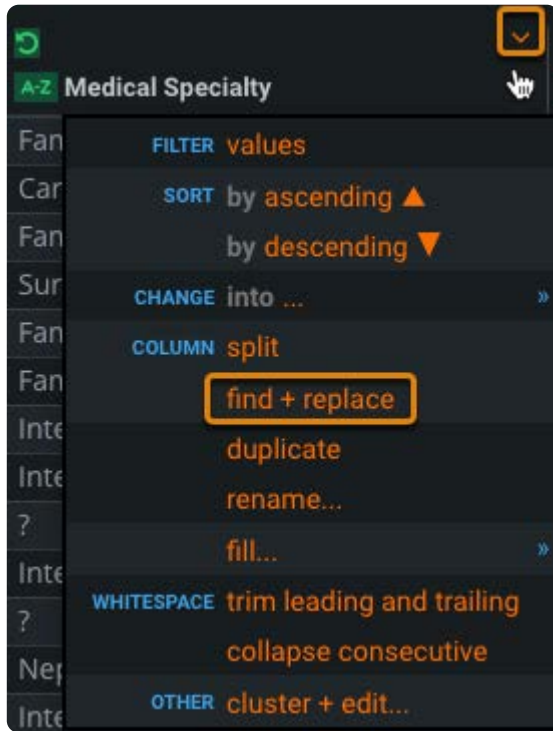
Medical Specialty	Number Lab Procedures	Number Procedures	Number Me
Family/GeneralPractice	12	0	
Cardiology	29	0	
Family/GeneralPractice	54	0	
Surgery-Cardiovascular/Thoracic	45	2	
Family/GeneralPractice	60	0	
Family/GeneralPractice	43	0	
InternalMedicine	38	0	
InternalMedicine	50	3	
?	66	1	
InternalMedicine	59	0	
?	50	0	
Nephrology	17	1	
InternalMedicine	64	0	
?	46	0	
?	51	0	
?	49	3	

You can also find and replace across multiple or all columns. See [Multiple find and replace operations](#) for details.

Find and replace

To find and replace text in a single column:

1. Locate the column where you want to find and replace values.
2. Hover over the column menu icon  and click **find + replace**.



3. In the **Find** field, specify a value to find.

Alternatively, double-click a cell that has the value you want to find; the fields for find and replace are then automatically populated with the value from that cell.

4. In the **using** field, select how to match the text specified in the **Find** field:

- **Contains:** The match can be anywhere in the cell value.
- **Starts With:** The match must be at the beginning of the cell value.
- **Equals:** The match must be exact.
- **Ends With:** The match must be at the end of the cell value.

5. In the **and replace** field, select what portion of the match to replace:

- **the entire cell:** Replaces the entire contents of the cell.
- **the match:** Replaces only the matched portion of the cell.

Data Prep generates a copy of the original column that reflects the changes, for example:

Find + Replace on **Medical Specialty** Personal sort not present

using **Contains** and replace **the match** with **<BLANK>** Ignore case

A-Z Medical Specialty	A-Z Medical Specialty	123 Number Lab Procedures	123 Number Procedures	123 Number Me
Family/GeneralPractice	→ Family/GeneralPractice	12	0	
Cardiology	→ Cardiology	29	0	
Family/GeneralPractice	→ Family/GeneralPractice	54	0	
Surgery-Cardiovascular/Thoracic	→ Surgery-Cardiovascular/Thoracic	45	2	
Family/GeneralPractice	→ Family/GeneralPractice	60	0	
Family/GeneralPractice	→ Family/GeneralPractice	43	0	
InternalMedicine	→ InternalMedicine	38	0	
InternalMedicine	→ InternalMedicine	50	3	
?	→	66	1	
InternalMedicine	→ InternalMedicine	59	0	
?	→	50	0	
Nephrology	→ Nephrology	17	1	
InternalMedicine	→ InternalMedicine	64	0	
?	→	46	0	
?	→	51	0	
?	→	49	3	

6. Click **Save** at the top to accept the changes.

Examples

Example 1: original cell value = "123456" Find: "123" and replace with Match "321" and the result is: "321456"

Example 2: original cell value = "123456" Find: "123" and replace Entire Cell with "321" and the result is: "321"

Tip

If the **Steps** tool is turned off, the **Save** button is displayed in the **Find + Replace** pane. If you enable the **Steps** tool, the **Save** button is located at the top of the **Steps** pane.

Multiple find and replace operations

If you want to continue making replacements and save all of the find and replacement transformations to the current step, click the plus (+) to continue.

Find + Replace on **Medical Specialty** Personal sort not present Filters Save Cancel

using **Contains** and replace **the match** with **<BLANK>** Ignore case +

Each additional find and replace operation is iterative—a subsequent find operates on the results of the previous find and replace operation and all of the transformations are applied in order, from top to bottom as they are displayed in the **Find + Replace** pane.

Examples of multiple find and replace operations

In the following example, the first operation finds all instances of "memorial" and then replaces the entire cell—not just the matched text—with a null value:

Find + Replace On Hospital Name Personal data not present Filters Save Cancel

Find using and replace with Ignore case

Find using and replace with Ignore case

Hospital Name	Hospital Name	Address	City	State
SOUTHEAST ALABAMA MC	SOUTHEAST ALABAMA Medical Center	1108 ROSS CLARK CIRCLE	DOTHAN	AL
MARSHALL MC SOUTH	MARSHALL Medical Center SOUTH	2505 U S HIGHWAY 431 NORTH	BOAZ	AL
ELIZA COFFEE MEMORIAL HOSPITAL	→ null	205 MARENGO STREET	FLORENCE	AL
MIZELL MEMORIAL HOSPITAL	→ null	702 N MAIN ST	OPP	AL
CRENSHAW COMMUNITY HOSPITAL	CRENSHAW COMMUNITY HOSPITAL	101 HOSPITAL CIRCLE	LUVERNE	AL
ST VINCENT'S EAST	ST VINCENT'S EAST	50 MEDICAL PARK EAST DRIVE	BIRMINGHAM	AL
DEKALB REGIONAL MC	DEKALB REGIONAL Medical Center	200 MED CENTER DRIVE	FORT PAYNE	AL
SHELBY BAPTIST MC	SHELBY BAPTIST Medical Center	1000 FIRST STREET NORTH	ALABASTER	AL
CALLAHAN EYE HOSPITAL	CALLAHAN EYE HOSPITAL	1720 UNIVERSITY BLVD	BIRMINGHAM	AL
HELEN KELLER MEMORIAL HOSPITAL	→ null	1300 SOUTH MONTGOMERY AVENUE	SHEFFIELD	AL

Another example:

find "Detroit" and replace with "San Francisco" find "San Francisco" and replace with "San Jose" The result is that "Detroit" is transformed into "San Jose"

Note that you can rearrange the sequence of your transformations at any time by dragging and dropping a find and replace operation to another position in the panel.

This dotted line separates your list of transformations. You can click any transformation in the list and drag it to a different place in the order.

Find + Replace On Hospital Name Personal data not present Filters Save Cancel

Find using and replace with Ignore case

Find using and replace with Ignore case

Hospital Name	Hospital Name	Address	City	State
SOUTHEAST ALABAMA MC	SOUTHEAST ALABAMA Medical Center	1108 ROSS CLARK CIRCLE	DOTHAN	AL
MARSHALL MC SOUTH	MARSHALL Medical Center SOUTH	2505 U S HIGHWAY 431 NORTH	BOAZ	AL
ELIZA COFFEE MEMORIAL HOSPITAL	→ null	205 MARENGO STREET	FLORENCE	AL
MIZELL MEMORIAL HOSPITAL	→ null	702 N MAIN ST	OPP	AL
CRENSHAW COMMUNITY HOSPITAL	CRENSHAW COMMUNITY HOSPITAL	101 HOSPITAL CIRCLE	LUVERNE	AL
ST VINCENT'S EAST	ST VINCENT'S EAST	50 MEDICAL PARK EAST DRIVE	BIRMINGHAM	AL
DEKALB REGIONAL MC	DEKALB REGIONAL Medical Center	200 MED CENTER DRIVE	FORT PAYNE	AL
SHELBY BAPTIST MC	SHELBY BAPTIST Medical Center	1000 FIRST STREET NORTH	ALABASTER	AL

To save all find and replace transformations for this step, click **Save**. Or to remove any find and replace operation for this step, click the minus button for that operation.

Important notes for find and replace


- By default, capitalization is ignored in the search. For example, the text cat matches the text CaT. If capitalization is important, deselect the Ignore case checkbox.
- By default, you are limited to 250 find and replace operations in a single project step. If you exceed this number, an error message is displayed in the **Step** tool and you cannot continue find and replace transformations on the column. Contact your system administrator if you need to increase the limit.
- The **highlight** feature on the grid is only enabled for a single **Find + Replace** transformation; it is disabled if you add another **Find + Replace** transformation.

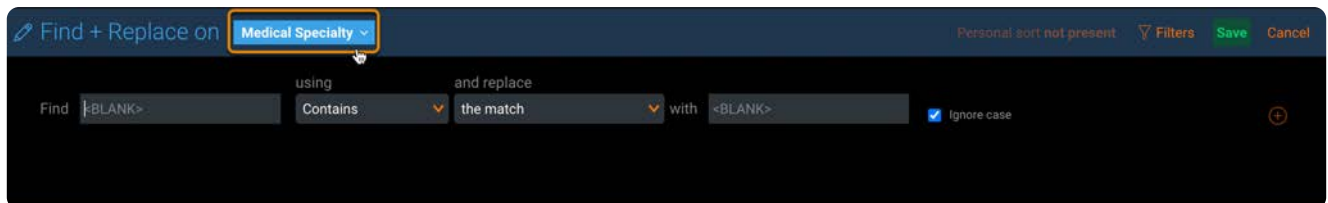
Find and replace across multiple columns

If you need to perform a search and replace across your entire dataset or a specific set of columns, you use the *advanced Find + Replace* pane. Here are some examples where the advanced function is useful:

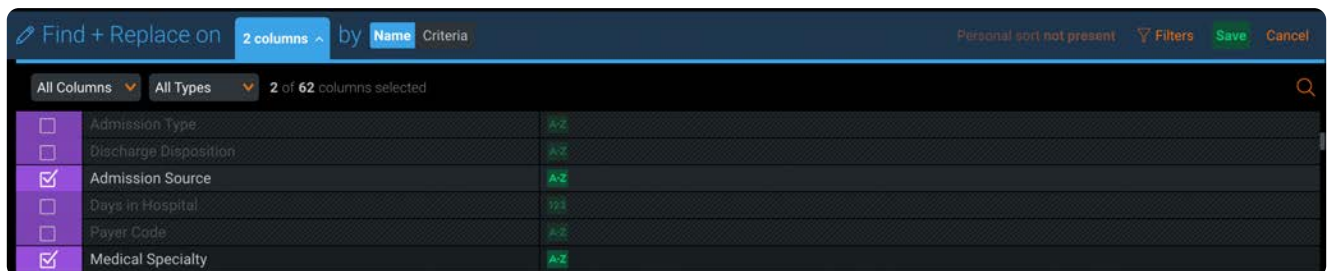
- The dataset has both "incorporated" and "Inc" everywhere. You want to standardize the entire dataset to have only the "Inc" value.
- The dataset has "incorporated" everywhere, and for the most part that's accurate. But you need to change the value to "Inc" for *some specific columns* in the dataset.
- You've pulled two datasets into your project—one has "NA" and the other uses blanks to represent non-applicable values. You want to change all of the "NA" values into blanks.

To find and replace across multiple columns:

1. Hover over any column's menu icon  and click **find + replace**.
2. Click the column name that appears in the **Find + Replace** pane.



3. In the advanced **Find + Replace** pane that displays, click the check box next to each column that you want to include in the find and replace operation.



The rest of the steps for find and replace across multiple columns are the same as the steps for find and replace for a single column. See [Find and replace](#).

Finding and replacing by Name or Criteria

In the advanced **Find + Replace** pane, you can select multiple columns by either *Name* or *Criteria*.

Find and replace by Name

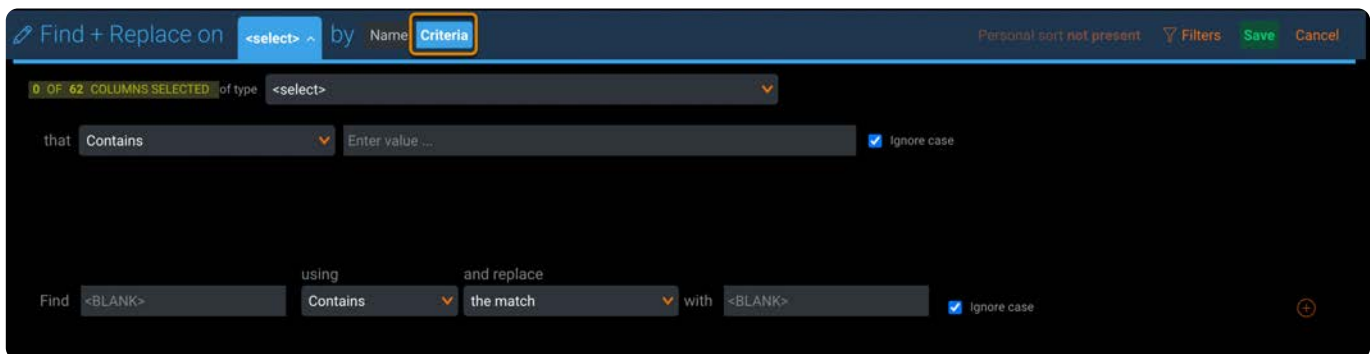
Finding and replacing by *Name* applies the replace operation *only* to the specific columns you select.

To select columns by *Name*:

- Click the check box adjacent to the column(s) that you want to select.
- Click the top-most check box to select all columns.
- Use the Columns and Types filters at the top of the panel to quickly filter down to the columns you want to select for the operation.
- Use the search function to locate a column by name.

Find and replace by Criteria

Finding and replacing by *Criteria* applies the replace operation to *any column* that meets the criteria you specify.



For example, if you have String type columns in your dataset and you specify the replace operation for String type columns, then all existing columns of this type in your dataset—and any new String type columns that are introduced to the dataset prior to this step—will be dynamically replaced.

To select columns based on criteria:

- Optionally specify the data type of the column—Boolean, DateTime, Number or String.
- Optionally specify the pattern for the column name—contains, starts with, equals or ends with.

Notice the header message updates to indicate the number of columns you have selected based on that criteria. You may later notice the number of selected columns increases or decreases if new data is brought into an earlier step that introduces or removes columns that meet your criteria.

Note

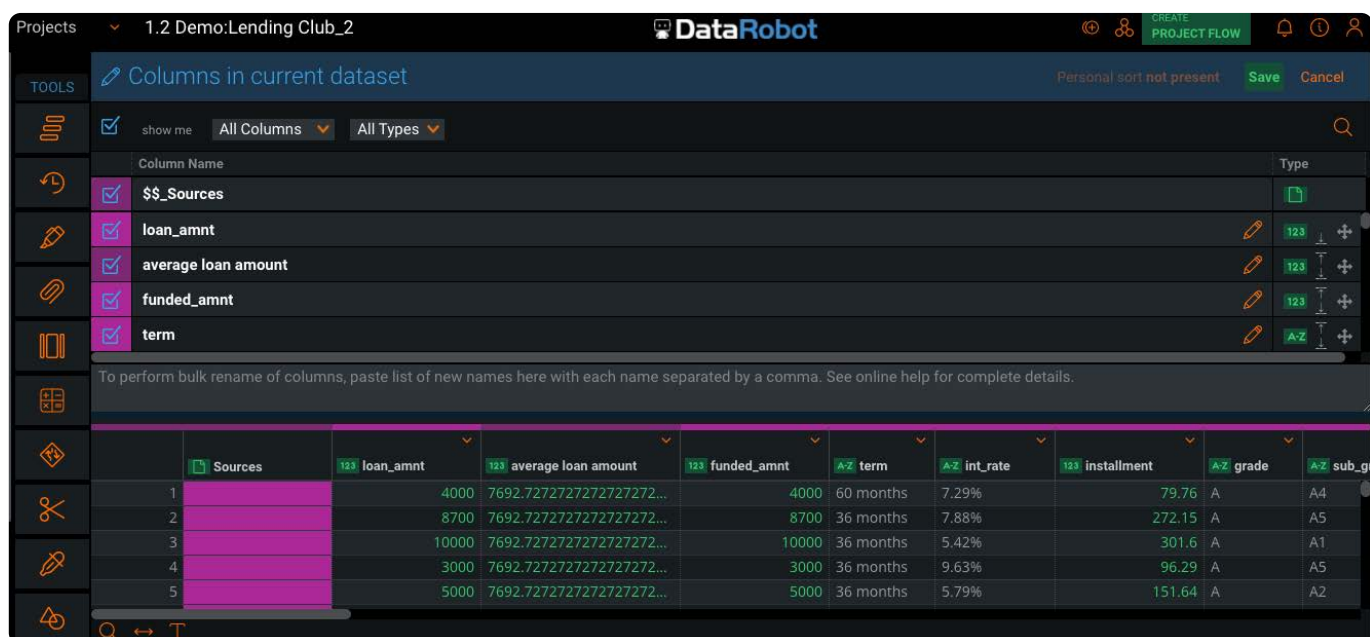
If you switch between the **Name** and **Criteria** options before saving the replace operation, Data Prep retains your selections and provides a **Restore last selection** link that returns you to your initial selection method.

Hide columns

There may be times when you want to simply hide, but not remove, columns in a project:

- **To create a personal view:** You want to hide columns from the grid so they don't clutter your view but you need to keep them in your data because they are used for calculating project steps. For example, if your work in a project is primarily with numeric columns, you could create a personal view of the data that hides all column types other than numeric.
- **To publish an AnswerSet that includes only specific columns:** You may want to publish a custom AnswerSet that only displays some columns from the data. First create a lens on the step where you want to create the AnswerSet, then hide the columns you do not want to show in the AnswerSet. Save the lens and your selections are saved. When you publish from the lens, only the columns you selected to show are published in the AnswerSet.

This example shows the **Columns in current dataset** pane with the `loan_amnt` and `funded_amnt` columns hidden:



Note

To remove columns permanently rather than hiding them temporarily, see [Remove columns](#).

Hide columns

To hide columns:

1. Click **Columns** in the upper right.
2. Click the eye icon on the left for any column you want to hide.

The selected eye icons are grayed out and the selected columns no longer appear in the **Preview Display**.

Considerations

Hiding a column does not affect filters or steps in the project. As soon as you begin to make edits in the project, you automatically exit personal view mode. This means columns you hide here for your personal view are still included and displayed in the live data—with the exception of using the Lens tool because when you publish from the lens, only the columns set to show are published in the AnswerSet.

Options for hiding columns

- You can set options to filter the columns you want to hide. You can filter for **All Columns** and **All Types**. Use these filter menus to view only the columns that you have either selected to show or hide, and to control which column types (String, Number, DateTime) are being displayed in the **View Columns** pane.
- You can select a contiguous group of columns to hide. To do so, click the eye icon adjacent to the column you want to hide, then hold the Shift key and click the eye icon for the column that brackets the group you want to select.
- You can search for a column in your data. To do so, click the magnifying glass icon in the upper right of the pane and begin typing a column name.
- You can bring a column into focus on in the **Preview Display**. To do so, hold down the Shift key as you hover over a column name. That column then comes into view and is highlighted in the **Preview Display**.

Tip

If you share a project and also want to share your personal view of the data, use the **Lens** tool to capture your view. The lens captures your personal view of the **Preview Display** and anyone else working in the project can see your personal view by clicking on the lens in the **Steps** pane.


Filter data

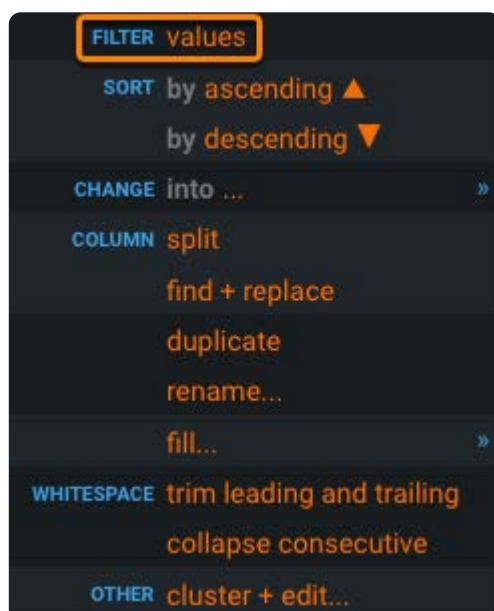
A Data Filtergram™ is both a filter that helps you transform your data and a histogram that helps you visualize your data. Filtergrams allow you to visualize your data before, during, and after every transformation. You can use Filtergrams to:

- Explore your data. When you use filters in this way, you can view your filtered selections in the **Preview Display**. Your selections are not saved to steps in the Project. Rather, they let you quickly explore your data. When you're done with your exploration, you can remove the filters and continue with transformation operations in your Project.
- **Remove rows**. If your data prep work involves removing rows, you need to first use filters to identify the rows you want to remove.
- Change the data in just a particular set of rows. In this case, you use filters to isolate the rows where you want to make changes. Then you can apply your changes to just those rows. For example, if you have a column with first names and you want to change all instances of "Anna" to "Ann," then you filter on the column to display only the "Anna" values. Then you apply a **Change into Custom Value** column operation to convert "Anna" to "Ann."
- Publish just a particular set of rows. In this case, use filters to isolate the rows you want to publish. Then you can add a lens to create the publishing point.

Create a Data Filtergram

To create a Data Filtergram:

1. Locate the column you want to filter.
2. Hover over the column menu icon  and click **Filter values**.



Note that you can have Filtergrams for multiple columns open so that you can preview the results of those filters on the **Preview Display**.

There are five types of Filtergrams that correspond to each type of column data:

- Text
- Numeric
- Date and Time
- Boolean
- Source

The following sections describe how to use filters of each type.

Text Filtergrams

The Text Filtergram pane displays a list of each distinct text value that appears in your dataset. The bars running left-to-right provide a histogram view of the relative number of times each value occurs. The total number of unique values is listed on the bottom left side of the pane. From the list, you can select values to dynamically display in your dataset.

The screenshot shows a 'Last Name' filtergram for a 'TEXT' column. It displays a list of 207 unique values with horizontal bars representing their frequency. The top of the pane shows 'SHOW 4 SELECTED ITEMS' and a search icon. The bottom of the pane shows '207 UNIQUE VALUES IN LIST', 'CLEAR', and 'INVERT' buttons, along with a filter type selector set to 'TYPE' and buttons for '+ OTHER', '+ BLANK', and '+ ERROR'.

Last Name	Count
Kelly	25
Jackson	25
Roberts	24
Richardson	24
Adams	24
Rice	23
Morales	23
Mason	23
Morgan	22
Elliott	22
Sanchez	21
McDonald	21
Armstrong	19

Callouts:

- Length of bars provide histogram view of relative occurrence
- Click to select values for filtering your dataset.
- Number of unique values in list

See [Work with Text Filtergrams](#) for actions you can take with Text Filtergrams.



If you have non-text values in the column, this color-coded bar provides an additional histogram to indicate the relative occurrence of each value type in the column:
 green = type Text
 gray = blank cells
 red = both Other (non-Text) types and cell Errors

a. Currently Selected (upper left): When you make selections from the list, the button's label changes to indicate your number of selections. Click the button to open a new pane that lists every selected value. From this pane, you can continue to refine the text values you want to filter on in your dataset. The actions you can take from this pane are described in the section [Work with Text Filtergrams](#).

b. Sort list order: by default, the list of text values displays from highest to lowest count. To reverse the order to display least to most, click the triangle in the upper right corner above the count column. You can also sort the list alphabetically by clicking the triangle located above the list. The triangle's orange color indicates which sort order (numeric or alpha) is currently applied to the dataset.

c. Color-coded filter bar: When moving your mouse anywhere over the Filtergram pane, the following buttons display:

- **Type:** When you mouse over this button, a ratio displays to express the count of currently selected text type rows relative to all text type rows in the dataset. If nothing is selected in the list, the ratio reflects the total count for text type rows relative to all rows in the dataset. When you click this button, you dynamically hide these text values in your dataset. This is useful if you have blanks, errors or other non-text values in this column and want to view only those data types.
- **Other:** If there are non-text values in the column, for example numbers, this button displays. When you mouse over this button, a ratio displays to express the count of currently selected rows with non-text values relative to all rows with non-text values in the dataset. If nothing is selected on the histogram, the ratio reflects the total count of rows with non-text values relative to all rows in the dataset. Click **Other** to toggle off and hide the other values from the current dataset view.
- **Blank:** If there are blanks in the column, this button displays. When you mouse over this button, a ratio displays to express the count of currently selected blank rows relative to all blank rows in the dataset. If nothing is selected on the histogram,

the ratio reflects the total count of blank rows relative to all rows in the dataset. Click **Blank** to toggle off and hide the blank cells from the current dataset view.

- **Error**: If there are errors in the column, this button displays. When you mouse over this button, a ratio displays to express the count of currently selected rows with cell errors relative to all rows with errors in the dataset. If nothing is selected on the histogram, the ratio reflects the total count of rows with cell errors relative to all rows in the dataset. Click **Error** to toggle off and hide the cell errors from the current dataset view.



If you have "other" values, blanks or errors in the column, notice there is also a + button adjacent to the button for that value. When you click the + button you add all values of that type, for example all "other" values, to the "Selected Items" list where you can continue working with them. See [Work with Text Filtergrams](#) for the filtering actions you can take from **Selected Items**.

Work with Text Filtergrams

The Text Filtergram pane provides tools that enable you to dynamically filter your data with great precision. This section describes the actions you can take.

The top screenshot shows a list of names with corresponding bar lengths. The names are: Kelly (25), Jackson (25), Roberts (24), Richardson (24), Adams (24), Rice (23), Morales (23), Mason (23), Morgan (22), Elliott (22), Sanchez (21), McDonald (21), and Armstrong (19). A blue circle 'a.' is next to the first four names, and a blue circle 'b.' is at the top left. A blue circle 'c.' is at the bottom right. The bottom screenshot shows the 'SELECTED VALUES' list with a 'SAVE' button. The selected values are: Richardson, Adams, Rice, and Morales. A blue circle 'b1.' is next to the list, a blue circle 'b2.' is next to the 'SAVE' button, a blue circle 'b3.' is at the top right, and a blue circle 'b4.' is at the bottom right. A blue circle 'b5.' is at the top left. A large blue arrow points from the top screenshot to the bottom one.

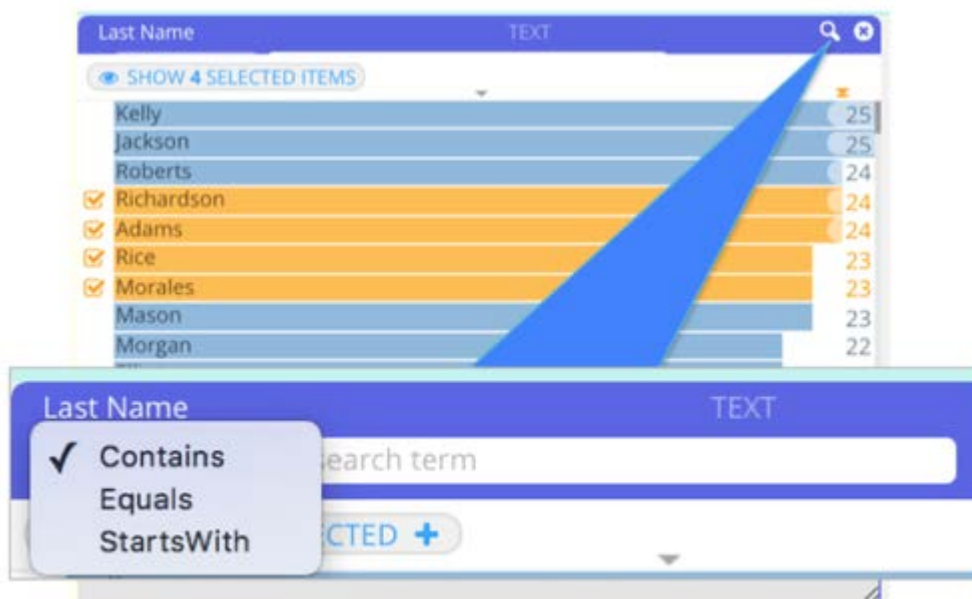
- a. Select values from the list to display in your dataset: Click any text value in the pane to dynamically filter your dataset to display only that value. To select more than one value: CTRL+click (Windows) or Command+click (Mac). To select a contiguous range: Shift+click. To remove any selection: Alt+click.

b. Show Selected Items: After selecting text values for filtering, click **Show Selected Items**. A new pane opens and displays every selected Range and Value currently displayed in your dataset. Note: a Range for text fields is defined by ASCII sort order. From this pane, you can take the following actions:

- **b1.** Manually remove any of the Ranges or Values in the dataset by clicking to remove the orange check mark adjacent to the Range or Value. Click to return that range or value to your filtered dataset. Note: a trash can icon appears at the top of the pane when you deselect a range or value. This allows you to discard the selection from your filters.
- **b2.** Manually update any of the Ranges or Values listed here by clicking the value you want to edit. The value then becomes editable. Enter a new value and click **Save**.
- **b3.** Manually add other Ranges or Values for filtering the data by clicking the **+** button. Provide a Min and Max for the Range, or provide a Value, and click **Save**. Your dataset dynamically updates to reflect your additions.
- **b4.** Click EXCLUDE to exclude a Range or Value from the dataset. This is particularly useful if you have already selected a range. From the range, you can exclude specific values (from within that range) to hide from the current dataset. Note that ranges and values marked with EXCLUDE are displayed with orange, dotted outlines in the histogram to remind you of their exclusion.
- **b5.** When you are finished working in this pane, click **Hide Selected Items** to return to the list view.

c. Clear and Invert: Clear removes all of your current filters. Invert displays all of your data *except* that which you have *selected to filter*.

You can also search for values. To do so, click the magnifying glass icon in the upper right to open a search field. Search for a value that contains, equals, or starts with any value you specify:



Numeric Filtergrams

When you open the Numeric Filtergram pane, there are two views that you can use for performing filtering operations on your data:

- **Show as Graph** (the default view) This histogram displays a numeric histogram representing the distribution of numeric values in your dataset. The horizontal line, or x-axis, represents the range of values that occur in your dataset within this

column. The height of each bar represents the count of each distinct value in the column. The scale along the left side, or y-axis, provides a measure for the count. The filtering actions you can take with the histogram are described below in the section [Work with Numeric Filtergrams](#).



