

## Publication 1

Samuel Kaski, Jarkko Venna, and Teuvo Kohonen. Coloring that reveals cluster structures in multivariate data. *Australian Journal of Intelligent Information Processing Systems*, 6:82–88, 2000.

© 2000 AJIIPS. Reprinted with the kind permission of the Australian Journal of Intelligent Information Processing Systems.

# Coloring that Reveals Cluster Structures in Multivariate Data

Samuel Kaski, Jarkko Venna, and Teuvo Kohonen

Helsinki University of Technology, Neural Networks Research Centre  
P.O. Box 5400, FIN-02015 HUT, Finland  
samuel.kaski@hut.fi

## Abstract

A method is introduced for assigning colors to displays of cluster structures of high-dimensional data, such that the perceptual differences of the colors reflect the distances in the original data space as faithfully as possible. The method consists of three parts: First the cluster structures are discovered with the Self-Organizing Map (SOM), and then a new nonlinear projection method is applied to map the cluster structures into the CIE Lab color space. Finally the cluster structures are visualized using the colors found by the projection. The projection method preserves best the local data distances that are the most important ones, while ensuring that the global order is still discernible from the colors, too. This allows the method to conform flexibly to the available color space. The output space of the projection need not necessarily be the color space. Projections onto, say, two dimensions can be visualized as well.

## 1 Introduction

In exploratory data analysis or interactive data mining one central goal is to illustrate high-dimensional data sets by overviews that are easily understandable but still preserve the essential properties of the data. The Self-Organizing Map (SOM) algorithm [3, 4] can be used for visualizing one of the most central properties of high-dimensional statistical data, namely its cluster structure, on graphical map displays.

In this paper we introduce a method for visualizing the cluster structures discovered by the SOM with colors. Unlike in commonly used visualizations of data types or clusters the (relative) colors are not chosen arbitrarily, but the perceptual differences of the colors are supposed to reflect the distance relations within the cluster structure as faithfully as possible.

The perceptual differences between different colors can be approximated by distances in suitably defined color spaces, for example the CIE Lab color space [1]. Therefore, any coloring of high-dimensional data items actually corresponds to a mapping of the items into a smaller-dimensional space, the color space. In this work we do not map the data items themselves but their representations, the models formed by the SOM algorithm. A nonlinear projection method that is suitable for the SOM models will be introduced in Sec-

This work was supported by the Academy of Finland.

tion 3, and applied to coloring cluster structures in Section 4.

## 2 Discovering Cluster Structures with the Self-Organizing Map

The SOM [3, 4] is an unsupervised neural network algorithm that can be used to visualize cluster structures in high-dimensional data. The SOM consists of a regular grid of neurons (*units*) that are all connected to the input of the SOM. Each of the SOM units contains a model vector  $\mathbf{m}_i \in \mathbb{R}^n$ , where  $n$  is the dimension of the data. When a data item  $\mathbf{x} \in \mathbb{R}^n$  is input to the SOM, the best matching model vector,  $\mathbf{m}_c$ , is searched for using the equation

$$c = \operatorname{argmin}_i \{\|\mathbf{x} - \mathbf{m}_i\|\}. \quad (1)$$

When the best matching model vector is found, the model vectors are updated using the following equation:

$$\mathbf{m}_i(t+1) = \mathbf{m}_i(t) + h_{ci}(t)[\mathbf{x}(t) - \mathbf{m}_i(t)], \quad (2)$$

where  $t = 0, 1, 2, \dots$  is the number of the learning step. The function  $h_{ci}(t) = h(\|\mathbf{r}_c - \mathbf{r}_i\|; t)$ , where  $\mathbf{r}_i \in \mathbb{R}^2$  and  $\mathbf{r}_c \in \mathbb{R}^2$  are the location vectors on the SOM grid for units  $i$  and  $c$ , is the neighborhood function. The neighborhood function controls how each unit is updated in each step. As the result of the SOM algorithm the model vectors become ordered: neighboring units in the grid represent similar data. Additionally the density of the model vectors reflects the density of the data in the input space.

