Efficient Clustering Of High Dimensional Data Partitioning Using Projected Clustering Algorithm

Mr.I.INFANT RAJ AP/CSE K.RAMAKRISHNAN COLLEGE OF ENGINEERING SAMAY APURAM, TRICHY rajsindhu098@gmail.com 9944299364 ABSTRACT

SHAFANA.K, RAJALAKSHMI.D II CSE B K.RAMAKRISHNAN COLLEGE OF ENGINEERING shafanakathar@gmail.com **8508275596**

Clustering techniques are used to partition the transaction data values. Similarity measures are used to analyze the relationship between the transactions. Vector based similarity models are suitable for low dimensional data values. High dimensional data values are clustered using subspace clustering methods.

Clustering high-dimensional data is a major challenge due to the curse of dimensionality. Projective clustering attempts to find projected clusters in subsets of the dimensions of a data space. Probability model describes projected clusters in high-dimensional data space. Model-based algorithm for fuzzy projective clustering that discovers clusters with overlapping boundaries in various projected subspaces. Model Based Projective Clustering (MPC) algorithm is used in the system.

The projective clustering techniques are used to cluster the high dimensional data. The model based projective clustering algorithm is a subspace clustering technique. Non-axis-subspaces are used with similarity analysis. Anomaly transactions are partitioned with projected clusters. The proposed system is designed to perform clustering on high dimensional spaces. Non access subspaces are included in the similarity analysis. Anomaly data values are verified with similarity under the clustering process. The subspace selection process is optimized.

1. Introduction

Data clustering has a wide range of applications and has been studied extensively in the statistics, data mining, and database communities. Many algorithms have been proposed in the area of clustering. One popular group of such algorithms, the model-based methods, have sparked wide interest because of their additional advantages, which give them the capacity to describe the underlying structures of populations in the data.

In model-based methods, data are thought of as originating from various possible sources, which are typically modeled by Gaussian mixture. The goal is to identify the generating mixture of Gaussians, that is, the nature of each Gaussian source, with its mean and covariance. Examples include the classical k-means and its variants. However, such methods would suffer from the curse of dimensionality problem for high dimensional data.

Many types of real-world data, such as the documents represented in the Vector Space Model (VSM) used in text mining and the microarray gene expression data of bioinformatics, consist of very high dimensional features. The data are inherently sparse in high-dimensional spaces, making the Gaussian function inappropriate in this case. Verleysen states that when the dimension increases, the percentage of the samples of a normalized multivariate Gaussian distribution falling around its center would rapidly decrease to 0. In other words, most of the volume of a Gaussian function is contained in the tails instead of near the center in high-dimensional space: the so called "empty space phenomenon".

Furthermore, in a high-dimensional space, clusters may exist in different subspaces comprised of different combinations of features. In many realworld applications, in fact, some points are correlated with a given set of dimensions, and others are correlated with different dimensions [11]. For example, in document clustering, clusters of documents on different topics are characterized by different subsets of keywords. The keywords for one cluster may not occur in the documents of other clusters. To address the above challenges, projective clustering has been defined to find clusters in different subspaces of the same data set.

A projected cluster is an ensemble of subsets of points, each of which is associated with a subset of attributes. Two different projected clusters are illustrated for a set of data points in 3-dimensional space. There are two clusters in this example; however, they are associated with two different lowdimensional subspaces. The first cluster corresponds to the data in group C_1 , which are close to each other when projected into the subspace consisting of the dimensions A_1 and A_2 , while the second one corresponds to the data in group C_2 projected onto the $A_1 - A_3$ plane.

A number of algorithms for finding such projected clusters have been proposed in the literature. They fall into two categories [2]. Those in the first category, which include PROCLUS, ORCLUS and FINDIT, are aimed at discovering the exact subspaces of different clusters. The algorithms in the second category cluster data points in the entire data space but assign different weighting values to different dimensions of clusters: examples include EWKM, FWKM and LAC. Most of the algorithms in the second category are of the k-means type, whose sequential structure is analogous to the mathematics of the EM algorithm. However, there is a general lack of underlying models on which these methods can be built.

