

Algorithmic challenges in Big Data analytics

Verónica Bolón-Canedo¹, Beatriz Remeseiro², Konstantinos Sechidis³,
David Martínez-Rego⁴ and Amparo Alonso-Betanzos¹ *

1- Departamento de Computación, Universidade da Coruña
Campus de Elviña s/n, 15071 A Coruña, Spain

2- Departament de Matemàtiques i Informàtica, Universitat de Barcelona
Gran Via de les Corts Catalanes 585, 08007 Barcelona, Spain

3- School of Computer Science, University of Manchester
M13 9PL Manchester, UK

4- Engineering Building, University College London
Malet Place, WC1E 7JE London, UK

Abstract. This session studies specific challenges that Machine Learning (ML) algorithms have to tackle when faced with Big Data problems. These challenges can arise when any of the dimensions in a ML problem grows significantly: a) size of training set, b) size of test set or c) dimensionality. The studies included in this edition explore the extension of previous ML algorithms and practices to Big Data scenarios. Namely, specific algorithms for recurrent neural network training, ensemble learning, anomaly detection and clustering are proposed. The results obtained show that this new trend of ML problems presents both a challenge and an opportunity to obtain results which could allow ML to be integrated in many new applications in years to come.

1 Introduction

Several years have passed after the term Big Data was coined and it is now clear that no practitioner or academic in the fields of Machine Learning and Computational Statistics can ignore the implications that this new reality has had in their discipline. Big Data challenges cannot be measured just in storage units but in social impact of Machine Learning research and constitutes a new reformulation of goals. While in the past ML restricted itself to well formulated problems such as classification, regression and clustering, being oblivious to the scalability and guarantees of its proposals; nowadays we are reformulating our goals moving towards a true Artificial Intelligence, where these classic well formulated problems are only a (fundamental) piece in the puzzle and ML solutions are required a level of guarantees only reserved to other engineering areas

*This research has been partially funded by project TIN-2015-65069-C2-1-R of Ministerio de Economía y Competitividad (MINECO) of the Spanish Government, and Xunta de Galicia project GRC2014/035, both partially funded by FEDER funds of the European Union. Beatriz Remeseiro acknowledges the support of MINECO under *Juan de la Cierva* Program (ref. FJCI-2014-21194).

in the past. Although Big Data is a significant area of study in itself, it does not substitute the value and importance of past Machine Learning research. On the contrary, it builds on past fundamental results and presents a new set of challenges that only arise in this very specific scenario. This is precisely the aim of the special session *Algorithmic Challenges in Big Data Analytics*, to isolate, study and tackle specific problems that only arise when any of the dimensions of a ML problem: a) size of training set, b) size of test set or c) dimensionality reaches a size that poses a different statistical or computational challenge when compared to its constrained counterpart.

The algorithms proposed in this section range from Recurrent Neural Network training to Similarity Pattern Search, showing the breath of areas where Big Data poses a challenge. This introduction is divided as follows. Section 2 presents a historic review of the emergence of Big Data analytics and its landscape. Section 3 presents a taxonomy of challenges that this special session would like to focus on. Section 4 presents a review of recent contributions in this area of study. Finally, Section 5 closes with a review of the main contributions of this session.

2 Big Data analytics

Storing data and the information derived from it has been present in human civilization from ancient times. But never until last decade has been as much easy, quick and inexpensive to find, copy, share and process data. Besides “datification” of most activities and processes has become possible due to the variety of sensors that have appeared, with constantly decreasing costs, and thus allowing us to represent digitally music, experiences as traveling or scientific advances as sequencing human genomes. As digital entities can be easily replicated, transmitted, modified, analyzed, etc., complete sectors as for example health or finances are being transformed in information and knowledge services. Big Data refers to the voluminous amount of data, structured or unstructured, that has the potential to be analyzed to derive patterns and other information that can add value to the business activities at hand.