The SOM can often be computed faster by using a batch version of the SOM algorithm instead of the online algorithm described above. The batch version consists of the following steps [5]:

1. Initialize the model vectors  $\mathbf{m}_i$ . (Suitable initialization methods can be found for example in [4]).
2. For each unit  $j$  compute the average of the data vectors that the unit  $j$  is the best match for. Denote this average with  $\bar{\mathbf{x}}_j$ .
3. Compute new values for the model vectors  $\mathbf{m}_i$  using the equation

$$\mathbf{m}_i = \frac{\sum_j n_j h_{ji} \bar{\mathbf{x}}_j}{\sum_j n_j h_{ji}}, \quad (3)$$

where  $j$  goes through all model vectors. The term  $h_{ij}$  is the neighborhood function of the SOM and  $n_j$  is the number of data vectors that the unit  $j$  is the best match for.

4. Repeat the steps two and three until the algorithm converges.

Once the SOM has been computed, the ordered SOM grid can be used as a groundwork for displaying the cluster structure of the data. Clusteredness, i.e. the density of the data, is reflected in the density of the model vectors, which in turn is reflected in the distances between neighboring model vectors. This idea has been utilized in the U-matrix method [8] in which the distances are visualized by gray levels. Figures 2a and 3a illustrate the cluster structure of two data sets. The first set consists of 39 statistical indicators that describe various aspects of poverty of each country in the world [2]. The second set consists of cepstral-coefficient vectors picked up from continuous Finnish speech of one speaker (test material available in [6]).

### 3 SOM-Based Non-Linear Projection that Preserves Local Distances

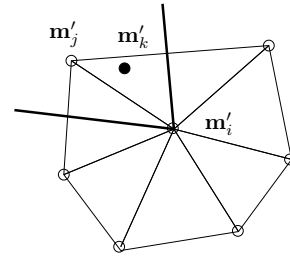
When projecting a set of high-dimensional data vectors into a lower-dimensional space it is in general impossible to preserve *all* of the pairwise distances between the data items. In the mapping of Sammon (cf. the end of this section), a compromise is made by minimizing an objective function that describes the *global error*, errors in all distances.

Fortunately, preservation of all distances is not equally important. When aiming at a projection that preserves the cluster structures, i.e. the density of the data, it is especially important to preserve the *local distances* that are indicators of the local density. In the SOM the model vectors are already ordered so that the distances between model vectors of neighboring map units are the local distances. This is the reason why distances between neighbors are used in the SOM-based clustering displays to indicate the clusteredness.

The requirements for a good projection can be summed up as follows:

- i The perceived color differences should reflect the distances between model vectors of neighboring units as well as possible.
- ii The SOM grid should still be ordered. Otherwise different non-neighboring areas on the map may attain the same color.

We introduce below a new nonlinear projection method that aims at fulfilling these requirements.



**Figure 1:** Definition of orderliness of the projection. The projections of the model vectors in the neighborhood of the unit  $i$  have been depicted with open circles, and neighbors have been connected with thin lines. The most obvious definition of orderliness might be that the map may not become “twisted”, i.e., no non-neighbor may enter the area bordered by the outer thin lines. This definition has disadvantages, however: it is not easily applicable, especially in higher-dimensional spaces, and not as obvious as here if the neighborhood is very irregular. We have therefore used a more straightforward definition: If a non-neighbor  $m'_k$  belongs to the “sector” of  $m'_j$ , illustrated with the thick lines, and is closer to  $m'_i$  than  $m'_j$  is, then  $m'_k$  is too close to  $m'_i$  and the projection is not ordered. The sector consists of those points for which the angle from  $m'_j$ , as seen from  $m'_i$ , is smaller than from the other neighbors of the middle point  $m'_i$ .

Denote the index set corresponding to the neighborhood of map unit  $i$  by  $N_i$  (here excluding the map unit  $i$  itself), and denote the distance between the model vectors  $i$  and  $j$  by  $d_{ij} = \|\mathbf{m}_i - \mathbf{m}_j\|$ . If the low-dimensional projection of  $\mathbf{m}_i$  is denoted by  $\mathbf{m}'_i$ , and the distance after the projection by  $d'_{ij} = \|\mathbf{m}'_i - \mathbf{m}'_j\|$ , then

