

Enhancing Visual Clustering using Adaptive Moving Self-Organizing Maps (AMSOM)

Gerasimos Spanakis and Gerhard Weiss

Department of Data Science and Knowledge Engineering,
Maastricht University,
6200MD, Maastricht, Netherlands
{jerry.spanakis,gerhard.weiss}@maastrichtuniversity.nl
<http://www.maastrichtuniversity.nl/DKE>

Abstract. Recent advancements in computing technology allowed both scientific and business applications to produce large datasets with increasing complexity and dimensionality. Clustering algorithms are useful in analyzing these large datasets but often fall short to provide completely satisfactory results. Integrating clustering and visualization not only yields better clustering results but also leads to a higher degree of confidence in the findings. Self-Organizing Map (SOM) is a neural network model which is used to obtain a topology-preserving mapping from the (usually high dimensional) input/feature space to an output/map space of fewer dimensions (usually two or three in order to facilitate visualization). Neurons in the output space are connected with each other but this structure remains fixed throughout training and learning is achieved through the updating of neuron reference vectors in feature space. Despite the fact that growing variants of SOM overcome the fixed structure limitation, they increase computational cost and also do not allow the removal of a neuron after its introduction. In this paper, a variant of SOM is presented called AMSOM (Adaptive Moving Self-Organizing Map) that on the one hand creates a more flexible structure where neuron positions are dynamically altered during training and on the other hand tackles the drawback of having a predefined grid by allowing neuron addition and/or removal during training. Experimental evaluation on different literature datasets with diverse characteristics improves SOM training performance, leads to a better visualization of the input dataset, and provides a framework for determining the optimal number and structure of neurons as well as the optimal number of clusters.

Keywords: Self-Organizing Maps, Clustering, Visualization, Unsupervised Learning

1 Introduction

Clustering is one of the basic data analysis tasks: It is a process of organizing data into similar groups, without any prior knowledge or training. With the increasing graphics capabilities of the available computers, researchers realized [2, 8] that integrating the visual component into the clustering process helps to improve the effectiveness of automated clustering algorithms. This synthesis of computational clustering methods and interactive visualization techniques not only yields better clustering results but allows

exploration and refinement of the clustering structure. Designing such embedded algorithms is tricky mainly due to two limitations: (a) automated clustering algorithms are sensitive to input parameters and results may significantly vary and (b) large data spaces often have skewed distributions which are difficult to be approximated.

The Self-Organizing Map (SOM) [26] is an unsupervised neural network model which effectively maps high-dimensional data to a low-dimensional space (usually two-dimensional). The low-dimensional space (also called output space) consists of a grid of neurons connected with each other, according to a specific structure (can be hexagonal, rectangular, etc.). This structure allows the topology preservation of input data (i.e., similar input patterns are expected to be mapped to neighboring neurons in the output grid) [23]. By this way, SOM manages to achieve dimensionality reduction, abstraction, clustering and visualization of the input data and this is the reason that it has been applied successfully to many different domains and datasets like financial data [11], speech recognition [24], image classification [29], document clustering [28, 37].

The SOM algorithm raises some issues and problems: (a) SOM's architecture is fixed and predefined in terms of number and arrangement of neurons. In case of largely unknown input data, it is difficult to determine a priori the correct structure that provides satisfactory results. There is some work in this area in order to how to add/remove neurons but none of current approaches adjusts neuron positions on the grid according to training progress. (b) Training a SOM comes with a large computation cost, especially in cases of large datasets and/or large maps. Many epochs might be needed in order for the SOM to converge and the map to reach a final state.

In this paper we study how an extension of the traditional SOM can effectively used for visual clustering and handles both issues described above: First, it allows neurons to change positions during training which provides better visualization and faster training time. Second, number of neurons can be adjusted (neurons can be either added or removed) according to dataset requirements and training progress. Due to this enhanced training scheme, the number of epochs required for training is significantly reduced. The rest of the paper is organized as follows. Section 2 presents background work on SOM, extensions on the traditional algorithm and their limitations. The proposed method is presented in Section 3 while experimental setup is described in Section 4. Finally, Section 5 concludes the paper.

2 Related work

2.1 SOM and Competitive Learning

The Self-Organizing Map (SOM) is a fully connected single-layer linear neural network. The SOM uses a set of neurons, often arranged in a 2-D rectangular or hexagonal grid, to form a discrete topological mapping of an input space, $\mathbf{X} \in R^D$. Input space consists of a set of vectors $\mathbf{x}_j \in R^D$:

$$\mathbf{x}_j = [x_{j1}, x_{j2}, \dots, x_{jD}]^T \quad (1)$$

\mathbf{w}_i is the weight vector associated to neuron i and is a vector of the same dimension (D) of the input space, M is the total number of neurons. Obviously, these weights represent the synaptic connections of each neuron i and can be denoted:

$$\mathbf{w}_i = [w_{i1}, w_{i2}, \dots, w_{iD}]^T \quad (2)$$

The fundamental principle of SOM is the soft competition between the nodes in the output layer; not only the node (winner) but also its neighbors are updated [27].

A SOM architecture can be found in Figure 1.

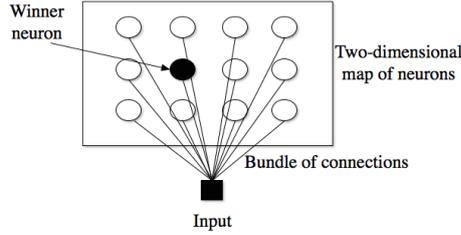


Fig. 1: The SOM (fully connected) architecture.

All the weights $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M$ are initialized to random numbers, in the range of the corresponding input characteristics. We also introduce a discrete time index t such that $\mathbf{x}(t), t = 0, 1, \dots$ is presented to network at time t and $\mathbf{w}_i(t)$ is the weight vector of neuron i computed at time t . The available input vectors are recycled during the training (or learning) process: a single pass over the input data is called an epoch.

On-line Training of SOM In the conventional “on-line” or “flow-through” method, the weight vectors are updated recursively after the presentation of each input vector. As each input vector is presented, the Euclidean distance between the input vector and each weight vector is computed:

$$d_i(t) = \|\mathbf{x}(t) - \mathbf{w}_i(t)\|^2 \quad (3)$$

Next, the winning or best-matching node (denoted by subscript c) is determined by:

$$c = \{i, \min_i d_i(t)\} \quad (4)$$

Note that we suppress the implicit dependence of c on discrete time t . The weight vectors are updated using the following rule:

$$\mathbf{w}_i(t+1) = \mathbf{w}_i(t) + \alpha(t) \cdot h_{ci}(t) \cdot [\mathbf{x}(t) - \mathbf{w}_i(t)] \quad (5)$$

where $\alpha(t)$ is the learning-rate factor and $h_{ci}(t)$ is the neighborhood function. The learning rate factor controls the overall magnitude of the correction to the weight vectors, and is reduced monotonically during the training phase. The neighborhood function controls the extent to which $\mathbf{w}_i(t)$ is allowed to adjust in response to an input most closely resembling $\mathbf{w}_c(t)$ and is typically a decreasing function of the distance on the 2-D lattice between nodes c and i . We use the standard Gaussian neighborhood function:

$$h_{ci}(t) = \exp\left(-\frac{\|\mathbf{r}_i - \mathbf{r}_c\|^2}{\sigma(t)^2}\right) \quad (6)$$

where \mathbf{r}_i and \mathbf{r}_c denote the coordinates of nodes i and c , respectively, on the output space (usually two-dimensional grid). The width $\sigma(t)$ of the neighborhood function decreases during training, from an initial value comparable to the dimension of the lattice to a final value effectively equal to the width of a single cell. It is this procedure which produces the self-organization and topology preserving capabilities of the SOM: presentation of each input vector adjusts the weight vector of the winning node along with those of its topological neighbors to more closely resemble the input vector. The converged weight vectors approximate the input probability distribution function, and can be viewed as prototypes representing the input data.

