



Growing Cell Structures – A Self-organizing Network for Unsupervised and Supervised Learning*

Bernd Fritzke[†]

TR-93-026

May 1993

Abstract

We present a new self-organizing neural network model having two variants. The first variant performs unsupervised learning and can be used for data visualization, clustering, and vector quantization. The main advantage over existing approaches, e.g., the Kohonen feature map, is the ability of the model to automatically find a suitable network structure and size. This is achieved through a controlled growth process which also includes occasional removal of units. The second variant of the model is a supervised learning method which results from the combination of the abovementioned self-organizing network with the radial basis function (RBF) approach. In this model it is possible - in contrast to earlier approaches - to perform the positioning of the RBF units and the supervised training of the weights in parallel. Therefore, the current classification error can be used to determine where to insert new RBF units. This leads to small networks which generalize very well. Results on the two-spirals benchmark and a vowel classification problem are presented which are better than any results previously published.

*submitted for publication

[†]International Computer Science Institute, Berkeley, CA, fritzke@icsi.berkeley.edu

Contents

1	Introduction	1
2	Unsupervised Growing Cell Structures	1
2.1	Problem Definition	1
2.2	Network Architecture	2
2.3	Network Dynamics	3
2.4	Removal of Cells	9
2.5	Approximation of the Voronoi Regions	9
2.6	Efficient Manipulation of High-dimensional Topologies	11
2.7	Network Visualization for High-dimensional Input Data	12
2.8	Alternative Insertion Criteria	16
3	Extension to Supervised Learning	19
3.1	Motivation	19
3.2	Radial Basis Functions	19
3.3	Supervised Growing Cell Structures	21
3.4	Simulation Examples	24
4	Discussion	31
	Notation	32

1 Introduction

Self-organizing neural network models, as proposed by Willshaw & von der Malsburg (1976) and Kohonen (1982), generate mappings from high-dimensional signal spaces to lower-dimensional topological structures. These mappings are able to preserve neighborhood relations in the input data and have the property to represent regions of high signal density on correspondingly large parts of the topological structure. This makes them interesting for applications in various areas ranging from speech recognition (Kohonen, 1988) and data compression (Schweizer et al., 1991) to combinatorial optimization (Favata & Walker, 1991). The fact that similar mappings can be found at various places in the brains of humans and animals indicates that preservation of topology is an important principle at least in natural “signal processing systems”.

It has been noted that the predetermined structure and size of Kohonen’s model imply limitations on the resulting mappings. A number of variations have been proposed concerning networks with variable topology or variable number of elements. The approach of Jokusch (1990) leads to networks with rather complicated structure. In the minimum-spanning-tree network described by Kangas, Kohonen & Laaksonen (1990) the preservation of neighborhood relations is done only to a small degree due the sparse connectivity of the network. The “neural gas” algorithm of Martinetz & Schulten (1991) seems to produce compact networks which preserve the neighborhood relations extremely good. It generates, however, in general networks with the same dimensionality as the input data so that no dimensionality reduction is performed. Other models allow a variable number of elements, but have predefined principal structure (e.g., rectangular array), namely the interpolative algorithm of Rodrigues & Almeida (1990) and the “learning expectation” method introduced by Xu (1990). A proposal to use random structures stems from Ritter (1991). Recently an interesting approach with a network growing on a grid has been introduced by Blackmore & Miikkulainen (1992).

The network presented in this contribution has a flexible as well as compact structure, a variable number of elements, and a k -dimensional topology whereby k can be arbitrarily chosen. Recently it was demonstrated that the new model improves over Kohonen’s feature map with respect to various important criteria (Fritzke, 1993a). We acknowledge, however, that the new model owes several ideas to Kohonen’s approach and that it is an extension of his work rather than a completely different formalism.

First we outline the network for unsupervised learning and introduce later on the extension of the model to supervised learning.

2 Unsupervised Growing Cell Structures

2.1 Problem Definition

Before we describe our network model, it seems appropriate to exactly define the kind of problems the network is supposed to solve. In the first place, we have a number of n -dimensional input signals obeying an unknown probability distribution $P(\xi)$. With $V = R^n$ we denote the vector space the input signals stem from.

Our objective is to generate a mapping from V onto a discrete k -dimensional topological

$$g) \mathbb{K} = \mathbb{J}$$

Ἐπισημαίνεται ἡ ἀποστολὴ τῆς ἐπιστολῆς.

Ἐπισημαίνεται ἡ ἀποστολὴ τῆς ἐπιστολῆς.

is chosen large at the beginning and decreases then, too. Our model follows the same basic strategy. There are, however, two important differences:

- The adaptation strength is constant over time. Specifically we use constant adaptation parameters ε_b and ε_n for the best matching unit and the neighboring cells, respectively.
- Only the best-matching unit and its direct topological neighbors are adapted.

These choices eliminate the need to define a cooling schedule for any of the model parameters.

In the following N_c denotes the set of direct topological neighbors of a cell c . Furthermore, we define for every cell c a local counter variable τ_c basically containing the number of input signals for which the cell has been best-matching unit. Since the cells are slightly moving around, more recent signals should be weighted stronger than previous ones. This is achieved by decreasing all counter variables by a certain fraction after each adaptation step. To enable this decay, the signal “counters” must be represented by real-valued variables.

An adaptation step in our model can be formulated as follows¹ (see also fig. 3):

1. Choose an input signal ξ according to the probability distribution $P(\xi)$.
2. Locate the best matching unit $s = \phi_w(\xi)$.
3. Increase matching for s and its direct topological neighbors

$$\Delta w_s = \varepsilon_b(\xi - w_s) \quad (3)$$

$$\Delta w_c = \varepsilon_n(\xi - w_c) \quad (\text{for all } c \in N_s) \quad (4)$$

4. Increment the signal counter of s .

$$\Delta \tau_s = 1 \quad (5)$$

5. Decrease all signal counters by a fraction α .

$$\Delta \tau_c = -\alpha \tau_c \quad (\text{for all } i \in A)$$

If we choose small values for ε_b and ε_n , then the cells move from their initial random positions to locations with a dynamic equilibrium between the changes in all directions. They do not stop moving completely since the adaptation parameters are not decreased (so this is not stochastic approximation).

Our objective is a structure with the synaptic vectors w_c distributed according to $P(\xi)$. This is achieved when every cell has the same probability of being best-matching unit for the current input vector. We do not know $P(\xi)$ explicitly, but with the local signal counters we can compute an estimate of $P(\xi)$, namely the relative frequency of input signals received by a certain cell.

The *relative signal frequency* of a cell c is

$$h_c = \tau_c / \sum_{j \in A} \tau_j. \quad (6)$$

¹Here, and throughout the whole paper, $\Delta x = y$ stands for $x^{\text{new}} = x^{\text{old}} + y$. This is to have a concise notation for incremental changes.

a) Initial Situation

Voluntary respiration.