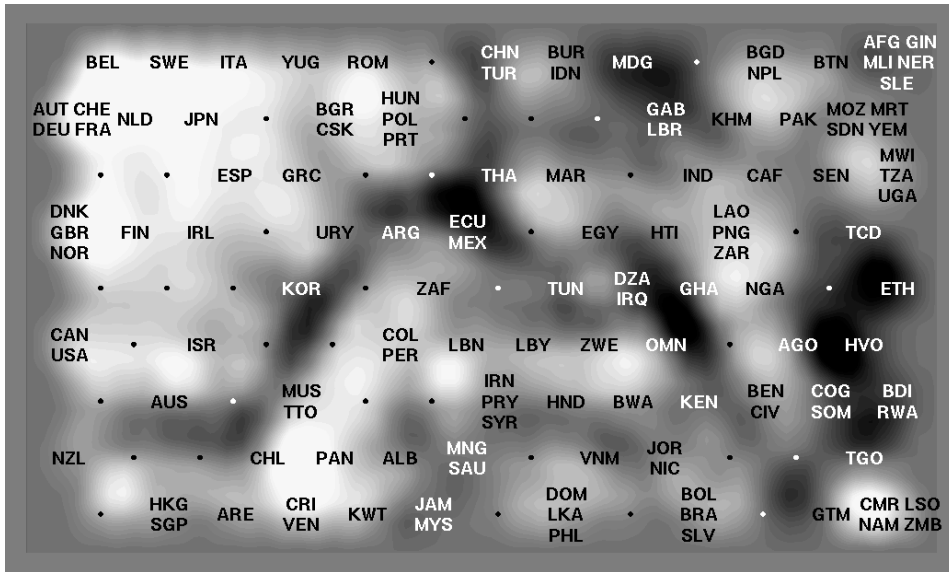
$$E_1 = \sum_i \sum_{j \in N_i} (d_{ij} - d'_{ij})^2 \quad (4)$$

is the mean-square cost function that measures distortions in the local distances. The projection method should find such values for the  $\mathbf{m}'_i$  that the cost function  $E_1$  is minimized.

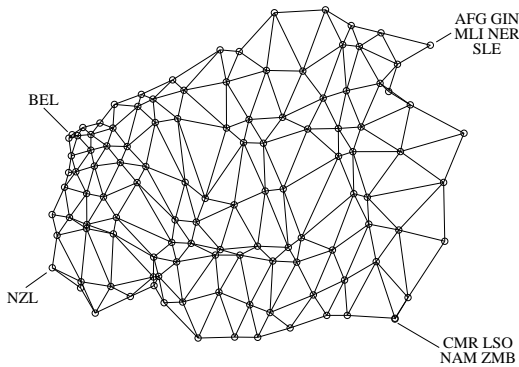
Unfortunately, a projection that preserves the local distances is not guaranteed to be globally ordered. The projection may “fold” over itself and become unintelligible.

It has turned out in our experiments that the projection becomes globally ordered if an only slightly restrictive extra constraint is imposed on the non-local distances: Intuitively speaking, no other model vector may be projected in between the model vectors belonging to neighboring points on the SOM grid. The detailed definition is given and illustrated in Figure 1. Note that the exact value of the non-local distances may be arbitrary as long as they are large enough.

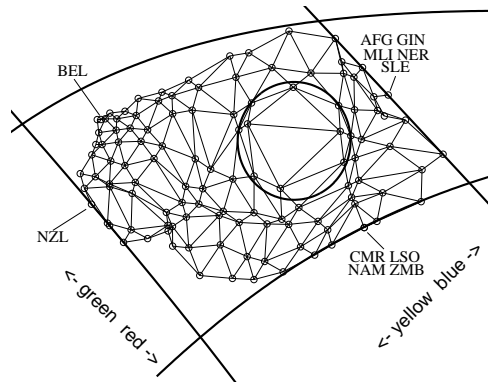
The orderliness can be guaranteed by adding a second term into the cost function (4). Denote by  $I_{ij}$  the set of indices of all such model vectors that do not satisfy the orderliness constraint for the pair of neighbors  $(i, j)$ . The total squared



a



b



c

**Figure 2:** The poverty map. **a** Structured diagram of the data set [9] chosen to describe the standard of living. The order of the abbreviated country names indicates the similarity of the standard of living of the countries, and the gray levels, obtained by smoothing the U-matrix, indicate the degree of clustering. Light shades correspond to a high degree of clustering and dark shades to gaps in the degree of clustering. Different types of welfare and poverty are visible as the clustered (light) areas on the map, separated by the dark “ravines”. The SOM grid points having no countries have been marked with dots. **b** Projection of the poverty map onto the two-dimensional space. The small circles denote the projections of the model vectors. Projections of model vectors of neighboring map units have been connected with thin lines and the countries mapped to the corner units have been marked in the figure. **c** The projection of the model vectors of the poverty map onto a constant-lightness cross section of the CIELab color space. The thick lines delimit the region representable by a typical CRT tube, and the circle in the middle encompasses colors that were considered too non-saturated (gray). Note that the axes have been rotated and mirrored to ease comparison with **a** and **b**.

