

Value and Relation Display for Interactive Exploration of High Dimensional Datasets

Jing Yang, Anilkumar Patro, Shiping Huang, Nishant Mehta, Matthew O. Ward and Elke A. Rundensteiner

Computer Science Department

Worcester Polytechnic Institute

Worcester, MA 01609

{yangjing,anil,shiping,nishantm,matt,rundenst}@cs.wpi.edu *

ABSTRACT

Traditional multi-dimensional visualization techniques, such as glyphs, parallel coordinates and scatterplot matrices, suffer from clutter at the display level and difficult user navigation among dimensions when visualizing high dimensional datasets. In this paper, we propose a new multi-dimensional visualization technique named a *Value and Relation* (VaR) display, together with a rich set of navigation and selection tools, for interactive exploration of datasets with up to hundreds of dimensions. By explicitly conveying the relationships among the dimensions of a high dimensional dataset, the VaR display helps users grasp the associations among dimensions. By using pixel-oriented techniques to present values of the data items in a condensed manner, the VaR display reveals data patterns in the dataset using as little screen space as possible. The navigation and selection tools enable users to interactively reduce clutter, navigate within the dimension space, and examine data value details within context effectively and efficiently. The VaR display scales well to datasets with large numbers of data items by employing sampling and texture mapping. A case study on a real dataset, as well as the VaR displays of multiple real datasets throughout the paper, reveals how our proposed approach helps users interactively explore high dimensional datasets with large numbers of data items.

CR Categories: H.5.2 [Information Interfaces and Presentation]: User Interfaces—Graphical user interfaces H.2.8 [Database Management]: Database Applications—Data mining

Keywords: Multi-dimensional visualization, pixel-oriented, multi-dimensional scaling, high dimensional datasets.

1 INTRODUCTION

High dimensional datasets are common in applications such as digital libraries, bioinformatics, simulations, process monitoring, and surveys. Automatic analysis tools are widely used for analyzing high dimensional datasets. For example, automatic dimension reduction approaches, such as Principal Component Analysis (PCA) [12] and Multi-dimensional Scaling (MDS) [16], are used to project high dimensional datasets into lower dimensional spaces. Subspace clustering algorithms, such as CLIQUE [1], are used to detect data clusters from high dimensional datasets.

However, due to the dimensionality curse [4], i.e., the lack of data separation in a high dimensional space, finding lower dimensional projections, data clusters and other trends from high dimensional datasets is much harder than it is from low dimensional datasets. Thus graphically presenting high dimensional datasets and allowing the user to apply his or her perceptual abilities and

domain knowledge to make sense of the data is an important approach to both analyzing high dimensional datasets and assessing and understanding the results of automatic analysis tools.

Traditional multi-dimensional visualization techniques, such as glyphs [2], parallel coordinates [10] and scatterplot matrices [6], do not scale well to high dimensional datasets. For example, a dataset containing 200 dimensions will generate star glyphs and parallel coordinates composed of 200 axes. The corresponding scatterplot matrix display would contain 40,000 plots. These large numbers of axes and plots not only clutter the screen but also make it difficult for users to navigate among different dimensions. They make it difficult for users to accomplish exploration tasks such as understanding relationships among dimensions and detecting data clusters and outliers.

In this paper, we propose a new multi-dimensional visualization technique named a *Value and Relation* (VaR) display. By explicitly conveying the relationships among the dimensions as well as data values, the VaR display helps users not only grasp the relationships among the dimensions and navigate within the dimension space, but also detect data clusters and outliers in subspaces composed of subsets of the dimensions. The VaR display uses pixel-oriented techniques [15] to utilize screen space efficiently. It also provides a rich set of navigation and selection tools to enable users to reduce clutter and interactively explore high dimensional datasets.

By graphically presenting each dimension of a high dimensional dataset as a glyph in a 2D space, the VaR display conveys relationships such as correlation among the dimensions through the positions of the glyphs. The positions of the glyphs are generated using Multi-dimensional Scaling (MDS) [16] according to the pairwise relationships among the dimensions. MDS is a technique that maps locations in high dimensional space to positions in a low dimensional space. It is widely used in visualization applications to convey relationships among data items within a multi-dimensional dataset. For example, [21] used MDS to map data items in a document dataset to a 2D space. It generated a Galaxies display as a spatial representation of relationships within the document collection. The VaR display uses MDS in a different way in that it maps dimensions rather than data items to a 2D space according to relationships among the dimensions. In Figure 1a, each dimension of the SkyServer dataset (361 dimensions, 50,000 data items) is mapped to a dot and positioned in the 2D space using MDS. Such a display is called a *star field* display. In this example the correlation among the dimensions is used to generate the positions through MDS. Thus closely related dimensions have positions adjacent to each other in this display. It reveals the correlation among the dimensions intuitively.

Besides the relationships among the dimensions, the VaR display conveys values of the data items using pixel-oriented techniques [15]. Pixel-oriented techniques are visualization methods that map values of data items to the color of pixels and arrange pixels to convey relationships. We use pixel-oriented techniques to map values of the data items within a single dimension to pixels and ar-

*This work was supported under NSF grant IIS-0119276.