was so far. The grey in

The cell \bar{a} was received the

columns represent signal co

a) Situation before an i

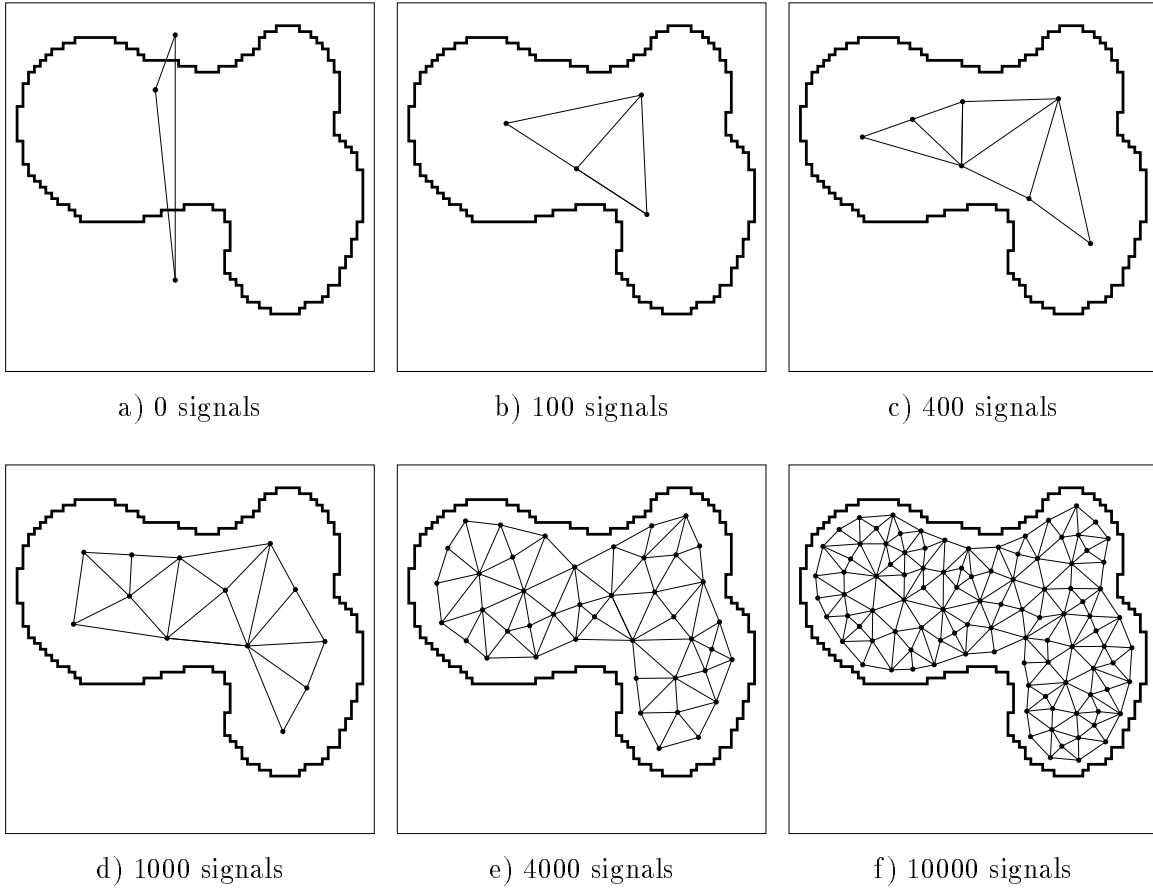
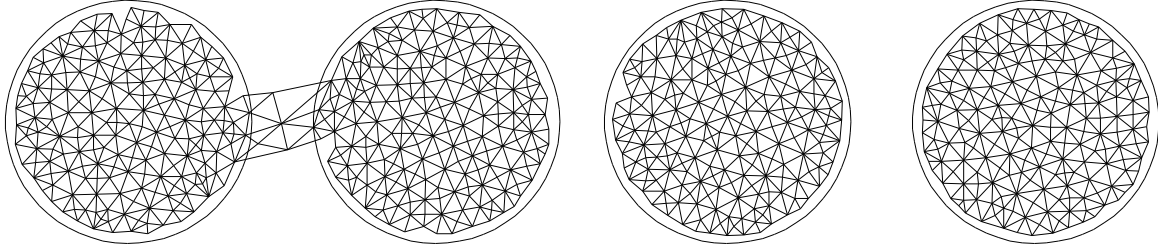


Figure 7: Development of a two-dimensional Growing Cell Structure. The underlying probability distribution $P(\xi)$ is in this case also two-dimensional and is uniform in a cloud-shaped area. Below every sub-picture the number of already received input signals (which is the number of adaptation steps) is shown. The employed simulation parameters are $\lambda = 100$, $\varepsilon_b = 0.06$, $\varepsilon_n = 0.002$, $k = 2$, $\alpha = 0.05$.



a) The growth process leads to a well-adapted network structure. Mostly short connections indicate good topology-preservation. Few synaptic vectors lie outside the relevant circular areas.

b) By removal of *superfluous* cells sub-structures can be formed. The positions of the synaptic vectors now indicate a nearly perfect modeling of the probability distribution and there are only short connections.

Figure 8: A Growing Cell Structures network with 400 cells has adapted to the probability distribution of the previous example. Simulation parameters: $\varepsilon_b = 0.06$, $\varepsilon_n = 0.002$, $\lambda = 100$, $\alpha = 0.05$, and $\eta = 0.09$ (for b only).

is this property which facilitates the extension of the model to a new supervised learning method as will be demonstrated in section 3.

2.4 Removal of Cells

In some cases, especially if $P(\xi)$ consists of several separate regions of positive probability density, a still better modeling can be achieved by removing "superfluous" cells. A cell can be regarded as superfluous if it has a position (synaptic vector) in a region of V with very low probability density. In general, $P(\xi)$ is unknown, but we can relate the relative signal frequency of a cell to the size of its "receptive field" (Voronoi field) to get a local estimate. Specifically, one can note that

$$\tilde{p}_c = h_c / |F_c| \quad (13)$$

is a local estimate of the probability density near w_c . By periodically removing cells with values of \tilde{p} below some threshold η we can model even structured distributions very accurately. Fig. 8 shows the simulation results for such a probability distribution without and with removal of neurons.

2.5 Approximation of the Voronoi Regions

The computation of the Voronoi tessellation is very difficult for dimensions $n > 2$. Thus we replaced the Voronoi region F_c in a first approach by an n -dimensional hypercube with a side length equal to the mean length \bar{l}_c of the edges emanating from c . Instead of $f_c = |F_c|$ we took

$$\tilde{f}_c = (\bar{l}_c)^n, \quad (14)$$

with \bar{l}_c computed by

$$\bar{l}_c = 1/\text{card}(N_c) \sum_{i \in N_c} \|w_c - w_i\|. \quad (15)$$

Although the estimation of the probability density is not anymore as accurate as before, it seems to be sufficient to reliably identify superfluous neurons. However, the appropriate value for the threshold depends strongly on the probability distribution. The reason for that is that a probability distribution which is non-zero in a large area of the input vector space has a lower density than a distribution which is concentrated more locally. But if one uses instead of \tilde{p} the normalized value \hat{p} , defined as

$$\hat{p} = \tilde{p} \sum_{c \in A} \tilde{f}_c, \quad (16)$$

which is computed by multiplying \tilde{p} with the total volume of all hypercubes, one gets sufficient independence. A threshold value of $\eta = 0.09$ is then appropriate in most cases. The check whether a cell i has a value $\hat{p}_i < \eta$ is performed after each insertion, and the cells fulfilling the condition are removed. Also the simulation leading to fig. 8b) has been performed with this method.

One should note, however, that the choice of an n -dimensional hypercube is only appropriate if the underlying data indeed spans the n -dimensional space. If, on the other hand, the data stems from a lower-dimensional subspace of R^n , it might be better to use a hypercube of that dimensionality.

To illustrate the point, let us consider an extreme example. Assume that our input data is 100-dimensional (which is not uncommon for some real problems), but stems from a two-dimensional sub-manifold of R^{100} (what we do not know). We might take a two-dimensional network to be able to visualize the data (see section 2.7). If now for one of our cells the mean edge length shrinks by five percent, then the volume of the corresponding 100-dimensional hypercube collapses to less than 0.6 percent of its previous size. Obviously, this does not reflect very well the change of the “receptive field” of the cell. In this case taking two-dimensional hypercubes would have been more appropriate.

