

On Computational Issues of Semi-supervised Local Fisher Discriminant Analysis

Masashi Sugiyama (sugi@cs.titech.ac.jp)

Department of Computer Science, Tokyo Institute of Technology
2-12-1 O-okayama, Meguro-ku, Tokyo 152-8552, Japan

Abstract

Dimensionality reduction is one of the important preprocessing steps in practical pattern recognition. *SEmi-supervised Local Fisher discriminant analysis* (SELF)—which is a semi-supervised and local extension of Fisher discriminant analysis—was shown to work excellently in experiments. However, when data dimensionality is very high, a naive use of SELF is prohibitive due to high computational costs and large memory requirement. In this paper, we introduce computational tricks for making SELF applicable to large-scale problems.

Keywords

dimensionality reduction, semi-supervised learning, semi-supervised local Fisher discriminant analysis, sparsity, generalized eigenvalue problem

1 Introduction

Reducing dimensionality of data is one of the important challenges in pattern recognition. However, a naive use of supervised dimensionality reduction methods in high-dimensional scenarios often results in overfitting. In such cases, the use of *unlabeled samples* could be helpful [1]. A semi-supervised dimensionality reduction method called *SEmi-supervised Local Fisher discriminant analysis* (SELF) [6] has been proposed and shown to work excellently. SELF is a regularized variant of a supervised dimensionality reduction method called *Local Fisher discriminant analysis* (LFDA) [5].

An advantage of SELF in addition to its good performance is that the globally optimal solution can be obtained analytically by solving a generalized eigenvalue problem. Dimensionality of the generalized eigenvalue problem depends only on dimensionality of the feature vectors, not on the number of samples. So SELF may be applicable to a dataset with a large number of samples as long as dimensionality of the feature vectors

Table 1: Computational issues of SELF. d and n are dimensionality and the number of samples.

d	n	Sparseness	Formulation
Moderate	Large	None	Primal
Large	Moderate	None	Dual
Large	Large	Sample	Primal
Large	Large	Kernel	Dual

is moderate. However, when the feature dimensionality is very high, solving the generalized eigenvalue problem may not be computationally tractable—even more critically, the matrix which needs to be eigendecomposed cannot be stored in memory. This critical limitation makes SELF inapplicable to high-dimensional real-world problems.

The purpose of this paper is to provide computationally efficient algorithms for SELF and make it applicable to high-dimensional problems. More specifically, we introduce the following computational tricks. The first is based on the dual formulation, where the size of the generalized eigenvalue problem does not depend on the feature dimensionality but only on the number of samples. Thus this dual formulation would be computationally efficient if the number of samples is moderate. The other method makes use of the *sparsity* of feature vectors—we show that the high-dimensional *dense* matrix which needs to be eigendecomposed in SELF can be expressed by the sum of a sparse matrix and low-rank matrices. This structure allows us to efficiently solve the primal generalized eigenvalue problem even when the feature dimensionality is very high. The same trick could be applied to the dual formulation (see Table 1). Through document classification experiments, we show the effectiveness of the proposed method.

2 Dimensionality Reduction

In this section, we formulate the dimensionality reduction problem and review existing methods.

2.1 Formulation

Let \mathbf{x} ($\in \mathbb{R}^d$) be a d -dimensional sample and let \mathbf{z} ($\in \mathbb{R}^r$) be a low-dimensional representation of \mathbf{x} , where r ($1 \leq r \leq d$) is dimensionality of the reduced space. We consider linear dimensionality reduction scenarios where an embedded representation \mathbf{z} of a sample \mathbf{x} is obtained by using a $d \times r$ transformation matrix \mathbf{T} as

$$\mathbf{z} = \mathbf{T}^\top \mathbf{x},$$

where $^\top$ denotes the transpose of a matrix or a vector.

