

Combining Cluster Analysis and Continuous Isosurface Colouring in a Tool for Data Exploration

Johan Hagman
SSKKII, University of Göteborg
41298 Göteborg, Sweden
hagman@ling.gu.se

Abstract

The paper shows how two powerful techniques for supporting data exploration of multi-dimensional data can be combined in a tool for this purpose. The techniques, cluster analysis and graphic visualization, are briefly presented and discussed as modules of a prototype. Its performance is illustrated by experiments resulting in cluster configurations where the value distribution of “underlying” dimensions are visualized with colors or shades of grey.

1 Introduction

A system which integrates *cluster analysis* and *visualization* has been developed to meet an often felt need for this special kind of tool for data exploration in scientific, industrial and everyday applications. The prototype has preliminary been called the *Cluster Analyzer & Visualizer* — or CAV for short and this paper discusses some experiences from its development and present performance.

The epoch we are living in at present is getting less “industrial” and more “informational”. In some areas this process has been going on to such an extent that we are presented with difficulties in surveying all the information offered and we find it hard to navigate through it to get what we want or to discover unexpected things. The tool presented here is an attempt to combine two techniques of making especially multidimensional data more graspable — cluster analysis and data visualization. The development of these techniques has essentially been parallel with that of the computers and their usability depends heavily on the latter’s performance.

2 The Structure of CAV

The development of CAV has included several aspects:

- requirements presented by a number of “real-world” applications
- choice of clustering methodology, algorithms, and implementation environment
- experiments with various types of visualization
- steps toward an integration with human-computer interactive facilities, aiming at a powerful tool for information exploration and visualization.

In [Hagman 94a] these aspects are presented more in detail along with the terminology that will be used here.

2.1 Data import format

CAV accepts *pattern matrices* and *proximity matrices* in plain ASCII format either as tables with rows and columns or as lists (in case the table would be too big to handle with available editors). Empty cells are allowed for and the proximity matrices are supposed to be marked as to whether the values denote *similarity* or *dissimilarity*.

2.2 Data processing

The data are equalized *before* they are normalized. Since equalization works with *relative* values such as the proportion of the standard (or maximal) deviation to the mean value for each feature¹, the values do not have to be within the same range when comparing their *relative* spread. CAV checks both the mean deviation and the maximal deviation in order to detect single, extreme cases. Then, when the internal variation has been “compressed” or “blowed up” to be within a user-specified relative range, the values are normalized.

¹ Or ‘attribute’, ‘parameter’, or ‘column’.

Missing data do not affect the equalization. Mean and maximum values can be calculated anyway out of the data actually supplied. Whether this is correct is discussed in the literature. One suggestion is that missing data should be filled in by the mean value. Another, more elaborate proposal is that a pattern which is lacking a value for a certain feature could have this value calculated by comparing (seemingly ignoring other possibly absent data) with the n patterns most similar to that pattern.

2.2.1 Pattern matrix

The operations on the original matrix result in a normalized and equalized pattern matrix.

2.2.2 Proximity matrix

The pattern matrix underlies the calculation of the proximity matrix. For each pair of patterns the sum or product² of the differences of their corresponding feature values is calculated. Sums are then divided by the number of comparable features (i.e. where neither pattern had an empty cell), giving the proximity index for that pattern pair. Eventually this process yields a proximity matrix that could be used directly as input for future runs. It is also possible to analyze the pattern matrix vertically instead of horizontally. Then the features are compared as to their covariance through the patterns. Whatever is analyzed, patterns or features, the term 'item' will be used henceforth.

2.2.3 Dendrogram

The proximity matrix is passed on to a procedure which builds the corresponding binary dendrogram. At present, an agglomerative, hierarchical clustering is carried out which iteratively merges the rows and columns of the most similar pair of items while simultaneously joining the corresponding subtrees in a new node. The output of all three steps described in 2.2.1 through 2.2.3 may be downloaded to separate files for special examination.

2.2.4 Iterative procedure

The next phase is an "animated" clustering, currently on a 2D scatter plot. The algorithm is a variant of the iterative one presented in [Sammon 69]. In CAV the items are represented by small, mobile labels on a 2D surface where they move around, driven by their interrelational forces of attraction and repulsion, like the balls in billiards. They continue to move until a satisfactory low-stress configuration has been reached. The criterion for satisfaction is given by the user; a typical one is to let the configuration adjust its shape until no new minimal stress value has been encountered during the latest n movements, having n set to some suitable value. Different initial configurations have been tried for the same proximity matrix and with very rare exceptions they yield the same final result. The crucial thing is to avoid cases where items are too astray to make it on time to reach their "home cluster" before other clusters are formed in its way blocking it. To minimize this risk the dendrogram is used; the order of its leaves from left to right is a kind of 1D representation of their internal order and that is followed when the items are lined up diagonally over the screen³. This speeds up the 2D clustering process considerably and straightens the item trajectories.

A comparison of [Sammon 69]'s Euclidean-based projection [Jain & Dubes 88, p. 39] of the so-called '80X data' has also been made with a CAV Minkowski-based projection (FIGURES 1a and 1b). The [Sammon 69] method gives good results and today's computers, as compared with those of the '60s, offer possibilities of another magnitude concerning execution time and memory. Other methods that also involve an iterative reshaping of the configuration have been proposed; one of them is *simulated annealing* [Kirkpatrick 83], [Bell 90], [Chalmers & Chitson 92] and below *neural networks* represent yet another.

When the best configuration has been found, a validation of the clustering can be done. In a series of tests the program was fed with proximity matrices describing 2D distances between corners of various known geometrical shapes which later during the clustering process have gradually re-emerged on the display. Another example of this is given in section 3.1.

² Depending on whether *Manhattan distance* or *Euclidean distance* is preferred.

³ An alternative is to start by forming a circle.