Big Data addresses requirements such as the combination of different datasets that might be not related in principle, the processing of large amount of data, perhaps some of it unstructured, and the gathering of hidden information that adds value to the business all in a time sensitivity way. Big Data implies big challenges, as data production is increasing dramatically every year. But not only the volume and velocity (as speed of generation makes that sometimes real-time is not fast enough) are a defiance for Machine Learning algorithms, but also other V’s are part of the equation. In fact, Big Data, a field dedicated to the analysis, processing, and storage of large collections of data frequently originated from different sources [1], can be defined using seven V’s: Volume, Velocity, Variety, Variability, Veracity, Visualization and Value:

- Volume, which is substantial nowadays, in the order of Petabytes, and increasingly growing, imposing new data storage and processing demands.

This volume dimension demands also scalable data preparation, preprocessing and management.

- Velocity, as data is acquired frequently at very fast speed, setting the need for temporally competitive data processing algorithms.
- Variety, as data might be in different formats (textual, image, video, audio, etc.) and types (structured, semi-structured, unstructured). This characteristic results in the need for powerful data integration, transformation, processing and storage methods.
- Variability, that refers to data with constantly changing meanings, due to concept changes, drifts, or in those particular cases in which data gathering relies on language processing.
- Veracity, that reflects the need for quality assessment in the data being acquired, thus implying data preprocessing challenges for removing noisy data, invalid data, etc.
- Visualization, that challenges the data scientists for the need of adequate representation of the data to permit visual perception of the meaning of information derived from the application of data analysis techniques, aiding identifying hidden patterns, anomalies, correlations, etc.
- Value, which refers to the usefulness of the data for the organization. Value is closely related with veracity and depends also on the temporal restrictions of the business field, as the longer the time needed for deriving useful information, the smaller the value.

3 Challenges

In this new scenario in which datasets are becoming larger everyday, Machine Learning researchers find more challenges that need to be faced. In this section we comment on some of them.

3.1 Millions of features

Analogous to Big Data, the term Big Dimensionality has been coined [2] to refer to the unprecedented number of features, which arrive at levels that render the existing Machine Learning methods not useful anymore. For example, the maximum number of features in the datasets posted in the popular UCI Machine Learning Repository [3] is more than 3 million, and in the case of the LIBSVM Database [4], the maximum dimensionality is more than 29 million.

It might seem that having such a huge number of features (and thus, information) is desirable for researchers, but it is not the case. On the one hand, this ultrahigh dimensionality requires massive memory, which implies a high computational cost. And, on the other hand, having such a large number of features

might deteriorate the generation ability of learning algorithms, which is known as *curse of dimensionality* [5].

A common solution to reduce the high dimensionality of the data is to apply a feature selection method before learning and therefore obtaining a smaller subset of practical and useful features. However, state-of-the-art feature selection methods are now confronted by key challenges that potentially have negative repercussions on performance. As an example, Zhai et al. [2] pointed out that it was necessary more than a day of computational effort for the popular methods SVM-RFE and mRMR to crunch the data for a psoriasis single-nucleotide polymorphism (SNP) dataset composed of *just* half a million features.

3.2 Real-time processing

Nowadays, data is being collected at an unprecedented fast pace, and society demands that it has be processed rapidly. We are surrounded by social media networks and portable devices which produce vast amounts of data in real time, so we also need sophisticated methods that are capable of dealing with it and provide valuable information.

Classical batch learning algorithms cannot deal with continuously flowing data streams, which require online approaches. Although some advances have been made in the field [6], there are still challenges that need to be solved such as dynamic feature spaces that would initially be empty but would add features as new information arrives (e.g., documents in their text categorization application).

