

Boosting Performance of Incremental IDR/QR LDA - from Sequential to Chunk

Yiming Peng

A thesis submitted to Auckland University of Technology
in fulfillment of the requirements
for the degree of Master of Computer and Information Sciences

July, 2011



School of Computing and Mathematical Sciences

Primary Supervisor: Prof. Alvis Fong
Secondary Supervisor: Dr. Shaoning Pang

Declaration

I hereby declare that this submission is my own work and that, to the best of my knowledge and belief, it contains no material previously published or written by another person nor material which to a substantial extent has been accepted for the qualification of any other degree or diploma of a University or other institution of higher learning, except where due acknowledgment is made in the acknowledgments.

Printed name: Yiming Peng

Signature:

Date:

Abstract

Training data in the real world is often presented in random chunks. Yet existing sequential incremental IDR/QR LDA (sIncLDA) can only process data one instance after another. This thesis proposes a new chunk incremental IDR/QR LDA (cIncLDA) capable of processing multiple data instances at one time. sIncLDA updates the reduced within-class scatter matrix \mathbf{W} by a QR decomposition of the centroid matrix for each newly-arrived data instance. It is assumed that the updated $\mathbf{Q}' \approx \mathbf{Q}$ for any data instance from an existing class and the updated $\mathbf{W}' \approx \mathbf{W}$ for any data instance from a new class. In practice, the assumption in sIncLDA leads to significant loss of the discriminative information from approximating \mathbf{Q} and \mathbf{W} when the number of classes is large. By utilizing a new method that accurately updates \mathbf{W} , the proposed cIncLDA can better preserve the discriminative information contained in \mathbf{W} . The limitation of sIncLDA is hence resolved. Experimental comparisons have been conducted on six facial datasets with diverse class numbers ranging from 40 to 1010. The result indicates that our algorithm achieves a competitive accuracy to batch QR/LDA and is consistently higher than sIncLDA. It is noted in the report that the computational complexity of our algorithm is more expensive than sIncLDA for single data processing (i.e., sequential manner); however, the efficiency of our algorithm surpasses sIncLDA as the chunk size increases for multiple instances processing (i.e., chunk manner).

Acknowledgments

I would like to thank all people who have helped and inspired me during my master study.

I especially want to thank my supervisors, Prof. Alvis Fong and Dr Shaoning Pang, for their guidance during my research and study at Auckland University of Technology. Their perpetual energy and enthusiasm for research have motivated all their students, including me. In addition, Prof. Alvis Fong and Dr. Shaoning Pang were always accessible and willing to help his students with their research. As a result, research life became smooth and rewarding for me. Also, many thanks to the professional editor, Heather Moodie, for correcting my thesis carefully and patiently.

All my lab buddies at the DMLI Lab made it a convivial place to work. In particular, I would like to thank Lei Song and Lei Zhu for their friendship and help during my thesis.

My deepest gratitude goes to my family for their unflagging love and support throughout my life; this thesis would be simply impossible without them. I am indebted to my father, Sige Peng, for his care and love. As a typical father, he worked industriously to support the family and spared no effort to provide the best possible environment for me to grow up and attend school. He has never complained in spite of all the hardships in his life. I cannot ask for more from my mother, Qi Zhang, as she is simply perfect. I have no suitable word that can fully describe her everlasting love for me. I remember her constant support when I encountered difficulties and I remember, most of all, her delicious dishes. Also, many appreciations to my brother Lin Peng and his wife Li Chen, for their encourage, help and love. I would also like to thank my girlfriend Yuanhua Sun for supporting me while I was writing the thesis and for ensuring that I did not lose sight of other aspects in life not related to academia.

