

A Visual Method for High-dimensional Data Cluster Exploration

Ke-Bing Zhang¹, Mao Lin Huang², Mehmet A. Orgun¹, and Quang Vinh Nguyen³

¹Department of Computing, Macquarie University, Sydney, NSW 2109, Australia
{kebing, mehmet}@science.mq.edu.au

²Faculty of Engineering and Information Technology, University of Technology, Sydney
NSW 2007, Australia
maolin@it.uts.edu.au

³School of Computing and Mathematics, University of Western Sydney, NSW 1797, Australia
vinh@scm.uws.edu.au

Abstract. Visualization is helpful for clustering high dimensional data. The goals of visualization in data mining are exploration, confirmation and presentation. However, the most of visual techniques serviced for cluster analysis are focused on cluster presentation rather than cluster exploration. Several techniques are proposed to explore cluster information by visualization, but most of them heavily depend on the individual user's experience. Inevitably, it incurs subjectivity and randomness in the clustering process. In this paper, we employ the statistical features of datasets as predictions to estimate the number of clusters by a visual technique, HOV³. This approach avoids the randomness and subjectivity of the user during the process of cluster exploration by other visual techniques. As a result, it provides an effective visual method for cluster exploration.

Keywords: Cluster Exploration, Visualization, Statistics

1 Introduction

Cluster analysis is an important technique of knowledge acquisition in data mining. To address the requirements of different applications, a large number of clustering algorithms have been developed [9, 3]. However, those algorithms are not very effective to cope with arbitrarily shaped clusters. In addition, cluster analysis is a highly iterative process. But the most of existing clustering methods are too automated to exclude the domain experts' knowledge in the intermediate process of clustering. As a consequence, they are not always effective to cluster datasets with a large number of variables and/or huge-sized datasets in real world applications. In a high dimensional space, traditional clustering algorithms tend to break down in terms of efficiency as well as accuracy because data does not cluster well anymore [1].

In order to solve those problems, Shneiderman [19] proposed that, to present data as a visual plot so that the interesting features could be seen by a human researcher. He pointed out that, visualization can be very powerful and effective in revealing trends, highlighting outliers, showing clusters, and exposing gaps in high-dimensional data analysis. Therefore, the use of visualization to explore and understand high-

The datasets exemplified in this paper are available from <http://archive.ics.uci.edu/ml/>

dimensional datasets is becoming an efficient way to combine human intelligence with the immense brute force computation power available nowadays [16].

Clustering is an exploratory activity [9]. It is an iterative process under the guidance of user domain knowledge. In most cases of the preprocessing stage of clustering, it is hard for the user to estimate the proper cluster number [3]. Visualization is very helpful for the user to do that. However, cluster exploration by visualization mostly depends on the individual user's experience. Thus, subjectivity, randomness and impreciseness may be introduced into the cluster exploration process. As a result, cluster analysis based on imprecise results may be inefficient and ineffective. On the other hand, cluster exploration based on the user's random interaction is arbitrary and it may not be easy to interpret where the grouped results come from.

In this paper, based on the projection of HOV^3 [22], we introduce the statistical features of datasets as the predictions of HOV^3 to guide the user on cluster exploration, because the statistical summaries objectively reflect the features of datasets. As a result, it provides the user an effective method on determining cluster numbers in the preprocessing stage of cluster analysis.

The rest of this paper is organized as follows. Section 2 briefly introduces visual cluster analysis and cluster visualization techniques as the background of this research. Section 3 provides a short introduction of the HOV^3 technique as the preliminary knowledge of this study. Section 4 presents the statistics-guided visual approach for cluster detection by HOV^3 , and demonstrates its effectiveness by experiments on several datasets. Finally, Section 5 summarizes the contributions of this paper.

