

A distributed approach for classification using distance metrics

L. Morán-Fernández, V. Bolón-Canedo and A. Alonso-Betanzos *

Laboratory for Research and Development in Artificial Intelligence (LIDIA),
Computer Science Dept., Universidade da Coruña, 15071 A Coruña, Spain

Abstract.

To cope with the huge quantity of data that fast development of sensing, networking and inexpensive data storage has come, many distributed approaches have been developed during the last years. The main reason is that, when dealing with large datasets, most existing data mining algorithms do not scale well, and their efficiency may significantly deteriorate. Thus, we present a distributed approach by samples in which the original dataset will be divided into several nodes or processors. For classifying a new test sample, first we compute the distance to the data on each node, and then it will be classified by the model learned from the “closest” data. The proposed method has proved to be useful, demonstrating important savings in runtime and satisfactory performance.

1 Introduction

Traditionally, data mining algorithms have been designed to run in a centralized computing environment. However, with the coming age of network technologies, during the last few years many distributed methods have been developed instead of the centralized approaches. Over 2.5 quintillion bytes of data are created daily, produced by everything from photos uploaded to social media websites, to scientific data from surveys and simulations, to the Curiosity rover currently exploring Mars. Big companies, such as Facebook, collects 500 terabytes everyday, including 2.5 billion pieces of content, 2.7 billion “likes”, and 300 million photos. Google processes billion requests per day meanwhile Amazon draws data from 152 million customers’ purchases to help users decide on items to buy. In this scenario, where data size increases beyond capacity, traditional machine learning techniques and, more specifically, data mining algorithms, do not scale well —memory demands and impracticable running times—, damaging performance and efficiency. Besides, data is often distributed across geographical and organizational boundaries, and it is not legal or economic to gather it in a single location. Therefore, a possible solution might be to distribute the data into several nodes, run a data mining algorithm on each partition and then combine the results.

Although distributed learning is a fairly new field, there exist in the literature several works to scale up large datasets. Lazarevic et al. [1] developed a general framework for distributed boosting to integrate efficiently specialized classifiers learned over very large and distributed homogeneous datasets that cannot be merged at a single location. Tsoumakas et al. [2] presented a framework for constructing a global predictive model

*This research has been financially supported in part by the Spanish Ministerio de Economía y Competitividad (research project TIN2015-65069-C2-1-R), by European Union FEDER funds and by the Consellería de Industria of the Xunta de Galicia (research project GRC2014/035).

from local classifiers which scales up efficiently and achieves high predictive accuracy with respect to a large number of distributed datasets. Zhang et al. [3] proposed a general framework for designing online learning algorithms for large-scale distributed data mining algorithms, which can flexibly balance among computational complexity, communication cost and prediction accuracy.

This paper will be focused on the most common approach to distribute the data: horizontal partitioning (i.e. by samples). We will present a methodology in which the training dataset will be divided into several nodes. Then, each classification algorithm is trained on each of these subsets of data, generating a corresponding model. At this point, the test samples will be classified by one of the generated models according to a certain metric distance. This distributed approach is proposed primarily for being able to learn in a situation in which a huge dataset is partitioned into several disjoint nodes for a more efficient analysis, although it can also be used for learning from several disjoint data locations when the data cannot be merged together.

2 Distributed learning approach based on distance measures

The methodology that we propose in this work consists of a distributed framework for classification by partitioning the data horizontally, that is, by samples. In the horizontal partitioning, the dataset is divided into several nodes that have the same features as the original dataset, each containing a subset of the original samples. Basically, our proposed framework can be summarized in the three following stages: (1) partition of the training dataset into several nodes, (2) calculate the distance between the test samples and the different nodes and (3) classify each test sample using the model generated by the partition with the minimum distance to that test sample.