In this paper, we will present a new modelbased method for projective clustering. The first contribution is the proposal of a probability model to describe projected clusters in a high-dimensional space. In contrast to existing models for highdimensional data clustering, our extended Gaussian model is designed for projective clustering, and by analysis is able to explain the general assumptions used in popular projective methods. Second, we derive an objective function for projective clustering based on the probability model and propose an EMtype, parameter-free algorithm, named MPC, for optimizing the objective function. The performance of MPC has been evaluated on synthetic data sets and some widely used real-world data sets, and the experimental results show its effectiveness. The method presented in this paper is very different from the one in our previous work [4]. Although the basic density function of the projected cluster is reused, the probability model for projected clusters has been changed. This results in a different algorithm which is no more dependent on any user-defined parameter for updating the dimension weights. The new algorithm has been much better motivated, analyzed, and experimentally evaluated.

2. Related Work

2.1 Techniques for High-Dimensional Clustering

Techniques for dimensionality reduction have been used in high-dimensional data clustering. Feature transformation techniques, such as PCA and SVD, attempt to summarize the data set in a smaller number of new dimensions created via linear combination of the original attributes, while feature selection methods select only the most relevant attributes for the clustering task. Because these traditional techniques are performed in the entire data space, they may encounter difficulties when clusters are found in different subspaces. Local Dimensionality Reduction (LDR) attempts to create a new set of dimensions for each cluster. The such method include the difficulties with determination of dimensionality for each subspace associated with the clusters. Additionally, LDR often has high computational complexity.

Biclustering also referred to as coclustering, has been proposed for simultaneous clustering on the

data points and dimensions of high-dimensional data. One of its typical applications is in the analysis of gene expression data, where the task is to find subgroups of genes and subgroups of conditions such that the genes exhibit highly correlated activities for every condition.

Finally, two related terms occur in the literature: subspace clustering and projective clustering. According to Parsons et al., projective clustering algorithms constitute a particular category of the subspace clustering techniques. However, different views are put forward elsewhere in the literature: see for instance [3]. We adopt the taxonomy and make a distinction between the two terms based on the ideas behind them. The idea of subspace clustering is to identify all dense regions in all subspaces, whereas in projective clustering the main focus is on discovering clusters that are projected onto particular spaces. In the subspace clustering field, CLIQUE was the pioneering approach, followed by a number of algorithms such as ENCLUS and MAFIA and SUBCLU. The major concern of this paper is projective clustering. In the following pages, we will focus only on such techniques.

2.2 Projective Clustering Methods

Projective clustering is typically based on feature weighting. Each dimension of each cluster is assigned a weighting value, indicating to what extent the dimension is relevant to the cluster. Usually, the weighting values of a given dimension may be different for different clusters. Based on the way the weights are determined, projective clustering algorithms can be divided into two categories: hard subspace clustering and soft subspace clustering.

In the first category, the dimensions are assigned weights with values of either 0 or 1, resulting in hard feature weighting for the subspaces. PROCLUS, which is based on the traditional kmedoids approach, is a representative algorithm using this weighting scheme. PROCLUS samples the data, then selects a set of medoids and iteratively improves the clustering, with the goal of minimizing the average within cluster dispersion. For each medoid, a set of dimensions is chosen whose average distances to the medoid are small compared to statistical expectation. After the subspaces have been identified, an average Manhattan segmental distance is used to assign points to medoids. PROCLUS requires users to provide the average number of relevant dimensions per cluster, which is usually unknown to users.