2 Background

2.1 Visual Cluster Analysis

Visual cluster analysis is a combination of visualization and cluster analysis. It is believed that the combined strength of visualization and data mining would enrich both approaches and enable more successful solutions [20]. However, the data that clustering processed is usually high dimensional. It is not easy to visualize multidimensional data on 2D or 3D space and still give a "genuine" visual interpretation. Because mapping higher dimensional data onto lower dimensional space inevitably introduces ambiguities, overlapping and even bias. Thus, choosing a technique to fit visualizing clusters of high dimensional data is the first and crucial task of visual cluster analysis.

In practice, instead of providing a quantitative guidance on cluster exploration, most of the cluster visualization techniques are typically used as an observational mechanism to assist the user in having intuitive comparisons and better understanding of clustering results. Several approaches are proposed to help the user on cluster exploration.

For example, Multidimensional scaling (MDS) maps multidimensional data as points into 2D Euclidean space, where the distances between data points reflect the

similarity/dissimilarity of them [14]. However, the relative high computational cost of MDS (polynomial time complexity $O(N^2)$) limits its usability on very large datasets. PCA is a commonly used multivariate analysis technique [10], mainly used for reducing the dimensionality of high dimensional data by extracting the representative variables. But PCA is sensitive to deal with the non-linear data structure. It is not suitable for the exploration of unknown data. A Grand Tour based visual technique is proposed to visualize cluster structure [5], but this technique visualizes 3 clusters. To deal with more than 3 clusters with a more sophisticated Grand Tour technique, more assistance is required.

OPTICS uses a density-based technique to detect cluster structure and visualizes them in “Gaussian bumps” [2]. It is an intuitive method to assist the user to observe cluster structures. But its non-linear time complexity makes it neither suitable to deal with very large data sets, nor suitable to provide the contrast between clustering results.

Huang *et. al* [7, 8] proposed several approaches to assist users in identifying and verifying the validity of clusters in visual form. Their techniques work well in cluster identification, but are unable to evaluate the cluster quality very well. On the other hand, these techniques are not well suited to the interactive investigation of data distributions of high-dimensional data sets.

CVAP [21] is a recently proposed prototype which several integrated clustering algorithms and cluster validation methods. It is a convenient toolkit to assist the user on the selection of clustering scheme for the application of small-sized datasets. However, CVAP is only for displaying the clustering and cluster validation results, rather than the purpose for directly evolving the user into the cluster exploration process.

2.2 Star Coordinates

The projection of Star Coordinates [11] has only linear time complexity, which is significant for interactive cluster visualization of very large datasets. VISTA [4] and HOV³ [22] extend Star Coordinates by additional features to mitigate the problem of overlapping and ambiguities caused by projecting high dimensional data onto 2D space. The visual approach reported in this paper has been developed based on the projection of HOV³. For the sake of completeness, we briefly introduce the Star Coordinates technique here.

Star Coordinates plots a 2D plane into n equal sectors with n coordinate axes, where each axis represents a dimension and all axes share the initials at the centre of a circle surface on the 2D space [11]. Star Coordinates first normalizes data in each dimension into a unit interval $[0, 1]$. Then the values of all axes are mapped to an orthogonal X-Y coordinate which shares the centre point with Star Coordinates on the 2D space. Thus, an n -dimensional data item is represented as a point in the X-Y 2D plane by Star Coordinates. Based on this projection, several interactions, such as axis scaling, axis rotation, data point filtering, etc. are provided in Star Coordinates to change the data distribution of a dataset in order to detect cluster characteristics and render clustering results with the interactions.

However, it is not easy to give an explanation of the grouping results produced by the user’s random interactions in Star Coordinates and VISTA, also the grouping

results are usually not repeatable. On the other hand, in Star Coordinates space, the user's interactions cannot change the data distribution too much when the dimensionality of the dataset is very high (a hundred or more dimensions, which is very common in data mining). This is because the alteration of the data distribution by applying interactions to an axis is much less than that of lower dimensional data in the Star Coordinates space. As a result, in very high dimensional space, it is not effective anymore to separate clusters or explore grouping clues by the interactions of Star Coordinates and VISTA.