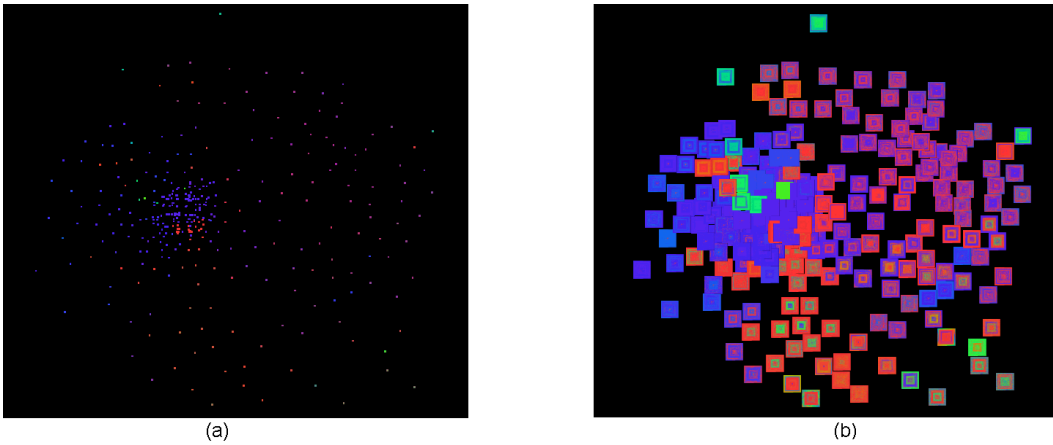


Figure 1: The VaR Display. (a) A star field display where each dimension is mapped to a dot and positioned using MDS according to the correlation among the dimensions. (b) The dots in (a) are replaced by glyphs that present values of the data items to form a VaR display. The dataset is the SkyServer dataset (361 dimensions, 50,000 data items), which was extracted from the Sloan Digital Sky Server (SDSS) data [8].

range them into a *glyph* (“subwindow”, as termed in other papers on pixel-oriented techniques [13]). For each dimension a glyph can be generated. We replace the dots in the star field display by their respective glyphs to produce the VaR display. Figure 1b shows the VaR display of the SkyServer dataset generated by replacing the dots in Figure 1a by glyphs. The textures of the glyphs reveal data patterns in the dimensions.

We provide a rich set of navigation and selection tools for the VaR display. Navigation tools help users reduce clutter in the display and interactively explore the dataset. They include interactions such as overlap reduction, zooming and panning, distortion, comparing, and refining. Selection allows human-driven dimension reduction, i.e., users select subsets of dimensions from the VaR display. Then a space composed of the selected dimensions can be further explored using the VaR display as well as alternative multi-dimensional visualization techniques. Automatic and manual selection tools of the VaR display make selection both flexible and easy to use.

The VaR display can be used for the following purposes:

Visually Exploring High Dimensional Datasets: The VaR display allows users to interactively explore high dimensional datasets with large numbers of data items. It visually reveals both the data item relationships and dimension relationships within a high dimensional dataset.

Guiding Automatic Data Analysis: The VaR display can assist users in (1) assessing and understanding the result of some types of automatic data analysis algorithms and (2) manually tuning the parameters used in those algorithms for better results. For example, by visually presenting the relationships among the dimensions and the values within the dimensions, the VaR display helps users understand the result of an automatic dimension reduction approach. The VaR display can also help users assess the result of an automatic subspace clustering algorithm by visually presenting the clusters.

Human-Driven Dimension Reduction: The VaR display allows users to interactively select dimensions of interest and further explore these dimensions using VaR displays as well as other multi-dimensional visualization techniques. For example, a user can select a group of closely related dimensions from the VaR display, project the dataset into the subspace composed of the selected dimensions, and view the projection using parallel coordinates [10].

2 VaR DISPLAY GENERATION

The following steps are necessary for generating a VaR display for a high dimensional dataset:

1. Build a distance matrix that captures the relationship (such as correlation) between each pair of dimensions in the dataset.
2. Apply MDS on the distance matrix to get a set of positions in a 2D space, where each position corresponds to a dimension.
3. Create a glyph for each dimension that reveals data patterns in that dimension.
4. Place the glyphs in their corresponding positions calculated in step 2.

In step 3, a glyph is created for each dimension of the displayed dataset. It can be generated using pixel-oriented techniques by mapping values of all data items on that dimension to pixels. Each value is represented by the color of a pixel. The pixels are arranged into a glyph using a particular layout scheme such as a spiral layout. Among different glyphs, pixels corresponding to values of the same data items are arranged in the same positions in the glyphs so that users can link values of a data item across different dimensions.

In a VaR display, the positions of the glyphs reveal the relationships among the dimensions. The patterns of the glyphs reveal data patterns in the dimensions. Thus relationships among dimensions can be examined in detail by comparing the patterns of the glyphs.

In the rest of the paper, we study the correlation among the dimensions as a concrete example of relationships among dimensions. However, all the discussion and interaction tools can be readily extended to other types of relationships among dimensions.

There are many possible ways to instantiate the four steps of the VaR display generation. The goal is to find a solution that can provide users with the largest amount of information. The major optimization problems in these steps are how to get a good set of positions for the glyphs and to arrange the pixels within the glyphs to reveal information of interest to users. These two problems are discussed in the following sections.

2.1 Glyph Position Optimization in VaR Display

To get a good set of glyph positions, we first need to identify the factors that affect the glyph positions. The glyph positions generated

by MDS are based on a distance matrix that records the correlation between each pair of dimensions in the dataset. Given that MDS techniques have been widely studied and are mature techniques, we expect that the positions generated by MDS convey the distance matrix with an acceptable quality. However, there are many different correlation measures [11, 3]. Different distance matrices will lead to different sets of positions generated by MDS with the same MDS parameter settings. Thus an appropriate correlation measure needs to be selected.

Second, one needs to make clear what are good glyph positions. We argue that a good VaR display should have the following properties: it helps users locate similar and dissimilar dimensions. According to this argument, the distances between glyphs should vary as much as possible in a good VaR display. In other words, the variance, i.e., the average squared deviation from the mean, of the non-diagonal elements in the distance matrix should be as large as possible. Thus the glyph position optimization problem can be expressed as follows: Find a correlation measure so that the variance of all non-diagonal elements of the distance matrix reaches a maximum. We name the variance of the non-diagonal elements of the distance matrix the *variance criteria* of the glyph position optimization problem.

