

# Clustering data streams with weightless neural networks

Douglas O. Cardoso<sup>1</sup>, Priscila M. V. Lima<sup>2</sup>, Massimo De Gregorio<sup>3</sup>, João Gama<sup>4,\*</sup> and Felipe M. G. França<sup>1</sup>

1- PESC/COPPE, Universidade Federal do Rio de Janeiro - Brazil

2- DEMAT/ICE, Universidade Federal Rural do Rio de Janeiro - Brazil

3- Istituto di Cibernetica "Eduardo Caianiello", CNR - Italy

4- LIAAD-INESC, Universidade do Porto - Portugal

**Abstract.** Producing good quality clustering of data streams in real time is a difficult problem, since it is necessary to perform the analysis of data points arriving in a continuous style, with the support of quite limited computational resources. The incremental and evolving nature of the resulting clustering structures must reflect the dynamics of the target data stream. The WiSARD weightless perceptron, and its associated DRASiW extension, are intrinsically capable of, respectively, performing one-shot learning and producing prototypes of the learnt categories. This work introduces a simple generalization of RAM-based neurons in order to explore both weightless neural models in the data stream clustering problem.

## 1 Introduction

Clustering is the process of grouping objects into different groups, such that the common properties of data in each subset is high, and between different subsets is low. The data stream clustering problem is defined as: *to maintain a continuously consistent good clustering of the sequence observed so far, using a small amount of memory and time*. The issues are imposed by the continuous arriving of data points, and the need to analyze them in real time. These characteristics require incremental clustering, maintaining cluster structures that evolve over time. Moreover, the data stream may evolve over time, and new clusters might appear, other disappears, reflecting the dynamics of the stream.

Weightless neural models are networks of artificial neurons based on Random Access Memories (RAMs). The pioneering WiSARD (Wilkie, Stonham & Aleksander's Recognition Device)  $n$ -tuple classifier [5] is a system formed by RAM-based discriminators. A RAM-discriminator consists of a set of  $N$  one-bit word RAMs having  $n$  addresses bits (RAM nodes), and a summing device ( $\Sigma$ ). A RAM-discriminator receives a binary input pattern having  $n \cdot N$  bits ("input retina"), which are connected by means of a pseudo-random mapping to all  $N$  RAM address lines. Each RAM-discriminator is initialized with "0"s at all its contents. Training a discriminator with a particular category is done via writing "1" at the positions addressed by the target input patterns biunivocally transformed by the pseudo-random

---

\* Thanks to the Project Knowledge Discovery from Ubiquitous Data Streams (PTDC/EIA/098355/2008).

mapping. On the other hand, classification by the overall system is performed as follows: as a target pattern is given as input, each RAM-discriminator produces a response to that input, i.e., RAM locations are read and fed to the summing device ( $\Sigma$ ). The set of all discriminator responses are evaluated via a comparison between such responses and a relative confidence  $c$  of the highest response is computed (e.g., the difference  $d$  between the highest response and the second highest response, divided by the highest response). Generalization and noise tolerance capabilities, analogously to the classical artificial neural networks based on synaptic strength, are found in RAM-based neural systems. An introduction to the most important weightless neural networks can be found in [10].

This work presents a preliminary exploration of the WiSARD weightless neural model in the problem of clustering data streams. The agility of the WiSARD is explored on the training of RAM discriminators with data points presented in a sequential style. It is shown how a new DRASiW (associated training view of RAM discriminators) policy is able to reproduce the dynamics of cluster transformations in real time. Synthetic data streams from the MOA framework [11] are used to demonstrate the novel approach.

The remainder of this paper is organized as follows. Section 2 presents a brief overview about data stream clustering. The WiSARD weightless neural model and DRASiW, the associated training view architecture, are explained in Section 3. Section 4 presents WiSARD and DRASiW setups in the experimentation with data stream clustering. Section 5 contains conclusion and final remarks.

## 2 Clustering data streams

Major clustering approaches in data stream cluster analysis include: Partitioning algorithms: construct a partition of a set of objects into  $k$  clusters that minimize some objective function (e.g., the sum of squares distances to the representative centroid). Examples include  $k$ -means [1], and  $k$ -medoids [4]; Micro-clustering algorithms: divide the clustering process into two phases, where the first phase is online and summarizes the data stream in local models (micro-clusters) and the second phase generates a global cluster model from the micro-clusters. Examples of these algorithms include BIRCH [3] and CluStream [2].