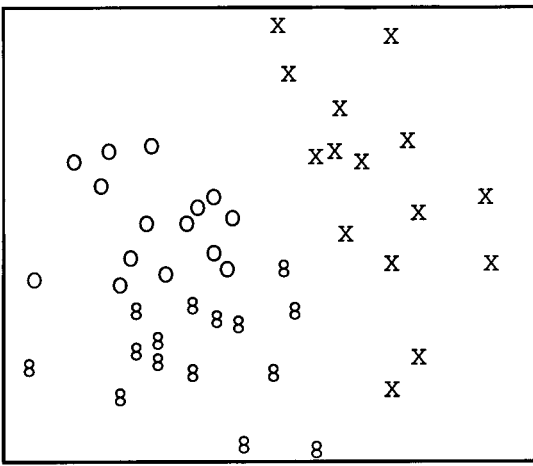


FIGURE 1a Euclidean projection of the '8OX' data by [Sammon 69] from [Jain& Dubes 88].

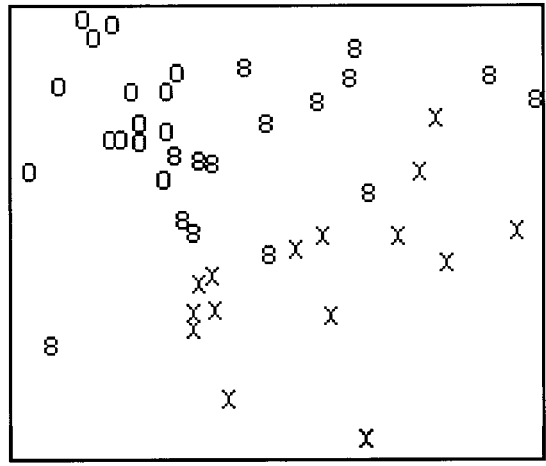


FIGURE 1b Manhattan projection of the '8OX' data by CAV.

2.3 Visualization

We proceed to the next module, which is a more advanced visualization where shades of grey or colour are added to the 2D cluster configuration. Both the grey shades and the colours constitute interval scales with a finite number of degrees which, however, depending on the hardware, may be millions. Grey shades increase with the values from black to white. The colour scale also runs from black to white but passes through blue, purple, red, orange, and yellow. Note that colour intervals may be differently defined; frequently (mostly) the interval is an inverted variant of the visible spectrum with the increasing sequence purple, blue, green, yellow, orange, and red. This could be claimed to be a "natural" interval but even the choice of the aforementioned interval could be supported by nature, i.e. e.g. an analogy of metal at different temperatures (assuming a bluish metal) which explains the absence of green. Anyway, [Brodie & al. 92] state that:

As yet there are no widely accepted standard colour scales or range scales (such as logarithmic) for data exploration in visualization. A need is perceived for empirical studies leading to the derivation of formal guidelines and reliable standards in this area. [p. 119]

For a further discussion of colours, see also [Tufté 83].

Shades or colours (henceforth: 'intensity'⁴) could be used to enhance the groupings of the clusters so that for each picture element (i.e. $n*n$ pixel square, $n \geq 1$) of the picture the density of items in that part of the configuration is represented by a value on the intensity interval. To do this, the minimum and maximum densities have to be calculated in order to calibrate the intensity scale and thereby use it optimally. When the impression of groups in this way gets enhanced, each item contributes to its "background" shading/colouring by its mere presence. All items contribute equally to the intensity distribution throughout the picture. This is illustrated by the FIGURE 2.4a.

Now, it is also possible to associate different values to the items — values, for instance, corresponding to one of the underlying features (if the items correspond to the *rows* of an underlying pattern matrix). In such a visualization, the intensity would reflect the distribution of that feature throughout the picture. Examples of this are the FIGURES 3x and 3y. Items with missing data for the chosen feature can have their background intensity decided by interpolating the intensities of the surrounding items. This could also be a way of *estimating* a missing value, probably better than the assignment of the mean value, mentioned above, and fully comparable with the method using the nearest neighbours.

3 Experiments

One of the many pattern matrices processed by CAV was a table⁵ of the consumption per capita of three types of alcoholic beverage (spirits, wine, and beer) in 32 countries. The resulting 3D \Rightarrow 2D reduction and cluster configurations are showed in FIGURES 3.1a through 3.3.

⁴ Note that this refers to the values of the underlying data features and not to the brightness on the screen.

⁵ Extracted from [SCB 91]

Another pattern matrix, provided by [Hofstede 91], characterizes some fifty countries according to the four dimensions 'individuality', 'masculinity', 'power distance', and 'uncertainty avoidance'. The cluster configuration of this 4D⇒2D reduction is shown as part of the FIGURES 4.1 and 4.2.

CAV has also been fed directly with proximity matrices. "Semantic spaces" for part of the content of a Swedish thesaurus have been generated where the words were clustered according to their co-appearance in definitions. The arrival of computers has made possible many interesting experiments with automatic (re)construction of thesauri [Morris & Hirst 91], [Crouch & Yang 92]. In addition, we are currently experimenting with clustering newspaper articles of different genres based on the presence of certain word stems.

A special type of proximity matrix is the *asymmetric* one, i.e. where the relation $A \rightarrow B$ is *not* the same as the relation $B \rightarrow A$. An example of this is the speakers' asymmetric sequencing in a discourse and such a matrix was analyzed and visualized in e.g. [Allwood & Hagman 94]. Moreover, similar experiments have been done with the same data as are partly presented in [Eikmeyer 92] (i.e. confusion and association matrixes) — in collaboration with this author. One of those experiments is the following.

3.1 Regeneration of a Map from a Road Distance Table

In this experiment a table of road distances between 11 cities (in FIGURE 2.1) was fed into CAV as a *dissimilarity* matrix.