It is impossible to get an optimal solution to this problem since there are infinite possible correlation measures. Thus we use a heuristic approach to calculate a distance matrix based on the fact that in a large-scale dataset, two dimensions might be closely related in subsets of the data items rather than in all data items.

In this approach, data values are first normalized within each dimension for invariance against scaling [3]. Second, for each pair of dimensions, the data items are divided into subsets within which data items have similar value differences between the two dimensions. In particular, an equal-width histogram of the value differences of all data items between the two dimensions is used. It naturally divides the data items into subsets by their value difference and records the number of data items in each subset.

Then the top k (k is decided according to the variance criteria) subsets with the largest populations are picked out. The sum of their populations is proportional to the correlation between the two dimensions. For example, if 50% data items have value differences of 0, while another 40% data items have value differences of 0.5 between two dimensions, the correlation between them is 0.5 if k equals to 1 while the correlation is 0.9 if k equals to 2 in a 0 to 1 correlation scale. Given a k , the correlation between each pair of dimensions can be calculated. Thus a distance matrix can be generated. To decide the optimal k , a distance matrix is built for each possible k . The variance of the non-diagonal elements of each distance matrix is calculated and the distance matrix with the largest variance will be chosen.

The above approach can be described using the following steps:

1. Normalize the values within each dimension;
2. For each pair of dimensions, build an equal-width histogram of the value differences of all data items between the two dimensions. The number of bins in each histogram is $numBins$, which is a constant that can be set by the user;
3. For $i = 1$ to $numBins$ calculate a distance matrix $Matrix_i$ and its non-diagonal elements variance Var_i :
 - For each pair of dimensions calculate a distance according to the histogram: sort the bins in the histogram according to the number of data items falling into them in a decreasing order. Take the first i bins. The correlation of the dimensions is the percentage of data items falling into these i bins. Their distance is one minus the correlation.

- Build a distance matrix $Matrix_i$ using the calculated distances. Calculate its variance criteria Var_i .

4. Find the maximum variance Var_k from Var_i ($i = 1, \dots, numBins$). Output $Matrix_j$.

For a large-scale dataset stored in a low-speed memory, CPU computational cost can be ignored with regard to I/O cost. Thus the time complexity of this algorithm is analyzed as follows: step 1 is a common step needed for most correlation measures and needs two scans of the dataset. Step 2 requires one scan of the dataset with each data item accessed to build the histograms. But it can be combined into the second scan of step 1. Thus its cost can be ignored. The dimensionality is usually much smaller than the number of data items in a large-scale dataset. $NumBins$ is a constant usually chosen to be much smaller than the number of data items. Thus the histograms can be stored in high-speed memory. Therefore step 3 and step 4 are computations without I/O cost. So the total cost of the algorithm in terms of external data access is two scans of the dataset, which is similar to most common correlation calculations.

We have run a series of experiments to compare the variance criteria calculated using the presented algorithm with that calculated using a global Euclidean similarity measure [3]. The experiments were run on 10 real datasets whose dimensionalities range from 4 to 361 and the numbers of data items range from 256 to 95130. The results showed that the variance calculated using the presented algorithm is 45% to 95% larger than that calculated using the global Euclidean similarity measure for the tested dataset in all experiments. Thus the presented algorithm generated more qualitative distance matrices than using the global Euclidean similarity measure. The reason is that the presented algorithm captures similarity between two dimensions in a finer granularity than the global Euclidean similarity measure, thus is good for large-scale datasets.

2.2 Pixel Arrangement in Glyphs

Pixel arrangement is an important issue for pixel-oriented techniques [13]. Proposed solutions include spiral arrangements, space-filling curves, and axes techniques, among many others [13]. As an initial layout, we use a simple spiral arrangement [14]. Given an order of the data items, the pixels are placed from the center of a square to the outside of the square spirally according to that order. More pixel arrangement schemes will be explored in the future.

One option is to employ the original order of the data for time series data or other datasets where order is meaningful. Besides that, the VaR display allows the user to reorder the data items by their values in one dimension. Such a dimension is called a *base dimension*. Figure 2 shows that selection of the base dimension greatly affects the information conveyed by a VaR display. Patterns existing in dimensions closely related to the base dimension are more explicitly presented than those existing in other dimensions.

In the VaR display, a base dimension is automatically selected by default when the VaR display is initially presented to the user. The selection criteria is that the base dimension should be a dimension that has the largest number of closely related dimensions so that patterns of the largest number of dimensions are better conveyed in the initial view of a VaR display.

Since “closely related” is a subjective measure, a heuristic approach is used to find the initial base dimension. This initial dimension is chosen such that it has the smallest total distance to all other dimensions in the distance matrix. The user can select another dimension as the base dimension using the manual pixel reordering tool provided by the VaR display (see Section 3.1).

