

# Distance-based large margin classifier suitable for integrated circuit implementation

L.C.B. Torres, C.L. Castro, F. Coelho, F. Sill Torres and A.P. Braga

A new learning method for classification problems that is suitable for integrated circuit implementation is presented. The method, which outperforms current approaches in many data sets, is based on a structural description of the learning set represented by a planar graph. The final classification function is composed of a hierarchical mixture of local experts, which yields a large margin classifier for the whole learning set. Since it is based only on distance calculations, on-chip learning can also be executed. The method is also appropriate for online and incremental learning, since model parameters are obtained directly from the data set, without need of user interaction for learning.

**Introduction:** Machine learning design usually requires user interaction to set model structure and learning parameters, and may also involve the implementation of complex optimisation algorithms. In addition, learning from data involves trading-off bias and variance [1, 2], which may also require setting regularisation parameters and performing cross-validation on data. These have been regarded as major limitations for the implementation of autonomous learning machines such as artificial neural networks (ANNs) on the integrated circuit (IC) level [3]. In the 1990s and early 2000s, ANN accelerator boards and ICs appeared on the market motivated mainly by the slow convergence rates of learning algorithms. Although ETANN [4], Intel's neural chip launched in 1989, had adaptive synapses, a host computer was needed to run the optimisation algorithms and to provide the user interface. In recent years, a new interest in embedded learning systems has been motivated by the rise of the Internet of Things and by the high throughput of data collection systems. IBM has recently launched a new high-performance neural chip [5], which is a further step towards next generation online adaptive systems.

The old limitations, however, still exist, since current high-performance learning algorithms, which treat learning as a trade-off problem, demand user interaction and complex computations, which are not feasible to be directly implemented in ICs. In this Letter, we present a new learning algorithm that does not depend on user-defined parameters to deal with the bias and variance dilemma [2]. The method takes into consideration only the structure of the data set in order to minimise the training set error and maximise the separation margin between classes. Since only distance computations are required, the new method is particularly suitable for implementation in application specific ICs (ASICs) and field-programmable gate arrays (FPGAs).

**Data set structure and support vectors:** Large margin classifiers became popular in the past two decades after the description of support vector machines (SVMs) [2]. In SVMs, margin maximisation is accomplished in feature space after kernel mapping and linearisation occurs. The resulting quadratic programming problem yielded by minimising the square norm of the weights ( $\|w\|^2 = w^T w$ ) subject to the error function constraint ( $J_e(w) \leq \epsilon$ ) has a single minimum, subject to pre-established kernel and regularisation parameters [2].

In SVM's learning, support vectors (SVs) are an outcome of optimisation. However, the existence of SV is regarded in this Letter as a property of the problem, inherent to margin definition. According to this perspective, the maximum margin solution can be obtained by exploring the statistical properties and geometry of the learning set. This is supported by the demonstration that the SVM's solution can be approximated geometrically by finding the maximum margin separator in relation to the convex hulls of the two classes [6].

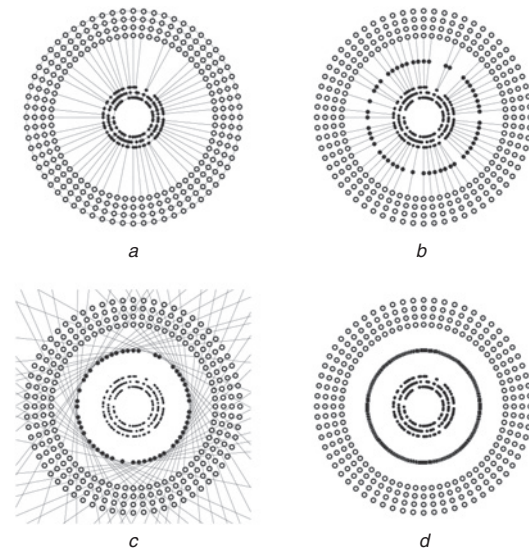
The distance-based large margin classifier presented next is based on a planar graph model of the data set, the Gabriel graph [7]. The resulting discrimination function is also supported by margin vectors, the structural SVs (SSVs), which are analogous to the SVM's SVs, but are obtained here directly from the structure of the data set. As will be shown next, what differs the present method from previous ones is that its parameters are unique and can be obtained directly from the data set structure.