2.2 Principal Component Analysis (PCA)

The most well-known dimensionality reduction method would be PCA. Given samples $\{\mathbf{x}_i\}_{i=1}^n$, PCA seeks a transformation matrix \mathbf{T} such that the variance of the data samples in the embedding space is maximized. Let $\mathbf{S}^{(t)}$ be the *total scatter matrix*:

$$\begin{aligned} \mathbf{S}^{(t)} &:= \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top \\ &= \frac{1}{2} \sum_{i,j=1}^n W_{i,j}^{(t)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top, \end{aligned}$$

where $\boldsymbol{\mu} := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ is the mean of all samples and $\mathbf{W}^{(t)}$ is the $n \times n$ matrix with $W_{i,j}^{(t)} := 1/n$. Then a PCA solution is given by the leading eigenvectors of

$$\mathbf{S}^{(t)} \boldsymbol{\varphi} = \lambda \boldsymbol{\varphi}.$$

2.3 Fisher Discriminant Analysis (FDA) for Dimensionality Reduction

A popular supervised dimensionality reduction technique is *Fisher discriminant analysis* (FDA) [3]. Suppose that we have n' labeled samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^{n'}$, where $y_i (\in \{1, \dots, c\})$ is a class label associated with the sample \mathbf{x}_i and c is the number of classes. Let n'_m be the number of labeled samples in class m .

Let $\mathbf{S}^{(b)}$ and $\mathbf{S}^{(w)}$ be the *between-class scatter matrix* and the *within-class scatter matrix*:

$$\begin{aligned} \mathbf{S}^{(b)} &:= \sum_{m=1}^c n'_m (\boldsymbol{\mu}_m - \boldsymbol{\mu})(\boldsymbol{\mu}_m - \boldsymbol{\mu})^\top, \\ \mathbf{S}^{(w)} &:= \sum_{m=1}^c \sum_{i:y_i=m} (\mathbf{x}_i - \boldsymbol{\mu}_m)(\mathbf{x}_i - \boldsymbol{\mu}_m)^\top, \end{aligned}$$

where $\boldsymbol{\mu}_m := \frac{1}{n'_m} \sum_{i:y_i=m} \mathbf{x}_i$ is the mean of samples in class m .

FDA seeks a transformation matrix \mathbf{T} such that between-class scatter in the embedding space \mathbb{R}^r is maximized and within-class scatter in the embedding space minimized—a solution is given by the leading generalized eigenvectors of

$$\mathbf{S}^{(b)} \boldsymbol{\varphi} = \lambda \mathbf{S}^{(w)} \boldsymbol{\varphi}.$$

$\mathbf{S}^{(b)}$ and $\mathbf{S}^{(w)}$ are related to the total scatter matrix $\mathbf{S}^{(t)}$ as

$$\mathbf{S}^{(t)} = \mathbf{S}^{(b)} + \mathbf{S}^{(w)}.$$

This can also be confirmed from the fact that $\mathbf{S}^{(b)}$ and $\mathbf{S}^{(w)}$ are expressed in the pairwise form as follows [5]:

$$\mathbf{S}^{(b)} = \frac{1}{2} \sum_{i,j=1}^{n'} W_{i,j}^{(b)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top,$$

$$\mathbf{S}^{(w)} = \frac{1}{2} \sum_{i,j=1}^{n'} W_{i,j}^{(w)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top,$$

where $\mathbf{W}^{(b)}$ and $\mathbf{W}^{(w)}$ are the $n' \times n'$ matrices with

$$W_{i,j}^{(b)} := \begin{cases} 1/n' - 1/n'_{y_i} & \text{if } y_i = y_j, \\ 1/n' & \text{if } y_i \neq y_j, \end{cases}$$

$$W_{i,j}^{(w)} := \begin{cases} 1/n'_{y_i} & \text{if } y_i = y_j, \\ 0 & \text{if } y_i \neq y_j. \end{cases}$$

Then we have

$$\mathbf{W}^{(t)} = \mathbf{W}^{(b)} + \mathbf{W}^{(w)}.$$

2.4 Local Fisher Discriminant Analysis (LFDA)

LFDA [5] is a supervised dimensionality reduction technique which is a local extension of FDA.

Let $A_{i,j} (\in [0, 1])$ be the *affinity* between \mathbf{x}_i and \mathbf{x}_j ; $A_{i,j}$ is large if \mathbf{x}_i and \mathbf{x}_j are ‘close’ and $A_{i,j}$ is small if \mathbf{x}_i and \mathbf{x}_j are ‘far apart’. We assume that the affinity is symmetric, i.e., $A_{i,j} = A_{j,i}$. There are several different manners of defining the affinity; among them, we use the *local scaling heuristic* [7]:

$$A_{i,j} = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\sigma_i \sigma_j}\right).$$

$\sigma_i = \|\mathbf{x}_i - \mathbf{x}_i^{(k)}\|$ represents the local scaling around \mathbf{x}_i , where $\mathbf{x}_i^{(k)}$ is the k -th nearest neighbor of \mathbf{x}_i . A heuristic choice of $k = 7$ was shown to be useful through extensive simulations [7, 5].