From the above it should be evident that it would be very helpful to know the *true* dimensionality of the data, meaning the smallest dimensionality t , such, that a t -dimensional sub-manifold of V can be found containing all (or most) input data. Then t -dimensional hypercubes could be used to estimate the size of the Voronoi regions in our model. Unfortunately, it is in general difficult to figure out the value of t , especially because the mentioned sub-manifold does not have to be linear but could be arbitrarily twisted (e.g. a curved surface in R^3). Therefore, even a principal-component analysis of the data does in general not reveal their true dimensionality, but gives only (or at least) an upper bound.

As long as there is no simple method to determine the true data dimensionality t , one has to define an estimate \tilde{t} of it. In the following we give some general rules for choosing such an estimate which do work well for all problems we encountered so far:

- Always set $\tilde{t} \leq n$.
- If the different components of the input vectors are known to be stochastically independent from each other use $\tilde{t} = n$.
- If there are known dependencies among the components set \tilde{t} to the number of independent variables.

- Always set \tilde{t} smaller or equal to the number of input vectors. This rule applies only to the rather unusual case that the total number of vectors is smaller than their dimensionality n .
- Finally one can perform a principal component analysis of the data. Then \tilde{t} should be set to the number of principal eigenvalues of the covariance matrix of the data.

However, in most cases it is not necessary to do the principal component analysis since our method is not very sensitive to the choice of \tilde{t} . The only case one should avoid is, to choose a value of \tilde{t} which is *much* too high. This can happen only if the data is very high-dimensional and there are strong dependencies among the components. In such a case it can happen that most insertions occur in one region of the structure. This is due to the fact that a newly inserted cell then gets attributed nearly all the signals of its neighbors because the change of their Voronoi fields is overestimated (see above). A simple remedy for this problem is to choose a lower value for \tilde{t} and to perform another simulation.

In conclusion, we choose in the following in each case a value for \tilde{t} and approximate the volume of the Voronoi regions by

$$\tilde{f}_c = (\bar{l}_c)^{\tilde{t}} \quad (17)$$

with \bar{l}_c being the mean edge length (see eqn. 15). It should be stressed again, however, that the choice of \tilde{t} seems not at all to be a critical step. For a given set of data usually many different estimates work well.

2.6 Efficient Manipulation of High-dimensional Topologies

The implementation of k -dimensional Growing Cell Structures is somewhat more complicated than the implementation of the Kohonen feature map (for which usually a rectangular array of processing units is chosen). Therefore, it seems appropriate to give some hints how this can be done with relatively small effort.

Any implementation of the model must support the two structural update operations:

- Insertion of a neuron
- Deletion of a neuron

These operations have to be performed such that the resulting structure consists exclusively of k -dimensional hypertetrahedrons again.

The general structure of the network can be represented as an undirected graph which is a standard data type consisting of nodes and of edges between pairs of nodes.³ The nodes correspond to neurons and the edges to topological neighborhood relations. Although such a data structure is already sufficient in principle, a considerable search effort is needed to make consistent update operations. The problem is that the removal of a neuron might require that also other neurons and connections are removed to make the structure consistent again. Simple heuristics as, e.g.,

³Our current implementation of the model is based on LEDA (see Mehlhorn & Näher, 1989), a publicly available library of data types and algorithms. LEDA contains in particular a very elaborated data type “graph”.

FIGURE 8: SIMPLE PENT

THE ADJACENT EDGES AND THE

TO BE REMOVED. THIS IS DO

a) GROWING CELL STRUCTURE

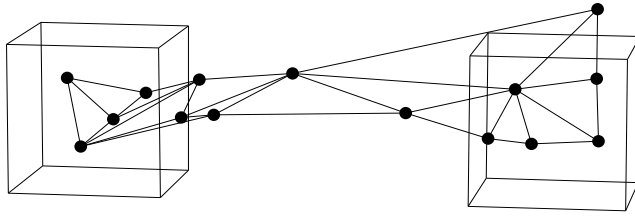
FIGURE 10: COLLECT LEWIS

IN WHICH PARTICIPATES:

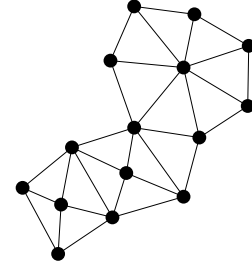
ANDERSON (two-dimensional plane)

AND CONSIDER

a) GEOMETRIC CELL STRUCTURES.



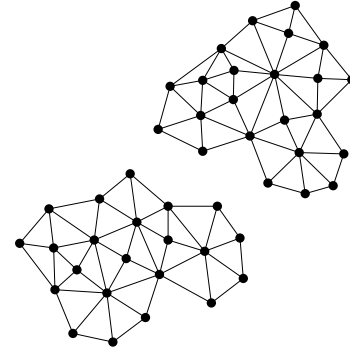
a) Projection into input vector space. Shortly after the start of the simulation all cells are connected since no deletions took place yet.



b) Embedding into the two-dimensional space



c) Projection into input vector space. Deletion of superfluous cells has led to two separate structures.



d) Embedding into the two-dimensional space. The two sub-structures can be recognized easily.

Figure 11: Example for the embedding method. The probability distribution is uniform in two separated cubes. The network is two-dimensional. Figure a) and b) as well as figure c) and d) show the same state of the simulation, respectively. Through the embedding it is easily possible to detect the splitting of the network as can be seen from fig. d).

$$f_b = \begin{cases} 0 & \text{if } 3d < e \\ d/5 & \text{if } 2d < e \leq 3d \\ d/2 & \text{if } d < e \leq 2d \\ d & \text{if } 0 < e \leq d \\ 0 & \text{if } 0 = e \end{cases} \quad (18)$$

Attracting force f_n of two connected discs with center distance e :

$$f_n = \begin{cases} 0 & \text{if } e < d \\ (e - d)/2 & \text{(otherwise)} \end{cases} \quad (19)$$

These two forces have to be balanced against each other. We usually multiplied f_b by 0.2 and f_n by 1.0. In some cases, however, different values might be more appropriate.

An example of the results obtained by the described method is shown in fig. 11. A

animal		d	o	h	u	v	e	n	k	e	l	k	e	x	g	f	d	o	c	g	i	r	h	z	b	e	c	o	a	w
is	small	1	1	1	1	1	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
	medium	0	0	0	0	0	0	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	big	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1
has	2 legs	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	4 legs	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	hair	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	hooves	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1
	mane	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	1	1	1	1	1	0	0	0
	feathers	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
likes to	hunt	0	0	0	0	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0
	run	0	0	0	0	0	0	0	0	1	1	0	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	0	0	0
	fly	1	0	0	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	swim	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

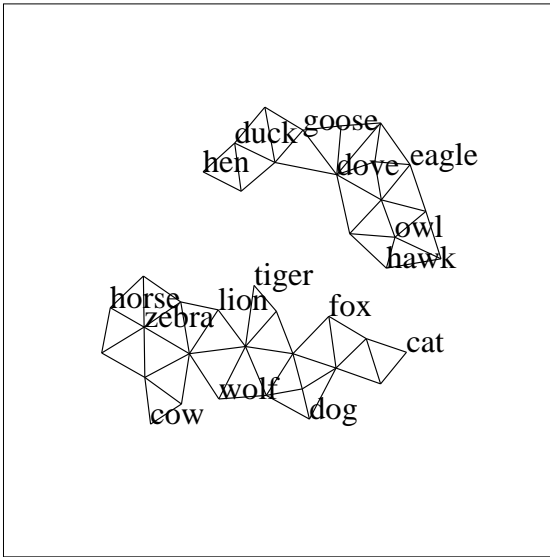
Table 1: Animal names and binary attributes (after Ritter & Kohonen, 1989): If an attribute applies for an animal the corresponding table entry is 1, otherwise 0.

still larger benefit can be gained by having the embedding when the input data is so high-dimensional that we can not visualize the network in input vector space anymore. This is in many real applications the case since often we have data consisting of many more than three components.

