

## **Publication 7**

Clustering of the Self-Organizing Map

Juha Vesanto and Esa Alhoniemi  
In *IEEE Transactions on Neural Networks*, Volume 11,  
Number 3, pp. 586–600, 2000.

# Clustering of the Self-Organizing Map

Juha Vesanto and Esa Alhoniemi, *Student Member, IEEE*

**Abstract**—The self-organizing map (SOM) is an excellent tool in exploratory phase of data mining. It projects input space on prototypes of a low-dimensional regular grid that can be effectively utilized to visualize and explore properties of the data. When the number of SOM units is large, to facilitate quantitative analysis of the map and the data, similar units need to be grouped, i.e., clustered. In this paper, different approaches to clustering of the SOM are considered. In particular, the use of hierarchical agglomerative clustering and partitive clustering using  $k$ -means are investigated. The two-stage procedure—first using SOM to produce the prototypes that are then clustered in the second stage—is found to perform well when compared with direct clustering of the data and to reduce the computation time.

**Index Terms**—Clustering, data mining, exploratory data analysis, self-organizing map.

## I. INTRODUCTION

**D**ATA mining processes can be divided to six sequential, iterative steps:

- 1) problem definition;
- 2) data acquisition;
- 3) data preprocessing and survey;
- 4) data modeling;
- 5) evaluation;
- 6) knowledge deployment.

Each step is essential: The problem defines what data are used and what is a good solution. Modeling makes it possible to apply the results to new data. On the other hand, data modeling without good understanding and careful preparation of the data leads to problems. Finally, the whole mining process is meaningless if the new knowledge will not be used [1].

The purpose of survey is to gain insight into the data—possibilities and problems—to determine whether the data are sufficient and to select the proper preprocessing and modeling tools. Typically, several different data sets and preprocessing strategies need to be considered. For this reason, efficient visualizations and summaries are essential. In this paper, we focus on clusters since they are important characterizations of data.

The self-organizing map (SOM) [2] is especially suitable for data survey because it has prominent visualization properties. It creates a set of prototype vectors representing the data set and carries out a topology preserving projection of the prototypes from the  $d$ -dimensional input space onto a low-dimensional grid. This ordered grid can be used as a convenient visualization surface for showing different features of the SOM (and thus of the data), for example, the cluster structure [3].

However, the visualizations can only be used to obtain qualitative information. To produce summaries—quantitative descriptions of data properties—interesting groups of map units must be selected from the SOM. The most obvious such a group is the whole map. While its properties are certainly interesting, even more useful summaries can be prepared if the SOM (and thus the data) actually consist of two or more separate regions, and these are studied separately. Another option would be to consider all map units individually, but in the case of a large map, this could result in far too many summaries. Thus, to be able to effectively utilize the information provided by the SOM, methods to give good candidates for map unit clusters or groups are required. It should be emphasized that the goal here is not to find an optimal clustering for the data but to get good insight into the cluster structure of the data for data mining purposes. Therefore, the clustering method should be fast, robust, and visually efficient.

The clustering is carried out using a two-level approach, where the data set is first clustered using the SOM, and then, the SOM is clustered. The most important benefit of this procedure is that computational load decreases considerably, making it possible to cluster large data sets and to consider several different preprocessing strategies in a limited time. Naturally, the approach is valid only if the clusters found using the SOM are similar to those of the original data. In the experiments, a comparison between the results of direct clustering of data and clustering of the prototype vectors of the SOM is performed, and the correspondence is found to be acceptable.

## II. CLUSTERING

### A. Definitions

A clustering  $Q$  means partitioning a data set into a set of clusters  $Q_i$ ,  $i = 1, \dots, C$ . In crisp clustering, each data sample belongs to exactly one cluster. Fuzzy clustering [4] is a generalization of crisp clustering where each sample has a varying degree of membership in all clusters. Clustering can also be based on mixture models [5]. In this approach, the data are assumed to be generated by several parametrized distributions (typically Gaussians). Distribution parameters are estimated using, for example, the expectation-maximization algorithm. Data points are assigned to different clusters based on their probabilities in the distributions. However, neither fuzzy clustering nor mixture models based clustering are considered in this study because the goal was to evaluate clustering of the SOM using a few simple standard methods.

A widely adopted definition of optimal clustering is a partitioning that minimizes distances within and maximizes distances between clusters. However, this leaves much room for variation: within- and between-clusters distances can be defined

Manuscript received June 15, 1999; revised November 13, 1999.

The authors are with Neural Networks Research Centre, Helsinki University of Technology, Helsinki, Finland (e-mail: Juha.Vesanto@hut.fi; Esa.Alhoniemi@hut.fi).

Publisher Item Identifier S 1045-9227(00)04039-X.

TABLE I  
WITHIN-CLUSTER DISTANCES  $S(Q_k)$  AND  
BETWEEN-CLUSTERS DISTANCES  $d(Q_k, Q_l)$ ;  $x_i, x_{i'} \in Q_k, i \neq i',$   
 $x_j \in Q_l, k \neq l, N_k$  IS THE NUMBER OF SAMPLES IN CLUSTER  $Q_k$  AND  
 $c_k = 1/N_k \sum_{x_i \in Q_k} x_i$

Within-cluster distance	$S(Q_k)$
average distance	$S_a = \frac{\sum_{i,i'} \ \mathbf{x}_i - \mathbf{x}_{i'}\ }{N_k(N_k-1)}$
nearest neighbor distance	$S_{nn} = \frac{\sum_i \min_{i'} \{\ \mathbf{x}_i - \mathbf{x}_{i'}\ \}}{N_k}$
centroid distance	$S_c = \frac{\sum_i \ \mathbf{x}_i - \mathbf{c}_k\ }{N_k}$
Between-clusters distance	$d(Q_k, Q_l)$
single linkage	$d_s = \min_{i,j} \{\ \mathbf{x}_i - \mathbf{x}_j\ \}$
complete linkage	$d_{co} = \max_{i,j} \{\ \mathbf{x}_i - \mathbf{x}_j\ \}$
average linkage	$d_a = \frac{\sum_{i,j} \ \mathbf{x}_i - \mathbf{x}_j\ }{N_k N_l}$
centroid linkage	$d_{ce} = \ \mathbf{c}_k - \mathbf{c}_l\ $

in several ways; see Table I. The selection of the distance criterion depends on application. The distance norm  $\|\cdot\|$  is yet another parameter to consider. In this paper, Euclidean norm is used because it is widely used with SOM. In addition, the  $k$ -means error criterion is based on it.

Many algorithms utilize local criteria in clustering data. For example,  $S_{nn}$  and  $d_s$  in Table I are based on distance to nearest neighbor. However, the problem is that they are sensitive to noise and outliers. Addition of a single sample to a cluster can radically change the distances [6]. To be more robust, the local criterion should depend on collective features of a local data set [7]. Solutions include using more than one neighbor [8] or a weighted sum of all distances. It has been shown that the SOM algorithm implicitly uses such a measure [9].

### B. Algorithms

The two main ways to cluster data—make the partitioning—are hierarchical and partitive approaches. The hierarchical methods can be further divided to agglomerative and divisive algorithms, corresponding to bottom-up and top-down strategies, to build a hierarchical clustering tree. Of these agglomerative algorithms are more commonly used than the divisive methods and are considered in this paper.

Agglomerative clustering algorithms usually have the following steps:

- 1) Initialize: Assign each vector to its own cluster.
- 2) Compute distances between all clusters.
- 3) Merge the two clusters that are closest to each other.
- 4) Return to step 2 until there is only one cluster left.

In other words, data points are merged together to form a clustering tree that finally consists of a single cluster: the whole data

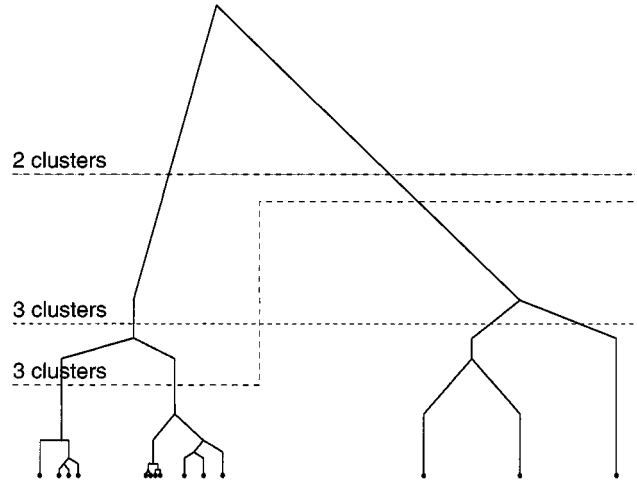


Fig. 1. Dendrogram of a set of 14 points in 1-D space. A partitioning can be obtained by cutting the dendrogram at a certain level, for example, at the level where there are only two clusters left, because there is a large jump in the dendrogram. An alternative way is to cut the dendrogram at different level for each branch.

set. The clustering tree (dendrogram) can be utilized in interpretation of the data structure and determination of the number of clusters.

However, the dendrogram does not provide a unique clustering. Rather, a partitioning can be achieved by cutting the dendrogram at certain level(s); see Fig. 1. The most typical solution is to cut the dendrogram where there is a large distance between two merged clusters. Unfortunately, this ignores the fact that the within-cluster distance may be different for different clusters. In fact, some clusters may be composed of several subclusters; to obtain sensible partitioning of the data, the dendrogram may have to be cut at different levels for each branch [10]. For example, two alternative ways to get three clusters are shown in Fig. 1.

Partitive clustering algorithms divide a data set into a number of clusters, typically by trying to minimize some criterion or error function. The number of clusters is usually predefined, but it can also be part of the error function [11]. The algorithm consists of the following steps.

- 1) Determine the number of clusters.
- 2) Initialize the cluster centers.
- 3) Compute partitioning for data.
- 4) Compute (update) cluster centers.
- 5) If the partitioning is unchanged (or the algorithm has converged), stop; otherwise, return to step 3.

If the number of clusters is unknown, the partitive algorithm can be repeated for a set of different number of clusters, typically from two to  $\sqrt{N}$ , where  $N$  is the number of samples in the data set. An example of a commonly used partitive algorithm is the  $k$ -means, which minimizes error function

$$E = \sum_{k=1}^C \sum_{\mathbf{x} \in Q_k} \|\mathbf{x} - \mathbf{c}_k\|^2 \quad (1)$$

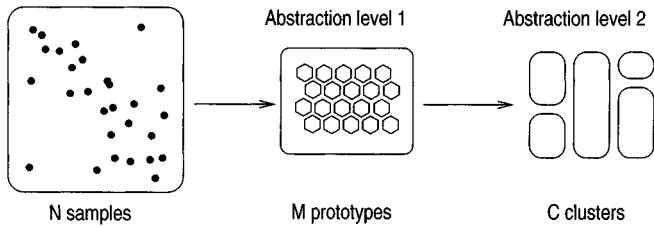


Fig. 2. First abstraction level is obtained by creating a set of prototype vectors using, e.g., the SOM. Clustering of the SOM creates the second abstraction level.

where  $C$  is the number of clusters, and  $c_k$  is the center of cluster  $k$ .