Let $\mathbf{S}^{(lb)}$ and $\mathbf{S}^{(lw)}$ be the *local* between-class scatter matrix and the *local* within-class scatter matrix:

$$\mathbf{S}^{(lb)} := \frac{1}{2} \sum_{i,j=1}^{n'} W_{i,j}^{(lb)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top,$$

$$\mathbf{S}^{(lw)} := \frac{1}{2} \sum_{i,j=1}^{n'} W_{i,j}^{(lw)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top,$$

where $\mathbf{W}^{(\text{lb})}$ and $\mathbf{W}^{(\text{lw})}$ are the $n' \times n'$ matrices with¹

$$W_{i,j}^{(\text{lb})} := \begin{cases} A_{i,j}(1/n' - 1/n'_{y_i}) & \text{if } y_i = y_j, \\ 1/n' & \text{if } y_i \neq y_j, \end{cases}$$

$$W_{i,j}^{(\text{lw})} := \begin{cases} A_{i,j}/n'_{y_i} & \text{if } y_i = y_j, \\ 0 & \text{if } y_i \neq y_j. \end{cases}$$

When $A_{i,j} = 1$ for all i, j (i.e., no locality), the above scatter matrices are reduced to the ones used in original FDA. Thus, LFDA could be regarded as a ‘localized’ variant of FDA since the effect of data pairs in the same class but having small affinity values are deemphasized in the scatter matrices.

LFDA seeks a transformation matrix \mathbf{T} such that local between-class scatter in the embedding space is maximized and local within-class scatter in the embedding space is minimized—a solution is given by the leading generalized eigenvectors of

$$\mathbf{S}^{(\text{lb})}\boldsymbol{\varphi} = \lambda\mathbf{S}^{(\text{lw})}\boldsymbol{\varphi}.$$

LFDA is advantageous over original FDA in the following two respects. FDA is known to perform poorly when within-class multimodality or outliers exist [3]. On the other hand, LFDA can overcome this weakness by evaluating within-class scatter locally. The reduced dimensionality r is at most $c - 1$ in FDA [3], which is a critical limitation when the number of classes is small. On the other hand, LFDA can be generally applied to dimensionality reduction into *any* dimensional spaces, thanks to the affinity factor $A_{i,j}$.

2.5 Semi-supervised LFDA (SELF)

LFDA (and any other supervised methods) tends to suffer from overfitting when the number of samples is small. For mitigating the overfitting problem, we assume the availability of unlabeled samples; from here on, we consider the *semi-supervised* setup [1] where, in addition to the labeled samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^{n'}$, unlabeled samples $\{\mathbf{x}_i\}_{i=n'+1}^n$ are available.

With the help of unlabeled samples, the overfitting problem of LFDA could be mitigated by combining it with PCA—SELF smoothly bridges LFDA and PCA [6]. More specifically, a SELF solution is given by the leading generalized eigenvectors of

$$\mathbf{S}^{(\text{rlb})}\boldsymbol{\varphi} = \lambda\mathbf{S}^{(\text{rlw})}\boldsymbol{\varphi}, \tag{1}$$

where $\mathbf{S}^{(\text{rlb})}$ and $\mathbf{S}^{(\text{rlw})}$ are *regularized* local between-class scatter matrix and *regularized*

¹A more compact expression would be

$$W_{i,j}^{(\text{lb})} = \bar{A}_{i,j}/n' - W_{i,j}^{(\text{lw})},$$

$$W_{i,j}^{(\text{lw})} = \delta_{y_i, y_j} \bar{A}_{i,j}/n'_{y_i},$$

where $\delta_{i,j}$ denotes Kronecker’s delta and $\bar{A}_{i,j} = A_{i,j}$ if $y_i = y_j$ and $\bar{A}_{i,j} = 1$ otherwise.

local within-class scatter matrix defined by

$$\begin{aligned}\mathbf{S}^{(\text{rlb})} &:= (1 - \beta)\mathbf{S}^{(\text{lb})} + \beta\mathbf{S}^{(\text{t})}, \\ \mathbf{S}^{(\text{rlw})} &:= (1 - \beta)\mathbf{S}^{(\text{lw})} + \beta\mathbf{I}_d.\end{aligned}$$

\mathbf{I}_d is the d -dimensional identity matrix. β ($\in [0, 1]$) is a trade-off parameter that controls the ‘degree’ of unsupervisedness—when $\beta = 0$, SELF is fully supervised and is reduced to LFDA; when $\beta = 1$, SELF is fully unsupervised and is reduced to PCA; otherwise, SELF inherits the properties of both LFDA and PCA.