3.3 Visualization and interpretability

Data is everywhere, continuously increasing, and heterogeneous. In the last few years we are witnessing a kind of Diogenes syndrome referring to data. Given that the cost of storage is being reduced, organizations collect and store huge amounts of data, but they are unable to extract from them useful information. There is a clear need to gather data in a meaningful way, so as to evolve from a data-rich/knowledge-poor scenario, to a data-rich/knowledge-rich scenario. The trend now is to use sophisticated techniques that involve deep neural networks, kernel methods and ensembles of different classifiers. However, although these methods provide impressive prediction accuracies, their outputs are blind and very little insight is available about their inner workings, which complicates the task of the decision makers. Therefore, the challenge now is to achieve user-friendly visualization of the results to facilitate their interpretability and understanding, so in this way Machine Learning algorithms can have impact on consequential real-world applications.

4 Recent contributions

In the era of Big Data, the urgent need for analyzing big databases has revolutionized Machine Learning [7], data mining [8] and statistics [9]. Proof of this

is the growing number of recent textbooks that are dedicated to analyzing Big Data aimed at both researchers and practitioners [10, 11, 12]. Here we present a brief review of recent techniques that can be used for efficient prediction (classification/regression), clustering and dimensionality reduction on massive datasets.

Scaling-up Machine Learning algorithms to deal with large datasets has recently become an important research topic. According to Tong [13], there are two main ways of scaling-up: one is by adapting the algorithm to handle Big Data in a single machine, and the other is by scaling-up prediction algorithms by parallelism. Wang et al. [14] identify three main strategies for analyzing big-data: (1) Sub-sampling, where we repeat the analyses several times in subsets of data sampled through advanced sampling techniques such as Bag of little bootstraps [15] a scalable bootstrap procedure for massive data; (2) Divide and conquer, where we split the data into blocks, process each of them independently and, at the end, combine the results from each block [16]; and, (3) Online updating for streaming data [17].

For all of the popular prediction methods of Machine Learning, different ways for dealing with massive datasets have been suggested. For example, Muja and Lowe [18] present efficient methods for k-nearest neighbor (k-nn) classification by using space partitioning structures, such as the randomized k-d trees or the priority search k-means tree. A different solution is suggested by Deng et al. [19]; in their work, the authors firstly conduct a k-means clustering to split the data into several parts, and then they apply k-nn classification in each one. Recently, Maillo et al. [20] suggested a solution on k-nn classification based on MapReduce, while in [21] Spark was used. A different set of works focuses on scaling-up support vector machines (SVM). For example, Hsieh et al. [22] suggest a divide-and-conquer approach, where they use clustering to partition the kernel SVM problem into smaller subproblems that can be solved independently and efficiently. Rebstroff et al. [23] follow a different approach and present a quantum support vector machine, which achieves a logarithmic complexity in the size of dimensions and the number of training examples. Chen and Xie [16] propose a split-and-conquer approach under the generalized linear regression models (i.e. to solve regression or classification problems), where they use subsets of data to estimate the parameters parallel and then they use weighted combination rules to produce an overall estimate. Another recent approach for both regression and classification in Big Data is [24], where the authors introduce a method for stochastic variational inference for Gaussian process models. While Bem-Haim and Tom-Tov [25] propose a distributed algorithm for building decision trees, which is designed for big and stream data. Finally, during the last few years, there are many new works that explore the challenges and the applications of deep learning in Big Data analytics [26, 27], and more precisely, of deep convolutional networks when the data are images or videos [28].

Another important challenge is clustering big amounts of data. For example, Du et al. [29] suggest an approach for clustering by combining k-nn and principal component analysis, while Boutsidis et al. [30] suggest a method to scale-up clustering by a random projection algorithm for k-means. For clustering Big

Data streams, Hassani and Seidl [31] presents a thorough review of a wide variety of methods.

Finally, as we mentioned earlier, a main challenge of dealing with Big Data is to efficiently handle the big dimensionality [2]. The two main ways of reducing dimensionality is through feature extraction or through feature selection [32]. For feature selection, there is a vast amount of recent works for clustering, regression and classification [6]. For example, Tan et al. [33] present an adaptive feature scaling schema for feature selection and they demonstrate its performance on datasets with tens of million of examples. For information theoretic approaches, a recent work by Ramírez-Gallego et al. [34] presents an extension of the famous minimum-redundancy-maximum-relevance (mRMR), which improves mRMR by transforming the quadratic complexity of mRMR into an efficient greedy process. Finally, Halko et al. [35] present an algorithm for the principal component analysis of datasets that are too large to be stored in random-access memory.