FINDIT, which uses a distance measure called the Dimension-Oriented Distance (DOD), is similar in structure to PROCLUS. As a hierarchical clustering algorithm, HARP automatically determines the relevant attributes of each cluster without requiring user-defined parameters. HARP is based on the assumption that two data points are likely to belong to the same cluster if they are very similar to each other along many dimensions. DOC also defines the subspace as a subset of attributes on which the projection of points in a partition is contained within a segment. DOC computes projected clusters using a randomized algorithm to minimize a certain quality function. MINECLUS improves on DOC by transforming the problem of finding the projected clusters into the problem of mining the frequent item set.

PROCLUS and the other algorithms mentioned above search for axis-aligned subspaces for the clusters, while some other methods search more general subspaces, termed nonaxis-aligned, where the new features are linear combinations of the original dimensions. ORCLUS is a generalization of PROCLUS that can discover clusters in arbitrarily oriented subspaces. Bv covariance matrix diagonalization, ORCLUS selects the eigenvectors corresponding to the smallest eigenvalues of the matrix of the set of points. ORCLUS inherits the weaknesses of PROCLUS mentioned above. KSM, a k-means type projective clustering algorithm. determines the non-axis-aligned subspaces by SVD computations, while EPCH performs non-axisaligned projective clustering by histogram construction.

Instead of identifying hard subspaces for clusters, the algorithms in the second category assign weights in the range [0, 1]. Since the weights can be any real number in [0, 1], we can call these soft projective clustering algorithms. Typically, the weight value for a dimension in a cluster is inversely proportional to the dispersion of the values from the center in the dimension of the cluster. In other words, a high weight indicates a small dispersion in a dimension of the cluster. Virtually all of the existing algorithms in this category are based on the following general assumptions: 1) the data project along a significant dimension onto a smaller range of values than on the other dimensions; 2) the data are more likely to be uniformly distributed along each irrelevant dimension. We will examine the capabilities of our projective clustering model, presented below, with respect to these two general assumptions.

A number of soft projective clustering algorithms have been reported recently. In [8], an algorithm making use of particle swarm optimization is presented. Since a heuristic global search strategy is used, the near-optimal feature weights could be obtained by this algorithm; however, it would run more slowly than other algorithms. To build an efficient soft projective clustering algorithm, the kmeans type structure has been widely adopted. Based on the classical k-means clustering process, an additional step for computing the weighting values is added in each iteration in these algorithms, which include EWKM, FWKM, LAC and FSC [5], etc. Algorithm 1 shows a typical structure for these algorithms.

Input: the dataset and the number of clusters K;

Output: the partition C and the associated weights W; Begin

Find the intial cluster V and set W with equal v values;

Report

1. Re-group the dataset into C according to V and W:

2. Re-compute V according to C;

3. Re-compute W according to C;

Until convergence is reached;

end

From Algorithm 1, the common projective clustering algorithm can be thought of as an EMbased process for estimating the unknown parameters C, V, and W of a model F(C, V, W) from which the data originate. However, the underlying F(C, V, W) is generally neglected in the above methods. The lack of such a model makes derivation of more effective clustering algorithms difficult. This has led us to work on projected cluster modeling, since we are convinced this type of modeling process allows us to benefit from the full potential of cluster analysis: for example, in describing the underlying mechanism that generates the cluster structure and addressing cluster validity problems.

In a typical model-based clustering analysis, one tries to find a mixture of multivariate distributions to approximate the data. Due to the empty space phenomenon and the property of projective clustering, as mentioned above, cluster modeling on high-dimensional data is a difficult problem. In one of the few attempts to use modelbased high-dimensional data clustering, Hoff [7] proposed a model of "clustering shifts in mean and variance" based on a nonparametric mixture of sequences of independent normal random variables. The model is learned by a Markov chain Monte Carlo process; however, its computational cost is prohibitive. Harpaz et al. [10] presented a nonparametric density estimation modeling technique, where the data are described as a mixture of linear manifolds. A Bayesian approach is used to identify groups of points that fit or are embedded in lower dimensional linear manifolds. The low dimensional subspaces associated with the individual clusters are computed by PCA. The problems with this method lie in its inflexibility in determining the dimensionality of the subspaces, and its inefficient clustering process.

