

Distance Dimension Reduction on QR Factorization for Efficient Clustering Semantic XML Document Using the QR Fuzzy C-Mean (QR-FCM)

Hsu-Kuang Chang^{1,2} and I-Chang Jou¹

¹ Institute of Engineering Science and Technology,
National Kaohsiung First University of Science and Technology, Kaohsiung, Taiwan

² Department of Information Engineering, I-Shou University, Kaohsiung, Taiwan
hkchang@isu.edu.tw, icjou@ccms.nkfust.edu.tw

Abstract. The rapid growth of XML adoption has urged for the need of a proper representation for semi-structured documents, where the document semantic structural information has to be taken into account so as to support more precise document analysis. In order to analyze the information represented in XML documents efficiently, researches on XML document clustering are actively in progress. The key issue is how to devise the similarity measure between XML documents to be used for clustering. Since XML documents have hierarchical structure, it is not appropriate to cluster them by using a general document similarity measure. Dimension reduction plays an important role in handling the massive quantity of high dimensional data such as mass semantic structural documents. In this paper, we introduce distance dimension reduction (DDR) based on the QR factorization (DDR/QR) or the Cholesky factorization (DDR/C). DDR generates lower dimensional representations of the high-dimensional XML document, which can exactly preserve Euclidean distances and cosine similarities between any pair of XML documents in the original dimensional space. After projecting XML documents to the lower dimensional space obtained from DDR, our proposed method QR fuzzy c-mean to execute the document-analysis clustering algorithms (we called the QR-FCM). DDR can substantially reduce the computing time and/or memory requirement of a given document-analysis clustering algorithm, especially when we need to run the document analysis algorithm many times for estimating parameters or searching for a better solution.

Keywords: QR factorization, singular value decomposition, distance dimension reduction, PEWF, PEIDF, PESSW, fuzzy C-mean, QR-FCM.

1 Introduction

XML document which is a semi-structured data, have hierarchical structure. Therefore, rather than using the similarity measure of the general document clustering techniques as is, a new similarity measure that considers the semantic and structural information of XML document must be investigated. However, some XML clustering methods used

the similarity measure that takes only the structural information of XML documents into account. Hwang proposed a clustering method that extracts a typical structure of the maximum frequent pattern using *PrefixSpan* algorithm [1] on XML documents [2, 3]. However, since such a typical structure extracted from XML document is not only the structure that represents the XML document itself, but also it cannot be the representative of the whole documents corpus, there is an accuracy issue of similarity. Lian summarized XML documents into *S-graph* which is a structural graph and proposed the calculation method of distance between *S-graphs* to be used for clustering [4]. However, they have no consideration for semantic information on XML documents as they focus on structural information only. Since dimension reduction is one of the fundamental methods for data analysis, there have been a lot of studies on effective and efficient dimension reduction algorithms. There are linear dimension reduction algorithms including principal component analysis (PCA) [5] and multidimensional scaling (MDS) [6]. There are also nonlinear dimensional reduction algorithms (NLDR) including Isomap [7], locally linear embedding (LLE) [8], [9], Hessian LLE [10], Laplacian eigenmaps [11], local tangent space alignment (LTSA) [12] and distance preserving dimension reduction based on the singular value decomposition (DPDR/SVD) [13]. These dimension a variety of areas such as biomedical image recognition, biomedical text data mining, and biological data analysis.

In this paper, we introduce distance dimension reduction (DDR) based on the QR factorization or the Cholesky factorization. DDR can produce t dimensional representations where t is the rank of the original XML data set, which exactly preserve Euclidean L_2 -norm distances as well as cosine similarities between any pair of XML documents in the original m -dimensional space when m is much larger than the number of XML documents n , i.e. $m > n$. It projects the original XML data set into a much lower dimensional space without any loss of the pair-wise Euclidean distance information. Then, other XML documents analysis algorithms can be executed so that we can substantially reduce their computing times and memory requirements without any quality loss of their results.

The rest of this paper is organized as follows. In Section 2, we introduce the preparation XML documents on vector space model. In Section 3, we introduce DDR based on the QR factorization and DDR based on the Cholesky factorization after showing theorems of distance and cosine similarity preserving properties. Section 4 presents experimental results illustrating properties of the proposed DDR methods. Summary is given in Section 5.

2 Preparation for Semantic-Based XML Documents

In this section, we first introduce pre-processing steps for the incorporation of hierarchical information in encoding the XML tree's paths. It is based on the preorder tree representation (PTR) [14] and will be introduced after a brief review of how to generate an XML tree from an XML document. To do so, we have to first go through

Table 1. Preprocessing steps for XML document

Step 1. <u>Conversion</u>	Convert the XML document to tree Format. The values of the elements in the tree are not considered here, and only the structural information will be passed on to the subsequent steps.
Step 2. <u>Path Extraction</u>	Traverse the elements from the root to each leaf node of the tree. Record the sequence and hierarchical information for each path.
Step 3. <u>Similar Element Identification and Transformation</u>	Rename the terms of the paths with unique identifiers. Replace every element with a unique identifier in increasing order.
Step 4. <u>Nested and Duplicated Path Removal</u>	Remove any nested and duplicated path. The nested and duplicated paths in the tree are not considered here, and only the unique ones will be passed on to the next step.
Step 5. <u>Path Elements Encoding</u>	Based on the structural summary by Step 4, the path elements are encoded.

the following five preprocessing steps for XML documents. They are illustrated in Table 1.