A powerful idea in clustering from data streams is the concept of *cluster feature* — CF. A cluster feature, or *micro-cluster*, is a compact representation of a set of points. A CF structure is a triple  $(P, LS, SS)$ , used to store the sufficient statistics of a set of points:  $P$  is the number of data points;  $LS$  is a vector, of the same dimension of data points, that store the linear sum of the  $P$  points;  $SS$  is a vector, of the same dimension of data points, that store the square sum of the  $P$  points.

The properties of cluster features are:

- Incrementality: If a data point  $x$  is added to the cluster, the sufficient statistics are updated as follows:  $LS = LS + x$ ,  $SS = SS_A + x^2$ ,  $P = P + 1$ ;
- Additivity: If two clusters,  $A$  and  $B$ , are merged, the sufficient statistics of the resulting cluster  $C$  are: if  $A_1$  and  $A_2$  are disjoint sets, merging them is equal to the sum of their parts. The additive property allows us to merge sub-clusters incrementally.

A CF entry has sufficient information to calculate the norms  $L1$  and  $L2$ .

### 3 DRASiW: reflections of WiSARD's discriminators

The information stored by the RAM during the training phase is used to deal with previous unseen patterns. When one of these is given as input, the RAM memory contents addressed by the input pattern are read and summed by  $\Sigma$ . The number  $r$  thus obtained, which is called the discriminator response, is equal to the number of RAMs that output "1";  $r$  reaches the maximum value  $N$  if the input pattern belongs to the training set. Intermediate values of  $r$  express a kind of "similarity measure" of the input pattern with respect to the patterns in the training set. The WiSARD's classification behavior is modulated by the size of  $n$  (number of address lines of each RAM node): lower values of  $n$  favor generalization skills, while higher values of  $n$  makes the WiSARD to respond with higher specificity (see [5]).

DRASiW is an extension to the WiSARD model provided with the ability of producing pattern examples, or prototypes, derived from learned categories [7][8]. RAM-discriminators are modified in what their memory locations may hold and, correspondingly, in their training algorithm. Similarly to PLN nodes, introduced by Aleksander [6], such change allows one to store  $h$ -bit words in memory locations, and such can be exploited in the generation of "mental" images of learned pattern categories, i.e., to be able to produce prototypes [9].

The training algorithm of RAM-discriminators is generalized in the following way: memory location contents that are addressed by input patterns are incremented (+1). At the end of the training phase, values at the memory contents will vary between 0 and  $Y$  (where  $Y$  is the number of training patterns). The bidirectional behavior one wants to obtain from a RAM-discriminator  $D$  must satisfy the following conditions: (a) in one direction (WiSARD),  $D$  has to perform the usual classification process of RAM-discriminators; (b) in the opposite direction (DRASiW),  $D$  has to provide, given the name of class  $C$  as input, an example of  $C$ . The solution herein outlined involves the construction of grey level (rather than black and white) images (see Fig. 1), in an internal retina having the same dimensions of the input field, by exploiting the information held in the modified RAM memory locations.



Fig. 1: "Mental" images produced by DRASiW in different applications: [7][8][9].

A new use of the mental image formation process was introduced in [9]: improving the classification abilities of the WiSARD model. It is possible that some of the "mental" images produced by DRASiW are contour-saturated images, predominantly composed by dark grey/black pixels (right picture in Fig. 1). In such a scenario, WiSARD's discriminators generate ambiguous responses, i.e., draws between true and false winner responses. This is due to the saturation that (i) is quickly reached for a training set with a relevant number of class examples and (ii) occurs in an inverted order of the size of RAM neurons: smaller RAM neurons get saturated (most RAM positions written) sooner as the number of patterns used in the training phase increases.