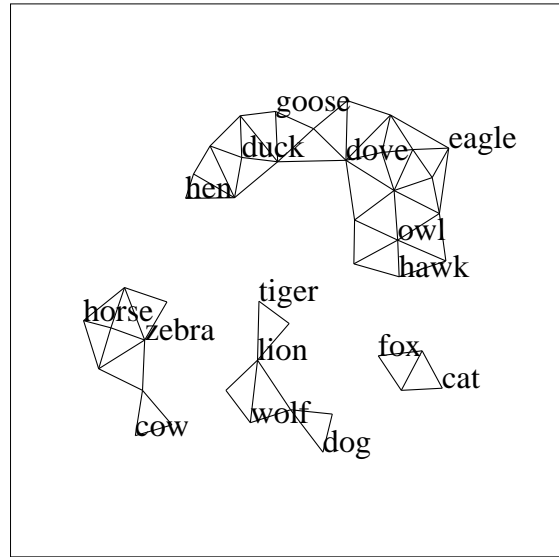
Ritter and Kohonen have introduced an illustrative example of high-dimensional data. It consists of the description of 16 animals by binary property lists (see table 1). The thirteen properties together with a 1-out-of- n coding of the name of the animal led to 29-dimensional vectors. These vectors were fed into a two-dimensional Kohonen feature map consisting of 10×10 neurons. After the end of the self-organization process it was tested where each of the input vectors was represented on the map. It came out that Kohonen's method had found an interesting projection positioning similar animals generally at neighboring locations on the map. It was, e.g., possible to partition this "semantotopic" map into three connected regions containing all birds, herbivores, and carnivores, respectively (see fig. 12).

We tested the Growing Cell Structures with the same data and constructed during the self-organization a two-dimensional embedding of the network with the method just described. Two different stages of a specific simulation are shown in fig. 13. When comparing the results with those of Ritter and Kohonen the main advantage of our model lies in the fact that it automatically finds meaningful partitions of the data, while Ritter and Kohonen had to identify those partitions by themselves.

In general this technique makes it possible to visualize and cluster high-dimensional data which might be useful in many application areas as e.g. process control or pattern recognition.

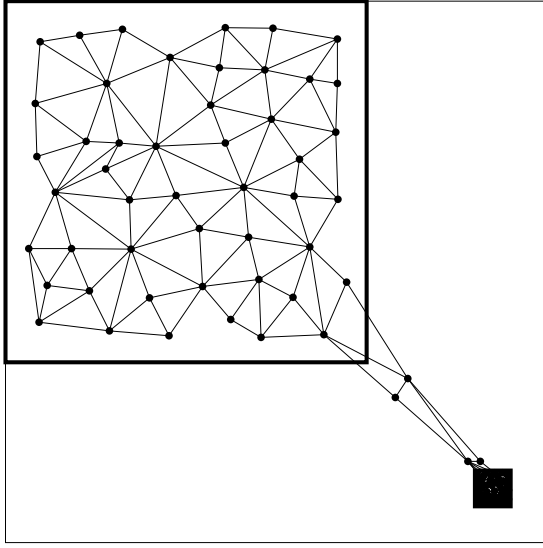


a) The birds have been divided from the mammals. Among the birds the peaceful ones and the birds of prey are at different positions. Also in the mammal cluster similar animals can generally be found at neighboring positions.

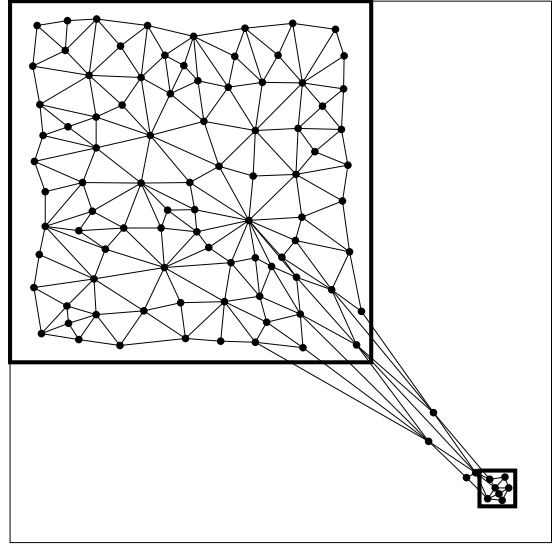


b) The mammal cluster has been split into three other clusters. One contains the large and peaceful animals (horse, zebra, cow), the second contains animals which like to run (tiger, lion, wolf, dog) and the third cluster contains animals which like to hunt, but avoid excessive running (cat, fox).

Figure 13: Semantotopic Growing Cell Structures. The data used stems from Ritter and Kohonen (see table 1). The data is ordered as by Kohonen's model but beyond that the ability of the Growing Cell Structures to form sub-structures makes it possible to partition the data in clusters of mutually similar items.



a) The original version of the Growing Cell Structures leads to a solution with approximately fifty percent of the reference vectors in the 10×10 field as well as in the 1×1 field. The mean square error is 0.00095.



b) The error-minimizing variant of the Growing Cell Structures positions most of the reference vectors in the 10×10 field. The mean square error is 0.00054.

Figure 14: Minimization of quantization error. The probability distribution consists of a 10×10 field and a 1×1 field. Fifty percent of the input signal come from either of these areas. After letting the networks grow until size 100 the mean square error was determined by 1000 test signals.

the signal counter of the best-matching unit as we did earlier, we change it through

$$\Delta\tau_s = \|\xi - w_s\|^2 \quad (20)$$

which effectively replaces eqn. 5. By using this measure as insertion criterion new neurons are inserted not anymore near those neurons getting the most input signals but rather near those neurons the input signals of which are very different from their reference vectors. The resulting network structures differ especially for probability distributions with a non-uniform probability density (see, e.g., fig 14). Recently this particular insertion criterion has been used to develop a new method for vector quantization (see Fritzke, 1993b). For this application the consistency requirements for the structures have been loosened by allowing also separate cells (without any neighbors) to exist. The method is able to generate codebooks of exceptionally good quality.

Another useful example for the resource is discussed in the next section where we report first results on a new supervised network based on the Growing Cell Structures.

3 Extension to Supervised Learning

3.1 Motivation

Self-organizing networks perform unsupervised learning. Frequently they generate ordered mappings of the input data onto some low-dimensional topological structure. In other cases they are used to partition the input data into subsets (or clusters) such that data items inside one subset are similar but items from different subsets are dissimilar.

In many situations, however, one has given input as well as corresponding output data. The problem is then to learn the underlying relation from a limited number of examples. For sake of concreteness let us in the following assume that our data consists of a number of pairs

$$(\xi_i \in R^n, \zeta_i \in R^m)$$

whereby ξ_i is the input and ζ_i is the desired output of the i -th pair.

Supervised learning methods are in these cases used to train networks to generate the desired output when they are presented with the input part of a specific data pair. Although this is not very useful *per se*, it is hoped that after finishing the training the network will be able to generate “reasonable” output values also for unknown input data. This is often denoted as *generalization*. It is a commonplace today that to achieve good generalization the number of free parameters of the network must be kept small. Otherwise there is the danger of “over-fitting” which denotes a situation where the network still improves on the training data, but already has a decreasing performance on the test data. Typical applications areas for supervised learning include pattern classification or function approximation.

In the following we demonstrate how the self-organizing model we presented in this paper can be extended to a supervised learning procedure. The result is a method which resembles the well-known radial basis function network (RBF) but eliminates some serious drawbacks of this approach.