3. Aprobability Model For Projective Clustering

The attributes of a non-axis-aligned subspace are typically combinations of the dimensions of the original data space. Since they are difficult to interpret, often making the clustering results less useful for many real applications, such as document clustering, only projected clusters in axisaligned subspaces are formalized in the following presentation.

3.1. Basic Notation and Definitions

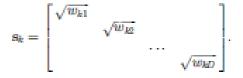
The notation used throughout the paper is summarized in Table 1. Given a data set $DB = \{x_1, x_2, \ldots, x_N\}$ containing K clusters, $x_i \in \mathbb{R}^D$ for $i = 1, 2, \ldots, N$ are called data points in the D-dimensional space. It is assumed that the data set has been normalized such that each $x_{ij} \ge [0, 1]$ where $j = 1, 2, \ldots$.,D. The membership degree of x_i with regard to the kth cluster c_k , where $k = 1, 2, \ldots, K$, is denoted as u_{ki} , subject to the following constraints:

$$o \le u_{ki} \le 1; \sum_{k=1}^{n} u_{ki} = 1, i = 1, 2, \dots, N.$$
(1)

The cluster c_k is associated with a weight vector $w_k = \langle w_{k1}, w_{k2}, \ldots, w_{kD} \rangle$, satisfying

$$\begin{cases} \sum_{j=1}^{D} w_{kj} = 1, \quad k = 1, 2, \dots, K \\ 0 \le w_{kj} \le 1, \quad k = 1, 2, \dots, K; j = 1, 2, \dots, D. \end{cases}$$
(2)

Here, the weight w_{kj} is defined to measure the relevance of the jth dimension to c_k . The greater the relevance, the larger the weight. Furthermore, we introduce a $D \times D$ matrix s_k , which is defined as



For a given c_k , the assignment of w_{k1} , w_{k2} , . . , w_{kD} can be regarded as a soft feature selection procedure for the space in which ck exists [6]. We thus use such a matrix to stand for the subspace associated with a cluster.

3.2 Probability Model

It is important to note that the Gaussian mixture is a fundamental hypothesis that many model-based clustering algorithms make regarding the data distribution model [9]. In this case, data points are thought of as originating from various possible sources, and the data from each particular source is modeled by a Gaussian. However, Gaussian functions are not appropriate in high-dimensional space due to the curse of dimensionality.

1	5
$X_i = \langle x_{i1}, x_{i2}, \dots, x_{iD} \rangle$	Ith data point R ^D ,
	i=1,2,,N

$DB=\{x_1, x_2, \dots, x_N\}$	The data set
K	Number of clusters
c_1, c_2, \dots, c_k	K clusters of DB
u _{ki}	Membership degree of
	x_i in $c_k, k=1, 2,, K$
$U=\{u_{ki}\}k\times N$	Membership matrix,
	where k=1,2,,k
	and i=1,2,N
$v_k = < v_{k1}, v_{k2,,k2} > 0$	Cluster center vector
	of c _k
$V = \{ v_{kj} \} k \times D$	Cluster center matrix,
	where k=1,2,k
	and j=1,2,,D
$w_k = \langle w_{k1}, w_{k2,,k}, w_{kD} \rangle$	A weight vector
	associated with c_k
W={ w_{kj} }k×D	Weight matrix, where
	k=1,2,,k and
	j=1,2,,D