As discussed above, the issues of arbitrary exploration and/or complicated visual representation of cluster structure make those techniques inefficient and time consuming on cluster exploration of large and high dimensional data. As Seo and Shneiderman [18] mentioned that "A large number of clustering algorithms have been developed, but only a small number of cluster visualization tools are available to facilitate researchers' understanding of the clustering results". Thus developing an effective visualization technique to assist the user during cluster exploration and detection is the main aim of this research.

3 HOV³

To remedy the randomness and arbitrariness of visualization on cluster analysis, Zhang *et al.* mathematically generalized the Star Coordinates model by the Euler formula and proposed their visual approach HOV³ to detect clusters [22]. According to the Euler formula: $e^{ix} = \cos x + i \sin x$, where $z = x + i.y$, and i is the imaginary unit. Let $z_0 = e^{2\pi i/n}$, such that $z_0^1, z_0^2, z_0^3, \dots, z_0^{n-1}, z_0^n$ (with $z_0^n = 1$) divide the unit circle on the complex 2D plane into n equal sectors. Then Star Coordinates can be simply written as:

$$P_j(z_o) = \sum_{k=1}^n [(d_{jk} - \min_k d_k) / (\max_k d_k - \min_k d_k) \cdot z_o^k] \quad (1)$$

where $\min_k d_k$ and $\max_k d_k$ represent the minimal and maximal values of the k th coordinate respectively. In any case equation (1) can be viewed as mapping from $R^n \rightarrow C^2$.

Conversely, instead of using random exploration of cluster information by axis scaling or axis rotation in Star Coordinates/VISTA, HOV³ quantifies the user's priori knowledge/estimation of a studied dataset as a measure vector to precisely guide the user on the exploration of group information. A measure vector M in HOV³ represents the corresponding axes' weight values. Then given a non-zero measure vector M in R^n , and a family of vectors P_j , the projection of P_j against M , according to formula (1), the HOV³ model is presented as:

$$P_j(z_o) = \sum_{k=1}^n [(d_{jk} - \min_k d_k) / (\max_k d_k - \min_k d_k) \cdot z_o^k \cdot m_k] \quad (2)$$

where m_k is the k th variable of measure M .

It can be observed that, equation (2) is a standard form of linear transformation of n variables, where m_k is the coefficient of k th variable of P_j .

4 Cluster Exploration by HOV³

4.1 The Idea

In analytic geometry, the difference of two vectors A and B can be expressed by their inner product $A \cdot B$, its geometrical meaning is that, the data distribution plotted by vector A against vector B (vice versa). The inner product between a dataset and a measure vector in HOV³ can be geometrically viewed as a data distribution plotted by a set of vectors against the measure vector in the HOV³ space, as shown in the equation (2).

Predictive knowledge discovery is an important knowledge acquisition method, which utilizes the existing knowledge to deduce, infer, reason and establish predictions, and verify the validity of the predictions. As mentioned above, the user can quantify his/her priori knowledge of a studied dataset as the guidance on the exploration of group information. Thus the statistical summaries of a dataset can be directly employed as the statistical predictions (measure vectors) of the dataset in HOV³, since the statistical summaries reflect the nature comparisons of data objectively [19]. Also, it is easy to interpret the grouping results of a dataset plotted by statistical predictions in HOV³.

Based on this idea, we propose a statistics-guided cluster exploration approach by HOV³. The detailed description of our approach is presented in the follows.

4.2 The Algorithm and the Features

4.2.1 The Algorithm

Table 1. The Algorithm of Statistics-guided Cluster Exploration by HOV³

Algorithm: Statistics-guided Cluster Exploration by HOV³	
Input: D : a dataset; M : statistical measures of D ;	
Output: G : data distribution of D or subsets of D ;	
1:	cluster exploration \leftarrow true;
2:	$p \leftarrow D$;
3:	$m_i \leftarrow$ a statistical measure of p ; ($m_i \in M$)
4:	$m \leftarrow m_i$;
5:	while (cluster exploration)
6:	$G \leftarrow Hc(D, m_i)$;
7:	if (G well grouped?)
8:	if (stop exploration?)
9:	cluster exploration \leftarrow false;
10:	break;
11:	endif
12:	endif
13:	if (new statistical measure of p ?)
14:	$m_i \leftarrow$ a statistical measure of p ;
15:	$m \leftarrow m_i$;
16:	endif
17:	$m_i \leftarrow m \cdot m_i$;
18:	endwhile

We formalized our idea of using statistical predictions to explore clusters by HOV³ into the algorithm in table 1. The detailed explanation of our algorithm is given next.