3.2 Radial Basis Functions

Radial Basis Function networks (Moody & Darken, 1988) consist of a layer L of units with Gaussian activation functions⁵ and an output layer of m linear summation units (see fig. 15). We assume again data pairs $(\xi_i \in R^n, \zeta_i \in R^m)$ of input and desired output.

Each Gaussian unit c has an associated vector $w_c \in R^n$ indicating the position of the Gaussian in input vector space and a standard deviation σ_c . For a given input datum ξ the activation of a unit c is described by

$$D_c(\xi) = \frac{f_c(\xi)}{\sum_{i \in L} f_i(\xi)} \quad (21)$$

whereby

$$f_c(\xi) = \exp\left(-\frac{\|\xi - w_c\|^2}{\sigma_c^2}\right). \quad (22)$$

⁵In general every activation function could be used which is only in a limited and local area of the input vector space considerably different from zero.

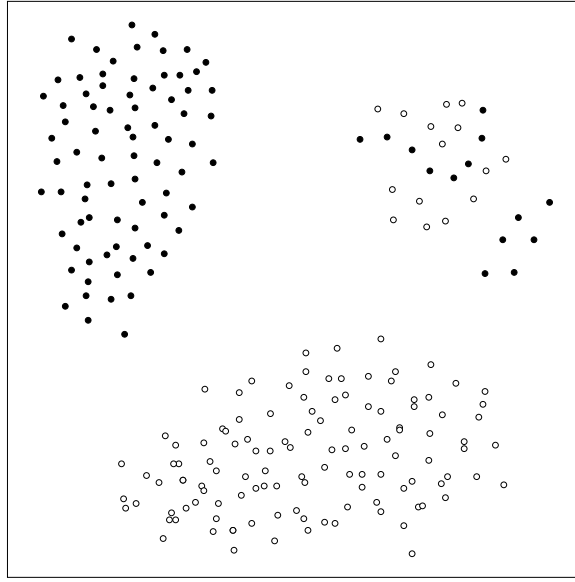


Figure 16: Classification problem with two classes: Given the shown example points find a good method to map all points in the square to one or the other class. (Alternatively one could also consider rejection of points, for which neither class seems to be appropriate.)

optimal. Consider a simple classification problem with two classes where most of the data vectors lie in two well-separated clusters, but the remaining vectors of both classes are scattered in several small clusters which are pretty close to each other (see fig. 16). In this case k-means would position most of the available Gaussians on the two large clusters. A much better choice, however, would be to cover the large clusters with only few Gaussians (having a large standard deviation) and to use the rest to cover the more complicated region containing the small clusters. Generally, relatively more Gaussians should be positioned at those locations where it is difficult to differentiate between the classes. These locations, however, are not known a priori.

3.3 Supervised Growing Cell Structures

In a fairly obvious way one can extend the Growing Cell Structures to a supervised radial basis function network (see fig. 17):

- For every cell c the reference vector w_c defines the center of a Gaussian activation function.
- The standard deviation σ_c of the Gaussian is defined as the mean length of all edges emanating from c (this is comparable to the heuristic proposed by Moody e.a).
- A number of m linear output units are defined and the Gaussian units are completely connected to them by weighted connections. This can be realized by associating with every cell c an output weight vector $w_c^{\text{out}} = (w_{1c}, w_{2c}, \dots, w_{mc})$. Thereby, w_{ic} denotes the weight of the connection from cell c to output unit i .

The change of weights (according to the delta rule) is defined by

$$\Delta w_{ic} = \eta(\zeta_i - o_i) D_c \quad (\text{for all } i \in \{1, \dots, m\}) (\text{for all } c \in A), \quad (26)$$

whereby η is the learning rate.

Finally, we update the resource variable of the current best-matching unit s by adding to it the overall squared error between actual output $o = o_1, \dots, o_m$ and desired output $\zeta = \zeta_1, \dots, \zeta_m$:

$$\Delta \tau_s = \|\zeta - o\|^2 \quad (27)$$

This replaces eqn. 5 where we incremented the resource variable τ respectively eqn. 20 where we summed up the quantization error.

If the current task is a classification problem (as opposed to a continuous input/output mapping), we can alternatively use the classification error. In this case the resource would be updated according to

$$\Delta \tau_s = \begin{cases} 0 & \text{if } \xi \text{ is classified correctly} \\ 1 & \text{otherwise} \end{cases} \quad (28)$$

Networks built with the classification error as insertion criterion tend to be still very small when they start classifying all training examples correctly. This is due to the fact that new cells are only inserted in those regions of the input vector space where still misclassifications occur. On the other hand, learning does practically halt when no misclassifications occur anymore even if the “raw” mean square error of the network is still rather large⁶. In some cases this can lead to poor generalization for unknown patterns. It, therefore, seems advisable to use a weighted combination of classification and mean square error. It has to be pointed out, however, that this is merely a matter of fine-tuning. From our experience the networks generate usually satisfying mappings in all areas where training vectors are available, no matter which combination of the two kinds of error is used.

Whenever a new cell r is inserted, it gets a vector $w_r^{\text{out}} = (w_{1r}, w_{2r}, \dots, w_{mr})$ of weighted connections to the m output units. Instead of initializing these vectors with zero or random values, they are obtained through a redistribution very similar to that used for the resource variable of the new cell (compare eqn. 10 and 11):

$$\Delta w_c^{\text{out}} = \frac{|F_c^{(\text{new})}| - |F_c^{(\text{old})}|}{|F_c^{(\text{old})}|} w_r^{\text{out}} \quad (\text{for all } c \in N_r). \quad (29)$$

whereby $|F_c|$ is the n -dimensional volume of F_c . Finally, the initial output weight vector of the new cell is defined as

$$w_r^{\text{out}} = - \sum_{c \in N_r} \Delta w_c^{\text{out}}, \quad (30)$$

In doing this redistribution the new cell is given output weights such, that it will activate the output units in a way similar to its “mean” neighbor. Since the neighboring Gaussians

⁶This particular behavior is also a characteristic of the original perceptron learning rule introduced by Rosenblatt (1958).

overlap considerably, the overall output behavior of the network is not changed very much. In future adaptation steps, however, the new unit can develop different weights and contribute so to better error reduction in this area of the input vector space. The complete algorithm for Supervised Growing Cell Structures is shown in fig. 18.

3.4 Simulation Examples

Example 1: A simple classification problem

We used the described supervised version of the Growing Cell Structures to construct a classifier for the data shown in fig. 16. The network was chosen to be two-dimensional. Since the data had to be classified into two classes, two output units were used. The combined growth and learning process was continued until the MSE for the training data fell below some bound. The resulting network (see fig. 19a) was used to map 200×200 points inside the the square region to either one or the other class (see fig. 19b).

One can observe that the size of the triangles and, therefore, the standard deviation of the Gaussians is considerably smaller in the region with the four small clusters (upper right). The reason is that the classification in this area is difficult and, therefore, many classification errors occur during training. This leads to insertions in this area. The resulting decision regions demonstrate that the final network classifies all training vectors correctly. Moreover, it seems to do a rather good job on classifying the other points inside the depicted region.

Example 2: The Two Spirals

A well-known benchmark in the connectionist community is the so called two-spiral problem. It consists of 194 two-dimensional vectors lying on two interlocked spirals which are the classes in this case (see fig. 20a). The task is to construct a classifier being able to distinguish between the two classes. This benchmark is interesting since, due to the low data dimensionality, it is possible to visualize the decision regions of the network during and after training. Moreover, it seems to be a rather difficult task for typical feed-forward networks, e.g., multi-layer perceptrons with sigmoidal activation functions. Lang & Witbrock (1989) were unable to solve the problem with a standard multi-layer network and had to use additional connections to achieve convergence. Fahlman & Lebiere (1990) used a constructive algorithm called Cascade-Correlation to solve the problem. The resulting decision regions of this network are shown in fig. 20b. One can note that the Cascade-Correlation algorithm is able to learn the training data, but the decision regions show several artifacts. In many cases points between two training vectors of a specific class are classified as belonging to the other class. This occurs especially in the outer parts of the spirals where the example patterns of one class are further apart from each other than from the representants of the other class. The resulting “cuts” in the spiral can be interpreted as poor generalization. In absence of other evidence it seems more natural to assume that those intermediate points belong to the same class. The decision regions produced by the network of Lang and Witbrock look similar.