Partitive methods are better than hierarchical ones in the sense that they do not depend on previously found clusters. On the other hand, partitive methods make implicit assumptions on the form of clusters. For example,  $k$ -means tries to find spherical clusters.

To select the best one among different partitionings, each of these can be evaluated using some kind of validity index. Several indices have been proposed [6], [12]. In our simulations, we used the Davies–Bouldin index [13], which uses  $S_c$  for within-cluster distance and  $d_{cc}$  for between clusters distance. According to Davies–Bouldin validity index, the best clustering minimizes

$$\frac{1}{C} \sum_{k=1}^C \max_{l \neq k} \left\{ \frac{S_c(Q_k) + S_c(Q_l)}{d_{cc}(Q_k, Q_l)} \right\} \quad (2)$$

where  $C$  is the number of clusters. The Davies–Bouldin index is suitable for evaluation of  $k$ -means partitioning because it gives low values, indicating good clustering results for spherical clusters.

Many clustering methods, for example, agglomerative clustering using the complete linkage  $d_{co}$ , require the clusters to be compact and well separated. In real data, this is rarely the case. Rather, the gaps between clusters are obscured by noise, the clusters overlap, and there are outliers. When the data set size increases, computational load also becomes a problem, especially if some of the distance measures or indices are computationally expensive, like average distance  $S_a$  and average linkage  $d_a$ .

### C. Two-Level Approach

The approach used in this paper (clustering of the SOM rather than clustering the data directly) is depicted in Fig. 2. First, a large set of prototypes—much larger than the expected number of clusters—is formed using the SOM or some vector quantization algorithm. The prototypes can be interpreted as “proto-clusters,” which are in the next step combined to form the actual clusters. Each data vector of the original data set belongs to the same cluster as its nearest prototype. Similar multiple-level approaches to clustering have been proposed earlier, e.g., in [9]. While extra abstraction levels yield higher distortion, they also effectively reduce the complexity of the reconstruction task [14].

The primary benefit of the two-level approach is the reduction of the computational cost. Even with relatively small number of samples, many clustering algorithms—especially hierarchical

TABLE II  
TRAINING PARAMETERS OF THE SOM'S. PARAMETERS  $\sigma_1(0)$  AND  $\sigma_2(0)$  ARE INITIAL NEIGHBORHOOD WIDTHS FOR THE ROUGH AND FINE-TUNING PHASES, RESPECTIVELY. LAST COLUMN CONTAINS THE TOTAL TRAINING TIME

data set	map size	$\sigma_1(0)$	$\sigma_2(0)$	elapsed time
I	$19 \times 17$	10	2	45 s
II	$16 \times 13$	8	2	20 s
III	$24 \times 19$	12	2	569 s

TABLE III  
EXECUTION TIMES OF DIFFERENT PARTITIONINGS. TIMES FOR  $k$ -MEANS ARE FOR 29 PARTITIONINGS  $k = 2, \dots, 30$  USING 100 RUNS WITH DIFFERENT INITIALIZATIONS. COMPUTATION TIMES FOR CONSTRUCTING THE SOM'S ARE SHOWN IN THE FIRST COLUMN

data set	single linkage	complete linkage	average linkage	$k$ -means (100 runs)
I	27 min	27 min	1.9 h	84 min
II	2.3 min	2.3 min	9.4 min	32 min
III	3.2 h	3.1 h	13 h	9.0 h
SOM of I (45 s)	1.9 s	1.9 s	5.5 s	6.7 min
SOM of II (20 s)	1.2 s	1.2 s	3.1 s	3.8 min
SOM of III (9.5 min)	11 s	11 s	19 s	10 min

ones—become intractably heavy. For this reason, it is convenient to cluster a set of prototypes rather than directly the data [15].

Consider clustering  $N$  samples using  $k$ -means. This involves making several clustering trials with different values for  $k$ . The computational complexity is proportional to  $\sum_{k=2}^{C_{\max}} Nk$ , where  $C_{\max}$  is preselected maximum number of clusters. When a set of prototypes is used as an intermediate step, the total complexity is proportional to  $NM + \sum_k Mk$ , where  $M$  is the number of prototypes. With  $C_{\max} = \sqrt{N}$  and  $M = 5\sqrt{N}$ , the reduction of computational load is about  $\sqrt{N}/15$ , or about six-fold for  $N = 10\,000$ . Of course, this is a very rough estimate since many practical considerations are ignored.

The reduction is even greater for agglomerative algorithms since they cannot start from  $\sqrt{N}$  but have to start with  $N$  clusters and work their way down from there. If agglomerative clustering is used during the second step, the two-level strategy is a way to use partitive clustering to avoid those  $N - 5\sqrt{N}$  first steps. Table III shows computation times for both direct and two-level clustering strategies in our experiments. Even with these relatively small data sets, the difference is remarkable especially in the agglomerative methods.

Another benefit is noise reduction. The prototypes are local averages of the data and, therefore, less sensitive to random variations than the original data. In [16], partitive methods (i.e., a small SOM) greatly outperformed hierarchical methods in clustering imperfect data. Outliers are less of a problem since—by definition—there are very few outlier points, and therefore, their impact on the vector quantization result is limited. On the other hand, this may also be a problem if the outliers are actually the interesting areas. A recently proposed clustering method—Chameleon—also utilizes a two-level strategy to get more reliable estimates of cluster similarities [8].

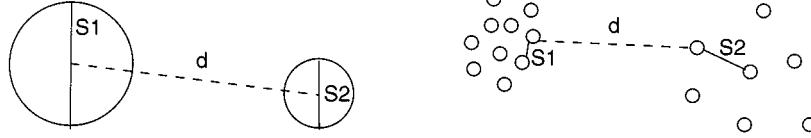


Fig. 3. Different criteria for finding interesting cluster merges: Davies–Bouldin (left) and gap (right). Merge is interesting when  $d > S1 + S2$ , where  $S1$  and  $S2$  are within-cluster distances, and  $d$  is the distance between clusters.

### III. CLUSTERING OF THE SOM

#### A. SOM Training

The SOM consists of a regular, usually two-dimensional (2-D), grid of map units. Each unit  $i$  is represented by a prototype vector  $\mathbf{m}_i = [m_{i1}, \dots, m_{id}]$ , where  $d$  is input vector dimension. The units are connected to adjacent ones by a neighborhood relation. The number of map units, which typically varies from a few dozen up to several thousand, determines the accuracy and generalization capability of the SOM. During training, the SOM forms an elastic net that folds onto the “cloud” formed by the input data. Data points lying near each other in the input space are mapped onto nearby map units. Thus, the SOM can be interpreted as a topology preserving mapping from input space onto the 2-D grid of map units.

The SOM is trained iteratively. At each training step, a sample vector  $\mathbf{x}$  is randomly chosen from the input data set. Distances between  $\mathbf{x}$  and all the prototype vectors are computed. The best-matching unit (BMU), which is denoted here by  $b$ , is the map unit with prototype closest to  $\mathbf{x}$

$$\|\mathbf{x} - \mathbf{m}_b\| = \min_i \{\|\mathbf{x} - \mathbf{m}_i\|\}. \quad (3)$$

Next, the prototype vectors are updated. The BMU and its topological neighbors are moved closer to the input vector in the input space. The update rule for the prototype vector of unit  $i$  is

$$\mathbf{m}_i(t+1) = \mathbf{m}_i(t) + \alpha(t)h_{bi}(t)[\mathbf{x} - \mathbf{m}_i(t)] \quad (4)$$

where

$t$	time;
$\alpha(t)$	adaptation coefficient;
$h_{bi}(t)$	neighborhood kernel centered on the winner unit:

$$h_{bi}(t) = \exp\left(-\frac{\|\mathbf{r}_b - \mathbf{r}_i\|^2}{2\sigma^2(t)}\right) \quad (5)$$

where  $\mathbf{r}_b$  and  $\mathbf{r}_i$  are positions of neurons  $b$  and  $i$  on the SOM grid. Both  $\alpha(t)$  and  $\sigma(t)$  decrease monotonically with time. There is also a batch version of the algorithm where the adaptation coefficient is not used [2].

In the case of a discrete data set and fixed neighborhood kernel, the error function of SOM can be shown to be [17]

$$E = \sum_{i=1}^N \sum_{j=1}^M h_{bj} \|\mathbf{x}_i - \mathbf{m}_j\|^2 \quad (6)$$

where  $N$  is number of training samples, and  $M$  is the number of map units. Neighborhood kernel  $h_{bj}$  is centered at unit  $b$ , which is the BMU of vector  $\mathbf{x}_i$ , and evaluated for unit  $j$ . If neigh-

borhood kernel value is one for the BMU and zero elsewhere, this leads to minimization of (1)—the SOM reduces to adaptive  $k$ -means algorithm [18]. If this is not the case, from (6), it follows that the prototype vectors are not in the centroids of their Voronoi sets but are local averages of all vectors in the data set weighted by neighborhood function values.

The SOM algorithm is applicable to large data sets. The computational complexity scales linearly with the number of data samples, it does not require huge amounts of memory—basically just the prototype vectors and the current training vector—and can be implemented both in a neural, on-line learning manner as well as parallelized [32]. On the other hand, the complexity scales quadratically with the number of map units. Thus, training huge maps is time consuming, although the process can be speeded up with special techniques; see, for example, [33] and [34]. For example, in [33], a SOM with million units was trained with 6.8 million 500-dimensional data vectors.

If desired, some vector quantization algorithm, e.g.,  $k$ -means, can be used instead of SOM in creating the first abstraction level. Other possibilities include the following.

- Minimum spanning tree SOM [19], neural gas [20], growing cell structures [21], and competing SOM's [22] are examples of algorithms where the neighborhood relations are much more flexible and/or the low-dimensional output grid has been discarded. Their visualization is much less straightforward than that of the SOM.
- In addition, several such growing variants of the SOM have been proposed where the new nodes do have a well-defined place on low-dimensional grid, and thus, the visualization would not be very problematic [23]–[27].

The SOM variants were not used in this study because we wanted to select the most commonly used version of the SOM. However, the principles presented in this paper could naturally be applied to the prototype vectors of these variants as well.