TABLE 1: Notation Used throughout the Paper

In order to learn the underlying structure of clusters in a high-dimensional space, we will examine the distribution on each dimension. Consider the projections of the data points of the cluster k onto the jth dimension. It is reasonable to describe the projections using a 1D Gaussian function. The probability density function is

$$G(y_j|\mu_{kj};\sigma_k) = \frac{1}{\sqrt{2\pi\sigma_k}} \exp\left(-\frac{1}{2\sigma_k^2}(y_j - \mu_{kj})^2\right),$$

where μ_{kj} and $\boldsymbol{\sigma}_k$ denote the mean and covariance of the Gaussian. The above expression thus becomes

$$G(x_j | v_{kj}, w_{kj}; \sigma_k) = \frac{1}{\sqrt{2\pi\sigma_k}} \exp\left(-\frac{w_{kj}}{2\sigma_k^2}(x_j - v_{kj})^2\right).$$
 (5)

The major difference between (5) and the standard Gaussian is the introduction of the weighting value wkj, indicating the contribution of the jth dimension to c_k . Curves of (5) with different values of w_{kj} and a fixed σ_k . As we can see, the smaller the weighting value, the more uniformly distributed the data points. With a large weighting value, the data points would distribute within a small range. Note that the characteristic of this extended Gaussian meets the general requirements of projective clustering.

The probability model is created based on the following two assumptions. First, it is assumed that the distribution of points on each of the dimensions spanning the subspace is independent of the others. Although this assumption may not be realistic in some applications, it is a common assumption in many qualitative models, which allows us to approximate a joint distribution of the set of uncorrelated variables by the product of their marginals. Second, it is assumed that variations of points are independent of each other. Because

$$\int G(x_j|v_{kj}, w_{kj}; \sigma_k) dx_j = \frac{1}{\sqrt{w_{kj}}},$$

we then suppose the N inputs x_1, x_2, \ldots, x_N are independently and identically distributed from the following mixture density population:

$$F(\mathbf{x}; \Theta) = \sum_{k=1}^{K} \alpha_k \prod_{j=1}^{D} \sqrt{w_{kj}} G(x_j | v_{kj}, w_{kj}; \sigma_k)$$

with

$$\sum_{k=1}^{K} \alpha_k = 1, \ \alpha_k \ge 0, \ k = 1, 2, \dots, K,$$
(6)

where $\boldsymbol{\theta} = \{(\boldsymbol{\propto}_k, v_k, w_k, \boldsymbol{\sigma}_k) | 1 \le k \le K) \text{ is }$

the set of parameters, and \propto_k denotes the mixing weight of the kth component of the model.

3.3 Clustering Criterion

By applying the probability model to clustering, the goal is to estimate $\boldsymbol{\theta}$ from the given data set. Supposing $\hat{\boldsymbol{\theta}} = \{(\boldsymbol{\infty}_k, \mathbf{v}_k, \mathbf{w}_k, \boldsymbol{\sigma}_k) | 1 \le k \le K)$ is an estimator of $\boldsymbol{\theta}$, the distance between $F(\mathbf{x}, \boldsymbol{\theta})$ and $F(\mathbf{x}, \boldsymbol{\theta})$ can be measured by the following Kullback-Leibler divergence function:

$$R(\hat{\Theta}) = \int F(\mathbf{x}; \Theta) \ln \frac{F(\mathbf{x}; \Theta)}{\hat{F}(\mathbf{x}; \hat{\Theta})} d\mathbf{x}.$$

The equation can be decomposed into two terms. The first, $\int F(x; \theta) \ln F(x; \theta) dx$, is a constant that is irrelevant to $\hat{\theta}$; therefore, the following objective criterion needs to be maximized:

$$\begin{split} \uparrow Q_1(\hat{\Theta}) &= \int F(\mathbf{x};\Theta) \ln \hat{F}(\mathbf{x};\hat{\Theta}) \mathrm{d}\mathbf{x} \\ &= \sum_{k=1}^{K} \int p(k|\mathbf{x}) F(\mathbf{x};\Theta) \ln \hat{F}(\mathbf{x};\hat{\Theta}) \mathrm{d}\mathbf{x} \end{split}$$