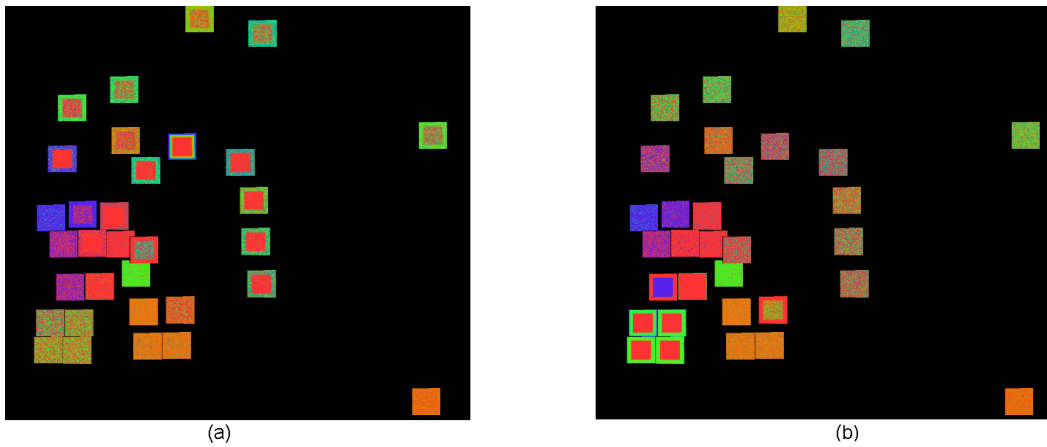


Figure 2: Ordering of Pixels. Pixels in the glyphs are ordered by values of data items in different dimensions in (a) and (b). (a) It is clearly visible that the dimensions in the top of the display are closely related since their glyphs have similar patterns. (b) It is clearly visible that the dimensions in the bottom left of the display are closely related. The dataset is the Census-Income-Part dataset (42 dimensions, 20,000 data items), which is a subset of the Census-Income dataset extracted from the Census Bureau database [17].

3 INTERACTIVE TOOLS IN THE VaR DISPLAY

A rich set of navigation and selection tools has been developed for the VaR display. Navigation tools help users reduce clutter of the display and learn information about the dataset. Automatic and manual selection tools allow users to perform human-driven dimension reduction by selecting subsets of dimensions for further exploration using the VaR display as well as other multi-dimensional visualization techniques. In the following sections, details of each navigation and selection tool will be presented.

3.1 Navigation Tools

Different from all the other pixel-oriented techniques, where each pixel is assigned a unique position on the screen, our VaR display allows overlaps among the glyphs. Overlaps emphasize close relationships among the dimensions because glyphs overlap only if their dimensions are closely related. However, overlaps can prevent a user from seeing details of an overlapped glyph. We provide the following operations to overcome this problem:

Showing Names: By putting the cursor on the VaR display, the dimension names of all glyphs under the cursor position are shown in a message bar. Thus a user can be aware of the existence of glyphs hidden by other glyphs.

Layer Reordering: With a mouse click, a user can force a glyph to be displayed in front of the others. In this way he/she can view details of a glyph originally overlapped.

Manual Relocation: By holding the control key, a user can drag and drop a glyph to whatever position he/she likes. In this way a user can separate overlapping glyphs.

Extent Scaling: Extent scaling allows a user to interactively decrease the sizes of all the glyphs proportionally to reduce overlaps, or to increase them to see larger glyphs. Figure 3b gives an example of extent scaling.

Dynamic Masking: Dynamic masking allows users to hide the glyphs of unselected dimensions from the VaR display. In Figure 5, the glyphs of unselected dimensions are hidden using dynamic masking.

Automatic Shifting: This operation automatically reduces the overlaps among the glyphs by slightly shifting the positions of the glyphs. Figure 3c gives an example of automatic shifting using a simple distortion algorithm for reducing glyph overlaps borrowed from [20]. There are many more advanced overlap reducing algorithms that can be used, as those listed in [20].

Other navigation tools provided by the VaR display include:

Distortion: Users can interactively enlarge the size of some glyphs while keeping the size of all other glyphs fixed. In this way users are allowed to examine details of textures of the enlarged glyphs within the context provided by the other glyphs. Figure 3d gives an example of distortion.

Zooming and Panning: Users can zoom in, zoom out and pan the VaR display. For example, in order to reduce overlaps, sometimes the size of the glyphs has to be set very small when there are a large number of dimensions. Zooming into the display will enlarge the glyphs so that the user can have a clear view of the texture of the glyphs.

Manual Pixel Reordering: Users can click the middle button on a glyph to reorder the data items according to their values on the dimension corresponding to that glyph. That dimension is called the base dimension, as discussed in Section 2.2. Glyphs will have different textures for different base dimensions. Thus different patterns of the dataset will be revealed as users reorder the data items by different base dimensions. Figure 2 was generated using manual pixel reordering.

Comparing: It is important to allow a user to compare the values of the data items in one dimension with those in other dimensions so that the relationship between a dimension and other dimensions can be revealed in a more intuitive manner. We allow users to switch to a comparison mode. In comparison mode, except the glyph of the base dimension, the pixels of all other glyphs will be colored according to the differences between the values of the base dimension and their dimensions. Figure 4 shows an example of the comparison operation.

Refining: A refined VaR display can be generated for selected dimensions where only the glyphs of the selected dimensions