If you have non-numeric values in the column, this color-coded bar provides an additional histogram to indicate the relative occurrence of each value type in the column:
 green = type Numeric
 gray = blank cells
 red = both Other (non-Numeric) types and cell Errors

- **Show as List** Click this tab to hide the numeric histogram and display the unique occurrence of every value in the column. From the list, you select values to dynamically display in your dataset. The filtering actions you can take with the list are described in the section [Work with Numeric Filtergrams](#).



When moving your mouse anywhere over the Filtergram pane, the following buttons display:



a. Currently Selected button (top left corner): When you make selections on the histogram, the button's label changes to indicate your number of selections. Click the button to open a new pane that lists every selected Range and Value currently selected in the histogram. From this pane you can continue to refine the numeric values you want to filter on in your dataset. If you already know the ranges and values that you want to filter on in your dataset, you can click this button instead of using the histogram. From the new pane you can enter values and ranges to begin your filtering operations. The actions you can take from this pane are described in the section [Work with Numeric Filtergrams](#).

b. Log button (bottom left corner): Toggles on a logarithmic scale (log) view of your data. By default, a linear view of the data is displayed in the Filtergram. However, if you have a large numerical range in which one or a few points in the data are much larger than the bulk of the data, the log view adjusts for skewness in your data.

c. Show Zoom/Pan (bottom left corner): Toggles on and off an overview tool used for viewing the relative position of values and ranges that you have zoomed into on the numeric histogram. Zooming and panning operations are described below in the section [Show numeric columns as graphs](#).

d. When moving your mouse anywhere over the Filtergram pane, the following buttons display:

- **Type:** When you mouse over this button, a ratio displays to express the count of currently selected numeric type rows relative to all numeric type rows in the dataset. If nothing is selected on the histogram, the ratio reflects the total count for numeric type rows relative to all rows in the dataset. When you click this button, you dynamically hide these numeric values in your dataset. This is useful if you have blanks, errors or other non-numeric values in this column and want to view only those data types.
- **Other:** If there are non-numeric values in the column, for example text values, this button displays. When you mouse over this button, a ratio displays to express the count of currently selected rows with non-numeric values relative to all rows with non-numeric values in the dataset. If nothing is selected on the histogram, the ratio reflects the total count of rows with non-numeric values relative to all rows in the dataset. Click **Other** to toggle off and hide the other values from the current dataset view.
- **Blank:** If there are blanks in the column, this button displays. When you mouse over this button, a ratio displays to express the count of currently selected Blank rows relative to all Blank rows in the dataset. If nothing is selected on the histogram, the ratio reflects the total count of Blank rows relative to all rows in the dataset. Click this button to toggle off and hide the blank cells from the current dataset view.

Error: If there are errors in the column, this button displays. When you mouse over this button, a ratio displays to express the count of currently selected rows with cell errors relative to all rows with errors in the dataset. If nothing is selected on the histogram, the ratio reflects the total count of rows with cell errors relative to all rows in the dataset. Click this button to toggle off and hide the cell errors from the current dataset view.



If you have "other" values, blanks or errors in the column, notice there is also a + button adjacent to the button for that value. When you click the + button you add all values of that type, for example all "other" values, to the "Selected Items" list where you can continue working with them. See [Show numeric columns as graphs](#) for the filtering actions you can take from **Selected Items**.

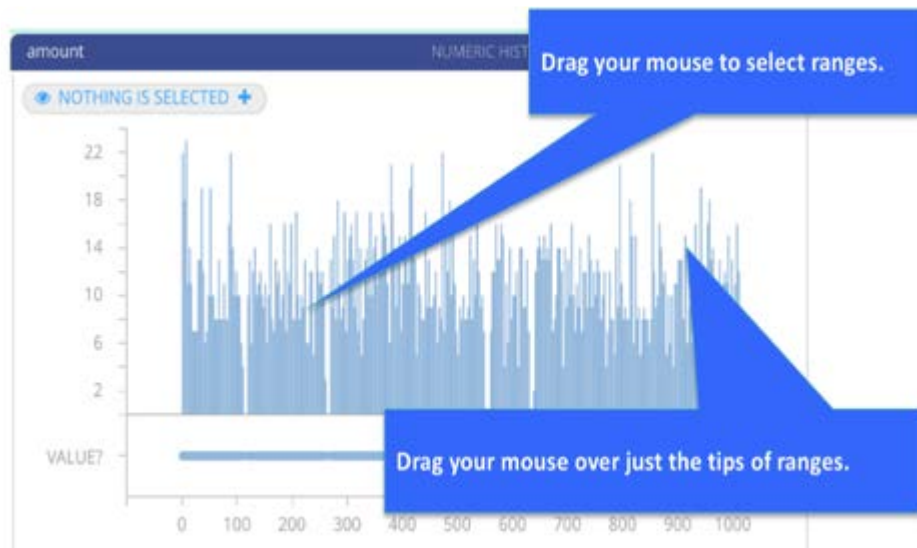
Work with Numeric Filtergrams

The Numeric Filtergram provides tools that enable you to dynamically filter your data with great precision. This section describes the actions you can take.

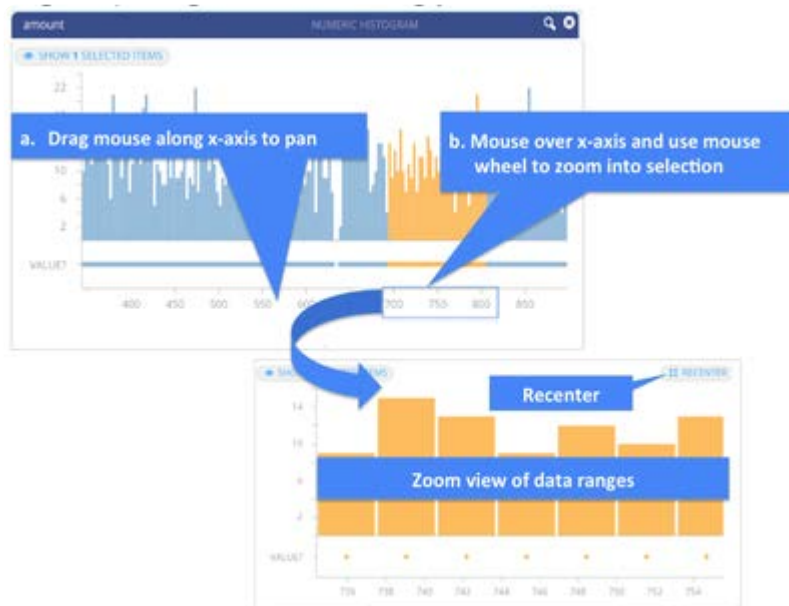
Show numeric columns as graphs

- Select ranges to view on the histogram: click and drag your mouse across a range of values, and notice that your dataset updates to reflect your selection. To select additional, non-contiguous ranges in the histogram, use the following key commands while dragging your mouse: CTRL+click (Windows) or Command+click (Mac). To remove any selection or portion of a selection, use the following key command while dragging your mouse: Alt+click.

You can also drag your mouse over the tips of ranges in the histogram to display only those values in your dataset. The y-axis is helpful in determining the relative values of peaks in your data:



- Begin exploring and transforming your data:



a. Pan across the histogram: Move your mouse over the values on the x-axis. Notice your cursor changes from the pointer to the click and drag symbol. Click and drag across the values on the x-axis to pan across the histogram and your selections. Click **Recenter** to return the histogram to its default view. Note that your selections are retained.

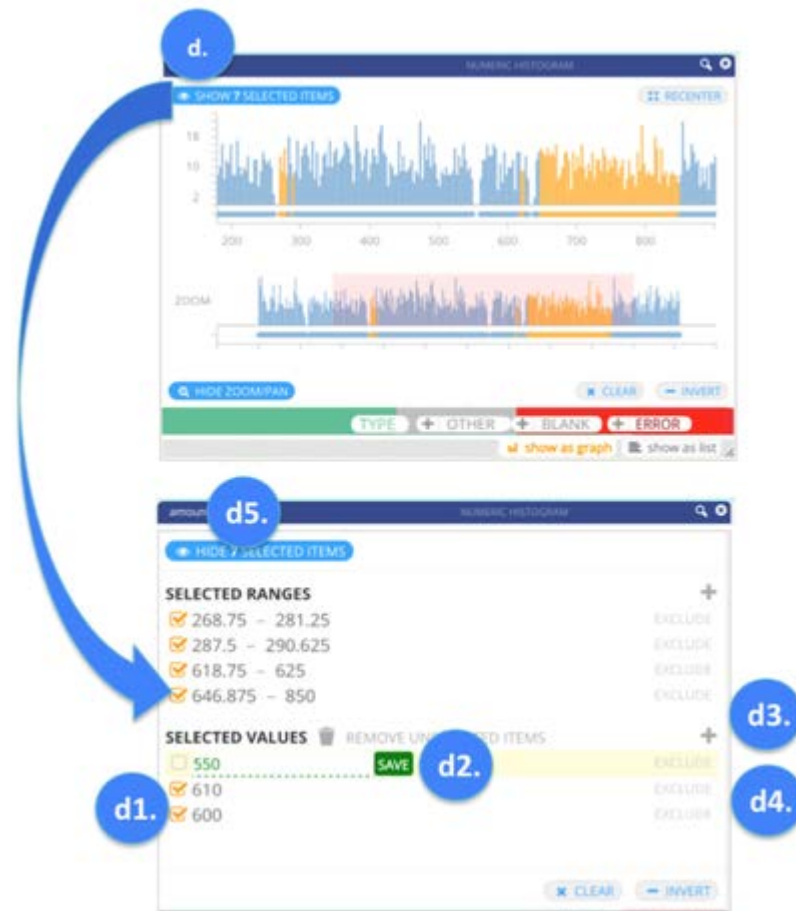
b. Zoom into a selected range: Move your mouse over the values on the x-axis and use the mouse wheel to zoom into the selection. For a Mac, drag two fingers down to zoom into your selection; drag two fingers up to zoom out. You can continue selecting ranges while zoomed in. To return the histogram to its default view, click **Recenter**. Note that your selections are retained after the histogram is re-centered.



c. View the relative position of a zoomed range: After zooming into a selected range, you can simultaneously view where that range is located relative to the entire histogram. Click **Zoom/Pan** to open a second, overview tool histogram below. The overview tool displays the entire range of data and highlights your zoomed selection with a red box. From the overview tool histogram, you can also:

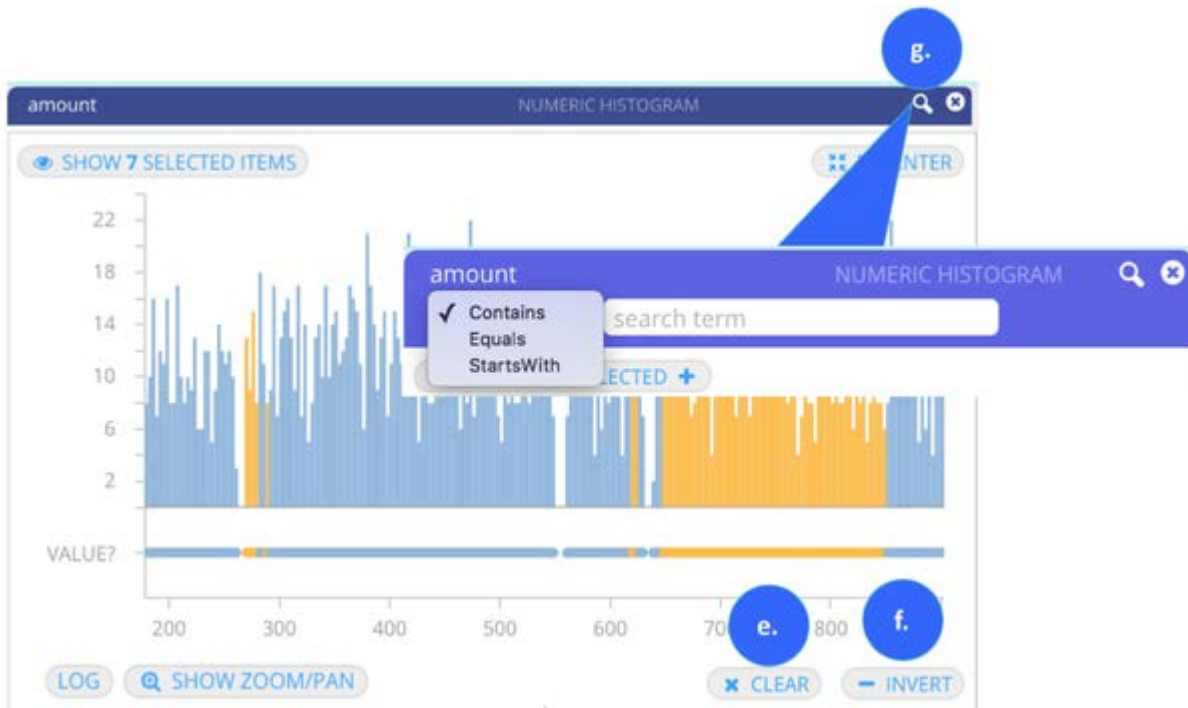
- Drag the red box to view, at the same zoom range, other values that you have *not* selected.

Work interactively with both the overview tool and the top, primary histogram. Remember that actions you take on the overview tool histogram are only reflected on the primary histogram. To dynamically filter your data, you must make your selections on the top, primary histogram.



d. Show Selected Items: After filtering for selected data ranges, click **Show Selected Items**. A new pane opens where you specify the exact ranges and values you want to view in your dataset. From this pane, you can take the following actions:

- **d1.)** Manually remove any of the Ranges or Values in the dataset by clicking to remove the orange check mark adjacent to the Range or Value. Click to return that range or value to your filtered dataset. Note: A trash can icon appears at the top of the pane when you deselect a range or value. This allows you to discard the selection from your filters.
- **d2.)** Manually update any of the Ranges or Values listed here by clicking the value you want to edit. The value then becomes editable. Enter a new value and click **Save**.
- **d3.)** Manually add Ranges or Values for filtering the data by clicking the **+** button. Provide a Min and Max for the Range, or provide a Value, and click **Save**. An entry is created for the Value or Range. Click the check box for an entry and your dataset dynamically updates to reflect your selection(s).
- **d4.)** **EXCLUDE** a Range or Value from the dataset. This is particularly useful if you have already selected a range. From the range, you can **EXCLUDE** specific values (from within that range) to hide from the current dataset. For example, select range 1-2000 to view in your dataset. Then **EXCLUDE** value 195. The dataset displays everything between 1 and 2000, except for value 195. Note that ranges and values marked with **EXCLUDE** are displayed with orange, dotted outlines in the histogram to remind you of their exclusion.
- **d5.)** Click **Hide selections** to toggle off the pane and return to the histogram where your selections are highlighted in the histogram.



e. Invert your selection(s): Displays *all* of your data *except* that which you have selected to filter.

f. Clear: Removes all of your current filters.

g. Search for values: Click the magnifying glass icon in the upper right corner to open a search field. You can search for a value that contains, equals, or starts with any value you specify.

Show numeric columns as lists

The List view allows you to see the count of every numeric value in the column. Viewing your data in this way is particularly useful when you want to quickly select and filter on specific numeric values. Note that selections you make in the List are reflected in the histogram when you click the **Show as Graph** tab. This section describes the actions you can take.



a. Sort list order: By default the list of values displays from lowest to highest count. To reverse the order to display highest to lowest, click the triangle in the upper right corner above the count column. You can also sort the list numerically by clicking the triangle located above the list. The triangle's orange color indicates which sort order (occurrence or numeric) is currently applied to the dataset.

b. Select values to dynamically display in your dataset: Click to select values from the List. To select more than one value: CTRL+click (Windows) or Command+click (Mac). To select a contiguous, multi-row range:Shift+click. To remove any selection, use the following key command: Alt+click. After making selections, click **Currently Selected** in the top left corner to open a new pane where you can continue refining your filtering operations. If you already know the exact ranges and values you want to use for filtering your dataset, you can skip making selections from the List. Instead, click **No Items Selected** to open a new pane where you enter exact values and ranges.

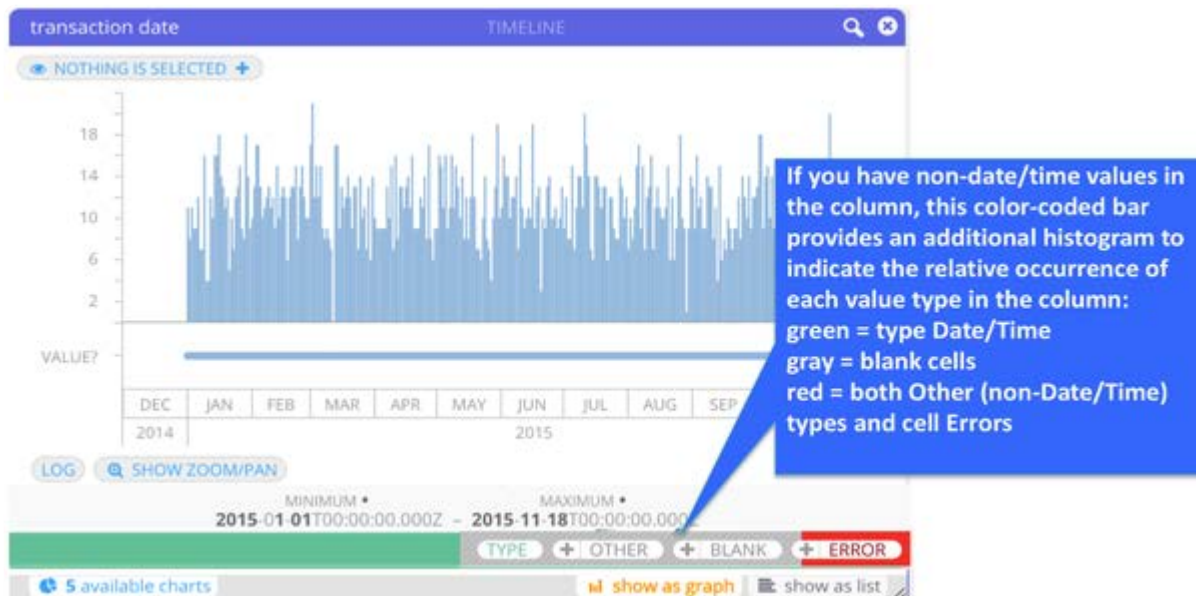
c. Refine your selections: After making selections, click **Show Selected Items** in the top left corner to open a new pane where you can continue refining your filtering operations. The button's label updates to reflect the number of items you have currently selected.

If you already know the exact ranges and values you want to use for filtering your dataset, you can skip making selections from the List. In this case the button's label is **No Items Selected**. Click the button to open a new pane where you enter exact values and ranges. The filtering operations you can perform from the pane are described in steps d1-d5 in [Show numeric columns as graphs](#). You can also **Invert** and **Clear** your selections, and perform a search for specific values from the List. See steps e-g of the same section for details.

Date and Time Filtergrams

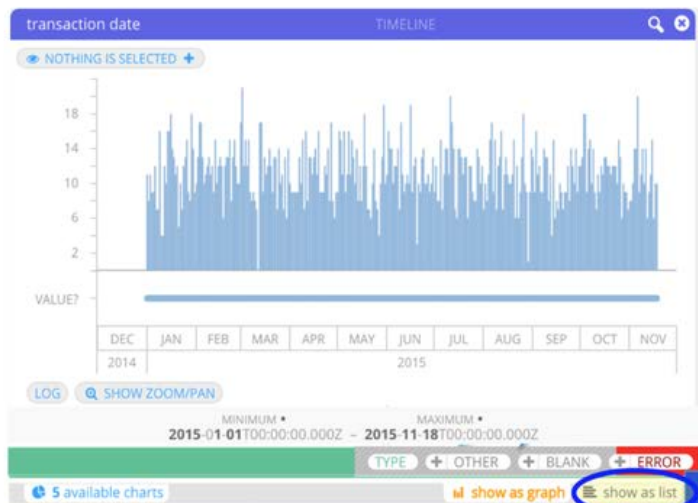
When you open the Date/Time Filtergram pane, there are two views that you can use for performing filtering operations on your data:

- **Show as Graph** (the default view)



This histogram displays the distribution of date/time values in your dataset. The horizontal line, or x-axis, represents the range of date values that occur in your dataset within this column. The height of each bar represents the count of each distinct date value in the column. The scale along the left side, or y-axis, provides a measure for the count. The filtering actions you can take with the histogram are described below in the section [Show date and time columns as graphs](#).

- **Show as List** Click this tab to hide the date/time histogram and display the unique occurrence of every date/time value in the column. From the List, you select values to dynamically display in your dataset. The filtering actions you can take with the List are described below in the section [Show date and time columns as lists](#).



transaction date TIMELINE

NOTHING IS SELECTED +

2015-08-29T00:00:00.000Z	1
2015-06-20T00:00:00.000Z	3
2015-01-10T00:00:00.000Z	4
2015-01-11T00:00:00.000Z	4
2015-05-27T00:00:00.000Z	4
2015-09-13T00:00:00.000Z	4
2015-01-21T00:00:00.000Z	5
2015-02-18T00:00:00.000Z	6
2015-03-30T00:00:00.000Z	6
2015-04-29T00:00:00.000Z	6
2015-05-22T00:00:00.000Z	6
2015-07-15T00:00:00.000Z	6
2015-07-22T00:00:00.000Z	6

321 UNIQUE VALUES IN LIST

TYPE + OTHER + BLANK + ERROR

5 available charts show as graph show as list

Work with Date and Time Filtergrams

When moving your mouse anywhere over the Filtergram pane, the following buttons display:



a. Currently Selected (top left corner): When you make selections on the histogram, the button's label changes to indicate your number of selections. Click the button to open a new pane that lists every selected Range and Value currently selected in the histogram. From this pane you can continue to refine the date/time values you want to filter on in your dataset. If you already know the ranges and values that you want to filter on in your dataset, you can click this button instead of using the histogram. From the new pane you can enter dates and ranges of dates to begin your filtering operations. The actions you can take from this pane are described below in the sections [Show date and time columns as graphs](#) and [Show date and time columns as lists](#).

b. Log (bottom left corner): Toggles on a logarithmic scale (log) view of your data. By default, a linear view of the data is displayed in the Filtergram. However, if you have a large date range in which one or a few points in the data are much larger than the bulk of the data, the log view adjusts for skewness in your data.

c. Show Zoom/Pan (bottom left corner): Toggles on and off an overview tool used for viewing the relative position of values and ranges that you have zoomed into on the histogram. Zooming and panning operations are described in the section [Show date and time columns as graphs](#).

d. When moving your mouse anywhere over the Filtergram pane, the following buttons display:

- **Type:** When you mouse over this button, a ratio displays to express the count of currently selected date/time type rows relative to all date/time type rows in the dataset. If nothing is selected on the histogram, the ratio reflects the total count for date/time type rows relative to all rows in the dataset. When you click this button, you dynamically hide the date/time values in your dataset. This is useful if you have blanks, errors or other non-date/time values in this column and want to view only those data types.

- **Other:** If there are non-date/time values in the column, for example numbers, this button displays. When you mouse over this button, a ratio displays to express the count of currently selected rows with non-date/time values relative to all rows with non-date/time values in the dataset. If nothing is selected on the histogram, the ratio reflects the total count of rows with non-date/time values relative to all rows in the dataset. Click **Other** to toggle off and hide the other values from the current dataset view.
- **Blank:** If there are blanks in the column, this button displays. When you mouse over this button, a ratio displays to express the count of currently selected blank rows relative to all blank rows in the dataset. If nothing is selected on the histogram, the ratio reflects the total count of blank rows relative to all rows in the dataset. Click **Blank** to toggle off and hide the blank cells from the current dataset view.
- **Error:** If there are errors in the column, this button displays. When you mouse over this button, a ratio displays to express the count of currently selected rows with cell errors relative to all rows with errors in the dataset. If nothing is selected on the histogram, the ratio reflects the total count of rows with cell errors relative to all rows in the dataset. Click **Error** to toggle off and hide the cell errors from the current dataset view.



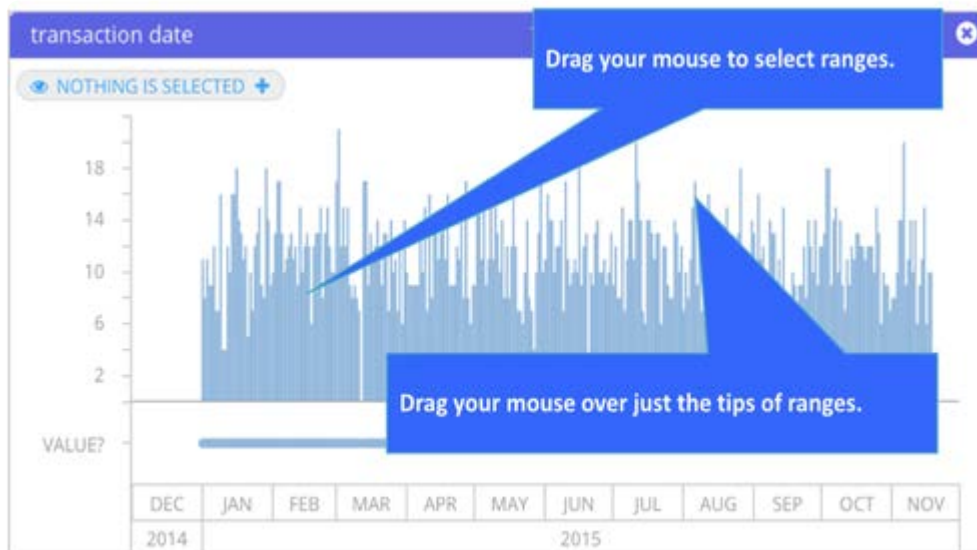
If you have "other" values, blanks or errors in the column, notice there is also a + button adjacent to the button for that value. When you click the + button you add all values of that type, for example all "other" values, to the "Selected Items" list where you can continue working with them. See below for the filtering actions you can take from **Selected Items**.

The Date/Time histogram provides a number of powerful tools that enable you to dynamically filter your data with great precision. This section describes the actions you can take.

Show date and time columns as graphs

- Select ranges to view on the histogram: Click and drag your mouse across a range of values, and notice that your dataset updates to reflect your selection. To select additional, non-contiguous ranges in the histogram, use the following key commands while dragging your mouse: CTRL+click (Windows) or Command+click (Mac). To remove any selection or portion of a selection, use the following key command while dragging your mouse: Alt+click.

You can also drag your mouse over just the tips of ranges in the histogram to display only those values in your dataset. The y-axis is helpful in determining the relative values of peaks in your data:



- Begin exploring and transforming your data:



a. Pan across the histogram: Move your mouse over the values on the x-axis. Notice your cursor changes from the pointer to the click and drag symbol. Click and drag across the values on the x-axis to pan across the histogram and your selections. Click **Recenter button** to return the histogram to its default view. Note that your selections are retained.

b. Zoom into a selected range: Move your mouse over the values on the x-axis and use the mouse wheel to zoom into the selection. For a Mac, drag two fingers down to zoom into your selection; drag two fingers up to zoom out. You can continue selecting ranges while zoomed in. To return the histogram to its default view, click **Recenter**. Note that your selections are retained after the histogram is re-centered.



c. View the relative position of a zoomed range: After zooming into a selected range, you can simultaneously view where that range is located relative to the entire histogram. Click the Zoom/Pan button to open a second, overview tool histogram below. The overview tool displays the entire range of data and highlights your zoomed selection with a red box. From the overview tool histogram you can also:

- Drag the red box to view, at the same zoom range, other values that you have *not* selected.
- Work interactively with both the overview tool and the top, primary histogram. Remember that actions you take on the overview tool histogram are only reflected on the primary histogram. To dynamically filter your data, you must make your selections on the top, primary histogram.



d. Show Selected Items: After filtering for selected ranges, **Show Selected Items**. A new pane opens and displays every selected Range and Value currently displayed in your dataset. From this pane, you can EXCLUDE a Range or Value from the dataset. This is particularly useful if you have already selected a range. From the range, you can EXCLUDE specific values (from within that range) to hide from the current dataset. For example, you have selected the following range of dates to display in your dataset: 03/01/2015 - 03/15/2015. You then EXCLUDE the date 03/10/2015. The dataset displays everything in the range except for 03/10/2015.

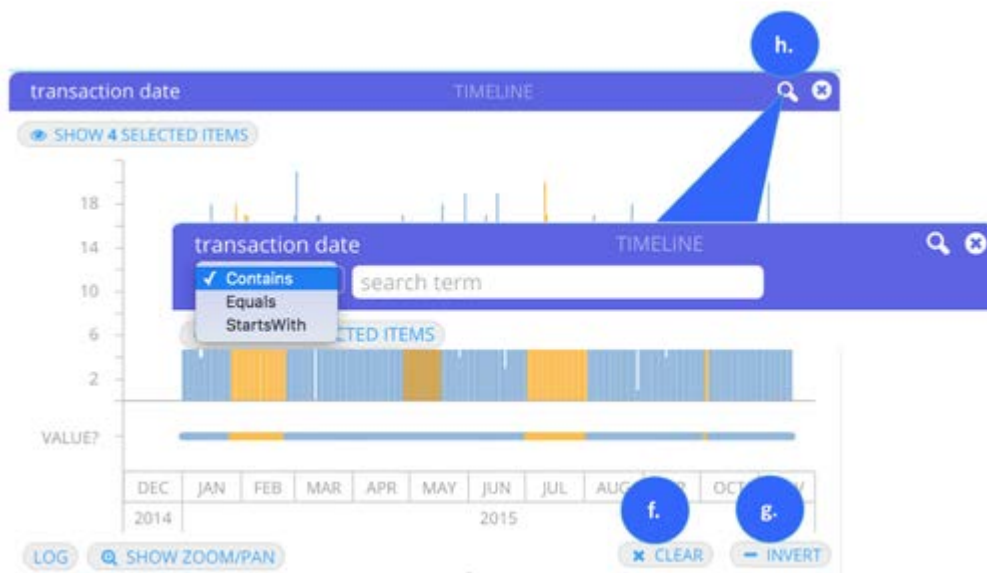
Note that ranges and values marked with EXCLUDE are displayed with orange, dotted outlines in the histogram to remind you of their exclusion.

A trash can icon appears at the top of the pane when you deselect a range or value. This allows you to discard the selection from your filters.

Click **Hide selections** to toggle off the pane and return to the histogram where your selections and exclusions are highlighted in the histogram.



e. Available charts: Click this tab to choose from four additional filters—Month of year, Day of month, Day of week, Hour of day—that you can use to filter your date/time data with more precision. After selecting a filter, click the "5 available charts" tab again to pin that filter to your view. While working with any of these filters, notice that your dataset is dynamically updated to reflect your selections.



f. Clear: Removes all of your current filters.

g. Invert your selection(s): Displays *all* of your data *except* that which you have selected to filter.

h. Search for values: Click the magnifying glass icon in the upper right corner to open a search field. You can search for a value that contains, equals, or starts with any value you specify.