The first step of our methodology is to randomly partition horizontally the training dataset into a number of disjoint nodes. Since the partition is being made by samples, it is necessary to bear in mind that random distributions of the data might imply that some of the classes are not represented exhaustively in all nodes. To solve this problem, we divide the data maintaining the original class proportions in the training dataset. Using the above framework, as we divide the dataset into p disjoint subsets, the memory requirement will decrease by a factor of $1/p$, which is a substantial reduction. Then, for each test sample, we calculate the distance between these test instances and the different nodes in which the original dataset has been divided. Among the existing distance metrics, we used Hamming and Hellinger since they measure differences between distributions, described as follows.

- The **Hamming distance** (d_{Ham}) between two vectors is the number of coefficients in which they differ [4]. Given an $m \times n$ -by- n data matrix X , which is treated as $m \times (1\text{-by-}n)$ row vectors x_1, x_2, \dots, x_{mx} , and $m \times n$ -by- n data matrix Y , which is treated as $m \times (1\text{-by-}n)$ row vectors y_1, y_2, \dots, y_{my} , the distance between the vector x_s and y_t is defined as follows:

$$d_{Ham} = (\#(x_{sj} \neq y_{tj})/n)$$

Rows of X and Y correspond to samples, and columns correspond to features.

- The **Hellinger distance** (d_{Hel}) is a measure of distributional divergence [5]. Recently, Hellinger metric has been receiving significant attention in the machine learning area since it detects failures in classifier performance due to shifts in data distributions [6]. It can be expressed as follows, where x_s and x_t are discrete probability distributions.

$$d_{Hel}^2 = \frac{1}{2} \sum_{j=1}^n (\sqrt{x_{sj}} - \sqrt{y_{tj}})^2$$

At this point, we classify each test sample using the model generated by the partition with the minimum distance to that test instance. If several partitions obtain the same distance to a test sample, we take the least complex partition. In order to calculate the complexity of the partitions, we use the data complexity measure Fisher's multiple discriminant ratio for C classes [7], which is defined as:

$$F1 = \frac{\sum_{i=1, j=1, i \neq j}^c p_i p_j (\mu_i - \mu_j)^2}{\sum_{i=1}^c p_i \sigma_i^2}$$

where μ_i , σ_i^2 y p_i are the mean, variance, and proportion of the i th class, respectively. In practice the inverse of F1 is preferred, such that a small complexity value represents an easy problem. The distributed method is applied to the training dataset in several rounds or iterations (10 in this work). Repeating the process in several rounds ensures that we have gathered enough information.

An example. Imagine we divide the original dataset in 5 partitions (see Figure 1). Each partition will have the same f features as the primary dataset and $m/5$ samples, where m is the total number of instances of the original dataset. Then, the distance (Hamming or Hellinger) between each test sample and the 5 partitions is calculated. In this example, the minimum distance is obtained by the partitions 1 and 3. Therefore, the model generated by the third partition will be used to classify this test sample, since it presents the lowest complexity value $F1$.

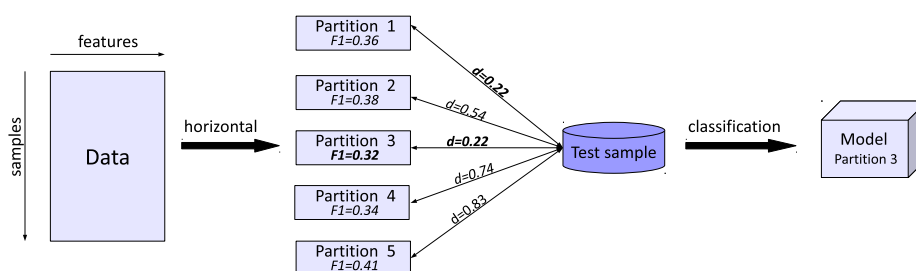


Fig. 1: An example of the proposed distributed methodology.

3 Experiments

In this section we present and discuss experimental results in terms of classification accuracy and runtime over five datasets, comparing the centralized standard approach

and the distributed method proposed in this work.