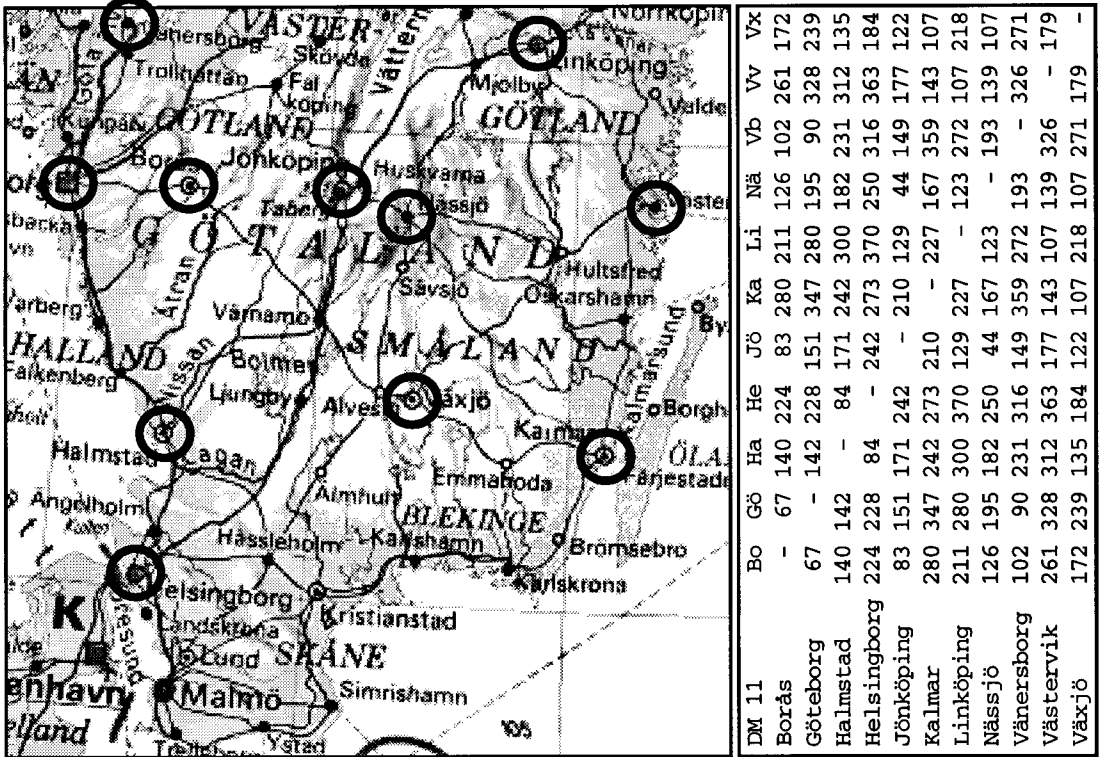


FIGURE 2.1 (above left)
Map of southern Sweden with 11 cities encircled.

TABLE 2.2 (above right)
Road distance table (in kilometers), used as a dissimilarity matrix.

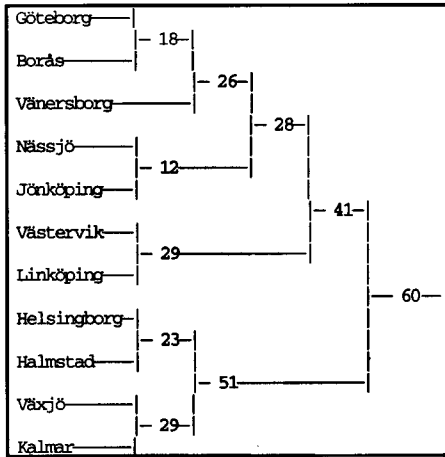


FIGURE 2.3a (left)
Dendrogram reflecting city groupings. Each number indicates the distance between the two subtrees expressed as percentage of the maximum intercity distance in Figure 2.2.

TABLE 2.2 shows this matrix and FIGURE 2.3a shows the resulting dendrogram after a single-linkage hierarchical clustering based on Manhattan distances. After letting the forces of attraction and repulsion move the city labels around on the 2D “animated” clustering display for a couple of minutes, the low-stress configuration shown in the final picture of FIGURE 2.4a is reached. Note however that, of course, this configuration has no *dimensions* or explicit *axes*; its form is completely determined by the intercity *relations*. The fact that even the orientation of this picture so closely coincides with FIGURE 2.1 is accidental; it could as well have been rotated or mirror-inverted.

3.2 Cluster Analysis and Neural Networks — a Comparison

The dissimilarity matrix in FIGURE 2.1 was also used as input to a self-organizing Kohonen map⁶. Output of this approach is shown in TABLES 2.3b and 2.4b. TABLE 2.3b shows the similarities between a city’s distance pattern to all other cities, where the mean similarity distance is assumed to be 100 with a standard deviation of 27. A low similarity distance indicates that the cities have a highly similar distance pattern. The cities have been put together into five groups for which the group internal similarity is significantly high (i.e. the similarity distance is below $100-27=73$). Compare this table with the dendrogram in FIGURE 2.3a.

	He	Ha	Bo	Gö	Vb	Jö	Nä	Vx	Ka	Vv	Li	
He	47		104	109	117	87	90	119	126	123	121	
Ha			103	108	117	90	92	123	130	126	124	
Bo			50	62		79	87	126	131	128	127	
Gö						57	85	92	132	136	134	132
Vb							77	83	120	122	120	119
Jö						50		82	86	83	82	
Nä								77	79	79	77	
Vx								52		106	109	
Ka										98	103	
Vv											50	
Li												

FIGURE 2.3b City groupings (framed cells) as suggested by the neural network in this modified version of TABLE 2.2. The mean distance here is 100 (Courtesy of H.-J. Eikmeyer).

The data were represented in the neural network by a 2D 10x10-cell Kohonen map and one of the best output results is shown in FIGURE 2.4b, where the city abbreviations in the cells reflect the cities’ “area of dominance” and distance relations⁷. As well as in FIGURE 2.4a, even in this configuration there are no “natural axes” as could be claimed for FIGURE 2.1 but, nevertheless, the similarity among all these three figures becomes clearer when we turn FIGURE 2.4b 90° clockwise and then rotate it along its vertical Vb—Ka/He axis.