Batch Training of SOM The SOM update given by Equation (5) is “on-line” in the sense that the weight vectors are updated after the presentation of each input record. In the batch SOM algorithm (proposed in [25]), the weights are updated only at the end of each epoch according to:

$$\mathbf{w}_i(t_f) = \frac{\sum_{t'=t_0}^{t'=t_f} \tilde{h}_{ci}(t') \cdot \mathbf{x}(t')}{\sum_{t'=t_0}^{t'=t_f} \tilde{h}_{ci}(t')} \quad (7)$$

where t_0 and t_f denote the start and finish of the present epoch, respectively, and $\mathbf{w}_i(t_f)$ are the weight vectors computed at the end of the present epoch. Hence, the summations are accumulated during one complete pass over the input data. The winning node at each presentation of new input vector is computed using:

$$\tilde{d}_i(t) = \|\mathbf{x}(t) - \mathbf{w}_i(t_0)\|^2 \quad (8)$$

$$c = \{i, \min_i \tilde{d}_i(t)\} \quad (9)$$

where $\mathbf{w}_i(t_0)$ are the weight vectors computed at the end of the previous epoch. The neighborhood functions $\tilde{h}_{ci}(t)$ are computed using Equation (6), but with the winning nodes determined from Equation (9). This procedure for computing the neighborhood function is identical to the Voronoi partitioning. As is in the on-line method, the width of the neighborhood function decreases monotonically over the training phase.

A more concrete explanation of the batch algorithm is given by the following Equation:

$$\mathbf{w}_i = \frac{\sum_j n_j \cdot h_{ji} \cdot \tilde{x}_j}{\sum_j n_j \cdot h_{ji}} \quad (10)$$

where n_j is the number of input items mapped into node j and the index j runs over the nodes in the neighborhood of node i . The basic idea is that for every node j in the grid, the average \tilde{x}_j of all those input items $x(t)$ is formed that have node j (i.e., vector \mathbf{w}_j) as the closest node. The above Equation is used for updating the node weight

vectors and this is repeated for a few times, always using the same batch of input data items to determine the updated \tilde{x}_j .

The batch SOM offers several advantages over the conventional on-line SOM method. Since the weight updates are not recursive, there is no dependence upon the order in which the input records are presented. In addition to facilitating the development of data-partitioned parallel methods, this also eliminates concerns [31] that input records encountered later in the training sequence may overly influence the final results. The learning rate parameter $\alpha(t)$ does not appear in the batch SOM algorithm, thus eliminating a potential source of poor convergence [9] if this parameter is not properly specified.

The mathematical theory of the SOM is very complicated and only the one-dimensional case has been analyzed completely [15], since the SOM belongs to the ‘ill posed’ problems in mathematics. The SOM can also be looked at as a ‘nonlinear projection’ of the probability density function of high-dimensional input data onto the two-dimensional display.

Usually, the input is mapped onto a 1- or 2-dimensional map. Mapping onto higher dimensions is possible as well, but complicates the visualization. The neurons connected to adjacent neurons by a neighborhood relationship define the structure of the map. The two most common 2-dimensional grids are the hexagonal grid and the rectangular grid and are shown in Figure 2.

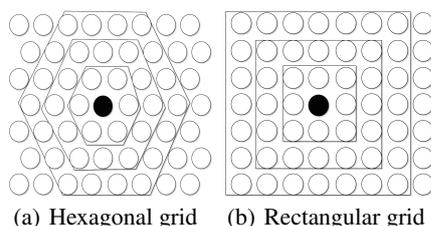


Fig. 2: Typical SOM grids with different neighborhoods around the winner neuron.

The neighborhood function defines the correlation between neurons. The simplest neighborhood function is called *bubble*; it is constant over the neighborhood of the winner neuron and zero otherwise. The neighborhood of different sizes in rectangular and hexagonal maps can be seen in Figure 2. A more flexible definition is the *gaussian neighborhood function* defined by Equation (6).

The number of neurons, the dimensions of the map grid, the map lattice and shape must be specified before training. The more neurons the grid has, the more flexible the mapping becomes but the computation complexity of the training phase increases as well. The choice of the map structure and size is both related to the type of problem and the subjective choice of the user.

2.2 Flexible Structure in Neural Networks and SOM

The norm in artificial neural nets is that classic techniques involve simple and often fixed network topologies trained via stimulus-based methods such as backpropagation. However, there are cases in which the structural design of the network is strongly influenced by the environment and by utilizing constructive and pruning algorithms. Both these algorithmic categories deliver a network which is gradually adjusted in response to training data. There are many approaches which apply these algorithms in classic neural networks [6, 20, 22, 32, 43].

Also, there are many variations of SOM that allow a more flexible structure of the output map which can be divided into two categories: In the first type, we include growing grid (GG) [17], incremental GG [5], growing SOM (GSOM) [1] all coming with different variants. GG is the only variant which allows growing a new node from the interior of the grid (but this is a whole row or column of nodes). In the rest cases, new nodes are generated by a boundary node, despite the fact that the highest error could have been generated by an internal node. The idea is that the error will be propagated to the exterior to guarantee that growing can only be from the boundaries but this process can lead to a map structure with not perfect topology preservation. Therefore, map size becomes very wide after a limited number of insertions, with some additional nodes, which have no effect. MIGSOM [3] allows a more flexible structure by adding neurons internally and from the boundary but still does not offer the ability to remove neurons if necessary.

In the second type of growing variants, the rectangular grid is replaced with some connected nodes. We denote growing cell structures (GCSs) [16], growing neural gas (GNG) [18] and growing where required [30]. These works add just the necessary nodes at the same time, to fine-tune the optimal map size. Nevertheless, GCS and GNG are facing many difficulties for visualizing high-dimensional data. Visualization in these cases is guaranteed only with low-dimensional data.

Limitations in growing and visualization led to hierarchical variants of the previous model like the Growing Hierarchical SOM (GHSOM) [36]. With GHSOM you can get an idea of the hierarchical structure of the map, but the growing parameter of the map has to be defined beforehand. Other approaches (like TreeGNG [12] or TreeGCS [21]) use dendrograms for representation but due to this tree structure they lose the topological properties.

Disadvantages of these approaches are: (a) the high computational cost due to the fact that structure starts from a very basic architecture and has to grow in order to reach an acceptable structure for the data and (b) the fact that after adding neurons there is not the possibility of removing a neuron if performance is not improving.