The two-level approach in clustering can only work if the prototypes reflect the properties of the data. The clusters found using the prototypes should be similar to those that would have been found directly from the data. In vector quantization, it has been shown that the density of the prototype vectors is proportional to  $\text{const} \cdot p(\mathbf{x})^{d/(d+r)}$ , where  $p(\mathbf{x})$  is the probability density function (p.d.f.) of the input data,  $d$  is dimension, and  $r$  is distance norm [28], [29]. For the SOM, connection between the prototypes and the p.d.f. of the input data has not been derived in general case. However, a similar power law has been derived in the 1-D case [30]. Even though the law holds only when the number of prototypes approaches infinity and neighborhood width is very large, numerical experiments have shown that the computational results are relatively accurate even for a

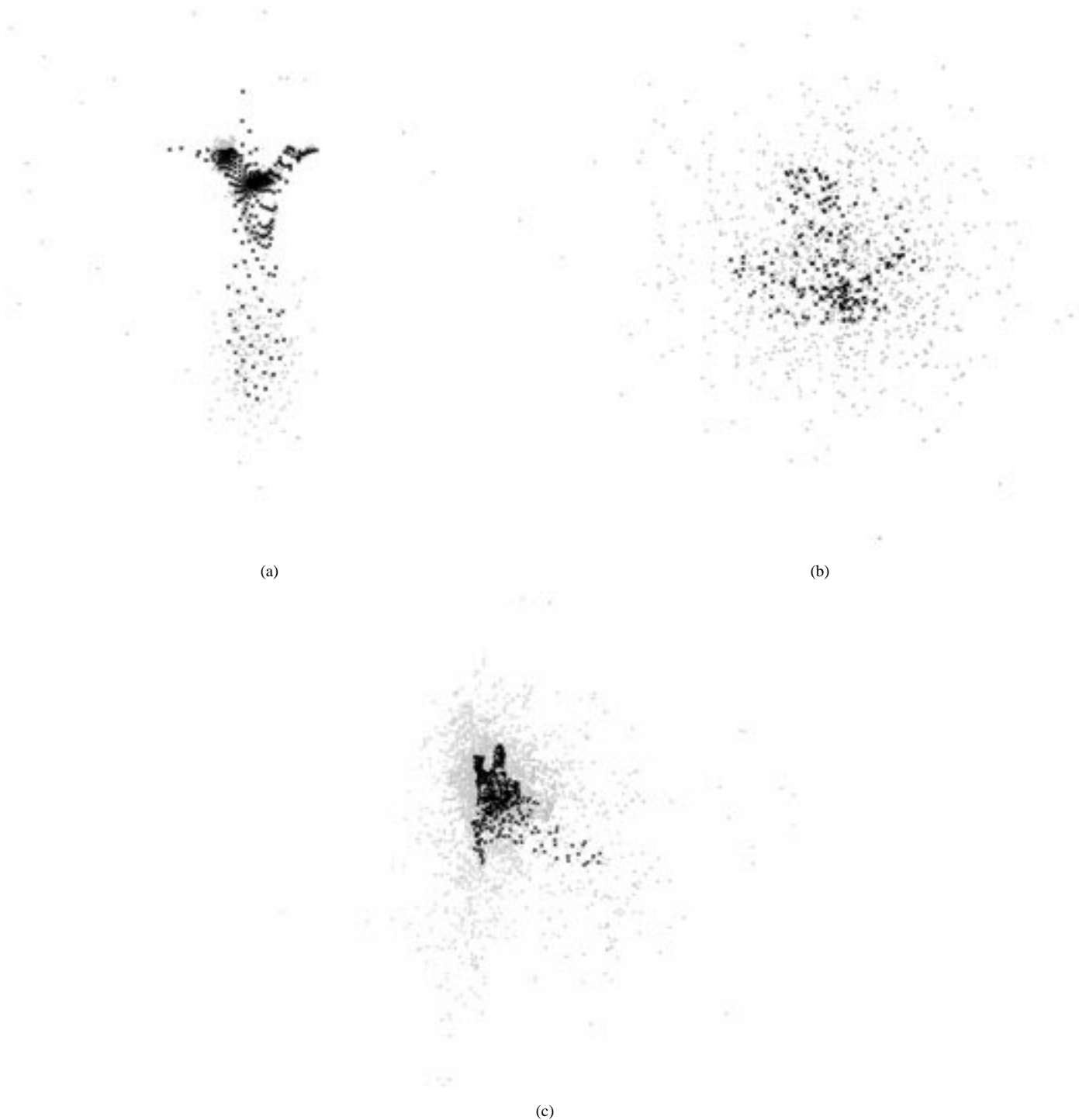


Fig. 4. Illustration of (a) data set I, (b) II, and (c) and III. Two-dimensional data set I is directly plotted in two dimensions, whereas data sets II and III are projected to the 2-D subspace spanned by the two eigenvectors with largest eigenvalues. The gray dots are data points, and the black crosses are prototype vectors of the SOM trained with corresponding data.

small number of prototypes [31]. Based on close relation between the SOM and  $k$ -means, it can be assumed that the SOM roughly follows the density of training data when not only the number of map units but also the final neighborhood width are small, as was the case in our experiments.

### B. Visual Inspection

An initial idea of the number of clusters in the SOM, as well as their spatial relationships, is usually acquired by visual

inspection of the map. The most widely used methods for visualizing the cluster structure of the SOM are distance matrix techniques [35], [36], especially the unified distance matrix (U-matrix). The U-matrix shows distances between prototype vectors of neighboring map units. Because they typically have similar prototype vectors, U-matrix is actually closely related to the single linkage measure. It can be efficiently visualized using gray shade [37]; see, for example, Figs. 7(a), 11(a), and 12(a).

Another visualization method is to display the number of hits in each map unit. Training of the SOM positions interpolating map units between clusters and thus obscures cluster borders. The Voronoi sets of such map units have very few samples (“hits”) or may even be empty. This information can be utilized in clustering the SOM by using zero-hit units to indicate cluster borders [38].

Generic vector projection methods can also be used. As opposed to the methods above, these are generally applicable to any set of vectors, for example the original data set. The high-dimensional vectors are projected to a low dimension so that intersample distances are preserved as faithfully as possible. Thus, the high-dimensional data set can be visualized while still preserving its essential topological properties. Examples of such nonlinear projection methods include multidimensional scaling techniques [39], [40], Sammon’s mapping [41], and curvilinear component analysis [42]. A special technique is to project the prototype vectors into a color space so that similar map units are assigned similar colors [43], [44].

Of course, based on the visualization, one can select clusters manually. However, this is a tedious process and nothing guarantees that the manual selection is done consistently. Instead, automated methods are needed.

### C. SOM Clustering

In agglomerative clustering, the SOM neighborhood relation can be used to constrain the possible merges in the construction of the dendrogram [45]. In addition, knowledge of interpolating units can be utilized both in agglomerative and partitive clustering by excluding them from the analysis. If this is used together with the neighborhood constraint in agglomerative clustering, the interpolative units form borders on the map that the construction of the dendrogram must obey. It may even be that the interpolating units completely separate some areas from the rest of the map.

The dendrogram of a map with  $M$  units contains  $M - 1$  different nonsingleton clusters, which is far too many to facilitate summarization. Additionally, most of these clusters are completely uninteresting: Many of them differ only by addition of a single or a few map units. Therefore, the dendrogram has to be pruned. The objective is to find the interesting groups. This can be done by testing for each merge of the dendrogram whether the involved clusters are sufficiently different from each other. If they are, both are added to the list of interesting groups. Interesting merges can be defined, for example, as follows (see Fig. 3).

- The Davies–Bouldin index (2) can be calculated for the two clusters. If the index is greater than one, there is an interesting merge.
- If the gap between the two clusters  $d_s(Q_k, Q_l)$  is greater than the sum of mean distances between points in the two clusters  $S_{nn}(Q_k) + S_{nn}(Q_l)$ , there is an interesting merge.

To be able to evaluate the criteria above, all singleton clusters have to be ignored.

Note that the distance measures used in the gap criterion are very sensitive to noise, especially if the number of samples in a cluster is small. However, the problem can be alleviated some-

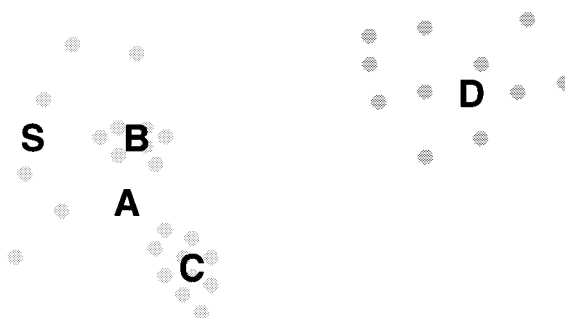


Fig. 5. High-level cluster  $A$  consists of subclusters  $B$  and  $C$  and a set of separate samples  $S$ . The subclusters  $B$  and  $C$  represent focus on local details, but in larger context, they, together with  $S$ , form the high-level cluster  $A$ .

what by requiring a wider gap if the sample size is small. Let  $x$  be a random sample of size  $N$ , whose mean can be estimated as  $\mu = E\{x\} = \sum x/N$ . It is easy to show that its variance is  $\text{Var}\{\mu\} = \text{Var}\{x\}/N$ . In many distribution models—for example, the gamma distribution—the variance is directly proportional to expectation  $\text{Var}\{x\} = a_1 E\{x\}$ . Therefore, let  $\text{Var}\{\mu\} = a_1 E\{x\}/N$ . Now, let  $x$  be the distance between each sample and its nearest neighbor in cluster  $Q_k$  of size  $N_k$  samples. A robust upper limit for nearest neighbor distance  $S_{nn}(Q_k)$  can be acquired by  $S_{nn}^*(Q_k) = E\{x\} + a_2 \text{Var}\{x\} = E\{x\} + a_2 a_1 E\{x\}/N_k = E\{x\}(1 + a/N_k)$ . In order to make also the between-clusters distance  $d_s(Q_k, Q_l)$  more robust, the mean of four shortest distances between the samples can be used instead of only the shortest distance.

The pruning produces a more clear dendrogram but still does not provide a unique partitioning. This is not a problem if the objective is only to find a set of interesting groups. If necessary, a final partitioning can be selected by hand using some kind of interactive tool [10]. By viewing the clustering hierarchy down from top, one can concentrate on the big, firmly separated clusters and discard the low-level groups. In our experiments, however, an automated procedure for finding the partitioning from the dendrogram was used; see Section IV-C for details.

## IV. EXPERIMENTS

### A. Data Sets

In the experiments, one real-world and two artificial data sets were used. Plots of the three data sets are presented in Fig. 4.

- Data set I (“clown”) consisted of 2220 2-D samples. The data set had the following seven clusters:
  - cluster with three subclusters (right eye);
  - spherical cluster (left eye);
  - elliptical cluster (nose);
  - nonspherical cluster (U-shaped: mouth);
  - large and sparse cluster (body).

In addition, some outliers were added to see how they affect the clustering. As it happens, the outliers are mostly concentrated on the left and seem to form a cluster of their own.

- Data set II was generated to see how the algorithms perform when there is only one cluster: a spherical Gaussian distribution in eight dimensions (1000 samples).