By taking advantage of DRASiW's prototype generation capability, one can avoid ambiguous discriminator responses. Consider the introduction of an integer variable threshold  $b$ ,  $b \geq 1$ , over all RAM neurons contents at all discriminators. At the start of a pattern test,  $b=1$  and, if one observes a draw between discriminator responses,  $b$  is incremented and the output of  $\Sigma$  units are re-calculated taking into account only RAM neuron contents above  $b$ . A straightforward convergence policy, called *bleaching*, is to have  $b$  incremented until just one of the discriminators producing a winner response. Notice that this process is directly related to the way pattern examples are produced from "mental" images in the DRASiW internal retina. Furthermore, based on the "mental" image threshold idea to generate a prototype, it was shown in [9] that, by re-evaluating the set of RAM-discriminator responses upon the detection of a draw, higher rates of correct classification could be obtained.

#### 4 DRASiW-generated clusters

In order to deal with mutable data clusters formed from a target data stream, the WiSARD/DRASiW approach under exploration here assumes the following: (i) a one-to-one association between clusters and discriminators; (ii) a two dimensional synthetic data stream produced by the MOA [11] framework is taken as example for this qualitative study; (iii) a maximum number of coexistent clusters, known *a priori*, where clusters could disappear, but could not appear; (iv) (ii) is a noiseless data set.

A simple generalization related to the contents of the RAM-nodes is introduced in the WiSARD/DRASiW model: instead of a  $h$ -bits counter, each location now holds a list of training time stamps, i.e., when a "1" is written. This list is continuously updated so that each discriminator always has up-to-date information of the cluster it represents. It is now possible to change the way a RAM-node answers to an address read: given an address to be read and a threshold, the node answers "1" iff there is at least one entry greater than the threshold in the address' time stamps list. In other words, the node answers "1" only if the address was written after a given time-threshold.

In order to maintain the order semantics of the input values, a transformation was used to concatenate the Binary Reflected Gray Code (BRGC) bit representation of each feature value. The 0 to  $2^{n+1}/3$  range is used, so the Hamming distance between any two consecutive values is one (1) and between the first and the last range values are maximum. The DRASiW-based data stream clustering approach is performed in two phases: (i) startup phase, when each RAM-discriminator is initialized/trained with points of the cluster it will represent; such examples must be previously clustered and should be as recent as possible; (ii) online phase; for each arriving point a tournament happens to define the discriminator that will learn such example.

Notice that the CF concept is not being explicitly explored in the present architecture. Also noteworthy is the fact that the online phase does not suffer any interruption after startup. These features suggest advanced architectures composed by two concatenated WiSARD/DRASiW network layers; the first layer dedicated to the creation of micro-clusters, and the second layer, fed by the first one, producing the desired data stream clustering. The following pseudo-code describes a tournament carried out during the online phase:

```

candidates = <set of all discriminators>
threshold = 0
Loop
  bestAnswer = -1
  winners = {}
  For each candidate in the candidates' set
    Given the point as input and the threshold, if the
    candidate's answer > bestAnswer
      winners = {current candidate}, an unitary set
      bestAnswer = current candidate's answer
    else, if the candidate's answer = bestAnswer
      add the current candidate to the winners set
  If the winners set is unitary
    Return its only element
  Else, if bestAnswer = 0
    choose randomly one winner element of and return it
  increment threshold
  candidates = winners

```

After processing each arriving point the time window is shifted and all time stamps (RAM-node writing entries) outside it are discarded. In order to produce a DRASiW representation of the ongoing data clusters, a cloud of prototype points is produced in the following way: (i) a single address from each RAM-node  $R$  is

randomly selected according to probability  $s(R_p) / \sum_{i=0}^{2^n-1} s(R_i)$ , where  $s(R_p)$  is the size

of the timestamp list hosted at address  $p$ ; (ii) upon having an address selected from each node, the resulting binary vector is unmapped into a prototype point which would be fully recognized by the discriminator. By generating a population of these prototype points, it is possible to produce the “mental” image of the clusters learnt by the discriminators. A set of data stream clustering dynamics produced by the proposed WiSARD/DRASiW architecture can be seen in [13]. Figure 2 presents a snapshot of a dynamics of four clusters from a synthetic data stream (from MOA [11]).

