

A distributed learning algorithm based on two-layer artificial neural networks and genetic algorithms

Diego Peteiro-Barral, Bertha Guijarro-Berdiñas, Beatriz Pérez-Sánchez and Óscar Fontenla-Romero *

University of A Coruña, Faculty of Informatics - Dept of Computer Science
Campus de Elviña s/n, 15071 A Coruña - Spain

Abstract. In many real-world applications of machine learning, the amount of data is now beyond the capability of learning algorithms because they cannot process all available data in a reasonable time. Moreover, most large datasets are naturally distributed or they are being stored in a distributed manner. A promising line of research in order to deal with large and/or distributed data is distributed learning. We propose a new distributed learning algorithm based on two-layer artificial neural networks and genetic algorithms. The results obtained show that our method performs better than other distributed learning algorithms.

1 Introduction

Nowadays, machine learning algorithms have to deal most often with large and/or distributed data. Nevertheless, practically all existing implementations of algorithms operate with the training set entirely in main memory, and thus they cannot deal with neither large nor distributed data. Firstly, in order to overcome the issue of learning from large volumes of data, preprocessing techniques (e.g. subsampling) are frequently used, but reducing the size of training sets often decreases the performance of learning algorithms [1]. On the other hand, in order to overcome the issue of learning from distributed data, gathering data in a single location is the most common practice. However, this is often unrealistic or ineffective due to the fact that the necessary storage capacity might not be affordable (the cost of storing a single dataset is much larger than the sum of the costs of storing smaller parts of the dataset); and secondly, the necessary bandwidth to efficiently transmit the data to a single location might not be available (note also that it is common to have frequently updated databases and the required communication could be a continuous overhead) [2]. In addition, even when communication cost was not too high, it is often the case that sensitive data cannot be moved around distributed locations due to privacy issues. One of the most promising lines of research in order to deal with large and/or distributed data is distributed learning, since

- Large volumes of data can be scattered across different locations.

*This research was funded by the Xunta de Galicia by projects 2008/060 and PGIDT-08TIC012105PR, and the Ministerio de Ciencia e Innovación by project TIN2009-10748.

- Distributed learning makes possible learning from distributed data by minimizing communication costs.

The current trend of reducing speed of processors in favor of multicore processors and computer clusters leads to a suitable context for developing new distributed algorithms. However, few algorithms have been proposed so far in the literature and, even more, some of them focus only on exploiting parallel processing in order to speed up learning, but they do not take into account issues related to naturally distributed data as privacy-preserving computation or data skewness.

Our previous contribution in the field of distributed learning is the algorithm DEvoNet (Distributed Evolved Networks) [3]. DEvoNet is a fast and accurate algorithm based on local learning and further model integration. As local classifiers, it uses single-layer (no hidden layer) artificial neural networks (ANNs) and, as model integration method, genetic algorithms (GAs). DEvoNet performs well on many datasets, but nonlinear problems cannot be properly represented by a single-layer ANN. Algorithm DEvoNet-2L is developed using two-layer ANNs in order to overcome this limitation.

This paper is structured as follows: Section 2 describes the proposed distributed learning algorithm DEvoNet-2L, Section 3 presents the experimental results obtained, and Sections 4 and 5 include discussion and conclusions.

2 DEvoNet-2L

Following DEvoNet, in the first place, a set of classifiers is trained locally on distributed data and, subsequently, the classifiers are integrated.

2.1 Local learning

As local classifiers, DEvoNet-2L uses two-layer ANNs trained with the efficient learning algorithm SBLLM [4]. This learning method considers an ANN as composed of two subnetworks (see Fig. 1). The weights of layers 1 and 2 are independently computed by minimizing the loss functions $Q^{(1)}$ and $Q^{(2)}$, respectively. Functions Q are based on the mean squared error (MSE) computed *before* the nonlinear activation functions g_k and f_j rather than *after* them, as in regular learning algorithms.