3 Expanding the idea of self-organization in neuron locations

During the classic SOM algorithm neuron positions remain unchanged and the grid is fixed from the beginning till the end of the training. This facilitates the process of learning (since neighborhood structure is known beforehand) but is restricting regarding the final result and ways of visualizing it. We propose a different and more flexible

scheme in regard to position vectors \mathbf{r}_i of neurons, which allows a more adaptive form of the neuron grid and acts as an extension to the batch learning algorithm.

Starting from an already grown map size, AMSOM can adapt both its size and structure in order to better represent the data at a specific level of detail. After a specific number of steps, neurons are analyzed to see whether the level of representation is sufficient or adjustments are needed: removal and/or addition of neurons. Initially, connections between neurons are determined based on the grid structure but as training advances, these can change and adjust according to the way that neuron positions are also changed during the process. The algorithm flow is described in Figure 3 and more details about the steps are presented in the following subsections.

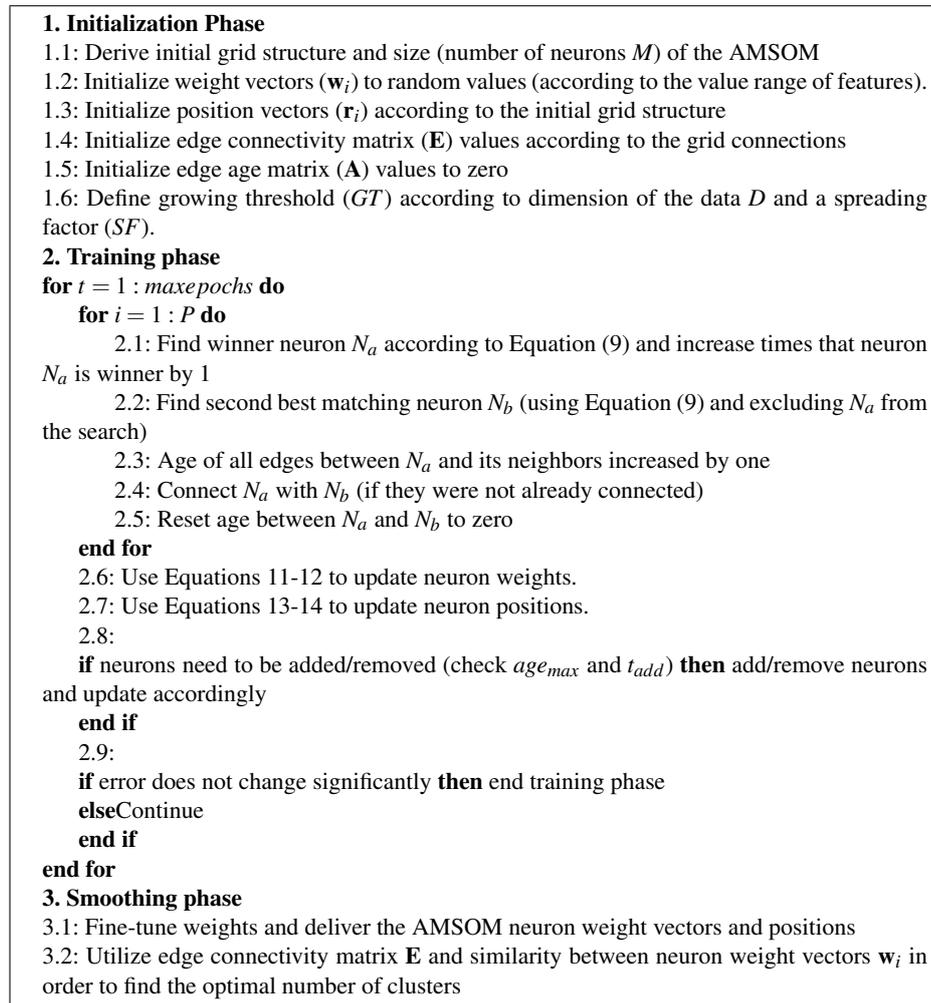


Fig. 3: AMSOM algorithm overview.

3.1 Phase I: AMSOM Initialization

Grid Structure and Size. The first step of AMSOM algorithm is to define the initial grid structure (as the classic SOM). This process facilitates training time in contrast to starting from a small-size structure and building on that as other approaches do [40]. It is also in agreement with the neural development which suggests that nearly all neural cells used through human lifetime have been produced in the first months of life [13]. This overproduction of neuron cells is thought to have evolved as a competitive strategy for the establishment of efficient connectivity [10].

Having this in mind, the initial structure of SOM is determined. Several empirical rules [34] suggest that the number of neurons should be $5 \cdot \sqrt{N}$ where N is the number of patterns in the dataset. In this case, the two largest eigenvalues of the training data are first calculated, then the ratio between side lengths of the map grid is set to the ratio between the two maximum eigenvalues. The actual side lengths are finally set so that their product is close to the number of map units determined according to [40] rule. The eigenvalues ratio shows how well the data is flattened and elongated [14]. At this point a more precise determination of the number of neurons is not essential, since this number will be fine tuned during the training process. Initially, neurons are connected with their neighbors following the idea of Figure 2 using a rectangular or hexagonal grid. For example, if the algorithm suggests that the initial grid of the AMSOM should be 5×4 (let's suppose rectangular), every neuron has 4 neighbors (except the marginal ones). Figure 4 demonstrates two different topologies, a rectangular and a hexagonal one with the corresponding connections between neurons.

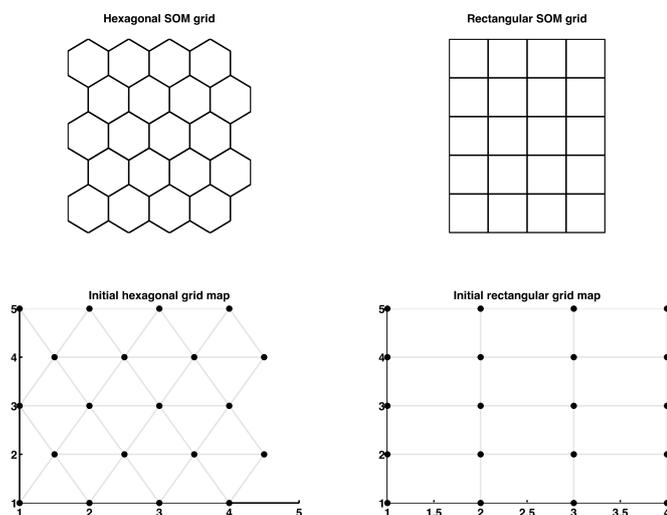


Fig. 4: Initial grid example (hexagonal & rectangular).

Vector, Matrix and Parameters Initialization. For each neuron the following are defined and initialized accordingly:

- Neuron vector (weight vector, \mathbf{w}_i): It is the same as the classic SOM (see Equation (2)) and shows the representation of the neuron in the feature (input) space. Initialization of neuron vectors is random according to literature standards.
- Neuron position (position vector, \mathbf{r}_i): Depending on the output space (mostly it is two-dimensional), it's a vector that shows the position of the neuron. Initial position vectors are equal to the positions of the neurons in the grid, i.e., in Figure 4 one can see the coordinates of neurons according to the structure (hexagonal or rectangular).