From five steps preprocess, now XML document is modeled as a XML tree $T=(V,E)$. T is connected tree with $V=\{ v_1, v_2, \dots \}$ as a set of vertices and $v_1 \in V, v_2 \in V, (v_1, v_2) \in E$ as a set of edges. One distinguished vertex $r \in V$ is designated the root, and for all $v \in V$, there is a unique path from r to v . As an example, Figure 1 depicts a sample XML tree containing some information about collection of books. The *book* consists of *intro* tags, each comprising *title*, *author* and *date* tags. Each *author* contains *fname* and *lname*, each *date* includes *year* and *month* tags. Figure 1 left shows only the first letter of each tag for the simplicity.

XML document has a hierarchical structure and this structure is organized with tag paths. Each tag path represents document characteristics that can predict the contents of XML document. Strictly speaking, it shows the semantic structural characteristics of XML document. In this paper, we propose a method for calculating the similarity using all tag paths of XML tree representing the semantic structural information of XML document.

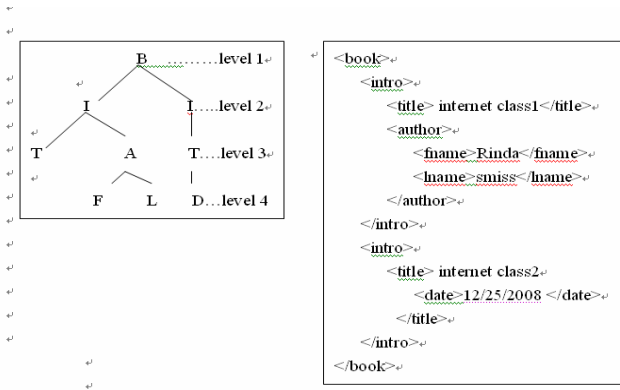


Fig. 1. An example of XML Document

3 Path Element Vector Space Model (PEVSM)

Vector model represents a document as a vector whose elements are the weights of the path elements within a document. In calculating the weight of each path element within a document, Term Frequency and IDF (Inverse Document Frequency) method is used [15]. We define PESSW (Path Element Structural Semantic Weight) that calculate the weight of a path element in a XML document. The PESSW is PEWF (Path Element Weighted Frequency) multiplied by PEIDF (Path Element Inverse Document Frequency). PESSW_{ij} of ith path element in the jth document is shown in equation (1).

$$PESSW_{ij} = PEWF_{ij} \times PEIDF_{ij} \quad (1)$$

PEWF_{ij} is shown in equation(2).

$$PEWF_{ij} = freq_{ij} \times \frac{1}{x^n} \quad (2)$$

PEIDF_{ij} is shown in equation (3),

$$PEIDF_{ij} = \log \frac{N}{DF_j} \quad (3)$$

Table 2 shows PEWF, PEIDF, and PESSW on an example trees in Figure 1.

Table 2. An example of PTWF, PTIDF and PESSW

Path	PEWF			PEIDF			PESSW		
	doc ₁	doc ₂	doc ₃	doc ₁	doc ₂	doc ₃	doc ₁	doc ₂	doc ₃
Element									
/B/LT/D	1.0	0.0	0.0	1.1	0.0	0.0	1.1	0.0	0.0
/B/LA/F	1.0	1.0	0.0	0.41	0.41	0.0	0.41	0.41	0.0
/B/LA/L	1.0	1.0	0.0	0.41	0.41	0.0	0.41	0.41	0.0
/M/LA/L	0.0	0.0	1.0	0.0	0.0	1.1	0.0	0.0	1.1
/B/LT/	2.0	1.0	0.0	0.41	0.41	0.0	0.81	0.41	0.0
/B/LA	1.0	1.0	0.0	0.41	0.41	0.0	0.41	0.41	0.0
/B/L/	2.0	1.0	0.0	0.41	0.41	0.0	0.81	0.41	0.0
/B	1.0	1.0	0.0	0.41	0.41	0.0	0.41	0.41	0.0
/M/LT	0.0	0.0	1.0	0.0	0.0	1.1	0.0	0.0	1.1
/M	0.0	0.0	1.0	0.0	0.0	1.1	0.0	0.0	1.1
/M/L	0.0	0.0	2.0	0.0	0.0	1.1	0.0	0.0	2.2
/T/D	0.5	0.0	0.0	1.10	0.0	0.0	0.5	0.0	0.0
/LA/F	0.5	0.5	0.0	0.41	0.41	0.0	0.21	0.21	0.0
/LA/L	0.5	0.5	0.5	0.0	0.0	0.0	0.0	0.0	0.0
/L/T	1.0	0.5	0.5	0.0	0.0	0.0	0.0	0.0	0.0
/LA	0.5	0.5	0.5	0.0	0.0	0.0	0.0	0.0	0.0
/L	1.0	0.5	1.0	0.0	0.0	0.0	0.0	0.0	0.0
/T/D	0.33	0.0	0.0	1.1	0.0	0.0	0.33	0.0	0.0
/LA/F	0.33	0.33	0.0	0.41	0.41	0.0	0.14	0.14	0.0
/LA/L	0.33	0.33	0.33	0.0	0.0	0.0	0.0	0.0	0.0
/T	0.67	0.33	0.33	0.0	0.0	0.0	0.0	0.0	0.0
/LA	0.33	0.33	0.33	0.0	0.0	0.0	0.0	0.0	0.0
D	0.25	0.0	0.0	1.1	0.0	0.0	0.27	0.0	0.0
F	0.25	0.25	0.0	0.41	0.41	0.0	0.1	0.1	0.0
L	0.25	0.25	0.25	0.0	0.0	0.0	0.0	0.0	0.0