**Structural information extracted from a planar graph:** Structural information is obtained from the corresponding Gabriel graph ( $G_G$ ) [7],

which depends only on pattern-to-pattern distances within the learning set. Considering the data set  $\mathcal{D} = \{(x_i, y_i) | i = 1, \dots, N\}$ , with  $y_i \in \{+1, -1\}$  and  $x_i \in \mathbb{R}^n$ , the graph  $G_G$  of  $\mathcal{D}$  with vertices  $V = \{x_i \in \mathcal{D} | i = 1, \dots, N\}$  has edge  $E$  with vertices  $x_i$  and  $x_j$  if and only if  $\delta(x_i, x_j)^2 \leq [\delta(x_i, x_k)^2 + \delta(x_j, x_k)^2]$ ,  $\forall x_k \in V$  and  $i \neq j \neq k$ , where  $\delta(\cdot, \cdot)$  is the square Euclidean distance operator.

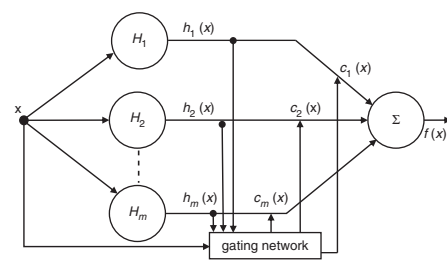
Support edges (SEs) are those edges with vertices  $x_i$  and  $x_j$  that belong to different classes, so they are located in the separation margin between classes. All those vertices  $x_i$  and  $x_j$  that form a SE are selected as SSVs. A local hyperplane  $H_1$  in the middle point of the SE corresponds to a maximum margin classifier in relation to  $x_i$  and  $x_j$ .

**Combination of large margin classifiers:** Local hyperplanes  $H_1$  are not expected to separate all the  $N$  patterns in  $\mathcal{D} = \{(x_i, y_i) | i = 1, \dots, N\}$ , since they are based on a single SE. However, the combination of all hyperplanes yields a large margin classifier that is based on the spatial information of all patterns in  $\mathcal{D}$ , as shown in the example of Fig. 1.



**Fig. 1** Combination of large margin classifiers yielded by SEs

- a  $G_G$  with original data set
- b SEs with middle points
- c Local hyperplanes
- d Final classifier



**Fig. 2** SEs yielding local experts of hierarchical mixture model

Fig. 1a shows the original data set, formed by two concentric classes, with the corresponding  $G_G$ . Fig. 1b shows the data set with the obtained SEs and their middle points represented as dots. Fig. 1c shows all the classifiers corresponding to all middle points of all SEs and Fig. 1d the original data set and a continuous line indicating the projection in the plane of the separation surface due to the final combined classifier. In addition to separating the two classes, the resulting classifier has also maximised the separation margin, as can be seen in Fig. 1. As will be shown next, what differs the present method from previous ones is that its parameters are unique and can be obtained directly from the data set structure.

**Hierarchical mixture of experts:** The final classifier results from hierarchical mixture of experts [2] with a gating network that weights the relevance of each classifier (expert) in the mixture for a given input

pattern  $\mathbf{x}$ . The first layer of the mixture corresponds to the local experts, the classification functions  $h_1(\mathbf{x}), \dots, h_m(\mathbf{x})$ , yielded by all  $m$  individual hyperplanes  $\{H_1, \dots, H_m\}$ . As shown schematically in Fig. 2, the responses of the local experts are weighted by the gating network. Weighting parameters  $c_i(\mathbf{x})$  are also obtained according to distance calculations, as will be shown next.

Let  $H_i$  be a local hyperplane yielded by SE with vertices  $(\mathbf{x}_i, \mathbf{x}_j)$  such that  $y_i = -1$  and  $y_j = +1$ . The classification outcome due to hyperplane  $H_i$  for an arbitrary pattern  $\mathbf{x}$  is given by  $h_i(\mathbf{x}) = \text{sign}(\mathbf{x}^T \mathbf{w}_i - b_i)$ , where  $\text{sign}(\cdot)$  is the sign function and  $\mathbf{w}_i = (\mathbf{x}_i - \mathbf{x}_j)$ , since it corresponds to the hyperplane that is perpendicular to the straight line connecting  $\mathbf{x}_i$  and  $\mathbf{x}_j$ . Similarly, the bias term  $b_i$  is obtained by the expression  $b_i = [(1/2)(\mathbf{x}_i + \mathbf{x}_j)] \mathbf{w}_i^T$ . The middle point of each SE is also obtained directly as  $\mathbf{p}_i = (\mathbf{x}_i + \mathbf{x}_j)/2$ . Therefore, the parameters of the local experts depend only on the vertices  $\mathbf{x}_i$  and  $\mathbf{x}_j$  and are then obtained by direct calculation. The weighting parameter  $c_i(\mathbf{x})$  of  $H_i$  for an arbitrary pattern  $\mathbf{x}$  is calculated according to (1)

$$c_i(\mathbf{x}) = \exp \left( - \frac{\left\{ \max[\delta(\mathbf{x}, \mathbf{p}_k)] \right\}^2}{d(\mathbf{x}, \mathbf{p}_i)} \right), \quad \forall k = 1, \dots, m \quad (1)$$