4.2.2 The Features

There are two significant features of the use of statistical predictions to explore clusters by HOV³: *Enhanced separation of data groups* and *quantitatively guided exploration*. The projection of HOV³ is simply written as $G \leftarrow Hc(D, m)$ [23], where D is the processing dataset, m is a measure vector, and G is the distribution of D projected by HOV³.

- **Enhanced group separation**

It is proved that if there are several data point groups that can be roughly separated by applying a measure vector m in HOV³ to a dataset, then multiple applications of the projection in HOV³ with the same measure vector to the dataset would lead to the groups being more condensed, i.e., have a good separation of the groups [24].

This feature is achieved by step 6 and step17 in the **while** loop (steps 5-18) of the algorithm, as shown in Table 1. The enhanced group feature is significant for cluster exploration by HOV³ with statistical predictions, since clearly separated groups cannot be usually observed by applying a measure vector to a dataset in HOV³ once.

- **Quantitatively Guided Exploration**

The HOV³ technique provides a quantitative mechanism to visually detect cluster clue by measure vectors. Definitely, the statistical summaries of a dataset are quantitative depiction of the dataset. They objectively reflect the nature comparisons of the dataset. Thus introducing them as the predications in HOV³ avoid the randomness and subjectivity of the user during the cluster exploration by visualization.

To highlight these two features and demonstrate the effectiveness of our approach, we provide several examples with brief explanation of our algorithm in the next subsection.

4.3 The Examples

4.3.1 Parkinson's disease dataset

Parkinson's disease dataset has 23 attributes and 195 instances. The original data distribution of Parkinson's disease dataset is shown in Fig. 1, where we cannot recognize any groups of the dataset. Then we choose the standard deviation of the dataset $pstd=[0.24096, 0.18676, 0.25056, 0.15401, 0.13764, 0.14296, 0.14786, 0.14293, 0.17215, 0.16013, 0.19555, 0.16314, 0.12977, 0.19553, 0.12865, 0.17987, 0.43188, 0.24253, 0.22046, 0.13688, 0.18776, 0.17029, 0.18665]$ as a statistical prediction to explore the clusters of the dataset. Its projected data distribution is illustrated in Fig. 2, where data points are roughly separated, but we still cannot distinguish groups clearly (3, or 4 groups?) there.

According to the enhanced separation feature of HOV³ [24], we adopt two times inner product of $pstd$ as a statistical prediction to try again. The newly projected result is shown in Fig.3, where the data points are separated into two mains groups, based on the user's observation. We have also used three times mean value of Parkinson's

dataset as the statistical prediction to plot the dataset. Its data distribution is listed in Fig.4. It can be viewed that, clearly, there are two groups in Fig.3 and Fig.4.

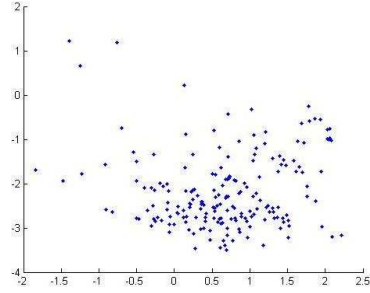


Fig.1 projecting data distribution by HOV^3 in MATLAB of Parkinson's disease dataset without any measurement

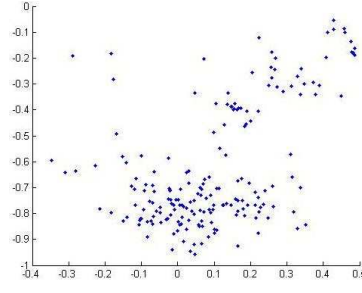


Fig.2 the data distribution projected by HOV^3 in MATLAB of Parkinson's disease dataset with its standard deviation, *pstd* as a statistical prediction