Since the structure of the grid is subject to changes during training, we need to keep track of the neighbors of each neuron. There is the possibility that some neurons which were connected in the initial grid become disconnected after some time (or vice versa). In order to keep track of these changes we introduce the orthogonal and symmetrical matrices \mathbf{E} and \mathbf{A} (both size $M \times M$) where $E(p, q)$ shows if neurons p and q are connected (0 translates to no connection, 1 translates to connected neurons) and $A(p, q)$ shows the age of edge (as implied by $E(p, q)$) between neurons p and q : This will be used in order to determine which neurons had incidental connections to other neurons or strong connections as training moves forward. When $A(p, q)$ is 0 that means that neurons p and q were closest neighbors at current epoch but any other value (i.e., 2) implies that neurons p and q were closest neighbors some epochs before (i.e., 2). An example of matrices \mathbf{E} and \mathbf{A} is seen in Figure 5.



Fig. 5: Example of matrices A and E describing connections between AMSOM neurons.

In this example, neurons number (M) is 4 and connectivity matrix E shows how neurons are connected to each other (as implied by the graph). Age matrix A shows how many epochs an edge has ‘survived’: Connection between neuron #1 and #2 has age 2 whereas connection between neuron #2 and #4 has age 0. Notice that age 0 can either mean that neurons are not connected, like neurons #1 and #4 or that neurons are connected at this current epoch (so their connection is ‘recent’), like neurons #2 and #4.

Also, at this stage the growing threshold GT of the map is defined as a function of data dimension (D) and a spread factor (SF) defined by the user. Formula used is $GT = -\ln(D) \times \ln(SF)$ (from [1]). Generally, a SF value of 0.5 always yields good results but its fine tuning is up to the user requirements and the dataset structure.

3.2 Phase II: Training

Weight and Position Updating For the weight learning of neurons, the SOM batch algorithm is utilized, as it was given in Equations 7-10, which are repeated here for clarity.

$$\mathbf{w}_i(t+1) = \frac{\sum_j n_j(t) \cdot h_{ji}(t) \cdot \tilde{\mathbf{x}}_j(t)}{\sum_j n_j(t) \cdot h_{ji}(t)} \quad (11)$$

$$h_{ji}(t) = \exp\left(-\frac{\|\mathbf{r}_j - \mathbf{r}_i\|^2}{\sigma(t)^2}\right) \quad (12)$$

where:

- $\mathbf{w}_i(t+1)$ marks neurons i updated weight (at epoch $t+1$),
- t marks current epoch and $t+1$ marks the next epoch,
- $n_j(t)$ marks the number of patterns that are assigned to neuron j ,
- $h_{ji}(t)$ marks the neighborhood function and is a measure of how close are neuron j and neuron i ,
- $\tilde{\mathbf{x}}_j(t)$ is the mean feature vector of all x that are assigned to neuron j at epoch t ,
- $\mathbf{r}_j, \mathbf{r}_i$ are the position vectors (in the output space) for neurons j and i ,
- $\sigma(t)$ is the adaptation factor, decreasing through training

Building on top of this, at the end of each epoch, the neuron position vectors are adjusted in a similar manner to the SOM training algorithm. In more detail, at the end of each epoch and after the neuron weight vectors update is over, the distances between the neuron vectors (\mathbf{w}_i) are computed. These distances show how close neurons are (in the input space) and can be used as a measure in order to update neuron positions (in the output space). This is achieved through the following Equations:

$$\mathbf{r}_i(t+1) = \mathbf{r}_i(t) + \alpha(t) \cdot \frac{\sum_j n_j(t) \cdot \delta_{ji}(t) (\mathbf{r}_j(t) - \mathbf{r}_i(t))}{\sum_j n_j(t) \cdot \delta_{ji}(t)} \quad (13)$$

$$\delta_{ji}(t) = \exp\left(-\frac{\|\mathbf{w}_j - \mathbf{w}_i\|^2}{\gamma \times \sigma(t)^2}\right) \quad (14)$$

where:

- $t, n_j(t)$ were defined in Equations 11 and 12,
- $\alpha(t)$ denotes the learning rate at epoch t and controls the rate that positions of neurons are moving,
- $\delta_{ji}(t)$ is a neighborhood function denoting how close neurons j and i are (during time t and is based on their distance in the input space (i.e., distance computed based on their vectors w_i),
- γ is a parameter that controls the neighborhood shrinking as a fraction of σ which was used in Equation (12)

Notice the similarity of δ_{ji} with h_{ji} : both are neighborhood functions and are used to determine how close two neurons are but the first one does so using their distances in the feature (input) space while the latter does so using their distances in the output space (map).

Equation (13) will adjust neuron's i position vector according to the neurons which proved winners for more patterns in its neighborhood and less (or even none) according to neurons which were winners for few patterns (or none). This process enhances the concept of neighborhood around the neurons that attract more patterns and also allows to cover any empty spaces in the data representation. It is expected to improve the training speed, since position updating will lead to more accurate position vectors that will be used for the next training epoch and leads to more insightful representations of the neurons in the output space.

Learning rate $\alpha(t)$ can also be set to a small value 0.01 since the neighborhood function controls well the percentage of change in the position vectors. It was selected to update the position vectors with this hybrid on-line-batch SOM rule, due to the fact that output space is much smaller (in most SOM applications) than the input space, so in many cases minor adjustments (than major repositioning of the neurons) are more necessary in order to guarantee satisfactory training but also representation. Also note that the parameter γ which controls neighborhood shrinking for position can also control how fast the map will be updated and how neurons are going to affect each other.

Adding and Removing Neurons During the weight updating process, for each input (pattern) the best matching neuron is determined (N_a) and also the second best matching (N_b). At this step the age of all edges between N_a and its neighbors is increased. Afterwards, N_a is connected to N_b . If both of the neurons were already connected then their age is reset to zero. This is another step that implements the competitive learning rule, since for each new pattern, a new edge connecting the two closest neurons is drawn. This process is repeated for all patterns as they are presented to the AMSOM. Finally, at the end of the epoch for each incident edge between neurons (i, j) , if $A(i, j) \geq age_{max}$, then this edge is removed. age_{max} can be set to a value not small enough (so as to avoid many disconnections) but also not big enough (so as to avoid having a fully connected grid). In our experiments this value was 30. If either of the implicated neurons becomes isolated from the remainder of the structure, then it is removed from the grid. The aim here is to remove edges that are no longer useful because they are replaced by younger edges that are created during the AMSOM training. That is the reason that each time two neurons are connected by an edge, then its age is reset to zero. By this process, neurons that were connected incidentally -especially at the beginning of the training when the map is still under forming- are disconnected after some epochs. This process has two distinct advantages: (a) self-organization and competitive learning will allow after some epochs the removal of redundant number of neurons and (b) adjustment of connections between neurons so as to enhance topological properties of the dataset. An example of a removal of a neuron is shown in Figure 6 along with the necessary adjustments to matrices **A** and **E**.