Let d_x and d_y be two vectors that represent a XML document doc_x and doc_y . Cosine similarity is defined as the angle between two vectors and quantified by equation (4) and (5).

$$\cos \theta = \frac{d_x \cdot d_y}{|d_x| \cdot |d_y|}, \text{ that is} \tag{4}$$

$$\text{sim}(doc_x, doc_y) = \frac{d_x \cdot d_y}{|d_x| \cdot |d_y|} = \frac{\sum_{k=1}^t w_{kx} \times w_{ky}}{\sqrt{\sum_{k=1}^t w_{kx}^2} \times \sqrt{\sum_{k=1}^t w_{ky}^2}} \tag{5}$$

where $d_x = (w_{1x}, w_{2x}, \dots, w_{tx})$ and $d_y = (w_{1y}, w_{2y}, \dots, w_{ty})$, t is the total number of path elements in d_x, d_y respectively[16].

3.1 Distance Dimension Reduction (DDR) via the QR Factorization

Let us deal with n XML documents whose dimension is m s.t. $m \succ n$. We compute the QR factorization of the XML document matrix $D \in \mathfrak{R}^{m \times n}$:

$$D = QR = Q \begin{pmatrix} R_1 \\ 0 \end{pmatrix} = (Q_1 \quad Q_2) \begin{pmatrix} R_1 \\ 0 \end{pmatrix} = Q_1 R_1,$$

where $Q \in \mathfrak{R}^{m \times m}$ is an orthogonal matrix and $R_1 \in \mathfrak{R}^{n \times n}$ is an upper triangular matrix. Then, $Q_1 \in \mathfrak{R}^{m \times n}$ can be considered as a dimensionality transformation matrix when $m > n$ and the lower dimensional representation $\hat{x} \in \mathfrak{R}^{n \times 1}$ of a vector $x \in \mathfrak{R}^{m \times 1}$ can be computed as $\hat{x} = Q_1^T x$. Thus, the lower dimensional representation d_i of each XML document $\hat{d}_i = Q_1^T d_i = r_i$, where d_i is the i -th column of D and r_i is the i -th column of R_1 .

Theorem 1. QR-Factorization

Let $(\vec{d}_1, \vec{d}_2, \dots, \vec{d}_n)$ be a basis of D a subspace of \mathfrak{R}^m . For $j=2,3,\dots,n$, where we write $\vec{d}_j = \vec{d}_j^{\parallel} + \vec{d}_j^{\perp}$, with respect to span $(\vec{d}_1, \vec{d}_2, \dots, \vec{d}_{j-1})$.

Then, $\vec{q}_1 = \frac{1}{\|\vec{d}_1\|} \vec{d}_1$

$$\vec{d}_2^{\perp} = \vec{d}_2 - \vec{d}_2^{\parallel} = \vec{d}_2 - (\vec{q}_1 \cdot \vec{d}_2) \vec{q}_1, \text{ and } \vec{q}_2 = \frac{1}{\|\vec{d}_2^{\perp}\|} \vec{d}_2^{\perp},$$

$$\vec{d}_3^{\perp} = \vec{d}_3 - \vec{d}_3^{\parallel} = \vec{d}_3 - (\vec{q}_1 \cdot \vec{d}_3) \vec{q}_1 - (\vec{q}_2 \cdot \vec{d}_3) \vec{q}_2, \text{ and } \vec{q}_3 = \frac{1}{\|\vec{d}_3^{\perp}\|} \vec{d}_3^{\perp},$$

....., and

$$\vec{d}_j^{\perp} = \vec{d}_j - \vec{d}_j^{\parallel} = \vec{d}_j - (\vec{q}_1 \cdot \vec{d}_j) \vec{q}_1 - \dots - (\vec{q}_{j-1} \cdot \vec{d}_j) \vec{q}_{j-1}, \text{ and } \vec{q}_j = \frac{1}{\|\vec{d}_j^{\perp}\|} \vec{d}_j^{\perp}$$

....., and

$$\vec{d}_m^{\perp} = \vec{d}_m - \vec{d}_m^{\parallel} = \vec{d}_m - (\vec{q}_1 \cdot \vec{d}_m) \vec{q}_1 - \dots - (\vec{q}_{m-1} \cdot \vec{d}_m) \vec{q}_{m-1}, \text{ and } \vec{q}_m = \frac{1}{\|\vec{d}_m^{\perp}\|} \vec{d}_m^{\perp}$$

The Gram-Schmidt process represents a change of basis from the old basis $D = (\vec{d}_1, \vec{d}_2, \dots, \vec{d}_n)$ to a new, orthonormal basis $Q = (\vec{q}_1, \vec{q}_2, \dots, \vec{q}_n)$ of vector

V. If R is the change basis matrix form D to Q, then

$$D=QR$$

$$(\vec{d}_1, \vec{d}_2, \dots, \vec{d}_n) = (\vec{q}_1, \vec{q}_2, \dots, \vec{q}_n) \begin{bmatrix} r_{11} & r_{12} & \dots & \dots & r_{1n} \\ 0 & r_{22} & \dots & \dots & r_{2n} \\ 0 & 0 & r_{33} & \dots & r_{3n} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & r_{nn} \end{bmatrix}$$