Show date and time columns as lists

The List filter view allows you to see the count of every date/time value in the column. Viewing your data in this way is particularly useful when you want to quickly select and filter on specific dates. Note that selections you make in the List are reflected in the histogram when you click the **Show as Graph** tab. This section describes the actions you can take.



a. Sort list order: By default, the list of dates displays from lowest to highest occurrence. To reverse the order to display highest to lowest, click the triangle in the upper right corner above the count column. You can also sort the list chronologically by clicking the triangle located above the date/time values. The triangle's orange color indicates which sort order (occurrence or chronological) is currently applied to the dataset.

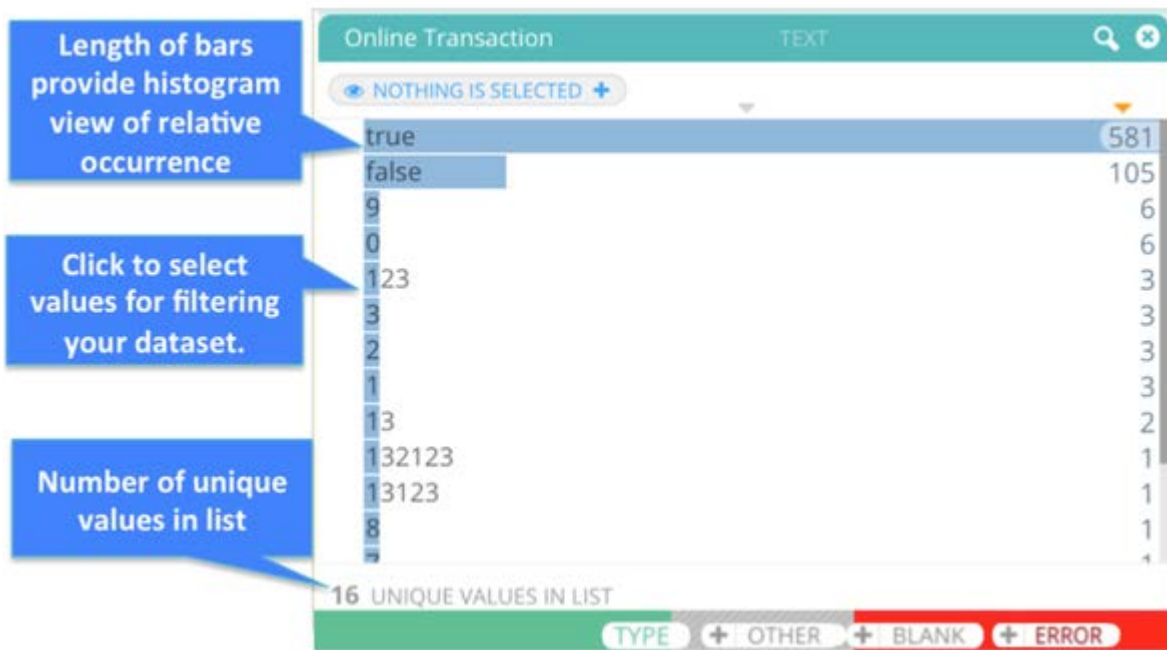
b Select values to dynamically display in your dataset: Click to select values from the List. To select more than one value: CTRL+click (Windows) or Command+click (Mac). To select a contiguous, multi-row range: Shift+click. To remove any selection, use the following key command: Alt+click. After making selections, click **Currently Selected** in the top left corner to open a new pane where you can continue refining your filtering operations. If you already know the exact ranges and dates you want to use for filtering your dataset, you can skip making selections from the List. Instead, click **No Items Selected** to open a new pane where you enter exact values and ranges.

c. Refine your selections: After making selections, click **Show Selected Items** in the top left corner to open a new pane where you can continue refining your filtering operations. The button's label updates to reflect the number of items you have currently selected.

The filtering operations you can perform from the pane are described in step d of the [Show date and time columns as graphs](#) section. You can also **Invert** and **Clear** your selections, and perform a search for specific values from the List. See steps f-g in the same section.

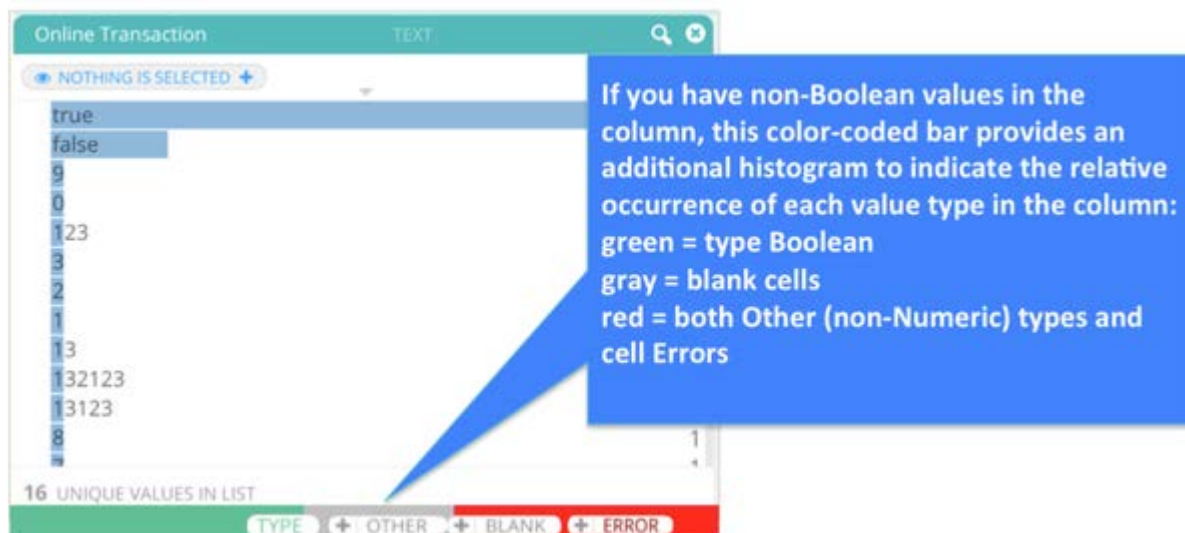
Boolean Filtergrams

The Boolean Filtergram shows the count of Boolean values in your dataset and filters out any other values from it.



The bars running left-to-right provide a histogram view of the relative number of times each value occurs. The total number of unique values is listed in the bottom left side of the pane. From the list, you can select values to dynamically display in your dataset.

Following is an overview of the Boolean Filtergram display:



Work with Boolean Filtergrams

When moving your mouse anywhere over the Filtergram pane, the following buttons display:

- **Type:** When you mouse over this button, the count of Boolean values in this column is listed. When you click this button, you dynamically hide the Boolean values in the column.
- **Other:** If there are non-Boolean values in the column, this button displays. When you mouse over this button, the count of non-Boolean values in this column is listed. When you click this button, you dynamically hide the non-Boolean values in the column. Alternatively, click the **+** button to add all "other" values to the Filtergram list. You can then filter out and hide specific "other" values from the current dataset view by using the following key commands while clicking the "other" value you want to hide: Alt+CRTL (Windows) or Alt+Command (Mac).
- **Blank:** If there are blanks in the column, this button displays. When you mouse over this button, the count of blank values in this column is listed. When you click this button, you dynamically hide the blank values in the column.
- **Error:** If there are errors in the column, this button displays. When you mouse over this button, the count of cell errors in this column is listed. When you click this button, you dynamically hide the cell errors in the column.

Use **Clear** and **Invert** to manage your selections in the Filtergram list. Note that **Invert** displays *all* of the values in the current dataset view *except* the values you have selected in the list.

Source Filtergrams

The Sources Filtergram allows shows how the rows from lookup and appended datasets are participating with the base dataset in your Project.

Source Filtergram example

You have an outer join in which there are rows from a lookup dataset that do not match the base dataset. The Sources Filtergram allows you to see how many rows from the base and lookup datasets are participating in the join. Additionally, a count for the unmatched rows from each dataset source is displayed.

Following is an overview of the Source Filtergram display:

The screenshot shows a 'Sources' list with the following data:

Source	Count
[Dispaly Technology Catalogue_2015.xlsx v.1] + [04.xlsx v.1] + [02.xlsx v.1] + [03.xlsx v.1]	659
[Dispaly Technology Catalogue_2015.xlsx v.1] + [04.xlsx v.1] + [03.xlsx v.1]	140
[Academic_DT_Inv_2015 (melanie).xlsx v.1] + [Dispaly Technology Catalogue_2015.xlsx v.1]	140
[Dispaly Technology Catalogue_2015.xlsx v.1] + [03.xlsx v.1]	73
[Academic_DT_Inv_2015 (melanie).xlsx v.1] + [Dispaly Technology Catalogue_2015.xlsx v.1]	16
[Academic_DT_Inv_2015 (melanie).xlsx v.1] + [Dispaly Technology Catalogue_2015.xlsx v.1]	8
[Dispaly Technology Catalogue_2015.xlsx v.1]	1

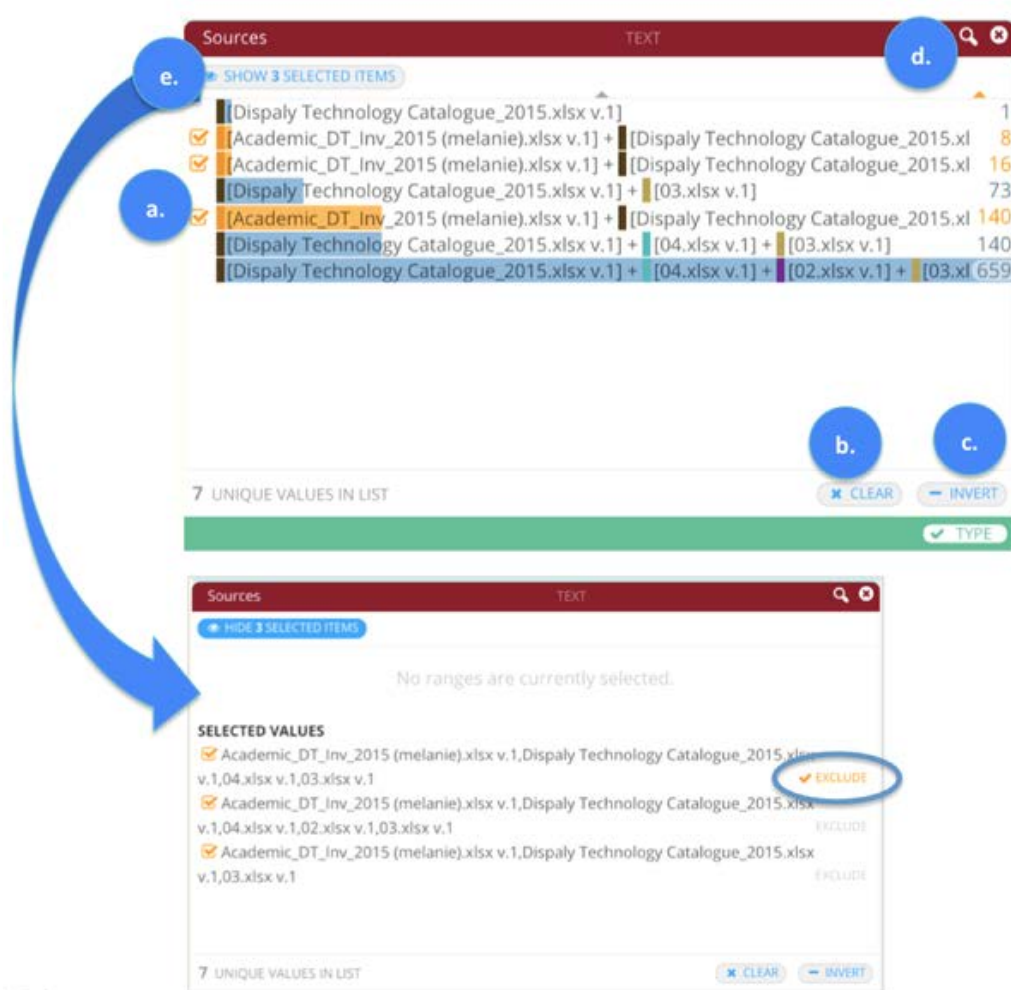
Callouts in the image:

- Click to sort alphabetically**: Points to the triangle above the list.
- Click to sort by count**: Points to the triangle above the count column.
- Datasets are color-coded to match their corresponding column header colors**: Points to the color-coded source names.

Additional UI elements: 'Sources', 'TEXT', 'NOTHING IS SELECTED +', '7 UNIQUE VALUES IN LIST', and a 'TYPE' button.

By default, the list of Sources displays from highest to lowest count. To reverse the order to display least to most, click the orange triangle in the upper right corner above the count column. You can also sort the list alphabetically by clicking the triangle located above the list. The triangle's orange color indicates which sort order (numeric or alpha) is currently applied to the list.

Work with Source Filtergrams



a. Select Sources to display in your dataset: Click any Source to dynamically filter your dataset to display only that Source. To select more than one Source: CTRL+click (Windows) or Command+click (Mac). To select a contiguous, multi-row range: Shift+click.

b. Clear: Remove all Source filter selections.

c. Invert your selections: Displays all Sources except those you have selected.

d. Search for Source files: Click the magnifying glass icon in upper right corner to open a search field. You can search for a Source file that contains, equals, or starts with any text value you specify.

e. Refine your filter selections: Click **Show Selected Items**. A new pane opens and displays the currently selected Sources. From this pane, you can take the following action:

- **EXCLUDE** a Source. This is a toggle that hides the associated Source from your dataset until you click **EXCLUDE** again. Note that Sources marked with **EXCLUDE** are displayed with an orange dash in the initial Filtergram pane to remind you of their exclusion.
- A trash can icon appears at the top of the pane when you deselect a Source. This allows you to discard the Source from your filters.
- When you are finished working in this pane, click **Hide Selected Items** to return to the initial Filtergram view.

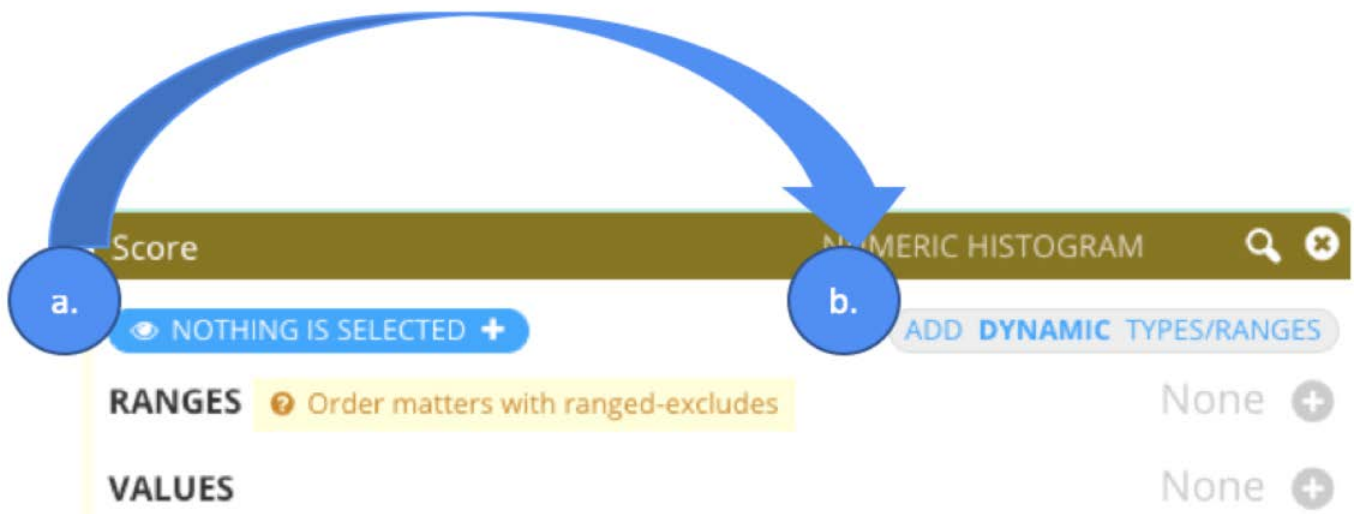
Dynamic Ranges

The Dynamic Percentiles feature in Filtergrams provides you with the powerful option to specify percentiles of your selected values. For example, if you have an inventory dataset that has a column for products sold per week per region, you can use Dynamic Percentiles to filter select the top 5% highest selling products for each region. Your percentile selections are also dynamically applied to newer versions of datasets that are automatically updated in the Library through [Automatic Project Flows \(APF\)](#). For example, you can use the APF feature with Dynamic Percentiles to automatically produce an AnswerSet each week to identify the top 5% highest selling products for each region per week.

The Dynamic Filtering options are applicable for Date/Time, String, and Numeric type columns. To open the Dynamic filtering options pane:

(a.) Click **Currently Selected** in the top left corner of your Filtergram.

(b.) Click **Add Dynamic Types/Ranges**.



The Dynamic Ranges pane opens:

The screenshot shows a 'Score' column in a 'NUMERIC HISTOGRAM' view. The interface is divided into 'DYNAMIC SELECTIONS' and 'STATIC SELECTIONS'. The 'DYNAMIC SELECTIONS' section includes a 'VALUES THAT ARE' filter with options for 'VALID', 'INVALID', 'BLANK', and 'ERROR'. Below this is a 'DYNAMIC RANGES' filter with a 'REMOVE UNSELECTED ITEMS' button. A number line shows percentile ranges from 0% to 100%. Two dynamic ranges are shown: one from 0% to 10% and another from 75% to 100%. The '0% - 10%' range is selected, and its 'ORDERED BY' is set to 'VALUE'. The '75% - 100%' range is also selected, and its 'ORDERED BY' is set to 'COUNT'. The 'STATIC SELECTIONS' section has 'RANGES' and 'VALUES' filters, both currently set to 'None'. At the bottom, there are 'MINIMUM' and 'MAXIMUM' values of 67 and 67.0625, and buttons for 'CLEAR' and 'INVERT'. At the very bottom, there are 'show as graph' and 'show as list' options.

1 First, select and deselect which types of values from your dataset to include in the percentiles.

2 Then, click the plus sign to add more Dynamic Ranges for the selected values.

3 Drag both sides of the orange band to select your desired percentile range. The arrow keys on your keyboard can be used to get exact increments on the number line.

4 Select one or more of the Dynamic Ranges that you have created. By selecting more than one, they will be simultaneously displayed.

5 Toggle selection by:

- Value = percentile based on actual values in column.
- Count = percentile based on frequency of occurrence per value.

6 Click to toggle between setting your Dynamic Ranges, and viewing them on a graph or list.

Work with Dynamic Selections

1. Select or deselect the types of values you want to include in your percentile ranges: Valid, Invalid, Blank, Error:

- Valid: Values that are of the same type as the column type—for example numbers in a numeric type column.
- Invalid: Values that are not of the same type as the column type—for example alpha characters in a numeric type column.
- Blank: Use if there are blanks in the column.
- Error: Use if there are errors in the column.

2. Click **plus** to add as many different dynamic ranges as desired.

Note

Each of these ranges will be applied to only the types of values that you select in step 1.

3. Set your Dynamic Range: Drag the limits on the number line to set your desired value.

Tip

Use your left and right keyboard arrow keys to adjust for exact values on the number line.

4. (Optional) Add more range selections (following steps 1-3 above). When you create multiple range selections, each additional selection is treated as an AND operation. For example, if you have created two range selections—one set to select the top 10th percentile, and the other to select the bottom 10th percentile—then the values that correspond to these percentiles are highlighted in the filtergram and displayed accordingly on the data grid.

Note

You can always set your dynamic filtering options to work in conjunction with the other filtering operations described in this document for each column type.

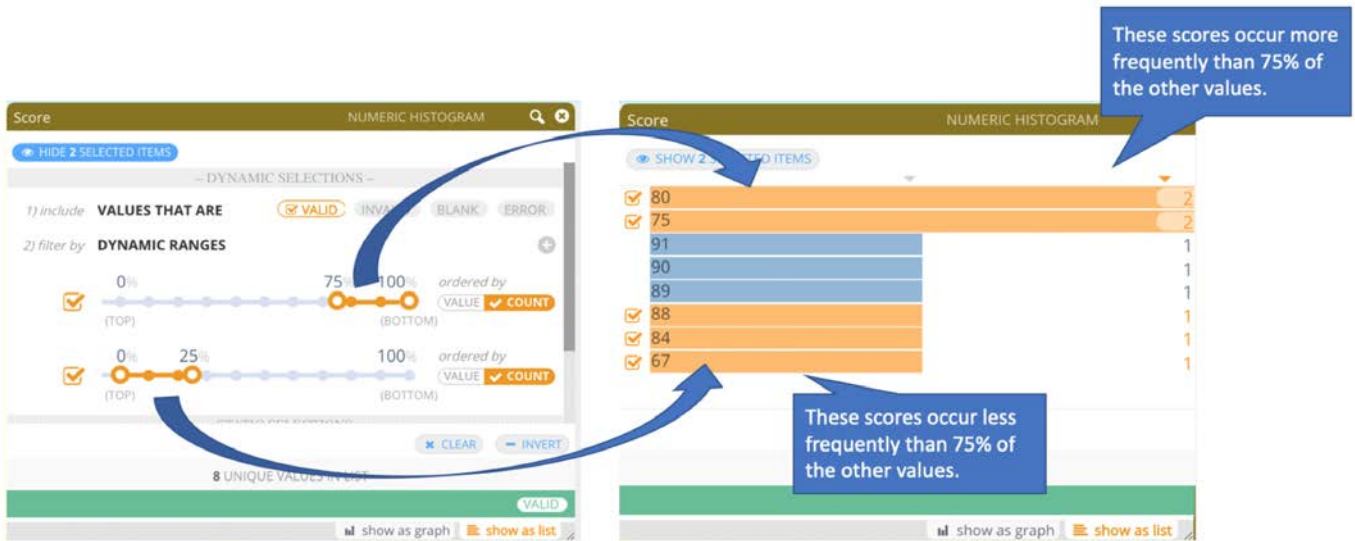
5. Order by **Value** or by **Count**: when ordered by **Value**, the percentile is based on the *actual values* in the column. When you order your data by **Count**, the percentile is based on the *frequency of occurrences* per value. For example, you have the following dataset with participant ages and scores on an exam.

	Sources	A-Z Name	123 Age	123 Score
1		Dennis	32	91
2		Sara	30	90
3		Jess	35	89
4		Mark	25	88
5		Rick	36	84
6		Grace	31	80
7		Jack	35	80
8		Yana	34	75
9		Erin	26	75
10		Max	36	67

If you want to determine how the participants scored in relation to one other, order by **Value**:



If you want to determine the validity and usefulness of the exam itself or show patterns in the exam results, order by **Count**:



6. Click **Currently Selected** to view your highlighted percentiles in the graph or list view.

Behavior of multiple Filtergrams across columns

Filtergrams for multiple columns may be open and dynamically filtered at the same time. Note that the resolution for multiple filtergrams is from left to right—whichever dynamic ranges are set on the left filtergram affect the resulting values that can be selected by the adjacent filtergram on the right:

Filtering the 75th-100th percentile in the Age column produces four values.

Creating a dynamic range for the Score column will filter from the remaining values.

The result shows those who are in the top Age percentile AND the bottom percentile of Scores within that range.

	Name	Age	Score
1	Jess	35	89
2	Rick	36	84
3	Jack	35	80
4	Max	36	67

Name	Age	Score
Max	36	67

If you exit out of a Filtergram on the left, then the dynamic ranges being applied for that column are removed.

Note


- If the [Interactive Mode](#) feature is enabled, then your selections are applied to the entire dataset.
- Dynamic Ranges is a feature that must be enabled. If you do not see this button in your Project, contact your Data Prep System Administrator.

Detect and transform date formats

Data Prep detects date formats in a column and lets you transform the cells within the column to ISO date format. Date formatting issues are common problems you'll find with the raw data that you need to prep for your training datasets. Often the values are of type string or numeric and you need to cast them as a date type. For example, you might have a format like 20210101 that you want to transform into a format like 01/01/2021.

Format dates

To format date in a column of your dataset:

1. Locate the column where you want to format dates.
2. Hover over the column menu icon  and click **Change into > date**.



3. Click **Detect Date Format**.

Change start date

into date using

Value is required

Detect Date Format

empl	A-Z team matrix	A-Z location	123 years education	A-Z major	123 age	A-Z start date	123 an
	D52 - eng	New York	12	Computer Science	32	3/27/21	
	D17 - sales	CA	16	Business	28	5/2/20	
	D24 - sales	California	18	Another Engineering	46	2/25/1999	
	D22-sales	NY	16	Social Science	35	7/09/2011	
	D55 - exec	California	16	Business	48	11/2/20	
	D17 - sales	Calif.	18	Humanities	29	4/6/2008	
	D12 - Ops	New York	16	Computer Science	38	3/5/05	
	D51 - eng	CA	16	Another Engineering	34	2/12/2019	
	D52-eng	Amsterdam	18	Computer Science	42	4/16/2005	

The most common format appears as the custom format.

Change start date

into date using Custom M/d/yyyy

Detect Date Format Validate First 1000 Rows

First 1000 rows will be used to detect date format. If more than one format is fo

location	123 years education	A-Z major	123 age	A-Z start date	start date
	12	Computer Science	32	3/27/21	0021-03-27T00:00:00.000Z
	16	Business	28	5/2/20	0020-05-02T00:00:00.000Z
	18	Another Engineering	46	2/25/1999	1999-02-25T00:00:00.000Z
	16	Social Science	35	7/09/2011	2011-07-09T00:00:00.000Z
	16	Business	48	11/2/20	0020-11-02T00:00:00.000Z
	18	Humanities	29	4/6/2008	2008-04-06T00:00:00.000Z
	16	Computer Science	38	3/5/05	0005-03-05T00:00:00.000Z
	16	Another Engineering	34	2/12/2019	2019-02-12T00:00:00.000Z
im	18	Computer Science	42	4/16/2005	2005-04-16T00:00:00.000Z

If a valid format is not detected, an error message displays.

4. Click **Validate First 1000 Rows**.

Data Prep validates the data against the specified format and verifies if there are inconsistencies. If validation fails for any cells in a column, an error message specifies which columns failed.

If there are validation errors, perform the date conversion on each invalid date value. Repeat the operation until all values are converted to the specified date format.

5. Click **Save** on the top right.

After you save the selection, notice the column type is converted to a date type as indicated by the timer symbol.

A-Z location	123 years education	A-Z major	123 age	start date
New York	12	Computer Science	32	0021-03-21
CA	16	Business	28	0020-05-01
California	18	Another Engineering	46	1999-02-20
NY	16	Social Science	35	2011-07-01

You can also detect and validate more than one column.

Change 2 columns by Name Criteria Personal sort not present

All Columns All Types 2 of 3 columns selected

<input type="checkbox"/>	Column Name	Type
<input checked="" type="checkbox"/>	GROUP	A-Z
<input type="checkbox"/>	StartDate	A-Z
<input checked="" type="checkbox"/>	EndDate	A-Z

into date using yyyy-MM-ddT'HH:mm:ssZ

Detect Date Format Validate First 1000 Rows

The following columns failed validation: EndDate, GROUP

Sources	A-Z GROUP	GROUP	A-Z StartDate	A-Z EndDate	EndDate
1	A	A	1900-01-01	1949-12-31	1949-12-31
2	A	A	1949-01-01	1959-12-31	1959-12-31
3	A	A	1960-01-01	1969-12-31	1969-12-31
4	A	A	1970-01-01	1979-12-31	1979-12-31
5	A	A	1980-01-01	1989-12-31	1989-12-31
6	A	A	1990-01-01	1999-12-31	1999-12-31
7	A	A	2000-01-01	2009-12-31	2009-12-31

Split columns

Data Prep provides features to help you split columns during data prep. Splitting a column involves taking the values in a single column and distributing them into one or more new columns in the same row. The split is performed either by selecting a character string on which to separate the column or by giving a character count (number of characters) along which the split should occur.

Match

The screenshot shows the 'Split' configuration interface in Data Prep. The 'Match' tab is selected, with a separator '@' and two columns: 'Email (1)' and 'Email (2)'. Below the configuration is a table with columns: First Name, New Prospect ID, Source, Workflow Status, Email, and Email (1). The table contains 20 rows of data, with the 'Email' column being split into 'Email (1)'.

APPLICANTS (FINAL)	APPLICANTS (FINAL)	APPLICANTS (FINAL)	APPLICANTS (FINAL)	APPLICANTS (FINAL)	APPLICANTS (FINAL)
First Name	New Prospect ID	Source	Workflow Status	Email	Email (1)
ffrey	Q76UGKKJW2WC4ZJ	LinkedIn	Hire	jshah@yahoo.cm	→ jshah
alerie	8QEIOCY1AYWID1LR	Intellisource	New	vwilson@gmail.com	→ vwilson
ay	PASDYM4BH933WYOS	LinkedIn	Hired Elsewhere	rdeaser@gmail.com	→ rdeaser
egan	Z7YYJAXYOFJTREA4	LinkedIn	Hire	myu@gmail.com	→ myu
cott	XMRVAVUHUMUML2C7	Simply Hired	Filed for Later	sthota@gmail.com	→ sthota
llie	X4JNBCMPH4MYDHSV	LinkedIn	Filed for Later	igopina@gmail.com	→ igopina
se	XDYTCILRDNKGGGGT	our website	Filed for Later	janupin@gmail.com	→ janupin
am	FXQ45G6B1LSH6M8V	LinkedIn	New	slin@cal.berkeley.edu	→ slin
finerva	JNVFXBR32EMIGCR	Andiamo	Phone Screen	mvelayu@yahoo.com	→ mvelayu
ob	KSJTOQRDIGWWAQQM	Simply Hired	New	blin@yahoo.com	→ blin
ferrill	PRRQHBRH3LQFQEVV	Jivaro	Schedule On-site Interview	mshah@gmail.com	→ mshah
arveen	VJFB83WPUIHLG0TT	Simply Hired	Phone Screened	pvipa@gmail.com	→ pvippa
heik	RY458MYFNSU041OJ	Intellisource	Filed for Later	skuikar@ccs.neu.edu	→ skuikar
icky	UFWVALNEDY9LRRG2	Simply Hired	New	vsingh@gmail.com	→ vsingh
wipan	GOCMFGEB1OMO62V6	LinkedIn	Not Qualified	asimmon@gmail.com	→ asimmon
ictoria	RRFOL5NG0VANSUTO	Simply Hired	Filed for Later	vsarawg@gmail.com	→ vsarawg
fargaret	QQTkRLBHYKSTBL2U	Andiamo	New	mmehta@gmail.com	→ mmehta
oug	LTSKSPUE5GTPDI0L	Simply Hired	Phone Screen	dli@gmail.com	→ dli
mv	ABGZHILOOSEHNV00	Elevate	Schedule On-site Interview	aaopach@email.com	→ aaopach

In the **Separator** field, enter one or more characters on which the split should occur. This will serve as the pattern against which the match will occur. The characters specified will not become part of any of the resulting column(s)—they will be excluded from the distribution to the new column(s).