Also, there is the possibility that after some epochs (t_{add}), new neurons are added. The criterion is based on the training progress and when an addition happens, then

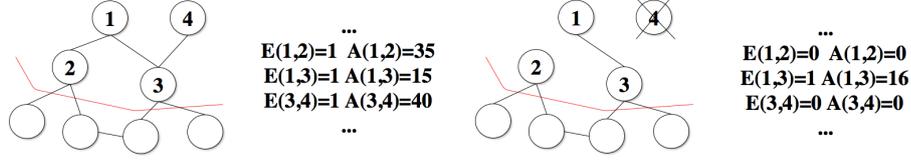


Fig. 6: The process of removing neurons in a part of AMSOM: With age_{max} set to 30, neuron 4 is disconnected from neuron 3 and neuron 1 is disconnected from neuron 2 (notice that matrices A and E are updated accordingly). Neuron 4 is left with no connections so it is removed.

new neurons can be added only after a number of epochs (t_{add}) in order to allow weight adaptation of the map, before evaluating current structure. First step is to spot the neuron N_u with the largest quantization error. A new neuron will be added, if its quantization error is higher than GT , where GT is the growing threshold of the map: A high value for GT will result in less spread out map and a low GT will produce a more spread map. If the quantization error satisfies the above condition then its Voronoi region is considered to be under-represented and therefore a new neuron has to be added to share the load of the high-error-valued neuron.

Regarding the new neuron that will be added, we follow the biological process of ‘cell division’ [33]. By this way the neuron with the highest quantization error is ‘splitted’ to two new neurons (instead of just adding one new neuron somewhere at random with no connections at all). Both new neurons preserve the same connectivity (and also they are connected to each other) with the original neuron, thus we achieve a preservation of behavioral link between the parent and the offspring. Regarding the exact position of the two neurons the following process is followed: Neuron with the largest error among N_u ’s neighbors is spotted (let it be N_v). One neuron will preserve N_u ’s position and the other one will be placed in the middle between N_u and N_v . In detail, weights and positions of the two new neurons (u_1 and u_2) are calculated using the following Equations:

$$\mathbf{w}_{u1} = (1 + \beta) \times \mathbf{w}_u \quad (15)$$

$$\mathbf{w}_{u2} = -\beta \times \mathbf{w}_u \quad (16)$$

$$\mathbf{r}_{u1} = \mathbf{r}_u \quad (17)$$

$$\mathbf{r}_{u2} = \frac{\mathbf{r}_u + \mathbf{r}_v}{2} \quad (18)$$

where \mathbf{w}_u refers to the weight vector of neuron u (neuron that is splitted) and β is a mutation parameter which can take either a fixed or random value according to a certain distribution rule (following [33]). In any case, value of β has to be chosen small in order to avoid a large change both in network topology but also in the weight vectors. In this paper, β takes a random value according to a Gaussian distribution with a mean of zero and variance of one. New neurons retain the same connectivity to other neurons as the parent neuron but age weights are zeroed. The process of adding a new neuron (along with any changes in matrices \mathbf{E} and \mathbf{A}) is described in Figure 7.

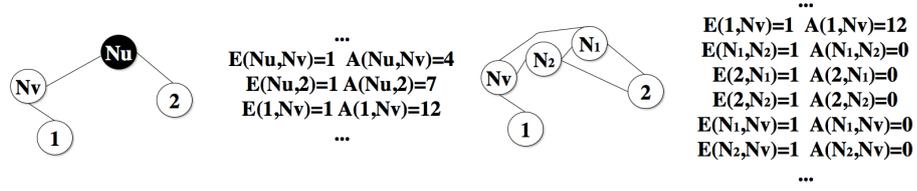


Fig. 7: The process of adding new neurons in a part of AMSOM: N_u is highlighted as the neuron with the highest error and N_v is the neuron among its neighbors with the largest error. Neurons N_1 and N_2 are added instead of N_u , matrices E and A are updated accordingly and weight/position vectors are determined by Equations 15-18.

It has to be pointed out that there is the possibility that a neuron would be removed from a region of the map and to be added in another region (removing and adding neurons are consecutive processes). This comes to agreement with several theories in neural organization, suggesting that cortical regions can adapt to their input sources and are somewhat interchangeable or ‘reusable’ by other modalities, especially in vision- or hearing-impaired subjects [41].

Architecture Adaptation and Termination Criterion As it is described before, initial structure of AMSOM is adapted through learning and training in order to find what is optimal for the number of neurons, their weights and their connections. The adaptation process starts by training the initial structure of AMSOM. When the criteria of adding or removing neurons are satisfied, then the network is adapted. In order to maintain (as possible) the initial structure (i.e., rectangular or hexagonal or any other lattice selected), after this adaptation process we re-evaluate all connections of all neurons and make sure that each neuron has at most Q neighbors (where Q is decided in the beginning, i.e., in the case of rectangular lattice, $Q = 4$): This can be ensured by checking edge matrix E after each epoch and if a neuron is found to have more than Q connections then only the Q -‘recent’ are kept (utilizing age of edges in matrix A). This process is presented in Figure 8.

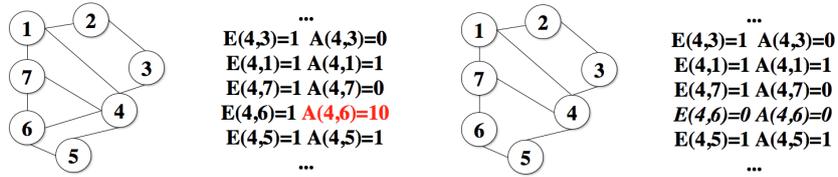


Fig. 8: Maintaining the structure of AMSOM: With $Q = 4$ (i.e., a rectangular grid) neuron 4 is connected to five neurons, so its connection with neuron 6 (oldest connection) is removed.

By this training scheme, AMSOM adapts simultaneously the structure of the map (number of neurons and connections) and the weight vectors. Removing and adding neurons occur when different criteria are met, so they can be applied in any sequence, depending on when the criteria are satisfied. By applying these operations repeatedly, AMSOM is expected to find a near-optimal structure and representation of a given dataset.

Finally, like every SOM algorithm, AMSOM has an upper limit of epochs that training takes place. This number is set to 1000 but there is also a premature termination criterion depending on the mean quantization error change between two consecutive epochs. Thus, if $mqe(t) - mqe(t - 1) < \epsilon_1$ where ϵ_1 is a small value (like $1E - 06$) then the map has reached the desired size (according to the *GT* provided) and training is terminated.

3.3 Phase III: AMSOM Finalization and Clustering of the map

Final phase of AMSOM happens when learning is complete and structure of the network is not any more changing. No neurons are added or removed at this phase and no connections between neurons are added or removed but weight and position vector adaptation is continued with a lower rate. Purpose of this process is to smooth out any quantization error and fine tune weights and positions of the neurons, especially for neurons added at the latter epochs. For this purpose, neighborhood function (both for Equations 11 and 13 is constrained only to the immediate neighborhood and learning rate $\alpha(t)$ in Equation (13) is set to 0.001 (even smaller than in phase II). Phase III is concluded when there is no significant change in change in mean quantization error (i.e., when $mqe(t) - mqe(t - 1) < \epsilon_2$), where ϵ_2 is set to a smaller value than ϵ_1 (like $1E - 10$).