This equation is called the QR-factorization of the matrix D. Using Gram-Schmidt algorithm, we know that

$$\begin{aligned} \vec{d}_j &= \vec{d}_j^{\parallel} + \vec{d}_j^{\perp} \\ &= (\vec{q}_1 \cdot \vec{d}_1)\vec{q}_1 + \dots + (\vec{q}_i \cdot \vec{d}_j)\vec{q}_i + \dots + (\vec{q}_{j-1} \cdot \vec{d}_j)\vec{q}_{j-1} + \|\vec{d}_j^{\perp}\|\vec{q}_j \\ &= r_{11}\vec{q}_1 + \dots + r_{ij} + \dots + r_{j-1,j}\vec{q}_{j-1} + r_{jj}\vec{q}_j \end{aligned}$$

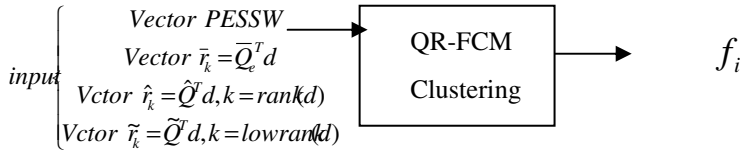
It follows that $r_{ij} = \vec{q}_i \cdot \vec{d}_j$ if $i < j$, $r_{jj} = \|\vec{d}_j^{\perp}\| \geq 0$; and $r_{ij}=0$ if $i > j$. The last equation

implies that R is an upper triangular matrix. The first diagonal entry is $r_{11} = \|\vec{d}_1\|$,

$$\text{since } \vec{d}_1 = \|\vec{d}_1\|\vec{q}_1 \cdot (\because \vec{q}_1 = \frac{1}{\|\vec{d}_1\|}\vec{d}_1)$$

3.2 Our Proposed QR-FCM Algorithm

From the QR factorization as the previous section described, we have the originated document vector PESSW, document vector \vec{r} (economic QR factoring), document vector \hat{r} (rank of PESSW), and document vector \tilde{r} (low rank of QR) based on the XML documents, which taken as the QR-FCM input data and then go through the clustering.



$$F(I)=[f_1, f_2, \dots, f_c], \text{ where } f_i = \sum_{j=1}^N \mu_{ij} P_j = \frac{1}{N} \sum_{j=1}^N \mu_{ij} .$$

Fuzzy c-means (FCM) is a method of clustering which allows one piece of data to belong to two or more clusters. This method (developed by Dunn in 1973 and improved by Bezdek in 1981) is frequently used in pattern recognition. It is based on minimization of the following objective function:

$$J_m = \sum_{i=1}^N \sum_{j=1}^C \mu_{ij}^m \|d_i - c_j\|^2, \quad 1 \leq m \leq \infty$$

where m is any real number greater than 1, u_{ij} is the degree of membership of d_i in the cluster j , d_i is the i th of d -dimensional measured data, c_j is the d -dimension center of the cluster, and $\|\cdot\|$ is any norm expressing the similarity between any measured data and the center.

Fuzzy partitioning is carried out through an iterative optimization of the objective function shown above, with the update of membership u_{ij} and the cluster centers c_j by:

$$\mu_{ij} = \frac{1}{\sum_{k=1}^C \left(\frac{\|d_i - c_j\|}{\|d_i - c_k\|} \right)^{\frac{2}{m-1}}}, \quad c_j = \frac{\sum_{i=1}^N \mu_{ij}^m . d_i}{\sum_{i=1}^N \mu_{ij}^m}$$

This iteration will stop when $\max_{ij} \{|\mu_{ij}^{(k+1)} - \mu_{ij}^{(k)}|\} < \xi$, where ξ is a termination criterion between 0 and 1, whereas k are the iteration steps. This procedure converges to a local minimum or a saddle point of J_m . The QR-FCM algorithm is composed of the following steps:

1. Input the number of clusters c , the weighting exponent m , and error tolerance ϵ .
2. Initialize the cluster centers $c = \{c_i\}$, for $1 \leq i \leq c$.
3. Initialize $U=[u_{ij}]$ matrix, $U^{(0)}$
4. At k -step: calculate the centers vectors $C^{(k)}=[c_i]$ with $U^{(k)}$

$$c_j = \frac{\sum_{i=1}^N \mu_{ij}^m \cdot d_i}{\sum_{i=1}^N \mu_{ij}^m}$$

5. Update $U^{(k)}$, $U^{(k+1)}$

$$\mu_{ij} = \frac{\left(\frac{1}{\|d_j, v_i\|_2} \right)^{\frac{1}{m-1}}}{\sum_{i=1}^c \left(\frac{1}{\|d_j, c_i\|_2} \right)^{\frac{1}{m-1}}}, \quad \mu_{ij} = \frac{1}{\sum_{k=1}^c \left(\frac{\|d_i - c_j\|}{\|d_i - c_k\|} \right)^{\frac{2}{m-1}}}$$

6. If $\|U^{(k+1)} - U^{(k)}\| < \xi$ then STOP; otherwise return to step 4.

4 Experiment Result

4.1 Within-Group-Variance and Between-Group-Variance Using Membership μ

The purpose of any clustering technique [17], [18], [19], [20], [21] is to evolve a $K \times n$ partition matrix $(U)=D$ of a data set D ($D = \{d_1, d_2, \dots, d_n\}$) in \mathfrak{R}^N , representing its partitioning into a number, say K , of clusters ($C_1; C_2; \dots; C_K$). The partition matrix $(U) = D$ may be represented as $U=[\mu_{kj}]$, $k=1, \dots, K$, and $j= 1, \dots, n$, where u_{kj} is the membership of pattern d_j to clusters C_k . In crisp partitioning of the data, the following condition holds: $u_{kj}=1$ if $d_j \in C_k$; otherwise, $u_{kj}=0$. Right now, we use the membership u as separation measuring degree to figure out the variance of the within-group-variance and between-group-variance among the XML documents in the different cluster. First, we define the useful definitions as follows.