With

$$p(k|\mathbf{x}) = \frac{\hat{\alpha}_k \prod_{j=1}^{D} \sqrt{\hat{w}_{kj}} G(x_j | \hat{v}_{kj}, \hat{w}_{kj}; \hat{\sigma}_k)}{\hat{F}(\mathbf{x}; \hat{\Theta})}, \quad 1 \le k \le K \quad (7)$$

where p(k|x) is the posterior probability of an input x from the kth probability density function, given x. Substituting for $\hat{F}(x; \hat{\theta})$ according to (7) in $Q_1(\hat{\theta})$, we obtain

$$\uparrow Q_1(\hat{\Theta}) = \sum_{k=1}^{K} \int p(k|\mathbf{x}) F(\mathbf{x};\Theta) \\ \times \ln \frac{\hat{\alpha}_k \prod_{j=1}^{D} \sqrt{\hat{w}_{kj}} G(x_j | \hat{v}_{kj}, \hat{w}_{kj}; \hat{\sigma}_k)}{p(k|\mathbf{x})} d\mathbf{x}.$$
(8)

By the law of large numbers, given a data set DB, maximizing (8) is equivalent to the maximum likelihood learning of $\boldsymbol{\theta}$ from all the inputs x_1, x_2, \ldots ., x_N . Therefore, using (5) to replace $G(x_j|\hat{v}_{kj}, \hat{w}_{kj}, \hat{\theta}_k)$, the objective criterion can be further transformed into

$$Q_{2}(\hat{\Theta}) = \frac{1}{N} \sum_{k=1}^{K} \sum_{i=1}^{N} p(k|\mathbf{x}_{i}) \\ \times \left(\frac{1}{2} \sum_{j=1}^{D} \left(\frac{\hat{u}_{kj}}{\hat{\sigma}_{k}^{2}} (x_{ij} - \hat{v}_{kj})^{2} - \ln \frac{\hat{w}_{kj}}{2\pi \hat{\sigma}_{k}^{2}} \right) - \ln \frac{\hat{\alpha}_{k}}{p(k|\mathbf{x}_{i})} \right).$$
(9)

For an input x_i, the posterior probability p(k|x_i) is thought of as the fuzzy membership u_{ki} in clustering. Given that $\frac{1}{N}$ and $\sum_{k=1}^{K} \sum_{i=1}^{N-1} \frac{1}{2} u_{ki} (\sum_{j=1}^{D} \ln 2\pi)$ are constants irrelevant to $\hat{\theta}$, the resulting clustering objective

irrelevant to θ , the resulting clustering objective function can be obtained as

$$\downarrow J(U, V, W, Z) = \sum_{k=1}^{K} \sum_{i=1}^{N} \left(\frac{u_{ki}}{2} \sum_{j=1}^{D} \left(\frac{w_{kj}}{\sigma_k^2} (x_{ij} - v_{kj})^2 - \ln \frac{w_{kj}}{\sigma_k^2} \right) - u_{ki} \ln \frac{\alpha_k}{u_{ki}} \right)$$
(10)

subject to the constraints of (1), (2), and (6). Here, Z = $\{\alpha_1, \alpha_2, \dots, \alpha_K, \sigma_1, \sigma_2, \dots, \sigma_K\}$.

4. Amodel-Based Algorithm For Projective Clustering

This section presents our algorithm, MPC for projective clustering by minimizing (10) subject to the constraints of (1), (2), and (6), which is a constrained nonlinear optimization problem. Using the Lagrangian multiplier technique, this can be transformed into an unconstrained optimization problem

$$\min J_1(U, V, W, Z) = J(U, V, W, Z)$$

$$+\sum_{k=1}^{K} \lambda_{k} \left(\sum_{j=1}^{D} w_{kj} - 1 \right) + \xi \left(\sum_{k=1}^{K} \alpha_{k} - 1 \right) \\ + \sum_{i=1}^{N} \zeta_{i} \left(\sum_{k=1}^{K} u_{ki} - 1 \right),$$
(11)

where $\lambda_k (k = 1, 2, ..., K)$, ξ , and $\zeta_i (i = 1, 2, ..., N)$ are the Lagrange multipliers corresponding to the constraints defined in (1), (2), and (6).

