

LANN-DSVD: A privacy-preserving distributed algorithm for machine learning

Oscar Fontenla-Romero, Bertha Guijarro-Berdiñas, Beatriz Pérez-Sánchez
and Marcelo Gómez-Casal *

Universidade da Coruña
Centro de Investigación en Tecnologías de la Información y
las Comunicaciones (CITIC)
Elviña, 15071 - Spain

Abstract. In the Big Data era new challenges have arisen in machine learning related with the Volume (high number of samples or variables), the Velocity, etc. making many of the classic and brilliant methods not applicable. One main concern derives from Privacy issues when data is distributed and cannot be shared among locations. Herein, we present LANN-DSVD a non iterative algorithm for One-Layer Neural Networks that allows distributed learning guaranteeing privacy. Moreover, it is non iterative, parameter-free and provides incremental learning, thus making it very suitable to manage huge and/or continuous data. Results demonstrate its competitiveness both in efficiency and efficacy.

1 Introduction

Over recent years, with the appearance of new possibilities such as the Internet of Things, new challenges have arisen for the machine learning field. The more obvious challenge is related to the scalability and efficiency of machine learning algorithms and its ability to deal with huge datasets. However, some other interesting challenges arise related to certain external restrictions that those algorithms must accomplish as imposed by the conditions under which data is available for training. One of these conditions occurs when data is distributed. In order to overcome the issue of learning from distributed data, one could think on the simple solution of gathering data at a single location. However, this solution can be unrealistic or ineffective: the total storage capacity needed to store a single big dataset might not be affordable or the spatial complexity of the algorithms could impede to process the whole training dataset due to time or memory restrictions. Moreover, the necessary bandwidth to efficiently transmit the data to a single location might not be available: the Velocity of data, one of the “V” big data properties, could cause to have frequently updated databases and the required communication could be a continuous unmanageable overhead. In addition, even when communication cost was not too high or the spatial complexity was not a problem, it is every day more often the case that sensitive data cannot be moved around distributed locations due to

*This research has been supported in part by the Spanish Government’s Secretaría de Estado de Investigación (Grant TIN2015-65069-C2-1-R), the Xunta de Galicia (Grants GRC2014/035, ED431G/01 and ED341D R2016/045) and European Union FEDER funds.

privacy issues. In all these cases, distributed learning is the perfect solution, as it is a natural way of scaling up learning algorithms and to deal with natural distributed data, and thus it has become an active and promising research line for big data learning. Moreover, the popularization of multicore processors and computer clusters led to a suitable context for the appearance of new distributed algorithms [1]. Nevertheless, most of them are focused only on exploiting parallel processing in order to speed up learning or overcome overflow memory problems. Very often they do not take into account issues related to naturally distributed data and, more specifically, to the management of sensitive data through privacy-preserving learning. As a consequence, they work by moving data between locations during learning to obtain the final aggregate model. However, in many applications domains, data might belong to different, perhaps competing, organizations that want to exchange knowledge and share the benefit of a machine-learning model that aggregates all their data but without the need of exchanging raw private data among them [1, 2]. In this paper, we focus our attention on this specific context, in which distributed data is not allowed to be moved between locations, to propose a new privacy-preserving distributed learning algorithm: LANN-DSVD. This algorithm is based on the LANN-SVD algorithm [3][4], a previous contribution of the authors. The cornerstone of this method is its mathematical formulation that makes the optimization problem separable, thus allowing to learn in parallel from each partition.

The paper is structured as follows: Section 2 describes the proposed distributed learning algorithm and the background needed for comprehension, Section 3 presents the experimental results that demonstrate its learning accuracy, and Section 4 includes conclusions.

2 LANN-DSVD

In our context of interest, we can define a distributed computing environment as a set of locations physically separated but probably connected by a communication network. In this context, a common strategy to obtain a valid machine learning model is local learning and further model integration. LANN-DSVD takes advantage of this learning paradigm, avoiding moving raw data across the locations and thus allowing privacy-preserving learning as well as minimizing communications, as it is described below.

2.1 Local learning

When working with large datasets, it is essential to employ low complexity algorithms due to time and memory restrictions. Moreover, in this context, local models must allow to be integrated somehow to obtain the final model. With this conditions in mind, as local learners we employ single-layer (no hidden layers) Artificial Neural Networks (ANN) trained with LANN-SVD [3][4]. This a very efficient learning algorithm as it computes the weights of the network analytically and also it is able to learn in an incremental, and hence distributed, way. Consider a single-layer ANN. For the sake of comprehension, in what follows,

we will consider only one output. The equation relating the inputs $\mathbf{X} \in \mathbb{R}^{m \times n}$ and outputs $\mathbf{y} \in \mathbb{R}^{n \times 1}$ is given by:

$$\mathbf{y} = f(\mathbf{X}^T \mathbf{w})$$

f being the nonlinear activation function of the neurons and \mathbf{w} the weights. Let $\mathbf{d} \in \mathbb{R}^{n \times 1}$ be the vector with the desired outputs for training. Assuming that f is invertible, instead of computing the error at the output of the network between \mathbf{y} and \mathbf{d} , LANN-SVD minimizes the mean-squared-error before f , i.e. between $\mathbf{X}^T \mathbf{w}$ and $\bar{\mathbf{d}} = f^{-1}(\mathbf{d})$. This error function does not contain local minima and its global solution for \mathbf{w} can be obtained very efficiently by solving a system of linear equations defined by:

$$\mathbf{XFFX}^T \mathbf{w} = \mathbf{XFF}\bar{\mathbf{d}} \quad (1)$$

$\mathbf{F} = \text{diag}(f'(\bar{d}_1), f'(\bar{d}_2), \dots, f'(\bar{d}_n))$ being a diagonal matrix formed by the derivative of the f function in the components of the $\bar{\mathbf{d}}$ vector. Using the Singular Value Decomposition (SVD) of $\mathbf{XF} = \mathbf{USV}^T$ and after some transformations [3] an equivalent expression is obtained to calculate the optimal weights:

$$\mathbf{w} \approx \mathbf{U}(\mathbf{SS}^T)^\dagger \mathbf{U}^T \mathbf{XFF}\bar{\mathbf{d}} \quad (2)$$

where \dagger stands for the Moore-Penrose pseudoinverse. The advantage is that the size of the system of linear equations in (2) depends on $r \leq \min(m, n)$ and, therefore, the computational complexity of this calculation relies on the smaller value between the number of data n and the number of variables m of the problem making the training algorithm suitable for both situations.

2.2 Model integration

In this section, we describe the proposed algorithm LANN-DSVD (*Linear learning algorithm for Artificial Neural Networks by Distributed Single Value Decomposition*). Suppose we have K locations. Once classifiers are trained on their corresponding (local) subset of data $\mathbf{X}_k; k = 1 \dots K$, the aim is to integrate these local models using some combination method to obtain a single (global) ANN representing the union of knowledge stored distributively. This can be achieved if we can obtain, from the corresponding local elements, the elements implied in Equation 2 as if they were calculated using the entire dataset $\mathbf{X} = [\mathbf{X}_1 | \mathbf{X}_2 | \dots | \mathbf{X}_K]$.

The solution can be found by identifying two different parts in the right hand side of Equation 2. The first part $\mathbf{U}(\mathbf{SS}^T)^\dagger \mathbf{U}^T$ depends on matrices \mathbf{S} and \mathbf{U} obtained from the SVD of matrix \mathbf{XF} which ideally would contain information on the entire data set. At each location we have already computed the SVD of $\mathbf{X}_k \mathbf{F}_k = \mathbf{U}_k \mathbf{S}_k \mathbf{V}_k^T; k = 1, \dots, K$ (the calculation of matrix \mathbf{V}_k can be avoided as it is not needed). Using the results by Iwen et al. [5] we can calculate the terms \mathbf{U} and \mathbf{S} corresponding to SVD of the global \mathbf{XF} matrix from the local results as $SVD(\mathbf{XF}) = SVD([\mathbf{U}_1 \mathbf{S}_1 | \mathbf{U}_2 \mathbf{S}_2 | \dots | \mathbf{U}_K \mathbf{S}_K])$.

The second part of Equation (2), $\mathbf{XFF}\bar{\mathbf{d}}$, can also be obtained from local calculations. As can be demonstrated [4], this term can be calculated incrementally as $\mathbf{XFF}\bar{\mathbf{d}} = \mathbf{M}_1 + \mathbf{M}_2 + \dots + \mathbf{M}_K$ where $\mathbf{M}_k = \mathbf{X}_k \mathbf{F}_k \mathbf{F}_k \bar{\mathbf{d}}_k$.

Putting the above results together, to obtain the global optimal weights \mathbf{w} only local matrices \mathbf{U}_k , \mathbf{S}_k and \mathbf{M}_k are needed. Therefore, no raw data is transmitted among locations thus guaranteeing privacy. In the interest of reproducible research the Matlab code is available at <http://mloss.org/software/>.

3 Experimental Results

In this section, we present some results to demonstrate the behavior in accuracy and efficiency of LANN-DSVD. With this aim, it will be compared to a Deep Neural Network (DNN), a main Big Data paradigm nowadays, trained in both distributed and batch setups using the implementation provided by Tanaka et al. [6]. To apply DNN to a distributed environment a ring algorithm was followed: node 1 computes the first DNN over \mathbf{X}_1 , afterwards, this DNN₁ is forwarded to node 2 to continue training over \mathbf{X}_2 and so on, until node K computes the final DNN_K.

Table 1 contains the data sets used in the experiments, treated as binary classification problems, that were selected to check different situations: small and large data sets with more samples than variables and viceversa.

To simulate a distributed environment the first three data sets were splitted into 5 non-overlapped blocks while the last two data sets were splitted into 10. The adequate DNN topologies and hyperparameters were experimentally obtained using grid search for the first three data sets whereas for MiniBoone and Susy data sets the employed architecture was the one in the research by Baldi et al. [7]. Table 1 shows the topologies that obtained the best results.