Definition 1. The closeness (denseness) of the XML documents in the i^{th} cluster (C_i)

$$C_i = \frac{\sum_{j \in S_i} \mu_{ij}}{|S_i|} \text{ where } I_{i,j} = \begin{cases} 1 & \text{if } \mu_{i,j} = \max_{1 \leq k \leq C} \mu_{k,j} \\ 0 & \text{if } \mu_{i,j} \neq \max_{1 \leq k \leq C} \mu_{k,j} \end{cases}$$

, μ_{ij} is the membership calculated from the QR-FCM (QR fuzzy C-mean) and

$$S_i = \{j \mid I_{i,j} = 1\}.$$

Definition 2. Within-Group-Variance (WGV) of the cluster

$$WGV(\mu, V; D) = \sum_{i=1}^C \sum_{j \in S_i} (C_i)^{-1} \|d_j, v_i\|_2 \tag{6}$$

C_i is the denseness within the cluster which is defined as the membership μ_{ij} with the i th cluster. μ_{ij} is calculated from QR-FCM (QR fuzzy C-mean). μ_{ij} represents the membership of the XML document d_j within the cluster i . If the all membership within the i th cluster is large, then each XML document is close to the central document. That means each XML document is close to the central document within the cluster, each close to the central document, high denseness, low variance. C_i is the denseness determined by the μ_{ij} membership. Beside, from the distance of XML document and central document, $\|d_j, v_i\|_2$, defined the XML document diversity within the cluster. We defined the within-group-variance as combine both of these two criteria, the bigger C_i the more dense, and the smaller $\|d_j, v_i\|_2$ the better result for the within-group-variance, so we have $\min\{WGV\}$.

Definition 3. The contribution of the separation XML document d_j among λ cluster and

m cluster defined as $\mu_{\lambda,j} * \mu_{m,j}$.

Example 1

$$m = C_1 \begin{bmatrix} d_1 & d_2 & d_3 \\ \mu_{11} & \mu_{12} & \mu_{31} \end{bmatrix} = C_1 \begin{bmatrix} d_1 & d_2 & d_3 \\ 0.8 & 0.1 & 0.4 \end{bmatrix},$$

$$C_2 \begin{bmatrix} \mu_{12} & \mu_{22} & \mu_{32} \end{bmatrix} = C_2 \begin{bmatrix} 0.2 & 0.9 & 0.6 \end{bmatrix}$$

$$\mu_{11} * \mu_{21} = 0.16, \mu_{12} * \mu_{22} = 0.09, \mu_{13} * \mu_{23} = 0.2$$

From the membership matrix m , the smaller the $\mu_{\lambda,j} * \mu_{m,j}$, the bigger contribution of the separation for the λ cluster and m cluster.

Example 2

$$m = \begin{matrix} C_1 \\ C_2 \\ C_3 \end{matrix} \begin{bmatrix} d_1 & d_2 \\ \mu_{11} & \mu_{12} \\ \mu_{21} & \mu_{22} \\ \mu_{31} & \mu_{32} \end{bmatrix} = \begin{matrix} C_1 \\ C_2 \\ C_3 \end{matrix} \begin{bmatrix} d_1 & d_2 \\ 0.8 & 0.5 \\ 0.1 & 0.3 \\ 0.1 & 0.2 \end{bmatrix} \quad \mu_{11} * \mu_{21} = 0.08 \quad \mu_{12} * \mu_{22} = 0.15$$

Shows the separation contribution of the XML document d_1 with cluster 1 and cluster 2 (C_1 and C_2) and the separation contribution of the XML document d_2 with cluster 1 and cluster 2 (C_1 and C_2).

Definition 4. Between-Group-Variance (BGV) among two clusters is defined as the following

$$BGV(\mu, v) = \frac{1}{\binom{C}{2}} \sum_{\lambda=1}^{C-1} \sum_{m=\lambda+1}^C \left(\frac{\sum_{j \in S_\lambda \cup S_m} \mu_{\lambda,j} * \mu_{m,j}}{|S_\lambda| + |S_m|} \right)^{-1} \|v_\lambda, v_m\|_2 \tag{7}$$

where $S_\lambda = \{j | I_{\lambda,j} = 1\}$ and $S_m = \{j | I_{m,j} = 1\}$, v_λ, v_m are the vector of the central document on the λ^{th} and the m -th cluster separately.