Contents

1	Introduction	1
1.1	Background	1
1.2	Research Objectives	4
1.3	Research Contributions	4
1.4	Thesis Structure	5
2	Literature Review	6
2.1	Linear Discriminant Analysis	6
2.1.1	Classic Fisher LDA	7
2.1.2	Small Sample Size Problem	8
2.1.3	LDA Extensions to the SSS Problem	9
2.1.4	Summary	12
2.2	Incremental Linear Discriminant Analysis	12
2.2.1	Chunk Incremental Learning	12
2.2.2	One-pass Incremental Learning	13
2.2.3	Incremental Linear Discriminant Analysis	14
2.2.4	Summary	22
2.3	Motivations for the Presented Research	24
3	Existing Sequential Incremental IDR/QR LDA	26
3.1	Batch QR/LDA	26
3.2	Sequential Incremental IDR/QR LDA	27
3.3	Discussion	29
4	Proposed Chunk Incremental IDR/QR LDA	30
4.1	Updating the Centroid Matrix and its QR-decomposition	30
4.2	Updating the Reduced Within-Class Scatter Matrix	32

4.3	Updating the Reduced Between-Class Scatter Matrix	34
4.4	The Pseudocode of the Proposed Algorithm	34
4.5	Time Complexity Analysis	34
5	Experiment Evaluation	38
5.1	Data Description	38
5.2	Experimental setup	39
5.3	Results Evaluation	40
5.3.1	Similarity to Ground Truth Eigenspace	40
5.3.2	Class Separability	40
5.3.3	Computational Efficiency	43
5.4	Summary	46
6	Conclusions and Directions for Future Research	47
6.1	Conclusions	47
6.2	Directions for Future Research	48
A	Proof of Proposition 1	50
B	Proof of Proposition 2	52
C	Proof of Proposition 3	54
	References	56

Chapter 1

Introduction

1.1 Background

Incremental learning (IL) addresses real-world situations when data are frequently presented in either sequential (i.e., one sample after another) or chunk manner (i.e., a subset of samples presented at a time) (Giraud-Carrier, 2000; Pang, Ozawa & Kasabov, 2005b). Given an existing learning model Ω , the incremental learning of Ω is an updating process as follows:

$$\Omega' = \mathcal{F}_{sq}(\Omega, \tilde{\mathbf{x}}),$$

or

$$\Omega' = \mathcal{F}_{ck}(\Omega, \tilde{\mathbf{X}}), \tilde{\mathbf{X}} = \{\tilde{\mathbf{x}}_i\}_{i=1}^s,$$

where $\tilde{\mathbf{x}}$ denotes a newly presented data sample, and $\{\tilde{\mathbf{x}}_i\}_{i=1}^s$ represents a chunk of s data samples. \mathcal{F}_{sq} and \mathcal{F}_{ck} perform respectively the sequential and chunk incremental learning to calculate the updated model Ω' over any newly presented data. The \mathcal{F}_{ck} is generally computationally less efficient than \mathcal{F}_{sq} , when the chunk size $s = 1$ (Pang et al., 2005b; Cauwenberghs & Poggio, 2000). This thesis studies incremental linear discriminant analysis (InclDA) with a focus on how an \mathcal{F}_{ck} is able to match and surpass the \mathcal{F}_{sq} on learning efficiency without sacrificing effectiveness.

Linear Discriminant Analysis (LDA) seeks a linear projection of data that best discriminates two or more classes by the Fisher criterion or its equivalents (Fisher, 1936; Fukunaga, 1990; McLachlan, 2004). In principle, the computing of batch LDAs in different disciplines can be formulated as an n -tuple function. The clas-

