



ELSEVIER

Contents lists available at ScienceDirect

Pattern Recognition

journal homepage: www.elsevier.com/locate/pr

A Rayleigh–Ritz style method for large-scale discriminant analysis



Lin Zhu*, De-Shuang Huang

School of Electronics and Information Engineering, Tongji University, 4800 Caoan Road, Shanghai 201804, China

ARTICLE INFO

Article history:

Received 5 September 2012
 Received in revised form
 6 September 2013
 Accepted 5 October 2013
 Available online 20 October 2013

Keywords:

Dimensionality reduction
 Linear Discriminant Analysis (LDA)
 Generalized Eigenvalue Decomposition (GEVD)
 Rayleigh–Ritz procedure

ABSTRACT

Linear Discriminant Analysis (LDA) is one of the most popular approaches for supervised feature extraction and dimension reduction. However, the computation of LDA involves dense matrices eigendecomposition, which is time-consuming for large-scale problems. In this paper, we present a novel algorithm called Rayleigh–Ritz Discriminant Analysis (RRDA) for efficiently solving LDA. While much of the prior research focus on transforming the generalized eigenvalue problem into a least squares formulation, our method is instead based on the well-established Rayleigh–Ritz framework for general eigenvalue problems and seeks to directly solve the generalized eigenvalue problem of LDA. By exploiting the structures in LDA problems, we are able to design customized and highly efficient subspace expansion and extraction strategy for the Rayleigh–Ritz procedure. To reduce the storage requirement and computational complexity of RRDA for high dimensional, low sample size data, we also establish an equivalent reduced model of RRDA. Practical implementations and the convergence result of our method are also discussed. Our experimental results on several real world data sets indicate the performance of the proposed algorithm.

© 2013 Elsevier Ltd. All rights reserved.

1. Introduction

With the advancements in data collection and storage technologies, there has been an exponential increase in the availability and usage of large, high-dimensional datasets. These data can directly be represented as vectors in high-dimensional vector spaces. Obviously, operating directly on such high-dimensional vector space is ineffective and may lead to high computational and storage demands as well as poor performance. A typical way to circumvent the “curse of dimensionality” problem [1] and other undesired properties of high dimensional spaces is to use dimensionality reduction techniques. The goal of dimensionality reduction is to map a set of high-dimensional samples into a lower dimensional space while preserving the intrinsic structure in the data.

Up to now, researchers have developed a variety of dimensionality reduction methods. Based on how to utilize the label information, these algorithms can be broadly divided into three classes, i.e., unsupervised, supervised and semi-supervised. The first class is supervised which typically includes linear discriminant analysis (LDA) [2,3], maximization of the geometric mean of all divergences (MGMD) [4], max–min distance analysis (MMDA) [5], etc. The second class is unsupervised which include principal component analysis

(PCA) [6] and locality preserving projections (LPP) [7,8], etc. The third class is semi-supervised [9,10] which can use unlabeled data for promoting supervised methods. Beyond the commonness in mathematical formulation [11] shared by these algorithms, Zhang et al. proved that most popular dimensionality reduction algorithms, unsupervised or supervised, can be explained as instances of a ubiquitously two-stage framework named “patch alignment” [12].

One of the most popular supervised dimensionality reduction algorithms is linear discriminant analysis (LDA). It has been widely used in many applications such as microarray data classification; face recognition and gait recognition, etc [2,13–15].

Computationally, LDA amounts to solving a generalized symmetric semi-definite eigenvalue decomposition (GEVD) problem. A straightforward implementation can thus be very time-consuming for large datasets [16]. So far a number of methods [3,16–22] have been explored to solve LDA with improved scalability. The key idea of these approaches is to transform the generalized eigenvalue problem into a least squares formulation, which can be solved efficiently using existing algorithms such as LSQR [21,23]. Two representative algorithms are least squares linear discriminant analysis (LS-LDA) [18,19] and spectral regression discriminant analysis (SRDA) [16]. Recently, Sun introduced 2SrLDA [21], a generalization of LS-LDA which relaxes the equivalence conditions and incorporates a regularization term.

Alternatively, one can use general-purpose GEVD algorithms to solve LDA. It is known that for solving small to medium size GEVD

* Corresponding author. Tel.: +86 551 3493367.