It is practically useful to down-weight the eigenvectors according to their associated eigenvalues [6]; more specifically, the transformation matrix is given by

$$\mathbf{T}^{(\text{SELF})} = (\sqrt{\lambda_1}\boldsymbol{\varphi}_1 | \cdots | \sqrt{\lambda_r}\boldsymbol{\varphi}_r),$$

where $\{\boldsymbol{\varphi}_k\}_{k=1}^d$ are the generalized eigenvectors associated with the generalized eigenvalues $\lambda_1 \geq \cdots \geq \lambda_d$ of Eq.(1); $\{\boldsymbol{\varphi}_k\}_{k=1}^d$ are normalized as

$$\boldsymbol{\varphi}_k^\top \mathbf{S}^{(\text{rlw})} \boldsymbol{\varphi}_k = 1 \quad \text{for } k = 1, \dots, d.$$

SELF was shown to be useful even when only a small number of labeled samples are available [6].

3 Computational Issues of SELF for Large-scale and High-dimensional Datasets

Here, we address computational issues of SELF.

3.1 Primal Formulation for Moderate Dimensionality

Dimensionality of the generalized eigenvalue problem (1) depends only on dimensionality of feature vectors, not on the number of samples. Thus SELF may be applicable to a dataset with a large number of samples as long as dimensionality of the feature vectors is moderate. However, when the feature dimensionality is very high, solving the generalized eigenvalue problem may not be computationally tractable—even more critically, the matrix which needs to be eigendecomposed cannot be stored in memory. This critical limitation makes SELF inapplicable to high-dimensional problems.

In the rest of this section, we show how the computational problem of SELF could be alleviated.

3.2 Dual Formulation for Moderate Sample Size

When d is very large but n is moderate, the SELF solution can be computed efficiently using the dual formulation.

For obtaining the dual problem, let us first express the primal problem (1) in terms of the sample matrix

$$\mathbf{X} = (\mathbf{x}_1 | \cdots | \mathbf{x}_n)^\top.$$

Let $\mathbf{W}^{(\text{rlb})}$ be the $n \times n$ matrix with

$$W_{i,j}^{(\text{rlb})} := \begin{cases} (1 - \beta)A_{i,j}(1/n' - 1/n'_{y_i}) + \beta/n & \text{if } y_i = y_j, \\ (1 - \beta)/n' + \beta/n & \text{if } y_i \neq y_j, \\ \beta/n & \text{otherwise,} \end{cases}$$

and let $\mathbf{D}^{(\text{rlb})}$ and $\mathbf{D}^{(\text{lw})}$ be the $n \times n$ diagonal matrices with

$$D_{i,i}^{(\text{rlb})} := \sum_{j=1}^n W_{i,j}^{(\text{rlb})},$$

$$D_{i,i}^{(\text{lw})} := \sum_{j=1}^n W_{i,j}^{(\text{lw})}.$$

Let

$$\mathbf{L}^{(\text{rlb})} := \mathbf{D}^{(\text{rlb})} - \mathbf{W}^{(\text{rlb})},$$

$$\mathbf{L}^{(\text{rlw})} := (1 - \beta)(\mathbf{D}^{(\text{lw})} - \mathbf{W}^{(\text{lw})}) + \beta(\mathbf{X}^\top \mathbf{X})^\dagger,$$

where \dagger denotes the Moore-Penrose generalized inverse. Then the primal problem (1) can be expressed as follows [6]²:

$$\mathbf{X} \mathbf{L}^{(\text{rlb})} \mathbf{X}^\top \boldsymbol{\varphi} = \lambda \mathbf{X} \mathbf{L}^{(\text{rlw})} \mathbf{X}^\top \boldsymbol{\varphi}. \quad (2)$$

