

Dimension Reduction in Text Classification with Support Vector Machines

Hyunsoo Kim
Peg Howland
Haesun Park

Department of Computer Science and Engineering
University of Minnesota
200 Union Street S.E., 4-192 EE/CS Building
Minneapolis MN 55455, USA

HSKIM@CS.UMN.EDU
HOWLAND@CS.UMN.EDU
HPARK@CS.UMN.EDU

Editor: Nello Christianini

Abstract

Support vector machines (SVMs) have been recognized as one of the most successful classification methods for many applications including text classification. Even though the learning ability and computational complexity of training in support vector machines may be independent of the dimension of the feature space, reducing computational complexity is an essential issue to efficiently handle a large number of terms in practical applications of text classification. In this paper, we adopt novel dimension reduction methods to reduce the dimension of the document vectors dramatically. We also introduce decision functions for the centroid-based classification algorithm and support vector classifiers to handle the classification problem where a document may belong to multiple classes. Our substantial experimental results show that with several dimension reduction methods that are designed particularly for clustered data, higher efficiency for both training and testing can be achieved without sacrificing prediction accuracy of text classification even when the dimension of the input space is significantly reduced.

Keywords: dimension reduction, support vector machines, text classification, linear discriminant analysis, centroids

1. Introduction

Text classification is a supervised learning task for assigning text documents to pre-defined classes of documents. It is used to find valuable information from a huge collection of text documents available in digital libraries, knowledge databases, the world wide web (WWW), and company-wide intranets, to name a few. Several characteristics have been observed in vector space based methods for text classification (20; 21), including the high dimensionality of the input space, sparsity of document vectors, linear separability in most text classification problems, and the belief that few features are irrelevant. It has been conjectured that an aggressive dimension reduction may result in a significant loss of information, and therefore, result in poor classification results (13).

Assume that training data (\mathbf{x}_i, y_i) with $y_i \in \{-1, +1\}$ for $1 \leq i \leq n$ are given. The dual formulation of soft margin support vector machines (SVMs) with a kernel function K and control parameter

C is

$$\begin{aligned} \max_{\alpha_i} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j), \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, n. \end{aligned} \quad (1)$$

The kernel function

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle,$$

where \langle, \rangle denotes an inner product between two vectors, is introduced to handle nonlinearly separable cases without any explicit knowledge of the feature mapping ϕ . The formulation (1) shows that the computational complexity of SVM training depends on the number of training data samples which is denoted as n . The dimension of the feature space does not influence the computational complexity of training or testing due to the use of the kernel function.

However, an often neglected fact is that the computational complexity of training depends on the *dimension of the input space*. This is clear when we consider some typical kernel functions such as the linear kernel

$$K(\mathbf{x}, \mathbf{x}_i) = \langle \mathbf{x}, \mathbf{x}_i \rangle,$$

the polynomial kernel

$$K(\mathbf{x}, \mathbf{x}_i) = [\langle \mathbf{x}, \mathbf{x}_i \rangle + \beta]^d,$$

where d is the degree of the polynomial, and the Gaussian RBF (radial basis function) kernel

$$K(\mathbf{x}, \mathbf{x}_i) = \exp(-\gamma \|\mathbf{x} - \mathbf{x}_i\|^2),$$

where γ is a parameter to control. The evaluation of the kernel function *depends on the dimension of the input data*, since the kernel functions contain the inner product of two input vectors for the linear or polynomial kernels or the distance of two vectors for the Gaussian RBF kernel. Let α_i^* denote the optimal solution for (1). The optimal separating hyperplane $f(\mathbf{x}, \alpha^*, b)$ also requires evaluation of the kernel function since

$$f(\mathbf{x}, \alpha^*, b) = \sum_{\mathbf{x}_i \in SV} \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

where SV denotes the set of support vectors, b is a bias given by

$$b = -\frac{\min_{y_i=1} \langle w^*, \phi(\mathbf{x}_i) \rangle + \max_{y_i=-1} \langle w^*, \phi(\mathbf{x}_i) \rangle}{2}$$

and

$$w^* = \sum_{i=1}^l y_i \alpha_i^* \phi(\mathbf{x}_i).$$

Therefore, more efficient testing as well as training is expected from dimension reduction.

Throughout the paper, we will assume that the document set is represented in an $m \times n$ term-document matrix $A = (a_{ij})$, in which each column represents a document, and each entry a_{ij} represents the weighted frequency of term i in document j (1; 2). The clustering of data is assumed to be performed previously.

In the next section, we review Latent Semantic Indexing (LSI) (2; 1), which uses the truncated singular value decomposition (SVD) as a low-rank approximation of A . Although the truncated SVD provides the closest approximation to A in Frobenius or L_2 norm, LSI ignores the cluster structure while reducing the dimension of the data. In contrast, in Section 3, we review several dimension reduction methods that are especially effective for classification of clustered data: two methods based on centroids (16; 12), and one method which is a generalization of linear discriminant analysis (LDA) using the generalized singular value decomposition (GSVD) (10). With dimension reduction, computational complexity can be dramatically reduced for all classifiers including support vector machines and k-nearest neighbor classification. For k-nearest neighbor classification (kNN), the distances of vector pairs need to be computed when finding k nearest neighbors. Therefore, one can significantly reduce computational complexity by dimension reduction.

In many document data sets, documents can be assigned to more than one cluster upon classification. To handle this problem more effectively, we introduce a threshold based extension of several classification algorithms in Section 4. Our numerical experiments illustrate that the cluster-preserving dimension reduction algorithms we employ reduce the data dimension without any significant loss of information. In fact, in many cases, they seem to have the effect of noise reduction, since prediction accuracy becomes better after dimension reduction when compared to that in the original high dimensional input space.