5 Contributions to the special session

The special session *Algorithmic Challenges in Big Data Analytics* has received contributions from different research groups, presenting approaches to deal with both theoretical and applied aspects of the main topic. Each accepted paper is briefly introduced in the following.

Partition-wise Recurrent Neural Networks (pRNNs) are effective RNNs that have their feature space partitioned (discretized) and learn different parameters for different partitions of each feature. Their use is proposed by Jian et al. [36] to detect fishing activities in the ocean by means of point-based trajectory classification and using Automatic Identification System (AIS) data. The authors compare the performance of Gated Recurrent Unit (GRU) and Long Short Term Memory (LSTM) networks using this partitioning approach against that of GRUs and LSTMs not using it on the problem at hand. The experimental results demonstrated that the former group of approaches (pRNNs) outperform the latter one.

Eiras-Franco et al. [37] explore an approach to efficiently construct an approximate k Nearest Neighbors Graph (kNNG) with a *divide and conquer* method based on locality sensitive hashing. The kNNG is a key component for a wide variety of Machine Learning algorithms across a variety of tasks, and its complexity makes necessary the development of efficient scalable solutions to extend these methods to larger datasets. The proposed approach consists in dividing the dataset into disjoint subsets based on locality sensitive hashing, in such a way that a graph is built upon each subset separately and then all these graphs are combined. Their method was compared with the preferred algorithm for obtaining kNNG, known as NN-Descent, and the results demonstrated the adequacy of their approach which provides good performance on distributed architectures.

The negative correlation (NC) ensemble learning approach is studied by Reeve and Brown [38] in a regression context. The authors derived a closed form solution for the minimizer of the NC loss when each function of the en-

semble is a linear map of a fixed basis function, allowing them to give an exact formula for the effective degrees of freedom of this type of NC ensembles. The paper provides an interesting theoretical analysis of the NC method which leads to a better understanding of its behavior with respect to diversity parameter as well as to an efficient way of tuning the parameter. In addition, the experimentation showed that the degrees of freedom formula gives rise to an efficient tuning of the diversity parameter on large dataset, succeeded by plugging the calculated degrees of freedom into the Generalized Cross Validation (GCV) formula.

Fernandez et al. [39] presented a method to reduce the number of random projections used in a Scaled Convex Hull (SCH) one-class classification algorithm, making it more suitable for larger dataset. Their approach introduces a projection pruning phase, in which projections are ranked by the maximum mutual information between each pair of projections. Then, projections with a small maximum mutual information are discarded, getting rid of the redundant projections. The experimentation included a comparison of their approach with a method that randomly removes projections and a method which ranks projections in reverse order based on mutual information. The obtained results proved that the proposed method outperforms the other two strategies, reducing the testing time whilst maintaining the predictive performance.

A distributed approach to split a dataset into several nodes is proposed by Morán-Fernández et al. [40], allowing to deal with huge quantities of data with no degradation in performance. The algorithm horizontally splits the training dataset among many nodes and uses two distance metrics (the Hamming distance and the Hellinger distance) to find the node with the highest similarity. In case of a tie, the complexity is used to decide the best candidate node among the multiple matches obtained. In this manner, the test samples are classified by the model learned from the closest data. The experimental results demonstrated that their approach allows to drastically decrease the running time without affecting the accuracy of the classifier.