deviation from orderliness can then be measured with

$$E_2 = \sum_{i,j} \sum_{k \in I_{i,j}} (d'_{ij} - d'_{ik})^2. \quad (5)$$

The total cost function to be minimized is then

$$E' = E_1 + \lambda_2 E_2. \quad (6)$$

Here  $\lambda_2$  is a parameter that determines how strongly deviations from orderliness affect the projection. We recommend initializing the projection in an orderly fashion to avoid local minima, and setting  $\lambda_2$  to a relatively large value, say,  $\lambda_2 = 100$ .

The projection of the SOMs of Figures 2a and 3a into a two-dimensional space, obtained by minimizing  $E'$  in (6), has been shown in Figures 2b and 3b. We optimized the cost function by a simple stochastic gradient descent: At each optimization step one model vector was selected randomly and moved into the direction of the negative gradient of  $E'$  by an amount that decreased gradually during the optimization. If the gradient step would have increased the value of the cost function, the step size was temporarily halved until either the value of the cost function decreased or the maximum number of allowed iterations was reached.

It can be seen from Figures 2b and 3b that the long distances between the projected model vectors correspond to black areas in Figures 2a and 3a and short distances to light areas, respectively.

**Relation to Sammon's Mapping.** The traditional multidimensional scaling methods like the Sammon's mapping [7] could in principle be used instead of the method presented here. They do not, however, produce as flexible mappings as our method, since they try to represent *all* of the pairwise distances of the original data. The local differences will then necessarily be represented less accurately in Sammon's mapping, except in special cases. Moreover, if only a restricted area of the output space is available like in the coloring task discussed in Section 4 below, then it may be more difficult to fit Sammon's projection to the area because the mapping is "stiffer".

*Note:* Sammon's mapping could be applied to the model vectors of the SOM as well. The cost function of Sammon's mapping would then be

$$E_S = \sum_{i \neq j} \frac{1}{d_{ij}} (d_{ij} - d'_{ij})^2,$$

(omitting a constant normalizing term). Sammon's mapping would consider distances between *all* model vector pairs  $(i, j)$ , weighted by  $1/d_{ij}$ , whereas for the present purpose we consider only the distances between model vectors of neighboring SOM units. Therefore, Sammon's mapping would still be less flexible.

## 4 Coloring by Projection into the Color Space

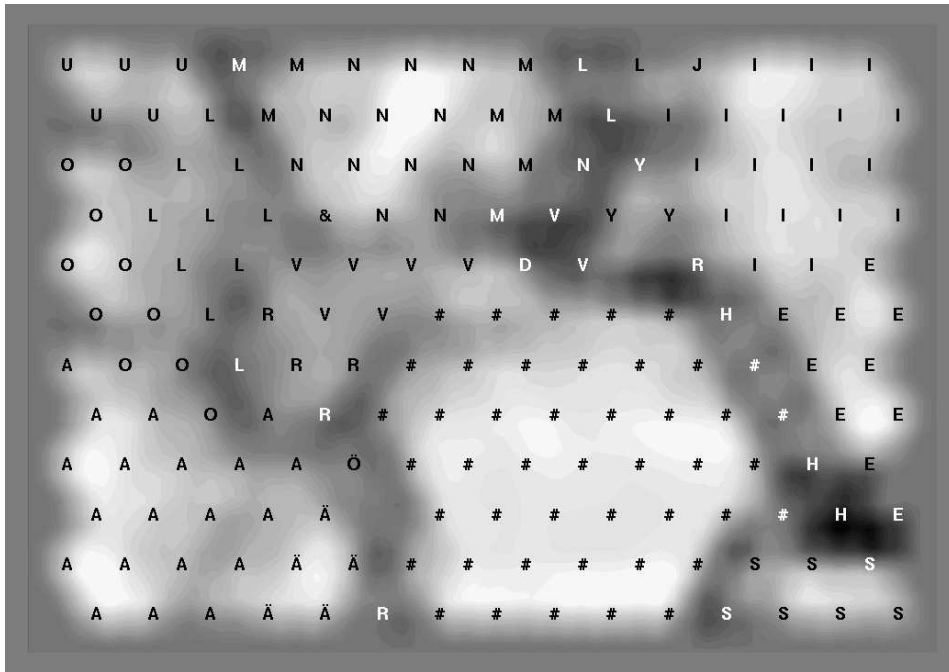
Coloring corresponds to projection into a color space. There exist several color spaces that could be used as targets for projection. One of the color spaces most often used in visualization is RGB which is based on CRT display mechanics. Each component in the RGB corresponds to the intensity of one of the phosphors on a CRT screen. RGB color space has two drawbacks however. First it is device dependent, i.e. a color specified in the RGB color space looks different when viewed on different displays. The second drawback is the perceived nonuniformity of color differences in different parts of the color space. The perceived color difference between a pair of colors that have the same Euclidean distance in the RGB color space varies widely depending on the location of the pair in the color space. Luckily there is a group of color spaces, the uniform color spaces, that are defined so that perceptual differences in colors correspond to Euclidean distances in the space as well as possible, at least when small color differences are considered. We have chosen to use a color space called CIELab [1], which is a well known uniform color space recommended by the Commission Internationale de L'Éclairage (CIE).