In order to effectively compute the number of clusters that are discovered by the AMSOM algorithm, we take into account the graph structure that is created during training, which is represented by matrix \mathbf{E} . The result is a segmented map that represents the clusters. Literature methods for SOM clustering involve the use of a clustering algorithm like *K*-means or Hierarchical Agglomerative Clustering (HAC) [39] or Spectral Clustering [38] but disadvantage of these approaches is that they require additional (sometimes time consuming) steps after the SOM training and additional parameters to be determined (e.g. number of clusters *K* or a distance measure to be used [7, 19, 42]). These approaches accurately extract the clusters when they are well separated but that is not the case when cluster structure is not that direct. Finally, other approaches utilize graph theoretic methods for partitioning the map but these approaches also rely on extra steps after the end of SOM training [35]. On the other hand, our proposed approach has two distinct advantages over traditional SOM clustering algorithms: (a) The graph construction is inherent to the algorithm, since matrices \mathbf{E} and \mathbf{A} are computed during the training process and (b) involves a simple heuristic process with only one parameter which is automatically determined.

The map to be clustered is represented by an undirected adjacency graph $G(V, \mathbf{E})$; where V represents the set of neurons after the end of training and \mathbf{E} is the edge adjacency matrix as formed during training but we take into account the similarity between

adjacent neurons in terms of the input space (weight vectors). The exact steps of the algorithm are the following:

- Given a trained AMSOM map and matrix \mathbf{E} , compute the distance between any adjacent neurons i and j ($E(i, j) = 1$ for adjacent neurons and $dist(i, j) = \|w_i - w_j\|^2$ as used in Equation 14),
- For each adjacent neuron the edge is considered adjacent when $dist(i, j) \leq v$, where v is a threshold,
- For each edge inconsistency ($dist(i, j) > v$), a null connection is considered in position (i, j) of the graph ($E(i, j)$ is set to 0), otherwise we retain 1,
- A different code is assigned to each connected neuron set.

The threshold v is automatically set to the mean value of distance between all neurons $(i, j) \in V$, so it is not mandatory to be determined, unless there are specific user requirements: a larger value of v will lead to less clusters whereas a smaller value will lead to more clusters. The result is a partitioned map, which indicates the number of clusters. Obviously, after this process the input patterns can be projected into the map (by finding the best matching neuron) and the label of the neuron can be assigned to each input.

4 Experiments

AMSOM performance has been tested with several literature datasets in order to evaluate both map quality (in terms of topology preservation) and the number of epochs needed to converge. Quantization Error (QE) and Topographic Error (TE) were used as intrinsic measures of evaluation (for more details readers are encouraged to read [4]). All (non-textual) datasets were provided by the UCI repository ¹, except the CLUSTER dataset which is a simple and random but large two-dimensional dataset with four groups. In order to test the scalability of the approach for larger datasets, we also tested AMSOM to a large textual dataset ². This dataset involves 4 classes from the RCV1 dataset and counts 9625 documents with 29992 discrete features (which correspond to different words, since we are using the bag-of-words model). All datasets used with their characteristics are presented in Table 1.

Each dataset is shuffled and split to training, testing and validation set (60%, 20% and 20% respectively). Each experiment was performed 20 times and the results presented here are average over these runs (deviations were small and are not presented here). Results for AMSOM QE and TE (compared to classic SOM) along with the number of neurons used by each model are presented in Table 1. From this Table it is obvious that AMSOM's performance is much better than classic SOM. AMSOM starts from the same number of neurons as classic SOM but by removing and adding neurons when necessary reaches a number which is suitable to represent the dataset. Both QE and TE are improved using AMSOM algorithm and this improvement is more significant in TE because of the neuron position changing which allows the map to better adjust to the dataset.

¹ <http://archive.ics.uci.edu/ml/>

² <http://www.cad.zju.edu.cn/home/dengcai/Data/TextData.html>

Table 1: Quality of AMSOM compared to classic SOM and number of neurons for different datasets

Dataset name	Characteristics			QE		TE		# of neurons	
	Instances	Features	Classes	AMSOM	SOM	AMSOM	SOM	AMSOM	SOM
CLUSTER	1000	3	4	0.108	0.1090	0.028	0.063	121	154
IRIS	150	4	3	0.1047	0.3930	0.009	0.013	40	66
WINE	178	13	3	1.7394	1.8830	0.008	0.017	42	66
IONOSPHERE	351	35	2	2.5697	2.9418	0.0026	0.0057	78	91
CANCER	699	9	2	0.7941	0.9456	0.0145	0.0286	103	132
GLASS	214	10	7	0.9797	1.1178	0.0041	0.0093	43	72
RCV1-4	9625	29992	4	1.6937	2.3251	0.0864	0.1232	431	500

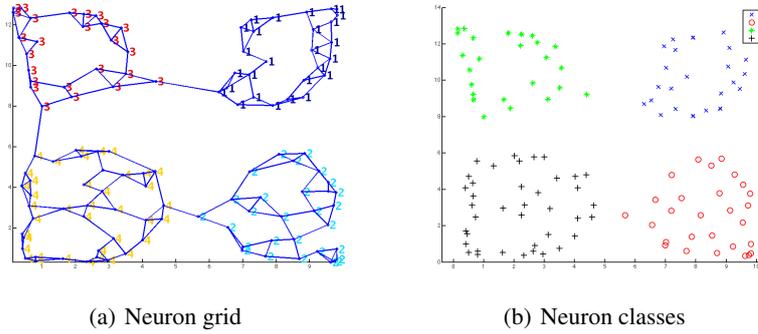


Fig. 9: Visualization results for CLUSTER dataset (4 classes).

Visualization results for four of the small-scale datasets are presented in Figures 9 through 12. In these figures final positions of the neurons and their final positions are represented. For each neuron the process described in Section 3.3 was followed in order to determine the optimal number of clusters but also the class that each neuron belongs to. For the simple CLUSTER dataset it is obvious that the four classes are identified and the grid structure can effectively represent their relations. For the IRIS dataset one class is completely identified whereas the other two (which are more similar) are also highlighted. Also notice that neurons that belong to the same class are mostly connected with each other on the grid and only some spontaneous connections between classes exist. Ability of AMSOM when it comes to multiple class can be seen in GLASS dataset where it successfully manages to highlight all seven classes and separate the corresponding neurons accordingly. Finally, for the more demanding dataset of IONOSPHERE (see the relatively higher QE), AMSOM manages to differentiate in a great degree the two classes. Neuron grids in all figures also reveal that the percentage of dead units (neurons that do not represent any pattern) is significantly small, which is an improvement to the classic SOM algorithm (fewer inactive neurons). Finally, visualization of the results on the RCV1-4 dataset is presented in Figure 13(a). Graph structure is not presented due to the many connections which make the graph not readable in such a format. Moreover, we also conducted a series of experiments in order to

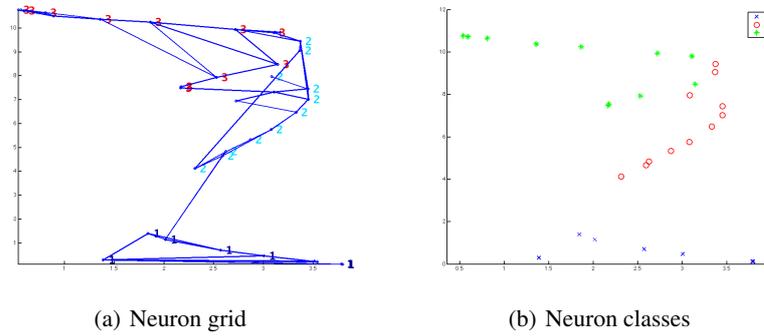


Fig. 10: Visualization results for IRIS dataset (3 classes).