Since $\mathbf{X}^\top \boldsymbol{\varphi}$ belongs to the range of \mathbf{X}^\top , it can be expressed by using some vector $\boldsymbol{\alpha}$ ($\in \mathbb{R}^n$) as

$$\mathbf{X}^\top \boldsymbol{\varphi} = \mathbf{X}^\top \mathbf{X} \boldsymbol{\alpha} = \mathbf{K} \boldsymbol{\alpha},$$

where \mathbf{K} is the $n \times n$ matrix with

$$K_{i,j} := \mathbf{x}_i^\top \mathbf{x}_j.$$

Then multiplying Eq.(2) by \mathbf{X}^\top from the left-hand side yields the following dual eigenvalue problem:

$$\mathbf{K} \mathbf{L}^{(\text{rlb})} \mathbf{K} \boldsymbol{\alpha} = \lambda \mathbf{K} \mathbf{L}^{(\text{rlw})} \mathbf{K} \boldsymbol{\alpha}. \quad (3)$$

Note that $\mathbf{K} \mathbf{L}^{(\text{rlw})} \mathbf{K}$ has a simpler expression:

$$\mathbf{K} \mathbf{L}^{(\text{rlw})} \mathbf{K} = (1 - \beta) \mathbf{K} (\mathbf{D}^{(\text{lw})} - \mathbf{W}^{(\text{lw})}) \mathbf{K} + \beta \mathbf{K},$$

²This eigenvalue problem would be indeterminate due to rank deficiency. However, this is not a problem since we only work in the range of \mathbf{X}^\top .

so the Moore-Penrose generalized inverse of $\mathbf{X}^\top \mathbf{X}$ does not have to be computed. Let $\{\boldsymbol{\alpha}_k\}_{k=1}^d$ be the eigenvectors associated with the eigenvalues $\lambda_1 \geq \dots \geq \lambda_d$ of Eq.(3), which are normalized as

$$\boldsymbol{\alpha}_k^\top \mathbf{K} \mathbf{L}^{(\text{rlw})} \mathbf{K} \boldsymbol{\alpha}_k = 1 \quad \text{for } k = 1, \dots, d.$$

Then the embedded representation \mathbf{z} of a sample \mathbf{x} can be computed in terms of $\{\boldsymbol{\alpha}_k\}_{k=1}^d$ as

$$\mathbf{z} = (\sqrt{\lambda_1} \boldsymbol{\alpha}_1 | \dots | \sqrt{\lambda_d} \boldsymbol{\alpha}_d)^\top (\mathbf{x}_1^\top \mathbf{x}, \dots, \mathbf{x}_d^\top \mathbf{x})^\top.$$

Since Eq.(3) is an n -dimensional eigenvalue problem which is independent of d , it may be solved efficiently even when d is large as long as n is moderate.

Note that the dual formulation does not directly include the feature vectors, but only through their inner product $\mathbf{x}^\top \mathbf{x}'$. This means that SELF can be non-linearized by replacing the inner product $\mathbf{x}^\top \mathbf{x}'$ with a reproducing kernel $K(\mathbf{x}, \mathbf{x}')$ [6]. Furthermore, the use of kernels for structured data such as sequences, trees, and graphs allows us to reduce dimensionality in the structured domain.

3.3 Primal Formulation with Sparse Data

If \mathbf{X} is sparse (i.e., only a small number of elements are non-zero), we may solve the primal eigenvalue problem efficiently even when n and d are both large.

$\mathbf{S}^{(\text{rlb})}$ and $\mathbf{S}^{(\text{rlw})}$ can be expressed as

$$\begin{aligned} \mathbf{S}^{(\text{rlb})} &= \mathbf{X} \mathbf{D}^{(\text{rlb})} \mathbf{X}^\top - \mathbf{X} \mathbf{W}^{(\text{rlb})} \mathbf{X}^\top, \\ \mathbf{S}^{(\text{rlw})} &= (1 - \beta) \mathbf{X} \mathbf{D}^{(\text{lw})} \mathbf{X}^\top - (1 - \beta) \mathbf{X} \mathbf{W}^{(\text{lw})} \mathbf{X}^\top + \beta \mathbf{I}_d. \end{aligned}$$

If \mathbf{X} is sparse, then $\mathbf{S}^{(\text{rlw})}$ is also sparse since $\mathbf{D}^{(\text{lw})}$, $\mathbf{W}^{(\text{lw})}$, and \mathbf{I}_d are all sparse. On the other hand, $\mathbf{D}^{(\text{rlb})}$ is sparse but $\mathbf{W}^{(\text{rlb})}$ is dense, so $\mathbf{S}^{(\text{rlb})}$ is not sparse. However, $\mathbf{S}^{(\text{rlb})}$ has nice structure as explained below.