Baum & Lang (1991) proposed a constructive method and also tested it with the two-spiral problem. Their approach employs an “oracle” that can tell for every point in the plane the desired class. Queries to the oracle are then used to position the hyperplanes

Initialize cell structure A with one k -dimensional simplex at random positions in $V = R^n$.
Create m linear output units.
Create a weighted connection w_{ic} from each cell $c \in A$ to each output unit i , ($i \in \{1, \dots, m\}$)
Associate every cell (vertex of the simplex) with a Gaussian function.
while (classification error not low enough)
repeat λ times
Choose I/O-pair $(\xi, \zeta) \in (R^n \times R^m)$ from training data
Determine best-matching unit s for ξ .
Increase matching for s and its direct neighbors.
Compute activation D_c for every cell $c \in A$ (see eqn. 24).
Compute the vector $o = (o_1, \dots, o_m)$ of all output unit activations (see eqn. 25).
Perform one delta-rule learning step for the weights (see eqn. 26)
Increase resource variable of s through $\Delta\tau_s = \ \zeta - o\ ^2$
Determine cell q with maximum resource value
Insert a new cell r between q and the direct neighbor f with maximum distance in input vector space
Redistribute resource values and weight vectors among r and its direct neighbors according to eqn. 10 - 11 and 29 - 30, resp.

Figure 18: Supervised Growing Cell Structures algorithm

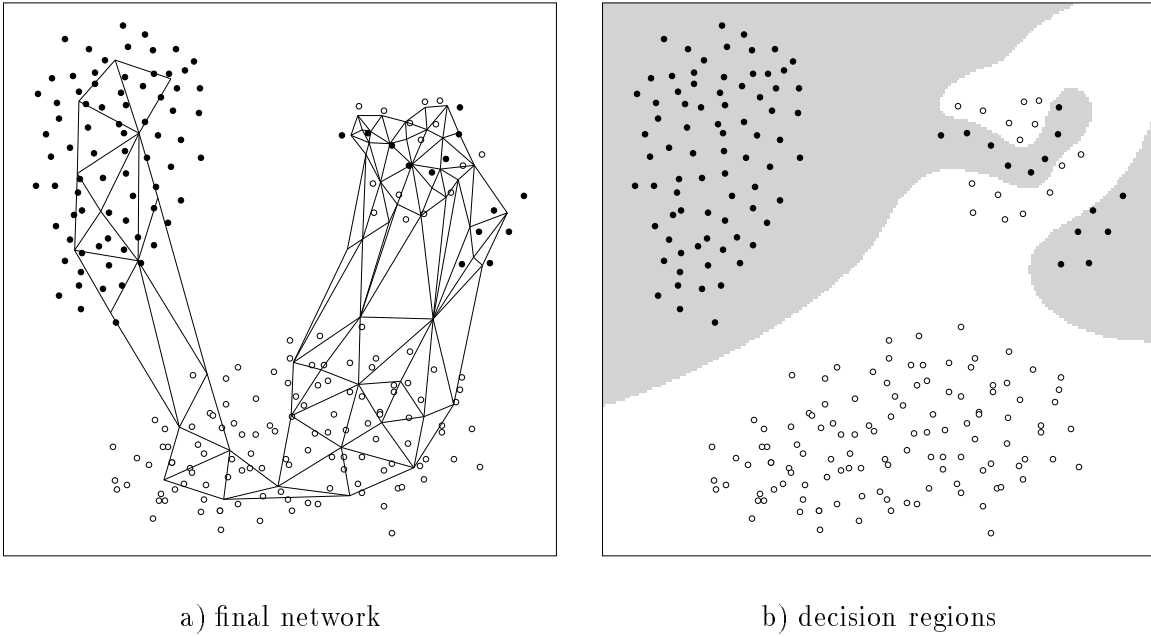


Figure 19: Supervised Growing Cell Structures. Network and decision regions for the data shown in fig. 16 Simulation parameters: $\lambda = 240$, $\varepsilon_b = 0.1$, $\varepsilon_n = 0.006$, $k = 2$, $\tilde{t} = 2$, $\alpha = 0.005$, $\eta = 0.15$, no removal of cells.

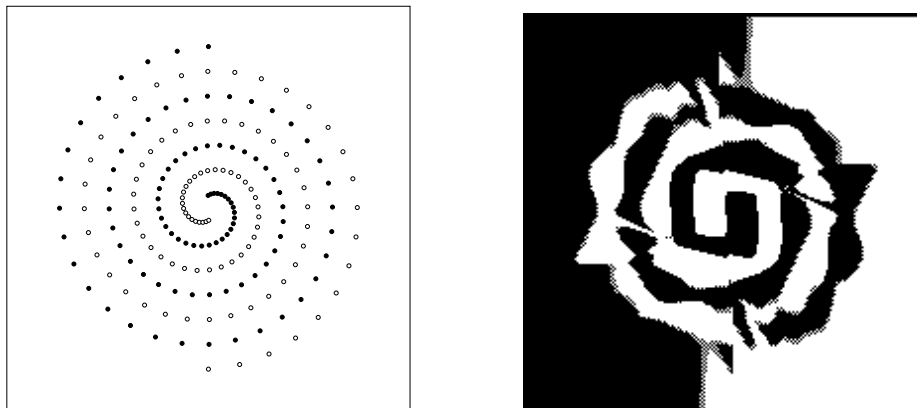


Figure 20: Two spiral problem and learning results of a constructive network.

network model	number of epochs	reported in
Backpropagation	20000	Lang & Witbrock (1989)
Cross Entropy BP	10000	Lang & Witbrock (1989)
Cascade-Correlation	1700	Fahlman & Lebiere (1990)
Growing Cell Structures	180	(this paper)

Table 2: Training epochs necessary for the two spiral problem

corresponding to certain hidden units. An explicit test set of 576 points has been defined consisting of three points between each pair of adjacent same-class training points. Therefore, training and test points together form two spirals with a four times higher point density than the training set alone. For their best model Baum and Lang report an average of 29 errors on the test set.

We generated a two-dimensional Growing Cell Structure to solve the two-spiral problem. The network and the corresponding decision regions are shown in fig. 21. In this case the decision regions form two well-separated spirals with very smooth borders. In fact, the decision regions exhibit a strong similarity to the oracle defined by Baum and Lang. Data points in between training vectors of one class are mapped onto that class and, therefore, the network makes no errors at all on the mentioned test set of Baum and Lang. Even in the outer regions of the spiral the decision regions follow the example vectors accurately. The local density of cells is rather uniform and does not follow the density of the training vectors which is higher near the center of the spirals. This is not surprising since near the center fewer units *per training point* are needed to facilitate correct classification.