3.3 Colouring and shading

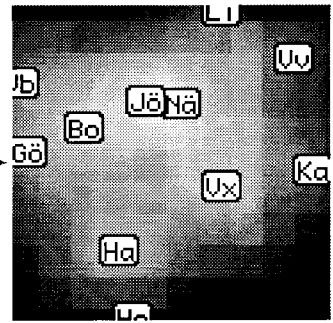
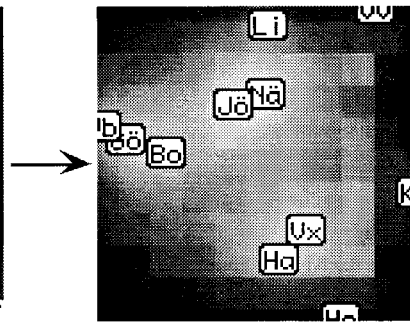
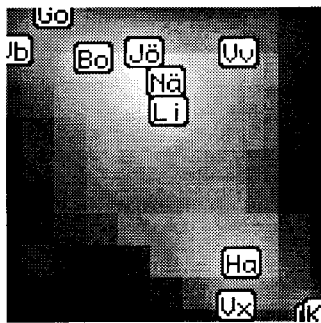
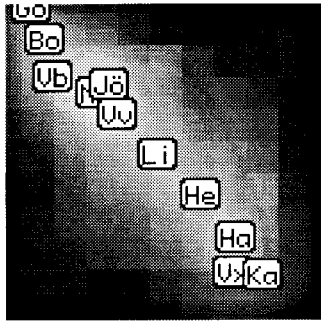
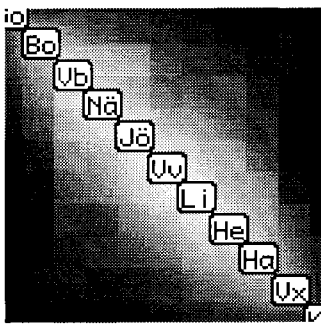
The intensity shading or colouring of the 2D configurations has invited to many experiments. The principle of assigning intensity to each point of the picture was described above. In one type of visualizations⁸ the *isosurfaces* (areas of the same intensity) are indicated only by their boundaries, the *isolines*. Often height is marked in this way on topographic maps. When a geometric third dimension as height is visualized it often has to be shown in perspective to render the effect of 3D. Such “3D” visualizations could also be combined with intensity hues that either show the same feature as does the height dimension (i.e. the same colours on the same levels) or that show yet another feature with intensity fields running over hills and valleys independently of these.

All the FIGURES 3.1a through 3.3 are based on the same set of data, mentioned above, and their cluster configuration is identical. In FIGURES 3.1a and 3.1c the value distribution of the

⁶ By H.-J. Eikmeyer and his colleagues at the Dept of Linguistics, University of Bielefeld. For related reading, see e.g. [Ritter & Kohonen 89] and [Eikmeyer 92].

⁷ For an interesting analogy: cf. [Lin 92].

⁸ See [Brodie & al. 92] for a classification of visualization types.



Bo	Bo	Bo	BoGo	Go	GoHa	Ha	Ha	Ha	Ha
Bo	Bo	BoGo	Go	Go	Go	GoHa	Ha	Ha	HaHe
Vb	BoVb	BoGo	Go	Go	Go	GoHa	HaHe	He	He
Vb	Vb	NÄVb	GoNÄ	GoNÄ	GoJo	Jo	HeJo	He	He
Vb	Vb	NÄVb	NÄ	JoNÄ	Jo	Jo	Jo	He	He
Vb	NÄVb	NÄ	NÄ	NÄ	Jo	Jo	Jo	JoKa	Ka
LiVb	LiNÄ Vb	NÄ	NÄ	JoNÄ	Jo	Jo	JoKa	Ka	Ka
Li	Li	LiNÄ	NÄVv	NÄVv	JoVv	Jo	JoKa Vx	KaVx	Ka
Li	Li	LiVv	Vv	Vv	Vv	VvVx	Vx	Vx	KaVx
Li	Li	LiVv	Vv	Vv	Vv	VvVx	Vx	Vx	Vx

FIGURE 2.4b Kohonen map based on TABLE 2.2 (H. -J. Eikmeyer)

FIGURE 2.4a Stepwise reshaping in five consecutive cluster configurations generated by CAV. Starting from a diagonal line-up of the leaves of the dendrogram (FIGURE 2.3a), the 'stress' decreases as the city labels (centered at the actual city coordinates) move according to the internal forces of attraction and repulsion among the cities. The 'stress' is the discrepancy between the distances in the dissimilarity matrix (TABLE 2.2) and the current configuration.

features 'wine consumption' and 'spirits consumption', respectively, is visualized⁹. Note, by the way, that the items in this case already have a natural configuration suggested by their geographical extension and that this as well could have been chosen as the configuration. Then, however, we would suffer the disadvantage of having large areas to which no values are associated (e.g. the areas between the European countries and JAPAN, South AFRICA, or PERU). To illustrate a situation of data exploration and the possibility to colour the configuration according to other features which have not contributed to the clustering of the data, FIGURES 3x and 3y are added. There, in the same configuration the reported rates of two different causes of death are visualized.

The cluster configurations of the [Hofstede 91] data in FIGURES 4.1, and 4.2 have been briefly commented upon. In Figure 4.1 the value distribution of the feature dimension 'individuality' is visualized with intensity less smoothed than it is in FIGURE 4.2 which shows the value distribution of the feature 'masculinity'.

⁹ See [Hagman 94b] for colour versions of these figures.