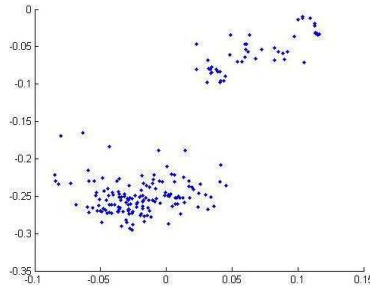


Fig.3 The data distribution projected by HOV^3 in MATLAB of the dataset in Fig.1 with two times of *pstd* as the prediction

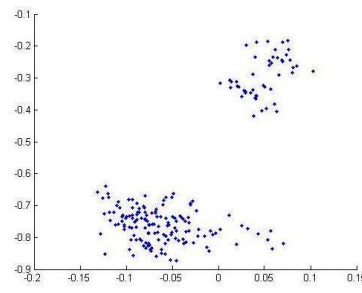


Fig.4 The data distribution projected by HOV^3 in MATLAB of Parkinson dataset with three times of mean values of the dataset as a prediction

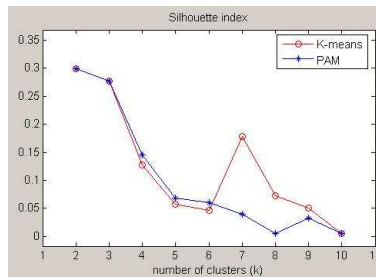


Fig.5 The quality indicated by Silhouette index of clustering results of Parkinson dataset produced by K-means and PAM clustering algorithms with cluster number 2 to 10.

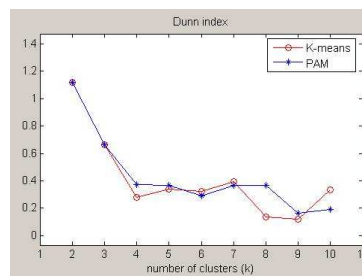


Fig.6 The quality indicated by Dunn index of clustering results of Parkinson dataset produced by K-means and PAM clustering algorithms with cluster number 2 to 10.

We believe that there are two well-separated clusters in Parkinson's dataset, based on above the experiments. The cluster exploration can be done iteratively until the user satisfies the grouping result by HOV³. He/she can terminate the cluster exploration process by his/her decision (steps 7-12) in the table 1.

To verify the validation of above the experiments produced by HOV³, we employed the CVAP system [21] to exam the quality of clustering results of Parkinson disease dataset by the K-means [15] and PAM [12] clustering algorithms with cluster number 2 to 10 respectively. Then we checked the quality of those clustering results by the cluster validation methods Shiouette index [17] and Dunn index [6]. The higher Shiouette and Dunn indices indicate the better quality of clustering result. The quality tests of those clustering results are illustrated in Fig.5 and Fig.6. It is clear that number 2 is the optimal cluster number of Parkinson dataset for K-means and PAM clustering. This example shows that statistics-guided cluster exploration by HOV³ provides an effective visual method to assist the user on acquisition of the cluster number in the preprocessing stage of clustering.

4.3.2 Wine dataset

We have also applied our approach to *wine* dataset, which has 13 attributes and 178 instances. Fig.7 and Fig.8 present the original data distribution of wine and the data distribution projected by HOV³ in MATLAB with three times standard deviation of dataset *wine* respectively. Clearly, there are three well-separated groups in Fig.8. Then we cluster these three groups (C_H).