From the Example 2 on the Definition 3, we know that $\frac{\sum_{j \in S_\lambda \cup S_m} \mu_{\lambda,j} * \mu_{m,j}}{|S_\lambda| + |S_m|}$ represents the separation contribution for all XML documents in the cluster λ and m . Combine the representing separation contribution using membership with between-group distance to define a between-group-variance (BGV), the smaller value of

$\frac{\sum_{j \in S_\lambda \cup S_m} \mu_{\lambda,j} * \mu_{m,j}}{|S_\lambda| + |S_m|}$, the more separation and the more separation contribution. The

maximized $\|v_\lambda, v_m\|_2$ is desired. So, $BGV(\mu, v)$ get larger get better. Finally, we

combine $BGV(\mu, v)$ with $WGV(\mu, v; D)$ to define a $WB(\mu, v; D)$ membership cluster validity indicator so called $WB(\mu, v; D)$ as follows.

$$\begin{aligned}
 WB(\mu, v; D) &= \frac{(B)etween-(G)roup-(V)ariance(\mu, v)}{(W)ithin-(G)roup-(V)ariance(\mu, v; D)} = \frac{BGV(\mu, v)}{WGV(\mu, v; D)} \\
 &= \frac{\frac{1}{\binom{C}{2}} \sum_{\lambda=1}^{C-1} \sum_{m=\lambda+1}^C \left(\frac{\sum_{j \in S_\lambda \cup S_m} \mu_{\lambda, j} * \mu_{m, j}}{|S_\lambda| + |S_m|} \right)^{-1} \|v_\lambda, v_m\|_2}{\sum_{i=1}^C \sum_{j \in S_i} (C_i)^{-1} \|d_j, v_i\|_2} \\
 &= \frac{\frac{1}{\binom{C}{2}} \sum_{\lambda=1}^{C-1} \sum_{m=\lambda+1}^C \frac{\|v_\lambda, v_m\|_2}{\frac{\sum_{j \in S_\lambda \cup S_m} \mu_{\lambda, j} * \mu_{m, j}}{|S_\lambda| + |S_m|}}}{\sum_{i=1}^C \sum_{j \in S_i} \frac{\|d_j, v_i\|_2}{C_i}} \tag{8}
 \end{aligned}$$

The maximized $WB(\mu, v; D)$ is better that means $\max_C WB$ is desired.

Definition 5. Precision, Recall, and F-Measure

Let X represent the set of XML documents and let $C = \{C_1, \dots, C_k\}$ be a clustering of X . Moreover, let $C^\Delta = \{C_1^\Delta, \dots, C_\ell^\Delta\}$ designate the human reference classification. Then the recall of cluster C_i with respect to class C_i^Δ , R_i , is defined as $|C_i \cap C_i^\Delta| / |C_i^\Delta|$. The precision of cluster C_i with respect to class C_i^Δ , PR_i , is defined as $|C_i \cap C_i^\Delta| / |C_i|$. The F -Measure combines the precision and recall measures from information retrieval [22]. The F -Measure combines both values as follows:

$$F - Measure = \frac{2}{\frac{1}{PR_i} + \frac{1}{R_i}} \tag{9}$$

and uses the formula to evaluate the QR-FCM clustering result on the following section.

4.2 Working on Real Data Sets

We have developed a prototype and performed extended evaluation of our framework for validity clustering XML documents. We tested the performance as well as the

quality of the clustering results using real data. The prototype testbed is a java-based software that can (a) generate synthetic XML documents or use existing ones, (b) extract feature (structural summaries) from XML documents, (c) norm distance calculate pair-wise structural distances between these summaries, (d) perform the QR-Fuzzy C-mean (QR-FCM). The goal of our work is to find documents with structural similarity, that is, documents generated from a common DTD. For any choice of distance metric, we can evaluate how closely the reported lower rank dimension document matrix corresponding to the actual originated XML. The experiments were conducted as follows. The following five DTDs were downloaded from ACM's SIGMOD Record homepage [23]: OrdinaryIssuePage.dtd (O in short), ProceedingsPage1999.dtd (P-1999 in short), ProceedingsPage2002.dtd (P-2002 in short), IndexTerm1999.dtd (IT in short), Ordinary2002.dtd (Ord-2002 in short) , and Ordinary2005.dtd (Ord-2005 in short). For another real data set we used the documents on ADC/NASA [24]: 70 XML documents from adml.dtd (Astronomical Dataset Markup Language DTD). Also we download the nigara data [24]: 150 XML documents from movie.dtd, department.dtd, club.dtd, and personnel.dtd. Based upon these sets of XML documents with similar characteristics, their accuracy of low rank dimension information retrieval were computed, analyzed and reported as follows. Table 1 first shows the results of QR fuzzy C-mean (QR-FCM) on the variant numbers of PESSW XML documents (50, 70, 90, and 120) from originated 3 DTDs, 4 DTDs and 5 DTDs which we called the heterogeneous XML documents, then Table 2 shows the results of QR fuzzy C-mean (QR-FCM) on the variant numbers of economic QR factorization $R = \overline{Q}^T d$ XML documents (50, 70, 90, and 120) from originated 3 DTDs, 4 DTDs and 5 DTDs , and Table 3 shows the results of QR fuzzy C-mean (QR-FCM) on the variant numbers of $\hat{r}_k = Q_k^T d$ (k=rank of PESSW) XML documents (50, 70, 90, and 120) from originated 3 DTDs, 4 DTDs and 5 DTDs and finally Table 4 shows the results of QR fuzzy C-mean (QR-FCM) on the variant numbers of $\tilde{r}_k = Q_k^T d$ (k=low rank of k) XML documents (50, 70, 90, and 120) from originated 3 DTDs, 4 DTDs and 5 DTDs. We also compute the value of F-measure on the clustering quality measure.