Note

The text separator value is case-sensitive.

Based on the separator you have chosen for the operation and the number of times that separator exists across the cells of this column, the **Min** and **Max** links allow you to quickly select the number of new columns to create after the split. In addition, the **Custom** link allows you to provide the exact number of columns you want to generate. Adjacent to each field, you can also click the plus and minus buttons to manually add and remove new column fields.

You also have the option to select **Right to left** for the split. By default, the split function parses left to right using the specified separator. The **Right to left** option allows you to parse the text in the column starting on the right. This is particularly useful in cases where you want to separate a string of text, for example a file name, from its directory path. In this case, specify the slash "/" as your separator, and then select the **Right to left** option.

Length

APPLICANTS (FINAL)	APPLICANTS (FINAL)	APPLICANTS (FINAL)	APPLICANTS (FINAL)	APPLICANTS (FINAL)	APPLICANTS (FINAL)
Email	Email (1)	Email (2)	Phone	Address	City
1	→ jshah	@yahoo.	4088382341	201 Trigg St.	Abingdon
om	→ vwils	on@gmai	870 489 3657	4 Fuller St	Alexandria Bay
:om	→ rdeas	er@gmai	650.207-0151	1000 Alderman Dr	Alpharetta
1	→ myu@g	mail.co	3125047493	120 E Oak St	Anderson
im	→ sthot	a@gmail	9139542404	14817 Oak Ln Ste Ph	Hialeah
om	→ jgopi	na@gmai	4087616535	1071 Pemberton Hill Rd Ste 202	Apex
:om	→ janup	in@gmai		825 E Wisconsin Ave	Appleton
y.edu	→ slin@	cal.ber	6505339136	2002 Summit Blvd 6th Floor	Atlantaa
.com	→ mvela	yu@yaho	469-835-4432	1310 Seaboard Industrial Blvd NW	Atlanta
1	→ blin@	yahoo.c	4257484459	1170 Peachtree St NE Ste 2400	Atlanta
om	→ mshah	@gmail.	650-933-6613	12770 Gateway Dr S	Seattle
om	→ pvipp	a@gmail	770-595-2574	2002 Summit Blvd 6th Floor	Ahtlanta
i.edu	→ skulk	ar@ccs.	1-617-820-4663	2455 Paces Ferry Rd SE	Atlanta
im	→ vsing	h@ymail	6044176561	303 Peachtree St	Atlanta
l.com	→ asirm	on@gmai	(408) 561-0608	250 Williams St NW Ste m100	Atlanta
com	→ vsara	wg@gmai	9197459142	1 Coca Cola Plz NW	Atlinta
com	→ mmeht	a@gmail	206-661-1881	303 Peachtree St	Atlanta
	→ dli@g	mail.co	3522266444	2455 Paces Ferry Rd SE	Atlahnta
.com	→ aadda	ch@gmai	925-658-0664	1 Kellogg Sa	Battle Creek

The **Lengths** field requires one or more numbers, separated by commas. The **Lengths** field is used to determine the number of characters in each of the new columns. Thus, the values "2,3,2" placed in the field will generate one column with the first two characters, create a second column with the subsequent three characters, and create a third column with the next two characters within the original column.

Unlike **Split by Separator**, this type of split does not drop any characters from the middle of the column. However, instead of any unaccounted-for (leftover) characters being placed together in the last column, they are excluded from the new columns completely. To avoid inadvertently truncating a column, it is recommended that the last number specified be large enough to account for the remaining column length.

Note that doing this will not generate a "whitespace" buffer to account for this value being larger than the length of the available character values. This will only ensure that all characters in the column are accounted for in the last column. Extra spaces specified by the **Lengths** parameter that go beyond the number of characters in the column are simply ignored.

Regular Expression (Regex)

If you are familiar with [regular expressions \(Regex\)](#), you can use this option to define a search pattern for locating strings where you want to enforce the split. In the following example, Regex is used to split on alpha characters in the string. The split then results in two new numeric columns:

Split Part Numbers by Match **Regex** Length

REGEX: `[a-z]` COLUMNS: SPLIT TO 0 (MIN), 2 (MAX) OR CUSTOM

OPTIONS: Ignore case, Capture mode

Part Prefix
Part Suffix Numeric
Part Suffix Alpha

Sources	Part Numbers	Part Prefix	Part Suffix Numeric	Part Suffix Alpha
1	123a456a	→ 123	456	
2	123a457b	→ 123	457	
3	123a458c	→ 123	458	
4	123a459d	→ 123	459	
5	123b235e	→ 123	235	
6	123b236f	→ 123	236	
7	123b236g	→ 123	236	
8	123b238h	→ 123	238	

The **Capture Mode** option allows you to extract strings that match your Regex pattern.

In the following example, Regex Capture mode is used to extract only the second set of numbers in the string. The split then results in a new numeric column:

Split Part Numbers by Match **Regex** Length

REGEX: `[0-9]{3}[a-z]([0-9]{3})[a-z]` COLUMNS: SPLIT TO 0 (MIN), 2 (MAX) OR CUSTOM

OPTIONS: Ignore case, Capture mode

Part Suffix

Sources	Part Numbers	Part Suffix
1	123a456a	→ 456
2	123a457b	→ 457
3	123a458c	→ 458
4	123a459d	→ 459
5	123b235e	→ 235
6	123b236f	→ 236
7	123b236g	→ 236
8	123b238h	→ 238

Fill columns

You use the Data Prep Fill operation to populate blank cells within a column based on the known values directly preceding or following the blank(s). To use the Fill function for a column, hover over the column operations menu, scroll to fill in the drop-down, and then select the action of choice. The Group and Sort by options are available for all Fill operations. Click **Save** to apply the filled in values to your Project.

- **Fill Up:** When a data point is missing, it will be filled in using the next available non-blank value.
- **Fill Down:** When a data point is missing, it will be filled in using the most recently seen non-blank value.
- **Fill Average:** When a data point is missing, it will be filled in with the average of the nearest previous and following non-blank values within a partition. Consecutive blanks will be filled with the same value because their input values are the same.
- **Fill Linear Fit:** When a series of data points are missing, they will be filled in with values fitting a line between the surrounding available values. This is different from fill average, as fill average assigns the same average value to all missing values. Fill linear fit calculates a linear average and adjusts the values based on the number of missing values.

Fill Up

123 Value	123 Value
5	5
	5
9	9
	9
	9
15	15
	15

Fill Average

123 Value	123 Value
5	5
	5
9	7
	9
	12
	12
15	15

The top row becomes 5 because the average is calculated from only one value.

All rows between the last known value (9) and next known value (15) are assigned the same Average value.

Fill Down

123 Value	123 Value
5	5
	5
9	9
	9
	9
15	15

Fill Linear Fit

123 Value	123 Value
	5
5	5
	7
9	9
	11
	13
15	15

Paxata assigns equal increments along the line between the last known value (9) and the next known value (15).

Note

Fill average and linear fit can only be applied to numeric column types.

Normalize using clustering

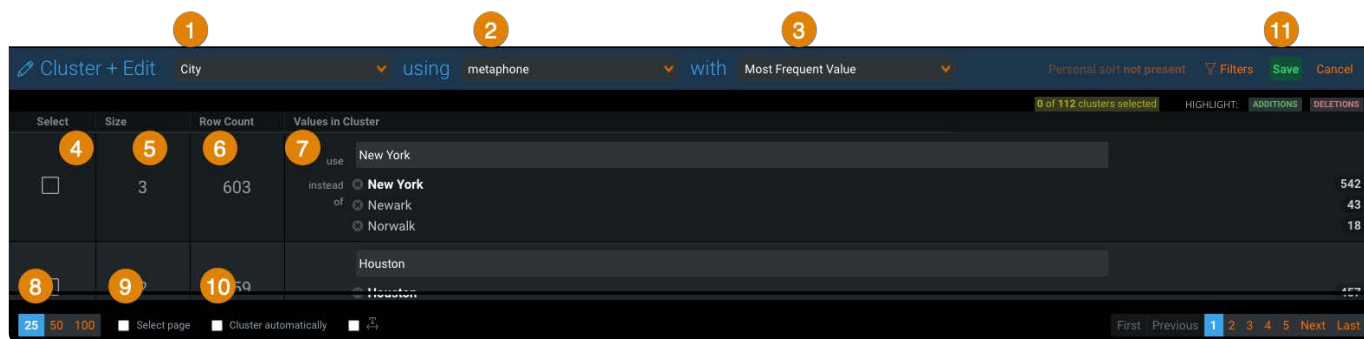
You can use the Data Prep **Cluster + Edit** operation to quickly normalize column data. It is especially useful in spotting inconsistencies and errors in a column.

When you execute **Cluster + Edit** on a column in your dataset:

- Data Prep searches all column values and groups closely matching values together in a cluster.
- Each cluster is listed in the **Cluster + Edit** pane, along with its cluster size (the number of unique values in the cluster) and its row count (the number of times each unique value occurs in the column.)

Based on the clustered data, Data Prep suggests a single replacement value to normalize all of the values in the cluster. You can accept the suggestion or specify another value to use for normalizing the cluster.

The following is the **Cluster + Edit** pane with descriptions of its components:



ELEMENT	ACTION
---------	--------

- 1** Column field
The column on which you want to perform the operation.
- 2** **using**
Use the dropdown menu to select an algorithm to use for the clustering operation. See [Clustering algorithms](#) for details.
- 3** **with**
Use the dropdown menu to select one of the available algorithms to use for the output option. The algorithm you choose determines the new value suggestion for the cluster. See [Output algorithms](#) for details.
- 4** **Select**
Click to select a cluster you want to update.
- 5** **Size**
Number of unique values in the cluster.


ELEMENT	ACTION
6 Row Count	Number of rows in the cluster.
7 Values in Cluster	Displays the value Data Prep suggests as a replacement for all values in the cluster. Data Prep's suggestion is based on the clustering algorithm you choose. To override Data Prep's suggestion, enter a different value. Click the 'X' next to a value to prevent it from being updated.
8 25 / 50 / 100	Select the number of clusters per page for bulk edits.
9 Select page	Lets you select clusters for bulk editing. The bulk edit operates on one page of clusters. Use the Page Size field to specify the number of clusters per page.
10 Cluster automatically	Performs the bulk edit on all rows in the dataset after you click Save .
11 Save	Saves changes.

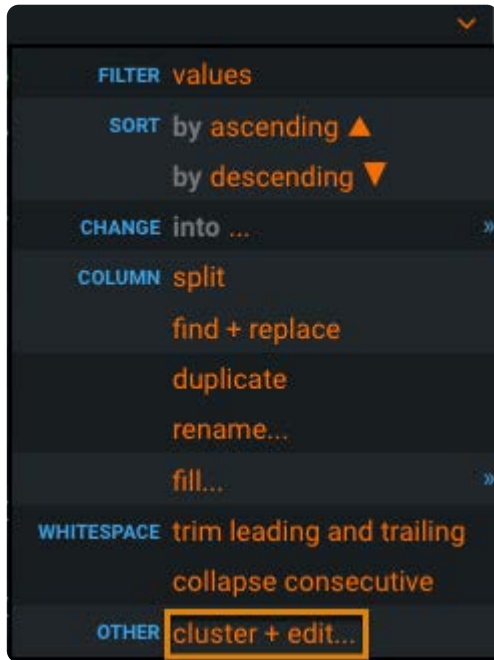
When to use Cluster + Edit

The following table describes common scenarios for using **Cluster + Edit**.

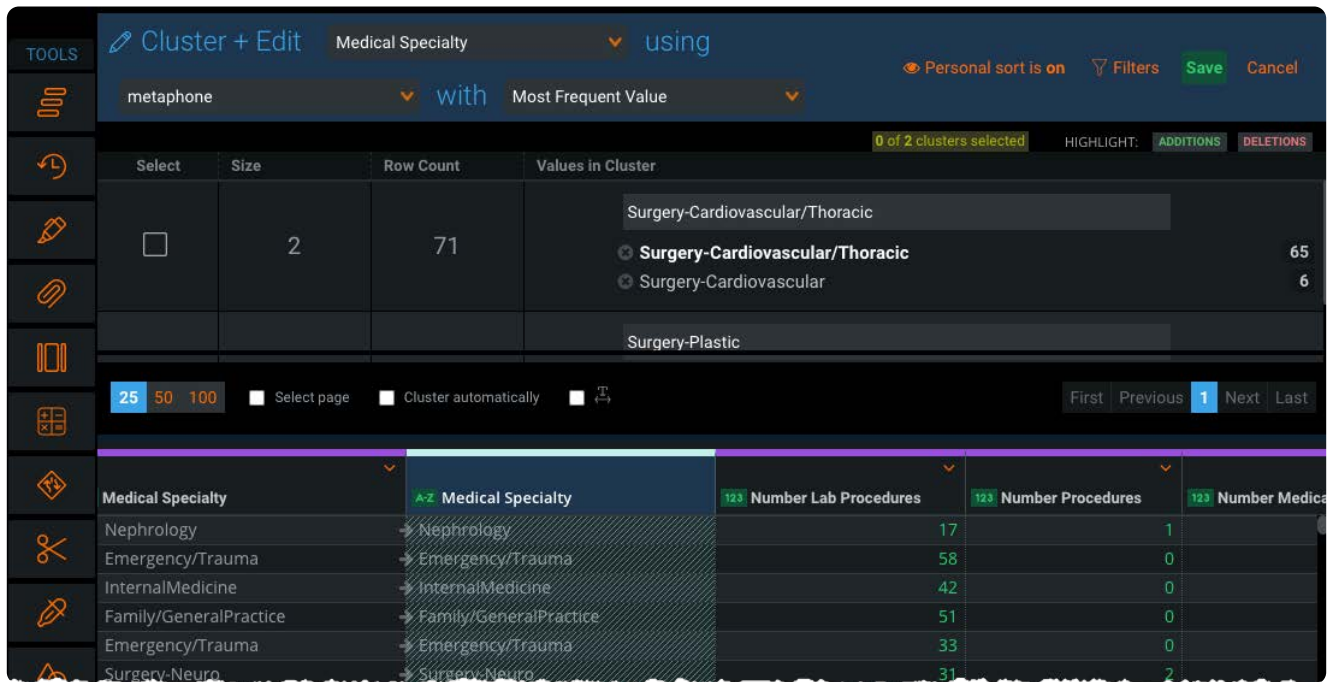
SCENARIO	DESCRIPTION	EXAMPLE
Corrections and inconsistencies	Bulk correction of data entry mistakes, spelling errors, and use of different abbreviations or shorthand conventions.	Acme Co., Acme Company, Acme Comp.
Reclassification	Reclassifying detailed values into an aggregate value.	"12oz soda" and "8oz soda" both become "soda."
Consolidation	Consolidating consistent, but different, values that emerge when the data from different systems is combined in a single column.	One data source refers to "soda" consistently and another data source refers to "tonic" consistently.

Perform Cluster + Edit on a column

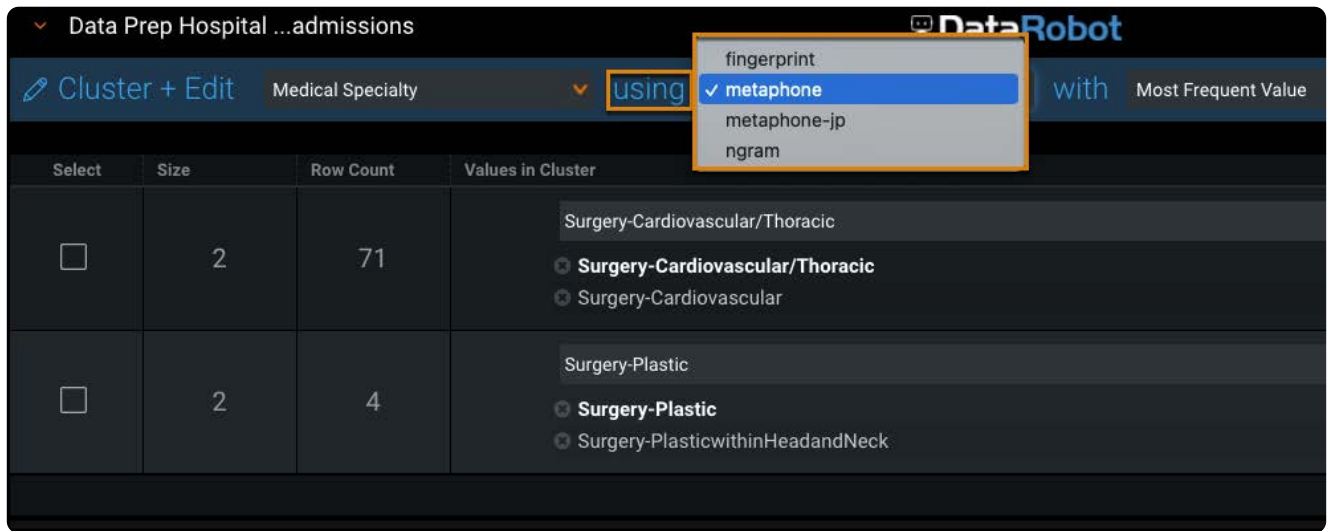
1. Locate the column you want to normalize.
2. Hover over the column menu icon  and click **Other > cluster + edit**.



3. The **Cluster + Edit** pane opens.

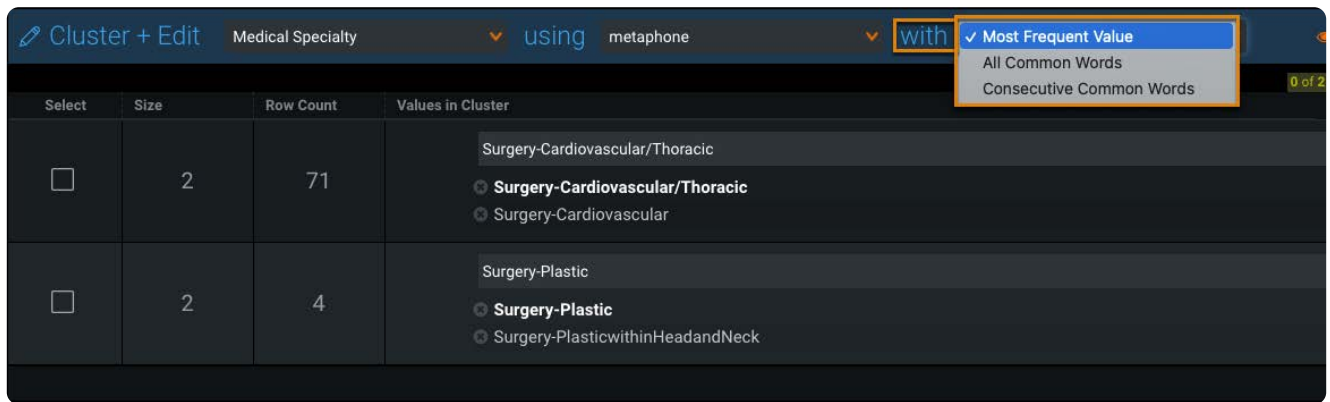


4. In the **using** field, use the dropdown menu to select an algorithm to use for the clustering operation.



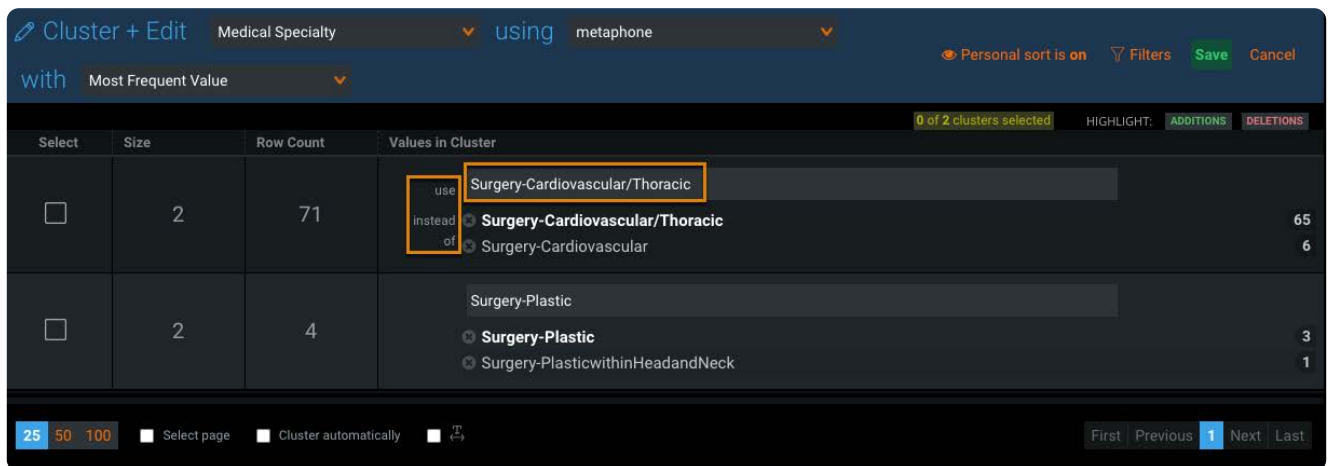
See [Clustering algorithms](#) for details.

- In the **with** field, use the dropdown menu to select an algorithm to use for the output option. The algorithm you choose determines the new value suggestion for the cluster.



See [Output algorithms](#) for details.

- To change the value Data Prep has suggested for the cluster, enter a term in the **Use** field.



- To prevent Data Prep from replacing a value, click the **X** next to the value.

Cluster + Edit Medical Specialty using metaphone

Personal sort is on Filters Save Cancel

with Most Frequent Value

0 of 2 clusters selected HIGHLIGHT: ADDITIONS DELETIONS

Select	Size	Row Count	Values in Cluster	
<input type="checkbox"/>	2	71	use Surgery-Cardiovascular/Thoracic instead of Surgery-Cardiovascular	65 6
<input type="checkbox"/>	2	4	Surgery-Plastic Surgery-Plastic Surgery-PlasticwithinHeadandNeck	3 1

25 50 100 Select page Cluster automatically

First Previous 1 Next Last

8. In the **Select** column, click the check box next to a cluster you want to update and click **Save**.

Cluster + Edit Medical Specialty using metaphone

Personal sort is on Filters Save Cancel

with Most Frequent Value

1 of 2 clusters selected HIGHLIGHT: ADDITIONS DELETIONS

Select	Size	Row Count	Values in Cluster	
<input checked="" type="checkbox"/>	2	71	Surgery-Cardiovascular/Thoracic Surgery-Cardiovascular/Thoracic Surgery-Cardiovascular	65 6
<input type="checkbox"/>	2	4	Surgery-Plastic Surgery-Plastic Surgery-PlasticwithinHeadandNeck	3 1

9. Continue making individual cluster edits.

Perform bulk editing on clusters

Use bulk editing to quickly normalize all clusters on a page.

1. Select 25, 50, or 100 clusters per page.

Cluster + Edit Medical Specialty using metaphone with

Most Frequent Value

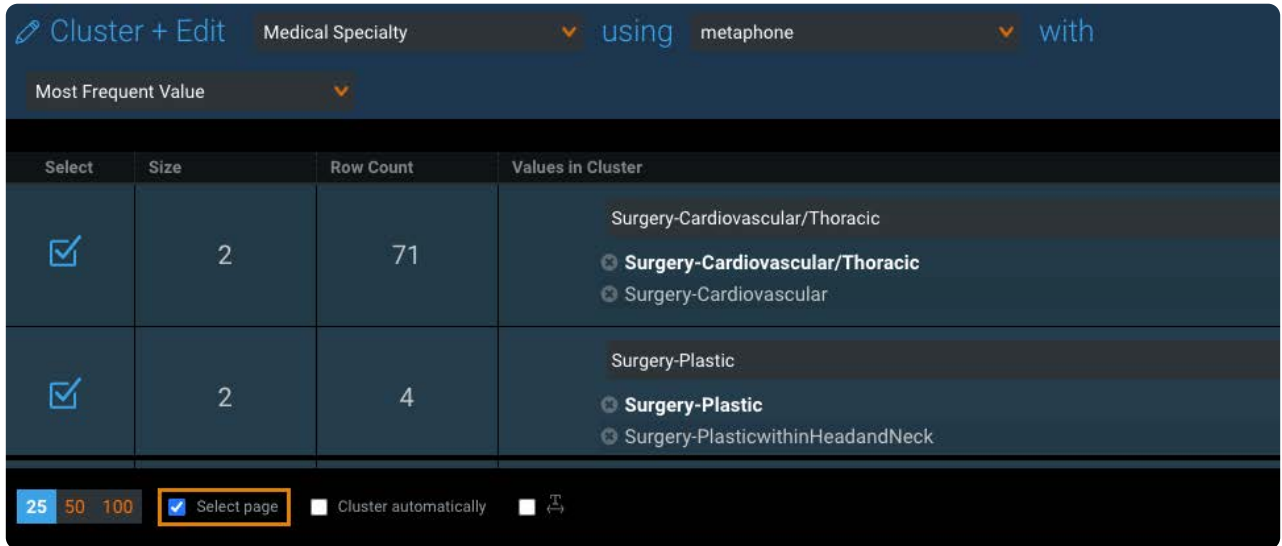
Select	Size	Row Count	Values in Cluster
<input type="checkbox"/>	2	71	Surgery-Cardiovascular/Thoracic Surgery-Cardiovascular/Thoracic Surgery-Cardiovascular
<input type="checkbox"/>	2	4	Surgery-Plastic Surgery-Plastic Surgery-PlasticwithinHeadandNeck

25 50 100 Select page Cluster automatically

The bulk edit operation is restricted to one page of clusters.

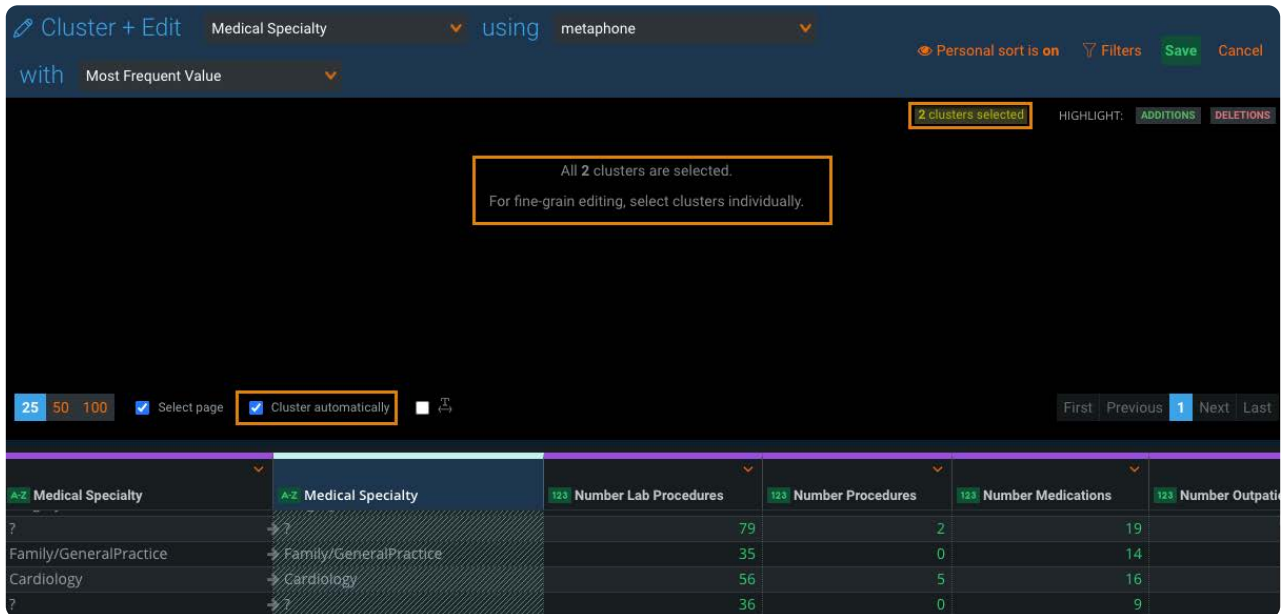
2. To perform bulk editing, do one of the following:

- Click **Select page** to select all clusters on the page.



Use this method if you want to review the suggested replacements and make edits before saving.

- Click **Cluster automatically** to select all of the clusters in the dataset.



Use this method if you are sure you want to accept all of the suggested replacement values.

3. Click **Save** to update the clusters.

For each cluster, all values change to the suggested value.

Tools for working with clusters

The following tools provide visual queues to better recognize how the suggested value for a cluster was derived.

Cluster + Edit Medical Specialty using metaphone

with Most Frequent Value

0 of 2 clusters selected HIGHLIGHT: ADDITIONS DELETIONS

Select	Size	Row Count	Values in Cluster	
<input type="checkbox"/>	2	71	Surgery-Cardiovascular/Thoracic Surgery-Cardiovascular/Thoracic Surgery-Cardiovascular	65 6
<input type="checkbox"/>	2	4	Surgery-Plastic Surgery-Plastic Surgery-Plastic withinHeadandNeck	3 1

25 50 100 Select page Cluster automatically

First Previous 1 Next Last

ELEMENT

DESCRIPTION

- 1** Fixed-width font

By default, cluster values display in a variable-width font. Click this option to display cluster values in a fixed-width font. The fixed-width option aligns all text characters, which allows you to more easily identify extra spaces within a cluster value and differentiate characters across the clusters.
- 2** Highlight tools

Highlighting lets you recognize how the suggested cluster replacement value was derived. The **Additions** tool highlights the characters that will be added for all common characters. The **Deletions** tool indicates where deletions will be made in order to derive the common characters. Deletions are condensed into a red X. You can enable the **Additions** and **Deletions** tools simultaneously.

Clustering algorithms

The **using** clustering algorithm defines the values that should be grouped together.

Note

For all clustering algorithms, blanks and nulls are not included when building a cluster.