sic Fisher LDA (Fukunaga, 1990) constructs the discriminant eigenspace based on the within-class scatter matrix \mathbf{S}_w , between-class scatter matrix \mathbf{S}_b , and the class label set \mathbf{C} . Thus a 3-tuple model $\mathbf{\Omega} = \{\mathbf{S}_w, \mathbf{S}_b, \mathbf{C}\}$ can be formed (Pang et al., 2005b). The GSVD/LDA (Ye, Janardan, Park & Park, 2004) applies generalized Singular Value Decomposition (SVD) on \mathbf{H} to obtain singular vector \mathbf{P} and singular value $\mathbf{\Gamma}$, then performs again the SVD on a submatrix of \mathbf{P} to solve the LDA optimization problem. Here, $\mathbf{H} = [\mathbf{H}_w \mathbf{H}_b]^T$ and $\mathbf{H} = \mathbf{P}\mathbf{\Gamma}\mathbf{\Pi}^T$. Thus, GSVD/LDA models $\mathbf{\Omega} = \{\mathbf{H}_w, \mathbf{H}_b, \mathbf{C}\}$. The Least Square LDA (Ye, 2007) transforms the LDA optimization into Multivariate Linear Regression (MLR), thus rears 3-tuple calculation $\mathbf{\Omega} = \{\mathbf{M}, \mathbf{Y}, \mathbf{C}\}$, where \mathbf{M} and \mathbf{Y} are the input centroid matrix and the output indicator matrix, respectively. The QR/LDA (Ye et al., 2005) employs QR decomposition to implement LDA optimization on size-reduced instead of full scatter matrices. Its eigenspace model is written as $\mathbf{\Omega} = \{\mathbf{M}, \mathbf{W}, \mathbf{B}\}$; here \mathbf{W} and \mathbf{B} represent reduced within-class and between-class scatter matrices respectively, and \mathbf{M} is the centroid matrix. Therefore, IncLDA can be explained as an n -tuple model updating process.

Thus, a general IncLDA model can be depicted in the example of existing sequential IDR/QR IncLDA (cIncLDA) (Ye et al., 2005) as

$$\{\mathbf{M}', \mathbf{W}', \mathbf{B}'\} = \mathcal{F}_{sq}(\{\mathbf{M}, \mathbf{W}, \mathbf{B}\}, \tilde{\mathbf{x}}).$$

In fact, for a preferable IL method, it should obey four general rules proposed by Polikar, Upda, Upda and Honavar (2001) as,

1. It should be able to learn new information from new data.
2. It should not require access to, and retain in memory, the original data.
3. It should preserve previously acquired information.
4. It should be able to accommodate new classes that may be introduced with new data.

In this thesis, we present two additional criteria for IncLDA, as (5) It should be able to process multiple data samples at one time, and single data sample can be processed as a special case; (6) It should be able to address the Small Sample Size problem.

IncLDA is essentially a model updating process, each component of the model needing to proceed through along a chunk of newly presented samples. As a result, the loss of discriminant information will become severe due to approximate projections at each update round. It is clear that the increase of update rounds is a main source of inefficiencies. To reduce update rounds at each IncLDA cycle, it is necessary to conduct chunk IncLDAs in which multiple samples can be processed at one time, instead of sequential IncLDAs.

Furthermore, most existing LDAs are derived from the Fisher criterion (Fukunaga, 1990) as

$$\phi = \operatorname{argmax}_{\phi^T \phi = \mathbf{I}} \frac{\operatorname{Tr}(\phi^T \mathbf{S}_b \phi)}{\operatorname{Tr}(\phi^T \mathbf{S}_w \phi)}, \quad (1.1)$$