Since all colours can be generated by combining the three elementary colours red, green, and blue, a special possibility is offered when one wants to visualize the value distribution of exactly three features on a 2D or 3D scatter plot or cluster configuration; each feature is simply associated with one of the three colours. This is described e.g. in [Ammermann & Cavalli-Sforza 84] where each colour was associated with one of the three *principal components* of a genetic analysis of the populations of Europe to try a hypothesis of prehistoric migration waves from Asia. The 2D representation in that study was a map over Europe and part of Asia and it was coloured in this way ...

... to take advantage of the power of the human eye to sum elementary colors and thus synthesize the information contained in the first three components, which together account for almost 60% of the information present in the original set of genetic data. It was heartening to find that the patterns obtained were in basic agreement with those expected under the demic hypothesis." [p. 105]

FIGURES 3.2a and 3.2c are monochromatic variants of FIGURES 3.1a and 3.1c, respectively, but with another "smoothness" of the colour shading. Each of these figures has been coloured with shades of just one colour: red and blue, respectively. FIGURE 3.2b visualizes the distribution of the beer consumption per capita in the configuration with shades of green. The synthesis of these three pictures is shown in FIGURE 3.3 and there, the transitions across the four consumption corners 'wine', 'beer', 'liquor', and '(officially) neither of the three' are clearly visualized (not without some aesthetic appeal) with millions of colour nuances of the spectrum which show the continuous transitions from one consumer category to the other.

4 Conclusion and Further Work

At present, the modules may be likened to the "turnkeys" [Brodie & al. 92] which can be run either independently or as integrated parts of a so-called application building system. Steps are taken to integrate all or part of the modules in a more powerful system for information exploration and visualization. The intent is to fuse CAV's possibilities with those offered e.g. by *dynamic queries* [Ahlberg & al 92], as well as other query devices [Ahlberg & Truvé 94], and *tight-coupled starfield displays* [Ahlberg & Shneiderman 94]. In this integrated system, the CAV parts will contribute with the possibilities of displaying data in cluster configurations and colouring this configuration universe or "starfield" in the ways presented here. FIGURES 4.1 and 4.2 are early sketches that illustrate what a simple version of this emerging tool could look like. These figures represent computer screens where the displays' current output is a cluster configuration of a data set from [Hofstede 91]. The only widgets shown are the four sliders, automatically created when reading the input file and assigned to one dimension, plus a regulator for the intensity representation of selected features.

Both the kind and the degree of the *interaction* of a system for data exploration is highly depending on the application and the expected type of users. As for university or industrial researchers, they are to a greater extent expected to know their data sufficiently to try different techniques for data processing and visualization. The extreme opposite type of application is the public service tools, e.g. automatic library guides or similar [Salton & al. 94]. In such a system a subset of all possible facilities may be chosen by the system developer, leaving to the user less (but sufficient) freedom — and maybe less problems. Since the requirements of interaction are not yet presented for the integrated tool for this information exploration and visualization, its interactive part is still to be completed. Thus, its final design will depend on application, user types, and, of course, hardware environment.

Acknowledgements

This work was enabled by support from NUTEK, grant no. 5321-93-2760.

I would like to thank Hans-Jürgen Eikmeyer, Dept of Linguistics, University of Bielefeld for our fruitful collaboration and my colleagues Christopher Ahlberg, Jens Allwood, and Staffan Truvé at SSKKII Research Centre (an intersection of the University of Göteborg & Chalmers University of Technology) for many valuable comments and suggestions.

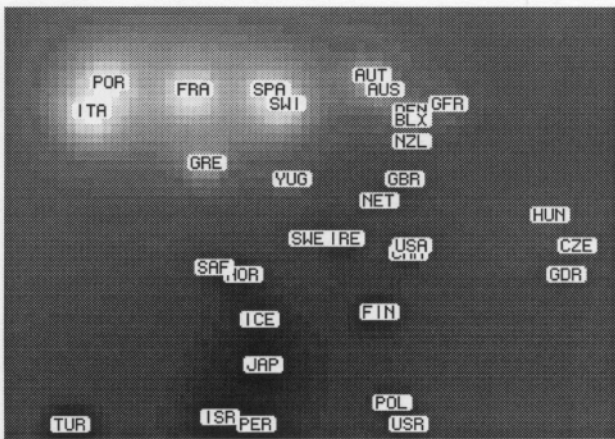


FIGURE 3.1a Cluster config. of the (reported) consumption of beer, wine, and spirits 1983-87, highlighting wine consumption.

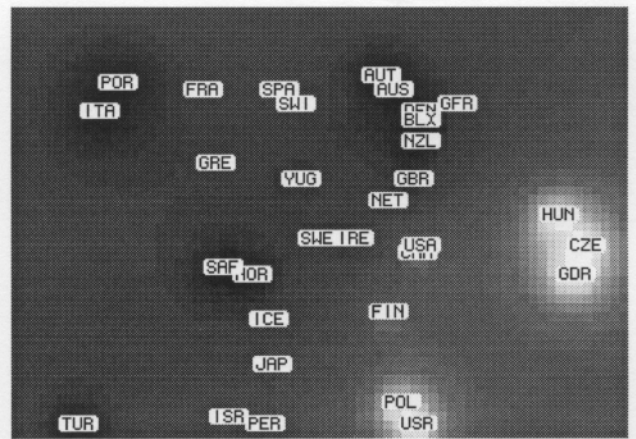


FIGURE 3.1c Same cluster as in FIGURE 3.1a. Highlights spirits consumption. (Values missing for GRE, ISR and PER)