E-mail addresses: zhomlynn@ustc.edu.cn, lizhonyx@163.com (L. Zhu).

problems, ‘Rayleigh–Ritz’ methods like Lanczos method and JDQZ are often the ideal choice [2,13–15]. In [21], Sun attempted to solve LDA based on Lanczos method. However, the experiments comparisons show that it is much slower than regression-based methods. This is not surprising because a standard GEVD solver would encounter severe computational bottleneck or even become impractical for large problems. For example, Rayleigh–Ritz type GEVD algorithms generally cannot guarantee finding the largest eigenpairs (which is required by LDA) via short-term recurrences, therefore more sophisticated strategy like trust-reign method has to be employed [24], which can make the computation inefficient. In addition, for symmetric eigenvalue problems, Rayleigh–Ritz methods like Lanczos algorithm can guarantee that the approximate solution at each iteration be computed efficiently using divide-and-conquer method [25]. However this advantage also doesn’t exist for general GEVD methods.

In this paper, we propose a novel Rayleigh–Ritz algorithm for LDA called Rayleigh–Ritz discriminant analysis (RRDA). Like JDQZ and Lanczos method, we adopt the ‘Rayleigh–Ritz’ procedure as a general framework. However, the subspace expansion and extraction strategy is specially designed to gradually optimize the ratio trace objective function of LDA. By exploiting the special structure of LDA problem, the above-mentioned issues of general purpose GEVD methods can also be solved. We demonstrate both theoretically and via numerical examples that this hybrid algorithm efficiently produces a solution of high precision.

Additionally, in many applications, data are highly dimensional, while the number of available training samples is usually much smaller. Such is the case for the image databases of facial recognition, gene expression data, as well as the text documents [26,27]. To reduce the storage requirement and computational complexity of the iteration for these high dimensional, low sample size data, we also establish an equivalent reduced model of RRDA which is of order n instead of order d , where d is the number of features and n is the number of samples. Practical implementations and the convergence result of RRDA are also discussed. Experiments show that RRDA is more efficient than 2SrLDA and SRDA, which are the state-of-the-art methods for solving LDA.

The rest of this paper is organized as follows. Section 2 outlines LDA and the Rayleigh–Ritz procedure. The proposed method is discussed in Section 3. Experimental results on several real world data sets are reported in Section 4. Finally, some concluding remarks are drawn in Section 5.

2. Background

In preparation for our description of the proposed method, we introduce some notations. Throughout this paper all matrices are boldface uppercase, and vectors are boldface lowercase. n is the number of samples, d is the data dimensionality, and c is the number of classes (or labels). The i th sample is denoted as $\mathbf{x}_i \in \mathbb{R}^{d \times 1}$, and $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ represents the data matrix. Without loss of generality, we assume that \mathbf{X} is partitioned into c classes as $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_c]$, where $\mathbf{X}_i \in \mathbb{R}^{d \times n_i}$ corresponds to the data points from the i th class. \mathbf{I}_q is the q -by- q identity matrix and \mathbf{e}_q is a vector of all ones with length q . The centering matrix is defined as $\mathbf{C}_q = \mathbf{I}_q - \mathbf{e}_q \mathbf{e}_q^T / q$. For a matrix \mathbf{A} , we let \mathbf{A}^+ be the Moore–Penrose inverse of \mathbf{A} , $\text{trace}(\mathbf{A})$ be the trace of \mathbf{A} , $\text{rank}(\mathbf{A})$ be the rank of \mathbf{A} and $\text{span}(\mathbf{A})$ be the range space of \mathbf{A} .

2.1. Linear discriminant analysis

The solution of LDA optimizes the ratio trace problem [21,28,29]:

$$\arg \max_{\mathbf{W}^T \mathbf{S}_t \mathbf{W} = \mathbf{I}} \text{trace}[(\mathbf{W}^T \mathbf{S}_t \mathbf{W})^{-1} \mathbf{W}^T \mathbf{S}_b \mathbf{W}] \quad (1)$$

where the total covariance matrix \mathbf{S}_t and the between-class covariance matrix \mathbf{S}_b are defined as

$$\begin{aligned} \mathbf{S}_t &= \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T \\ \mathbf{S}_b &= \frac{1}{n} \sum_{k=1}^c n_k (\boldsymbol{\mu}_k - \boldsymbol{\mu})(\boldsymbol{\mu}_k - \boldsymbol{\mu})^T \end{aligned} \quad (2)$$

where $\boldsymbol{\mu}_i = (1/n_i) \mathbf{X}_i \mathbf{e}_{n_i}$ denotes the centroid of the class i and $\boldsymbol{\mu} = (1/n) \mathbf{X} \mathbf{e}_n$ denotes the global centroid.