3.1 Setup

In order to evaluate empirically the proposed distributed framework, we used five datasets, described in Table 1 in terms of the number of training and test samples, the number of features, the number of classes and the imbalance ratio (IR), defined as the number of samples in the majority class divided by the number of samples in the minority class, where a high score indicated a highly unbalance dataset. The number of nodes in which each training dataset is partitioned was calculated trying to maintain a proportion between the number of samples and the number of features, with the constraint of dividing each dataset into at least three nodes.

Dataset	#Samples		#Features	#Classes	IR	#Nodes	Download
	Training	Test					
Connect4	45,038	22,519	42	3	7.95	45	[8]
Isolet	6238	1559	617	26	1.01	5	[8]
Madelon	1600	800	500	2	1.03	3	[8]
Ozone	1691	845	72	2	26.27	11	[8]
Spambase	3067	1534	57	2	1.45	5	[8]

Table 1: Characteristics of the datasets used in the experimental study.

Two classifiers were chosen to evaluate the performance of the framework: k -Nearest Neighbor (k NN) [9] and Support Vector Machine (SVM) using a linear kernel [10], since they are based on distances. Both classifiers are executed in the Weka [11] environment, using default values for the parameters.

3.2 Results

Table 2 displays the classification accuracy obtained by k NN and SVM classifiers both with the centralized and distributed approaches. For the distributed approaches, the average of the ten repetitions is shown, as well as the standard deviation. Statistical tests are not provided because the centralized approach has been validated using the hold-out technique. Notice that, since we have the datasets divided into training and test sets, the centralized approach obtains always the same result. In the case of the distributed approaches, the accuracy varies because on each repetition the partition of the training dataset. As can be seen, the results are very variable depending on the dataset. For some datasets (Madelon and Spambase) the highest accuracies are achieved by centralized approaches. For other datasets (Connect4, Isolet and Ozone), the distributed approach maintains or even improves in some cases with regard to the standard centralized process. However, the important conclusion is that by distributing the data there is not a significant degradation in accuracy.

In terms of runtime, we have to take into account the time required by the classification algorithms for both centralized and distributed approaches. Notice that in the distributed approach, considering that all the partitions can be processed at the same time, the time is the maximum of the times required by the classifier in each subset

		Connect4	Isolet	Madelon	Ozone	Spambase
kNN	C	57.97	54.30	57.40	95.62	78.46
	D-Hamming	54.55±1.72	57.88±0.65	56.58±1.82	95.52±0.28	73.85±1.17
	D-Hellinger	50.35±2.07	55.43±0.52	57.17±1.62	95.29±0.54	73.61±1.36
SVM	C	60.42	79.03	58.88	98.70	84.69
	D-Hamming	61.83±0.12	79.32±0.47	57.49±2.90	98.71±0.07	81.88±1.49
	D-Hellinger	54.40±0.16	79.18±0.39	55.62±1.14	98.60±0.14	81.82±0.99

Table 2: Test classification accuracy. C stands for centralized approach, while D refers to the distributed approaches.

generated in the partitioning stage. In these experiments, all the subsets were processed in the same machine, but the proposed algorithm could be executed in multiple processors. Table 3 shows the speed-up values, which indicate the performance improvement of the distributed approach with respect to the runtime of the centralized approach. For the distributed approaches, the runtime is the sum of the time required by the classifier, the time for calculating the distances (Hamming or Hellinger) and the partition complexity. As expected, the time required by the distributed methods is drastically reduced for all the datasets. It is worth mentioning the important reductions when the proportion between samples and features of the dataset grows (see Connect4).

		Connect4	Isolet	Madelon	Ozone	Spambase
kNN	Hamming	14.87	8.65	3.36	6.24	5.42
	Hellinger	13.78	6.60	3.38	6.56	5.77
SVM	Hamming	19.17	5.18	4.66	7.20	6.53
	Hellinger	17.87	4.49	4.56	7.38	6.82

Table 3: Datasets speed-up. Comparison of centralized and distributed approaches.

In light of the results, we can conclude that the distributed approaches performed successfully, since the runtime was considerably reduced while classification accuracy did not drop to inadmissible values.