The following algorithms are available in the application.

metaphone

The metaphone algorithm, the default selection, groups words together based on their English language pronunciation. It is classified as a “phonetic” algorithm because it is based on how similar or different the text would sound if spoken. This algorithm is particularly useful when working with manually entered data (where misspellings may occur) and data appended from multiple source systems (where minor variations may occur).

ngram

The ngram algorithm breaks the data in the column into a specified number (n) of characters. These “chunks” (or grams) of text are then compared based on the probability of what might follow each. Ngram algorithms are often used by search engines—as a user enters characters into the search bar, the engine examines the probability of what form the final search terms might take and makes suggestions as the user types.

fingerprint

The fingerprint algorithm groups similar values into a cluster where the only differences are punctuation, word order, and capitalization. Fingerprint algorithms are often used to match names, for example, "Adèle Smith" and "SMITH, ADELE".

Output algorithms

The **with** output option determines default replacement values for clusters. The output options attempt to make the best recommendation for the **New Value** replacement. The replacement value can always be manually edited to meet your specific business requirements.

The application provides the following algorithms.

Most Frequent Value

The Most Frequent Value output algorithm, the default selection, builds clusters using the value that occurs most frequently in the cluster.

All Common Words

The All Common Words output algorithm builds clusters using a string of matching words, starting at the beginning of the string, regardless of order. The frequency at which each string occurs then determines the **New Value**.

Example

Apple Computer Corporation

Apple Computer Inc

Apple Corporation Computer

Apple Computer

Apple Corp Computer

New Value: Apple Computer

The algorithm used to build your clusters affects the suggested value.

- Because **metaphone** attempts to preserve the semantic meaning of the words in your clusters, you may notice that some of the suggested values do not strictly reflect all common words in your clusters. This may be the case when punctuation is included in your clusters.
- The ngram algorithm must be used in order to include non-consecutive, common words in the cluster.

Consecutive Common Words

The Consecutive Common Words output algorithm builds clusters using the longest sequence of matching consecutive words, starting at the beginning of the string. Values that occur in less than 10% of the cluster are not included when determining the **New Value** recommendation. Note that most punctuation does not interrupt the sequencing for the match.

Example

Apple-Computer

Apple Computer

Apple ComputerAG

Apple Computer Corp

Apple Computer Corporation

Apple Computer Inc

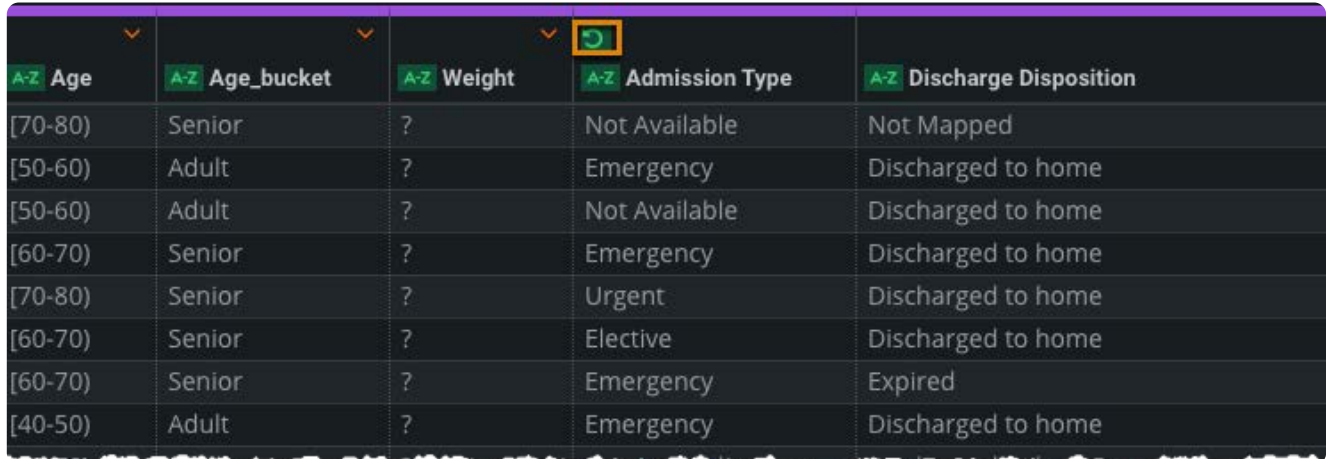
New Value: Apple Computer

View column lineage

Enable column lineage mode in the **Steps** pane to identify the project steps that resulted in the selected column.

To view the lineage of a column:

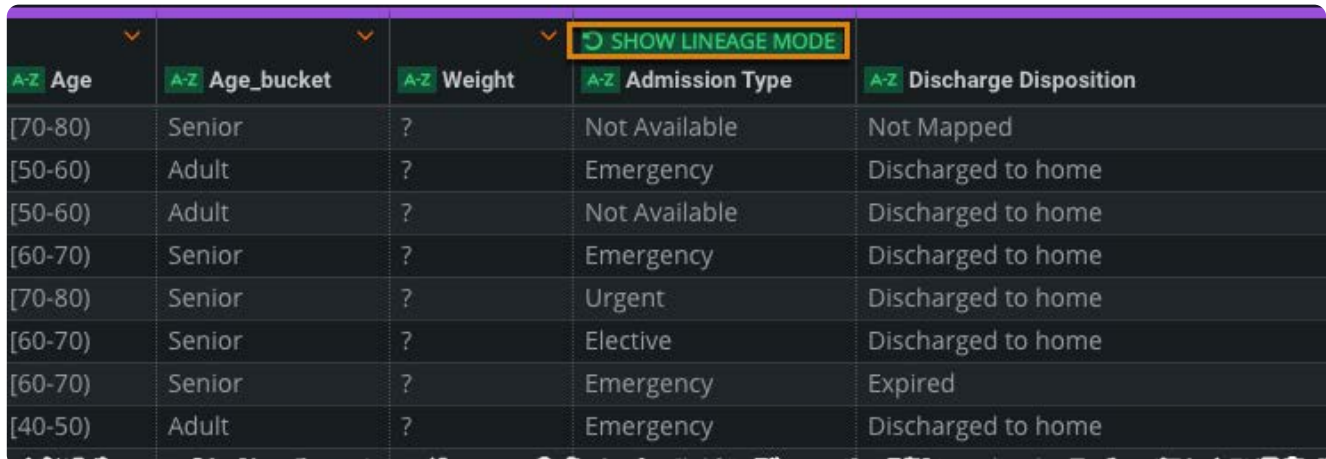
1. Hover over the  column icon.



A screenshot of a table with five columns: Age, Age_bucket, Weight, Admission Type, and Discharge Disposition. The 'Admission Type' column header has a small circular icon with a refresh symbol, which is highlighted by a yellow box. The table contains data rows with values like '[70-80)', 'Senior', '?', 'Not Available', and 'Not Mapped'.

Age	Age_bucket	Weight	Admission Type	Discharge Disposition
[70-80)	Senior	?	Not Available	Not Mapped
[50-60)	Adult	?	Emergency	Discharged to home
[50-60)	Adult	?	Not Available	Discharged to home
[60-70)	Senior	?	Emergency	Discharged to home
[70-80)	Senior	?	Urgent	Discharged to home
[60-70)	Senior	?	Elective	Discharged to home
[60-70)	Senior	?	Emergency	Expired
[40-50)	Adult	?	Emergency	Discharged to home

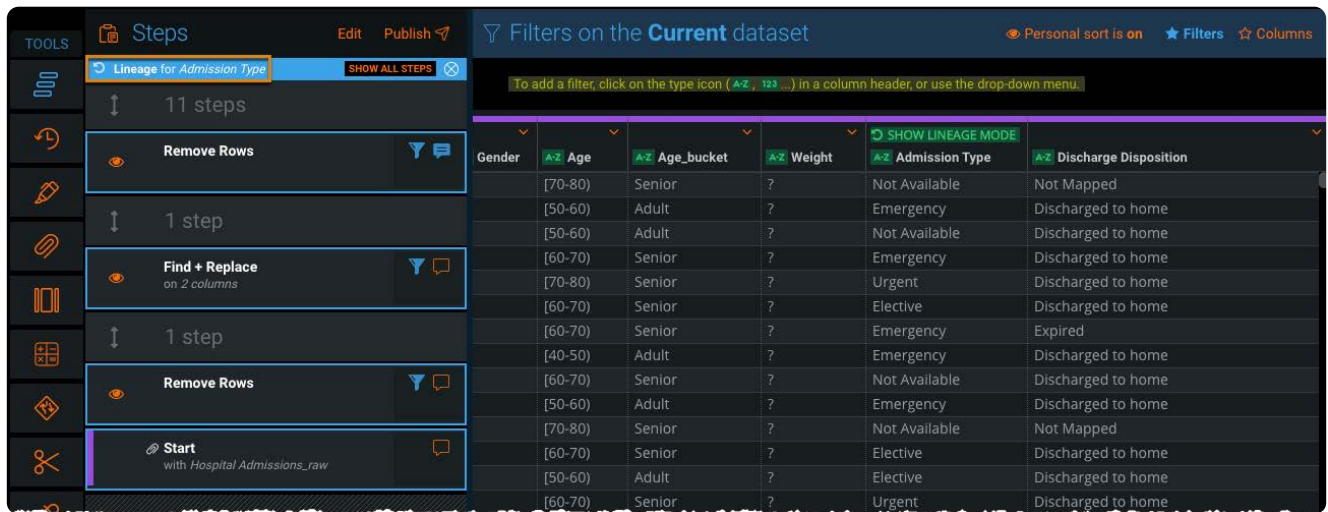
2. Click the **Show Lineage Mode** link that displays.



A screenshot of the same table as above, but now a link labeled 'SHOW LINEAGE MODE' is visible in the 'Admission Type' column header, highlighted by a yellow box. The rest of the table content remains the same.

Age	Age_bucket	Weight	SHOW LINEAGE MODE	Discharge Disposition
[70-80)	Senior	?	Not Available	Not Mapped
[50-60)	Adult	?	Emergency	Discharged to home
[50-60)	Adult	?	Not Available	Discharged to home
[60-70)	Senior	?	Emergency	Discharged to home
[70-80)	Senior	?	Urgent	Discharged to home
[60-70)	Senior	?	Elective	Discharged to home
[60-70)	Senior	?	Emergency	Expired
[40-50)	Adult	?	Emergency	Discharged to home

Data Prep outlines the step-level transformations that contributed to the selected column's state.

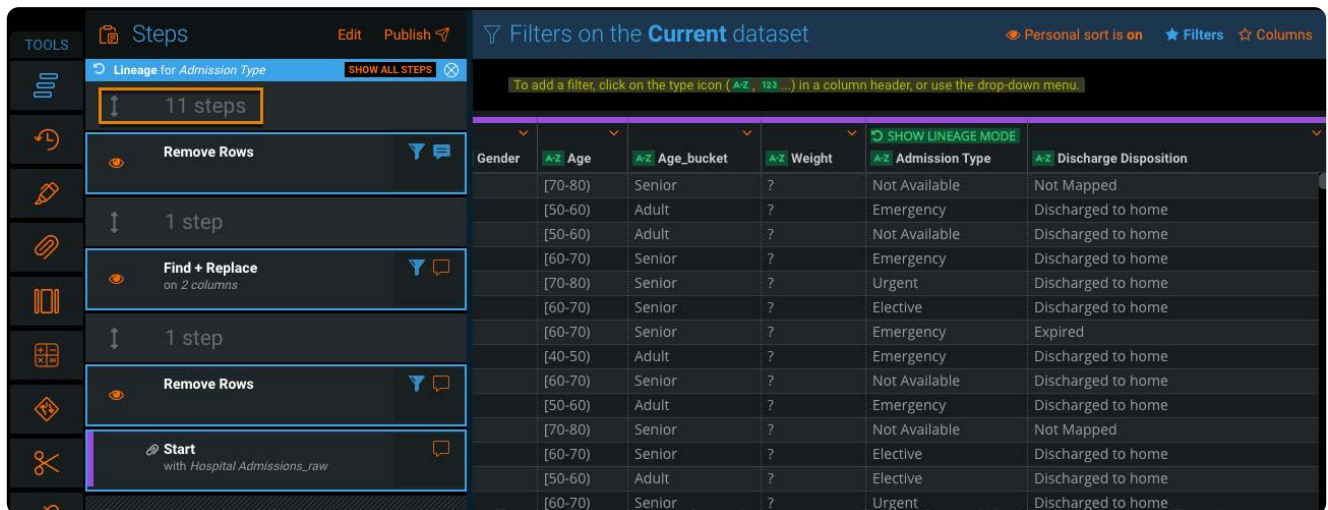


Use the outlines to identify the steps that affected the column or changed its data. If there are steps in the Editor that did not affect the column, those steps are grayed out, collapsed, and labeled to note the number of collapsed steps.

Lineage mode options

Following are options you can use when working in lineage mode.

- Click any grayed-out step to expand the associated collapsed steps.



- Click **Show all Steps** in the orange lineage mode header to expand all collapsed steps in your script.

The screenshot shows the 'Steps' editor on the left with 'Lineage for Admission Type' selected. The main data table has a header row with a green 'SHOW LINEAGE MODE' button. The data table contains the following rows:

Gender	Age	Age_bucket	Weight	Admission Type	Discharge Disposition
[70-80]	Senior	?	Not Available	Not Mapped	
[50-60]	Adult	?	Emergency	Discharged to home	
[50-60]	Adult	?	Not Available	Discharged to home	
[60-70]	Senior	?	Emergency	Discharged to home	
[70-80]	Senior	?	Urgent	Discharged to home	
[60-70]	Senior	?	Elective	Discharged to home	
[60-70]	Senior	?	Emergency	Expired	
[40-50]	Adult	?	Emergency	Discharged to home	
[60-70]	Senior	?	Not Available	Discharged to home	
[50-60]	Adult	?	Emergency	Discharged to home	
[70-80]	Senior	?	Not Available	Not Mapped	
[60-70]	Senior	?	Elective	Discharged to home	
[50-60]	Adult	?	Elective	Discharged to home	
[60-70]	Senior	?	Urgent	Discharged to home	

- Click X in the lineage mode header to close lineage mode.

The screenshot shows the same interface as above, but the 'SHOW LINEAGE MODE' header is now greyed out, indicating it has been closed. The data table content remains the same.

Note

Lineage mode closes automatically when you mute a step in the Steps Editor pane or begin making new transformations in the project.

Example

A project has the following steps:

1. Import a base dataset for customer contact information with a column for *int'l cell numbers*. In that column, all numbers follow this format: *+44-2071838750*.
2. Perform a split operation on the dash in the *int'l cell numbers* column to create two new columns: *country code* and *cell number*.
3. Rename the first newly created column: *country code*.
4. Perform a **Find + Replace** operation on the *country code* column to remove the preceding **+** character.
5. Rename the second newly created column: *cell number*.
6. Use the column tool to hide the original *int'l cell numbers* column.

When you enable Column Lineage mode for the *cell number* column, the second and fifth steps above are highlighted in the Steps Editor pane because those steps directly affect the data in the *cell number* column—the second step is the origin for the data and the fifth step is the new column name. All other steps are grayed out and collapsed because they do not affect the column.

Note

In addition to lineage mode, a column's header color provides a quick reference to indicate the original data source for the column's data. The color of the input step for the data source is used to identify all columns originating from that source. If there is no input data source for the column, for example, the column was created as the result of a compute column operation, then the column is color-coded with the project's color.

Automation and operationalization

You can automate and operationalize your projects and datasets using Data Prep. Using Automatic Project Flows (APFs), you can compute the entire sequence of data prep Steps across Data Prep projects, datasets, and AnswerSets to produce an end-to-end, automated output Flow for your data.

These pages describe:

TOPIC	DESCRIBES...
Automatic Project Flows	Intelligently operationalize curated data flows.
Automation	Automation is a legacy feature that allowed you to automate individual projects and datasets.

Automation

Note

Automation is a legacy feature that provided you the option to automate individual projects and datasets. [Automatic Project Flows \(APF\)](#), introduced in the 2019.1 release, allows you to intelligently operationalize curated data flows. The new APF feature computes the entire sequence of data prep steps across Data Prep projects, datasets and AnswerSets to produce an end-to-end, automated output Flow for your data. For customers who are currently using the 2018.2 Automation feature and are ready to upgrade their automated jobs to APF, contact your DataRobot representative for assistance.

There are two types of workload automation that reduce the number of repetitive tasks taken to produce AnswerSets: library automation and project automation.

Library automation

When you automate a data library dataset, you schedule it to automatically pull an update from its source based on a schedule you define. During the automation process, a dataset is updated with new versions of the data using the import and parse options specified when the file was originally uploaded into the data library. However, when you set up a dataset for automation, you have the option to modify those parse options.

Note

You cannot automate datasets on your local system.

Project automation

When you schedule a project for automation, you set it up to automatically publish an AnswerSet to the data library based on the schedule and parameters you define. The AnswerSet can also be exported to an external data source, for example AWS S3.

Note

Project lenses are essential for project automation because they define the publishing points to use for your automated jobs. In order to automate a project, you must have a lens defined for each point in the project where you want to publish data. You must have at least one lens defined in your project, otherwise no data can be published. For more information on lenses, see the article for [Project lenses](#).

After you configure automation schedules for data library files and/or projects, both are collectively referred to as automation “jobs.” The [Automation dashboard](#) provides you with details of all automation schedules and the status of all automation jobs.

Data library automation configuration page

To open the data library automation configuration page:

1. Open the data library.
2. Locate the file you want to automate.
3. Click the **More Actions** button that displays, then select the **Automation** option. The configuration page opens:

The screenshot shows the 'LIBRARY AUTOMATION' configuration page for a job named 'CDM_DRUG_ERAS.vw'. The page is divided into several sections:

- Job Name and Description:** Fields for 'JOB NAME' (CDM_DRUG_ERAS.vw) and 'JOB DESCRIPTION'.
- Schedules:** A section with a '+ Add' button and a message: 'No upcoming schedules configured at this time'.
- Notifications:** A section with two rows for 'EMAILS' and 'TEXTS', each with a field for 'email address' or 'text address' and a 'set' button.
- Importing From:** A section with 'DATA SOURCE' (M1 CDMS) and 'FILE PATH' (J:\SEARCH\ERAS\CDM_DRUG_ERAS.VW).
- Import Options:** A section with 'Separator' (FORMAT OF THE DATA) and 'Additional Options' (TEXT, VALUE SEPARATOR TYPE, SINGLE SEPARATOR, TEXT, LINE SEPARATOR TYPE, and LINE SEPARATOR).

Handwritten annotations in orange curly braces label these sections: 'Job Name and Description' on the left, 'Schedules' on the left, 'Importing from' on the left, 'Import Options' on the left, and 'Notifications' on the right.

Job name and job description

The dataset's name and description are listed in these fields. These are initial default values from when the file was originally imported into the data library. They can be changed here by entering new information into the fields.

Note

You may also notice a check box option for **Set me as the owner of this automated schedule**. This option only appears if you are not the person who initially set up this dataset for automation or not the person who currently owns its automated schedule. Ownership is significant because it provides a way to identify and audit users who are running automated jobs in the system. Typically, this option is used when automation responsibilities are transitioned to a new person in an organization. If you take ownership of an automation job, you must have all of the permissions that are required to perform every operation performed by the automation.

Schedules

Any upcoming schedules for the dataset are displayed here. The **Add** button allows you to set up new schedules. The **Deactivate** link in this pane allows you to indefinitely suspend all scheduled jobs for this dataset until you return and click the **Reactivate** button. To set up schedules, see [Set up a data library dataset for automation](#).

Notifications

Email notifications can be sent to notify users of either a successful upload into the data library or errors that have occurred. To set up notifications, see [Set up a data library dataset for automation](#).

Importing from

These are the connection parameters inherited from the most recent upload of the dataset. To change these connection parameters, manually upload a new version of the file to the data library with new parameters. Automation will then use the new connection parameters in its next scheduled update.

Import parsing options

For file-based datasets, the import parse options are displayed below the connection details. The import options are inherited from the most recent version of the dataset but are editable here.

Note

If you *manually* import another version of this dataset into the data library, the parse options you select for the manual upload will *not* be inherited from the automated version.

Set up a data library dataset for automation

You set up a dataset for automation by:

- Setting up schedules
- Setting up notifications
- Saving your automation configuration settings

Set up schedules

Click **Add** to set up a new time for the dataset to be updated by automation. The default setting for dataset automation frequency is to repeat on the time and day you specify. The **Repeat** toggle button lets you switch the automation to run **Once** at the time you specify.

To set up recurring updates:

1. Use the up and down arrows to adjust the time.

2. Toggle the **PM** or **AM** button to select the correct period.
3. Select the frequency: **week**, **day**, or **month**. The default is **week**. Click in the field to make a different selection.
4. Depending on your frequency selection, specify the day of the week or date in the month.
5. Click **Okay** to add the schedule. Your newly added schedule then appears. Click the pencil icon to edit it or the **X** button to delete it.

Note

- The time you select is based on your current time zone.
- The time, day, or date you select here must be in the future. For example, if it is currently 1pm on Monday and you set up the automation to run at 10am every Monday, the automation will not run today for this file.
- Datasets that are on your own local system cannot be automated.

To schedule a single update:

1. Click in the date field to open a calendar picker.
2. Use the up and down arrows to adjust the time.
3. Toggle the **PM** or **AM** button to select the correct period.
4. Click **Okay** to add the schedule. Your newly added schedule appears. Click the pencil icon to edit it or the **X** button to delete it.

Note

- The time you select for a schedule is based on your current time zone.
- When configuring a file's automation to run only once, do *not* set the job's start time too near to the current time. Your local computer's clock may not be precisely in sync with the web server that will process the job. If your local computer's clock is running behind the web server's clock, the time you specify for the job may have already passed on the web server. In this case, your job will not start.
- If you want to test one automated run of this dataset, use the **Add to Queue** feature instead of setting it up to run **Once**. For details on this feature, see [Save your automation configuration settings](#).
- Datasets that are on your own local system cannot be automated.

Review the following considerations when setting up a dataset for automation:

- If an automated project uses this dataset for input, you must ensure a safe buffer of time for this dataset update to finish uploading in the data library before the automated run of the project begins.
- The time you specify here is when this job will be added to the queue for uploading and not necessarily the start time for the automated import.

Set up notifications

Email notifications can be sent to notify users of either a successful upload into the data library or errors that have occurred. An error email provides a link to the file's log file where you can determine the cause of any errors.

To set up notifications:

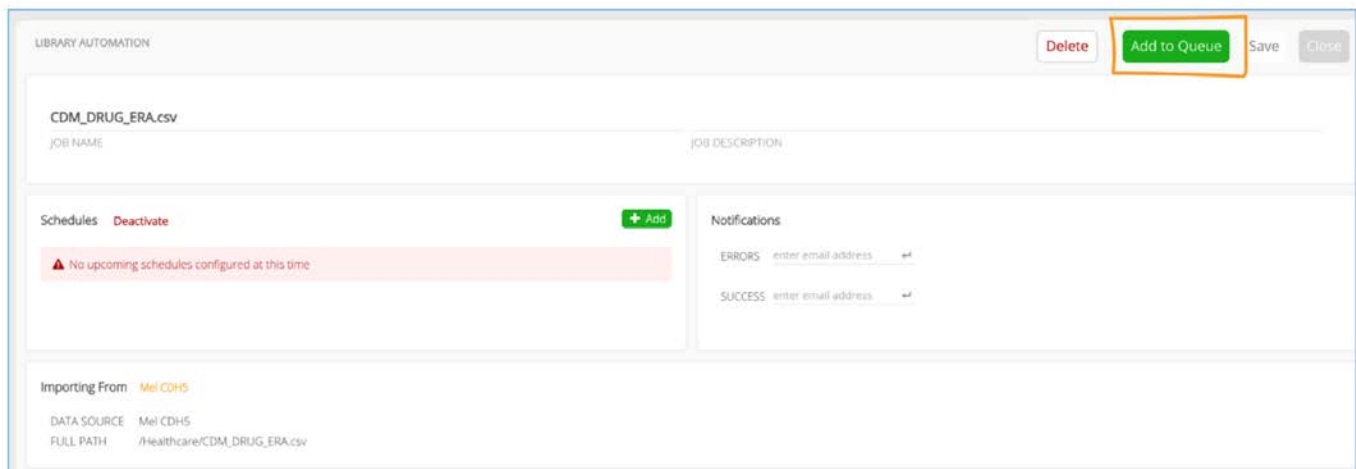
1. Click the dropdown menu to select which type of email notification to send: "Errors" or "Success".
2. Add the email address and press Enter.

Important considerations:

- An email address can only be added once for each notification type.
- Recipients must have the required system permissions to view the automation results.

Save your automation configuration settings

Click **Save** at the top of the configuration form to save all of your settings. After saving, notice the **Add to Queue** button that displays.

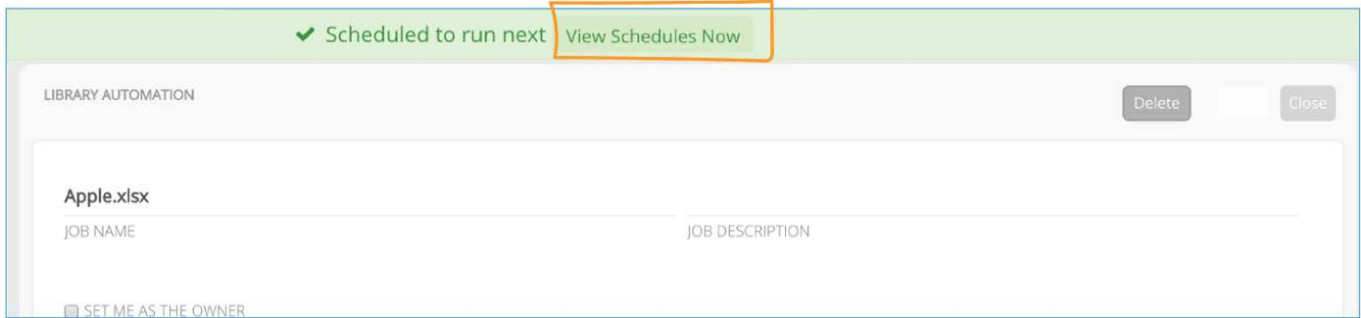


The screenshot shows a web interface for configuring automation. At the top right, there are buttons for 'Delete', 'Add to Queue' (highlighted with an orange box), 'Save', and 'Close'. Below this, the job name is 'CDM_DRUG_ERA.csv'. The interface is divided into 'Schedules' and 'Notifications' sections. The 'Schedules' section has a 'Deactivate' link and an '+ Add' button, with a message: 'No upcoming schedules configured at this time'. The 'Notifications' section has two rows: 'ERRORS' and 'SUCCESS', each with an 'enter email address' input field. At the bottom, there is an 'Importing From' section with 'Mel CDHS' and a table of details: DATA SOURCE: Mel CDHS, FULL PATH: /Healthcare/CDM_DRUG_ERA.csv.

The button allows you to add this automation job to the queue of upcoming jobs that will be run the next time automation starts. This option is useful if you want to test out this automation configuration without having to wait for its scheduled run time.

Tip

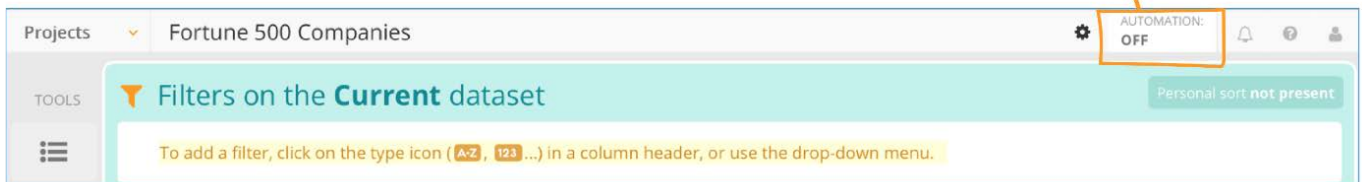
The Automation pane provides details of when the next automated run is scheduled to start. You can quickly navigate to the **Schedules** pane by clicking the **View Schedules Now** link that displays in the header after you click the **Add to Queue** button.



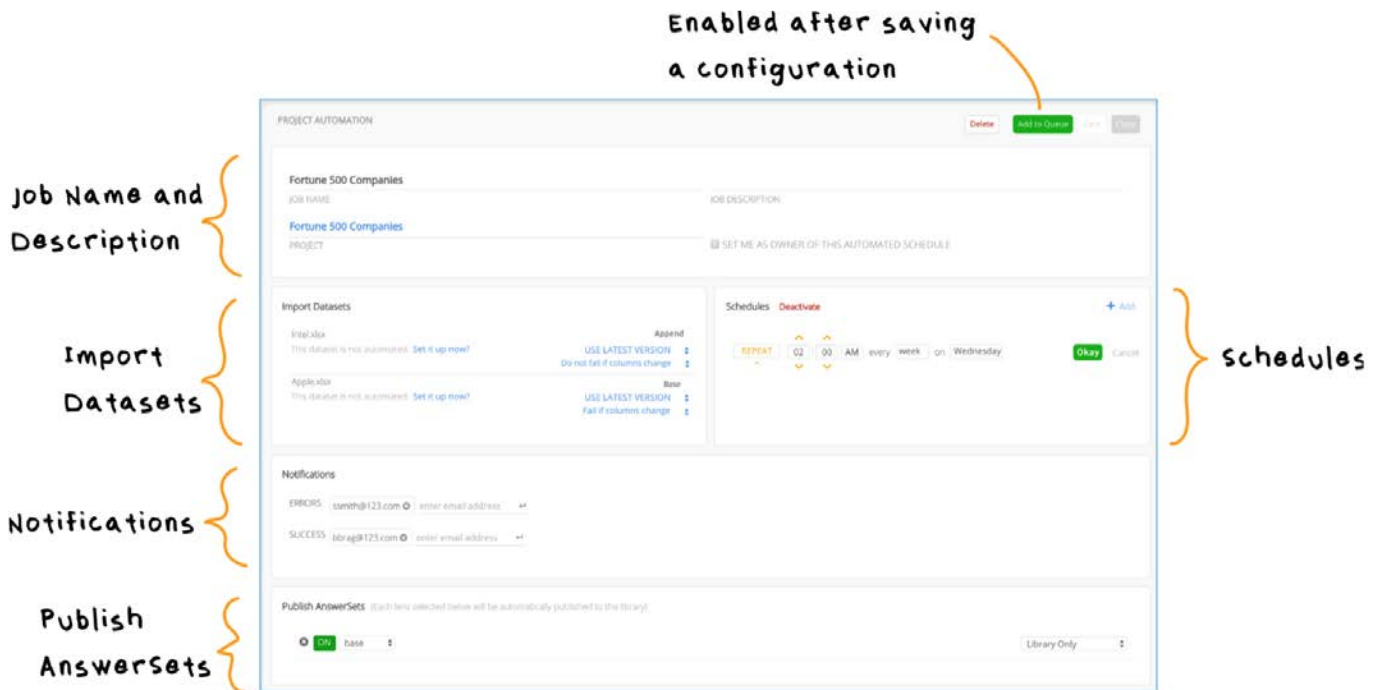
Project automation configuration page

To open the project automation page, open your project and click the automation status button:

Automation Status



The **Automation** configuration page opens:



Job name, job description, and project

The job's name and description are listed in these fields. This is the name and description provided when the project was originally created. It can be changed here for the automated version of your project by entering new information into the fields.