and two variations are $\phi = \operatorname{argmax}_{\phi^T \phi = \mathbf{I}} \frac{\operatorname{Tr}(\phi^T \mathbf{S}_t \phi)}{\operatorname{Tr}(\phi^T \mathbf{S}_w \phi)}$, $\phi = \operatorname{argmax}_{\phi^T \phi = \mathbf{I}} \frac{\operatorname{Tr}(\phi^T \mathbf{S}_b \phi)}{\operatorname{Tr}(\phi^T \mathbf{S}_t \phi)}$. The criterion requires that either \mathbf{S}_w or \mathbf{S}_t must be non-singular. However, it is not to be followed in many practical applications as the number of samples may be much smaller than the dimension of the sample space, namely $n \ll d$. This is the so-called ‘undersampled’ problems, or ‘small sample size’ (SSS) problem (Fukunaga, 1990). To address the problem, many LDA methods have been developed, such as QR-decomposition based LDA (QR/LDA)(Ye & Li, 2004; Ye et al., 2005), Generalized Singular Value Decomposition based LDA (GSVD/LDA)(Howland, Jeon & Park, 2003; Ye et al., 2004; H. Zhao & Yuen, 2008), Least Square LDA (LS/LDA)(Ye, 2007), Null Space LDA (NS/LDA)(Huang, Liu, Lu & Ma, 2002; Chu & Thye, 2010; Chen, Liao, Ko, Lin & Yu, 2000), PCA+LDA(Belhumeur, Hespanha & Kriegman, 1997), uncorrelated LDA (ULDA)(Jin, Yang, Hu & Lou, 2001; Ye, 2004), orthogonal LDA (OLDA) (Ye, 2005), etc. Accordingly, IncLDAs, which update some of those LDA extensions incrementally, should be able to address the SSS problem.

The existing sIncLDA satisfies all the criteria presented above, except criterion 5. Our work concentrates on developing a method that satisfies criterion 5 to boost its practical usage. To evaluate the proposed cIncLDA, its performance is experimentally compared with sIncLDA on some real-world datasets. Experimental results show that the proposed cIncLDA is more stable and effective on selected datasets than sIncLDA.

For convenience, we summarize most notations used in the thesis in Table 1.1. As a convention, “ \sim ” is used to represent additional information introduced by newly added data samples.

Notation	Descriptions
\mathbf{X}	data matrix
\mathbf{X}_i	data matrix of the i -th class
d	dimension of data matrix
n	number of data samples
n_i	number of data samples in the i -th class
\mathbf{C}	set of class labels
\mathbf{c}_i	class label of the i -th class
k	number of classes
\mathbf{M}	centroid matrix of data matrix
\mathbf{m}	general mean of data matrix
\mathbf{m}_i	mean of data matrix of the i -th class
\mathbf{W}	reduced within-class scatter matrix
\mathbf{B}	reduced between-class scatter matrix

Table 1.1: NOTATIONS

1.2 Research Objectives

This research aims to propose a new chunk incremental IDR/QR LDA to resolve the limitations of the sequential incremental IDR/QR LDA and boost its performance. The proposed method follows all incremental learning criteria and addresses the SSS problem. It is designed to process data presented in random chunks. In addition, it loses less discriminant information than sIncLDA, and the discriminant eigenspace model can be updated more accurately.

1.3 Research Contributions

The contributions of this thesis are:

1. A new evaluation criteria framework for IncLDAs has been proposed. It allows us to assess IncLDAs from a new point of view, not just by comparisons on performance, efficiency and so forth.
2. Two efficient matrix augmentation methods have been proposed. They allow us to easily accommodate information of new classes contained in chunk data through expanding the dimensions of matrices \mathbf{Q} and \mathbf{W} .

3. A new method for incremental updating of \mathbf{W} , more accurate than the existing approach, has been successfully developed in this thesis. This is achieved by relaxing a key assumption on \mathbf{Q} and \mathbf{W} that is vital to the existing algorithm but may be significantly violated as the number of classes increases.
4. Unlike Pang’s IncLDAs (Pang et al., 2005b) where only the chunk data in one class can be acquired each time, our method is capable of absorbing information obtained from newly added samples in whatever classes.

1.4 Thesis Structure

Chapter 2 investigates the fundamental concepts of LDA first, and reviews several popular LDA methods based on their principles and pseudocodes; After that, incremental extensions of those LDAs methods will be depicted in detail. Finally, based on what has been reviewed, motivations of this research is highlighted.

Chapter 3 recaptures and analyzes the existing batch QR/LDA and sequential incremental IDR/QR LDA from an embedding point of view. The principles of batch QR/LDA and sIncLDA are both recapitulated first. The limitations of the sIncLDA are summarized and discussed in depth at the end of the chapter.