2. Low-Rank Approximation Using Latent Semantic Indexing

LSI is based on the assumption that there is some underlying latent semantic structure in the term-document matrix that is corrupted by the wide variety of words used in documents and queries. This is referred to as the problem of polysemy and synonymy (6). The basic idea is that if two document vectors represent the same topic, they will share many associating words with a keyword, and they will have very close semantic structures after dimension reduction via SVD. Thus LSI/SVD breaks the original relationship of the data into linearly independent components (6), where the original term vectors are represented by left singular vectors and document vectors by right singular vectors. That is, if $l \leq \text{rank}(A)$, then

$$A \approx U_l \Sigma_l V_l^T$$

, where the columns of U_l are the leading l left singular vectors, Σ_l is an $l \times l$ diagonal matrix with the l largest singular values in nonincreasing order along its diagonal, and the columns of V_l are the leading l right singular vectors. Then $\Sigma_l V_l^T$ is the reduced dimensional representation of A , or equivalently, a new document $\mathbf{q} \in \mathbb{R}^{m \times 1}$ can be represented in the l -dimensional space as $\hat{\mathbf{q}} = U_l^T \mathbf{q}$.

This low-rank approximation has been widely applied in information retrieval (2). Since the complete orthogonal decomposition such as ULV or URV has computational advantages over the SVD including easier updating (22; 23; 24) and downdating (17), dimension reduction by these faster low-rank orthogonal decompositions has also been exploited (3). However, LSI ignores the cluster structure while reducing the dimension. In addition, since there is no theoretical optimum value for the reduced dimension, potentially expensive experimentation may be required to determine a reduced dimension l . As we report in Section 5, classification results after LSI vary depending upon the reduced dimension, classification method, and similarity measure employed. The experimental results confirm that when the data set is already clustered, the dimension reduction methods we present in the next section are more effective for classification of new data.

Algorithm 1 : Centroid algorithm for Dimension Reduction

Given a data set $A \in \mathbb{R}^{m \times n}$ with p clusters and a vector $\mathbf{q} \in \mathbb{R}^{m \times 1}$, this algorithm computes a p dimensional representation $\hat{\mathbf{q}} \in \mathbb{R}^{p \times 1}$ of \mathbf{q} .

1. Compute the centroid \mathbf{c}_i of the i th cluster, $1 \leq i \leq p$
 2. Set $C = [\mathbf{c}_1 \quad \mathbf{c}_2 \quad \cdots \quad \mathbf{c}_p]$
 3. Solve $\min_{\hat{\mathbf{q}}} \|C\hat{\mathbf{q}} - \mathbf{q}\|_2$
-

Algorithm 2 : Orthogonal Centroid algorithm for Dimension Reduction

Given a data set $A \in \mathbb{R}^{m \times n}$ with p clusters and a vector $\mathbf{q} \in \mathbb{R}^{m \times 1}$, this algorithm computes a p dimensional representation $\hat{\mathbf{q}}$ of \mathbf{q} .

1. Compute the centroid \mathbf{c}_i of the i th cluster, $1 \leq i \leq p$
 2. Set $C = [\mathbf{c}_1 \quad \mathbf{c}_2 \quad \cdots \quad \mathbf{c}_p]$
 3. Compute the reduced QR decomposition of C , which is $C = Q_p R$
 4. $\hat{\mathbf{q}} = Q_p^T \mathbf{q}$
-

3. Dimension Reduction Algorithms for Clustered Data

To achieve greater efficiency in manipulating data represented in a high dimensional space, it is often necessary to reduce the dimension *dramatically*. In this section, several dimension reduction methods that preserve the cluster structure are reviewed. Each method attempts to choose a projection to a reduced dimensional space that will capture the cluster structure of the data collection as much as possible.

3.1 Centroid-based Algorithms for Dimension Reduction of Clustered Data

Suppose we are given a data matrix A whose columns are grouped into p clusters. Instead of treating each column of the matrix A equally regardless of its membership in a specific cluster as in LSI/SVD, we want to find a lower dimensional representation Y of A so that the p clusters are preserved in Y . Given a term-document matrix, the problem is to find a transformation that maps each document vector in the m dimensional space to a vector in the l dimensional space for some $l < m$. For this, either the dimension reducing transformation $G^T \in \mathbb{R}^{l \times m}$ is computed explicitly or the problem is formulated as a rank reducing approximation where the given matrix A is to be decomposed into two matrices B and Y . That is,

$$A \approx BY \tag{2}$$

where $B \in \mathbb{R}^{m \times l}$ with $\text{rank}(B) = l$ and $Y \in \mathbb{R}^{l \times n}$ with $\text{rank}(Y) = l$. The matrix B accounts for the dimension reducing transformation. However, it is not necessary to compute the dimension reducing transformation G from B explicitly, as long as we can find the reduced dimensional representation of a given data item. If the matrix B is already determined, the matrix Y can be computed by solving

the least squares problem (8; 12; 16)

$$\min_{B,Y} \|BY - A\|_F. \quad (3)$$

Any given document $\mathbf{q} \in \mathbb{R}^{m \times 1}$ can be transformed to the lower dimensional space by solving the minimization problem

$$\min_{\hat{\mathbf{q}} \in \mathbb{R}^{l \times 1}} \|B\hat{\mathbf{q}} - \mathbf{q}\|_2. \quad (4)$$

Latent Semantic Indexing that utilizes the SVD (LSI/SVD) can be viewed as a variation of the model (2) with $B = U_l$ (16), where $U_l \Sigma_l V_l^T$ is the rank l truncated SVD of A . Then $\hat{\mathbf{q}} = U_l^T \mathbf{q}$ is obtained by solving the least squares problem

$$\min_{\hat{\mathbf{q}} \in \mathbb{R}^{l \times 1}} \|B\hat{\mathbf{q}} - \mathbf{q}\|_2 = \min_{\hat{\mathbf{q}} \in \mathbb{R}^{l \times 1}} \|U_l \hat{\mathbf{q}} - \mathbf{q}\|_2. \quad (5)$$