We can show that $\mathbf{W}^{(\text{rlb})}$ is expressed as

$$\mathbf{W}^{(\text{rlb})} = \mathbf{H} + \{(1 - \beta)/n'\} \boldsymbol{\eta} \boldsymbol{\eta}^\top + (\beta/n) \mathbf{1} \mathbf{1}^\top, \quad (4)$$

where \mathbf{H} is the $n \times n$ matrix with

$$H_{i,j} := \begin{cases} (1 - \beta) \{A_{i,j}(1/n' - 1/n'_{y_i}) - 1/n'\} & \text{if } y_i = y_j, \\ 0 & \text{otherwise.} \end{cases}$$

$\boldsymbol{\eta}$ is the n -dimensional vector with $\eta_i := 1$ for $1 \leq i \leq n'$ (i.e., \mathbf{x}_i is labeled) and $\eta_i := 0$ otherwise. $\mathbf{1}$ denotes the vector with all ones. Since the number of labeled samples is usually small in the semi-supervised setup (i.e., $n' \ll n$), \mathbf{H} would be sparse; furthermore, if samples are sorted according to the labels, \mathbf{H} becomes block-diagonal. Let us denote

the m -th ‘diagonal-block’ matrix by \mathbf{H}_m , which are dense but small. Then $\mathbf{S}^{(\text{rlb})}$ can be expressed as

$$\begin{aligned} \mathbf{S}^{(\text{rlb})} = & \mathbf{X} \mathbf{D}^{(\text{rlb})} \mathbf{X}^\top - \sum_{m=1}^c \mathbf{X}_m \mathbf{H}_m \mathbf{X}_m^\top \\ & - (1 - \beta) n' \boldsymbol{\xi} \boldsymbol{\xi}^\top - \beta n \boldsymbol{\mu} \boldsymbol{\mu}^\top, \end{aligned}$$

where $\boldsymbol{\xi} := \frac{1}{n'} \sum_{i=1}^{n'} \mathbf{x}_i$ is the mean of labeled samples; \mathbf{X}_m is a matrix of samples in class m , which is sparse. Thus $\mathbf{S}^{(\text{rlb})}$ is the sum of a sparse matrix and low-rank matrices. This structure is highly useful since $\mathbf{S}^{(\text{rlb})} \boldsymbol{\varphi}$ can be computed efficiently—for example, the *eigs* function in MATLAB[®] that sequentially solves a sparse eigenvalue problem can take advantage of this structure for improving computational efficiency.

3.4 Dual Formulation with Sparse Kernel

The above computational trick could be applied also to the dual formulation (3). Thus, when the kernel matrix \mathbf{K} is sparse, the dual eigenvalue problem may be solved efficiently even when both n and d are large.

4 Experiments

We have made SELF applicable to high-dimensional problems. Here, we apply it to real-world document classification tasks and evaluate its performance.

We use the *Technion Repository of Text Categorization*³ (TechTC), which contains 295 binary document classification tasks. Each task contains a few hundred documents with category labels. We use *term frequency/inverse document frequency* (TFIDF) features for classification [4], whose dimensionality ranges from thousands to tens of thousands (their dimensionality varies depending on the task since we removed entries which are zero for all documents in the dataset). The number of document samples is relatively small in this dataset, so we use the dual formulation for computing the SELF solution. Note that the original primal formulation cannot be employed for this experiment due to its high dimensionality.

We compare the performance of ‘Plain’ (without dimensionality reduction), LFDA, PCA, and SELF with $\beta = 0.5$. In each method, dimensionality of the reduced space r is chosen by 5-fold cross-validation from $\{1, \dots, 10\}$. For each dataset, we consider 3 configurations with different degrees of semi-supervisedness; given n document samples, we randomly choose 10%, 50%, 90% of them as labeled training samples and the rest are treated as unlabeled samples. We employ the 1-nearest neighbor method for evaluating the classification accuracy of the unlabeled samples. For each dataset and each training

³The dataset is available from ‘<http://techtc.cs.technion.ac.il/techtc300/techtc300.html>’. See [2] for its specification.