References

- [1] Thomas Erl, Wajid Khattak, and Paul Buhler. *Big Data Fundamentals: Concepts, Drivers & Techniques*. Prentice-Hall, 2016.
- [2] Yiteng Zhai, Yew-Soon Ong, and Ivor W Tsang. The Emerging Big Dimensionality. *IEEE Computational Intelligence Magazine*, 9(3):14–26, 2014.
- [3] Arthur Asuncion and David Newman. UCI Machine Learning Repository. <http://www.ics.uci.edu/mllearn/MLRepository.html>, 2007. Accessed: February 2017.
- [4] Chih-Chung Chang and Chih-Jen Lin. LIBSVM Data: Classification, Regression, and Multi-label. <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>. Accessed: February 2017.
- [5] Richard Bellman. Dynamic Programming, Princeton. *NJ: Princeton UP*, 18, 1957.
- [6] Verónica Bolón-Canedo, Noelia Sánchez-Marroño, and Amparo Alonso-Betanzos. Recent advances and emerging challenges of feature selection in the context of big data. *Knowledge-Based Systems*, 86:33–45, 2015.
- [7] Omar Y Al-Jarrah, Paul D Yoo, Sami Muhaidat, George K Karagiannidis, and Kamal Taha. Efficient machine learning for big data: A review. *Big Data Research*, 2(3):87–93, 2015.

- [8] Xindong Wu, Xingquan Zhu, Gong-Qing Wu, and Wei Ding. Data mining with big data. *IEEE Transactions on Knowledge and Data Engineering*, 26(1):97–107, 2014.
- [9] Michael I Jordan. On statistics, computation and scalability. *Bernoulli*, 19(4):1378–1390, 2013.
- [10] Shan Suthaharan. *Machine Learning Models and Algorithms for Big Data Classification: Thinking with Examples for Effective Learning*. Springer, 1st edition, 2015.
- [11] Jared Dean. *Big Data, Data Mining, and Machine Learning: Value Creation for Business Leaders and Practitioners*. Wiley and SAS Business Series. Wiley, 2014.
- [12] Bruce Ratner. *Statistical and Machine-Learning Data Mining: Techniques for Better Predictive Modeling and Analysis of Big Data*. CRC Press, Inc., 2nd edition, 2011.
- [13] Hanghang Tong. Big Data Classification. In *Data Classification: Algorithms and Applications*, pages 275–286. CRC Press, 2014.
- [14] Chun Wang, Ming-Hui Chen, Elizabeth Schifano, Jing Wu, and Jun Yan. Statistical methods and computing for big data. *Statistics and its interface*, 9(4):399–414, 2016.
- [15] Ariel Kleiner, Ameet Talwalkar, Purnamrita Sarkar, and Michael I Jordan. A scalable bootstrap for massive data. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(4):795–816, 2014.
- [16] Xueying Chen and Min-ge Xie. A split-and-conquer approach for analysis of extraordinarily large data. *Statistica Sinica*, pages 1655–1684, 2014.
- [17] Elizabeth D Schifano, Jing Wu, Chun Wang, Jun Yan, and Ming-Hui Chen. Online updating of statistical inference in the big data setting. *Technometrics*, 58(3):393–403, 2016.
- [18] Marius Muja and David G Lowe. Scalable nearest neighbor algorithms for high dimensional data. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(11):2227–2240, 2014.
- [19] Zhenyun Deng, Xiaoshu Zhu, Debo Cheng, Ming Zong, and Shichao Zhang. Efficient kNN classification algorithm for big data. *Neurocomputing*, 195:143–148, 2016.
- [20] Jesús Mailló, Isaac Triguero, and Francisco Herrera. A MapReduce-Based k-Nearest Neighbor Approach for Big Data Classification. In *IEEE Trustcom/BigDataSE/ISPA*, volume 2, pages 167–172, 2015.
- [21] Jesus Mailló, Sergio Ramírez, Isaac Triguero, and Francisco Herrera. kNN-IS: An Iterative Spark-based design of the k-Nearest Neighbors classifier for big data. *Knowledge-Based Systems*, 117:3–15, 2017.
- [22] Cho-Jui Hsieh, Si Si, and Inderjit Dhillon. A Divide-and-Conquer Solver for Kernel Support Vector Machines. In *Proceedings of the 31st International Conference on Machine Learning*, pages 566–574, 2014.
- [23] Patrick Reberntrost, Masoud Mohseni, and Seth Lloyd. Quantum Support Vector Machine for Big Data Classification. *Physical Review Letters*, 113:130503, 2014.
- [24] James Hensman, Nicolò Fusi, and Neil D Lawrence. Gaussian processes for Big Data. In *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence*, pages 282–290, 2013.
- [25] Yael Ben-Haim and Elad Tom-Tov. A streaming parallel decision tree algorithm. *Journal of Machine Learning Research*, 11(Feb):849–872, 2010.
- [26] Xue-Wen Chen and Xiaotong Lin. Big data deep learning: challenges and perspectives. *IEEE Access*, 2:514–525, 2014.
- [27] Maryam M Najafabadi, Flavio Villanustre, Taghi M Khoshgoftaar, Naeem Seliya, Randall Wald, and Edin Muharemagic. Deep learning applications and challenges in big data analytics. *Journal of Big Data*, 2(1):1, 2015.
- [28] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521:436–444, 2015.

