# Analysis of imputation bias
# for feature selection with missing data

Borja Seijo-Pardo[1], Amparo Alonso-Betanzos[1], Kristin Bennett[2],
Verónica Bolón-Canedo[1], Isabelle Guyon[3], Julie Josse[4], Mehreen Saeed[5]

1 . U. A Coruña, Spain. 3. RPI, New York, USA. 3. U. Paris-Saclay, France.
4. Ecole Polytechnique, France. 5. FAST, Lahore, Pakistan.

**Abstract**.   We study risk/benefit tradeoff of missing value imputation in the context of feature selection. We caution against using imputation methods that may yield false positives: features not associated to the target becoming dependent as a result of imputation. We also investigate situations in which imputing missing values may be beneficial to reduce false negatives. We use causal graphs to characterize when structural bias arises and introduce a de-biased version of the t-test.

## 1   Introduction

Sample bias is a well-known form of bias that plagues datasets large and small. In machine learning, this refers to obtaining (training) data drawn according to a distribution deviating from the "natural" distribution of the problem at hand. A variant of this problem is to have datasets with missing data. In this case "sample bias" corresponds entirely missing samples. We consider in this paper cases in which values of variables may be sporadically missing.

The difficulty of the missing data problem varies depending on the nature of the "missingness mechanism" [1, 2, 3]. In large datasets, samples containing values Missing Completely At Random (MCAR) can be discarded without biasing the data distribution. In this case, the missingness mechanism is unrelated to any study variable. A slightly more general and frequent case concerns data Missing At Random (MAR) for which the missingness mechanism depends solely on variables with complete information. Data that are neither MCAR nor MAR are referred to Missing Not At Random (MNAR). The literature abounds in algorithms addressing missing data with imputation (replacement of missing values), partial imputation, partial deletion, full analysis, and interpolation [4, 5, 6].

In this paper, we are interested in studying the potential bias introduced in **feature selection** when **handling missing data improperly**. It is generally very tempting to impute missing values to utilize standard feature selection methods. When large amounts of data are missing (e.g. 80%) this becomes almost "necessary" because full records can seldom be found at all, which renders multivariate selection methods very difficult to apply. However, we show that, particularly when large amounts of data are missing, imputation may *introduce* bias in data, even in the presumably "nicest" case of MCAR data. To the best of our knowledge, the **modified t-statistic** we propose and the **use of causal graph to evidence the notion of structural bias** are both new.[1]

---

[1]Source code available in `https://github.com/chalearn/missing-causal-relation.git`.

## 2   Motivation and problem setting

Feature selection may have several purposes [7]: **Prediction** (increasing or least deteriorating prediction performance) or **discovery** (explaining a phenomenon by identifying features associated with a target variable). While prediction and discovery are related and often addressed jointly, they differ in their emphasis on Type I and Type II error. **Type I errors (false positives)** affect more discovery tasks: Wrong associations identified as spurious by expert knowledge may discredit a feature selection method or, if not identified, can lead to harmful decisions (inefficient or detrimental new policies or treatments). **Type II errors (false negatives)** affect more the prediction task. It has been shown that predictors are very tolerant to large amounts of irrelevant variables, but their performance deteriorates a lot when key predictive variables are omitted [7].

To make our point more clearly, we consider **univariate feature ranking** methods. However, our analysis extends naturally to multivariate feature selection. We call "test statistic" any univariate feature ranking criterion. For example, the following t-statistic is commonly used for balanced binary classification problems with continuous features:

$$t_{orig} = \frac{\hat{\mu}_1 - \hat{\mu}_2}{\sigma_p \sqrt{2/n}} \quad \text{with} \quad \sigma_p^2 = \frac{\sigma_1^2 + \sigma_2^2}{2} \quad \text{and} \quad \sigma_i^2 = \sum_{k=1}^{n_i} (x_k - \hat{\mu}_i) , \quad (1)$$

$\hat{\mu}_i$ are the sample means, $\sigma_i$ the sample standard deviations, $\sigma_p$ the pooled within class standard deviation, and $n_1 = n_2 = n/2$ the number of samples per class.

We adopt standard methods of evaluation. The fraction of Type I errors (false positive rate) is assessed by the p-value of a statistical test (e.g. the t-test for the t-statistic). Corrections for multiple testing such as the Bonferroni correction, may be applied on top of our analysis. For non tabulated test statistics, we use "distractors" to emulate a null distribution (features obtained by randomly permuting the values of real features). The p-value is then estimated as the fraction of distractors whose test statistic exceeds the value obtained for the feature being tested. Type II errors are more difficult to assess in real data (for which no ground truth of the relevant features is known). We resort to using the prediction performances of a given predictor to quantify them indirectly.