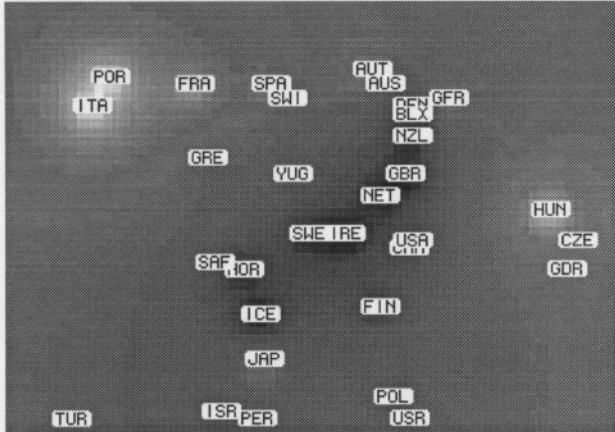


FIGURE 3x Clustered as FIGURES 3.1a & 3.1c. Deaths by cirrhosis 1982-86. (No values reported for PER, SAF, TUR, USA, and YUG)

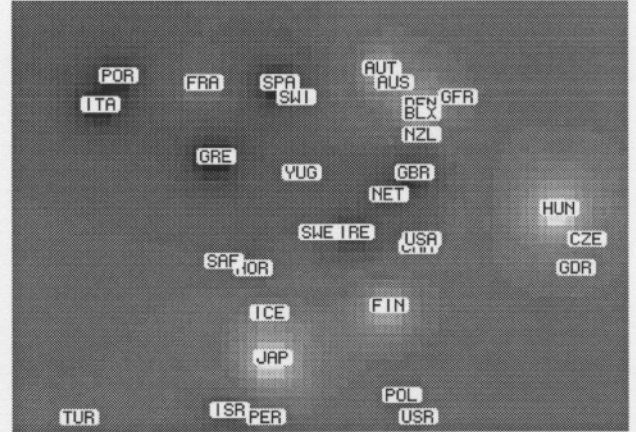


FIGURE 3y Clustered as FIG:s 3.1a & 3.1c. Shows suicide rates 1982-86. (No values for GDR, PER, SAF, TUR, USA, and YUG)