The project field provides a link that automatically opens the project you are setting up for automation. This link is particularly useful if you have multiple Versions of a project but are automating one specific Version of that project in this configuration form.

Note

You may also notice a check box option for **Set me as owner of this automated schedule**. This option only appears if you are not the person who initially set up this project for automation or not the person who currently owns its automated schedule. Ownership is significant because it provides a way to identify and audit users who are running automated jobs in the system. Typically, this option is used when automation responsibilities are transitioned to a new person in an organization. If you take ownership of an automation job, you must have all of the permissions that are required to perform every operation performed by the automation.

Import datasets

The datasets you have already imported into your project are listed here. If no datasets are listed, verify that you have saved the most recent set of changes to your project steps.

Your "Base" dataset is listed and any "Lookup" or "Append" datasets for the project are listed above it.

If this project uses a dataset that is set up for automation in the data library, the schedule is displayed here. When using an automated dataset, consider its automation schedule and allow a safe buffer of time for the new version to be published in the data library.

To set up a dataset for automation before you configure automation for this project, click the **Set it up now?** link adjacent to the dataset's name. You are taken to the data library scheduling page where you set up the automation parameters and schedule. See [Set up a data library dataset for automation](#).

Note

The **Use Latest Version** default setting refers to the version of the dataset that is used for this automation configuration. See [Set up a project for automation](#) for details on which version you should select for automating this project.

Schedules

Any upcoming schedules for the project are displayed here. The **Add** button allows you to set up new schedules. The **Deactivate** link in this pane allows you to indefinitely suspend all scheduled jobs for this dataset until you return and click the **Reactivate** button. To set up schedules, see [Set up a project for automation](#).

Notifications

Emails can be sent to notify users when automated projects finish updating or have errors. To set up notifications, see [Set up a project for automation](#).

Publish AnswerSets

Select a lens for publishing an AnswerSet. A lens is pinned to a step in your project and creates a publishing point that can be used by automation to publish an AnswerSet. You can save the setup for this project's automation without selecting a lens, but an automated run of this project will not succeed until you select a lens.

Automated projects are automatically published to the data library. However, automation can also be configured to export the published output to an external data source. See [Set up a project for automation](#).

Set up a project for automation

You set up a project for automation by:

- Importing datasets
- Setting up schedules
- Setting up notifications
- Selecting lenses and publishing destinations
- Saving your automation configuration settings

Import datasets

For each dataset used in your project, choose to use the **Latest Version** or **Current Version** for input:

- **Latest Version** uses the most up-to-date version of the dataset in the data library when the automated job is run.

Note

- Using **Latest Version** will result in a new Version of your project each time this automated configuration runs. When selecting the latest version, an additional option is available to specify if an automated run should fail because the latest version of the dataset has a different layout (schema)—for example new columns added, removed columns that are not used in the project's steps, different column types for existing columns, new order, etc.
- At least one of the datasets used for automating this project must be **Latest Version**. Otherwise, if no changes occur in the input datasets after an automated run of your project, the platform will not re-run this job.

- **Current Version** pins the dataset in its current state for all future automated runs. Using **Current Version** may be useful when a static dataset serves as a reference table for your project.

Note

If there have been no changes to the input datasets since the last project automation run, automation will not run again for the project until there are changes to the project. Therefore, at least one of the datasets used for automating this project must be **Latest Version**.

Set up schedules

Click **Add** to set up a new time for this project to run. The default setting for project automation frequency is to repeat on the time and day you specify here. The **Repeat** button is a toggle that you can click to switch the automation to run **Once** at the time you specify.

To set up recurring runs:

1. Use the up and down arrows or enter values in the fields to adjust the time.
2. Toggle the **PM** or **AM** button to select the correct period.
3. Select the frequency: **week**, **day**, or **month**. Note that **week** is the default. Click in the field to make a different selection.
4. Depending on your frequency selection, specify the day of the week or date in the month.
5. Click **Okay** to add the schedule.
6. Your newly added schedule displays. Click the pencil icon to edit it or the **X** button to delete it.

Note

- The time you select here is based on your current time zone.
- The time, day, or date you select here must be in the future. For example, if it is currently 1 PM on Monday and you set up the automation to run at 10 AM every Monday, then automation will not run today for this file.

To set up a single run:

1. Click in the date field to open a calendar picker.
2. Use the up and down arrows to adjust the time.
3. Toggle the **PM** or **AM** button to select the correct period.
4. Click **Okay** to add the schedule.
5. Your newly added schedule displays. Click the pencil icon to edit it or the **X** button to delete it.

Note

- The time you select here is based on your current time zone.
- When configuring a project's automation to run only once, do *not* set the job's start time too near to the current time. Your local computer's clock may not be precisely in sync with the web server that will process the job. If your local computer's clock is running behind the web server's clock, the time you specify for the job may have already passed on the web server. In this case, your job will not start.
- If you want to simply test one automated run of this project, use the **Add to Queue** feature instead of setting it up to run **Once**. For details, see [Save your project automation configuration settings](#).

Important considerations when setting up a project for automation:

- If a project's automation depends on input from an automated data library file or an AnswerSet published from another automated project, ensure a safe buffer of time for all input updates to finish before the automated run of the project begins.
- The time you specify in the automation set-up is when this project will be added to the queue for publishing an AnswerSet, and not necessarily the publishing start time.

Set up notifications

Emails can be sent to notify users when automated projects finish updating or have errors. An error email provides a link to the project's log file where you can determine the cause of any errors.

To set up notifications:

1. Click the dropdown menu to select which type of email notification to send: "Errors" or "Success".
2. Add the email address and press Enter.

Important considerations:

- An email address can only be added once for each notification type.
- Recipients must have the required system permissions to view the automation results.

Select lenses and publish destinations

To add a lens:

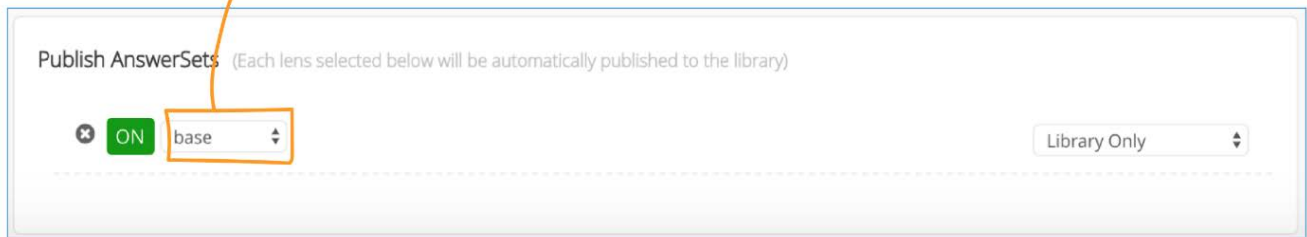
1. Click the green **Add** button:



A lens from the project is added to this automation configuration. By default, the lens that occurs earliest in your project steps is selected.

2. To change the default selection, click the dropdown menu and select a different lens that currently exists in your project.
3. To add additional lenses to be used for this automated run of your project, click the **Add** button and continue to select lenses.

The lens on the earliest step in your Project is selected as the default



To disable a lens:

- Click the green **On** button for the lens to toggle it off

To remove a lens:

- Click the **X** button for the lens.

If you need to add a new lens for automating this project, you will need to open the project and add the lens on the desired step. For more help on adding lenses to your project, see [Use lenses for publishing](#).

The default location for publishing this project when automation runs is **Library Only**. If you want to export the published output to an external data source, in addition to publishing it to the data library, click the dropdown menu and select **Library & Data Source**.

Publish AnswerSets (Each lens selected below will be automatically published to the library) + Add

ON base Toggle detail options ^ Library & Data Source ▾

Fortune 500 Companies.json

NAME

CREATE UNIQUE NAME for every automated publish

JSON

FORMAT

SFTP

DATA SOURCE NAME

Credentials

sftpuser

USERNAME ? PASSWORD ?

DIRECTORY PATH OR DATABASE NAME

- **Name:** The name that will be used for automated versions of this project.
- **Data Source Name:** Click the dropdown menu to select an available data source. !!! note Only the data sources that have been configured for export and that you have permissions to access are displayed in the dropdown menu for **Data Source Name**. Contact your System Administrator if you don't see the data source you want to select for export.
- **Directory Path or Database Name:** Provide the path or database on the data source where the export will be written.
- **Format:** Depending on the data source you select for export, the option to select a file format is also available. Any applicable parsing options are also presented.
- **Credentials:** The user credentials for writing to the selected data source are presented here. You can edit the credentials here.
- **Create unique name:** When enabled, automation appends an underscore and time stamp to the file or table name for each successive automated export so that any previous exports of this project are not overwritten on the data source.

Note

If you enable this option for a JDBC data source, ensure that your system administrator has also enabled the **Automatically Create Table** option in the JDBC Connector form. Otherwise, automation for this project will fail.

Save your project automation configuration settings

Click the **Save** button in the upper right pane to save all configurations you have made for automating this project. After saving the automation schedule, notice the **Add to Queue** button that displays:

PROJECT AUTOMATION Delete Add to Queue Save Close

Fortune 500 Companies

JOB NAME JOB DESCRIPTION

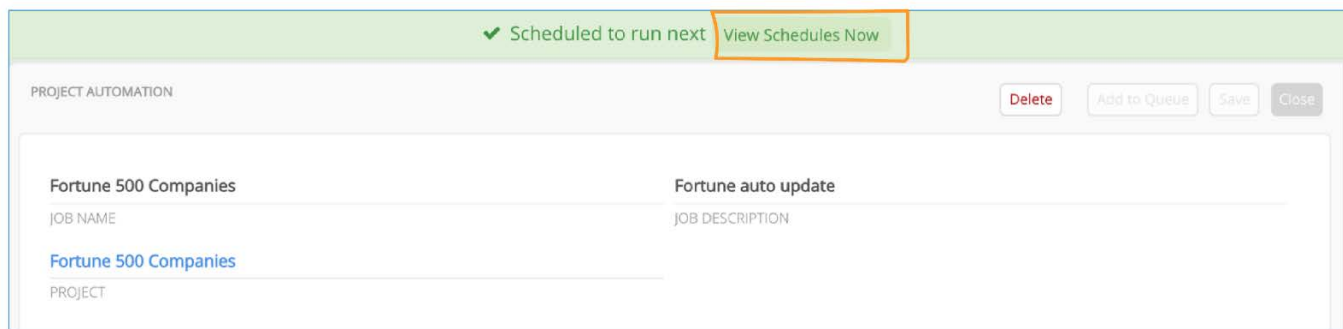
Fortune 500 Companies

PROJECT SET ME AS THE OWNER

This button allows you to add this automation job to the queue of upcoming jobs that will be run the next time automation starts. This option is useful if you want to test out this automated configuration without having to wait for its scheduled run time.

Note

The Automation pane provides details of when the next automated run is scheduled to start and you can quickly navigate to the Schedules pane by clicking the **View Schedules Now** link that displays in the header after you click the **Add to Queue** button.



Automation dashboard

The Automation dashboard provides details and history for all data library files and projects that are set up to be automated. This is where you:

- View your automation usage details.
- View and manage the schedules for automated jobs.
- View job execution history and statuses.
- Re-run failed jobs.

The dashboard is organized by [Schedules](#) and [Job Details](#).

Schedules

The **Schedules** page displays a list of all data library files and projects that are currently configured for automation. To view your automation usage details, mouse over the meters for additional information regarding the number of automated jobs you've already completed and the maximum number you can run for the day, week, or month.

The **Schedules** page can also be filtered in a variety of ways to display:

- **Active** or **Inactive** jobs that have had their automation schedules deactivated.
- Types of jobs—**Project only** or **Library only**.
- Automation jobs that you own.
- Job states—**Success**, **Complete with Error**, **Error**, or **Over Limit**. See [Definition of job states](#) for the meaning of each state.

- Jobs that last finished during a date range that you specify or that will be run in the next automated run based on the range you provide.

Filters for job listing

Automation usage meters

Project icon

Dataset icon

Add to queue of jobs that will be run the next time automation runs

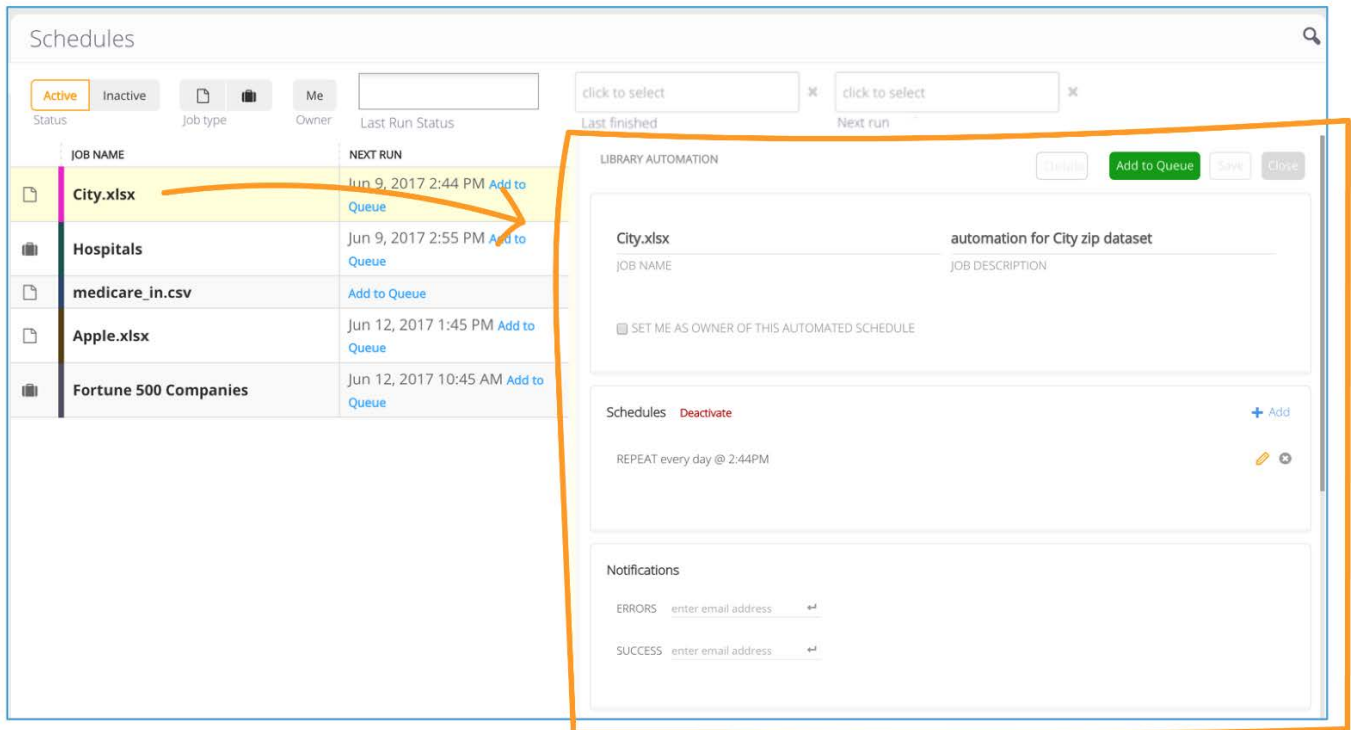
JOB NAME	SCHEDULE	LAST FINISHED	NEXT RUN	UPDATED AT
City.xlsx	REPEAT every day @ 2:44PM	Jun 9, 2017 10:01 AM Success	Jun 9, 2017 2:44 PM Add to Queue	Jun 9, 2017 10:01 AM by Sam
Hospitals	REPEAT every day @ 2:55PM	Jun 9, 2017 10:01 AM Complete with Error	Jun 9, 2017 2:55 PM Add to Queue	Jun 9, 2017 10:01 AM by Sam
medicare_in.csv	Inactive	Jun 5, 2017 3:01 PM Success	Add to Queue	Jun 5, 2017 3:01 PM by Abbey
Apple.xlsx	REPEAT every week on Monday @ 1:45PM	Jun 5, 2017 2:45 PM Success	Jun 12, 2017 1:45 PM Add to Queue	Jun 5, 2017 2:45 PM by Abbey
Fortune 500 Companies	REPEAT every week on Monday @ 10:45AM	Jun 5, 2017 2:45 PM Complete with Error	Jun 12, 2017 10:45 AM Add to Queue	Jun 5, 2017 2:45 PM by Abbey

You can rerun any job by clicking the **Add to Queue** link. This creates an internal schedule on-the-fly for the job, which triggers automation to run the job the next time the automation service wakes up to run regularly scheduled jobs. It's important to keep in mind that system resources must be available in order to run a queued job. For example, the number of threads allocated to run automation jobs must be sufficient. Otherwise, the job will remain in a queued state until resources become available to run it.

Note

- To determine the errors before re-running a job, go to the **Job Details** tab and open the **Results** page for that job run.
- **Add to Queue** for rerunning a job with errors does *not* count against your existing automation guardrail limits.

To make changes to a job's configuration settings, including deactivating it, click the job's name to open its configuration page. Then make and save the configuration changes.

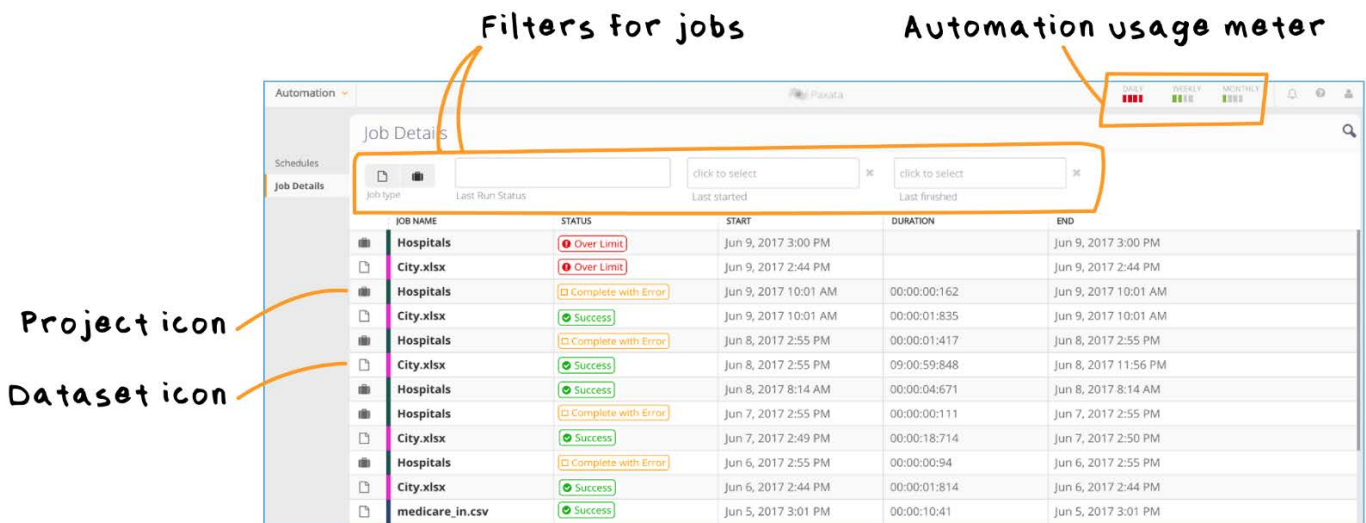


Job Details

The **Jobs Details** page provides an audit trail for every executed automated run—including automated jobs that have been deleted. To view your automation usage details, mouse over the meters for additional information regarding the number of automated jobs you've already completed and the maximum number you can run for the day, week, or month.

You can filter the **Job Details** page in a variety of ways to display:

- Types of jobs—**Project only** or **Library only**.
- Job states—**Success, Complete with Error, Error, Over Limit, Queued, Running**. See [Definition of job states](#) for the meaning of each state.
- Jobs that last started or last finished during a date range that you specify.



To display granular details for a job run, click the row for that job. The **Results** page for the job opens and displays a snapshot of the configuration settings used for this instance of the job run.

The screenshot displays the Automation interface. On the left, the 'Job Details' table lists various jobs with their statuses. An orange arrow points from the 'Fortune 100 Companies' row in the table to the 'Results' panel on the right. The 'Results' panel shows configuration settings for a specific job run, including the 'base dataset' and 'Import Datasets'.

JOB NAME	STATUS
Fortune 100 Companies	Success
Fortune 500 Companies	Complete with Error
Fortune 500 Companies	Complete with Error
Fortune 500 Companies	Complete with Error
Hospitals	Over Limit
City.xlsx	Over Limit
Hospitals	Complete with Error
City.xlsx	Success
Hospitals	Complete with Error
City.xlsx	Success
Hospitals	Success
Hospitals	Complete with Error
City.xlsx	Success
Hospitals	Complete with Error
City.xlsx	Success

Results START - JUN 9, 2017 4:50 PM END - JUN 9, 2017 4:50 PM [Download log](#)

base dataset [View Lens](#)

Library: Success [View AnswerSet](#)

PROJECT AUTOMATION [Delete](#) [Close](#)

Fortune 100 Companies

JOB NAME: Fortune 100 Companies JOB DESCRIPTION: Fortune 100 Companies PROJECT: Fortune 100 Companies

Import Datasets

City.xlsx Update REPEAT every day @ 2:44PM [USE LATEST VERSION](#)

Note

Because this is a snapshot, the job settings may have changed since this automated run.

If this is a project job, click the **View Lens** link to open the project to the lens that was used to publish the AnswerSet for this run. Click the **View AnswerSet** link to view the AnswerSet published by this run.

Opens the Project to the Lens that was used

Opens the AnswerSet

The screenshot shows a web interface for data results. At the top, it says "Results" with a start time of "START - JUN 9, 2017 4:50 PM" and an end time of "END - JUN 9, 2017 4:50 PM". There is a "Download log" link on the right. Below this, a section titled "base dataset" contains a "Library" tab, a "- Success" status, and a "View AnswerSet" link. To the right of this section is a "View Lens" link. A "PROJECT AUTOMATION" section has a "Close" button. Below that, a table titled "Fortune 100 Companies" is shown with columns for "JOB NAME" and "JOB DESCRIPTION". The table contains one row with the value "Fortune 100 Companies" under "JOB NAME" and "PROJECT" under "JOB DESCRIPTION". At the bottom, an "Import Datasets" section shows a file named "City.xlsx" with an update frequency of "Update REPEAT every day @ 2:44PM" and a "USE LATEST VERSION" link. The text "BASE" is also visible in the bottom right of the import section.

If this is a library job, click the **View Dataset** link to open the file in the data library.

Opens the dataset produced
by this job run

The screenshot shows a 'Results' page for a job run. At the top, it displays 'Results' with 'START - JUN 9, 2017 10:01 AM' and 'END - JUN 9, 2017 10:01 AM'. A 'Download log' link is in the top right. Below this, a 'Library' section shows '- Success' and a 'View Dataset' link. An orange arrow points from the handwritten text above to this 'View Dataset' link. The main content area is titled 'LIBRARY AUTOMATION' and contains a table with two columns: 'JOB NAME' and 'JOB DESCRIPTION'. The first row shows 'City.xlsx' as the job name and 'automation for City zip dataset' as the job description. There are 'Delete' and 'Close' buttons in the top right of the table area.

If any errors occurred during the run, they are displayed here. Download the log file for this job by clicking the **Download log** link:

The screenshot shows a 'Results' page for a job run that failed. It displays 'Results' with 'START - JUN 9, 2017 4:48 PM' and 'END - JUN 9, 2017 4:48 PM'. A 'Download log' link is in the top right. Below this, an error message is shown: 'Automation for the Project did not find any new dataset versions to run with'. A 'View Lens' link is next to the error. The main content area is titled 'PROJECT AUTOMATION' and contains a table with two columns: 'JOB NAME' and 'JOB DESCRIPTION'. The first row shows 'Fortune 500 Companies' as the job name and 'Fortune auto update' as the job description. Below this, there is a 'PROJECT' section with 'Fortune 500 Companies' as the project name. There are 'Delete' and 'Close' buttons in the top right of the table area.

Log for this
job run

Details regarding
failure

Definition of job states

The following are possible states for an automated job.

- **Running:** Job run is currently in progress.
- **Success:** Job successfully finished with no errors.
- **Error:** Job run failed.
- **Completed with errors:** Job run completed, but there were errors that prevented a complete run—for example, a job that successfully published to the data library but was unable to export to the specified data source will complete with this type of error.
- **Queued:** When a job is queued an internal schedule is created on-the-fly for it, which triggers automation to run the job the next time the automation service wakes up to run regularly scheduled jobs. However, it's important to keep in mind that system resources must be available in order to run a queued job. For example, the number of threads allocated to run automation jobs must be sufficient. Otherwise, the job will remain in a queued state until resources become available to run it.
- **Over limit:** When a job run exceeds the daily, weekly or monthly guardrail limits, then the job fails with an **Over limit** error.

Important notes:

- Automation guardrails are enforced at the tenant level.
- The *weekly* automation limit is defined as 00:00 Monday–23:59 Sunday.
- A job that ends in error is counted toward your limits, but a retry of the failed job (through **Add to queue**) is *not* counted.

Automatic Project Flows

Note

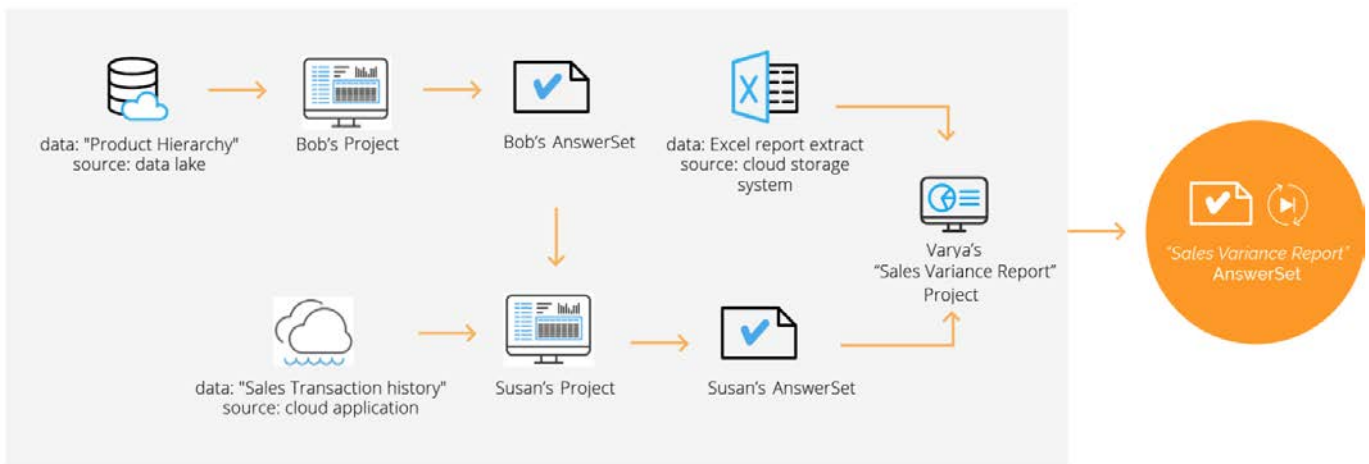
The APF feature must be enabled. If you do not see the **Create Project Flow** button on the top of the Project page, contact your Data Prep System Administrator.

The Data Prep Automatic Project Flows feature (APF) allows you to intelligently operationalize curated data flows. APF computes the entire sequence of data prep steps across projects, datasets, and AnswerSets to produce an automated end-to-end output Flow for your data. Set the Flow to run on a recurring time-based schedule or run it just once to produce an end result AnswerSet. Then manage all runs using APF's monitoring capabilities. Business Analysts and Data Engineers use APF to simplify complex data flows by breaking them into smaller groups of Data Prep projects.

APF lets you operationalize the data flow—each project performs a related or cohesive set of steps for improved readability and limited complexity. After creating projects, you can select the final project in the sequence as your *target* project. APF takes care of the rest—sequencing, preparing, and automating the entire end-to-end flow without requiring manual stitching.

APF helps teams share data and gather input from business and IT leaders. Team members can build Data Prep projects that depend on output AnswerSets created by others. Members complete their data prep work in their own Data Prep project, and then the entire sequence is operationalized from a single target project. APF takes care of the rest with no manual stitching required, regardless of who creates or owns the projects and AnswerSets. Members of the team can monitor the Flow and view a graph to see how their projects and AnswerSets contribute to the Flow's final output.

APF example



In this example, APF produces the end-state "Sales Variance Report" from a series of Data Prep projects and AnswerSets produced by multiple people.

Bob connects to the data lake for his "Product Hierarchy" data, preps and produces an AnswerSet that is shared with Susan who pulls in "Sales Transaction history" data from a Cloud application.

Susan preps this data and produces an AnswerSet, which she shares with Varya for the Sales Variance project that she maintains. In addition to the AnswerSet from Susan, Varya also combines data from an Excel report that she pulls in from a cloud storage system.

When Varya is finished with her data prep, she produces a "Sales Variance Report" AnswerSet. She needs to produce this report each week. She clicks **Create Project Flow** in her Sales Variance project and configures a time-based trigger for running the Flow. APF traverses back through the Flow of related projects, AnswerSets, and datasets to create the dependency chain required to produce the end-state AnswerSet. Varya then uses the APF Monitoring Interface to manage all subsequent runs of the Flow.

APF requirements

- Contributors must have permissions to all of the datasets and all of the projects in the Flow before creating a Flow, otherwise it will not run successfully.

Note

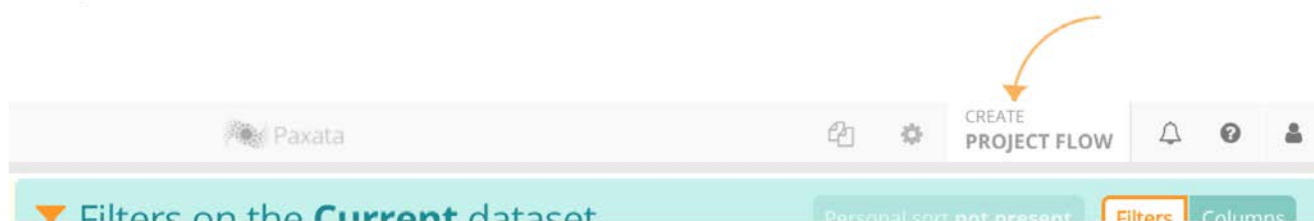
If a contributor has permissions to an AnswerSet, but not to the project from which that AnswerSet was produced, they can still create a Flow up to the point at which they ceases to have the read permission. This flexibility in Flow creation enables contributors to manage the operationalization of Flows for the portions they have permissions to access independently.