The whole three-dimensional space is not available for the coloring, however. A given display device like the CRT tube is able to present only a part of the space. We have restricted the available space even further since we wanted the colors to differ predominantly in one perceptual quality only, namely, the hue. Therefore, very non-saturated colors (more exactly, colors with a small chroma) are not allowed, and the lightness was set to a fixed value. Restricting the change to one perceptual component makes visualizations based on the coloring easier to interpret. After these restrictions a non-regularly shaped area in a two-dimensional cross section of the CIELab space remained (cf., e.g., Fig. 3c). To take these practical points into account we have to add two sub-requirements for the coloring method in addition to those mentioned in section 3:

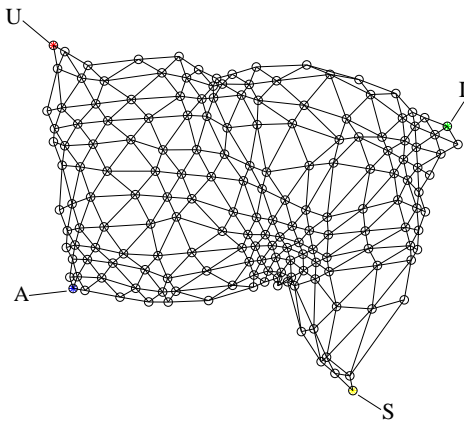
- iii The coloring should consist only of colors that belong to the available color space. All colors should be producible by the display device and preferably differ only in the hue.
- iv The available color space should be covered as well as possible.

We need an additional term in the cost function of the projection to keep the projected items within the available region. Before formulating the term we need to *scale* the distances of the CIELab space suitably, however, since the relative scale between the distances in the original data space and in the CIELab space may be arbitrary. A proper choice of the scaling factor, denoted by  $\kappa$ , will be discussed in more detail later in connection with Figure 4..

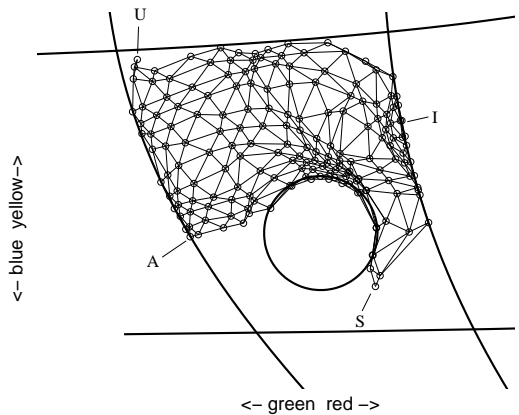
Denote, as previously, the  $i$ th model vector of the SOM by



a



b



c

**Figure 3:** The phonetic map. **a** Structured diagram of the data set consisting of cepstral-coefficient vectors picked up from continuous Finnish speech of one speaker [6]. The gray levels have been obtained by smoothing the U-matrix. The units have been labeled with the phoneme that occur most frequently within the node. **b** Projection of the phonetic map onto the two-dimensional space. **c** Projection of the model vectors of the phonetic map onto a constant-lightness cross section of the CIE Lab color space. For explanation of the figures see caption of Figure 2.