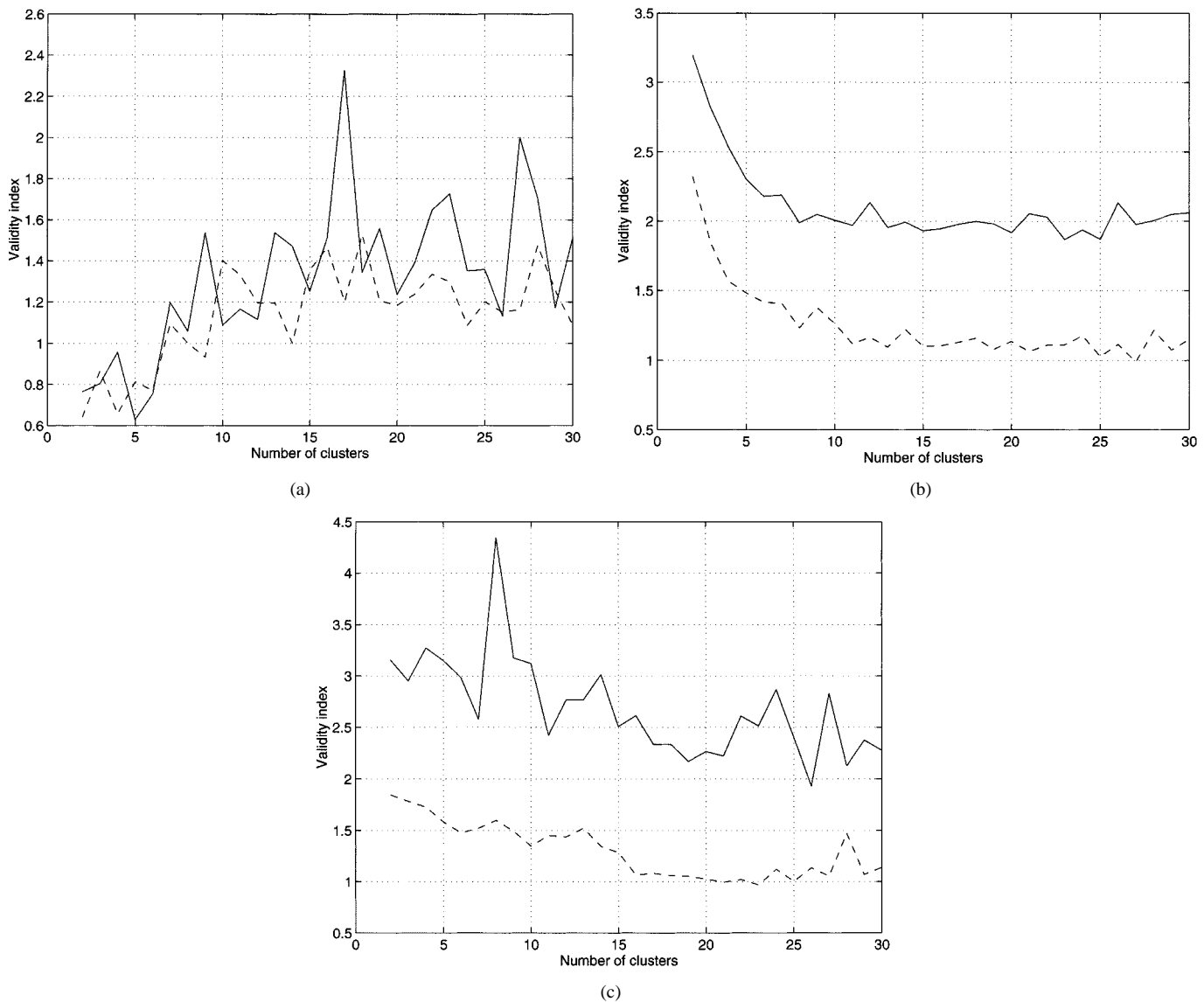


Fig. 6. Davies–Bouldin index for partitive clustering of data sets I–III (a)–(c) as function of the number of clusters. In each plot, the number of clusters runs on horizontal axis, and Davies–Bouldin index values lie on the vertical axis. Clustering results for the data are marked by solid line and results for the prototype vectors of the SOM by dashed line, respectively.

- Data set III consisted of 4205 75-D samples of the technology of pulp and paper mills of the world [46].

Especially in data set III, the original variables had very different scales. Therefore, as part of preprocessing, all variables in each data set were linearly scaled to have zero mean and unit variance.

### B. SOM's

A SOM was trained using the sequential training algorithm for each data set.<sup>1</sup> All maps were linearly initialized in the subspace spanned by the two eigenvectors with greatest eigenvalues computed from the training data. The maps were trained in two phases: a rough training with large initial neighborhood width and learning rate and a fine-tuning phase with small initial neighborhood width and learning rate; see Table II. The

<sup>1</sup>In the SOM training, freely available Matlab package SOM Toolbox was used. For more information, see URL <http://www.cis.hut.fi/projects/som-toolbox/>.

neighborhood width decreased linearly to 1; the neighborhood function was Gaussian [see (5)]. The training length of the two phases were 3 and 10 epochs, and the initial learning rates 0.5 and 0.05, respectively. The learning rate decreased linearly to zero during the training.

The prototype vectors of the SOM's have been superimposed on the plots of the data sets in Fig. 4. Note that the prototype vectors are well within the data sets and that this has removed most of the effect of outliers in data set I [Fig. 4(a)].

### C. Clustering

For each SOM, a visual inspection was performed. Then, agglomerative and partitive clustering algorithms were applied. Map units with empty Voronoi sets were excluded from the analysis. The neighborhood relations were not used.

In the agglomerative clustering, single, average, and complete linkages were used in the construction phase. The pruning was carried out using gap criterion with  $a = 3$  (see Section III-C).



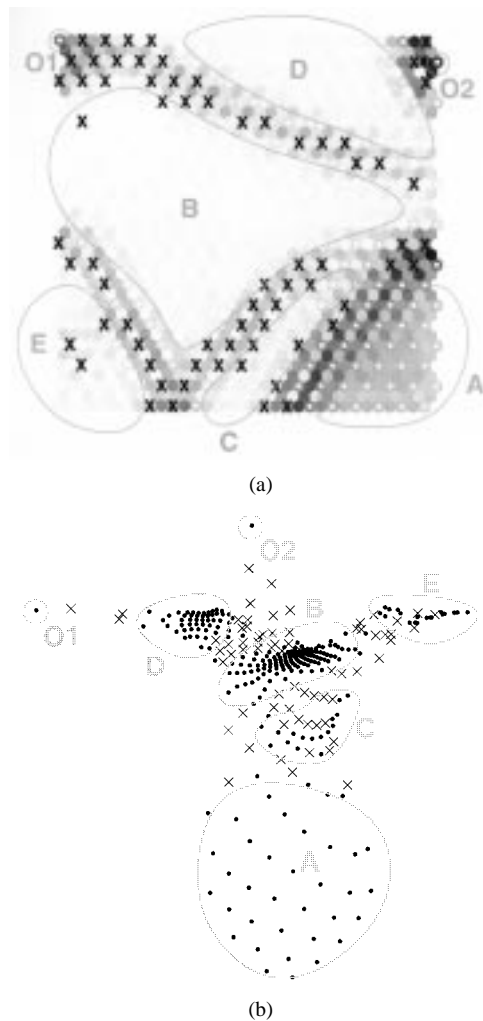


Fig. 7. (a) U-matrix and (b) prototype vectors of the SOM of data set I in input space. In the U-matrix, dark color indicates large distance between adjacent map units and light color small distance, respectively. Data clusters have been encircled in both figures, and the corresponding ones have been given the same labels. The interpolating map units have been marked with "X"s. Note that the interpolating map units completely separate the two outliers ("O1" and "O2") from the rest of the map.

Excluding the interpolating units proved to be very important; if this was not done, the interpolating units blurred out the gaps that the criterion was trying to detect.

As mentioned earlier, the list of interesting clusters does not give a specific partitioning. Looking at the pruned list down from top, the clusters form a hierarchy where each high-level cluster consists of some low-level clusters and a set of separate samples, see Fig. 5. The final partitioning was done by the following procedure.

- 1) A parent cluster was used if the subclusters represented less than half of its samples. Otherwise, the parent cluster was ignored, and the subclusters were used instead.
- 2) All remaining units (interpolative as well as those ignored above) were partitioned to nearest cluster.
- 3) Each sample in the original data was partitioned to the same cluster as its BMU.

Applying the same agglomerative clustering procedure directly to the data proved to be problematic. Because the original data sets were both much bigger and more noisy than the SOM

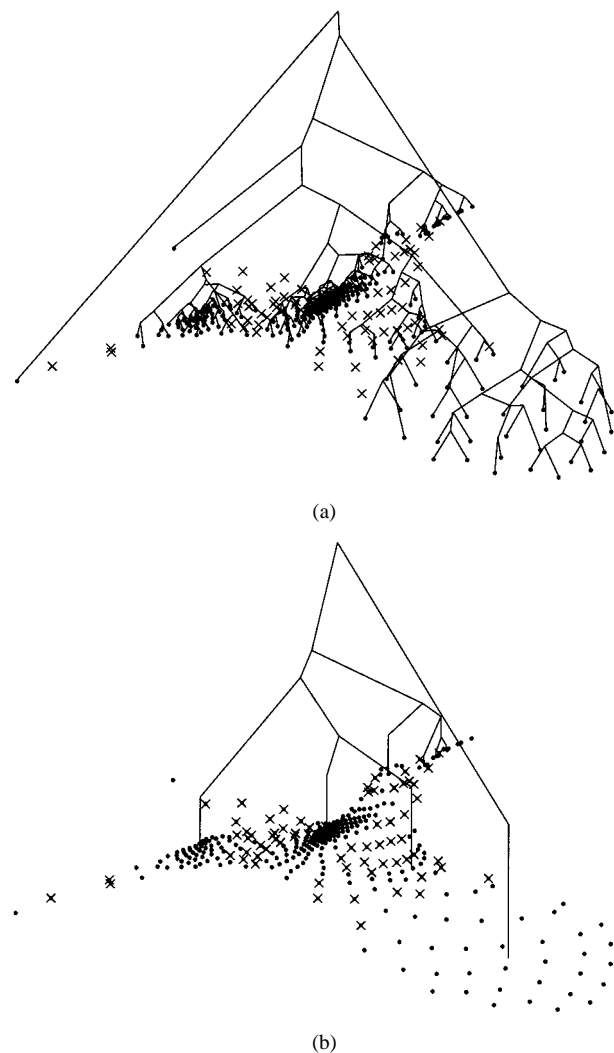


Fig. 8. (a) Full dendrogram and (b) pruned dendrogram of the SOM of data set I. The pruning is essential because it reduces the number of different clusterings into a reasonable level.

prototypes, using  $a = 3$  resulted in hundreds of small clusters. For data set I,  $a = 30$  gave reasonable results and was used for comparison. For data sets II and III, the gap criterion did not work at all. Instead, the comparisons in Section IV-G were done based on partitionings done by simply cutting the dendrograms at certain level.