Chapter 4 presents the proposed chunk incremental IDR/QR LDA. To improve computational efficiency, two new matrix augmentation methods are proposed. The proposed cIncLDA has been detailed with mathematical derivations and its pseudocode.

Chapter 5 examine the accuracy and efficiency of the proposed cIncLDA by comparing it with sIncLDA and batch QR/LDA. We have specifically examined its equality to the ground truth eigenspace, its class separability of the embedding eigenspace, and its execution time.

Chapter 6 concludes the work presented by the thesis and gives some directions for future research.

Chapter 2

Literature Review

In this chapter, the fundamental concepts of LDA will be investigated first, and popular LDA methods will be presented based on their principles and pseudocodes; After that, we study incremental extensions of those LDA methods in detail. Finally, based on what has been reviewed, motivations of this research will be highlighted.

2.1 Linear Discriminant Analysis

LDA has been widely demonstrated as an effective technique for dimension reduction and feature extraction (Ye & Li, 2004; Pang et al., 2005b; T.-K. Kim, Stenger, Kittler & Cipolla, 2011). Applications of LDA are growing in number, such as pattern recognition (Bishop, 2006; Duda, Hart & Stork, 2001; Fukunaga, 1990), pedestrian detection (X. Wang, Han & Yan, 2009), information retrieval (Kowalski, 1997; Frakes & Baeza-Yates, 1992), micro-array data analysis (Baldi & Hatfield, 2002; Dudoit, Fridlyand & Speed, 2002), text classification (Jain & Dubes, 1988), original texture analysis (Chan, Kittler & Messer, 2007) as well as face recognition (Jin et al., 2001; Swets & Weng, 1996). The aim of dimension reduction is to either find a linear combination of features, or select a subset of features from the original feature space (Ye et al., 2005). In this way, LDA can be simply defined as: LDA seeks a linear projection of data (in terms of a combination of features) that best discriminates two or more classes by the Fisher criterion or its equivalents (Fisher, 1936; Fukunaga, 1990; McLachlan, 2004).

2.1.1 Classic Fisher LDA

Most existing LDAs follow or extend the classic Fisher criterion. By employing within-class scatter matrix \mathbf{S}_w , between-class scatter matrix \mathbf{S}_b , and total scatter matrix \mathbf{S}_t , the class separability criterion can be formulated as follows (Fisher, 1936; Fukunaga, 1990):

$$\phi = \operatorname{argmax}_{\phi} \frac{\operatorname{Tr}(\phi^T \mathbf{S}_b \phi)}{\operatorname{Tr}(\phi^T \mathbf{S}_w \phi)}, \quad (2.1)$$

and two variations are

$$\phi = \operatorname{argmax}_{\phi} \frac{\operatorname{Tr}(\phi^T \mathbf{S}_t \phi)}{\operatorname{Tr}(\phi^T \mathbf{S}_w \phi)}, \quad (2.2)$$

$$\phi = \operatorname{argmax}_{\phi} \frac{\operatorname{Tr}(\phi^T \mathbf{S}_b \phi)}{\operatorname{Tr}(\phi^T \mathbf{S}_t \phi)}, \quad (2.3)$$

where

$$\mathbf{S}_b = \sum_{i=1}^k n_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T \quad (2.4)$$

denotes the scatter of the expected vectors around the general mean,

$$\mathbf{S}_w = \sum_{i=1}^k \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \mathbf{m}_i)(\mathbf{x}_{ij} - \mathbf{m}_i)^T \quad (2.5)$$