## 3   Imputation bias: an illustrative example

We selected a didactic example, carved out of the MNIST digit recognition problem [8], to illustrate "imputation bias" (Figure 1): a binary classification problem in which half the features are random distractors (permutations of real features). Features are ranked with the S2N filter (analogous to the t-statistic) and, for predictive modeling, ridge regression is performed with the top ranking features. We vary the fraction of missing values, selected completely at random. Two imputation methods are compared: Median and SVD. The SVD method, praised by many authors [9, 10, 11], capitalizes on correlations of lines and columns in the data matrix, while the median method treats features independently.
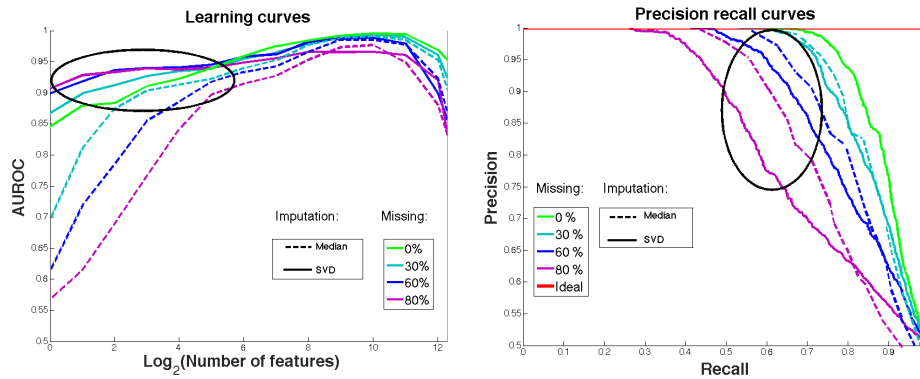
Fig. 1: **SVD biases feature discovery.** Feature selection for a binary classification problem (see text). **Left: Predictive power.** SVD performs suspiciously "too well" at low number of selected features: better results are obtained with MORE missing data! **Right: Discovery power.** For large fractions of missing values, SVD curves drop more quickly than those of median imputations. (Precision = fraction of "true features" retrieved of all features. Recall = fraction of "true features" retrieved of "all true features", a.k.a. True Positive Rate.)

The learning curves (Area under ROC curve as a function of number of features) show the superiority of SVD over median imputation (good prediction power). However, imputing with SVD before feature selection results in more false positives, see precision-recall curves (poor discovery power). Our interpretation is that SVD constructs falsely relevant features, which may be even more relevant than real features, by substituting missing values with combinations of ALL features. This is consistent with the results of the left figure.

## 4 Statistical bias

A first problem encountered with a naive usage of imputation can be traced to well-known statistical biases [12]. Applying Formula 1 to imputed data may yield false positive discoveries by under-estimating variance in data. The formula assumes independently and identically distributed samples, but imputed values break independence assumptions. A simple correction in the case of median imputation would be to divide the class variances by the number of observed values $n_{oi}$, not the total number of values, since imputed values carry no novel information. The analysis of bias introduced by SVD imputation is more complicated. To simplify, assume that we perform (single) imputation of the missing values of a feature of interest $S$ by linear regression of a fully observed helper variable $H$ (correlated both to $S$ and the target $T$). Using the imputed values in the calculation of $\sigma_1$ and $\sigma_2$ may under-estimate the variance for two reasons. Firstly, assuming the linear model $S = aH + noise$ is "correct", using single imputation corresponds to replacing missing values by their expected value $S = aH$, hence ignoring the noise term (which can be estimated by the RMSE residual of the fit $\sigma_r$). Secondly, the regression coefficient is evaluated only from a finite

(small) amount of observed data, hence subject to some uncertainty. In a simple (approximate) formula, we correct for both types of biases:

$$t_{modified} = \frac{\hat{\mu}_1 - \hat{\mu}_2}{\sqrt{2/n} \sqrt{\sigma_p^2 + \underbrace{f_m \sigma_r^2}_{\text{regression residual}} \underbrace{(1 + \alpha/n_o)}_{\text{regression coeff. uncertainty}}}} \quad (2)$$

Compared to Equation 1, $\sigma_p$ is the pooled within class standard deviation (estimated on all samples after imputation), $f_m$ is the fraction of missing values, $\sigma_r$ is the RMSE residual of the fit, $n_o$ the number of observed (non missing) values, and $\alpha$ a positive coefficient. See our supplemental material for details.[2]

## 5 Structural bias

A second problem encountered with a naive usage of imputation stems from wrong assumptions about the underlying data generating process. Simple structural causal model allow us to illustrate this point. Unlike other authors (e.g. [3]) we use causal models to analyze imputation mechanisms, NOT missingness mechanisms (assumed to be MCAR for simplicity). We state the problem as:
- $S$ is a feature of interest; some of its values are missing.
- $H$ is a feature **correlated to $S$ and $T$**; $H$ and $T$ are fully observed.
- $\Sigma$ is feature $S$ after imputation of missing values using H (helper).
- Does $\Sigma \perp T \Rightarrow S \perp T$ and $\Sigma \angle T \Rightarrow S \angle T$? ($\perp \doteq$ independent; $\angle \doteq$ dependent).