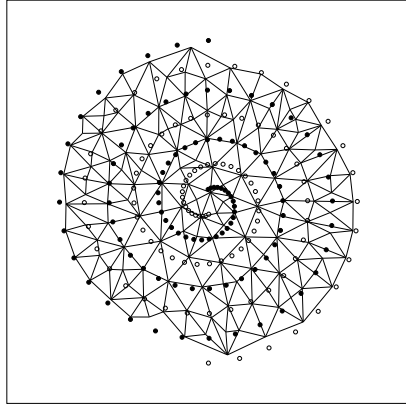
For every learning method an important practical aspect is the number of pattern presentations necessary to achieve a satisfying performance. In case of a finite training set a common measure is the number of cycles through all training patterns, also called *epochs*. We list in table 2 the number of epochs for the two-spiral problem for some earlier methods and for our approach. As can be seen the number of epochs required by the new method is about two orders of magnitude smaller than for standard backpropagation and nearly one order of magnitude smaller than for Cascade-Correlation.

Example 3: Speaker Independent Vowel Recognition

To explicitly investigate the generalization capability of our model, we performed experiments with a vowel recognition problem. The data used was collected by Deterding (1989), who recorded examples of the eleven steady state vowels of English spoken by fifteen speakers for a speaker normalization study. The vowel data (as well as the two-spiral data) is electronically available from the Carnegie-Mellon University connectionist benchmark collection (see Fahlman, 1993).

(An ASCII approximation to) the International Phonetic Association (I.P.A.) symbol and the word in which the eleven vowel sounds were recorded is given in table 3. The word was uttered once by each of the fifteen speakers, 7 of whom were female and 8 male.

The speech signals were low pass filtered at 4.7kHz and then digitized to 12 bits with



a) final network with 145 cells



b) decision regions

Figure 21: Performance of the Growing Cell Structures on the two-spiral benchmark. Simulation parameters: $\lambda = 240$, $\varepsilon_b = 0.1$, $\varepsilon_n = 0.006$, $k = 2$, $\tilde{t} = 2$, $\alpha = 0.005$, $\eta = 0.15$, no removal of cells.

vowel	word	vowel	word
i:	heed	O	hod
I	hid	C:	hoard
E	head	U	hood
A	had	u:	who'd
a:	hard	3:	heard
Y	hud		

Table 3: Words used in recording the vowels (from Robinson, 1989)

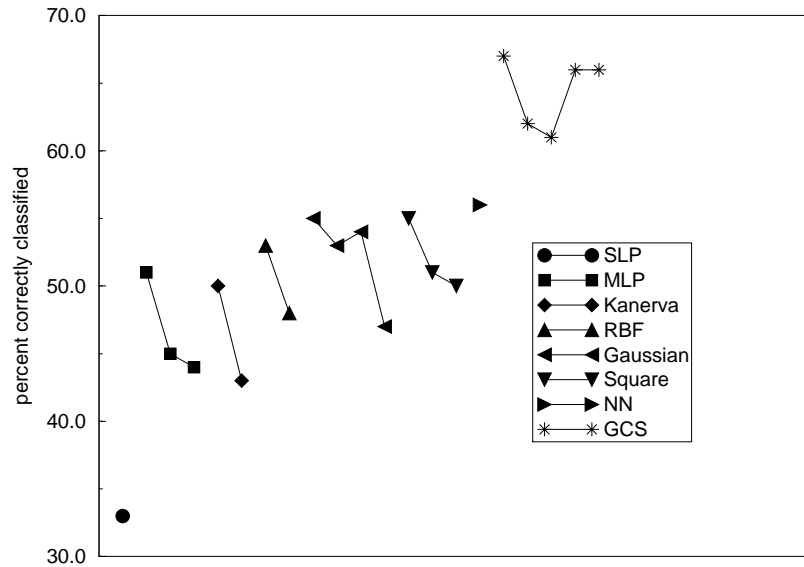


Figure 22: Percentage of correctly classified test patterns for the vowel recognition problem.

a 10kHz sampling rate. Twelfth order linear predictive analysis was carried out on six 512 sample Hamming windowed segments from the steady part of the vowel. The reflection coefficients were used to calculate 10 log area parameters, giving a 10 dimensional input space. A general introduction to speech processing and an explanation of this technique can be found in e.g. Rabiner & Schafer (1978). Each speaker, thus, yielded six frames of speech from eleven vowels. This gave 990 frames from the fifteen speakers.

Robinson used this data in his thesis (Robinson, 1989) to investigate several types of neural network algorithms. He used 528 frames from four male and four female speakers to train the networks and used the remaining 462 frames from four male and three female speakers for testing the performance.

The classifiers he examined were single-layer perceptrons, multi-layer networks with sigmoidal, Gaussian, and quadratic activation functions, a modified Kanerva model, radial basis networks, and also a conventional method, the nearest neighbor classifier. Due to the limited computational facilities available to Robinson, he did only one run for each of the different architectures. Every run was continued for about 3000 epochs (Robinson, 1993).

To get comparable results, we trained several Growing Cell Structure networks with the same data as Robinson and thereafter used his test data to evaluate the generalization capabilities of the networks. Since the input vector dimension was high-dimensional (10), we used also networks of a somewhat higher dimension than in the previous examples.

The results of Robinson and our results are shown in table 4. For easier comparison the percentage of correctly classified test patterns is shown graphically in fig. 22. It is evident from the simulations that our approach has the best results of the considered methods. The networks had to be trained only for about 80 epochs which compares rather well to the other methods. The ratio $3000/80 = 37.5$ is also approximately along the lines of our simulations for the two-spiral problem if one compares the number of epochs needed for cross entropy backpropagation and for our model (see table 2).

Classifier	number of hidden units	correctly classified	percent correct
Single-layer perceptron	-	154	33
Multi-layer perceptron	88	234	51
Multi-layer perceptron	22	206	45
Multi-layer perceptron	11	203	44
Modified Kanerva Model	528	231	50
Modified Kanerva Model	88	197	43
Radial Basis Function	528	247	53
Radial Basis Function	88	220	48
Gaussian node network	528	252	55
Gaussian node network	88	247	53
Gaussian node network	22	250	54
Gaussian node network	11	211	47
Square node network	88	253	55
Square node network	22	236	51
Square node network	11	217	50
Nearest neighbor	-	260	56
3-dimensional GCS	154	309	67
3-dimensional GCS	165	285	62
3-dimensional GCS	158	282	61
5-dimensional GCS	135	306	66
5-dimensional GCS	196	307	66

Table 4: Test results on vowel recognition problem. The table shows the network size, the number of correctly classified test patterns (out of 462), and the corresponding percentage. The upper box shows the results reported by Robinson in his thesis (Robinson, 1989). He got the best classification rate for the nearest neighbor method. The lower box shows the result of several nets generated by the Growing Cell Structures method. All of them have a higher rate of correctly classified test patterns than the nearest neighbor method (and all the other models examined by Robinson). We tried networks of dimensionality three and five. The parameter \tilde{t} was set equal to the network dimension in each case. The second run with a five-dimensional network was continued very long to see whether over-training effects could be produced which was not the case in that simulation. Also the different choices for \tilde{t} did not seem to influence the outcome of the algorithm very much.

4 Discussion

In the first part of the paper we introduced a new self-organizing network model. It has the following advantages over existing models:

- The network structure is determined automatically from the input data.
- The network size has not to be predefined. Instead, the growth process can be continued until a performance criterion is met.
- All parameter of the model are constant. It is, therefore, not necessary to define a decay schedule as in other models.
- The insertion of new units can be influenced such that the generated network estimates the probability density of the input signals, minimizes the quantization error or pursues still other goals.
- Since the final structure depends on the input data it can be used for data visualization and for clustering. In contrast, most other models have a fixed structure which does not provide any information of that kind.

In the second part of the paper we developed a combination of the self-organizing network with the radial basis function (RBF) approach. It provides a number of improvements over current network models with localized receptive fields (and also some other models):

- Number, diameter and position of RBF units are determined automatically through a growth process which can be stopped as soon as the network performance is good enough.
- Since positioning of RBF units and supervised training of connection weights is performed in parallel, the current classification error can be used to determine where to insert new RBF units. Previous approaches can only rely on clustering algorithms which often fail to find good positions for the RBF units with respect to classification accuracy.
- The networks are relatively small and generalize very well.
- The necessary number of training epochs seems to be one to two orders of magnitude smaller than for other approaches.