- Contributors must also have permissions to all datasets and projects in the Flow to manage them from the Monitoring Interface. Data Prep System Administrators provide these permissions.
- The target project does not include anything produced downstream in the defined Flow. In the previous example, if the "Sales Variance Report" AnswerSet consumes a project, the project is not included in the Flow—the target project is always the end point for a Flow.

Set up a Project Flow

To create a Project Flow:

- Open your target project, the project that will produce your end-state AnswerSet.
- Click **Create Project Flow** on the top right of the **Projects** page.



- Provide a name and optional description for the Flow, then click **Create**.

The screenshot shows a configuration form for a Project Flow. At the top right, there are two buttons: 'Cancel' (grey) and 'Create' (green). Below the buttons, the 'BASICS' tab is selected. The form has two main input fields: 'Project Flow Name' and 'Description'. The 'Project Flow Name' field contains the text 'My Project Flow' and has a small question mark icon to its right. The 'Description' field is currently empty and also has a question mark icon to its right.

The intelligent automation engine calculates the Flow dependencies and APF displays the **Project Flows** page where you configure APF. You can also access the **Project Flows** page when you edit an existing Project Flow.

See [Manage Flows](#) for common actions you can take for all Flows.

Configure APF

You configure APF by setting triggers and notifications on the **Project Flows** page. You can also adjust settings for the Flow's input and output datasets.

The **Project Flows** page has three tabs where you configure Flow settings:

- [General tab](#)
- [Inputs tab](#)
- [Outputs tab](#)

General tab

Use the **General** tab to update Project Flow details and to add triggers.

Project Flow Name

General Inputs Outputs BETA Graph Actions Discard Changes

Project Flow Name Description

WHEN TO RUN

Refresh **minutes** hourly daily weekly monthly yearly custom...

Every 1 Minutes

NOTIFICATIONS

If the Project flow has **Errors** then email

If the Project flow is **Success** then email

On the **General** tab, you can:

- Update the **Name** and **Description** of a Flow that you've created.
- Specify the triggers to run your Flow. The triggers are time- and frequency-based. You can also use the **custom** option to provide a cron expression for the trigger.
- Provide email addresses for run status. Separate each address with a comma.

Note

As soon as a Flow is created, a **Project ID Flow** displays on the **General** tab. This ID is used to identify the Flow for REST API calls and for troubleshooting the Flow.

Inputs tab

The **Inputs** tab provides a list of the datasets used in the Flow, the versions of those datasets used to create the Flow, and the projects in which each dataset is used.

Project Flow Name

STARTED COMPLETED DURATION RAN BY

General **Inputs** Outputs

DATASETS OPTIONS FOR DATASETS AS USED IN PROJECTS PROJECTS

Dataset 1
 Reimport dataset on run
 Pin to version of dataset used in Project **CONFIRM**
 Use latest version of dataset in Library **CANCEL**
 Fail if columns changed

Dataset 2
 Reimport dataset on run
 CONFIGURE REIMPORT OPTIONS **EDIT**
 Using latest version of dataset in Library

Dataset 3
 Reimport dataset on run
 Using latest version of dataset in Library **EDIT**

Finance ?
 Version: 5
 Created: Jul 19th 2019 9:52 AM
 By Anna A.
 VIEW PROJECT

Finance B ?
 2 PROJECTS
 SEE ALL PROJECTS

Annotations:

- Select version of dataset to use (points to Finance project)
- Automatically reimport the dataset as part of the Flow's run (points to Dataset 2 reimport checkbox)
- Click to drill into the list of Projects and optionally configure different versions of the dataset to use per Project (points to SEE ALL PROJECTS button)

On the **Inputs** tab, you can:

- Specify that a dataset is automatically reimported every time the Flow is run.

By default, all projects are configured to use the latest version of a dataset saved in the library. However, newer versions of a dataset may be available from the original data source before a new version of it is manually imported to the Data Prep library. In this case, you can configure a dataset to be automatically reimported from its original data source every time the Flow is run. Then this latest version will always be saved in the library. To enable this automatic update, click **Reimport dataset on run**. When the option is enabled, a **Configure Reimport Options** button also displays. The button opens the library import pane where you can change the data source path, query, or enter export parsing options. These options are saved with the dataset in the library and you only need to configure them if you want to change the current settings.

- Configure a dataset's version to use for the project.

By default, all projects are configured to use the latest versions of datasets saved in the library. However, you may want to change this default behavior, which can be done when you click **Edit** (in the **Options for Datasets as used in Projects** column):

- **Pin to version:** The dataset remains the exact version currently used by the project.
- **Fail if columns changed:** The dataset will fail to import into the project if the latest version coming in from the library has a different layout (schema), for example, if new columns are added, columns that are not used in the project's steps are removed, a column type is changed, the order of columns is changed, etc.
- If more than one project uses the same dataset as input for the Flow, this is noted in the projects column.

Click **See All Projects** to view all projects that use the dataset and to optionally configure different versions of the dataset to use per project. For example, you can specify that one project uses the latest version of the dataset from the library, while another project uses the exact version of the dataset currently saved in the associated version of the project.

Note

View metadata statistics for the dataset inputs by hovering over a dataset name in the **Datasets** column. The dataset's version, creation date, the user who added it to the library, and the number of columns and rows are displayed in a pop-up window.

Outputs tab

The **Outputs** tab provides a list of all the output AnswerSets that are published from the Flow.

Project Flow Name

General Inputs **Outputs**

PROJECT OUTPUT

Project Name Project Version # Publishing is disabled CONFIGURE LENS

Project Name Project Version # Publish to Library CONFIGURE LENS

Project Name Project Version # Publish to Library CONFIGURE LENS

Export Lens Project Library Library and Export Export Settings Cancel Confirm

Export location options

If exporting to external data source, export configuration is provided in this panel

Option to enable/disable publish of non-essential AnswerSet to Library

Click to open the publish options for the AnswerSet

All outputs are configured at the lens level because a publishing lens is always required to create a publishing point from a Data Prep project.

There are times when your Flow may include a project that has multiple lenses. Not all of those lenses are required to produce output AnswerSets. By default, only required lenses automatically publish AnswerSets that are saved in the library. If you'd like to enable the AnswerSets to be published even if they are not required for the Flow, you can enable them on the **Outputs** tab.

Note

Lenses that produce output AnswerSets and are required for the Flow can never be disabled.

In addition to adjusting the publishing options for non-essential AnswerSets, you can publish any lens output AnswerSet to an external data source, for example, a database or a cloud storage system. To specify a publish location in addition to the Data Prep library, click **Configure Lens** and open the **Exports** pane.

You can take the following actions on the **Outputs** tab:

- Disable a non-bridging lens to prevent it from publishing AnswerSets to the library.

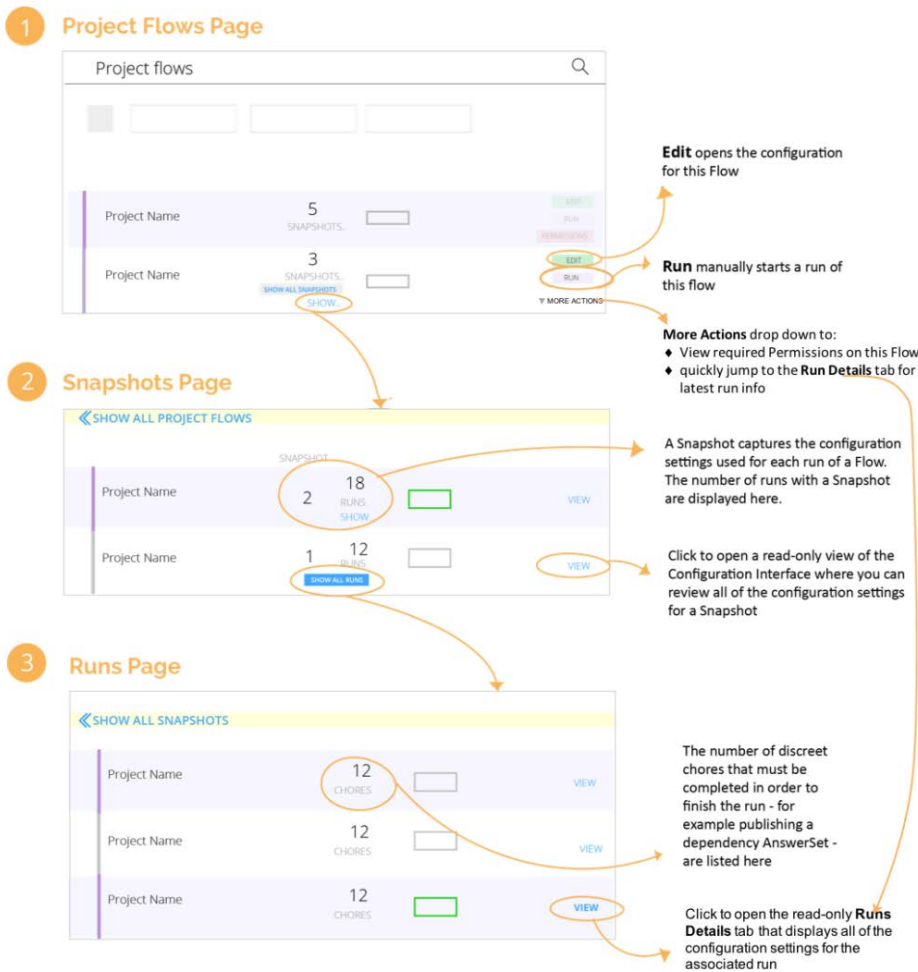
Click the slider adjacent to the lens to disable it.

- Export the published AnswerSet to a data source (in addition to the default library setting).

Click **Configure Lens** for the lens. The **Export** pane opens at the bottom of the page. By default, Data Prep publishes AnswerSets to the Data Prep library. To publish to an external data source, click the dropdown menu for the **Export Lens** field and select **Library and Export**. You can then specify the output location details and export parsing options for that AnswerSet.

Monitor Flows

APF lets you monitor the status of Flows. The key components for generating a Flow's output are Snapshots, Runs, and Chores. The following diagram illustrates how these components monitor Flows. See the following sections for details.



Project Flows page

The **Project Flows** page lists the Flows that you have permissions to view and edit, as well as the current status of the most recent run for each. On this page, you can:

- Edit the configuration details for the Flow. Click **Edit** to open the APF Configuration Interface where you can make adjustments to the configuration. See [Configure APF](#).
- Click **Run** to run the Flow manually. Starting a Flow manually is particularly useful if you need to test out a new Flow or a configuration change to the Flow and you don't want to wait for the time-based trigger to start it.
- Show the [Snapshots](#) for the Flow. Click **Show all Snapshots** to open the Snapshots pane.
- Click **More Actions > permissions** to update the permissions settings so that you can share this Flow with another person. Note that permissions are only visible to the user who created the Flow or to users with whom the creator has shared all of the permissions.
- Click **More Actions > view latest results** to go to the latest Flow. This will not display until there is at least one run of the Flow.

Snapshots page

The **Snapshots** page lists the Snapshots for a Flow. Every time a Flow is executed (called a "run" of the Flow), a Snapshot is created to capture the configuration settings used to create the output for the run. The runs continue with this Snapshot until any configuration changes are made to the Flow—for example, changes to the schedule, notifications, inputs, and output settings. Then a new Snapshot is created for the Flow. The new Snapshot captures the executed runs with the modified configuration settings. Snapshots allow you to audit the exact state of a Project Flow for each run.

Note

APF does not create a new Snapshot if datasets are configured to use the latest version from the library. See the [Inputs tab](#) for dataset configuration options.

On this page, you can:

- Click **View** to open a read-only view of the APF configuration settings for the Snapshot.
- Click **Show All Runs** to open the **Run List** page, which details every run for the Snapshot.

Runs page

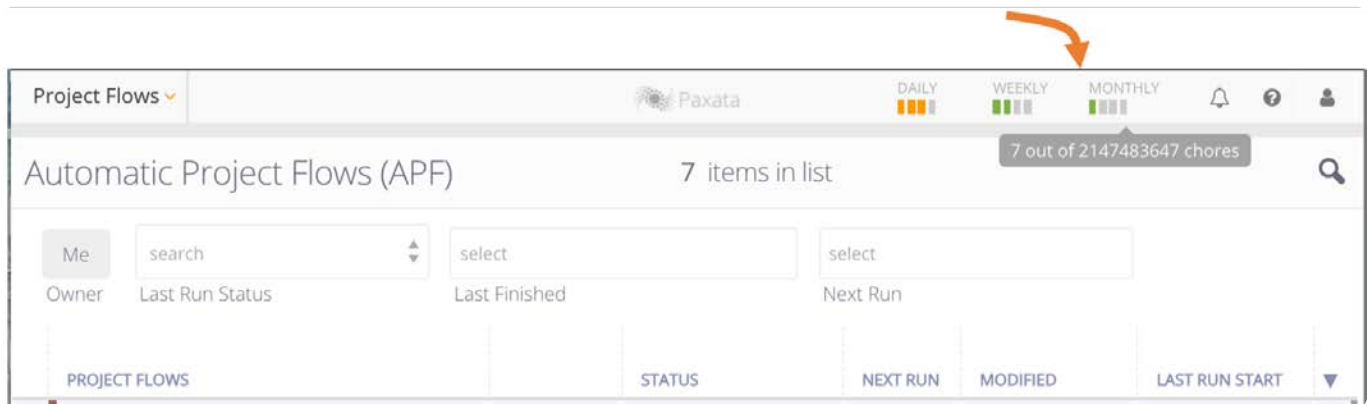
The **Run List** page captures all details for each individual run under a Snapshot. The number of discreet chores that must be completed to finish the run—for example, publishing a dependency AnswerSet—are listed on the page. Every time a Flow is run, a new run entry displays on this page.

To open a read-only view of the APF configuration settings associated with a run, click **View**.

Note

If there is no change to the data used to create the Flow—for example, all datasets used in the Flow remain exactly the same version as were used in the previous run—the APF engine conserves resources and does not rerun the Flow again until new data inputs are available.

The APF quotas meter displays at the top of the Flows page to indicate your usage. Hover over one of the counts—Daily, Weekly, or Monthly. A tooltip provides details of your current usage and limits.



The screenshot shows the 'Project Flows' page in the Paxata interface. At the top right, there are three quota meters for 'DAILY', 'WEEKLY', and 'MONTHLY'. An orange arrow points to the 'MONTHLY' meter, which has a tooltip showing '7 out of 2147483647 chores'. Below the quotas, the page title is 'Automatic Project Flows (APF)' with '7 items in list'. There are search and filter controls for 'Owner', 'Last Run Status', 'Last Finished', and 'Next Run'. At the bottom, a table header is visible with columns: 'PROJECT FLOWS', 'STATUS', 'NEXT RUN', 'MODIFIED', and 'LAST RUN START'.

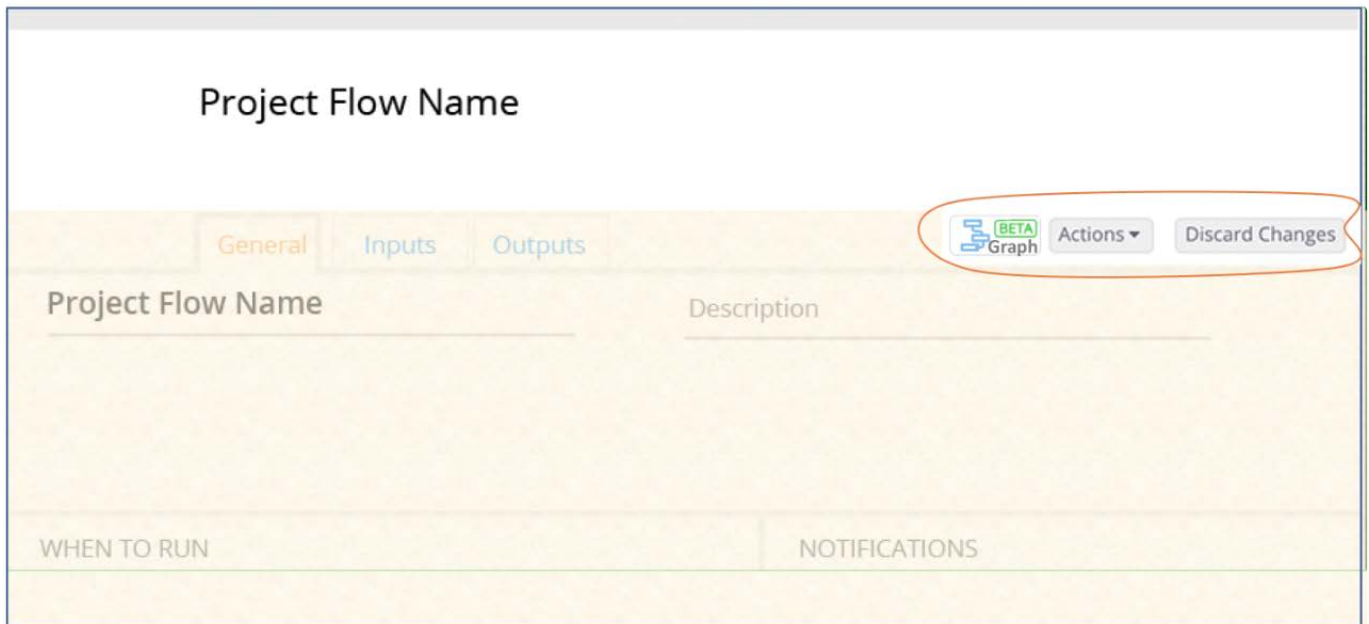
Quotas are based on Chore count, and Chores are defined as:

- The running of an individual project that is required to produce a Flow.
- An import (but not a publish) of any dataset or AnswerSet that is required to produce a Flow.

The sum of all Chores ultimately produces the output for your Flow. While a Flow is running, refresh your browser to update the quotas meter on the Flow's page. If you need your Chore count quotas increased, contact your DataRobot Data Prep Administrator or DataRobot Customer Success.

Manage Flows

Access tools for managing Flows on the top right of the **Project Flows** page.

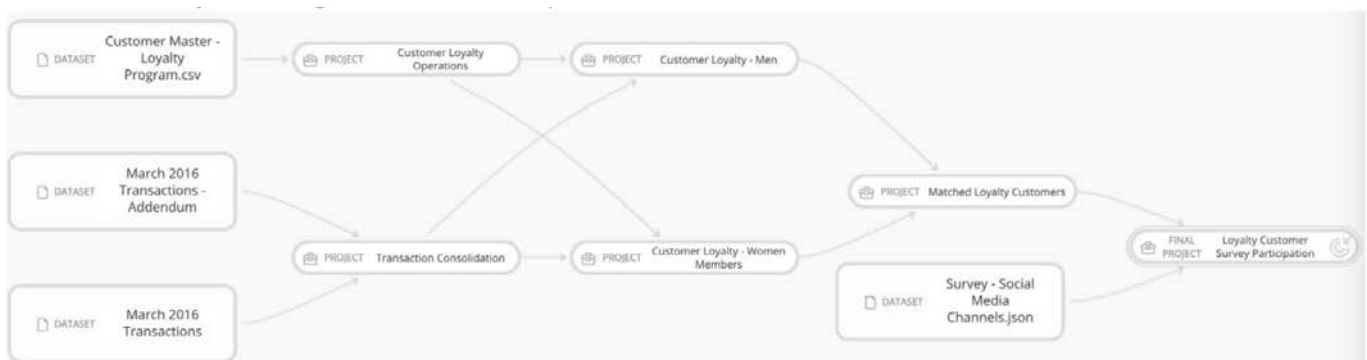


You can manage your saved Flow by:

- [Generating a visual graph for a Flow](#)
- [Running a Flow manually](#)
- [Deleting a Flow](#)
- [Updating Flows to use latest project versions](#)

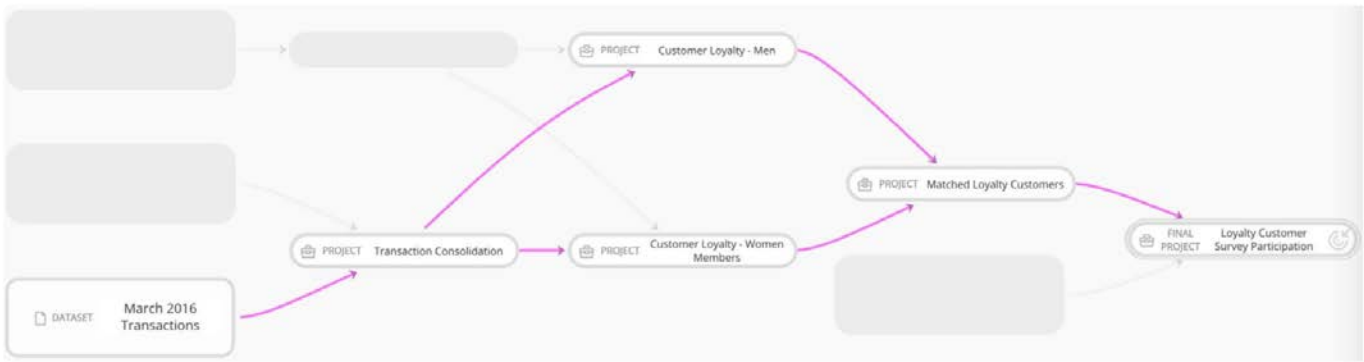
Generate a visual graph for a Flow

The **Graph** button generates an APF graph in a new browser window that displays the datasets and how they flow into the individual projects used to generate a Flow's final output AnswerSet.

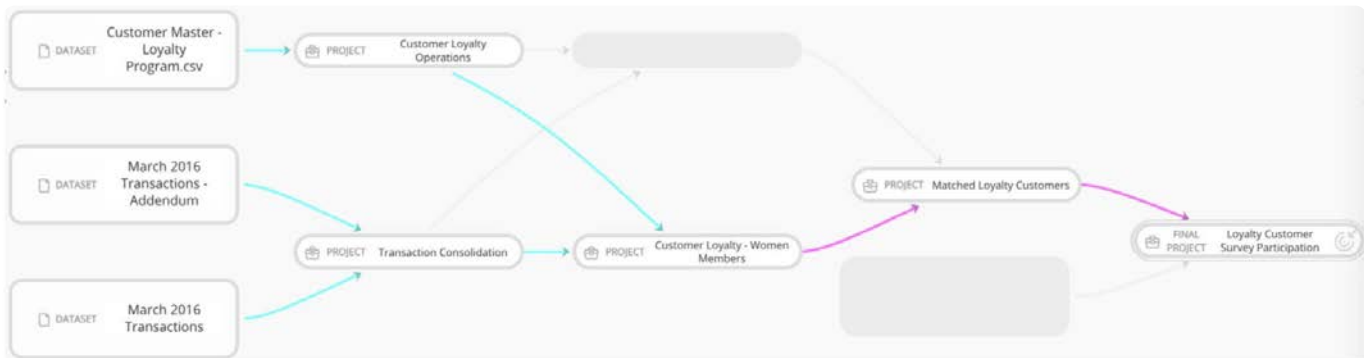


Hover over a dataset or project in the Flow to display the corresponding downstream lineage (in pink) and upstream dependencies (in blue).

For example, hovering over the dataset for **March 2016 Transactions** displays the following:



Hovering over an intermediate project in the Flow—in this example, **Customer Loyalty-Women Members**—the upstream dependencies display in blue and the downstream lineage displays in pink.



Notice in both examples that if datasets and projects do not participate in the portion of the Flow that you've selected, they are grayed out in the graph.

You may see a dotted line in a graph for some Flows. The dotted line indicates that an AnswerSet was published from a project in the Flow, and then later consumed again by the same or another project in the Flow. This is referred to as a looping input and is represented by the dotted line.

Run a Flow manually

There may be times when you want to manually kick off a run of a Flow without having to wait for its scheduled start time. This can be done from the **Actions** dropdown menu. Click **Run now**.

Delete a Flow

If you no longer want to keep a saved Flow, you can delete it. Click **Actions > Delete**. You are prompted to confirm your selection. Note that any AnswerSets that were published to the library as a result of running this Flow will not be deleted as a result of deleting the Flow.

Update a Flow with the latest project versions

Every time an action is taken in your project—for example, adding a step, removing a step, or rearranging steps—a new version of your project is created. Each version provides an audit trail of the changes you have made to your data during the course of your data prep work. When creating a Project Flow, the Flow is always pinned to the specific project versions at the time of the

Flow's creation. However, you can update a Flow to use the latest version of all projects. This can be done from the **Actions** dropdown menu. Select **Update Project Versions** and you are prompted to confirm your selection:

Update All Project Versions

This action updates all Project versions used in the Flow to the latest Project versions. There are **conditions that apply** to updating Project versions and no updates can be made to the Flow if the conditions are not met.

WARNING: Any unsaved changes in these Projects will be lost. If a Project is shared with someone else, ensure all changes have been saved before proceeding here.

Create a new APF instead of overwriting the existing one?

New APF name:

Cancel **Update All Projects**

You can choose to overwrite the existing APF or create a new one. If you choose to create a new APF, all triggers are copied to the new APF but they are disabled by default.

Note

The ability to update an existing APF must be enabled. If you do not see the **Update All Project Versions** window, contact your Data Prep System Administrator to enable this. If this feature is not enabled, a warning message displays and you can only update the versions if there were no significant changes to the project (e.g. no changes to project dataset of the lenses).

To update a specific project's version—instead of all projects in the Flow—on the **Outputs** tab, hover over the project for which you want to update the version, then click **Update Project Version** in the right-hand column.

APF terms

The following are terms specific to APF.

TERM	DEFINITION
Chore	A dataset import or a project execution. The dataset import chore performs a reimport of your dataset through a data source. The project execution chore addresses all other tasks required for the Flow, such as publishing an AnswerSet to the library or exporting an AnswerSet.

TERM	DEFINITION
Flow	A collection of projects that can be run as a unit. One or more frequency-based schedules can be associated with a Flow, which allows a Flow to run on a recurring basis.
Inputs	Datasets from the library that are required to run a Flow.
Outputs	The AnswerSets written to the library generated by the run of a Flow.
Run	The execution of each of the projects that are required by the target project. The run executes all of the steps from the upstream dependency projects, then writes the resulting AnswerSet(s) to the library.
Snapshot	The configuration settings captured for each run of a Flow. Your Data Prep Administrator must enable this feature in your application.
Target Project	The Data Prep project from which a Flow is created. Once a Flow is created, all upstream dependencies are automatically calculated by the APF engine.

Data Prep advanced topics

Data Prep provides features that protect your data, support interactive processing, and generate *ClicktoPrep* links from business intelligence (BI) and data visualization tools to your Data Prep steps and Filtergrams.

These pages describe:

TOPIC	DESCRIBES...
ClicktoPrep links	Generate links to Data Prep steps or Filtergrams that you can include in your business intelligence (BI) and data visualization tools.
Infrastructure and security	Learn about Data Prep's infrastructure and security features.
Interactive mode	Enable Interactive Mode which lets you perform data prep on a portion of you data so that you don't have to import large amounts of data into Data Prep.

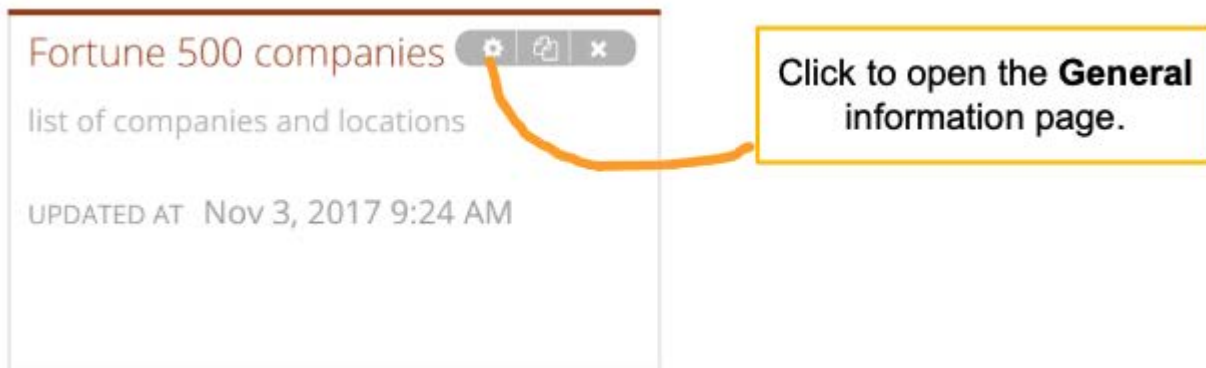
ClicktoPrep links

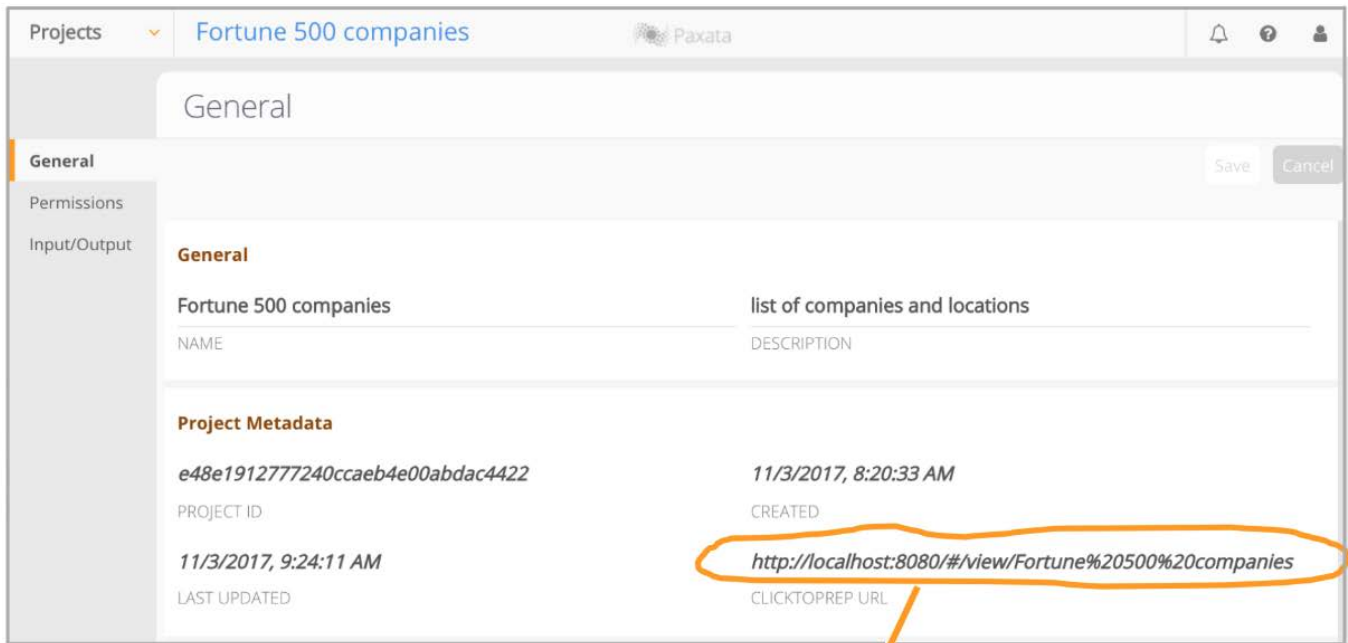
As a data preparation tool, Data Prep is frequently used in conjunction with business intelligence (BI) and data visualization tools, for example Tableau®, Power BI®, Qlik® and MicroStrategy®. For tools that support hyperlinks or URLs, you can generate a Data Prep ClicktoPrep link that takes you from your tools directly to any of the following:

- The last step in your Data Prep project
- A Data Prep project Filtergram
- A specific Data Prep project step that generated or modified the data

You can then open the link from your BI or visualization tool, make modifications to the Data Prep data, republish, and refresh the visualization or report to reflect your updated data.