A normalisation is also imposed so that  $\sum_{i=1}^m c_i(\mathbf{x}) = 1$ . The final classification is obtained as follows:  $f(\mathbf{x}) = \text{sign}[\sum_{i=1}^m h_i(\mathbf{x}) c_i(\mathbf{x})]$ . The coefficients  $c_i(\mathbf{x})$  of (1), which are combined with  $h_i(\mathbf{x})$ , yield larger values to those hyperplanes that are closer to  $\mathbf{x}$ . In practice, for circuit-level implementation, the final classification can be estimated solely by the nearest hyperplane.

**Results:** Experiments were performed with 13 real-world data sets from the UCI repository [8] and two gene expression problems: ‘Golub’ [9] and ‘BcrHess’ [10]. All these data sets have passed through the following preprocessing steps: noise filtering, removal of missing values and attribute rescaling to the range  $\{-1, 1\}$ . To ensure statistical relevance, the experiment was replicated using ten-fold cross-validation. The average performances (AUC) and standard deviations calculated after ten trials are listed in Table 1, along with some characteristics of the data sets in the last column.  $N$  and  $N_d$  correspond to the total amount of patterns and attributes, respectively.

**Table 1:** Results: average AUC and standard deviation. Best results are represented in bold.

Data set	New method	SVM-RBF	SVM-poly	$N/N_d$
<i>Australian Cr.</i>	0.85 ± 0.04	0.86 ± 0.04	<b>0.87 ± 0.04</b>	690/14
<i>Banknote Auth.</i>	0.98 ± 0.03	<b>1 ± 0</b>	<b>1 ± 0</b>	1372/4
<i>BcrHess</i>	<b>0.81 ± 0.12</b>	0.76 ± 0.11	0.77 ± 0.15	133/30
<i>B. Cancer W.P.</i>	0.96 ± 0.03	<b>0.97 ± 0.01</b>	0.96 ± 0.03	683/9
<i>Climate M.S.C.</i>	<b>0.84 ± 0.07</b>	0.53 ± 0.06	0.72 ± 0.11	540/18
<i>Fertility</i>	<b>0.59 ± 0.26</b>	0.5 ± 0	0.5 ± 0	100/9
<i>German Cr.</i>	0.67 ± 0.04	0.66 ± 0.07	<b>0.68 ± 0.05</b>	1000/24
<i>Golub</i>	0.77 ± 0.17	<b>0.8 ± 0.16</b>	0.78 ± 0.17	72/50
<i>Haberman’s S.</i>	<b>0.57 ± 0.09</b>	0.52 ± 0.06	0.5 ± 0.02	306/3
<i>ILPD</i>	<b>0.56 ± 0.09</b>	0.49 ± 0.02	0.5 ± 0	579/10
<i>liver disorders</i>	0.61 ± 0.1	0.67 ± 0.05	<b>0.72 ± 0.07</b>	345/6
<i>P. ind. diabetes</i>	<b>0.72 ± 0.04</b>	0.71 ± 0.05	0.71 ± 0.07	768/8
<i>Parkinsons</i>	<b>0.9 ± 0.15</b>	0.77 ± 0.11	0.81 ± 0.12	195/22
<i>Sonar. M against R.</i>	<b>0.88 ± 0.08</b>	0.84 ± 0.09	0.87 ± 0.08	208/60
<i>Stalog heart</i>	0.8 ± 0.08	<b>0.83 ± 0.07</b>	<b>0.83 ± 0.07</b>	270/13
Average rank	1.87	2.23	1.90	

The SVMs method was used to compare our method. Two configurations of SVM Kernel were tested. First, SVM-RBF with radial basis function, and second SVM-poly with polynomial function. Kernel and regularisation parameters for SVM-RBF and SVM-poly were set via a grid search with ten-fold cross-validation using the packages *Kernelab* and *Caret* available for the language *R* (<http://www.R-project.org/>).