Although the results obtained so far are very promising, it is necessary to investigate the performance of the network for larger problems than the ones presented here. Furthermore, it would be an improvement if one could find ways to automatically choose some of those parameters which still have to be set by the user. An interesting goal would be a model with no parameters except the properties of the desired classifier. This goal is, of course, still very distant but we hope that the proposed methods are a step in the right direction.

Acknowledgements

The author likes to thank Scott Fahlman for the permission to reproduce figure 20b and for maintaining the CMU Benchmark Collection.

Notation

A	Growing Cell Structures network, also denotes set of cells in the network
k	dimensionality of the Growing Cell Structures network A
V	n -dimensional input vector space
n	dimensionality of V
w_c	n -dimensional reference (synaptic, weight) vector of cell c
w	set of all reference vectors for cells in A
ϕ_w	mapping $V \rightarrow A$
λ	adaptation steps per insertion
ε_b	adaptation parameter for best-matching unit
ε_n	adaptation parameter for neighboring cells
η	threshold for cell removal
$P(\xi)$	probability distribution of input signals
α	decrease parameter for resource variables
N_c	set of direct neighbors of a cell c
F_c	Voronoi field of cell c
t	true data dimensionality
\tilde{t}	estimate for t
$ F_c $	n -dimensional volume of F_c
ξ	n -dimensional input signal
ζ	m -dimensional output signal
m	dimension of the output vector space for supervised learning
(ξ, ζ)	I/O-pair (for supervised learning)
o	m -dimensional vector of output unit activations
τ_c	resource variable of cell c (can contain, e.g., signals, quantization error, classification error)
L	layer of Gaussian units in RBF-networks
$D_c(\xi)$	activation of Gaussian unit c
h_c	relative signal frequency of cell c .
\tilde{p}_c	estimate of the probability density near w_c
\hat{p}_c	estimate of the normalized probability density near w_c
$\Delta x = y$	short-cut for $x^{\text{new}} = x^{\text{old}} + y$
w_c^{out}	m -dimensional vector of weights from cell c to the output units
w_{ic}	weighted connection from Gaussian unit c to output unit i
$\ \cdot\ $	Euclidean vector norm

References

- Baum, E. B. & K. E. Lang (1991), Constructing hidden units using examples and queries, in *advances in neural information processing systems 3*, R.P. Lippmann, J.E. Moody & D.S. Touretzky, eds., Morgan Kaufmann Publ., Inc, San Mateo, pp. 904–910.
- Blackmore, J. & R. Miikkulainen (1992), Incremental grid growing: encoding high-dimensional structure into a two-dimensional feature map, University of Texas at Austin, TR AI92-192, Austin, TX.
- Deterding, D. H. (1989), Speaker Normalisation for Automatic Speech Recognition, University of Cambridge, Ph.D. thesis.
- Fahlman, S. E. (1993), CMU Benchmark Collection for Neural Net Learning Algorithms, Carnegie Mellon University, School of Computer Science, [machine-readable data repository], Pittsburgh.
- Fahlman, S. E. & C. Lebiere (1990), The Cascade-Correlation Learning Architecture, in *Advances in Neural Information Processing Systems 2*, D.S. Touretzky, ed., Morgan Kaufmann, San Mateo, pp. 524–532.
- Favata, F. & R. Walker (1991), A study of the application of Kohonen-type neural networks to the travelling Salesman Problem, *Biological Cybernetics*, 64, pp. 463–468.
- Fritzke, B. (1993a), Kohonen feature maps and growing cell structures – a performance comparison, in *Advances in Neural Information Processing 5*, L. Giles, S. Hanson & J. Cowan, eds., Morgan Kaufmann Publishers, San Mateo, CA.
- Fritzke, B. (1993b), Vector quantization with a growing and splitting elastic net, (*to appear in the proceedings of ICANN-93*), Amsterdam.
- Jokusch, S. (1990), A neural network which adapts its structure to a given set of patterns, in *Parallel Processing in Neural Systems and Computers*, R. Eckmiller, G. Hartmann & G. Hauske, eds., Elsevier Science Publishers B.V., pp. 169–172.
- Kangas, J. A., T. Kohonen & T. Laaksonen (1990), Variants of Self-Organizing Maps, *IEEE Transactions on Neural Networks*, 1, pp. 93–99.
- Kohonen, T. (1982), Self-Organized Formation of Topologically Correct Feature Maps, *Biological Cybernetics*, 43, pp. 59–69.
- Kohonen, T. (1988), The Neural Phonetic Typewriter, *IEEE computer*, 21, pp. 11–22.
- Kohonen, T., K. Mäkisara & T. Saramäki (1984), Phonotopic maps, – insightful representation of phonological features for speech recognition, *Proc. 7th Int. Conf. on Pattern Recognition*, Montreal.
- Lang, K. J. & M. J. Witbrock (1989), Learning to tell two spirals apart, in *Proceedings of the 1988 Connectionist Models Summer School*, D. Touretzky, G. Hinton & T. Sejnowski, eds., Morgan Kaufmann, San Mateo, pp. 52–59.

- Martinetz, T. M. & K. J. Schulten (1991), A “neural-gas” network learns topologies, in *Artificial Neural Networks*, T. Kohonen, K. Mäkisara, O. Simula & J. Kangas, eds., North-Holland, Amsterdam, pp. 397–402.
- Mehlhorn, K. & S. Näher (1989), LEDA, a library of efficient data types and algorithms, Universität des Saarlandes, Fachbereich Informatik, TR A 04/89, Saarbrücken.
- Moody, J. & C. Darken (1988), Learning with Localized Receptive Fields, in *Proceedings of the 1988 Connectionist Models Summer School*, D. Touretzky, G. Hinton & T. Sejnowski, eds., Morgan Kaufmann, San Mateo, pp. 133–143.
- Rabiner, L. R. & R. W. Schafer (1978), *Digital Processing of Speech Signals*, Prentice Hall, Englewood Cliffs, NJ.
- Ritter, H. J. (1991), Learning with the self-organizing map, in *Artificial Neural Networks*, T. Kohonen, K. Mäkisara, O. Simula & J. Kangas, eds., North-Holland, Amsterdam, pp. 379–384.
- Ritter, H. J. & T. Kohonen (1989), Self-Organizing Semantic Maps, *Biological Cybernetics*, 61, pp. 241–254.
- Robinson, A. J. (1989), Dynamic Error Propagation Networks, Cambridge University, PhD Thesis, Cambridge.
- Robinson, A. J. (1993), (personal communication).
- Rodrigues, J. S. & L. B. Almeida (1990), Improving the learning speed in topological maps of patterns, *Proc. of INNC*, Paris.
- Rosenblatt, F. (1958), The perceptron: a probabilistic model for information storage and organization in the brain, *Psychological Review*, 65, pp. 386–408.
- Schweizer, L., G. Parladori, G. L. Sicuranza & S. Marsi (1991), A fully neural approach to image compression, in *Artificial Neural Networks*, T. Kohonen, K. Mäkisara, O. Simula & J. Kangas, eds., North-Holland, Amsterdam, pp. 815–820.
- Willshaw, D. J. & C. von der Malsburg (1976), How Patterned Neural Connections Can Be Set Up by Self-Organization, *Proceedings of the Royal Society of London B*, 194, pp. 431–445.
- Xu, L. (1990), Adding learning expectation into the learning procedure of self-organizing maps, *Int. Journal of Neural Systems*, 1, pp. 269–283.