We conduct first a listwise deletion (LWD) test of independence between $S$ and $T$ ignoring records with missing values. We then challenge the LWD test result by re-testing after imputation by regression with $H$. Two cases, depending on the outcome of the LWD test (Table 1), are exemplified in Figure 2 and Figure 3. We highlight our "nightmare case" in which the imputation mechanism reverses the causal arrow, which might lead to a false positive dependency $S \angle T$.

Table 1: **Challenging the results of listwise deletion (LWD test)**. Double arrows mean imputation. Directed arrows mean a causal relationship and bidirected arrows the presence of a latent common cause. Undirected edges mean any dependency (causal direction irrelevant). Stars are "wild cards" coding for "arrow or not arrow" e.g. $A \leftarrow *B$ means $A \leftrightarrow B$ or $A \leftarrow B$.

| LWD | Dependencies | Model graph | Imput. graph |
|---|---|---|---|
| $S \angle T$ | Null model ($S \perp T, S \angle H, H \perp T$) | $T \quad S * - * H$ | $T \quad \Sigma \Leftarrow H$ |
| | Alt. model ($S \angle T, S \angle H, H \perp T$) | $T* \rightarrow S \leftarrow *H$ | $T* \rightarrow \Sigma \Leftarrow H$ |
| $S \perp T$ | Null model ($S \perp T, S \angle H, H \angle T$) | $S* \rightarrow H \leftarrow *T$ | $\Sigma \Leftarrow H \leftarrow *T$ |
| | Alt. model ($S \angle T, S \angle H, H \angle T$) | $S - H$ | $\Sigma \Leftarrow H$ |
| | | $T$ | $T$ |

---

[2]Statistical properties of a univariate feature relevance estimator in the presence of missing data, K. Bennett, I. Guyon and B. Seijo-Pardo available in `http://www.lidiagroup.org/index.php/en/materials-en.html`.
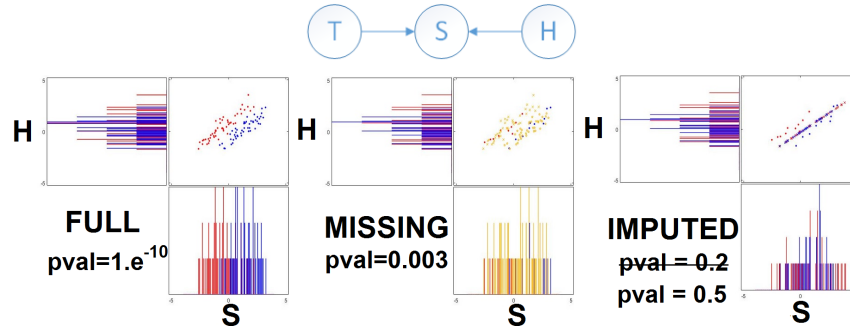
Fig. 2: **Imputation yields FALSE NEGATIVE.** We use the t-test to reveal dependencies between $S$ (source) and $T$ (target) (p-value $\leq 0.01$). $H$ (helper) is an "auxiliary" variable. In histograms and scatter plots of S and H, binary variable $T = \pm 1$ is color coded (red/blue) and missing values are represented in yellow. **Left:** 100 points randomly drawn following the data generating model $T \sim \text{Bern}\,(p = 1/2)\,; H \sim \mathcal{N}(0,1); S = T + H + noise$. The p-value reveal that $T$ **and** $S$ **are significantly DEPENDENT**. **Middle:** $S$ has 80% of values Missing Completely At Random (MCAR). H and T are fully known. The p-value indicates that $T$ **and** $S$ **remain significantly dependent** based on the remaining 20% of **complete data**. **Right:** We impute missing values by regressing $S$ on $H$. Imputation results in a loss of blue/red separation in the histogram of $S$. According to the p-value, **the dependency between $S$ and $T$ is no longer detectable**. Imputation using $H$, carrying no information about $T$, contributed *noise* w.r.t. detecting the dependency between $S$ and $T$.
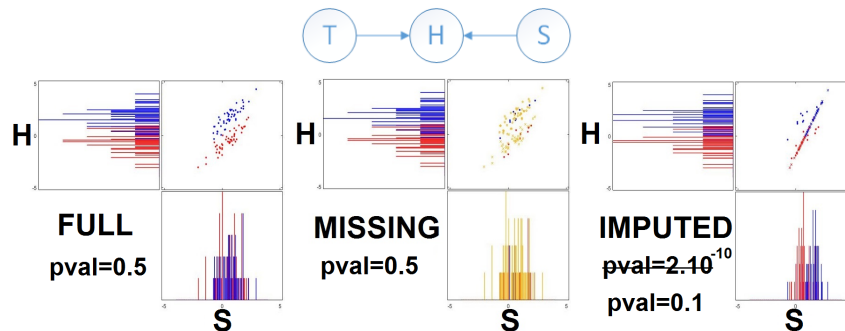