$\mathbf{m}_i$  and its projection by  $\mathbf{m}'_i$ , respectively. Denote then the point in the available space that is closest to  $\mathbf{m}'_i/\kappa$  by  $\bar{\mathbf{m}}_i$ . If  $\mathbf{m}'_i/\kappa$  is already in the available space, then  $\bar{\mathbf{m}}_i = \mathbf{m}'_i/\kappa$ . The additional term,  $E_3$ , in the cost function is then the squared distance between the actual projection and the closest available point, scaled to match the distances in the original data space,

$$E_3 = \sum_i \|\kappa \bar{\mathbf{m}}_i - \mathbf{m}'_i\|^2. \quad (7)$$

The total cost function that the projection shall minimize is

$$E = E_1 + \lambda_2 E_2 + \lambda_3 E_3, \quad (8)$$

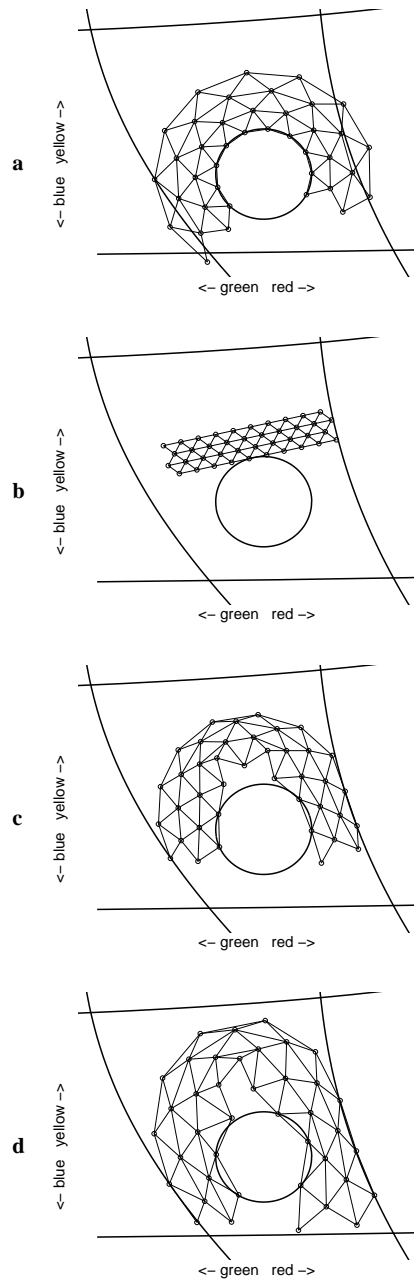
where  $\lambda_2$  and  $\lambda_3$  are suitably chosen parameters. We have minimized  $E$  by the stochastic gradient descent method described in Section 3. The cost function has many local minima, and thus the projection should be initialized properly, in an orderly fashion. A suitable sample initialization is shown in Figure 4a.

The parameter  $\lambda_2$  should be set to a large value, say,  $\lambda_2 = 100$  like in Section 3. We recommend setting the other parameter,  $\lambda_3$ , initially to zero to make the mapping flexible, and to increase its value gradually to, say,  $\lambda_3 = 10$  during the optimization to enforce the constraint.

The scaling factor  $\kappa$  in (7) is the only remaining free parameter in the cost function. Since the region of the CIELab space that corresponds to the available colors is very limited and irregularly shaped, a compromise must be made between a faithful representation of the distances and the color resolution. The scale should be chosen according to their relative importance.

If  $1/\kappa$  is small, like in Figure 4b, the distances are represented faithfully but only a small part of the available color space is utilized. As  $1/\kappa$  increases (Figs. 4c and d) the projection utilizes a larger portion of the available space but the distances between model vectors become somewhat more distorted because the projection must conform to the borders of the area. Fortunately, however, small changes in the scale usually only have small effects on the coloring. For example, there are only small distortions in the colorings corresponding to either of the projections in Figs. 4c and d although the projections seem very different (see <http://www.cis.hut.fi/sami/ajlips00/>).

The result of mapping the model vectors of the SOMs of Figures 2a and 3a into the CIELab color space is shown in Figures 2c and 3c. Compared with the non-constrained projections in Figures 2b and 3b the overall organization is similar but some compromises have clearly been needed especially in conforming to the disk of non-saturated colors in the middle. Nevertheless, it is clear, based on the results, that the projection method is able to make the necessary compromises.