The partitive clustering of both SOM's and the data sets was carried out using batch  $k$ -means algorithm. Because  $k$ -means is sensitive to initialization, it was run 100 times for each  $k$  with different random initializations. The best partitioning for each number of clusters was selected using error criterion in (1). Another possibility would have been to use some annealing technique to better avoid local minima of the error function [47]–[49].

To select the best clustering among the partitionings with different number of clusters, the Davies–Bouldin validity index (2) was used. In practice, though, it is better to use the index values as a guideline rather than absolute truth. As seen in Fig. 6, the index plots have several local minima. Each sharp local minimum in the validity index plot is certainly worth a look since it

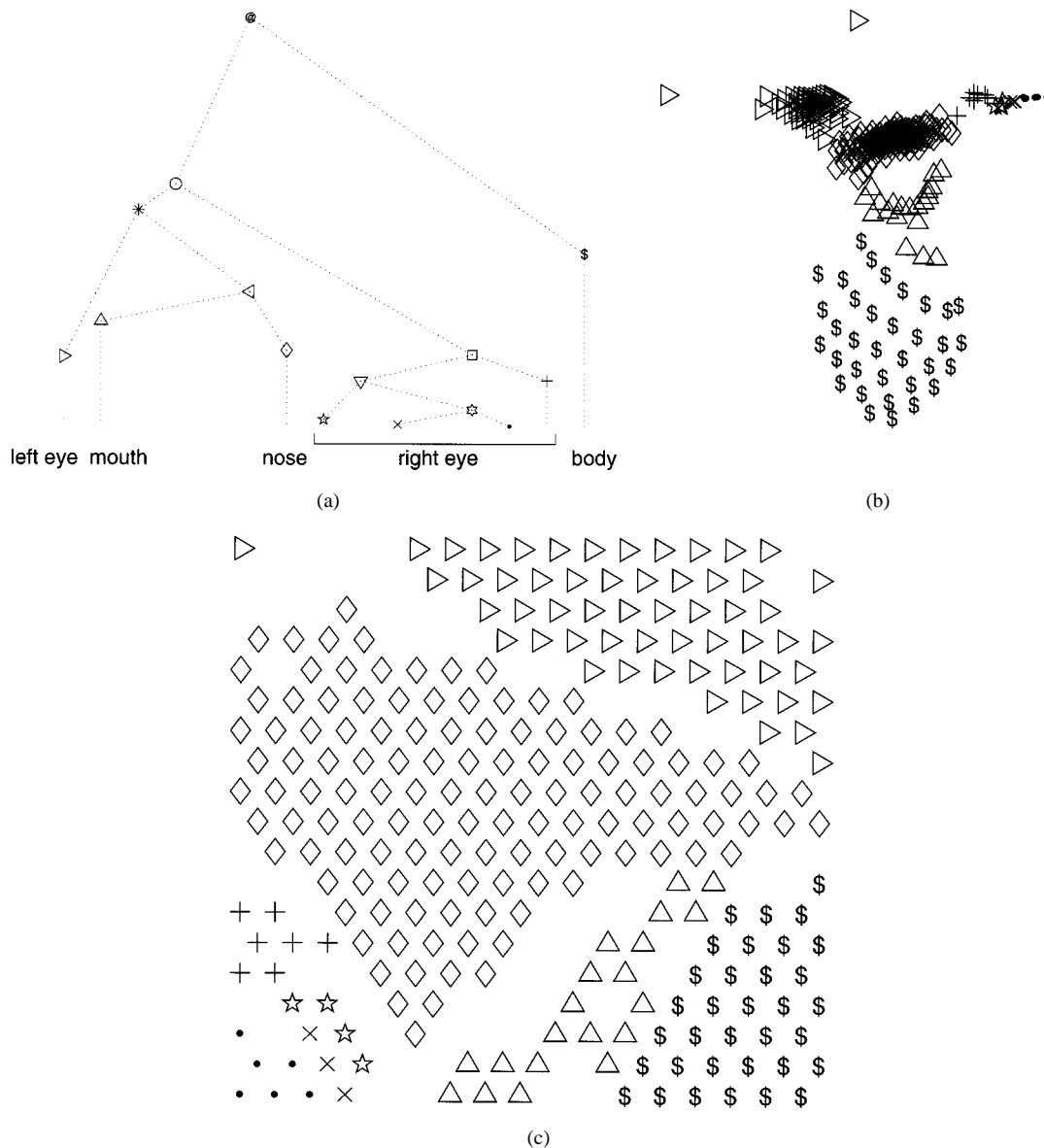


Fig. 9. Results of agglomerative clustering of the SOM of data set I. The interpolating units have been left out of the analysis. The pruned dendrogram (a) has 15 different clusters: Each node and leaf represents one “interesting” cluster. The same symbols have been used to indicate (b) cluster memberships of map units in the input space and (c) on the map grid.

indicates that the addition of one cluster enabled the algorithm to partition the data much better.

#### D. SOM of Data Set I

*Visual Inspection:* In Fig. 7, two visualizations of the SOM of data set I are shown: the U-matrix and the map units in the input space. Corresponding clusters in both figures have been encircled and labeled. From the U-matrix, one can easily detect seven distinct clusters, two of which (“O1” and “O2”) are very small and correspond to outliers. The interpolating units are positioned neatly between the clusters, although only the two outliers are completely separated from the rest of the map. There are also some interpolating units within the cluster corresponding to right eye (“E”), possibly indicating presence of two or more subclusters. Note that the relative size of the clusters is proportional to the number of samples in each cluster and not to the radius of the cluster in the input space.

*Agglomerative Clustering:* The full and pruned dendrograms of the SOM of data set I are shown in Fig. 8. The procedure produced 15 interesting clusters, whose hierarchy is shown in Fig. 9(a). The final partitioning was based on the pruned dendrogram: the results, excluding the interpolating units, are shown in Fig. 9(b) and (c).

The result is satisfactory. In the final partitioning, all the original clusters have been found: even the three subclusters of the right eye. However, there is one additional cluster in the right eye, so the performance is not perfect. The right eye causes problems in the SOM-based approach because it is in an extreme corner of the data distribution and the amount of data in that part of the input space is relatively small. If better accuracy is desired, a new clustering can be made using data in that area only. In addition, note that if the neighborhood constraint had been used in clustering, the outlying map units would have formed two small clusters of their own because

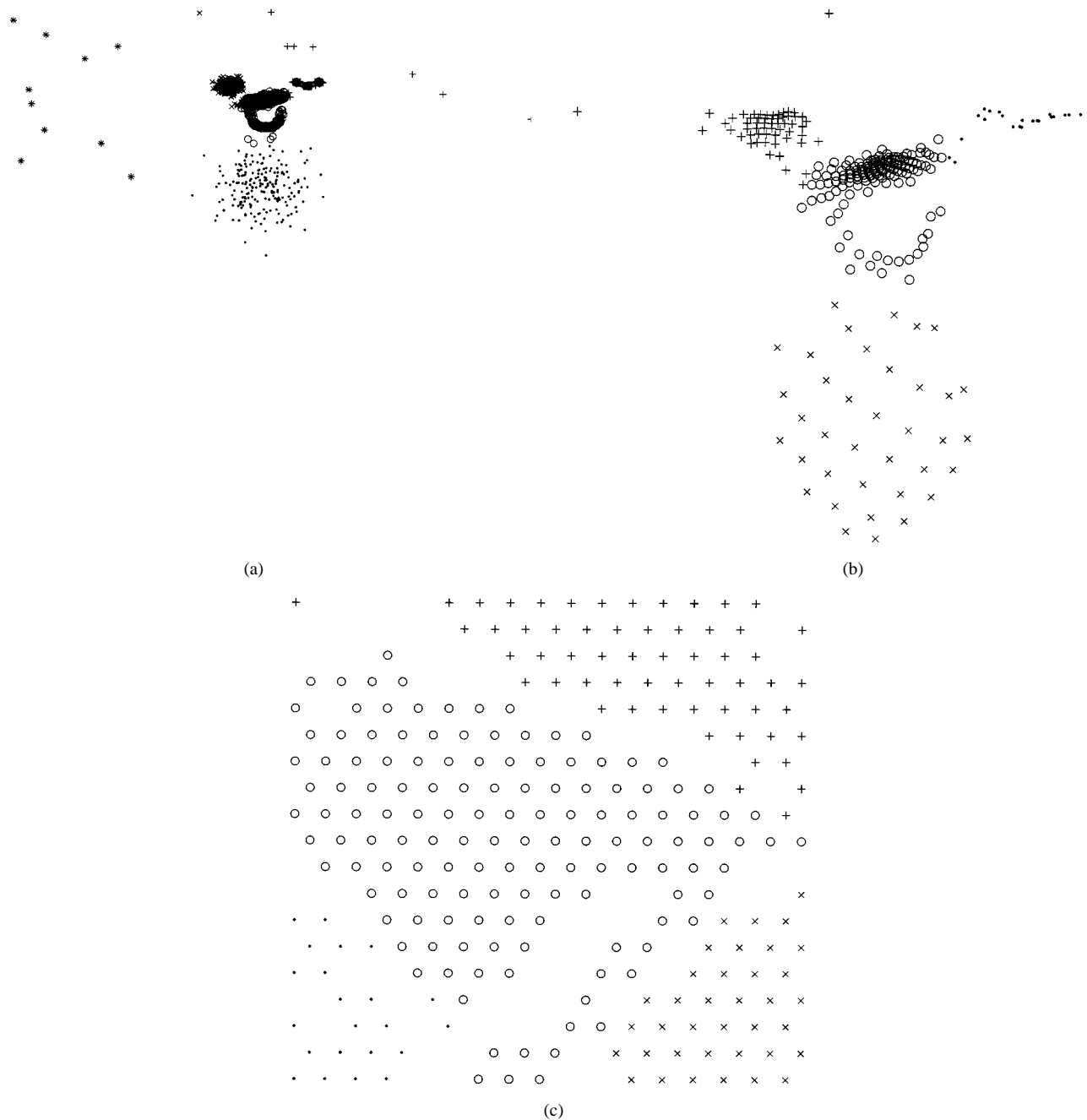


Fig. 10. Results of (a) partitive clustering of the data set I as well as (b) the corresponding SOM in the input space and (c) on the map grid.

they would have been completely separated from the rest of the map.

*Partitive Clustering:* According to the Davies–Bouldin index in Fig. 6(a), the number of clusters is small, say, below ten. The index computed from the clustering result of the data has a negative peak at five clusters, whereas the corresponding index of the prototypes of the SOM has similar peak at four clusters. Actually, at all clusters below 15, there is one cluster shift in the results. The reason for this off-by-one observation is that the direct clustering groups the outliers on the left to one cluster, whereas the two-level approach joins them to the left eye cluster. The clustering results using five clusters for the whole data and four clusters for SOM are presented in Fig. 10.

Both clustering the data directly and clustering of prototype vectors using  $k$ -means gave reasonable results, but the following problems were observed.