In the Centroid dimension reduction algorithm (see Algorithm 1), the i th column of B is the centroid vector of the i th cluster, which is the average of the data items in the i th cluster, for $1 \leq i \leq p$. This matrix B is called the centroid matrix. Then, any vector $\mathbf{q} \in \mathbb{R}^{m \times 1}$ can be represented in the p dimensional space as $\hat{\mathbf{q}}$, the solution of the least squares problem (4), where B is the centroid matrix. In the Orthogonal Centroid algorithm (see Algorithm 2), the p dimensional representation of a data vector $\mathbf{q} \in \mathbb{R}^{m \times 1}$ is given as $\hat{\mathbf{q}} = Q_p^T \mathbf{q}$ where Q_p is an orthonormal basis for the centroid matrix obtained from its QR decomposition.

The centroid-based dimension reduction algorithms are computationally less costly than LSI/SVD. They are also more effective when the data are already clustered. Although the centroid-based schemes can be applied only when the data are linearly separable, they are suitable for text classification problems, since text data is usually linearly separable in the original dimensional space (13). For a nonlinear extension of the Orthogonal Centroid method that utilizes kernel functions, see (18).

3.2 Generalized Discriminant Analysis based on the Generalized Singular Value Decomposition

Recently, a new algorithm has been developed for cluster-preserving dimension reduction based on the generalized singular value decomposition (GSVD) (10). This algorithm generalizes classical discriminant analysis, by extending its application to very high-dimensional data such as that encountered in text classification.

Classical discriminant analysis (7; 25) preserves cluster structure by maximizing the scatter between clusters while minimizing the scatter within clusters. For this purpose, the within-cluster scatter matrix S_w and the between-cluster scatter matrix S_b are defined. If we denote by N_i the set of column indices that belong to the cluster i , n_i the number of columns in cluster i , and \mathbf{c} the global centroid, then

$$S_w = \sum_{i=1}^p \sum_{j \in N_i} (\mathbf{a}_j - \mathbf{c}_i)(\mathbf{a}_j - \mathbf{c}_i)^T,$$

and

$$\begin{aligned} S_b &= \sum_{i=1}^p \sum_{j \in N_i} (\mathbf{c}_i - \mathbf{c})(\mathbf{c}_i - \mathbf{c})^T \\ &= \sum_{i=1}^p n_i (\mathbf{c}_i - \mathbf{c})(\mathbf{c}_i - \mathbf{c})^T. \end{aligned}$$

Algorithm 3 LDA/GSVD

Given a data matrix $A \in \mathbb{R}^{m \times n}$ with p clusters, this algorithm computes the columns of the matrix $G \in \mathbb{R}^{m \times (p-1)}$, which preserves the cluster structure in the reduced dimensional space, and it also computes the $p-1$ dimensional representation Y of A .

1. Compute $H_b \in \mathbb{R}^{m \times p}$ and $H_w \in \mathbb{R}^{m \times n}$ from A according to Eqns. (7) and (6), respectively.
2. Compute the complete orthogonal decomposition of $H = (H_b, H_w)^T \in \mathbb{R}^{(p+n) \times m}$, which is

$$P^T H Q = \begin{pmatrix} R & 0 \\ 0 & 0 \end{pmatrix}.$$

3. Let $t = \text{rank}(H)$.
4. Compute W from the SVD of $P(1:p, 1:t)$, which is $U^T P(1:p, 1:t)W = \Sigma_A$.
5. Compute the first $p-1$ columns of

$$X = Q \begin{pmatrix} R^{-1}W & 0 \\ 0 & I \end{pmatrix},$$

and assign them to G .

6. $Y = G^T A$
-

Since

$$\text{trace}(S_w) = \sum_{i=1}^p \sum_{j \in N_i} \|\mathbf{a}_j - \mathbf{c}_i\|_2^2$$

measures the closeness within the clusters, and

$$\text{trace}(S_b) = \sum_{i=1}^p \sum_{j \in N_i} \|\mathbf{c}_i - \mathbf{c}\|_2^2$$

measures the remoteness between the clusters, the goal is to minimize the former while maximizing the latter in the reduced dimensional space. Once again letting $G^T \in \mathbb{R}^{l \times m}$ denote the transformation that maps a column of A in the m dimensional space to a vector in the l dimensional space, the goal can be expressed as the simultaneous minimization of $\text{trace}(G^T S_w G)$ and maximization of $\text{trace}(G^T S_b G)$.

When S_w is nonsingular, this simultaneous optimization is commonly approximated by maximizing

$$J_1(G) = \text{trace}((G^T S_w G)^{-1} (G^T S_b G)).$$

It is well known that the global maximum is achieved when the columns of G are the eigenvectors of $S_w^{-1} S_b$ that correspond to the l largest eigenvalues (7; 25). In fact, when the reduced dimension $l \geq p-1$, $\text{trace}(S_w^{-1} S_b)$ is exactly preserved upon dimension reduction, and equals $\lambda_1 + \dots + \lambda_{p-1}$, where each $\lambda_i \geq 0$. Without loss of generality, we assume that the term-document matrix A is partitioned as

$$A = [A_1, \quad \dots, \quad A_p]$$

where the columns of each block $A_i \in \mathbb{R}^{m \times n_i}$ belong to the cluster i . Defining the matrices

$$H_w = [\mathbf{a}_1 - \mathbf{c}_1, \mathbf{a}_2 - \mathbf{c}_1, \dots, \mathbf{a}_n - \mathbf{c}_p] \in \mathbb{R}^{m \times n} \quad (6)$$

and

$$H_b = [\sqrt{n_1}(\mathbf{c}_1 - \mathbf{c}), \dots, \sqrt{n_p}(\mathbf{c}_p - \mathbf{c})] \in \mathbb{R}^{m \times p}, \quad (7)$$

then

$$S_w = H_w H_w^T \quad \text{and} \quad S_b = H_b H_b^T.$$

As the product of an $m \times n$ matrix with an $n \times m$ matrix, S_w will be singular when the number of terms m exceeds the number of documents n . In that case, classical discriminant analysis fails. However, if we rewrite the eigenvalue problem $S_w^{-1} S_b \mathbf{x}_i = \lambda_i \mathbf{x}_i$ as

$$\beta_i^2 H_b H_b^T \mathbf{x}_i = \alpha_i^2 H_w H_w^T \mathbf{x}_i,$$

it can be solved by the GSVD.