Table 3. QR Fuzzy C-mean (QR-FCM) on PESSW XML from variant DTDs

PESSW ^o	<u>3 Origin DTDs^o</u>				<u>4 Origin DTDs^o</u>				<u>5 Origin DTDs^o</u>				
	<u>Q</u>	<u>I</u>	<u>N</u>	<u>D</u>	<u>Q</u>	<u>I</u>	<u>N</u>	<u>M</u>	<u>Q</u>	<u>I</u>	<u>N</u>	<u>M</u>	<u>D</u>
# XML ^o	50 ^o	70 ^o	90 ^o	120 ^o	50 ^o	70 ^o	90 ^o	120 ^o	50 ^o	70 ^o	90 ^o	120 ^o	
Space ^o	125k ^o	397k ^o	484k ^o	619k ^o	144k ^o	427k ^o	564k ^o	713k ^o	142k ^o	497k ^o	518k ^o	689k ^o	
WB ^o	0.14 ^o	0.31 ^o	.21 ^o	.025 ^o	.09 ^o	.029 ^o	.021 ^o	.022 ^o	.08 ^o	.023 ^o	.02 ^o	.015 ^o	
PR ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	
R ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	
F-Measure ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	

(Q)rdinaryIssuePage.dtd, (I)ndex(T)ermPage.dtd, (N)asa.(M)ove.dtd, (D)ept.dtd

Table 4. QR Fuzzy C-mean (QR-FCM) on $\bar{r}_k = \bar{Q}_k^T d$ XML from variant DTDs

$\bar{r}_k = \bar{Q}_k^T d$ k=economic rank ^o	<u>3 Origin DTDs^o</u>				<u>4 Origin DTDs^o</u>				<u>5 Origin DTDs^o</u>				
	<u>Q</u>	<u>I</u>	<u>N</u>	<u>D</u>	<u>Q</u>	<u>I</u>	<u>N</u>	<u>M</u>	<u>Q</u>	<u>I</u>	<u>N</u>	<u>M</u>	<u>D</u>
# XML ^o	50 ^o	70 ^o	90 ^o	120 ^o	50 ^o	70 ^o	90 ^o	120 ^o	50 ^o	70 ^o	90 ^o	120 ^o	
Space ^o	23k ^o	45k ^o	75k ^o	93k ^o	23k ^o	44k ^o	76k ^o	166k ^o	23k ^o	133k ^o	75k ^o	133k ^o	
WB ^o	0.14 ^o	.031 ^o	.21 ^o	.025 ^o	.09 ^o	.029 ^o	.02 ^o	.022 ^o	.08 ^o	.023 ^o	.02 ^o	.015 ^o	
PR ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	
R ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	
F-Measure ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	

(Q)rdinaryIssuePage.dtd, (I)ndex(T)ermPage.dtd, (N)asa.(M)ove.dtd, (D)ept.dtd

Table 5. QR Fuzzy C-mean (QR-FCM) on $\hat{r}_k = \hat{Q}_k^T d$ XML from variant DTDs

$\hat{r}_k = \hat{Q}_k^T d$ k=rank(D) ^o	<u>3 Origin DTDs^o</u>				<u>4 Origin DTDs^o</u>				<u>5 Origin DTDs^o</u>				
	<u>Q</u>	<u>I</u>	<u>N</u>	<u>D</u>	<u>Q</u>	<u>I</u>	<u>N</u>	<u>M</u>	<u>Q</u>	<u>I</u>	<u>N</u>	<u>M</u>	<u>D</u>
# XML ^o	50 ^o	70 ^o	90 ^o	120 ^o	50 ^o	70 ^o	90 ^o	120 ^o	50 ^o	70 ^o	90 ^o	120 ^o	
Space ^o	4k ^o	18k ^o	23k ^o	35k ^o	6k ^o	19k ^o	28k ^o	48k ^o	7k ^o	40k ^o	28k ^o	40k ^o	
WB ^o	0.14 ^o	.031 ^o	.21 ^o	.025 ^o	.09 ^o	.029 ^o	.02 ^o	.022 ^o	.08 ^o	.023 ^o	.02 ^o	.015 ^o	
PR ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	
R ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	
F-Measure ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	1.0 ^o	

(Q)rdinaryIssuePage.dtd, (I)ndex(T)ermPage.dtd, (N)asa.(M)ove.dtd, (D)ept.dtd

Table 6. QR Fuzzy C-mean (QR-FCM) on $\tilde{r}_k = \tilde{Q}_k^T d$ XML from variant DTDs

$\tilde{r}_k = \tilde{Q}_k^T d$	3 Origin DTDs ^a				4 Origin DTDs ^a				5 Origin DTDs ^a				
	Q	IT	N		Q	IT	N	M	Q	IT	N	M	D
k= low rank													
# XML ^b	50 ^c	70 ^c	90 ^c	120 ^c	50 ^c	70 ^c	90 ^c	120 ^c	50 ^c	70 ^c	90 ^c	120 ^c	120 ^c
Space ^c	2k ^c	9k ^c	11k ^c	17k ^c	3k ^c	10k ^c	14k ^c	24k ^c	3k ^c	20k ^c	14k ^c	20k ^c	20k ^c
WB ^c	0.14 ^c	.031 ^c	.21 ^c	.025 ^c	.09 ^c	.029 ^c	.02 ^c	.022 ^c	.08 ^c	.023 ^c	.02 ^c	.015 ^c	
PR ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c
R ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c
F-Measure ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c	1.0 ^c

(Q)rdinaryIssuePage.dtd, (I)ndex(T)ermPage.dtd, (N)asa, (M)ove.dtd, (D)ept.dtd.