- The three small clusters close to each other in the right eye were not discovered as separate clusters but as a single cluster; as a matter of fact, they were not separated until over 20 clusters were used (not shown).
- The nonspherical cluster could not be properly recognized as one cluster.

Both difficulties originate from properties of the  $k$ -means algorithm: It not only seeks for spherical but also clusters with roughly equal number of samples. Therefore, the  $k$ -means should never be blindly used to cluster any data without careful

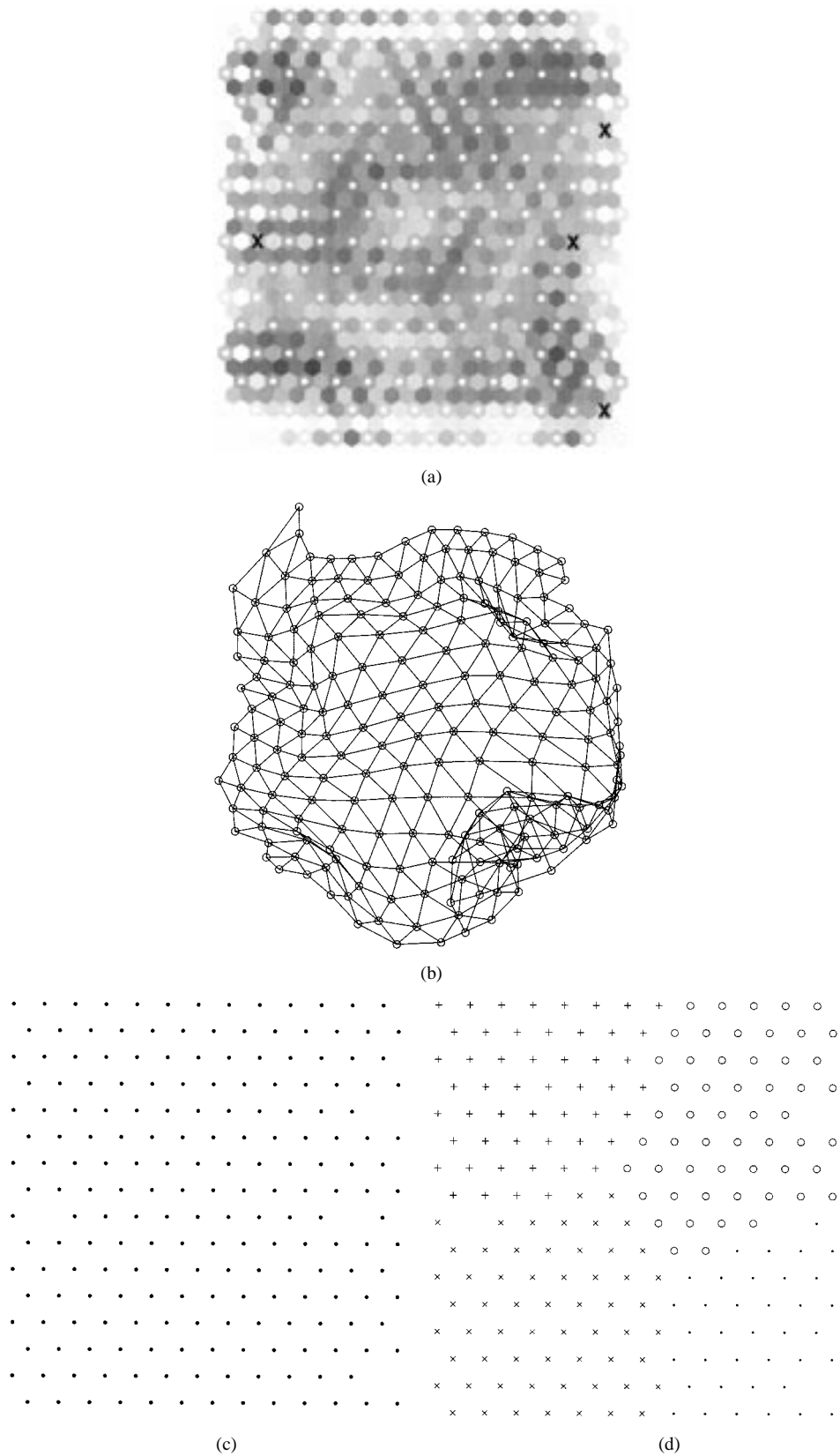


Fig. 11. Clustering of the SOM of data set II. (a) U-matrix. (b) Sammon's mapping. (c) Partitioning using agglomerative clustering. (d) Partitioning using partitive clustering with four clusters. In the Sammon's mapping, the map units have been indicated by circles and neighboring map units have been connected to each other with lines.

verification of the results. It should also be noted that as the number of clusters is increased, the number of samples in

clusters decreases, which makes algorithm more sensitive to outliers.

### E. SOM of Data Set II

U-matrix and Sammon's mapping of the SOM of data set II are presented in Fig. 11(a) and (b). According to them, there are no clusters; neither visualization reveals any cluster structure, and there are only four interpolating units randomly positioned around the map.

In Fig. 11(c), the partitioning obtained by the agglomerative clustering is shown. The result is perfect; except for interpolative units, which were left out of the analysis, all units belong to the same cluster. In the partitive clustering, the Davies–Bouldin index reduced to almost constant level as the number of cluster centers increased; see Fig. 6(b). Even though this does not indicate presence of only one cluster, it suggests that the data is somehow suspicious. Clustering the SOM of data set II using four clusters is shown in Fig. 11(d). As might be expected, the clustering nicely splits the SOM into four parts which are almost of equal size.

### F. SOM of Data Set III

U-matrix and Sammon's mapping of the SOM of data set III are presented in Fig. 12. From the U-matrix, one can clearly distinguish several—maybe about 20—separate areas. The interpolating map units (marked by "X"s on the U-matrix) are mostly positioned on the borders of the U-matrix. Sammon's mapping is badly folded, but in addition, there one can see several concentrations of points. The folding itself indicates that the actual dimension of the data set is much higher than two, which is the output dimension used in the projection.

In Fig. 13(a) and (b), the results of agglomerative clustering of SOM of data set III using average and single linkage are shown. The partitionings agree relatively well with the U-matrix in Fig. 12(a).

The shape of Davies–Bouldin index plot of data set III in Fig. 6(c) is somewhat similar to the one of data set II in Fig. 6(b). However, there are negative peaks at 6, 10, and 23 clusters. In Fig. 13(c) and (d), the partitioning of the SOM to 10 and 23 clusters is shown. It seems that compared with the U-matrix in Fig. 12, both give reasonable interpretations of the SOM at two accuracy levels.

### G. Comparison with Direct Clustering

*Computational Load:* To estimate the computational load, the execution times were measured. The tests were run in Matlab 5 environment in a SGI O2000 computer using one R10000 processor and enough memory to avoid swapping. While the times do not necessarily represent the true times when using optimized compiled code, they give a fair indication of the relative load of different algorithms. The computation times for SOM's and for dendrogram construction and  $k$ -means partitioning to 2–30 clusters (each using 100 runs with different initializations) are listed in Table III. The decrease in computation time is up to two orders of magnitude for agglomerative clustering using average linkage (e.g., 13 hours versus 10 min for data set III). For single and complete linkage, the decrease is smaller but still remarkable, especially for larger data sets (e.g., 27 versus 1 min for data set I). Exact comparison with partitive clustering is more difficult since the number of  $k$ -means runs plays a major role. Fur-

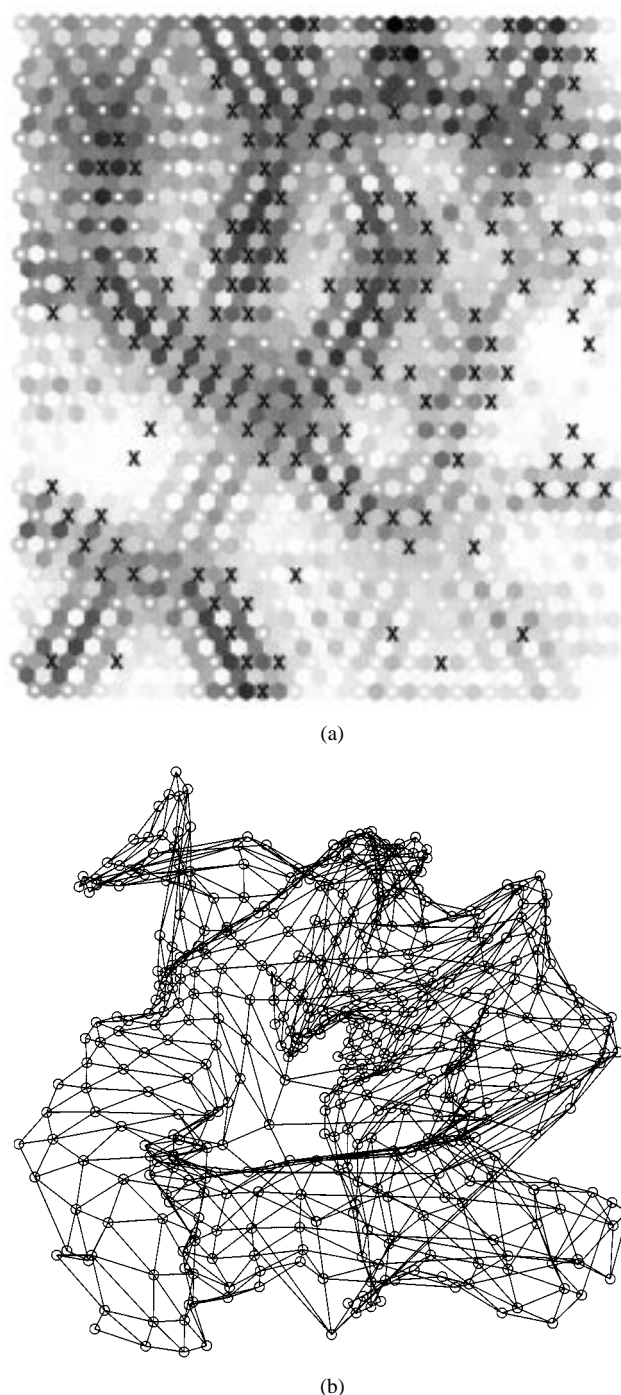


Fig. 12. (a) U-matrix and (b) Sammon's mapping of the SOM of data set III.

thermore, the sequential SOM algorithm was used as opposed to batch  $k$ -means. If batch SOM algorithm had been used, the SOM training times would have been much lower. However, for the testing procedure, the saving was about one order of magnitude (e.g., 84 versus 7 min for data set I).