The Friedman’s test [11] was applied assuming the null hypothesis  $H_0$  that all algorithms are equivalent. In our experiment, the statistic  $F_F$  was distributed according to the  $F$ -distribution with two degrees of freedom for the numerator and 28 for the denominator. Considering that the critical value  $F(2, 28)$  is 3.34 for  $\alpha = 0.05$ , the bottom row of Table 1 shows the average ranks  $[R(\mathcal{L})]$  achieved by all algorithms. The corresponding value of  $F_F$  is 0.6, and thus since  $F_F < F(2, 28)$ , the hypothesis  $H_0$

cannot be rejected. It is possible to conclude that our approach has outperformed the other methods for the data sets tested, according to the Friedman’s test.

**Conclusions:** An autonomous learning approach is described in this Letter, which is based on the principle that SEs embody a maximum margin classifier in its middle point, which also minimises the error in relation to its vertices. Individual classifiers are then combined in order to obtain a global large margin classifier that also minimises the error of the learning set. SEs are unique for a given data set and can be obtained directly from the structural information extracted from a planar graph, which is described uniquely by the pairwise distances within the learning set. The resulting classifier does neither require any parameter to be set in advance nor user interaction to balance bias and variance trade-off. In addition, mainly distance calculations are involved, which makes the classifier particularly suitable for circuit-level implementation.

The basic data structure requires an  $N \times (n + 1)$  matrix for storing the learning set and another  $m \times 2n$  matrix containing the parameters of the classifiers. This structure as well as the calculations for distance estimation during learning and classification are perfectly applicable for implementation in ASICs and FPGAs. As demonstrated in numerous works, both yield outstanding performance in speed, power and area for matrix operations with floating point numbers [12]. In addition, the presented method is well suited for autonomous and incremental learning, since the classifier can be updated without great impact on the current graph and hardware structures. In fact, a new learned pattern can be directly aggregated to the current data set without the need to recalculate the whole graph. It will only affect the current separation function if it results in a new SE, which can be simply added to the current set of SEs without disrupting the others.

**Acknowledgments:** The authors thank the CNPq and the CAPES for support.

© The Institution of Engineering and Technology 2015

Submitted: 12 May 2015 E-first: 30 October 2015

doi: 10.1049/el.2015.1644

L.C.B. Torres, C.L. Castro, F. Coelho, F. Sill Torres and A.P. Braga (Department of Electronic Engineering, Federal University of Minas Gerais, Belo Horizonte, Brazil)

✉ E-mail: apbraga@ufmg.br

## References

- Costa, M.A., Braga, A.P., and de Menezes, B.R.: ‘Convergence analysis of sliding mode trajectories in multi-objective neural networks learning’, *Neural Netw.*, 2012, **33**, pp. 21–31
- Haykin, S.: ‘Neural networks: a comprehensive foundation’ (Prentice-Hall, New Jersey, 1999, 2nd edn)
- He, M., Klein, J.O., and Belhaire, E.: ‘Design and electrical simulation of on-chip neural learning based on nanocomponents’, *Electron. Lett.*, 2008, **44**, (9), pp. 575–576
- Holler, M., Tam, S., Castro, H., et al.: ‘An electrically trainable artificial neural network (ETANN) with 10240 “floating gate” synapses’, *IEEE Neural Netw., IJCNN*, 1989, pp. 191–196
- Hsu, J.: ‘IBM’s new brain [News]’, *IEEE Spectrum*, 2014, **51**, (10), pp. 17–19
- Bennett, K.P., and Bredensteiner, E.J.: ‘Duality and geometry in SVM classifiers’. Proc. ICML, Stanford, CA, USA, July 2000, pp. 57–64
- Zhang, W., and Irwin, K.: ‘A study of the relationship between support vector machine and Gabriel graph’, *IEEE Neural Netw., IJCNN*, 2002, **1**, pp. 239–244
- ‘UCI machine learning repository’, 2013. [Online]. Available at <http://www.archive.ics.uci.edu/ml>
- Golub, T.R., Slonim, D.K., Tamayo, P., et al.: ‘Molecular classification of cancer: class discovery and class prediction by gene expression monitoring’, *Science*, 1999, **286**, (5439), pp. 531–537
- Hess, K.R., Anderson, K., Symmans, W.F., et al.: ‘Pharmacogenomic predictor of sensitivity to preoperative chemotherapy with paclitaxel and fluorouracil, doxorubicin, and cyclophosphamide in breast cancer’, *J. Clin. Oncol.*, 2006, **24**, (26), pp. 4236–4244
- Demšar, J.: ‘Statistical comparisons of classifiers over multiple data sets’. JMLR, 2006, vol. 7, pp. 1–30
- Dou, Y., Vassiliadis, S., Kuzmanov, G., et al.: ‘64-bit floating-point FPGA matrix multiplication’. ACM Int. Symp. on Field-programmable Gate Arrays, Monterey, CA, USA, February 2005, pp. 86–95