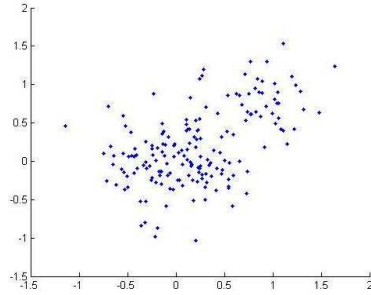


Fig.7 The original data distribution of *wine* dataset by HOV³ in MATLAB

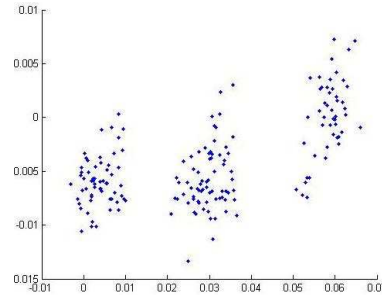


Fig.8 The data distribution projected by HOV³ in MATLAB of *wine* with its three times of stand deviation values

SOM (*Self-organizing Map*) is a neural network based clustering algorithm [13], which has been widely applied in machine learning and data mining. We applied the SOM to the *wine* dataset with cluster number 2-10, and employed the Shiouette index validation algorithm to verify the clustering results in CVAP. Fig.9 illustrates the curve of validation results produced by Shiouette index in CVAP, where we can observe that number 3 is the optimal cluster number of wine dataset.

The contrast of the clusters (C_H) projected by HOV³ and the clustering result (C_S) produced by the SOM clustering algorithm is summarized in Table 2. The weighted variance of the two clustering results is listed in the last row of the table. We can see that the quality of C_H is even slightly better than the quality of C_S based on the variance contrast. We believe that a domain expert could give a better and intuitive

explanation about this clustering result. This experiment also supports the effectiveness of our approach.

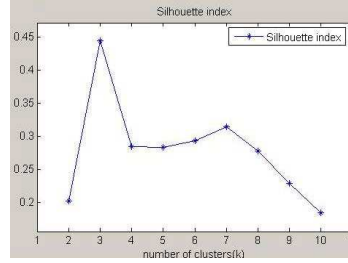


Fig.9 The quality indicated by Silhouette index of clustering results of wine dataset produced by SOM clustering algorithms with cluster number 2 to 10.

Table 2. The statistical contrast between the clusters ($k=3$) produced by HOV³ with three times standard deviation of wine dataset and the clusters produced by SOM clustering algorithms

C_H	%	Radius	Variance	Weighted Variance	C_S	%	Radius	Variance	Weighted Variance
1	26.966	102.286	0.125	3.37075	1	33.708	107.980	0.124	4.179792
2	39.888	97.221	0.182	7.259616	2	38.764	97.449	0.185	7.17134
3	33.146	108.289	0.124	4.110104	3	27.528	102.008	0.126	3.468528
14.74047					14.81966				

As the examples have demonstrated, visual projection based on the statistical prediction by HOV³ is a more purposeful and effective method for cluster exploration, and also it is easier to obtain a geometrical interpretation of the clustering results.

5 Conclusions

In this paper, we have proposed a statistics-guided visual approach to assist the user during cluster exploration, and demonstrated its effectiveness by experiments on several datasets. This approach adopts the statistical summaries of a high dimensional dataset as predictions to project the data by HOV³ so that the user can have an intuitive observation during cluster exploration. The use of statistical features of data avoids the weaknesses of randomness and arbitrary exploration of the existing visual methods employed in data mining. As a consequence, with the features of enhanced group separation and quantitatively guided exploration of this approach, the user can effectively identify the cluster number in the preprocessing stage of clustering.

References

1. Abul A. L., Alhaji R., Polat F., and Barker K.: Cluster Validity Analysis Using Subsampling, Proceedings of IEEE International Conference on Systems, Man, and Cybernetics, IEEE Press, Vol. (2) 1435-1440 (2003)
2. Ankerst M., Breunig M., Kriegel H.-P., Sander J.: OPTICS: Ordering Points To Identify the Clustering Structure, Proceedings of ACM SIGMOD'99, pp. 49-60 (1999)