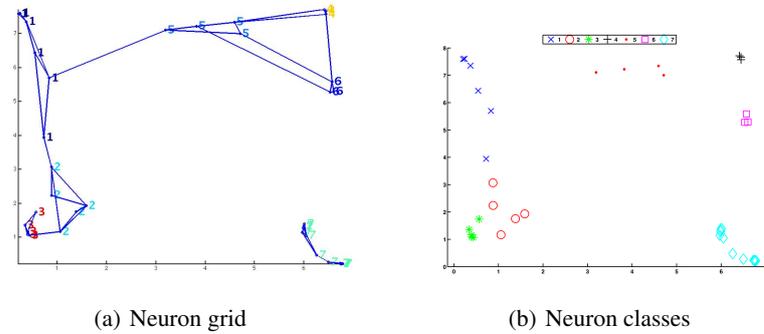


Fig. 11: Visualization results for GLASS dataset (7 classes).

determine handling unbalanced classes. We used uneven number of patterns per class and results were also promising. These can be found in Figure 13(b).

Regarding the spread factor (SF) which controls the growing threshold (GT), a value 0.5 was chosen for this series of experiments because for all datasets it yielded satisfactory results. In the general case that there is no prior knowledge on the data examined, a low value of SF (0-0.3) will allow highlighting of the most significant clusters.

Regarding γ parameter of Equation (13) it was found that it can effectively control the spreading or shrinking of neighborhood during position updating and by this way creating more isolated or more connected clusters. Several experiments were conducted (not presented here due to space limitations) and showed that small values of gamma (1 till 10) produce the best results for all datasets. The higher the γ , the better topographic preservation (reduced TE) but the quantization error (QE) rises. Also, high values of γ tend to increase the number of neurons that remain unused (dead units) whereas values close to 100 tend to approach the classic SOM algorithm (position updating is minimal).

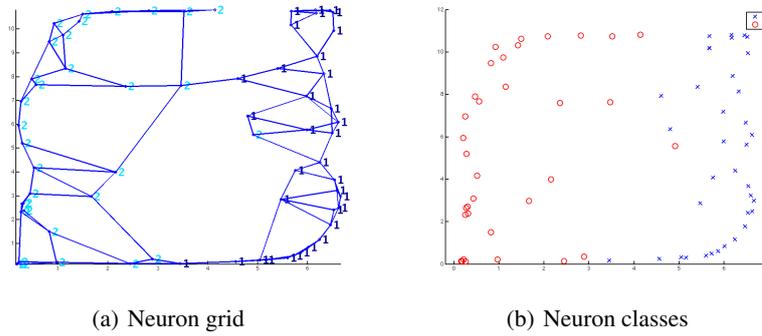


Fig. 12: Visualization results for IONOSPHERE dataset (2 classes).

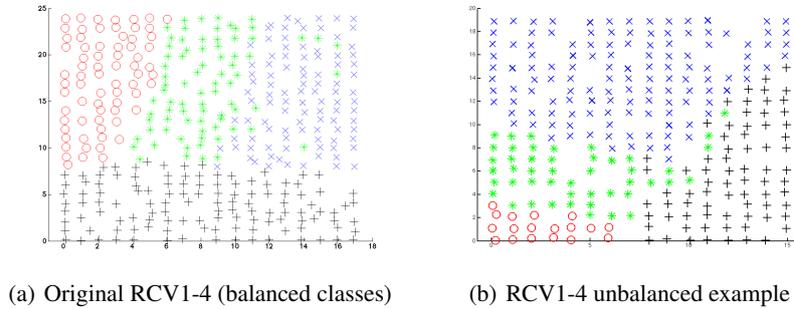


Fig. 13: Visualization results for RCV1-4 dataset (4 classes).

Two more parameters that need to be adjusted are age_{max} and t_{add} . For both parameters, 30 epochs were found to be optimal, which is sound given the fact that 30 epochs is enough time to check whether current structure performs well (reduced QE) or if adjustments are needed (adding/removing neurons).

Complexity of the developed algorithm is slightly increased due to the need for updating matrices \mathbf{A} and \mathbf{E} and also due to the more flexible structure. This overhead is counterbalanced by the faster training process (in all experiments there was a decrease in epochs number around 20%) since updating neuron positions clearly improves training time (requires less epochs) and for memory intensive tasks like the processing of textual datasets algorithm performed fast enough.

5 Conclusion

In this paper we presented AMSOM, a novel algorithm to perform visual clustering which extends SOM competitive learning in the output space. Neuron positions may change during this “double” training process and the number of neurons can be adjusted (addition or removal). These innovations allow a more flexible structure grid

which has the ability to represent the dataset more efficiently. Experimental results on different datasets (with different characteristics in regard to dimensionality, number of latent classes and origin of the data) show improvement in the performance of AMSOM against classic SOM algorithm. AMSOM produces better reference vectors by reducing the Quantization Error, topology is preserved through the neuron moving by significantly reducing the Topographic Error and the visualization result matches as much as possible the original dataset partitions. Also, AMSOM produces fewer nodes with no significant effect while at the same time it reduces required number of training epochs. Finally, AMSOM provides a framework to directly estimate the optimal number of clusters in the dataset with accurate quantitative and qualitative results.

AMSOM provides new insights on how to handle large volumes of otherwise incomprehensible data covering a wide range of human endeavor (science, business, medicine, healthcare, etc.). Obtained results highlight the effective use of competitive learning and self-organization in neural networks and demonstrate that AMSOM can be used with a big variety of datasets. Further work involves work in the following directions: firstly, evaluate the visualization result and its ability to facilitate discovery process, secondly, evaluate the algorithm in even larger datasets and more specifically explore whether statistical properties of the original data are preserved and finally, explore ways to improve performance in terms of time and space requirements.

References

1. Damminda Alahakoon, Saman K Halgamuge, and Bala Srinivasan. Dynamic self-organizing maps with controlled growth for knowledge discovery. *Neural Networks, IEEE Transactions on*, 11(3):601–614, 2000.
2. Gennady Andrienko, Natalia Andrienko, Salvatore Rinzivillo, Mirco Nanni, Dino Pedreschi, and Fosca Giannotti. Interactive visual clustering of large collections of trajectories. In *Visual Analytics Science and Technology, 2009. VAST 2009. IEEE Symposium on*, pages 3–10. IEEE, 2009.
3. Thouraya Ayadi, Tarek M Hamdani, and Adel M Alimi. MIGSOM: multilevel interior growing self-organizing maps for high dimensional data clustering. *Neural processing letters*, 36(3):235–256, 2012.
4. H-U Bauer, Michael Herrmann, and Thomas Villmann. Neural maps and topographic vector quantization. *Neural networks*, 12(4):659–676, 1999.
5. Justine Blackmore and Risto Miikkulainen. Incremental grid growing: encoding high-dimensional structure into a two-dimensional feature map. In *Neural Networks, 1993., IEEE International Conference on*, pages 450–455, 1993.
6. Maria Bortman and Mayer Aladjem. A growing and pruning method for radial basis function networks. *Neural Networks, IEEE Transactions on*, 20(6):1039–1045, 2009.
7. Dominik Brugger, Martin Bogdan, and Wolfgang Rosenstiel. Automatic cluster detection in kohonen’s som. *Neural Networks, IEEE Transactions on*, 19(3):442–459, 2008.
8. Pierrick Bruneau and Benoit Otjacques. An interactive, example-based, visual clustering system. In *2013 17th International Conference on Information Visualisation*, pages 168–173. IEEE, 2013.
9. Michele Ceccarelli, Alfredo Petrosino, and Roberto Vaccaro. Competitive neural networks on message-passing parallel computers. *Concurrency: Practice and Experience*, 5(6):449–470, 1993.