4.1. The Optimization Method

To achieve a local minimum of the objective function, the usual method is to use the partial optimization for each parameter in the function. Following this method, minimization of J_1 in (11) can be performed by optimizing U, V, W and Z in a sequential structure analogous to the mathematics of the EM algorithm. In each iteration, we first fix V = \hat{v} , W = \hat{w} , and Z = \hat{Z} , and solve U as \hat{U} to minimize $J_1(U, \hat{v}, \hat{w}, \hat{Z})$. Next, we fix $U = \hat{U}$, W = \hat{w} , and Z = \hat{Z} and solve V as \hat{v} to minimize $J_1(V, V, \hat{w}, \hat{Z})$. Then, $U = \hat{U}, V = \hat{v}$, and W = \hat{w} are fixed and the optimal Z, say \hat{Z} , is solved to minimize J1(\hat{U} , \hat{v} , \hat{w} , Z). Afterward, we fix U = \hat{U} , V = \hat{v} , and Z = \hat{Z} to obtain \hat{w} by minimizing $J1(\hat{U}, \hat{v}, W, \hat{Z})$. The four partial optimization problems are solved according to the following theorems.

4.2 The MPC Algorithm

The MPC algorithm, as outlined by Algorithm 2, performs projective clustering by minimizing the objective function of (10). Actually, this solution can also be regarded as an extension to the classical FCM algorithm by adding an additional step in each iteration to compute W for each cluster, an approach which is commonly adopted in existing soft subspace clustering algorithms such as [1].

Input: DB, K and a termination criterion which is a small positive number \in ;

Output: U, V and the associated weights W;

begin

Let *p* be the number of iteration, p=0

1. Initialization

1.1 Randomly choose K cluster centers. Denote V as $V^{(0)}$;

1.2 Set all the weights of W to
$$\frac{1}{D}$$
, and denote W

as W⁽⁰⁾;

1.3 Set all the
$$\alpha_k s$$
 to $\frac{1}{K}$ and $\sigma_k s$ to a non-

zero constant and denote them by $Z^{(0)}$;

2.1 Let $\hat{V} = V^{(p)}$, $\hat{W} = W^{(p)}$ and $\hat{Z} = Z^{(p)}$, compute $U^{(p+1)}$:

2.2 Let
$$\hat{U} = U^{(p+1)}$$
;

2.3 Let $\hat{V} = V^{(p+1)}$ s to compute $\hat{\alpha}k$ and $\hat{\sigma}k$ for k=1,2,...,K, respectively and obtain $Z^{(p+1)}$;

2.4 Let $\hat{Z} = Z^{(p+1)}$, to determine $\hat{\lambda}_k$ for k=1, 2,...., K; 2.5 Compute $W^{(p+1)}$; 2.6 p=p+1. until $|J(U^{(p)}, V^{(p)}, W^{(p)}, Z^{(p)}) - J(U^{(p-1)}, V^{(p-1)}, Z^{(p-1)})| < \epsilon$; 3. Output $U^{(p)}$ as U, $V^{(p)}$ as V and $W^{(p)}$ as W. end

It is important to note that MPC does not require user defined parameters for feature weighting, whereas most of the existing projective clustering algorithms do: for instance, 1 in PROCLUS, β in FWKM, γ in EWKM, etc. The only pending coefficient, say $\hat{\lambda}_k$, in the weight updating formula MPC can be determined by numerically solving. Step 2.4 of Algorithm 2 is designed for this purpose. All the variables except $\hat{\lambda}_k$ are given and thus can be considered as constants with respect to $\hat{\lambda}_k$. Consequently, we can resolve $\hat{\lambda}_k$ using a numerical method, such as the Newton-Raphson and bisection method.