To begin, you need the *base URL* for your project, which takes you to the last step in your Data Prep project. Copy your project's base URL from the project's **General** page.





Projects **Fortune 500 companies** Paxata

General

General

Permissions

Input/Output

General

Fortune 500 companies list of companies and locations

NAME DESCRIPTION

Project Metadata

e48e1912777240ccaeb4e00abdac4422 11/3/2017, 8:20:33 AM

PROJECT ID CREATED

11/3/2017, 9:24:11 AM http://localhost:8080/#/view/Fortune%20500%20companies

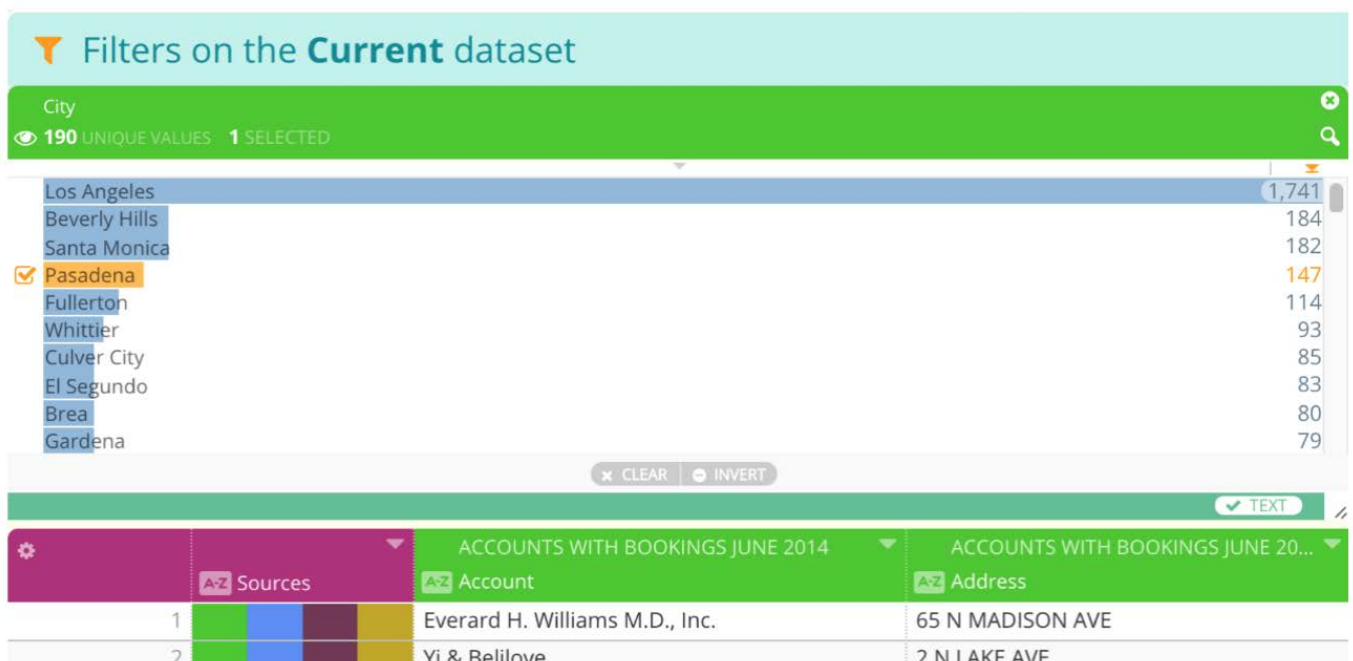
LAST UPDATED CLICKTOPREP URL

Drag your mouse over this URL and right-click to select copy.

Use this base URL to point from your visualization tool to the last step in your project. Or, continue to append the parameters described in the following sections to point to specific Filtergrams or steps in your Data Prep project.

ClicktoPrep links to project Filtergrams

This type of link opens a Data Prep project Filtergram for the dataset on the column and value you specify. For example, the following ClicktoPrep link navigates to the **Customers** project and opens a Filtergram on the column **City** filtered to **Pasadena**.



Filters on the **Current** dataset

City

190 UNIQUE VALUES 1 SELECTED

City	Count
Los Angeles	1,741
Beverly Hills	184
Santa Monica	182
<input checked="" type="checkbox"/> Pasadena	147
Fullerton	114
Whittier	93
Culver City	85
El Segundo	83
Brea	80
Gardena	79

CLEAR INVERT

TEXT

Sources	Account	Address
1	Everard H. Williams M.D., Inc.	65 N MADISON AVE
2	Yi & Belilove	2 N LAKE AVE

To create this type of link, see [Configure links to project Filtergrams](#).

ClicktoPrep links to project steps

This type of link opens a Data Prep project in edit mode, set to one of the following:

- The last step that edits or modifies a column.
- The last [StepType] that edits or modifies a column, for example a “Find and Replace” step on the column.

The following [StepTypes] are supported by ClicktoPrep.

StepType	Syntax
Import	AnchorTableStep
Append	AppendStep
Find and Replace	BulkEditStep
Cluster and Edit	ClusterEditStep
Duplicate Column	DuplicateColumnStep
Computed Column	ExpressionStep
Manage Columns—includes all of the following operations: Hide, Reorder, Delete columns (**see note below)	EditColumnsStep
Pivot—includes all of the following shaping operations: Deduplicate, Depivot, Transpose, Group and Pivot (**see note below)	PivotStep
Transform—includes all of the following transformations: capital case, lowercase, upper case, unescaped HTML, blanks, custom value, rename, whitespace for both numeric and string values (**see note below)	TransformStep

Note

If you have more than one of these [StepTypes] in your project, the URL points to the last one that occurs in your **Steps** pane.

Example 1

The following ClicktoPrep link navigates to the last step in your **Customer** project that affected or modified data in the **Division** column.

Steps: Save Cancel

Find + Replace Division using Equals StartsWith Contains Filters

find SOUTH replace Match with SOUTHEAST

OPTIONS
 Ignore case

StateZone	Division	Division	PromoZone
SOUTH	→ SOUTHEAST		ALABAMA - ALL OTHERS
SOUTH	→ SOUTHEAST		ALABAMA - ALL OTHERS
SOUTH	→ SOUTHEAST		ALABAMA - ALL OTHERS
SOUTH	→ SOUTHEAST		ALABAMA - ALL OTHERS
SOUTH	→ SOUTHEAST		ALABAMA - ALL OTHERS

24 COLUMNS + 670 ROWS AT STEP 6

Notice you are taken to the last Step edit that occurred for this column.

Example 2

The following ClicktoPrep link navigates to the **Customer** project where the last "Find + Replace" [**BulkEditStep**] step was made on the **Company** column.

Steps: Save Cancel

Find + Replace Company using Equals StartsWith Contains Filters

find RunWAY CAPITAL MANAGEM. replace Match with Runway Capital

OPTIONS
 Ignore case

Company	Company	Account
Curtiss Erickson	→ Curtiss Erickson	Curtiss Erickson
The Master Insurance Agency	→ The Master Insurance Agency	The Master Insurance
Jane Griffin	→ Jane Griffin	Jane Griffin
Discus Deb	→ Discus Deb	Discus Deb
RunWAY CAPITAL MANAGEMENT	→ Runway Capital	Runway Capital MA
RunWAY CAPITAL MANAGEMENT	→ Runway Capital	RunWAY CAPITAL MA
RunWAY CAPITAL MANAGEMENT	→ Runway Capital	RunWAY CAPITAL MA

44 COLUMNS + 4,576 ROWS AT STEP 8

Notice you are taken to the last Step in your Project where a Find + Replace occurred in the column.

To create this type of link, see [Configure links to project steps](#).

Requirements and considerations

Before creating Data Prep ClicktoPrep links, keep the following requirements and considerations in mind:

- Your visualization tool must support URLs.

- You must have the required Data Prep Resource-Level permissions to open any URL that directs you into a Data Prep project.

Configure links to project Filtergrams

The format and configuration parameters for ClicktoPrep project Filtergram links are described in this section. (For Tableau users, see [Create links to project Filtergrams in Tableau](#).)

The format is as follows.

```
https://[server]/#/view/[projectname]?filtercolumn=[column]&filtervalue=[value]
```

Tip

The characters before the question mark form the base URL. Copy the base URL from the **ClicktoPrep URL** field on the **General** page.

The following are requirements for project Filtergram links:

- Everything after the # in the URL is case sensitive.
- In order to create a link for a Data Prep project Filtergram, your visualization tool must support dynamic URLs.
- If your project name or a column name includes one or more spaces and you haven't copied it from the **General** page, replace (resolve) each space in the URL with %20 (the HTML encoding value for a space.) For example: `https://<server>/#/edit/Web%20Campaigns%20demo/Phone%20Number`

Note

Some visualization tools provide an option to automatically encode the spaces for you. Refer to [HTML URL Encoding Reference](#) for other special characters that require encoding.

Configure links to project steps

The format and configuration parameters for ClicktoPrep project step links are described in this section. (For Tableau users, see [Create links to Data Prep project steps in Tableau](#).)

For the last step that edits or modifies a column, the format is as follows:

```
https://[server]/#/edit/[projectname]/[columnname]
```

For the last [StepType] that edits or modifies a column, for example a "Find and Replace" step on the column, the format is as follows.

```
https://[server]/#/edit/[projectname]/[columnname]?filter=[StepType]
```

Tip

The characters before the question mark form the base URL. Copy the base URL from the **ClicktoPrep URL** field on the **General** page.

These [StepTypes] are supported:

- AnchorTableStep
- AppendStep
- BulkEditStep
- ClusterEditStep
- DuplicateColumnStep
- EditColumnsStep
- ExpressionStep
- PivotStep
- TransformStep

See [ClicktoPrep links to project steps](#) for an explanation of the [StepTypes].

Following are requirements for project step links:

- Everything after the # in the URL is case sensitive.
- If your project name or a column name includes one or more spaces and you haven't copied it from the **General** page as described above, you must replace (resolve) each space in the URL with %20 (the HTML encoding value for a space.) For example: `https://<server>/#/edit/Web%20Campaigns%20demo/Phone%20Number`

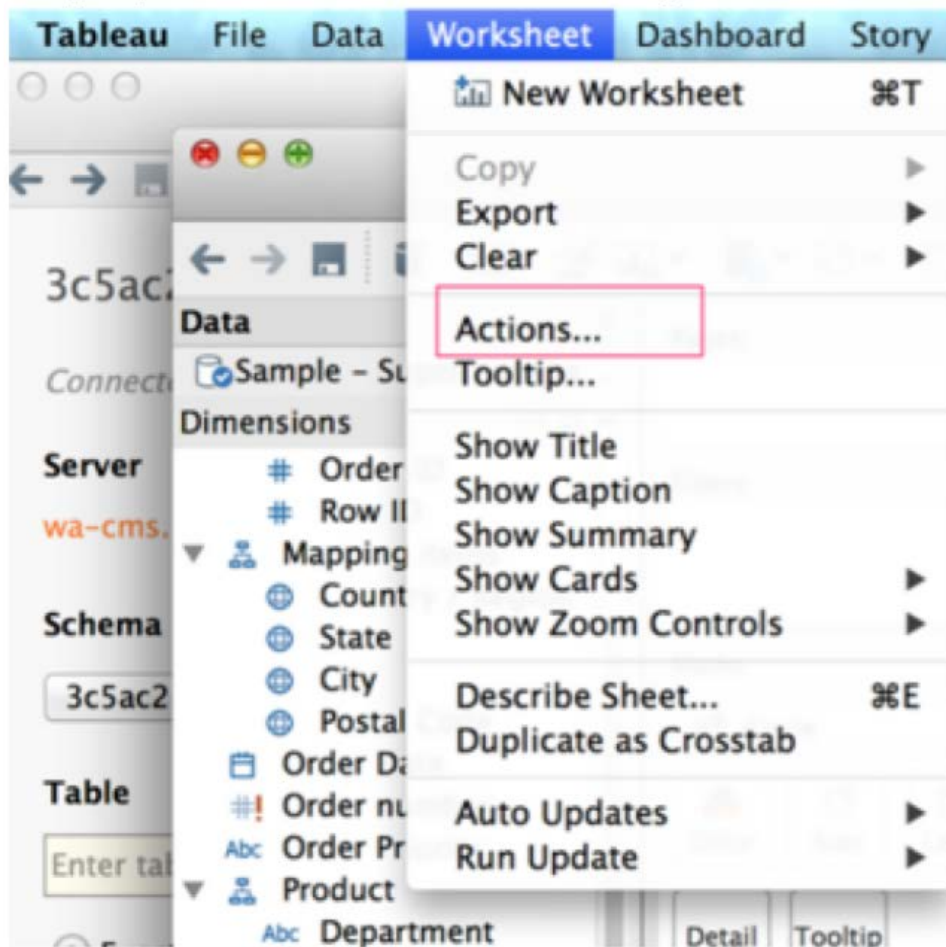
Note

Some visualization tools provide an option to automatically encode the spaces for you. Refer to [HTML URL Encoding Reference](#) for other special characters that require encoding.

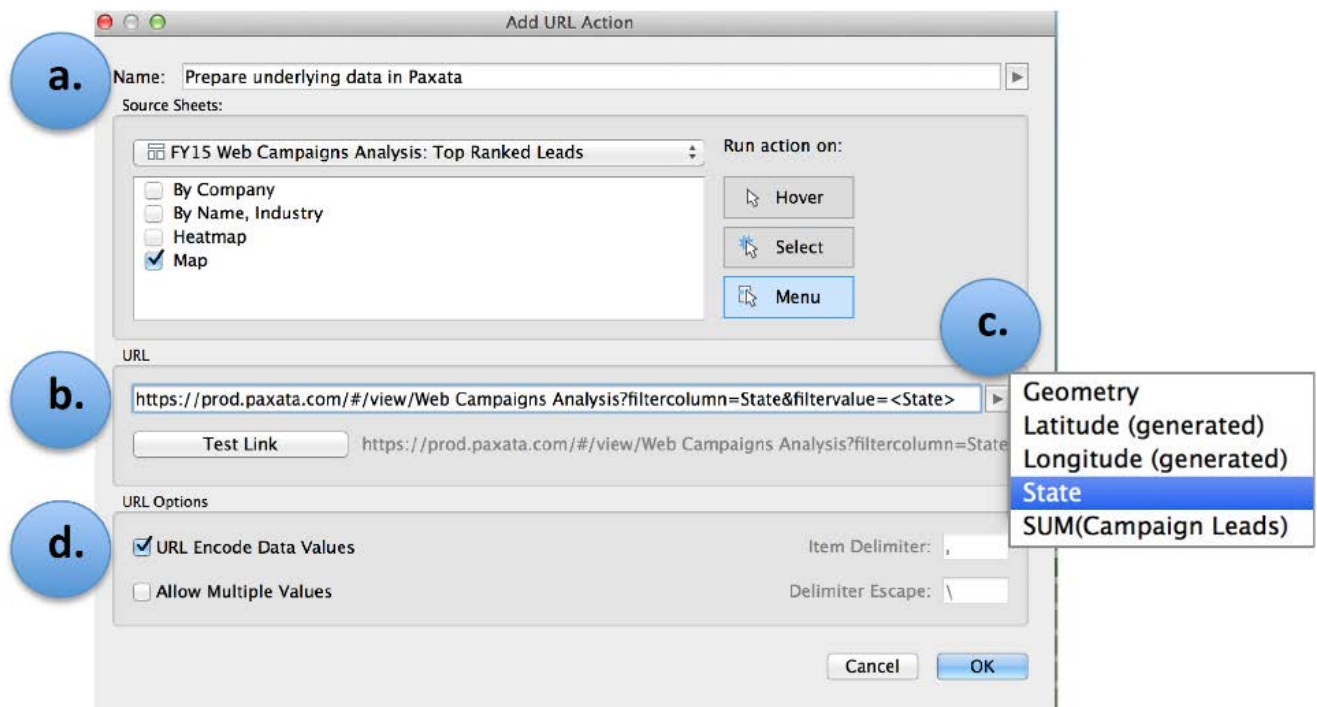
Create links to project Filtergrams in Tableau

For tools that support dynamic hyperlinks or URLs, you can link directly to your Data Prep project with a filter applied based on your selection tool. This section explains how to do this in Tableau.

1. Open your visualization in Tableau and go to: **Worksheet > Actions...**



2. Select **Add URL Action** and provide the following configuration information:



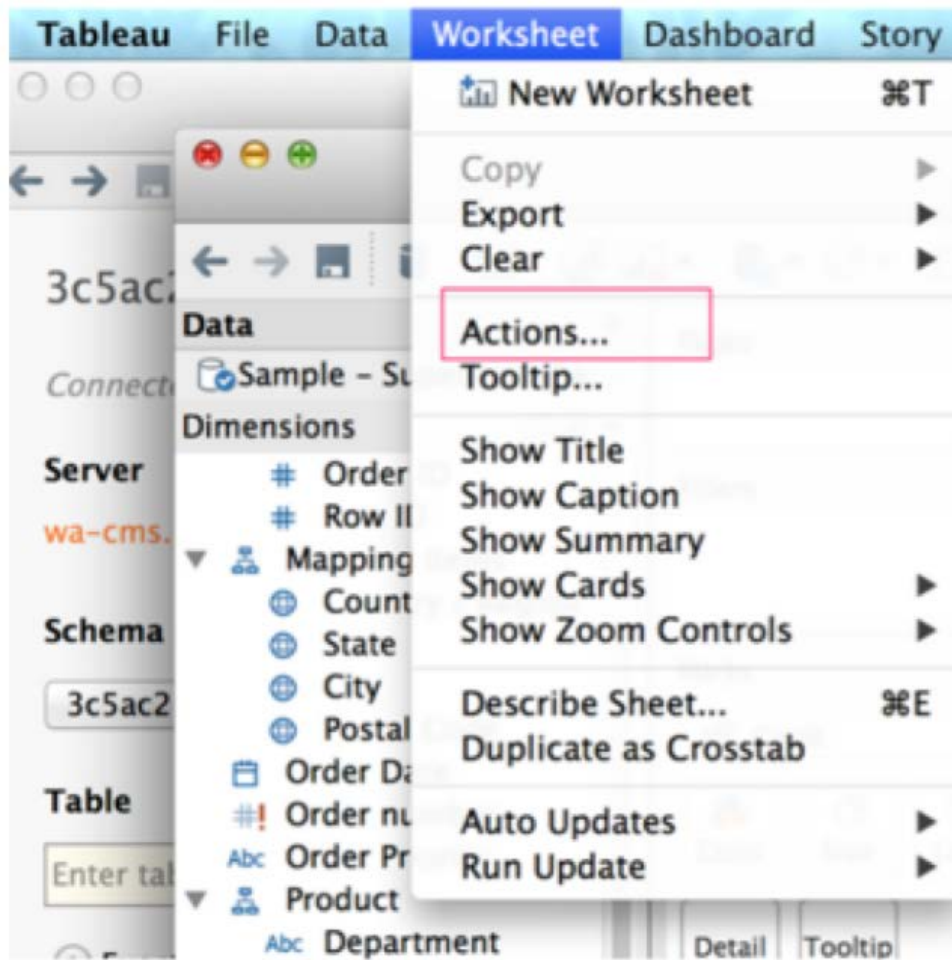
a. Name for the link displayed in Tableau.

- b. The URL is in this format: `https://[server]/#/view/[projectname]?filtercolumn=[column]&filtervalue=[value]` The characters before the question mark form the base URL. Copy the base URL from the **ClicktoPrep URL** field on the **General** page. Everything after the # in the URL is case sensitive.
 - c. Click to select the Tableau column that will dynamically receive data from your URL.
 - d. Ensure **URL Encode Data Values** is checked.
3. Click **Test Link** to verify the link works.
 4. Click **OK** to save.

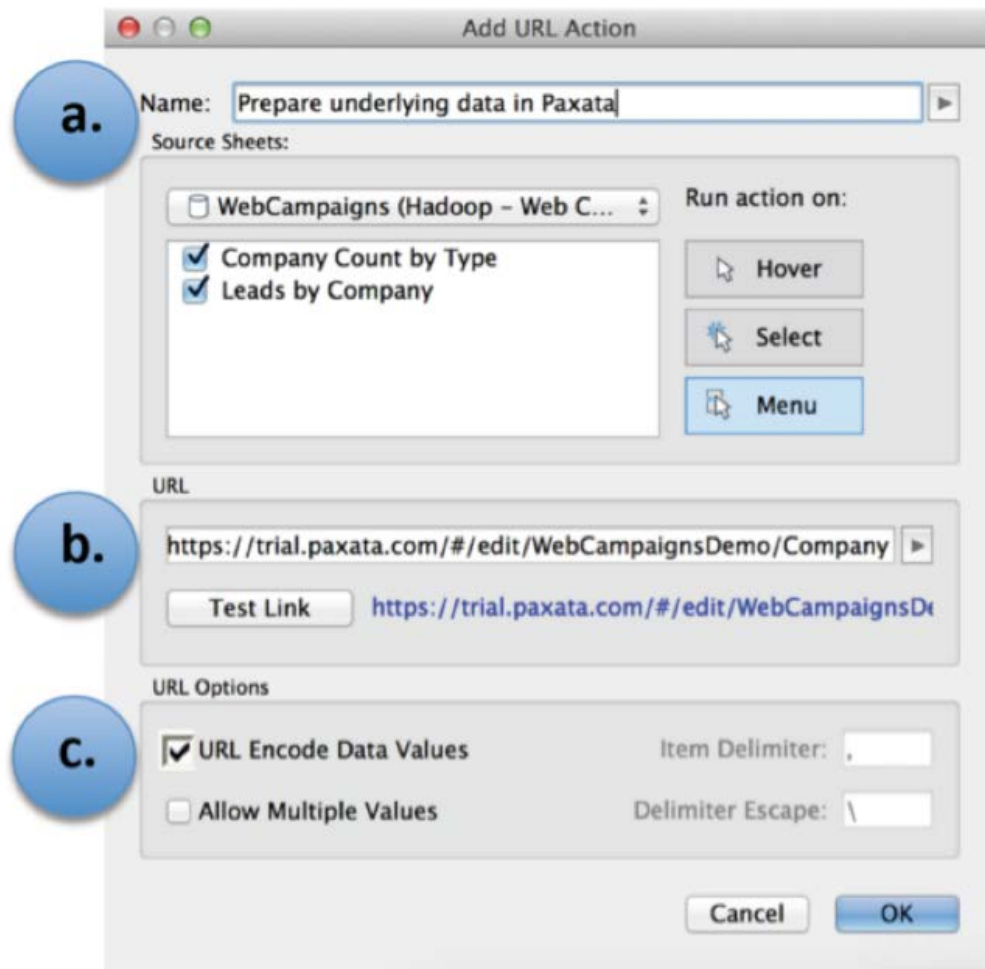
Create links to Data Prep project steps in Tableau

For tools that support dynamic hyperlinks or URLs, you can link directly to your Data Prep project with a filter applied based on your selection tool. This section explains how to do this in Tableau.

1. Open your visualization in Tableau and go to: **Worksheet > Actions...**



2. Select **Add URL Action** and provide the following configuration information:



a. Enter the name for the link displayed in Tableau.

b. Specify the URL using one of these formats:

- For the last step that edits or modifies a column: `https://[server]/#/edit/[projectname]/[columnname]`
- For the last [StepType] that edits or modifies a column: `https://[server]/#/edit/[projectname]/[columnname]?filter=[StepType]`

The characters before the question mark form the base URL. Copy the base URL from the **ClicktoPrep URL** field on the **General** page. Everything after the # in the URL is case sensitive.

See [ClicktoPrep links to project steps](#) for an explanation of the [StepTypes].

c. Click to select Tableau column that will dynamically receive data from your URL.

d. Ensure **URL Encode Data Values** is checked.

3. Click **Test Link** to verify the link works.

4. Click **OK** to save.

Data Prep infrastructure and security

This topic provides details on DataRobot Data Prep [infrastructure security](#) and [application security](#).

Infrastructure security

Access Controls

All ingress ports, whether on internal or external interfaces, are protected by security groups, which are automatically configured by Data Prep's system configuration tool. Customer/public facing ingress ports are: TCP 80/443.

Data Prep utilizes jump hosts for SSH access to production infrastructure and all production admins are access-controlled using multi-factor authentication.

The production accounts use strict IAM roles and only key employees with a verified business need receive administrative access.

Assessments

DataRobot does not allow customer-requested security scanning agents to be installed in our production SaaS environment. Data Prep leverages an on-demand cloud computing platform to perform vulnerability scans against the environment. Penetration testing of Data Prep is executed by a qualified third-party assessor and the results are integrated into the development workflow based on priority. Upon request, DataRobot can schedule vulnerability scans of our SaaS offering and coordinate the request with our cloud computing platform service.

Data Protection

Data Prep utilizes TLS and HTTPS to encrypt the data when in transit. Data Prep stores the data in an encrypted format when it is at rest to prevent access by unauthorized parties.

Application security

Password Policies

Native Data Prep accounts (defined as accounts that are not using LDAP or SAML) adhere to the following password requirements: the password must contain at least one number, one lowercase letter, one uppercase letter and one special character (!@#\$%^&*+=), and at least 8 characters.

Data Prep does not enforce account lockout policies or have any account lockout policy management capabilities for Native accounts.

For SAML authentication, the account policies and password requirements configured with the customer's SAML Identity Provider are enforced.

Production service accounts cannot be used for logins by any admin or user. The account is strictly used only to startup and run the Data Prep application. The account does not have any access to customer data or permissions within Data Prep.

Security Updates

Operating System security patches are applied to our Production SaaS environment after a security threat assessment/review. Careful testing is performed prior to applying any security updates so as to not compromise the integrity of our application or services. Application security updates to our SaaS offering are applied as soon as a fix is available.

Interactive mode

As your source datasets grow larger, they can impact your ability to efficiently import and work with your datasets in data preparation tools. To address this growing problem, Data Prep offers an Interactive Mode feature, which lets you work faster on a portion of your data—a portion size that you decide is right for your project needs. You can then efficiently and interactively prep that portion in a Data Prep project, without ever having to bring all of that data into the project.

Note

Your Data Prep Administrator must enable this feature in your application.

The major advantages of the Interactive Mode feature include the following:

- You don't need to wait for the entire dataset to load into your library before you can begin working with it in a Data Prep project. Instead, you define a portion size for datasets, and when that portion size is reached, the data is available for prep in a project while the remainder of the dataset continues loading in the library.
- When you've finished prepping your data in the project, you can easily apply the transformations to **all of the data** in the native datasets through the [Automatic Project Flows](#) feature.
- You can always reset the dataset portion that you want to work with in Interactive Mode. For example, after working in a project with a portion limit of 50k rows per dataset, you may realize you actually need larger portions from each dataset. Changing the portion size is a one-step operation for your Data Prep Administrator. Your project then dynamically recognizes your portion limits have changed and provides you with the option to refresh your datasets in order to pick up the new data.
- Your interactive experience in Data Prep projects is optimized because you only need to work with the defined portions of your datasets in your project.
- You have more flexibility in how you work with large projects. Data Prep projects in Interactive Mode have a row limit that defines the maximum number of rows that can be prepared within a project. This limit is set by your Data Prep Administrator and is useful because it allows the Administrator to ensure you have the optimal interactive experience based on available system resources.

Total number of rows you can work with interactively in the Project.

The screenshot shows the Data Prep interface. On the left, a 'Steps' panel lists: 'Append with medicare_in.xml', 'Edit Columns', 'Lookup from Hospitals by region.csv', and 'Start with Hospital General Information...'. Below the steps, it displays '27 COLUMNS + 9,861 ROWS AT STEP 4' and a 'REFRESH DATASETS' button. On the right, a table titled 'Filters on the Current dataset' shows columns: Provider ID, Hospital Name in region, Address, City, and State. The table contains 11 rows of hospital data. A yellow box highlights the 'REFRESH DATASETS' button and the table.

Number of rows per dataset you can use interactively. The Refresh Datasets option lets you choose which datasets to update with the latest versions of data.

By default, Interactive Mode is not enabled for your Data Prep projects and you will need to contact your Data Prep Administrator to enable it. Before enabling Interactive Mode, you should consider the following points:

- For existing projects, use the [Profiling](#) feature for the datasets in those projects. Profiling the datasets will give you fuller insights into the data and will inform your choice regarding the optimal portion size to select for your datasets.
- Existing projects with datasets whose row sizes now exceed the defined portion size will *not* be dynamically updated to remove any rows. Instead, when you open those projects, you will have the option to use the [Refresh Datasets](#) feature to enforce the row portion for each dataset. The row portions will only be applied if you elect to refresh those datasets.

After Interactive Mode is enabled, the following icon displays to indicate that you're operating on a portion of the dataset:

The screenshot shows the 'Steps' panel with three steps: 'Append with MedicareUpdated.xlsx', 'Lookup from medicare_in.csv', and 'Start with Hospital General Information.csv'. A yellow box labeled 'Portion icon' points to a small icon on the 'Lookup' step, which indicates that the data is being viewed on a portion of the dataset.

When you hover over the icon, the row portion value displays so that you can quickly discern the value enforced for all datasets.

Tip

To determine the *total number of rows* in a dataset, go to the library page where that total is displayed for each dataset. Additionally, the library page provides information specific to the Interactive Mode feature so that you can determine:

- The loading status of a dataset and when its interactive portion is available for use in a project.
- The AnswerSets that have been published from projects in Interactive Mode.

See [Data Prep library](#) for details.