The resulting algorithm, called LDA/GSVD, is summarized in Algorithm 3. It follows the construction of the Paige and Saunders (15) proof, but only computes the necessary part of the GSVD. The most expensive step of LDA/GSVD is the complete orthogonal decomposition of the composite H matrix in Step 2. When $\max(p, n) \ll m$, the SVD of $H = [H_b^T, H_w^T] \in \mathbb{R}^{(p+n) \times m}$ can be computed by first computing the reduced QR decomposition $H^T = Q_H R_H$, and then computing the SVD of $R_H \in \mathbb{R}^{(p+n) \times (p+n)}$ as

$$R_H = Z \begin{pmatrix} \Sigma_H & 0 \\ 0 & 0 \end{pmatrix} P^T.$$

This gives

$$H = R_H^T Q_H^T = P \begin{pmatrix} \Sigma_H & 0 \\ 0 & 0 \end{pmatrix} Z^T Q_H^T,$$

where the columns of $Q_H Z \in \mathbb{R}^{m \times (p+n)}$ are orthonormal. There exists orthogonal $Q \in \mathbb{R}^{m \times m}$ whose first $p+n$ columns are $Q_H Z$. Hence

$$H = P \begin{pmatrix} \Sigma_H & 0 \\ 0 & 0 \end{pmatrix} Q^T,$$

where there are now $m - t$ zero columns to the right of Σ_H . Since $R_H \in \mathbb{R}^{(p+n) \times (p+n)}$ is a much smaller matrix than $H \in \mathbb{R}^{(p+n) \times m}$, the required memory is substantially reduced. In addition, the computational complexity of the algorithm is reduced to $O(mn^2) + O(n^3)$ (8), since this step is the dominating part.

4. Classification Methods

To test the effect of dimension reduction in text classification, three different classification methods were used: centroid-based classification, k-nearest neighbor (kNN), and support vector machines (SVMs). Each classification method is modified by introducing some threshold values to perform classification correctly when a document has membership in multiple classes. In this section, we briefly review the three classification methods and discuss their modifications.

Algorithm 4 : Centroid-based Classification

Given a data matrix A with p clusters and p corresponding centroids, \mathbf{c}_i , $1 \leq i \leq p$, and a vector $\mathbf{q} \in \mathbb{R}^{m \times 1}$, this method finds the index j of the cluster in which the vector \mathbf{q} belongs.

- find the index j such that $\text{sim}(\mathbf{q}, \mathbf{c}_i)$, $1 \leq i \leq p$, is minimum (or maximum), where $\text{sim}(\mathbf{q}, \mathbf{c}_i)$ is the similarity measure between \mathbf{q} and \mathbf{c}_i . (For example, $\text{sim}(\mathbf{q}, \mathbf{c}_i) = \|\mathbf{q} - \mathbf{c}_i\|_2$ using the L_2 norm, and we take the index with the minimum value. Using the cosine measure,

$$\text{sim}(\mathbf{q}, \mathbf{c}_i) = \cos(\mathbf{q}, \mathbf{c}_i) = \frac{\mathbf{q}^T \mathbf{c}_i}{\|\mathbf{q}\|_2 \|\mathbf{c}_i\|_2},$$

and we take the index with the maximum value.)

4.1 Centroid-based Classification

Centroid-based classification, summarized in Algorithm 4, is one of the simplest classification methods. A test document is assigned to a class that has the most similar centroid. Using the cosine similarity measure, we can classify a test document \mathbf{q} by computing

$$\arg \max_{1 \leq i \leq p} \frac{\mathbf{q}^T \mathbf{c}_i}{\|\mathbf{q}\|_2 \|\mathbf{c}_i\|_2} \quad (8)$$

where \mathbf{c}_i is the centroid of the i th cluster of the training data. When dimension reduction is performed by the Centroid algorithm, the centroids of the full space become the columns $\mathbf{e}_i \in \mathbb{R}^{p \times 1}$ of the identity matrix. Then the decision rule becomes

$$\arg \max_{1 \leq i \leq p} \frac{\hat{\mathbf{q}}^T \mathbf{e}_i}{\|\hat{\mathbf{q}}\|_2 \|\mathbf{e}_i\|_2}, \quad (9)$$

where $\hat{\mathbf{q}}$ is the reduced dimensional representation of the document \mathbf{q} . This shows that classification can be performed by simply finding the index i of the vector $\hat{\mathbf{q}}$ with the largest component. Centroid-based classification has the advantage that the computation involved is extremely simple. We can also classify using the L_2 norm similarity measure by finding the centroid that is closest to \mathbf{q} in L_2 norm.