5. Projected Clustering with Outlier Analysis

The proposed system is designed to perform clustering on high dimensional spaces. Non access subspaces are included in the similarity analysis. Anomaly data values are verified with similarity under the clustering process. The subspace selection process is optimized. The system is designed to perform data clustering on high dimensional data values. The model based projective clustering is improved with anomaly analysis. The system also enhanced with attribute alignment process. The system is divided into six major modules. They are data cleaning process, subspace selection, subspace alignment, clustering with MPC, MPC with outliers and clustering with attribute and anomaly analysis.

The data cleaning module is designed to correct noise transactions. The sub space selection module is designed to select attribute subsets. The attribute alignment is performed under subspace alignment module. The clustering is performed with model based projective clustering technique. The outlier analysis is integrated with MPC model. The attribute and anomaly analysis is applied in the enhanced MPC model.

6. Conclusion

The projective clustering techniques are used to cluster the high dimensional data. The model based projective clustering algorithm is a subspace clustering technique. Non-axis-subspaces are used with similarity analysis. Anomaly transactions are partitioned with projected clusters. Cluster accuracy is improved in the system. Features space selection is optimized to handle non aligned attribute subspace. Outlier analysis is provided in clustering process. Cluster initialization is improved with subspace selection process.

REFERENCES

[1] Q. Wang, Y. Ye, and J.Z. Huang, "Fuzzy k-Means with Variable Weighting in High Dimensional Data Analysis," Proc. Ninth Int'l Conf. Web-Age Information Management (WAIM), pp. 365-372, 2008.

[2] L. Jing, M.K. Ng, and J.Z. Huang, "An Entropy Weighting k- Means Algorithm for Subspace Clustering of High-Dimensional Sparse Data," IEEE Trans. Knowledge and Data Eng., vol. 19, no. 8, pp. 1026-1041, Aug. 2007.

[3] G. Moise, J. Sander, and M. Ester, "Robust Projected Clustering," Knowledge Information System, vol. 14, no. 3, pp. 273-298, 2008.

[4] L. Chen, Q. Jiang, and S. Wang, "A Probability Model for Projective Clustering on High Dimensional Data," Proc. IEEE Int'l Conf. Data Mining (ICDM), pp. 755-760, 2008.

[5] G. Gao, J. Wu, and Z. Yang, "A Fuzzy Subspace Clustering Algorithm for Clustering High Dimensional Data," Proc. Int'l Conf. Advanced Data Mining and Applications (ADMA), pp. 271-278, 2006.

[6] C. Domeniconi et al., "Locally Adaptive Metrics for Clustering High Dimensional Data," Data Mining and Knowledge Discovery, vol. 14, pp. 63-97, 2007.

[7] P.D. Hoff, "Model-Based Subspace Clustering," Bayesian Analysis, vol. 1, no. 2, pp. 321-344, 2006.

[8] Y. Lu, S. Wang, S. Li, and C. Zhou, "Particle Swarm Optimizer for Variable Weighting in Clustering High-Dimensional Data," Proc. IEEE Swarm Intelligence Symp., pp. 37-44, 2009.

[9] M. Bouguessa, S. Wang, and H. Sun, "An Objective Approach to Cluster Validation," Pattern Recognition Letters, vol. 27, pp. 1419-1430, 2006.

[10] R. Harpaz and R. Haralick, "Linear Manifold Clustering in High Dimensional Spaces by Stochastic Search," Pattern Recognition Letters, vol. 40, pp. 2672-2684, 2007.

[11] Lifei Chen, Qingshan Jiang and Shengrui Wang," Model-Based Method for Projective Clustering "IEEE Transactions On Knowledge And Data Engineering, Vol. 24, No. 7, July 2012.