Fig. 3: **Imputation yields FALSE POSITIVE.** Same problem setting as in the previous figure, with a different data generating model for which $S$ **and** $T$ **are INDEPENDENT**. **Left:** We draw 100 points with the model $T \sim \text{Bern}\,(p = 1/2)\,; S \sim \mathcal{N}(0,1); H = T + S + noise$. No significant dependency between $S$ and $T$ is found according to the p-value of the T-test. **Middle:** $S$ has 80% of values Missing Completely At Random (MCAR). H and T are fully known. No change. **Right:** We impute missing values by regressing $S$ on $H$. The imputation model $H \Rightarrow S$ reverses the causal arrow $S \rightarrow H$. After imputation there is a blue/red separation in the histogram of $S$, which did not exist in the original data: $S$ **and** $T$ **become DEPENDENT** (strikethrough p-value). Our proposed correction to the t statistic (non strikethrough text) brings the p-value below the chosen significance level (0.01).

## 6   Conclusion

Imputing missing values before performing feature selection is tempting, particularly when there is a large fraction of missing values (above 80%). Yet this is precisely when it is important to be cautious. Replacing missing values may introduce bias in data, with adverse effects on type I errors (false positives) and type II errors (false negatives). This occurs even for the "nicest" types of missingness mechanisms: MCAR. The types of bias introduced are of two nature: statistical and structural. **Statistical bias** results in improper estimation of variance and/or co-variance between variables and can be corrected either analytically of by multiple imputation. We proposed for univariate feature selection of continuous features and a binary target variable a modified t-statistic, which takes into account the uncertainty of linear regression imputation analytically. It captures both the uncertainty due to the limited accuracy of the regression coefficients (estimated from a small data sample) and the residual of the fit. **Structural bias** is more insidious. It stems from the reversal of causal arrows by the imputation mechanism and can result in an increasing rate of false positives. For problems of prediction, this may not be a problem. But for problems of discovery, when a large fraction of variable values are missing, it is not advisable to use imputation methods such as regression or SVD, if one wants to avoid increasing the false discovery rate. Future work includes devising novel feature selection methods robust to missing data, without requiring imputation of missing values.

## References

[1]  D. B. Rubin. Inference and missing data. *Biometrika*, 63(3):581–592, 1976.

[2]  J.L. Schafer and J.W. Graham. Missing data: Our view of the state of the art. *Psychological Methods*, 7(2):147–177, 2002.

[3]  J. Pearl and K. Mohan. Recoverability and testability of missing data: Introduction and summary of results. *Available at SSRN 2343873*, 2013.

[4]  C. K. Enders. *Applied missing data analysis*. Guilford Press, 2010.

[5]  S. García, J. Luengo, and F. Herrera. *Data preprocessing in data mining*. Springer, 2015.

[6]  N.T. Longford. *Missing data and small-area estimation: Modern analytical equipment for the survey statistician*. Springer Science & Business Media, 2006.

[7]  I. Guyon and A. Elisseeff. An introduction to feature extraction. *Feature extraction*, pages 1–25, 2006.

[8]  I. Guyon, S. Gunn, A. Ben-Hur, and G. Dror. Result analysis of the NIPS 2003 feature selection challenge. In *NIPS*, pages 545–552, 2004.

[9]  M. Kurucz et al. Methods for large scale SVD with missing values. In *Proceedings of KDD cup and workshop*, volume 12, pages 31–38, 2007.

[10]  K. Moorthy et al. A review on missing value imputation algorithms for microarray gene expression data. *Current Bioinformatics*, 9(1):18–22, 2014.

[11]  O. Troyanskaya et al. Missing value estimation methods for DNA microarrays. *Bioinformatics*, 17(6):520–525, 2001.

[12]  A. Donders et al. A gentle introduction to imputation of missing values. *Journal of clinical epidemiology*, 59(10):1087–1091, 2006.