3. Berkhin, P.: A Survey of Clustering Data Mining Techniques, Kogan, Jacob; Nicholas, Charles; Teboulle, Marc (Eds.) Grouping Multidimensional Data, Springer Press 25-72 (2006)
4. Chen K. and Liu L.: VISTA: Validating and Refining Clusters via Visualization. Journal of Information Visualization Vol. 13 (4) 257-270 (2004)
5. Dhillon I. S., Modha D. S., and Spangler W. S.: Visualizing class structure of multidimensional data," the 30th Symposium on the Interface: Computing Science and Statistics, Vol.(30), 488–493 (1998)
6. Dunn, J.C.: Well Separated Clusters and Optimal Fuzzy Partitions. Journal of Cybern. Vol. (4) 95–104 (1974)
7. Huang Z., Cheung D. W., Ng M. K.: An Empirical Study on the Visual Cluster Validation Method with Fastmap, In the proceedings of DASFAA'01 84-91 (2001)
8. Huang Z., Lin T.: A visual method of cluster validation with Fastmap, In the proceedings of PAKDD'2000, 153-164 (2000)
9. Jain A., Murty M. N., and Flynn P.J.: Data Clustering: A Review. ACM Computing Surveys Vol. 31(3) 264-323 (1999)
10. Jolliffe I. T.: *Principal Component Analysis*, Springer Press (2002)
11. Kandogan E.: Visualizing multi-dimensional clusters, trends, and outliers using star coordinates. In the proceedings of ACM SIGKDD'2001, 107-116 (2001)
12. Kaufman L. and Rousseeuw P. J. 1990, Finding Groups in Data, An Introduction to Cluster Analysis, John Wiley and Sons, Brussels, Belgium (1990).
13. Kohonen T. 1995. Self-Organizing Maps. Berlin/Heidelberg, Germany Springer, vol. 30. (1995)
14. Kruskal J. B., and Wish, M.: Multidimensional Scaling, SAGE university paper series on quantitative applications in the social sciences, Sage Publications, CA. 07-011 (1978)
15. MacQueen J. B.: Some Methods for classification and Analysis of Multivariate Observations, In the proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability, Berkeley, University of California Press, 281-297 (1967)
16. Pampalk E., Goebel W., and Widmer G.: Visualizing Changes in the Structure of Data for Exploratory Feature Selection. In the proceedings of SIGKDD'03, Washington, DC, USA (2003)
17. Rousseeuw P.: Silhouettes: a Graphical Aid to the Interpretation and Validation of Cluster Analysis, Journal of Computing. Vol. 20. 53-65 (1987)
18. Seo J., Shneiderman B.: From Integrated Publication and Information Systems to Virtual Information and Knowledge Environments, Essays Dedicated to Erich J. Neuhold on the Occasion of His 65th Birthday. Lecture Notes in Computer Science Vol. 3379 (2005).
19. Shneiderman B.: Inventing Discovery Tools: Combining Information Visualization with Data Mining. Discovery Science 2001, Proceedings. Lecture Notes in Computer Science Volume: 2226. 17-28 (2001)
20. Westphal C., and Blaxton T.: Data Mining Solutions: Methods and Tools for Solving Real-World Problems, John Wiley and Sons (1999).
21. Wang K., Wang B. and Peng L.: CVAP: Validation for Cluster Analyses. Data Science Journal, Vol.8 (20) 88-93 (2009).
22. Zhang K-B., Orgun M. A., Zhang K.: HOV³, An Approach for Cluster Analysis, Proc. of ADMA 2006, XiAn, China, Lecture Notes in Computer Science series, Vol. 4093 317-328 (2006)
23. Zhang K-B., Orgun M. A., Zhang K.: A Visual Approach for External Cluster Validation, Proc. of IEEE Symposium on Computational Intelligence and Data Mining (CIDM2007), Honolulu, Hawaii, USA, April 1-5, 2007, IEEE Press, 576-582 (2007)
24. Zhang K-B., Orgun M. A., Zhang K.: A Prediction-based Visual Approach for Cluster Exploration and Cluster Validation by HOV3, Proceedings of PKDD 2007, Warsaw, Poland, Lecture Notes in Computer, LNAI 4702 Springer Press, 336-349 (2007)