*Clustering Accuracy:* The clustering results were compared using relative conditional entropy. Entropy quantifies the amount of uncertainty present in the state of a variable or, in this application, uncertainty of the cluster of the samples. Entropy of a clustering  $Q$  is defined as  $H(Q) = -\sum_{i=1}^C p_i \log(p_i)$ , where  $p_i$  is the probability of cluster  $i$ , and  $C$  is the number of clusters in  $Q$ . The entropy gets its maximum value  $\log(C)$  when

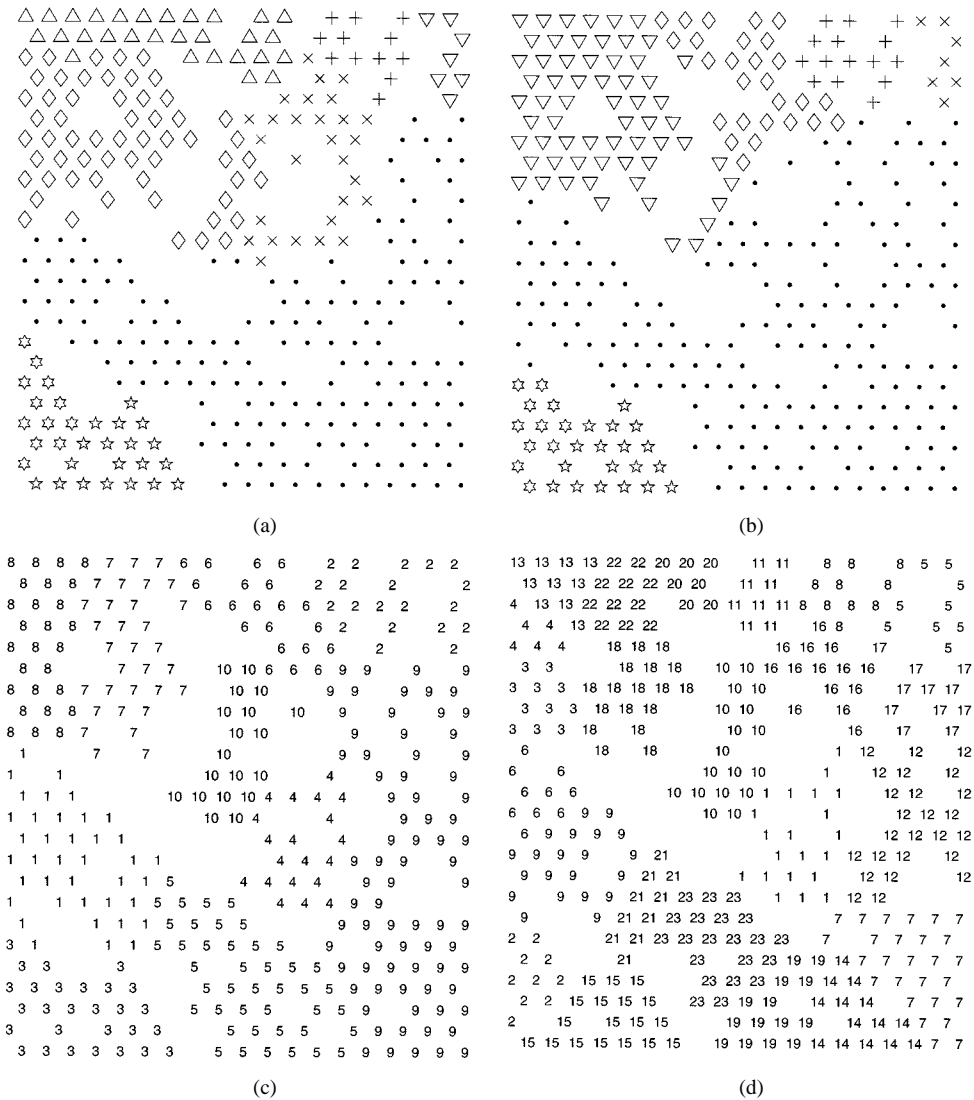


Fig. 13. Partitioning of the SOM of data set III. (a) With agglomerative clustering using average linkage and (b) single linkage. (c) With partitive clustering using 10 and (d) 23 clusters.

each cluster is equally probable. Relative entropy is defined  $H_r = H/\log(C)$  and has values between  $[0, 1]$ . As an example, a binary variable with probabilities  $p_1 = \{0.5, 0.1, 0.02, 0.01\}$  has relative entropies  $H_r = \{1, 0.47, 0.14, 0.08\}$ , respectively. The entropy measure used in the tests was relative conditional entropy  $H_r(\mathbf{X}|\mathbf{Y})$ , which describes uncertainty of clustering  $\mathbf{X}$  if clustering  $\mathbf{Y}$  is known

$$H_r(\mathbf{X}|\mathbf{Y}) = -\frac{1}{\log(C_X)} \sum_{j=1}^{C_Y} p_j \sum_{i=1}^{C_X} p_{ij} \log(p_{ij}), \quad (7)$$

where  $i$  and  $j$  denote clusters, and  $C_X$  and  $C_Y$  are the total number of clusters in  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively.

In Table IV, the direct and SOM-based clustering results of data set I are compared with the true known partitioning and with each other. In general, the match is very good, although the results are slightly worse for single and complete linkage. Fig. 14 shows the comparison results for partitive clustering of each data set and for agglomerative clustering of data set III.

TABLE IV  
RELATIVE CONDITIONAL ENTROPIES FOR DATA SET I.  $\mathbf{X}$  IS THE TRUE PARTITIONING OF THE DATA (OUTLIERS HAVE BEEN IGNORED),  $\mathbf{Y}$  IS THE DIRECT PARTITIONING OF THE DATA, AND  $\mathbf{Z}$  IS THE PARTITIONING OBTAINED BY CLUSTERING THE SOM

clustering	$H_r(\mathbf{X} \mathbf{Y})$	$H_r(\mathbf{X} \mathbf{Z})$	$H_r(\mathbf{Y} \mathbf{Z})$
average linkage	0.020	0.032	0.018
single linkage	0.024	0.031	0.34
complete linkage	0.0098	0.14	0.28
$k$ -means ( $k = 4$ )	0.22	0.22	0.029
$k$ -means ( $k = 7$ )	0.12	0.12	0.073

In all cases, knowledge of the SOM clustering result has significantly reduced the uncertainty. For data sets I and III, the relative conditional entropies indicate a good correspondence between the two clustering results, but for data set II, the results are poor. However, there is a natural explanation. Because data set II actually had no clusters, there is no reason why cluster centers from different clustering runs would match each other, especially with low number of clusters.

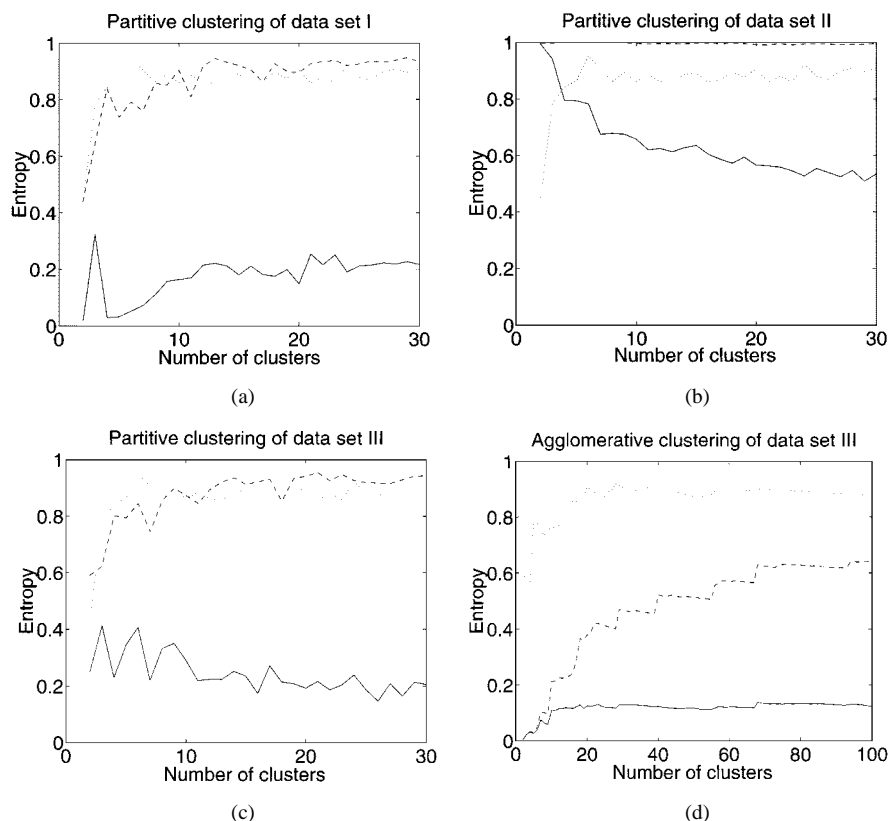


Fig. 14. Results for partitive clustering of (a) data set I, (b) II (c) III and (d) agglomerative clustering of data set III. Dashed line is  $H_r(\mathbf{Y})$ , dotted line  $H_r(\mathbf{Z})$ , and solid line  $H_r(\mathbf{Y}|\mathbf{Z})$ , where  $\mathbf{Y}$  and  $\mathbf{Z}$  are direct and SOM-based partitionings, respectively, with equal number of clusters. The agglomerative clustering in (d) was obtained using complete linkage. Results for single and average linkage were similar. Partitioning was selected by simply cutting the dendrogram at a certain level. From (d), the degrading effect of outliers on agglomerative clustering can be seen; the relative entropy  $H_r(\mathbf{Y})$  is very low when the number of clusters is small. This is because vast majority of the samples belongs to a single cluster and only a few samples—the outliers—belong to the others. Partitioning  $\mathbf{Z}$  based on the SOM is much more evenly distributed.

In the light of these results, correspondence between direct and two-level clustering can be considered acceptable.

## V. CONCLUSIONS

In our earlier work, it has become obvious that the SOM is an effective platform for visualization of high-dimensional data. However, to be able to fully understand contents of a data set, it is vital to find out if the data has cluster structure. If this is the case, the clusters need to be extracted to be able to fully exploit the properties of the data set by producing summary information. For each cluster, not only means (or medians) and ranges of single variables are interesting but also aspects like the following.

- Which variable(s) make the cluster different from the neighboring clusters?
- Which factors make the cluster different from the rest of the data?
- What is effective data dimension of the cluster?
- Is the cluster spherical or an outlier, and does it have sub-clusters?
- What are the dependencies between variables in the clusters?

The purpose of this paper was to evaluate if the data abstraction created by the SOM could be utilized in clustering of data. The SOM as the first abstraction level in the clustering has some

clear advantages. First, the original data set is represented using a smaller set of prototype vectors, which allows efficient use of clustering algorithms to divide the prototypes into groups. The reduction of the computational cost is especially important for hierarchical algorithms allowing clusters of arbitrary size and shape. Second, the 2-D grid allows rough visual presentation and interpretation of the clusters.

In the experiments, agglomerative and partitive ( $k$ -means) clustering algorithms were run both directly for data and for SOM trained using the data using three data sets. The experiments indicated that clustering the SOM instead of directly clustering the data is computationally effective approach. The clustering results using SOM as an intermediate step were also comparable with the results obtained directly from the data.