10. Jean-Pierre Changeux and Antoine Danchin. Selective stabilisation of developing synapses as a mechanism for the specification of neuronal networks. *Nature*, 264(5588):705–712, 12 1976.
11. Guido Deboeck and Teuvo Kohonen. *Visual explorations in finance: with self-organizing maps*. Springer Science & Business Media, 2013.
12. Kevin Doherty, Rod Adams, and Neil Davey. TreeGNG-hierarchical topological clustering. In *ESANN*, pages 19–24, 2005.
13. John E Dowling. *The Great Brain Debate: Nature Or Nurture?* Princeton University Press, 2007.
14. Pablo A Estévez, José C Príncipe, and Pablo Zegers. *Advances in Self-Organizing Maps: 9th International Workshop, WSOM 2012 Santiago, Chile, December 12-14, 2012 Proceedings*. Springer Science & Business Media, 2012.
15. Jean-Claude Fort. Soms mathematics. *Neural Networks*, 19(6):812–816, 2006.
16. Bernd Fritzke. Growing cell structures a self-organizing network for unsupervised and supervised learning. *Neural networks*, 7(9):1441–1460, 1994.
17. Bernd Fritzke. Growing grida self-organizing network with constant neighborhood range and adaptation strength. *Neural Processing Letters*, 2(5):9–13, 1995.
18. Bernd Fritzke et al. A growing neural gas network learns topologies. *Advances in neural information processing systems*, 7:625–632, 1995.
19. Maria Halkidi and Michalis Vazirgiannis. A density-based cluster validity approach using multi-representatives. *Pattern Recognition Letters*, 29(6):773–786, 2008.
20. Hong-Gui Han and Jun-Fei Qiao. A structure optimisation algorithm for feedforward neural network construction. *Neurocomputing*, 99:347–357, 2013.
21. Victoria J Hodge and Jim Austin. Hierarchical growing cell structures: TreeGCS. *Knowledge and Data Engineering, IEEE Transactions on*, 13(2):207–218, 2001.
22. Mohammad Islam, Abdul Sattar, Farnaz Amin, Xin Yao, and Kazuyuki Murase. A new adaptive merging and growing algorithm for designing artificial neural networks. *Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on*, 39(3):705–722, 2009.
23. Teuvo Kohonen. Self-organized formation of topologically correct feature maps. *Biological cybernetics*, 43(1):59–69, 1982.
24. Teuvo Kohonen. The 'neural' phonetic typewriter. *Computer*, 21(3):11–22, 1988.
25. Teuvo Kohonen. Things you haven't heard about the Self-Organizing Map. In *Neural Networks, 1993., IEEE International Conference on*, pages 1147–1156. IEEE, 1993.
26. Teuvo Kohonen. *Self-organizing Maps, vol. 30 of Springer Series in Information Sciences*. Springer Berlin, 2001.
27. Teuvo Kohonen. *Self-organization and associative memory*, volume 8. Springer, 2012.
28. Krista Lagus, Timo Honkela, Samuel Kaski, and Teuvo Kohonen. WEBSOM for textual data mining. *Artificial Intelligence Review*, 13(5-6):345–364, 1999.
29. Shin-ye Lu. Pattern classification using self-organizing feature maps. In *1990 IJCNN International Joint Conference on*, pages 471–480, 1990.
30. Stephen Marsland, Jonathan Shapiro, and Ulrich Nehmzow. A self-organising network that grows when required. *Neural Networks*, 15(8):1041–1058, 2002.
31. Filip Mulier and Vladimir Cherkassky. Learning rate schedules for self-organizing maps. In *Pattern Recognition, 1994. Vol. 2-Conference B: Computer Vision & Image Processing., Proceedings of the 12th IAPR International. Conference on*, volume 2, pages 224–228. IEEE, 1994.
32. Pramod L Narasimha, Walter H Delashmit, Michael T Manry, Jiang Li, and Francisco Maldonado. An integrated growing-pruning method for feedforward network training. *Neuro-computing*, 71(13):2831–2847, 2008.
33. Stevan V Odri, Dusan P Petrovacki, and Gordana A Krstonosic. Evolutional development of a multilevel neural network. *Neural Networks*, 6(4):583–595, 1993.

34. Young-Seuk Park, Juliette Tison, Sovan Lek, Jean-Luc Giraudel, Michel Coste, and Francois Delmas. Application of a self-organizing map to select representative species in multivariate analysis: A case study determining diatom distribution patterns across France. *Ecological Informatics*, 1(3):247 – 257, 2006. 4th International Conference on Ecological Informatics.
35. Do Phuc and Mai Xuan Hung. Using som based graph clustering for extracting main ideas from documents. In *Research, Innovation and Vision for the Future, 2008. RIVF 2008. IEEE International Conference on*, pages 209–214. IEEE, 2008.
36. Andreas Rauber, Dieter Merkl, and Michael Dittenbach. The growing hierarchical self-organizing map: exploratory analysis of high-dimensional data. *Neural Networks, IEEE Transactions on*, 13(6):1331–1341, 2002.
37. Gerasimos Spanakis, Georgios Siolas, and Andreas Stafylopatis. DoSO: a document self-organizer. *Journal of Intelligent Information Systems*, 39(3):577–610, 2012.
38. Kadim Taşdemir. Spectral clustering as an automated som segmentation tool. In *Advances in Self-Organizing Maps*, pages 71–78. Springer, 2011.
39. Juha Vesanto and Esa Alhoniemi. Clustering of the self-organizing map. *Neural Networks, IEEE Transactions on*, 11(3):586–600, 2000.
40. Juha Vesanto, Johan Himberg, Esa Alhoniemi, and Juha Parhankangas. *SOM toolbox for Matlab 5*. Citeseer, 2000.
41. Van J Wedeen, Douglas L Rosene, Ruopeng Wang, Guangping Dai, Farzad Mortazavi, Patric Hagmann, Jon H Kaas, and Wen-Yih I Tseng. The geometric structure of the brain fiber pathways. *Science*, 335(6076):1628–1634, 2012.
42. Sitao Wu and Tommy WS Chow. Clustering of the self-organizing map using a clustering validity index based on inter-cluster and intra-cluster density. *Pattern Recognition*, 37(2):175–188, 2004.
43. Shih-Hung Yang and Yon-Ping Chen. An evolutionary constructive and pruning algorithm for artificial neural networks and its prediction applications. *Neurocomputing*, 86:140–149, 2012.