The original form of centroid-based classification finds the nearest centroid and assigns the corresponding class as the predicted class. To allow an assignment of any document to multiple classes, we introduce the decision rule for centroid-based classification as

$$y(\mathbf{x}, j) = \text{sign}\{\text{sim}(\mathbf{x}, \mathbf{c}_j) - \theta_j^c\}, \quad (10)$$

where $y(\mathbf{x}, j) \in \{+1, -1\}$ is the classification for document \mathbf{x} with respect to class j (if $y > 0$ then the class is j , else the class is not j), $\text{sim}(\mathbf{x}, \mathbf{c}_j)$ is the similarity between the test document \mathbf{x} and the centroid vector \mathbf{c}_j for the class j , and θ_j^c is the class specific threshold for the binary decision for $y(\mathbf{x}, j)$ in centroid-based classification. In this way, document \mathbf{x} will be a member of class j if its similarity to the centroid vector \mathbf{c}_j for the class is above the threshold.

Algorithm 5 : k Nearest Neighbor (kNN) Classification

Given a data matrix $A = [\mathbf{a}_1, \dots, \mathbf{a}_n]$ with p clusters and a vector $\mathbf{q} \in \mathbb{R}^{m \times 1}$, this method finds the cluster in which the vector \mathbf{q} belongs.

1. Using the similarity measure $sim(\mathbf{q}, \mathbf{a}_j)$ for $1 \leq j \leq n$, find the k nearest neighbors of \mathbf{q} .
2. Among these k vectors, count the number belonging to each cluster.
3. Assign \mathbf{q} to the cluster with the greatest count in the previous step.

4.2 k-Nearest Neighbor Classification

The kNN algorithm, summarized in Algorithm 5, is one of the most commonly used classification methods. To correctly predict the membership of a document which belongs to multiple classes, we used the following modified decision rule for kNN (29):

$$y(\mathbf{x}, j) = \text{sign} \left\{ \sum_{\mathbf{d}_i \in kNN} sim(\mathbf{x}, \mathbf{d}_i) y(\mathbf{d}_i, j) - \theta_j^{kNN} \right\} \quad (11)$$

where kNN is the set of k nearest neighbors for document \mathbf{x} , $y(\mathbf{d}_i, j) \in \{+1, -1\}$ is the classification for document \mathbf{d}_i with respect to class j (if $y > 0$ then the class is j , else the class is not j), $sim(\mathbf{x}, \mathbf{d}_i)$ is the similarity between the test document \mathbf{x} and the training document \mathbf{d}_i , and θ_j^{kNN} is the class specific threshold for kNN classification.

4.3 Support Vector Machines

The optimal separating hyperplane of the one-vs-rest binary classifier can be obtained by conventional SVMs. We introduce the following decision rule for support vector machines as

$$y(\mathbf{x}, j) = \text{sign} \left\{ \sum_{\mathbf{x}_i \in SV} \alpha_i y_i K(\mathbf{x}, \mathbf{x}_i) + b - \theta_j^{SVM} \right\}, \quad (12)$$

where $y(\mathbf{x}, j) \in \{+1, -1\}$ is the classification for document \mathbf{x} with respect to class j , SV is the set of support vectors, and θ_j^{SVM} is the class specific threshold for the binary decision. This threshold is set so that a new document \mathbf{x} must not be classified to belong to class j when it is located very close to the optimal separating hyperplane, i.e. when the decision is made with a low reliability. We use the linear kernel $K = \langle \mathbf{x}, \mathbf{x}_i \rangle$, the polynomial kernel $K = [\langle \mathbf{x}, \mathbf{x}_i \rangle + 1]^d$, where d is the degree of the polynomial, and the Gaussian RBF (radial basis function) kernel $K = \exp(-\gamma \|\mathbf{x} - \mathbf{x}_i\|^2)$, where γ is a parameter that controls the width of the Gaussian function.

5. Experimental Results

Prediction results are compared for the test documents in the full space without any dimension reduction as well as those in the reduced space obtained by LSI/SVD, Centroid, Orthogonal Centroid, and LDA/GSVD dimension reduction methods. For SVMs, we optimized the regularization parameter C , polynomial degree d for the polynomial kernel, and γ for the Gaussian RBF (radial basis function) kernel for each full and reduced dimension data set.

classification methods	The rank- l approximation of LSI/SVD								Full
	$l=5$	$l=100$	$l=200$	$l=300$	$l=500$	$l=1000$	$l=1246$	$l=1247$	
centroid (L_2)	71.6	82.2	83.4	83.9	84.8	84.9	85.2	85.2	85.2
centroid (Cosine)	78.5	86.9	87.1	87.6	88.0	88.2	88.3	88.3	88.3
5NN (L_2)	77.8	68.8	55.4	49.2	63.8	76.9	79.0	79.0	79.0
15NN (L_2)	77.5	69.7	52.7	50.3	76.3	74.7	83.4	83.4	83.4
30NN (L_2)	77.5	64.3	47.8	58.0	80.8	73.2	83.8	83.8	83.8
5NN (Cosine)	77.8	82.2	79.1	79.6	79.4	78.7	77.8	77.8	77.8
15NN (Cosine)	80.2	83.1	82.5	83.6	82.9	82.5	82.5	82.5	82.5
30NN (Cosine)	79.8	83.4	83.8	84.1	84.2	84.1	83.8	83.8	83.8
SVM	79.1	87.6	88.4	88.5	88.6	89.2	89.7	89.7	88.9

Table 1: Text classification accuracy (%) using centroid-based classification, k-nearest neighbor classification, and SVMs, with LSI/SVD dimension reduction on the MEDLINE data set. The Euclidean norm (L_2) and the cosine similarity measure (Cosine) were used for the centroid-based and kNN classification.

The first data set that we used was a subset of the MEDLINE database with 5 classes. Each class has 500 documents. The set was divided into 1250 training documents and 1250 test documents. After stemming and stoplist removal, the training set contains 22095 distinct terms. For this data, each document belongs to only one class, and we used the original form of the three classification algorithms without introducing the threshold.