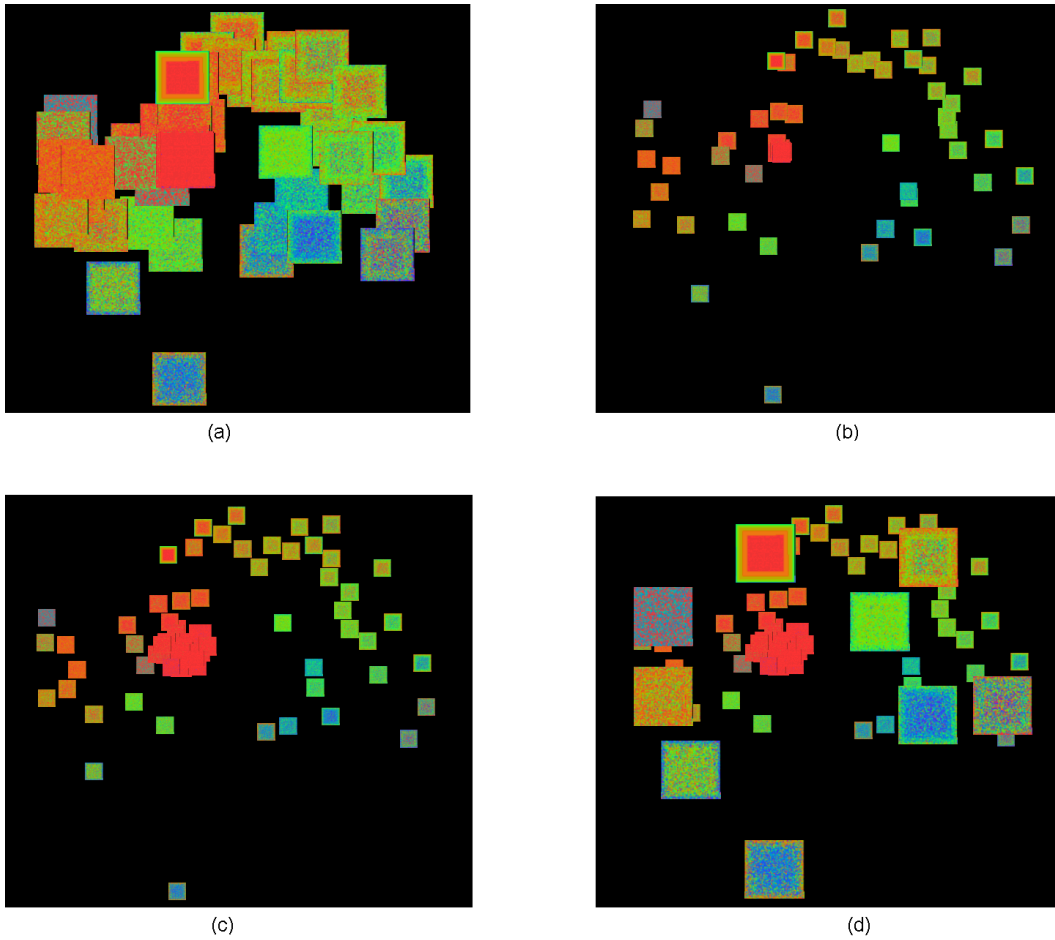


Figure 3: Extent Scaling, Automatic Shifting and Distortion. (a): A VaR display with seriously overlapped glyphs. (b) Overlap is reduced by decreasing the size of the glyphs. (c) Overlap of (b) is further reduced by automatic shifting. Notice that several glyphs appear in the center of the display which are previously non-visible in (b) due to overlaps. (d) Some glyphs are enlarged to examine detail within context. The dataset is the Ticdata2000 dataset (86 dimensions, 5,822 data items), which contains information on customers of an insurance company [18].

are shown, and their positions are relocated by MDS using relationships among only the selected dimensions. The glyph positions in the refined VaR display reflect the relationships among the selected dimensions more accurately than in the original VaR display since the effect on the positions of unselected dimensions is filtered.

3.2 Selection

Selection tools enable users to select dimensions of interest for further exploration using other multi-dimensional visualization techniques. They can also be used as a filter to reduce the number of glyphs displayed in a VaR display since we allow users to hide glyphs of unselected dimensions using dynamic masking (see Section 3.1). The selection tools we provide to users include an automatic selection tool for closely related dimensions, an automatic selection tool for well separated dimensions, and manual selection.

The automatic selection tool for related dimensions takes a user-assigned dimension and correlation threshold as input. Users can select the assigned dimension by clicking its glyph and adjust the threshold through a slide bar. The tool automatically selects all dimensions whose correlation measures to the input dimension are smaller than the threshold by traveling through the distance matrix. This tool enables the user to select a set of closely related dimen-

sions.

The automatic selection tool for separated dimensions takes a user-assigned dimension and correlation threshold as input and returns a set of dimensions that describe the major features of the dataset. The assigned dimension will be included in the returned set of dimensions. Between each pair of dimensions in the result set, the correlation measure is larger than the threshold. For any dimension that is not in the result set, there is at least one dimension in the result set such that the correlation measure between it and the unselected dimension is smaller than the threshold. Using this tool, a user is able to select a set of dimensions to construct a lower dimensional subspace revealing the major features of the dataset without much redundancy. Figure 5 shows an example of automatic selection for separated dimensions in a high dimensional dataset.

The following algorithm can be used for automatic selection of separated dimensions:

1. Set the assigned dimension as “selected” and all other dimensions as “unselected”.
2. Find all unselected dimensions whose distances to all existing selected dimensions are larger than the threshold. Mark them as “candidates”.
3. If there is no candidate dimension, go to step 4. Else, set one

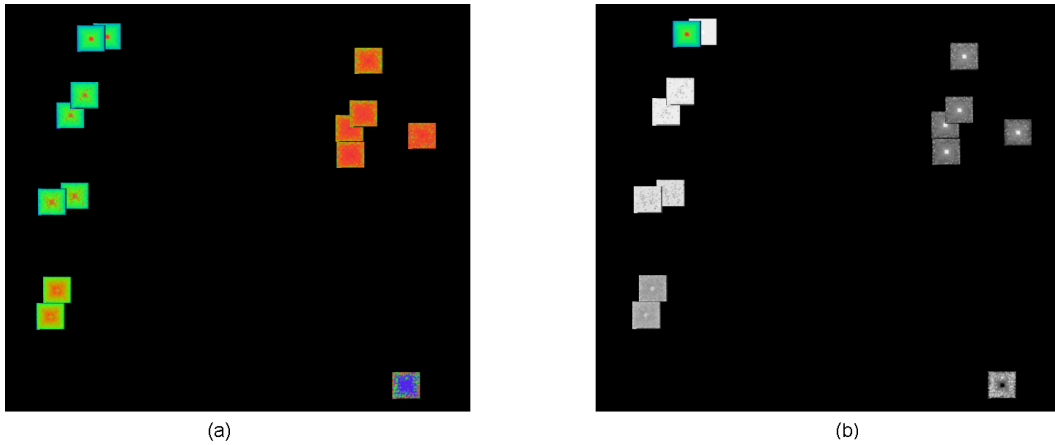


Figure 4: Comparing. (a): Original VaR display. The base dimension is in the top left of the display. (b): Comparison mode. The value differences between the base dimension and all other dimensions are visualized. The bigger the value difference, the darker a pixel is. Similar dimensions are clearly visible since they are brighter than dissimilar ones. The dataset is the AAUP salary dataset (14 dimensions, 1,161 data items).