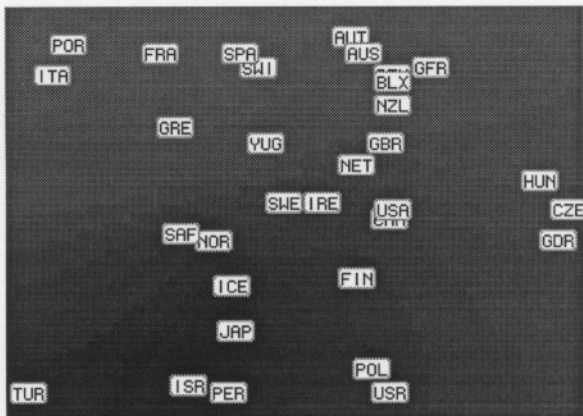


FIGURE 3.2a Monochr. (red) smoothed version of FIGURE 3.1.a

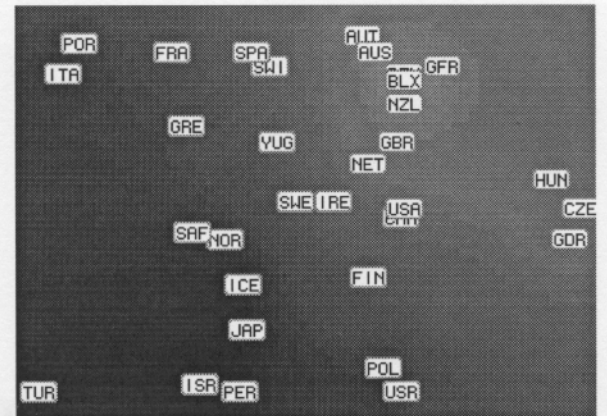


FIGURE 3.2b Monochr. (green) distrib. of beer consumption

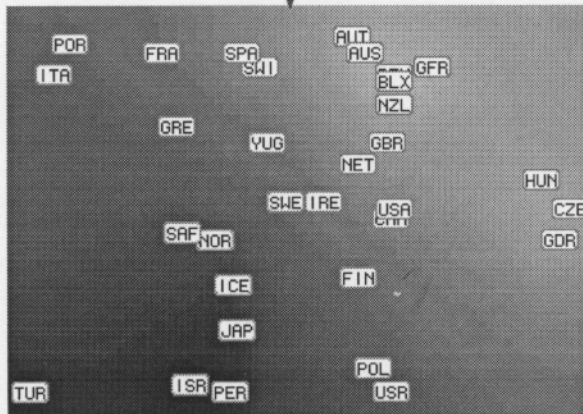


FIGURE 3.3 Polychrome synthesis of FIG:s 3.2a, 3.2b, and 3.2c

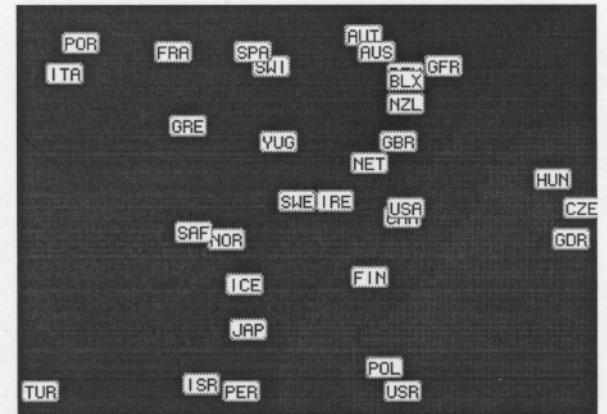


FIGURE 3.2c Monochr. (blue) smoothed version of FIGURE 3.1.c

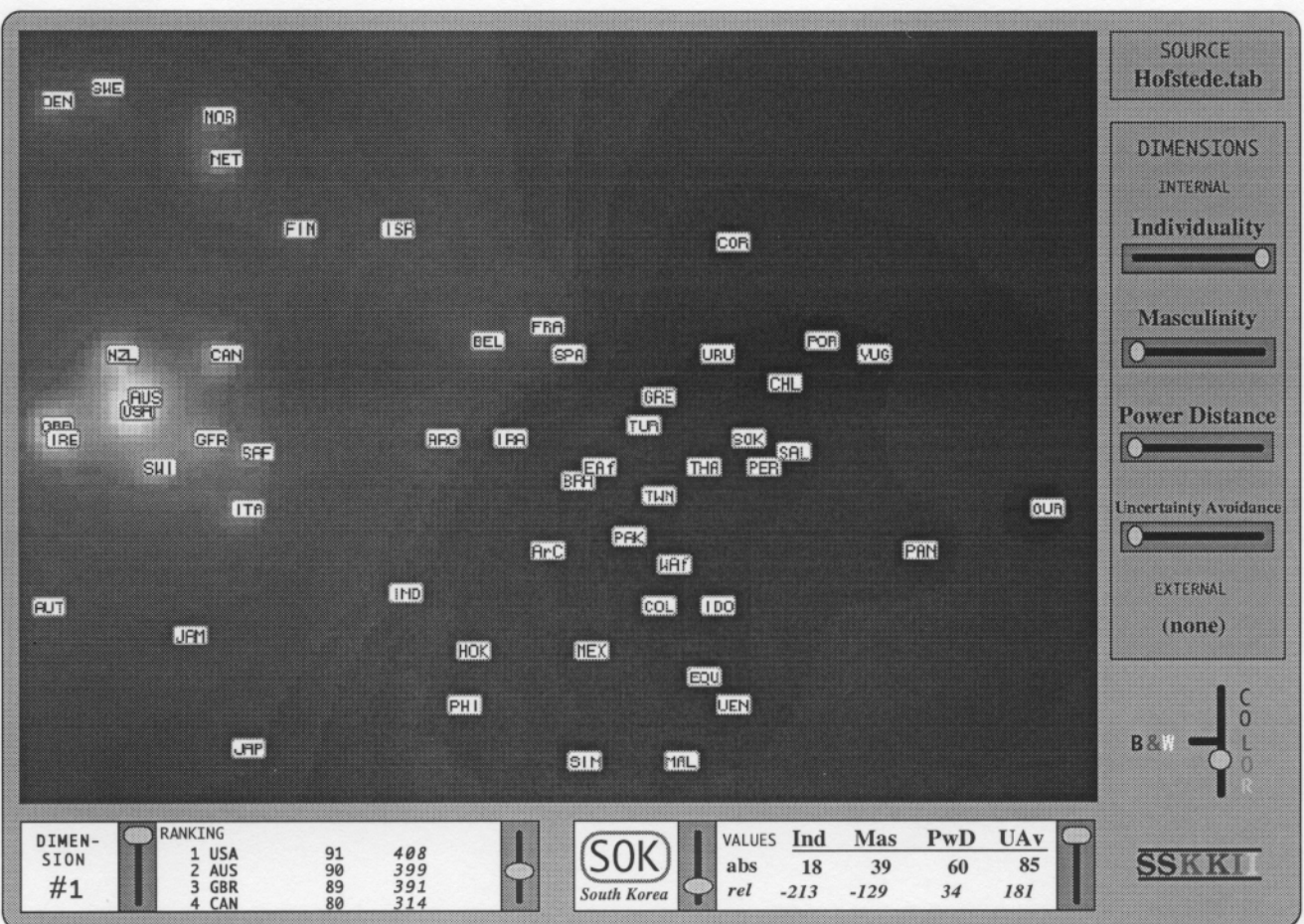


FIGURE 4.1 Example of interface of a tool for information visualization and exploration. The cluster configuration is the result of a CAV analysis of intercultural data from [Hofstede 91]. The current setting shows the value distribution of the feature 'individuality' with low colour-smoothing.