The second data set was the “ModApte” split of the Reuter-21578 text collection. We only used 90 classes for which there is at least one training and one test example in each class. It contains 7769 training documents and 3019 test documents. The training set contains 11941 distinct terms after preprocessing with stoplist removal and stemming. The Reuter data set contains documents that belong to multiple classes, so the classification methods utilize thresholds.

We used a standard weight factor for each word stem:

$$\phi_i(\mathbf{x}) = \frac{tf_i \log(idf_i)}{\kappa}, \quad (13)$$

where tf_i is the number of occurrences of term i in document \mathbf{x} , $idf_i = n/d$ is the ratio between the total number of documents n and the number of documents d containing the term, and κ is the normalization constant that makes $\|\phi\|_2 = 1$.

Table 1 reports text classification accuracy for the MEDLINE data set using LSI/SVD with a range of values for the reduced dimension. The smallest reduced dimension, $l = 5$, is included in order to compare with centroid-based and LDA/GSVD methods, which reduce the dimension to 5 and 4, respectively. Since the training set has the nearly-full rank of 1246, we include the reduced dimensions 1246 and 1247 at the high end of the range. For a training set of size 1250, the reduced dimension $l = 300$ is generous. However, we observe that kNN classification with L_2 norm similarity produces poor classification results for l values from 100 to 500. This is consistent with the common belief that cosine similarity performs better with unnormalized text data. Also, classification accuracy using 5NN lags that for higher values of k , suggesting that $k=5$ is too small for classes

kernel	Dimension reduction methods				
	Full	Centroid	Orthogonal	LDA/ GSVD4	LDA/ GSVD5
	22095×1250	5×1250	5×1250	4×1250	5×1250
linear (C=1.0)	88.1	88.9	85.9	86.5	86.6
linear (C=10.0)	88.9	88.5	88.3	86.7	86.7
linear (C=50.0)	88.9	87.7	88.8	87.1	87.1
linear ^{opt}	88.9	88.9	89.0	87.4	87.4
polynomial(d=2)	88.6	88.9	88.9	87.3	87.3
polynomial(d=3)	88.0	89.0	88.8	87.4	87.4
polynomial(d=4)	87.5	88.9	88.8	87.2	87.2
polynomial(d=5)	86.5	88.6	88.8	87.1	87.1
polynomial ^{opt}	88.6	89.0	88.9	87.4	87.4
RBF ($\gamma = 0.5$)	88.5	89.0	89.0	87.1	87.2
RBF ($\gamma = 1.0$)	87.6	89.2	89.0	87.3	87.2
RBF ($\gamma = 1.5$)	86.3	89.1	88.8	87.4	87.3
RBF ^{opt}	88.7	89.2	89.0	87.4	87.3

Table 2: Text classification accuracy (%) with different kernels in SVMs with and without dimension reduction on the MEDLINE data set. The regularization parameter C for each case was optimized by numerical experiments. Dimension of each training term-document matrix is shown. LDA/GSVD4 and LDA/GSVD5 represent the results from LDA/GSVD where the reduced dimensions are 4 and 5, respectively.

of size 250. It is noteworthy that even with LSI, which makes no attempt to preserve the cluster structure upon dimension reduction, SVM classification achieves very consistent classification results for reduced dimensions of 100 or greater, and the SVM accuracy exceeds that of the other classification methods.

Table 2 shows text classification accuracy (%) with different kernels in SVMs, with and without dimension reduction on the MEDLINE data set. Note that the linear^{opt} values are optimal over all the values of the regularization parameter C that we tried, and the RBF^{opt} values are optimal over all the γ values we tried. This table shows that the prediction results in the reduced dimension are similar to those in the original full dimensional space, while achieving a significant reduction in time and space complexity. In the reduced space obtained by the Orthogonal Centroid dimension reduction algorithm, the classification accuracy is insensitive to the choice of the kernel. Thus, we can choose the linear kernel in this case instead of the computationally more expensive polynomial or RBF kernel.

Table 3 shows classification accuracy obtained by all three classification methods – centroid-based, kNN with three different values of k , and the optimal result from SVM – for each dimension reduced data set and the full space. For the LDA/GSVD dimension reduction method, the classification accuracy with cosine similarity measure is lower with centroid-based classification as well as with kNN, while the results with L_2 norm are better. This is due to the formulation of trace optimization criteria in terms of the L_2 norm. With LDA/GSVD, documents from the same class in

classification methods	Dimension reduction methods				
	Full 22095×1250	Centroid	Orthogonal	LDA/ GSVD4	LDA/ GSVD5
		5×1250	5×1250	4×1250	5×1250
centroid (L_2)	85.2	88.0	85.2	88.7	88.7
centroid (Cosine)	88.3	88.0	88.3	83.9	83.9
5NN (L_2)	79.0	88.4	88.6	81.5	86.6
15NN (L_2)	83.4	88.3	87.8	88.7	88.6
30NN (L_2)	83.8	88.8	88.5	88.7	88.5
5NN (Cosine)	77.8	88.6	88.2	83.8	84.1
15NN (Cosine)	82.5	88.2	88.5	83.8	84.1
30NN (Cosine)	83.8	88.3	88.6	83.8	84.1
SVM	88.9	89.2	89.0	87.4	87.4

Table 3: Text classification accuracy (%) using centroid-based classification, k-nearest neighbor classification, and SVMs, with and without dimension reduction on the MEDLINE data set. The Euclidean norm (L_2) and the cosine similarity measure (Cosine) were used for centroid-based and kNN classification.

class	Dimension reduction				
	Full 22095×1250	Centroid	Orthogonal	LDA/ GSVD4	LDA/ GSVD5
		5×1250	5×1250	4×1250	5×1250
heart attack	92.4	94.4	94.4	92.4	92.4
colon cancer	84.8	84.8	86.0	83.2	83.2
glycemic	95.6	97.6	98.0	95.2	95.2
oral cancer	82.0	75.2	73.6	78.8	78.8
tooth decay	89.6	94.0	92.8	87.2	87.2
microavg	88.9	89.2	89.0	87.4	87.4

Table 4: Text classification accuracy (%) of the 5 classes and the microaveraged performance over all 5 classes on the MEDLINE data set. All results are from SVMs using optimal kernels.

the full dimensional space tend to be transformed to a very tight cluster or even to a single point in the reduced space, since the LDA/GSVD algorithm tends to minimize the trace of the within cluster scatter. This seems to make it difficult for SVMs to find a binary classifier with low generalization error.