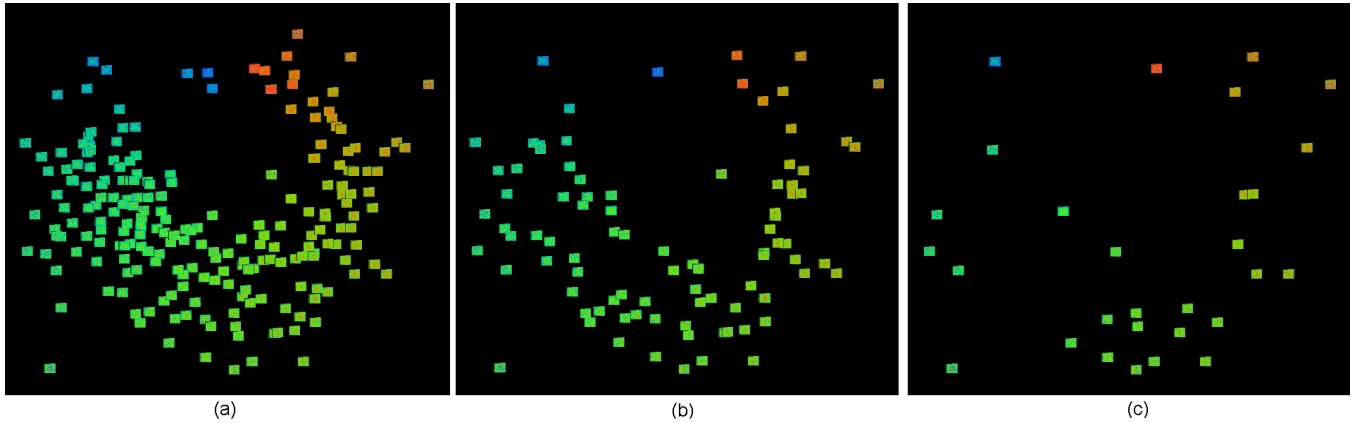


Figure 5: Automatic Selection of Separated Dimensions. Unselected dimensions are hidden using Dynamic Scaling. Selected dimensions in (a)(b)(c) are generated using automatic selection with the same assigned dimensions and an increasing correlation threshold. The dataset is the OHSUMED dataset (215 dimensions, 298 data items), which contains the word-counts of a medical abstract collection [9].

candidate dimension as “selected” and other candidate dimensions as “unselected”. Go back to step 2.

- Return all dimensions marked as “selected”.

Manual selection allows a user to manually select a dimension by clicking its corresponding glyph. The user can unselect a dimension by clicking the glyph again. The combination of manual and automatic selection makes the selection operation both flexible and easy to use.

4 SCALING TO LARGE NUMBERS OF DATA ITEMS

We have implemented a working prototype of the VaR display and its interaction tools in XmdvTool [19], a public-domain visualization system. In order to scale the VaR display to datasets with large numbers of data items, we have integrated sampling and texture mapping techniques. These techniques allow the VaR display to handle datasets with large numbers of data items efficiently.

The prototype stores datasets in an Oracle9i database server. It dynamically requests data from the server when needed. When generating a VaR display for a dataset containing a large number of data

items, we use a random sampling approach to reduce the response time for fetching data items from the server. In particular, the system keeps a default maximum number. When the number of data items contained in a dataset exceeds it, a uniform random sampling is performed on the dataset to only fetch the maximum number of data items. Users are allowed to interactively adjust the maximum number in order to reduce the response time or increase accuracy. Figure 6 shows two VaR displays of a dataset with and without sampling. It can be seen that the corresponding glyphs in the two displays have very similar patterns. However, more strict analysis of information loss caused by sampling and deeper research on sampling strategies to be used need to be performed in the future.

Secondly, in order to reduce the response time of user interactions for large-scale datasets, we store all glyphs as texture objects in OpenGL. Thus unless we need to regenerate the texture of the glyphs, each glyph can be refreshed, repositioned, or resized on the screen by simply redrawing the texture objects, mapping the texture objects to different positions on the screen, or mapping them to areas of different sizes. All these operations can be efficiently performed through OpenGL.