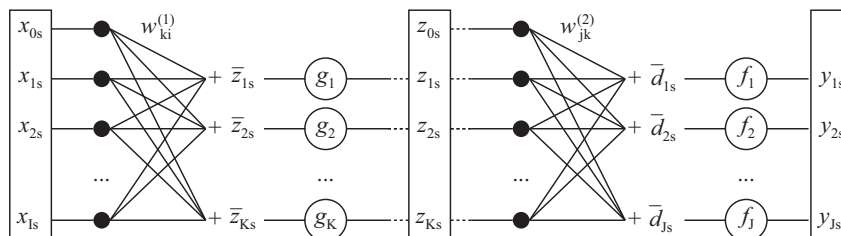


Fig. 1: Two-layer feedforward artificial neural network.

In the first place, considering layer 1 of the ANN, the set of equations relating inputs and outputs of the hidden layer is given by

$$z_{ks} = g_k \left(\sum_{i=0}^I w_{ki}^{(1)} x_{is} \right); k = 1, \dots, K; s = 1, \dots, S \quad (1)$$

where I , K and S are the number of inputs, hidden units and samples, respectively, $x_{0s} = 1$, and $w_{ki}^{(1)}$ is the weight of the connection between the i th input unit and the k th hidden unit. Based on Eq. 1, z_{ks} being the desired output for hidden neuron k and $\bar{z}_{ks} = g_k^{-1}(z_{ks})$, the loss function for layer 1 is defined as

$$Q^{(1)} = \sum_{s=1}^S \sum_{k=1}^K \left(\sum_{i=0}^I w_{ki}^{(1)} x_{is} - \bar{z}_{ks} \right)^2 \quad (2)$$

In the second place, considering layer 2 of the ANN, the set of equations relating inputs of the hidden layer and outputs is given by

$$y_{js} = f_j \left(\sum_{k=0}^K w_{jk}^{(2)} z_{ks} \right); j = 1, \dots, J; s = 1, \dots, S \quad (3)$$

where J is the number of outputs, $z_{0s} = 1$, and $w_{jk}^{(2)}$ is the weight of the connection between the k th hidden unit and the j th output unit. Based on Eq. 3, d_{js} being the the desired output for output neuron j and $\bar{d}_{js} = f_j^{-1}(d_{js})$, the loss function for layer 2 is defined as

$$Q^{(2)} = \sum_{s=1}^S \sum_{j=1}^J \left(\sum_{k=0}^K w_{jk}^{(2)} z_{ks} - \bar{d}_{js} \right)^2 \quad (4)$$

Notice that the weights in Eqs. 2 and 4 are not affected by the nonlinear activation functions and, as a result, the errors are linear with respect to the parameters of the ANN. When working with large datasets it is essential to keep the complexity of the algorithms as low as possible due to time and memory restrictions. The advantage of the loss functions Q is that the weights can be easily computed by solving a system of linear equations that are obtained by deriving $Q^{(1)}$ and $Q^{(2)}$ with respect to the weights and equating to zero, that is

$$\begin{aligned} \sum_{i=0}^I A_{pi}^{(1)} w_{ki}^{(1)} &= b_{pk}^{(1)}; p = 0, \dots, I; k = 1, \dots, K \\ \sum_{k=0}^K A_{qk}^{(2)} w_{jk}^{(2)} &= b_{qj}^{(2)}; q = 0, \dots, K; j = 1, \dots, J \end{aligned} \quad (5)$$

where

$$\begin{cases} A_{pi}^{(1)} = \sum_{s=1}^S x_{is} x_{ps} \\ b_{pk}^{(1)} = \sum_{s=1}^S \bar{z}_{ks} x_{ps} \end{cases} \quad (6)$$

and

$$\begin{cases} A_{qk}^{(2)} = \sum_{s=1}^S z_{ks} z_{qs} \\ b_{qj}^{(2)} = \sum_{s=1}^S \bar{d}_{js} z_{qs} \end{cases} \quad (7)$$

The solution to Eqs. 6 and 7 is unrelated to the order of the samples due to the commutative and associative properties of the sum, i.e. this method is able to learn incrementally since the coefficients $A_{pi}^{(1)}$, $b_{pk}^{(1)}$, $A_{qk}^{(2)}$ and $b_{qj}^{(2)}$ are simply computed as a sum of terms.