Table 4 shows text classification accuracy for the 5 classes using SVMs with and without dimension reduction methods on the MEDLINE data set. The colon cancer and oral cancer documents were relatively hard to classify correctly.

The REUTERS data set has many documents that are classified to more than 2 classes, whereas no document is classified to belong to more than one class in the MEDLINE data set. While we

classification methods	Full 11941×9579	Dimension reduction	
		Centroid 90×9579	Orthogonal Centroid 90×9579
centroid(L_2)	78.89	73.32	78.00
centroid(Cosine)	80.45	74.79	80.46
15NN	78.65	81.70	85.51
30NN	80.21	81.94	86.19
45NN	80.29	81.01	84.79
SVM	87.11	84.54	87.03

Table 5: Comparison of micro-averaged F_1 scores for 3 different classification methods with and without dimension reduction on the REUTERS data set. The Euclidean norm (L_2) and the cosine similarity measure (Cosine) were used for the centroid-based classification. The cosine similarity measure was used for the kNN classification. The dimension of the full training term-document matrix is 11941×9579 and that of the reduced matrix is 90×9579.

could handle relatively large matrices using a sparse matrix representation and sparse QR decomposition in the Centroid and Orthogonal Centroid dimension reduction methods, results for the LDA/GSVD dimension reduction method are not reported, since we ran out of memory while computing the GSVD. For this data set, we built a series of threshold-based classifiers, optimizing the thresholds to capture the multiple class membership. All class specific thresholds (θ_j^{kNN} , θ_j^c , θ_j^{SVM}) are determined by numerical experiments. Though we obtained precision/recall break even points by optimizing the thresholds, we report values of the F_1 measure (26) which is defined as

$$F_1 = \frac{2rp}{r+p}, \tag{14}$$

where r is recall and p is precision for a binary classification. Table 5 shows that the effectiveness of classification was preserved for the Orthogonal Centroid dimension reduction algorithm, while it became worse for the Centroid dimension reduction algorithm. This is due to a property of the Centroid algorithm that the centroids of the full space are projected to the columns of the identity matrix in the reduced space. This orthogonality between the centroids may make it difficult to represent the multiclass membership of a document by separating closely related classes after dimension reduction. The pattern of prediction measure F_1 for each class is also preserved by Orthogonal Centroid in Table 6. The macro-averaged F_1 and micro-averaged F_1 for the 10 most frequent classes are also presented.

6. Conclusion and Discussion

In this paper, we applied three methods, Centroid, Orthogonal Centroid, and LDA/GSVD, which are designed for reducing the dimension of clustered data. For comparison, we also applied LSI/SVD, which does not attempt to preserve cluster structure upon dimension reduction. We tested the effectiveness in classification with dimension reduction using three different classification methods:

class	Full 11941×9579	Dimension reduction	
		Centroid 90×9579	Orthogonal Centroid 90×9579
earn	98.25	97.49	96.60
acq	95.57	95.45	94.94
money-fx	75.78	77.97	79.44
grain	92.88	86.62	92.26
crude	88.11	86.49	87.70
trade	75.32	75.11	77.25
interest	77.99	78.13	83.21
ship	84.09	85.71	88.00
wheat	84.14	81.94	84.06
corn	87.27	74.78	89.47
microavg (top 10)	92.21	91.32	92.21
avg (top 10)	85.94	83.96	87.32
microavg(all)	87.11	84.54	87.03

Table 6: F_1 scores of the 10 most frequent classes and micro-averaged performance over all 90 classes on the REUTERS data set. All results are from SVMs using optimal kernels. The dimension of the full training term-document matrix is 11941×9579 and that of the reduced matrix is 90×9579.

SVMs, kNN, and centroid-based classification. For the three cluster-preserving methods, the results show surprisingly high prediction accuracy, which is essentially the same as in the original full space, even with very dramatic dimension reduction. They justify dimension reduction as a worthwhile preprocessing stage for achieving high efficiency and effectiveness. Especially for kNN classification, the savings in computational complexity in classification after dimension reduction are significant. In the case of SVM the savings are also clear, since the distance between two pairs of input data points need to be computed repeatedly with and without the use of the kernel function, and the vectors become significantly shorter with dimension reduction.

We have also introduced threshold based classifiers for centroid-based classification and SVMs in order to capture the overlap structure between closely related classes. Prediction results with the Centroid dimension reduction method became better compared to those from the full space for the completely disjoint MEDLINE data set, but became worse for the REUTERS data set. Since the Centroid dimension reduction method maps the centroids to unit vectors \mathbf{e}_i which are orthogonal to each other, it is helpful for the disjoint data set, but not for a data set which contains documents belonging multiple classes. We observed that prediction accuracy with the Orthogonal Centroid dimension reduction algorithm was preserved for SVMs as well as with centroid-based classification. The Orthogonal Centroid dimension reduction method maximizes the between cluster relationship using the relatively inexpensive reduced QR decomposition, compared to LDA/GSVD which also considers the within cluster relationship but requires a more expensive rank revealing decomposition such as the singular value decomposition (10; 11).