## REFERENCES

- [1] D. Pyle, *Data Preparation for Data Mining*. San Francisco, CA: Morgan Kaufmann, 1999.
- [2] T. Kohonen, *Self-Organizing Maps*. Berlin/Heidelberg, Germany: Springer, 1995, vol. 30.
- [3] J. Vesanto, "SOM-based data visualization methods," *Intell. Data Anal.*, vol. 3, no. 2, pp. 111–126, 1999.
- [4] J. C. Bezdek and S. K. Pal, Eds., *Fuzzy Models for Pattern Recognition: Methods that Search for Structures in Data*. New York: IEEE, 1992.
- [5] G. J. McLahlan and K. E. Basford, *Mixture Models: Inference and Applications to Clustering*. New York: Marcel Dekker, 1987, vol. 84.
- [6] J. C. Bezdek, "Some new indexes of cluster validity," *IEEE Trans. Syst., Man, Cybern. B*, vol. 28, pp. 301–315, 1998.

- [7] M. Blatt, S. Wiseman, and E. Domany, "Superparamagnetic clustering of data," *Phys. Rev. Lett.*, vol. 76, no. 18, pp. 3251–3254, 1996.
- [8] G. Karypis, E.-H. Han, and V. Kumar, "Chameleon: Hierarchical clustering using dynamic modeling," *IEEE Comput.*, vol. 32, pp. 68–74, Aug. 1999.
- [9] J. Lampinen and E. Oja, "Clustering properties of hierarchical self-organizing maps," *J. Math. Imag. Vis.*, vol. 2, no. 2–3, pp. 261–272, Nov. 1992.
- [10] E. Boudaillier and G. Hebrail, "Interactive interpretation of hierarchical clustering," *Intell. Data Anal.*, vol. 2, no. 3, 1998.
- [11] J. Buhmann and H. Kühnel, "Complexity optimized data clustering by competitive neural networks," *Neural Comput.*, vol. 5, no. 3, pp. 75–88, May 1993.
- [12] G. W. Milligan and M. C. Cooper, "An examination of procedures for determining the number of clusters in a data set," *Psychometrika*, vol. 50, no. 2, pp. 159–179, June 1985.
- [13] D. L. Davies and D. W. Bouldin, "A cluster separation measure," *IEEE Trans. Patt. Anal. Machine Intell.*, vol. PAMI-1, pp. 224–227, Apr. 1979.
- [14] S. P. Luttrell, "Hierarchical self-organizing networks," in *Proc. 1st IEE Conf. Artificial Neural Networks*, London, U.K., 1989, pp. 2–6.
- [15] A. Varfis and C. Versino, "Clustering of socio-economic data with Kohonen maps," *Neural Network World*, vol. 2, no. 6, pp. 813–834, 1992.
- [16] P. Mangiameli, S. K. Chen, and D. West, "A comparison of SOM neural network and hierarchical clustering methods," *Eur. J. Oper. Res.*, vol. 93, no. 2, Sept. 1996.
- [17] T. Kohonen, "Self-organizing maps: Optimization approaches," in *Artificial Neural Networks*, T. Kohonen, K. Mäkisara, O. Simula, and J. Kangas, Eds. Amsterdam, The Netherlands: Elsevier, 1991, pp. 981–990.
- [18] J. Moody and C. J. Darken, "Fast learning in networks of locally-tuned processing units," *Neural Comput.*, vol. 1, no. 2, pp. 281–294, 1989.
- [19] J. A. Kangas, T. K. Kohonen, and J. T. Laaksonen, "Variants of self-organizing maps," *IEEE Trans. Neural Networks*, vol. 1, pp. 93–99, Mar. 1990.
- [20] T. Martinez and K. Schulten, "A 'neural-gas' network learns topologies," in *Artificial Neural Networks*, T. Kohonen, K. Mäkisara, O. Simula, and J. Kangas, Eds. Amsterdam, The Netherlands: Elsevier, 1991, pp. 397–402.
- [21] B. Fritzke, "Let it grow—Self-organizing feature maps with problem dependent cell structure," in *Artificial Neural Networks*, T. Kohonen, K. Mäkisara, O. Simula, and J. Kangas, Eds. Amsterdam, The Netherlands: Elsevier, 1991, pp. 403–408.
- [22] Y. Cheng, "Clustering with competing self-organizing maps," in *Proc. Int. Conf. Neural Networks*, vol. 4, 1992, pp. 785–790.
- [23] B. Fritzke, "Growing grid—A self-organizing network with constant neighborhood range and adaptation strength," *Neural Process. Lett.*, vol. 2, no. 5, pp. 9–13, 1995.
- [24] J. Blackmore and R. Miikkulainen, "Visualizing high-dimensional structure with the incremental grid growing neural network," in *Proc. 12th Int. Conf. Machine Learning*, A. Prieditis and S. Russell, Eds., 1995, pp. 55–63.
- [25] S. Jockusch, "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, and G. Hauske, Eds. Amsterdam, The Netherlands: Elsevier, 1990, pp. 169–172.
- [26] J. S. Rodrigues and L. B. Almeida, "Improving the learning speed in topological maps of pattern," in *Proc. Int. Neural Network Conf.*, 1990, pp. 813–816.
- [27] D. Alahakoon and S. K. Halgamuge, "Knowledge discovery with supervised and unsupervised self evolving neural networks," in *Proc. Fifth Int. Conf. Soft Comput. Inform./Intell. Syst.*, T. Yamakawa and G. Matsumoto, Eds., 1998, pp. 907–910.
- [28] A. Gersho, "Asymptotically optimal block quantization," *IEEE Trans. Inform. Theory*, vol. IT-25, pp. 373–380, July 1979.
- [29] P. L. Zador, "Asymptotic quantization error of continuous signals and the quantization dimension," *IEEE Trans. Inform. Theory*, vol. IT-28, pp. 139–149, Mar. 1982.
- [30] H. Ritter, "Asymptotic level density for a class of vector quantization processes," *IEEE Trans. Neural Networks*, vol. 2, pp. 173–175, Jan. 1991.
- [31] T. Kohonen, "Comparison of SOM point densities based on different criteria," *Neural Comput.*, vol. 11, pp. 2171–2185, 1999.
- [32] R. D. Lawrence, G. S. Almasi, and H. E. Rushmeier, "A scalable parallel algorithm for self-organizing maps with applications to sparse data problems," *Data Mining Knowl. Discovery*, vol. 3, no. 2, pp. 171–195, June 1999.
- [33] T. Kohonen, S. Kaski, K. Lagus, J. Salojärvi, J. Honkela, V. Paatero, and A. Saarela, "Self organization of a massive document collection," *IEEE Trans. Neural Networks*, vol. 11, pp. XXX–XXX, May 2000.
- [34] P. Koikkalainen, "Fast deterministic self-organizing maps," in *Proc. Int. Conf. Artificial Neural Networks*, vol. 2, 1995, pp. 63–68.
- [35] A. Ultsch and H. P. Siemon, "Kohonen's self organizing feature maps for exploratory data analysis," in *Proc. Int. Neural Network Conf.* Dordrecht, The Netherlands, 1990, pp. 305–308.
- [36] M. A. Kraaijveld, J. Mao, and A. K. Jain, "A nonlinear projection method based on Kohonen's topology preserving maps," *IEEE Trans. Neural Networks*, vol. 6, pp. 548–559, May 1995.
- [37] J. Iivarinen, T. Kohonen, J. Kangas, and S. Kaski, "Visualizing the clusters on the self-organizing map," in *Proc. Conf. Artificial Intell. Res. Finland*, C. Carlsson, T. Järvi, and T. Reponen, Eds. Helsinki, Finland, 1994, pp. 122–126.
- [38] X. Zhang and Y. Li, "Self-organizing map as a new method for clustering and data analysis," in *Proc. Int. Joint Conf. Neural Networks*, 1993, pp. 2448–2451.
- [39] J. B. Kruskal, "Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis," *Psychometrika*, vol. 29, no. 1, pp. 1–27, Mar. 1964.
- [40] —, "Nonmetric multidimensional scaling: A numerical method," *Psychometrika*, vol. 29, no. 2, pp. 115–129, June 1964.
- [41] J. W. Sammon, Jr., "A nonlinear mapping for data structure analysis," *IEEE Trans. Comput.*, vol. C-18, pp. 401–409, May 1969.
- [42] P. Demartines and J. Héroult, "Curvilinear component analysis: A self-organizing neural network for nonlinear mapping of data sets," *IEEE Trans. Neural Networks*, vol. 8, pp. 148–154, Jan. 1997.
- [43] A. Varfis, "On the use of two traditional statistical techniques to improve the readability of Kohonen Maps," in *Proc. NATO ASI Workshop Statistics Neural Networks*, 1993.
- [44] S. Kaski, J. Venna, and T. Kohonen, "Coloring that reveals high-dimensional structures in data," in *Proc. 6th Int. Conf. Neural Inform. Process.*, 1999, pp. 729–734.
- [45] F. Murtagh, "Interpreting the Kohonen self-organizing map using contiguity-constrained clustering," *Patt. Recognit. Lett.*, vol. 16, pp. 399–408, 1995.
- [46] O. Simula, J. Vesanto, P. Vasara, and R.-R. Helminen, *Industrial Applications of Neural Networks*, L. C. Jain and V. R. Vemuri, Eds. Boca Raton, FL: CRC, 1999, ch. 4, pp. 87–112.
- [47] J. Vaisey and A. Gersho, "Simulated annealing and codebook design," in *Proc. Int. Conf. Acoust., Speech, Signal Process.*, 1988, pp. 1176–1179.
- [48] J. K. Flanagan, D. R. Morrell, R. L. Frost, C. J. Read, and B. E. Nelson, "Vector quantization codebook generation using simulated annealing," in *Proc. Int. Conf. Acoust., Speech, Signal Process.*, 1989, vol. 3, pp. 1759–1762.
- [49] T. Graepel, M. Burger, and K. Obermayer, "Phase transitions in stochastic self-organizing maps," *Phys. Rev. E*, vol. 56, pp. 3876–3890, 1997.



**Juha Vesanto** received the M.Sc. (Tech.) degree in technical physics from the Helsinki University of Technology (HUT), Helsinki, Finland, in 1997. He is currently pursuing the Ph.D. degree and doing research at the Neural Networks Research Centre, HUT.

His research interests focus on data mining (especially visualization) and analysis of multivariate data.



**Esa Alhoniemi** (S'97) received the M.Sc. (Tech.) and the Lic.Sc. (Tech.) degrees in computer science and engineering from the Helsinki University of Technology (HUT), Helsinki, Finland, in 1995 and 1998, respectively. He is pursuing the Ph.D. degree from HUT with a dissertation on applications of neural networks in the analysis of industrial processes.

Since 1995, he has been working as a Researcher with the Neural Networks Research Centre, HUT.