5 Conclusion and Future Work

The original XML documents $D_N=[d_1, d_2, \dots, d_N]$ are modeled on the vector space model according to the path element on each document, that is $D_N=PESSW$, then do QR factorization on the PE SSW we derived the $PESSW=QR$ ($D = \overline{Q}_k \overline{R}_k$), or $\overline{Q}_k^T D = \overline{R}_k$, and then take the rank on the $\hat{r}_k = \hat{Q}_k^T d$ with the rank of PE SSW , and finally $\tilde{r}_k = \tilde{Q}_k^T d$ with the low-rank of QR on PE SSW . We passed the 4 resulting vectors (PE SSW , \overline{R}_k , \hat{R}_k , and \tilde{R}_k) into the QR-FCM clustering algorithm to get the clustering result. From the clustering results on the section experiment, we found the same clustering result from the variant PE SSW , \overline{R}_k , \hat{R}_k and \tilde{R}_k vectors. We conclude that use the low-rank vector \tilde{R}_k instead of PE SSW original document, not only saving the space on the input vector but also spending less time to cluster on the documents. Based on the clustering, the next issues will be arisen such as how to index the clustered XML documents for speeding up query response, for example like R-tree and B⁺-tree, and how to manage the dynamic-updating XML document for adding XML document, adding a path element, and modifying a path element weight.

References

- [1] Pei, J., Han, J., Asi, B.M., Pinto, H.: PrefixSpan: Mining Sequential Pattern efficiently by Prefix-Projected Pattern Growth. In: Int. Conf. Data Engineering, ICDE (2001)
- [2] Hwang, J.H., Ryu, K.H.: XML A New XML clustering for Structural Retrieval. In: International Conference on Conceptual Modeling (2004)
- [3] Hwang, J.H., Ryu, K.h.: Clustering and retrieval of XML documents by structure. In: Gervasi, O., Gavrilova, M.L., Kumar, V., Laganá, A., Lee, H.P., Mun, Y., Taniar, D., Tan, C.J.K. (eds.) ICCSA 2005. LNCS, vol. 3481, pp. 925–935. Springer, Heidelberg (2005)

- [4] Lian, W., Wai-lok, D.: An Efficient and Scalable Algorithm for Clustering XML Documents by Structure. IEEE Computer Society Technical Committee on Data Engineering (2004)
- [5] Massay, W.F.: Principal components regression in exploratory statistical research. *J. Amer Statist. Assoc.* 60, 234–246 (1965)
- [6] Torgerson, W.S.: *Theory & Methods of Scaling*. Wiley, New York (1958)
- [7] Tenenbaum, J.B., de Silva, V., Langford, J.C.: A global geometric framework for nonlinear dimensionality reduction. *Science* 290(5500), 2319–2323 (2000)
- [8] Roweis, S.T., Saul, L.K.: Nonlinear dimensionality reduction by locally linear embedding. *Science* 290, 2323–2326 (2000)
- [9] Saul, L.K., Roweis, S.T.: Think globally, fit locally: Unsupervised learning of low dimensional manifolds. *Journal of Machine Learning Research* 4, 119–155 (2003)
- [10] Donoho, D.L., Grimes, C.E.: Hessian eigenmaps: locally embedding techniques for high-dimensional data. *Proc. Natl. Acad. Sci. USA* 100, 5591–5596 (2003)
- [11] Belkin, M., Niyogi, P.: Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Computation* 15(6), 1373–1396 (2003)
- [12] Zhang, Z., Zha, H.: Principal manifolds and nonlinear dimension reduction via tangent space alignment. *SIAM Journal of Scientific Computing* 26(1), 313–338 (2004)
- [13] Kim, H., Park, H., Zha, H.: Distance preserving dimension reduction for manifold learning. In: *Proceedings of the 2007 SIAM International Conference on Data Mining, SDM 2007* (2007)
- [14] Dalamagas, T., Cheng, T., Winkel, K.J., Sellis, T.: A Methodology for Clustering XML Documents by Structure. *Information Systems* 31(3), 187–228 (2006)
- [15] Gao, J., Zhang, J.: Clustered SVD strategies in latent semantic indexing. *Inf. Process. Manag.* 41(5), 1051–1063 (2005)
- [16] Berry, M.W., Shakhina, A.P.: Computing sparse reduced-rank approximation to sparse matrices. *ACM Trans. Math. Software* 31(2), 252–269 (2005)
- [17] Tou, J.T., Gonzalez, R.C.: *Pattern Recognition Principles*. Addison-Wesley, Reading (1974)
- [18] Jain, A.K., Dubes, R.C.: *Algorithms for Clustering Data*. Prentice Hall, Englewood Cliffs (1988)
- [19] Frigui, H., Krishnapuram, R.: A Robust Competitive Clustering Algorithm with Application in Computer Vision. *IEEE Trans. Pattern Analysis and Machine Intelligence* 21(1), 450–465 (1999)
- [20] Everitt, B.S.: *Cluster Analysis*, 3rd edn. Halsted Press (1993)
- [21] Maulik, U., Bandyopadhyay, S.: Genetic Algorithm Based Clustering Technique. *Pattern Recognition* 33, 1455–1465 (2000)
- [22] Ye, Y.Q.: Comparing matrix methods in text-based information retrieval. — Tech. Rep., School of Mathematical Sciences, Peking University (2000)
- [23] ACM SIGMOD Record home page, <http://www.acm.org/sigmod/record/xml>
- [24] <http://www.cs.wisc.edu/niagara/data/>