The better prediction accuracy using SVMs is due to low generalization error by maximizing the margin, and the capability to handle non-linearity by kernel choice. Although most classes of the Reuters-21578 data set are linearly separable (13), there seems to be some level of non-linearity. For non-linearly separable data, SVMs with appropriate nonlinear kernel functions would work as a better classifier. Another way to handle non-linearly separable data is to apply nonlinear extensions of the dimension reduction methods, including those presented in (18; 19). All of the dimension reduction methods presented here can also be applied to visualize the higher dimensional structure by reducing the dimension to 2- or 3-dimensional space.

We conclude that dramatic dimension reduction of text documents can be achieved, without sacrificing classification accuracy. For the document sets we tested, the Orthogonal Centroid method did particularly well at preserving the cluster structure from the full dimensional representation. That is, the prediction accuracies for Orthogonal Centroid rival those of the full space, even though the dimension is reduced to the number of clusters. The savings in computational complexity are significant using either kNN classification or SVM.

Acknowledgments

This material is based upon work supported by the National Science Foundation Grant No. CCR-0204109. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation (NSF). The authors would also like to thank University of Minnesota Supercomputing Institute (MSI) for providing the computing facilities.

References

- [1] M. W. Berry, Z. Drmac, and E. R. Jessup. Matrices, vector spaces, and information retrieval. *SIAM Review*, 41:335–362, 1999.
- [2] M. W. Berry, S. T. Dumais, and G. W. O’Brien. Using linear algebra for intelligent information retrieval. *SIAM Review*, 37:573–595, 1995.
- [3] M. W. Berry and R. D. Fierro. Low-rank orthogonal decompositions for information retrieval applications. *Numerical Linear Algebra with Applications*, 3(4):301–327, 1996.
- [4] Å. Björck. *Numerical Methods for Least Square Problems*. SIAM, Philadelphia, PA, 1996.
- [5] N. Cristianini and J. Shawe-Taylor. *Support Vector Machines and Other Kernel-based Learning Methods*. Cambridge University Press, 2000.
- [6] S. Deerwester, S.T. Dumais, G.W. Furnas, T.K. Landauer, and R. Harshman. Indexing by latent semantic analysis. *Journal of the Society for Information Science*, 41:391-407, 1990.
- [7] K. Fukunaga, *Introduction to Statistical Pattern Recognition*, Second ed., Academic Press, 1990.
- [8] G. H. Golub and C. F. Van Loan. *Matrix Computations, third edition*. Johns Hopkins University Press, Baltimore, 1996.

- [9] M. Heiler. *Optimization Criteria and Learning Algorithms for Large Margin Classifiers*. Diploma Thesis, University of Mannheim., 2002.
- [10] P. Howland, M. Jeon, and H. Park. Structure Preserving Dimension Reduction for Clustered Text Data based on the Generalized Singular Value Decomposition. *SIAM Journal of Matrix Analysis and Applications*, 25(1):165–179, 2003.
- [11] P. Howland and H. Park. Generalizing discriminant analysis using the generalized singular value decomposition, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(8): 995-1006, 2004.
- [12] M. Jeon, H. Park, and J. B. Rosen. Dimensional reduction based on centroids and least squares for efficient processing of text data. In *Proceedings for the First SIAM International Workshop on Text Mining*. Chicago, IL, 2001.
- [13] T. Joachims. Text categorization with support vector machines: Learning with many relevant features. In *Proceedings of the European Conference on Machine Learning*, pages 137–142, Berlin, 1998.
- [14] H. Lodhi, N. Cristianini, J. Shawe-Taylor, and C. Watkins. Text classification using string kernels. *Advances in Neural Information Processing Systems*, 13:563–569, 2000.
- [15] C. C. Paige and M. A. Saunders, Towards a generalized singular value decomposition, *SIAM Journal of Numerical Analysis*, 18, pp. 398–405, 1981.
- [16] H. Park, M. Jeon, and J. B. Rosen. Lower dimensional representation of text data based on centroids and least squares, *BIT Numerical Mathematics*, 42(2):1–22, 2003.
- [17] H. Park and L. Eldén. DOWDATING THE RANK-REVEALING URV DECOMPOSITION. *SIAM Journal of Matrix Analysis and Applications*, 16, pp. 138–155, 1995.
- [18] C. Park and H. Park. Nonlinear feature extraction based on centroids and kernel functions. *Pattern Recognition*, to appear.
- [19] C. Park and H. Park. Kernel discriminant analysis based on the generalized singular value decomposition. Technical report 03-017, Department of Computer Science and Engineering, University of Minnesota, 2003.
- [20] G. Salton, *The SMART Retrieval System*, Prentice Hall, 1971.
- [21] G. Salton and M. J. McGill, *Introduction to Modern Information Retrieval*, McGraw-Hill, 1983.
- [22] G. W. Stewart. An updating algorithm for subspace tracking. *IEEE Transactions on Signal Processing*, 40:1535–1541, 1992.
- [23] G. W. Stewart. Updating URV decompositions in parallel. *Parallel Computing*, 20(2):151–172, 1994.
- [24] M. Stewart and P. Van Dooren. Updating a generalized URV decomposition. *SIAM Journal of Matrix Analysis and Applications*, 22(2):479–500, 2000.

- [25] S. Theodoridis and K. Koutroumbas, *Pattern Recognition*, Academic Press, 1999.
- [26] C. J. van Rijsbergen. *Information Retrieval*. Butterworths, London, 1979.
- [27] V. Vapnik. *The Nature of Statistical Learning Theory*. Springer-Verlag, New York, 1995.
- [28] V. Vapnik. *Statistical Learning Theory*. John Wiley & Sons, New York, 1998.
- [29] Y. Yang and X. Liu. A re-examination of text categorization methods. In *22nd Annual International SIGIR*, pages 42–49, Berkeley, August 1999.