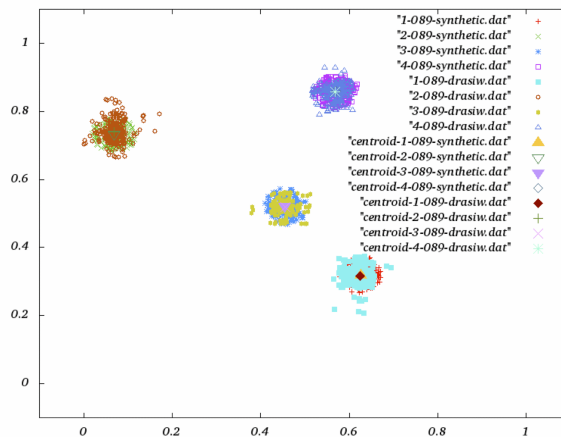


Fig. 2: Synthetic data stream superposed to cluster prototypes produced by DRASiW.

## 5 Conclusions

A qualitative exploration of both WiSARD and DRASiW weightless neural paradigms in the problem of clustering data streams was presented. The potential agility of one-shot training and the quite reduced amount of required memory of the proposed RAM-based neural architecture are compatible with data stream clustering requisites. However, neither a quantitative approach, via dealing with realistic data, nor a fully exploration of the clustering dynamics, e.g., cluster creation, were investigated (this was not exercised since no policies for the creation (or recycling) of discriminators were proposed). Such alternative is left for future work in which the AUTOWiSARD model [12], an unsupervised version of the WiSARD perceptron, could be employed.

## References

- [1] Farnstrom, F., Lewis, J., e Elkan, C. (2000). Scalability for clustering algorithms revisited. *SIGKDD Explorations*, 2(1):51–57.
- [2] Aggarwal, C., Han, J., Wang, J., e Yu, P. (2003). A framework for clustering evolving data streams. *Proc. of 29th International Conference on Very Large Data Bases*, pages 81–92. Morgan Kaufmann.
- [3] Zhang, T., Ramakrishnan, R., e Livny, M. (1996). Birch: an efficient data clustering method for very large databases. In *Proc. of ACM SIGMOD International Conference on Management of Data*, pages 103–114. ACM Press.
- [4] Guha, S., Meyerson, A., Mishra, N., Motwani, R., e O'Callaghan, L. (2003). Clustering data streams: Theory and practice. *IEEE Transactions on Knowledge and Data Engineering*, 15(3):515–528.
- [5] I. Aleksander, W. Thomas, and P. Bowden, WISARD, a radical new step forward in image recognition, *Sensor Rev.*, 4(3), 120-124, 1984.
- [6] I. Aleksander and W.W. Kan, A Probabilistic Logic Neuron Network for Associative Learning, *IEEE Proceedings of the First Int. Conf. on Neural Networks*, pp. 541-548, 1987.
- [7] M. De Gregorio, On the reversibility of multidiscriminator systems, Technical Report 125/97, Istituto di Cibernetica—CNR, Italy, 1997.
- [8] E. Burattini, M. De Gregorio, G. Tamburrini (2000). Mental imagery in explanation of visual object classification, *Proc. of the 6th Brazilian Symposium on Neural Networks (SBRN 2000)*, Rio de Janeiro, Brazil, pp. 137–143.
- [9] Grieco, B. P. A., Priscila M. V. Lima, Massimo De Gregorio, Felipe M. G. França (2010). Producing pattern examples from "mental" images. *Neurocomputing* 73(7-9):1057–1064.
- [10] I. Aleksander, M. De Gregorio, F. M. G. França, P. M. V. Lima, H. Morton (2009). A brief introduction to weightless neural systems, in: *Proceedings of the 17th European Symposium on Artificial Neural Networks (ESANN 2009)*, Bruges, Belgium, pp. 299–305.
- [11] A. Bifet, G. Holmes, R. Kirkby and B. Pfahringer (2010). MOA: Massive Online Analysis. *Journal of Machine Learning Research*. 11:1601-1604, [JMLR.org](http://www.jmlr.org).
- [12] I. Wickert and F. M. G. França, Validating an unsupervised weightless Perceptron (2002). *Proc. of the 9th international conference on neural information processing (ICONIP 2002)*, v. 2, pages 537-541.
- [13] [http://www.cos.ufrj.br/~dougascardoso/stream\\_wisard/](http://www.cos.ufrj.br/~dougascardoso/stream_wisard/)