The optimization problem (1) is solved by computing all the generalized eigenpairs [16]:

$$\lambda_i \mathbf{S}_t \mathbf{w}_i = \mathbf{S}_b \mathbf{w}_i, \quad \mathbf{w}_i \neq \mathbf{0}, \quad \lambda_i > 0 \quad (3)$$

and thus, the optimal \mathbf{W} consists of eigenvectors of $\mathbf{S}_t^{-1} \mathbf{S}_b$ corresponding to all positive eigenvalues, provided that \mathbf{S}_t is nonsingular.

In practical application, to deal with the singularity of \mathbf{S}_t , many generalizations of LDA have been proposed [15,17,27,28,30–39]. One widely used approach is called Regularized Discriminant Analysis (RDA) [33], which penalizes the total scatter matrix \mathbf{S}_t with a regularizer $\alpha \mathbf{I}_d$, $\alpha > 0$, and the solution is given by solving the GEVD problem

$$\lambda_i (\mathbf{S}_t + \alpha \mathbf{I}_d) \mathbf{w}_i = \mathbf{S}_b \mathbf{w}_i, \quad \mathbf{w}_i \neq \mathbf{0}, \quad \lambda_i > 0 \quad (4)$$

Since LDA is a special case of RDA with $\alpha=0$, we mainly discuss RDA in the sequel.

2.2. Rayleigh–Ritz method for GEVD problems

The Rayleigh–Ritz procedure serves as the basic framework of many algorithms for solving linear and nonlinear eigenvalue problems. Notable examples include the Lanczos method [40], JDQZ (including JD-related variants) [41–43] and GY [44].

In order to find the generalized eigenpairs of a given matrix pair $(\mathbf{S}_1, \mathbf{S}_2)$, the Rayleigh–Ritz procedure selects the approximate eigenvector from a search subspace $\text{span}(\mathbf{V}_k)$ that is expanded iteratively. Each iteration consists of two parts. In the first part, the ‘extraction’ part, the projected generalized eigenproblem

$$\mathbf{V}_k^T \mathbf{S}_1 \mathbf{V}_k \boldsymbol{\eta} = \lambda \mathbf{V}_k^T \mathbf{S}_2 \mathbf{V}_k \boldsymbol{\eta} \quad (5)$$

is solved and a solution $(\lambda, \boldsymbol{\eta})$ is selected. The Ritz value λ and Ritz vector $\mathbf{V}_k \boldsymbol{\eta}$ form an approximate generalized eigenvalue and generalized eigenvector, respectively. In the second part, the ‘expansion’ part, the search space $\text{span}(\mathbf{V}_k)$ is enlarged by adding a new basis matrix \mathbf{P}_k to it. This process is summarized in Algorithm 1. The idea is that, as the search subspace grows, the eigenpair approximations will converge to an eigenpair of the original problem. In order to keep computation costs low, we usually do not expand the search space to the whole space.

Algorithm 1. The Rayleigh–Ritz framework for solving GEVD Problems

- 1) **Start:** Matrix pair $(\mathbf{S}_1, \mathbf{S}_2)$, choose an initial non-zero matrix \mathbf{V}_1 .
- 2) **Output:** \mathbf{W}_{opt} as the eigenvectors of $(\mathbf{S}_1, \mathbf{S}_2)$ with positive eigenvalues.
- 3) **Iterate:** Until convergence, for $k=1, 2, \dots$ do:
 - 3.1) **Expansion:** Choose a subspace basis \mathbf{P}_k based on certain criteria, expand \mathbf{V}_{k-1} with \mathbf{P}_k to \mathbf{V}_k such that $\mathbf{V}_k = [\mathbf{V}_{k-1} \ \mathbf{P}_k]$.
 - 3.2) **Extraction:** Compute the positive eigenpairs $(\lambda_i, \boldsymbol{\eta}_i)$ of $(\mathbf{V}_k^T \mathbf{S}_1 \mathbf{V}_k, \mathbf{V}_k^T \mathbf{S}_2 \mathbf{V}_k)$ and the Ritz vectors $\mathbf{w}_i = \mathbf{V}_k \boldsymbol{\eta}_i$, $i=1, 2, \dots, c$. Let $\mathbf{W}_k = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots]$;
 - 3.3) Test for convergence. Stop if satisfied.
- 4) **End:** Output \mathbf{W}_k as the approximation of \mathbf{W}_{opt} .

Download English Version:

<https://daneshyari.com/en/article/532070>

Download Persian Version:

<https://daneshyari.com/article/532070>

[Daneshyari.com](https://daneshyari.com)