All experiments were carried out in the MatLab environment using a 64-bit CPU Quad Core (clock speed of 3.60GHz).

In order to obtain reliable results the Area Under the Curve (AUC) was calculated using 10-fold cross-validation. Table 2 shows the mean Area Under the Curve (AUC) and the standard deviation obtained. As can be noticed, only one result column is included for LANN-DSVD as it obtains exactly the same values as its batch version. This fact contrasts with what happens with DNN that does not guarantee the same results. Finally, we can observe that, despite of its simplicity, in all cases the ANN trained by LANN-DSVD achieved invariably competitive results when compare to a DNN.

A second comparison was made in terms of computational time. It is noticeable that for the Susy data set, the most significant case due its size, the learning process of the DNN needed 41,456 seconds in the case of the distributed setup and 195,301 seconds in the batch scenario whilst LANN-DSVD only needed, respectively, 22.2 and 21.4 seconds.

Data set	Samples	Variables	DNN topology
Lung [8]	181	12,533	[12533,10,2]
Smk [9]	187	19,993	[19993,10,2]
Ovarian [8]	253	15,154	[15154,10,2]
MiniBooNE [10]	130,064	50	[50,300,300,300,300,2]
Susy [10, 7]	5,000,000	18	[18,300,300,300,300,2]

Table 1: Data sets and DNN topology used for the experimental study.

Data set	AUC LANN-DSVD	AUC DNN-Inc	AUC DNN-Batch
Lung	1.000 \pm 0.000	0.983 \pm 0.053	0.983 \pm 0.053
Smk	0.795 \pm 0.075	0.640 \pm 0.117	0.719 \pm 0.064
Ovarian	1.000 \pm 0.000	0.984 \pm 0.029	0.997 \pm 0.010
MiniBooNE	0.955 \pm 0.002	0.842 \pm 0.025	0.836 \pm 0.014
Susy	0.836 \pm 0.000	0.756 \pm 0.002	0.798 \pm 0.001

Table 2: Mean test AUC \pm standard deviation for LANN-DSVD and DNN in an incremental and batch scenario.

4 Conclusions

Distributed learning is one active and promising line of research in order to deal with large and/or distributed data. In this paper, LANN-DSVD, a new distributed machine learning algorithm based on one-layer ANNs has been presented. This algorithm exhibits many interesting characteristics:

- It is fast and accurate.
- It is able to learn from distributed data while maintaining its privacy.
- Learning the distributed data sets can be done in parallel, in contrast with DNNs or other incremental learning algorithms that need to follow a ring scheme.
- Inside every location, if the amount of data is big, learning can be further distributed. Therefore, it is an interesting algorithm even if we only have one location but multi-core processors.
- It is parameter free, a highly desirable characteristic for large-scale learning environments, where tuning a model can be a very time consuming task.
- It allows incremental learning, which together with its parameter free characteristics, makes the algorithm very suitable to online learning from data streams.

These results encourage further research in the algorithm. As a near future work, the influence of several other Big Data Vs (Variety, Variability, Veracity...) will be studied.

References

- [1] D. Peteiro-Barral and B. Guijarro-Berdiñas. A survey of methods for distributed machine learning. *Progress in Artificial Intelligence*, 2:1–11, 2013.
- [2] O. Ohrimenko, F. Schuster, C. Fournet, A. Mehta, S. Nowozin, K. Vaswani, and M. Costa. Oblivious multi-party machine learning on trusted processors. In *25th USENIX Security Symposium*, pages 619–636, 2016.
- [3] O. Fontenla-Romero, B. Pérez-Sánchez, and B. Guijarro-Berdiñas. LANN-SVD: A non-iterative SVD-based learning algorithm for one-layer neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 2017. In Press.
- [4] O. Fontenla-Romero, B. Pérez-Sánchez, and B. Guijarro-Berdiñas. An incremental non-iterative learning method for one-layer feedforward neural networks. *Applied Soft Computing*, 2017. In Press.
- [5] M. A. Iwen and B. W. Ong. A distributed and incremental SVD algorithm for agglomerative data analysis on large networks. *SIAM Journal on Matrix Analysis and Applications*, 37:1699–1718, 2016.
- [6] M. Tanaka and M. Okutomi. A novel inference of a restricted boltzmann machine. In *IEEE International Conference on Pattern Recognition (ICPR2014)*, pages 1526–1531, 2014.
- [7] P. Baldi, P. Sadowski, and D. Whiteson. Searching for exotic particles in high-energy physics with deep learning. *Nature Communications*, 5:4308EP, 2014.
- [8] ELVIRA Biomedical Data Set Repository. <http://leo.ugr.es/elvira/DBCRepository/>, [Online; accessed November 2017].
- [9] Feature Selection Datasets at Arizona State University. <http://featureselection.asu.edu/datasets.php>, [Online; accessed November 2017].
- [10] M. Lichman. UCI machine learning repository. <http://archive.ics.uci.edu/ml>, [Online; accessed November 2017].