Both the above two approaches cause information loss in the VaR display. When random sampling is performed, data items not in the

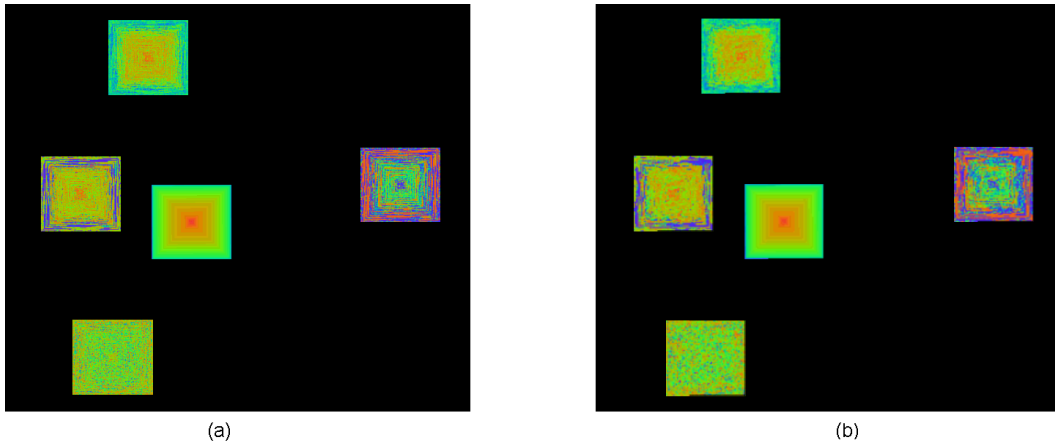


Figure 6: Approximation. (a): All 16,384 data items are displayed. The glyphs are shown in their original sizes. (b): A sample of 5,000 data items is displayed. The glyphs are magnified to the same sizes as in (a) using texture mapping. It can be seen that the corresponding glyphs in the two displays have very similar patterns. The dataset is the Out5d remote sensing dataset (5 dimensions, 16,384 data items).

sample are not visually presented to the user. When the texture objects are mapped to screen areas that are not exactly their original sizes, magnification or minification happens so that the pixels visualized are only approximations of the original pixels. However, information loss is exchanged for the reduction of response time and clutter in the display, which is very important for interactive visual exploration. Moreover, approximation is usually acceptable in a visualization system. Furthermore, users can always get the information accurately by setting the sampling threshold to a number no less than the number of data items contained in the dataset, and setting the size of the glyphs to exactly the size of the texture objects. System provided reset buttons allow users to do that easily.

5 CASE STUDY

A case study has been conducted on a real dataset, the Census-Income-Part dataset containing 42 dimensions and 20,000 data items. Its VaR display is shown in Figure 2. We accomplished the following tasks by interactively exploring the dataset through the VaR display:

- We are able to detect groups of closely related dimensions using three methods: (1) looking for glyphs clustered together in the VaR display; (2) looking for glyphs with similar patterns; (3) selecting dimensions closely related to a dimension of interest using the automatic selection tool for related dimensions. Using these three methods together helps us get results quickly and intuitively. Figure 2b shows a group of closed related dimensions in the bottom left of the display. By checking the dimension names we found that these are all dimensions recording people’s migration and moving status in the last year.
- We are able to find data clusters in a subset of the dimensions from similar patterns of the graphs. For example, in Figure 2b, within each glyph in the bottom left of the display, pixels in the center area have a different color from that in the outer area. We then determined that the data is divided into two clusters in those dimensions, which are the people who did not move in the last year and people who moved in the last year.
- We are able to find well separated dimensions of the dataset using three methods: (1) looking for glyphs evenly distributed

in the display; (2) looking for glyphs with significantly different patterns; (3) selecting well separated dimensions using the automatic selection tool for separated dimensions. Using these three methods together helps us get results quickly and intuitively.

- We are able to find dimensions with special patterns. For example, there were several dimensions with lots of values mapped to red in the VaR display. According to the color code we found that those dimensions contain a high rate of missing values. We can remove them from the display.

Through the case study we found that the VaR display and its navigation and selection tools could help users discover interesting patterns in a high dimensional dataset with a large number of data items effectively and efficiently.

6 RELATED WORK

Multi-dimensional Scaling (MDS) [16] is an iterative non-linear optimization algorithm for projecting multi-dimensional data down to a reduced number of dimensions. It is often used to convey relationships among data items of a multi-dimensional dataset. In our approach, MDS is used in a different way, namely to convey relationships among dimensions rather than data items.

Pixel-oriented visualization techniques [15, 13] are a family of multi-dimensional display techniques that map each data value to a pixel on the screen and arrange the pixels in such a way as to convey relationships. They generate condensed displays that may reveal clusters, trends, and anomalies. The VaR display is different from existing pixel-oriented visualizations because it uses positions of the subwindows (glyphs in the VaR display) to accurately convey the relationships among the dimensions. In addition, many interactions of the VaR display, such as extent scaling and comparing, have not previously been applied to pixel-oriented techniques.

The VHDR [24] and DOSFA approaches [23] explicitly convey the relationships among the dimensions of a high dimensional dataset using a dimension hierarchy. They allow users to interactively navigate and select dimensions from it. The VaR display is different in that it uses MDS to convey the relationships among the dimensions. In addition, the VaR display conveys values of data items, while VHDR and DOSFA do not.

Sampling has been used in pixel-oriented visualization systems. VisDB [7] allows users to interactively change the number of data

items displayed on the screen using sampling. The VaR display uses sampling in a similar way to limit the number of data items fetched in order to reduce I/O cost.

7 CONCLUSION

The major contributions of this paper are:

- A new method for the display of high dimensional datasets, the VaR display, has been proposed and developed. The VaR display not only conveys values of the data items to the users, but also explicitly conveys relationships among the dimensions of a high dimensional dataset.
- A rich set of navigation and selection tools for the VaR display has been implemented to allow users to interactively explore the dataset displayed. The navigation tools help users identify patterns hidden in a high dimensional dataset effectively. The selection tools enable users to interactively select dimensions of interest from the VaR display for further exploration.
- Criteria and algorithms for the distance matrix generation and the base dimension selection have been created for generating an informative VaR display among many possible ones.
- Sampling methods and texture mapping have been used to enable the VaR display to efficiently scale to datasets with large numbers of data items.