Study case. Figure 2 shows the accuracy and the computational cost over the Madelon dataset using 1-5 nodes. We can observe that the runtime decreases drastically using the distributed approach (SVM classifier with Hamming distance) over the standard method (1 node), but there is no significant reduction in the processing times using 2-5 nodes. According to classification accuracy, three nodes achieved the best performance among the distributed methods (only 1.39% inferior to the centralized approach).

4 Conclusions and future work

In the era of Big Data, as data sizes grow and models become more complex, it is natural to consider replacing centralized classification by distributed methods, as a way to reduce computational and memory costs. In this work, we present a methodology for distributed classification over horizontally partitioned data. First, the original dataset

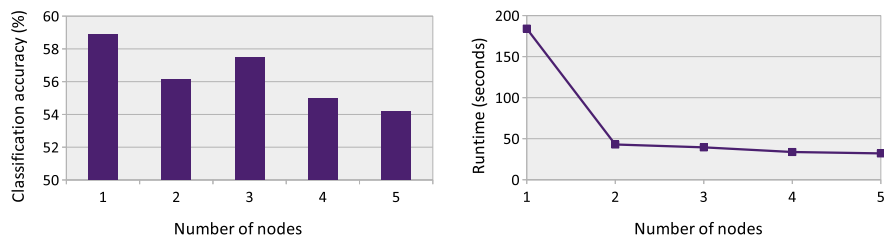


Fig. 2: Classification accuracy and runtime over the Madelon dataset.

will be divided in several nodes. Then, we calculate the distance —using the Hamming and Hellinger measures— between the test samples and the different partitions in which the original datasets has been divided. Finally, each test sample will be classified using the model generated by the partition with the minimum distance to that test instance. From the experiments carried out in five datasets, we can conclude that the novel procedure has demonstrated important savings in runtime and satisfactory performance. With respect to runtime, the behavior is excellent, being this fact the most important advantage of our method. As future research, we plan to extend this study exploring more data driven ways to split the data, as well as using other distance metrics.

References

- [1] Aleksandar Lazarevic and Zoran Obradovic. The distributed boosting algorithm. In *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 311–316. ACM, 2001.
- [2] Grigorios Tsoumakas and Ioannis Vlahavas. Effective stacking of distributed classifiers. In *Ecai*, volume 2002, pages 340–344, 2002.
- [3] Yu Zhang, Daby Sow, Deepak Turaga, and Mihaela van der Schaar. A fast online learning algorithm for distributed mining of bigdata. *ACM SIGMETRICS Performance Evaluation Review*, 41(4):90–93, 2014.
- [4] G Forney. Generalized minimum distance decoding. *IEEE Transactions on Information Theory*, 12(2):125–131, 1966.
- [5] Ayanendranath Basu, Ian R Harris, and Srabashi Basu. 2 minimum distance estimation: The approach using density-based distances. *Handbook of Statistics*, 15:21–48, 1997.
- [6] David A Cieslak and Nitesh V Chawla. A framework for monitoring classifiers’ performance: when and why failure occurs? *Knowledge and Information Systems*, 18(1):83–108, 2009.
- [7] Tin Kam Ho and Mitra Basu. Complexity measures of supervised classification problems. *IEEE transactions on pattern analysis and machine intelligence*, 24(3):289–300, 2002.
- [8] K. Bache and M. Lichman. UCI machine learning repository. University of California, Irvine, School of Information and Computer Sciences. [Online; accessed September 2016]. <http://archive.ics.uci.edu/ml/>.
- [9] D. W Aha, D. Kibler, and M. K. Albert. Instance-based learning algorithms. *Machine learning*, 6(1):37–66, 1991.
- [10] V. N. Vapnik and V. Vapnik. *Statistical learning theory*, volume 1. Wiley New York, 1998.
- [11] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten. The weka data mining software: an update. *ACM SIGKDD explorations newsletter*, 11(1):10–18, 2009.