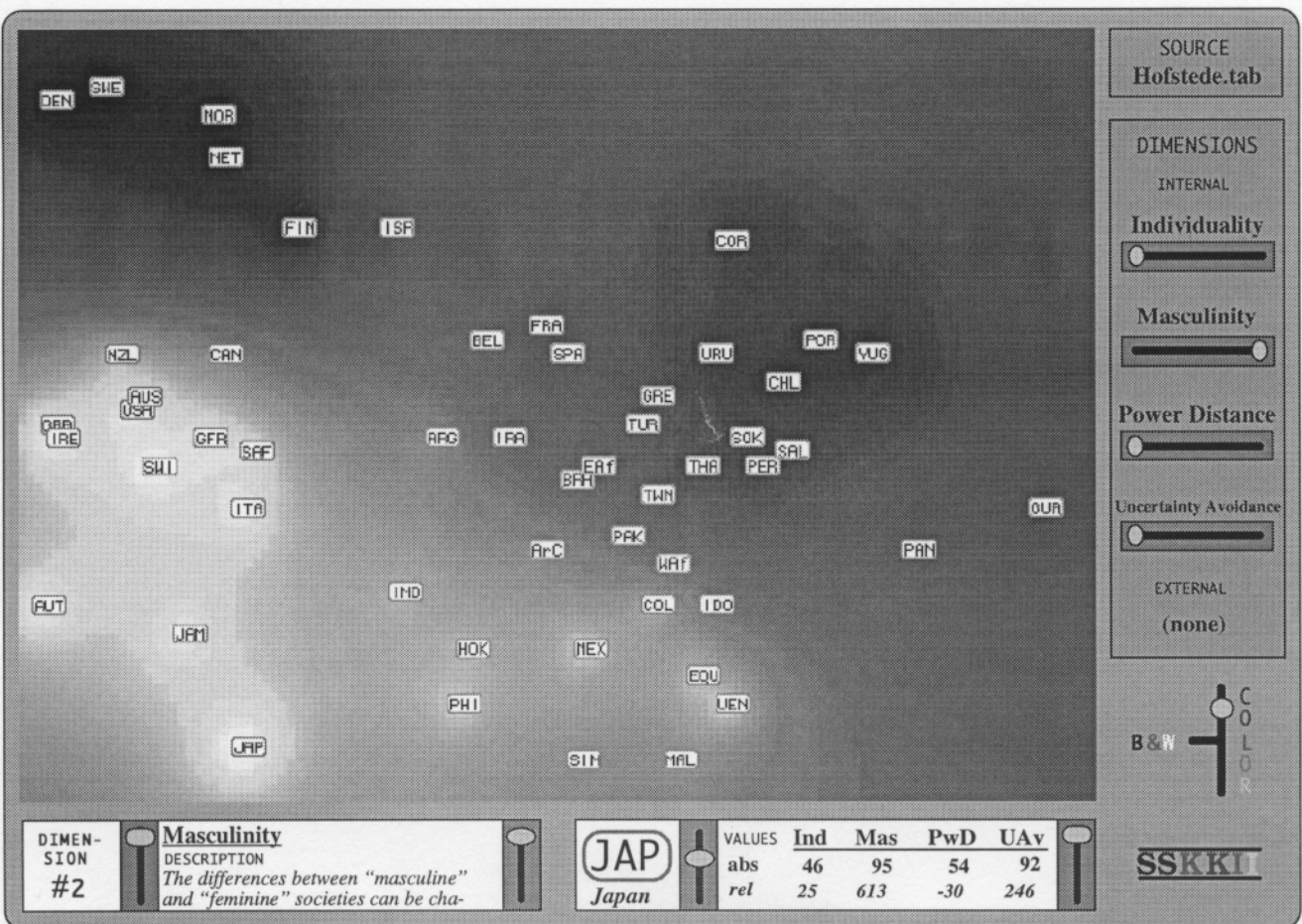


FIGURE 4.2 The same interface as in FIGURE 4.1 when indicating the value distribution of the feature 'masculinity' with high colour-smoothing.

References

- Ahlberg, C., C. Williamson, and B. Shneiderman (1992)
"Dynamic Queries for Information Exploration: An Implementation and Evaluation",
Proceedings ACM CHI'92: Human Factors in Computational Systems
- Ahlberg, C. and B. Shneiderman (1994)
"Visual Information Seeking: Tight Coupling of Dynamic Query Filters with Starfield Displays", In *Proceedings ACM CHI'94: Human Factors in Computational Systems*.
- Ahlberg, C. and S. Truvé, 1994. "Exploring Terra Incognita in the Design Space of Query Devices", submitted to UIST'94
- Allwood, J. and Hagman, J. (1994) "Some Simple Automatic Measures of Spoken Interaction", in *Proceedings from the joint 14th Scandinavian Conference of Linguistics & 8th Conference of Nordic and General Linguistics*, University of Göteborg
- Ammerman, A. J. and L. L. Cavalli-Sforza (1984)
The Neolithic Transition and the Genetics of Populations in Europe Princeton Univ. Pr.
- Bell, D. A., F. J. McEarlean, P. M. Stewart, and S. J. McClean (1990)
"Application of Simulated Annealing to Clustering Tuples in Databases", in *Journal of the American Society for Information Science*, 41/2 (pp. 98-110)
- Brodie, K.W, L. A Carpenter, R. A. Earnshaw, J. R. Gallop, R. J. Hubbard, A. M. Mumford, C. D. Osland, P. Quarendon (1992)
Scientific Visualization - Techniques and Applications, Springer-Verlag
- Chalmers, M. and P. Chitson (1992) "Bead: Explorations in Information Visualization" in *SIGIR '92, Proceedings from the 15th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval* ed. by N. Belkin, P. Ingwersen, and A. Mark Pejtersen IEEE Computer Society Press (pp. 330-337)
- Crouch, C. J., and B. Yang (1992)
"Experiments in Automatic Statistical Thesaurus Construction", in *SIGIR '92, Proceedings from the 15th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval* ed. by N. Belkin, P. Ingwersen, and A. Mark Pejtersen IEEE Computer Society Press
- Eikmeyer, H.-J. (1992)
"Structuring Phonological Space with Self-Organizing Maps", forthcoming in *Papers from the 13th Scandinavian Conference of Linguistics*, Roskilde (ed.) L. Heltoft
- Hagman, J. (1994a)
Information Visualization and Exploration — Potential and Requirements, SSKKII Technical Report #3, University of Göteborg / Chalmers University of Technology
- Hagman, J. (1994b)
The Cluster Analyzer & Visualizer — a Prototype for Data Exploration, SSKKII Technical Report #5, University of Göteborg / Chalmers University of Technology
- Hofstede, G. (1991) *Cultures and Organizations*, McGraw-Hill Book Company
- Jain, A. K. and Dubes, R. C. (1988)
Algorithms for Clustering Data, Prentice Hall Advanced Reference Series
- Kirkpatrick, S., C. D. Gelatt Jr, and M. P. Vecchi (1983)
"Optimization by Simulated Annealing", *Science*, vol. 220 (pp. 671-680)
- Lin, X. (1992)
"Visualization for the Document Space", in *Visualization '92, Conference Proceedings*, IEEE Computer Society Press, (pp. 274-281)
- Morris, J. and G. Hirst (1991)
"Lexical Cohesion Computed by Thesaural Relations as an Indicator of the Structure of Text" in *Computational Linguistics*, March '91, (pp. 21-48)
- Ritter, H. and T. Kohonen (1989)
"Self-Organizing Semantic Maps", in *Biological Cybernetics*, 61 (pp.241-254)
- Salton, G., J. Allan,, and C. Buckley
"Automatic Structuring and Retrieval of Large Text Files", in *Communications of the ACM*, February 1994 / Vol. 37, No.2 (pp. 97-108)
- Sammon, J. W. (1969)
"A nonlinear mapping for data structure analysis", *IEEE Transactions on Computers* C 18, (pp.401-409)
- SCB (1991) *Statistisk årsbok för Sverige*, SCBs Förlag, Stockholm
- Tufte, E. R. (1983)
The Visual Display of Quantitative Information, Graphics Press, Cheshire, Connect.