- [29] Mingjing Du, Shifei Ding, and Hongjie Jia. Study on density peaks clustering based on k-nearest neighbors and principal component analysis. *Knowledge-Based Systems*, 99:135–145, 2016.
- [30] Christos Boutsidis, Anastasios Zouzias, and Petros Drineas. Random projections for k-means clustering. In *Advances in Neural Information Processing Systems*, pages 298–306, 2010.
- [31] Marwan Hassani and Thomas Seidl. Clustering Big Data streams: recent challenges and contributions. *it-Information Technology*, 58(4):206–213, 2016.
- [32] Isabelle Guyon, Steve Gunn, Masoud Nikravesh, and Lofti A Zadeh. *Feature extraction: foundations and applications*, volume 207. Springer, 2008.
- [33] Mingkui Tan, Ivor W. Tsang, and Li Wang. Towards ultrahigh dimensional feature selection for big data. *Journal of Machine Learning Research*, 15(1):1371–1429, 2014.
- [34] Sergio Ramírez-Gallego, Iago Lastra, David Martínez-Rego, Verónica Bolón-Canedo, José Manuel Benítez, Francisco Herrera, and Amparo Alonso-Betanzos. Fast-mRMR: Fast Minimum Redundancy Maximum Relevance Algorithm for High-Dimensional Big Data. *International Journal of Intelligent Systems*, 32(2):134–152, 2017.
- [35] Nathan Halko, Per-Gunnar Martinsson, Yoel Shkolnisky, and Mark Tygert. An algorithm for the principal component analysis of large data sets. *SIAM Journal on Scientific computing*, 33(5):2580–2594, 2011.
- [36] Xian Jiang, Eirco N. de Souza, Xuan Liu, Behrouz Haji Soleimani, Xiaoguang Wang, Daniel L. Silver, and Stan Matwin. Partition-wise Recurrent Neural Networks for Point-based AIS Trajectory Classification. In *European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, 2017.
- [37] Carlos Eiras-Franco, Leslie Kanthan, Amparo Alonso-Betanzos, and David Martínez-Rego. Scalable approximate k-NN Graph construction based on Local Sensitivity Hashing. In *European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, 2017.
- [38] Henry Reeve and Gavin Brown. Freedom and diversity in regression ensembles. In *European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, 2017.
- [39] Diego Fernandez-Francos, Oscar Fontenla-Romero, Amparo Alonso-Betanzos, and Gavin Brown. Mutual information for improving the efficiency of the SCH algorithm. In *European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, 2017.
- [40] Laura Morán-Fernández, Verónica Bolón-Canedo, and Amparo Alonso-Betanzos. A distributed approach for classification using distance metrics. In *European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, 2017.