There are many open issues in the VaR display to be explored, such as:

- Are pixel-oriented techniques the only choice for generating glyphs in the VaR display? We will explore some alternatives, such as using 1-D histograms of the dimensions as the glyphs, or using 2-D scatterplots of a selected dimension and each of the other dimensions as the glyphs to allow users to compare the selected dimension with other dimensions more effectively.
- Is correlation (similarity) the only relationship among the dimensions that can be conveyed by the VaR display? Will a VaR display where dissimilar dimensions are close to each other help users locate informative subspaces more easily?
- What is the information loss caused by the sampling approach? What levels of sampling are enough and what kinds of sampling methods are appropriate for tasks such as revealing the major clusters of the datasets or emphasizing the outliers? We plan to explore these problems through two approaches. One approach is to compare the generated figures with and without sampling as in [22]. The other approach is to compare the distribution of the sampled data with the original data as in [5].

We also plan to explore different pixel arrangement approaches in glyph construction, develop interactions to allow users to compare and highlight values of the same data in different dimensions, and evaluate the effectiveness and efficiency of the proposed approach using more formal experiments and user studies.

Acknowledgments We gratefully thank Dr. Daniel A. Keim, who gave many valuable suggestions for this work.

REFERENCES

- [1] R. Agrawal, J. Gehrke, D. Gunopulos, and P. Raghavan. Automatic subspace clustering of high dimensional data for data mining applications. *Proc. ACM SIGMOD International Conference on Management of Data*, pages 94–105, 1998.
- [2] D.F. Andrews. Plots of high dimensional data. *Biometrics*, 28:125–136, 1972.
- [3] M. Ankerst, S. Berchtold, and D.A. Keim. Similarity clustering of dimensions for an enhanced visualization of multidimensional data. *Proc. IEEE Symposium on Information Visualization*, pages 52–60, 1998.
- [4] K. Beyer, J. Goldstein, R. Ramakrishnan, and U. Shaft. When is “nearest neighbor” meaningful? *Lecture Notes in Computer Science*, 1540:217–235, 1999.
- [5] S. Chaudhuri, R. Motwani, and V. Narasayya. Random sampling for histogram construction: how much is enough? *Proc. ACM SIGMOD International Conference on Management of Data*, pages 436–447, 1998.
- [6] W.S. Cleveland and M.E. McGill. *Dynamic Graphics for Statistics*. Wadsworth, Inc., 1988.
- [7] Database Systems Research Group, Institute of Computer Science, University of Munich. Visdb: A visual data mining and database exploration system. <http://www.dbs.informatik.uni-muenchen.de/dbs/projekt/visdb/visdb.html>.
- [8] J. Gray, A. Szalay, A. Thakar, P. Z. Zunszt, T. Malik, J. Raddick, C. Stoughton, and J. vandenBerg. The SDSS SkyServer - public access to the Sloan Digital Sky Server data. Technical Report MSR-TR-2001-104, Microsoft, 2001.
- [9] W. Hersh, C. Buckley, T. Leone, and D. Hickman. Ohsumed: An interactive retrieval evaluation and new large text collection for research. *Proc. ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 192–201, 1994.
- [10] A. Inselberg. The plane with parallel coordinates. *Special Issue on Computational Geometry, The Visual Computer*, 1:69–97, 1985.
- [11] R.A. Johnson and D.W. Wichern. *Applied Multivariate Statistical Analysis*. Prentice Hall International, Inc., second edition, 1988.
- [12] J. Jolliffe. *Principal Component Analysis*. Springer Verlag, 1986.
- [13] D.A. Keim. Designing pixel-oriented visualization techniques: Theory and applications. *IEEE Transactions on Visualization and Computer Graphics*, 6(1):1–20, January-March 2000.
- [14] D.A. Keim and H.-P. Kriegel. Visdb: Database exploration using multidimensional visualization. *IEEE Computer Graphics & Applications*, 14(5):40–49, 1994.
- [15] D.A. Keim, H.-P. Kriegel, and M. Ankerst. Recursive pattern: a technique for visualizing very large amounts of data. *Proc. IEEE Visualization '95*, pages 279–286, 1995.
- [16] J.B. Kruskal and M. Wish. *Multidimensional Scaling*. Sage Publications, 1978.
- [17] U.S. Census Bureau. The census bureau database. <http://www.census.gov/ftp/pub/DES/www/welcome.html>, 1997.
- [18] P. van der Putten and M. van Someren. Coil challenge 2000: The insurance company case. Technical Report Technical Report 2000-09, Leiden Institute of Advanced Computer Science, 2000.
- [19] M.O. Ward. Xmdvtool: Integrating multiple methods for visualizing multivariate data. *Proc. IEEE Visualization*, pages 326–333, 1994.
- [20] M.O. Ward. A taxonomy of glyph placement strategies for multidimensional data visualization. *Information Visualization*, 1(3-4):194–210, 2002.
- [21] J.A. Wise, J.J. Thomas, K. Pennock, D. Lantrip, M. Pottier, A. Schur, and V. Crow. Visualizing the non-visual: Spatial analysis and interaction with information from text documents. *Proc. IEEE Symposium on Information Visualization*, pages 51–58, 1995.
- [22] P. C. Wong, H. Foote, D. Adams, W. Cowley, and J. Thomas. Dynamic visualization of transient data streams. *Proc. IEEE Symposium on Information Visualization*, pages 97–104, 2003.
- [23] J. Yang, W. Peng, M.O. Ward, and E.A. Rundensteiner. Interactive hierarchical dimension ordering, spacing and filtering for exploration of high dimensional datasets. *Proc. IEEE Symposium on Information Visualization*, pages 105–112, 2003.
- [24] J. Yang, M.O. Ward, E.A. Rundensteiner, and S. Huang. Visual hierarchical dimension reduction for exploration of high dimensional datasets. *Eurographics/IEEE TCVG Symposium on Visualization*, pages 19–28, 2003.