2.2 Model integration

Due to the incremental learning property of the algorithm SBLLM, a single (global) ANN representing the union of knowledge stored by all distributed (local) ANNs may be obtained by summing their corresponding matrices of coefficients $A^{(1)}$, $b^{(1)}$, $A^{(2)}$ and $b^{(2)}$ (see Fig. 2).

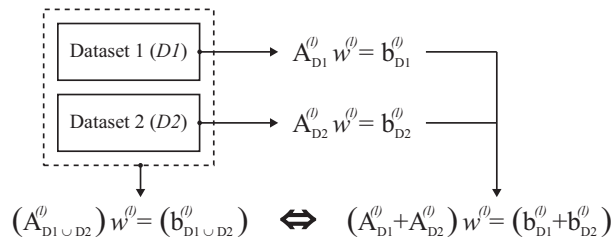


Fig. 2: Combination of ANNs using matrices of coefficients $A^{(l)}$ and $b^{(l)}$; $l = 1, 2$.

The global ANN computed following this method may be usable since it store the knowledge of all local ANNs. However, when talking about distributed learning, it is common that different partitions of data are not identical since naturally distributed real-world datasets have an inherent data skewness property, e.g. data related to diseases from hospitals around the world. For this reason, computing a global ANN by simply summing the matrices of coefficients of local ANNs, each trained on a data partition, may derive to a biased global ANN. In order to overcome this issue, a GA is used during the integration process of ANNs, where the initial population of the GA are such local ANNs. Basically, a GA is defined by the next three elements: *a*) the *fitness function* defines the optimality of a solution (individual of the population) and it is minimized during the execution of the GA. The standard classification error [5] was used in this work; *b*) the *crossover operator* defines how two individuals are combined. Due to the advantages of the learning algorithm used to train the ANNs, the crossover operator is simply defined as the sum of the matrices of coefficients $A^{(1)}$, $b^{(1)}$, $A^{(2)}$ and $b^{(2)}$ of two individuals, respectively; and *c*) the *mutation operator* defines how small random disturbances affects to an individual. In this work, a few elements of the matrices of coefficients $A^{(1)}$, $b^{(1)}$, $A^{(2)}$ and $b^{(2)}$ are altered with small normally-distributed probability.

The hypothesis here is that the crossover operator will find the best combination of local ANNs, and the mutation operator will reduce the impact of data skewness on the performance of the global ANN.

3 Experimental results

3.1 Materials and methods

The assessment of the proposed distributed learning algorithm was performed using the standard classification accuracy [5] computed on five well-known datasets, selected from the *UCI Machine Learning Repository* [6]. Data were divided using *holdout validation*, 90% of data were used for training while the remaining 10% were used for testing. Training data were scattered across 10 distributed locations in order to simulate a distributed environment. In order to ensure unbiased results, experimentation was performed 100 times.

The proposed algorithm was compared against the original single-layer ANN version of DEvoNet, as well as two well-known distributed learning algorithms:

- *Majority vote* [7] is a classic method for combining classifiers. A sample is classified as the class with the highest number of votes among classifiers.
- *Stacking* [8] combines classifiers by learning the way that their outputs correlates with their true class. Once local classifiers were trained, for every sample e of an independent dataset, the output $C_i(e)$ of the classifiers with respect to *all* classes along with the true class of the sample are used to form a sample m of a new dataset M . Each sample will be of the form $[C_1(e), C_2(e) \dots C_N(e)]$. Finally, a global classifier is trained from M .

Majority vote and Stacking are independent of the type of classifiers, in contrast to DEvoNet and DEvoNet-2L. In order to set up a fair comparative, two-layer ANNs trained with scaled conjugate gradient (SCG) [9] were used in majority vote and stacking. SCG shows a good tradeoff between computational requirements and accuracy and, in fact, SCG performs better than SBLLM on many not distributed datasets [4].