**Figure 4:** Effect of scale on the projection of an artificially generated, evenly spaced two-dimensional SOM into a constant-lightness cross section of the CIELab color space. **a:** Initial state. **b:** Optimal preservation of distances but low resolution: small scale  $1/\kappa$ . **c** and **d:** Progressively larger scale. The thick lines delimit the region representable by a typical CRT tube, and the circle in the middle encompasses colors that were considered too non-saturated (gray).

When the graphical SOM display was colored according to the projection (see <http://www.cis.hut.fi/sami/ajjips00/>), the perceptual differences in the hue of the neighboring map units corresponded well with the distances in the original data space, depicted as shades of gray in Figures 2a and 3a. Moreover, the hues became ordered globally; inside the different clustered areas the hues were relatively uniform, and different clustered areas attained different colors.

## 5 Coloring of the Original Data According to Its Cluster Structure

The coloring can be even more useful if the original data set can be visualized in some other manner besides the SOM display. Then each data item can be colored according to its color on the SOM display. The welfare and poverty structures, for instance, can be visualized in a straightforward manner on a geographic map display. The countries can then be colored according to their welfare or poverty type (see <http://www.cis.hut.fi/sami/ajjips00/>). The result is a display where countries having a similar welfare or poverty type have been colored similarly irrespective of their geographical location. Japan and Australia, for example, are fairly similar to the European countries and to the USA and Canada. Countries that belong to very different types than their neighbors pop out strongly, like Japan, Sri Lanka, and South Africa.

Such visualizations can be very useful if the data has a “natural” ordering like the geographical order here but, in fact, *any* order of the data items can be used. For example, if the countries were ordered simply according to the GNP per capita in a table and colored using a SOM, then countries in which the welfare or poverty type is different from the other countries having a similar value of GNP per capita would be clearly discernible based on sharp discontinuities in the coloring.

*Note:* In principle our method could be used to map the data set directly, instead of mapping the model vectors of a SOM computed from the data set. It would, however, be more difficult to define which distances are local enough so that they should be represented accurately. The SOM makes this decision automatically. If it is necessary to obtain a characteristic color for each data item then some suitable local fitting method may be used to complement the SOM-based mapping.

## 6 Conclusions

We have introduced an automatic method for coloring graphical SOM displays so that the perceptual properties of the coloring reflect the properties of the high-dimensional statistical data as closely as possible. In fact, the resulting coloring is almost tailored to the human color vision system which is very accurate in detecting differences between

the colors of neighboring areas, in this case the color differences between neighboring locations on cluster displays. The easily interpretable coloring makes it possible to visualize complex statistical structures automatically even for non-experts.

## References

- [1] “Colorimetry, 2nd Ed.” CIE Publication No. 15.2, 1986.
- [2] Kaski, S. and Kohonen, T., “Exploratory data analysis by the self-organizing map: Structures of welfare and poverty in the world,” in Refenes, A.-P. N., Abu-Mostafa, Y., Moody, J., and Weigend, A. (editors), *Neural Networks in Financial Engineering*. World Scientific, Singapore, pp. 498–507, 1996.
- [3] Kohonen, T., “Self-Organized Formation of Topologically Correct Feature Maps,” *Biological Cybernetics*, 43, pp. 59–69, 1982.
- [4] Kohonen, T., *Self-Organizing Maps*, Springer, Berlin, 1995, (Second, extended edition 1997).
- [5] Kohonen, T., “The self-organizing Map,” *Neurocomputing*, 21, pp. 1–6, 1998.
- [6] Kohonen, T., Hynninen, J., Kangas, J., Laaksonen, J., and Torkkola, K., “LVQ\_PAK: The Learning Vector Quantization Program Package,” Technical Report A30, Helsinki University of Technology, Laboratory of Computer and Information Science, 1996.
- [7] Sammon, Jr., J. W., “A nonlinear mapping for data structure analysis,” *IEEE Transactions on Computers*, C-18, pp. 401–409, 1969.
- [8] Ultsch, A., “Self-Organizing Neural Networks for Visualization and Classification,” in Opitz, O., Lausen, B., and Klar, R. (editors), *Information and Classification*, Springer-Verlag, Berlin, pp. 307–313, 1993.
- [9] World Bank, *World Development Report 1992*, Oxford Univ. Press, New York, NY, 1992.