Table 2: Document classification results. The mean and standard deviation of the misclassification rate over 295 datasets for 100 runs are described in the row of ‘Error’. For each dataset, the best method and comparable ones in terms of the mean misclassification rate over 100 runs based on the t -test at the significance level 5% are described in bold face. The number of times the corresponding method is judged to be the best over 295 datasets is described in the row of ‘Bests’. The mean value of the reduced dimensionality r chosen by 5-fold cross validation is described in the row of ‘ r ’.

n'/n		Plain	LFDA	SELF	PCA
0.1	Error	21.8 (2.6)	23.9 (2.6)	19.4 (2.4)	20.4 (2.0)
	Bests	91/295	17/295	239/295	122/295
	r	—	2.0	2.6	3.7
0.5	Error	21.0 (3.6)	14.3 (1.5)	13.2 (1.5)	16.5 (1.6)
	Bests	10/295	125/295	230/295	43/295
	r	—	3.3	4.1	5.5
0.9	Error	21.7 (3.4)	13.5 (2.9)	11.8 (3.2)	15.4 (3.5)
	Bests	15/295	149/295	253/295	98/295
	r	—	3.9	4.1	6.1

sample configuration, the experiments are repeated 100 times by changing the random choice of training samples.

The experimental results are summarized in Table 2. The table shows that all the dimensionality reduction methods tend to perform better than ‘Plain’ except LFDA for $n' = 0.1n$ due to strong overfitting. So dimensionality reduction overall contributes to improving the accuracy of document classification. The performance of LFDA is improved as the number of labeled samples increases, while the accuracy of ‘Plain’ does not improve. The performance of PCA also improves as the number of labeled samples increases; this is counter-intuitive at a glance since PCA is unsupervised. However, the labeled samples are used for choosing the reduced dimensionality r by cross-validation, so the performance tends to be improved if a large number of labeled samples are available. But the performance improvement of LFDA as the number of labeled samples is increased is more prominent than PCA, thanks to the supervised formulation.

SELF consistently outperforms LFDA and PCA for all cases, showing that SELF can effectively use information brought by unlabeled samples and improve the classification accuracy. We also tested SELF with β chosen by cross-validation; the results were almost the same as SELF with $\beta = 0.5$ (so we omit the detail), but it required more computation time.

Overall, combining LFDA and PCA by SELF is shown to be a useful dimensionality reduction method in practical document classification tasks; SELF was made applicable to such high-dimensional problems by the computational tricks introduced in this paper.

5 Conclusions

Accurately classifying high-dimensional patterns is an important challenge in pattern recognition. A semi-supervised dimensionality reduction method called SELF was demonstrated to work excellently, but its naive implementation does not allow us to apply SELF to high-dimensional problems due to high computational costs and large memory requirement. In this paper, we introduced computational tricks for making SELF applicable to high-dimensional problems. We demonstrated the usefulness of the proposed method through high-dimensional document classification simulations.

References

- [1] O. Chapelle, B. Schölkopf, and A. Zien, eds., *Semi-Supervised Learning*, MIT Press, Cambridge, 2006.
- [2] D. Davidov, E. Gabrilovich, and S. Markovitch, “Parameterized generation of labeled datasets for text categorization based on a hierarchical directory,” *The 27th Annual International ACM SIGIR Conference*, Sheffield, UK, pp.250–257, Jul. 25–29 2004.
- [3] K. Fukunaga, *Introduction to Statistical Pattern Recognition*, second ed., Academic Press, Inc., Boston, 1990.
- [4] T. Joachims, *Learning to Classify Text Using Support Vector Machines: Methods, Theory and Algorithms*, Kluwer Academic Publishers, Boston, 2002.
- [5] M. Sugiyama, “Dimensionality reduction of multimodal labeled data by local Fisher discriminant analysis,” *Journal of Machine Learning Research*, vol.8, pp.1027–1061, May 2007.
- [6] M. Sugiyama, T. Ide, S. Nakajima, and J. Sese, “Semi-supervised local Fisher discriminant analysis for dimensionality reduction,” *Advances in Knowledge Discovery and Data Mining*, ed. T. Washio, E. Suzuki, K.M. Ting, and A. Inokuchi, *Lecture Notes in Computer Science*, vol.5012, Berlin, pp.333–344, Springer, 2008.
- [7] L. Zelnik-Manor and P. Perona, “Self-tuning spectral clustering,” in *Advances in Neural Information Processing Systems 17*, ed. L.K. Saul, Y. Weiss, and L. Bottou, pp.1601–1608, MIT Press, Cambridge, MA, 2005.