3.2 Results

Table 1 shows the mean test accuracies and standard deviations, giving comparative results with respect to each of the four distributed learning algorithms.

	Forest	KDD99	Magic	Mushroom	Shuttle
Majority	71.27 ± 0.93	99.37 ± 0.19	83.66 ± 1.36	87.41 ± 1.23	97.02 ± 1.46
Stacking	74.73 ± 0.67	99.78 ± 0.03	84.12 ± 1.16	87.64 ± 0.86	98.34 ± 0.87
DEvoNet	75.55 ± 0.72	99.54 ± 0.02	78.56 ± 0.87	82.19 ± 0.91	89.50 ± 0.93
DEvoNet-2L	76.09 ± 0.51	99.87 ± 0.03	82.92 ± 0.82	86.94 ± 0.88	98.59 ± 1.03

Table 1: Comparative mean test accuracies (%) and standard deviations.

4 Discussion

DEvoNet-2L performs on average better than Majority vote, an improvement by 1.14%. With regards to DEvoNet, two-layer ANNs are capable of learning nonlinearly separable problems, while single-layer ANNs are not. As can be seen in Table 1, DEvoNet-2L overcomes DEvoNet in all cases, an improvement by 3.81%. Finally, DEvoNet-2L performs as well as Stacking with disagreement under 0.05%. However, notice that the complexity of dataset M in Stacking (see Section 3.1) is determined by the number of classifiers *times* the number of classes, such that it may be so high in real-world datasets with numerous locations and classes. Thus, DEvoNet-2L can solve problems that Stacking would not be able to solve due to a extremely large number of features in dataset M .

5 Conclusions

Distributed learning is one of the most promising lines of research in order to deal with large and/or distributed data. On the one hand, a large amount of data can be scattered across several locations, turning an impractical algorithm into a practical one. On the other hand, distributed learning avoids the necessity of gathering naturally distributed data into a single location. In this paper, a new distributed machine learning algorithm based on two-layer ANNs and GAs has been proposed. DEvoNet-2L is a fast and accurate algorithm, which is able to learn from large and/or distributed data while maintaining its privacy.

For future work, a wider comparison of existing distributed learning algorithms in terms of accuracy, complexity and scalability is needed.

References

- [1] J. Catlett. *Megainduction: machine learning on very large databases*. PhD thesis, School of Computer Science, University of Technology, Sydney, Australia, 1991.
- [2] G. Tsoumakas. Distributed Data Mining. *Database Technologies: Concepts, Methodologies, Tools, and Applications*, pages 157–171, 2009.
- [3] B. Guijarro-Berdiñas, D. Martínez-Rego, and S. Fernández-Lorenzo. Privacy-Preserving Distributed Learning Based on Genetic Algorithms and Artificial Neural Networks. *Distributed Computing, Artificial Intelligence, Bioinformatics, Soft Computing, and Ambient Assisted Living*, pages 195–202, 2009.
- [4] E. Castillo, B. Guijarro-Berdiñas, O. Fontenla-Romero, and A. Alonso-Betanzos. A very fast learning method for neural networks based on sensitivity analysis. *The Journal of Machine Learning Research*, 7:1159–1182, 2006.
- [5] S.M. Weiss and C.A. Kulikowski. *Computer systems that learn: classification and prediction methods from statistics, neural nets, machine learning, and expert systems*. Morgan Kaufmann, San Francisco, 1991.
- [6] A. Frank and A. Asuncion. UCI machine learning repository, 2010.
- [7] J. Kittler, M. Hatef, R.P.W. Duin, and J. Matas. On combining classifiers. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 20(3):226–239, 2002.
- [8] D.H. Wolpert. Stacked generalization. *Neural networks*, 5(2):241–259, 1992.
- [9] M.F. Møller. A scaled conjugate gradient algorithm for fast supervised learning. *Neural networks*, 6(4):525–533, 1993.