represents the expected vector of the general distribution, and

$$\mathbf{S}_t = \mathbf{S}_w + \mathbf{S}_b \quad (2.6)$$

is the covariance matrix of all samples without consideration of their class assignments. Given $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_k] \in \mathcal{R}^{d \times n}$, in which $\mathbf{X}_i = [\mathbf{x}_{i1}, \dots, \mathbf{x}_{in_i}] \in \mathcal{R}^{d \times n_i}$ denotes samples belonging to the i -th class, and $i = 1, \dots, k$. For the i -th class, \mathbf{m}_i and n_i denote the mean and the number of samples, respectively. The general mean and the total number of samples are represented by \mathbf{m} and n . **Note that**, the equivalence between (2.1) and those two alternative criteria (2.2) and (2.3) is conditional on the non-singular status of \mathbf{S}_w or \mathbf{S}_t . Thus, the Fisher LDA forms a 3-tuple model as $\Omega = \{\mathbf{S}_w, \mathbf{S}_b, \mathbf{C}\}$.

The pseudocode of the above algorithm is given as algorithm 1.

Algorithm 1 Fisher/LDA Algorithm

Input: Newly added training data matrix $\widetilde{\mathbf{X}}$, and its class label set $\widetilde{\mathbf{C}}$ **Output:** Optimal transform matrix ϕ , Fisher-LDA eigenspace model $\Omega = \{\mathbf{S}_w, \mathbf{S}_b, \mathbf{C}\}$.

- 1: Calculate within-class scatter matrix \mathbf{S}_w by (2.5);
 - 2: Calculate between-class scatter matrix \mathbf{S}_b by (2.4);
 - 3: Calculate matrix $\mathbf{Z} = \mathbf{S}_w^{-1}\mathbf{S}_b$;
 - 4: Perform eigen-decomposition of \mathbf{Z} : $\mathbf{Z} = \phi\Lambda\phi^T$.
-

2.1.2 Small Sample Size Problem

For many tasks, such as information retrieval and face recognition, LDA has already been showing promising results (Belhumeur et al., 1997; W. Zhao, Chellappa & Phillips, 1999; Chen et al., 2000; Yu & Yang, 2001; C. Liu & Wechsler, 2002; Lu, Plataniotis & Venetsanopoulos, 2003a, 2003b; Ye & Li, 2004). However, LDA suffers from the so-called Small Sample Size (SSS) problem (Raudys & Jain, 1990). The SSS problem is often encountered in high dimensional pattern recognition tasks where the number of training samples (i.e., n) is smaller than the dimensionality of the sample space (i.e., d), namely $n < d$. In such a case, the within-class scatter matrix \mathbf{S}_w or the total scatter matrix \mathbf{S}_t may not be singular. As a result, the Fisher/LDA may fail.

Simply put, there are three ways to solve the problem:

1. It can be solved by utilizing linear algebra techniques, such pseudo inverse, SVD, or GSVD, on the singular within-class scatter matrix. For example, Tian, Barbero, Gu and Lee (1986) proposed a solution to the problem by using pseudo inverse on the within-class scatter matrix. Some research (Hong & Yang, 1991; W. Zhao et al., 1999) gave a solution of adding a small perturbation to the within-class scatter matrix so that it became nonsingular.
2. It can also be solved by involving a subspace approach. For example, Belhumeur et al. (1997) applied PCA as a preprocessing step on the raw data to remove the null space of \mathbf{S}_w , and hereby the singularity of \mathbf{S}_w was well handled. Another solution used the null space of \mathbf{S}_w rather \mathbf{S}_w itself (Chen et al., 2000; Chu & Thye, 2010).
3. Regularization is another way to cope with the problem, adding a scaled iden-

tity matrix to diagonal elements of the within-class scatter matrix (Friedman, 1989; Lu, Plataniotis & Venetsanopoulos, 2005).

In summary, the third approach is the most efficient, since it only involves simple matrix addition. To address SSS problem of LDA, we need to extend the Fisher criterion by employing the above approaches. In the following text, three popular LDA extensions, GSVD/LDA (Ye et al., 2004; H. Zhao & Yuen, 2008), LS/LDA (Ye, 2007; L. Liu, Jiang & Zhou, 2009), and QR/LDA (Ye & Li, 2004; Ye et al., 2005), are depicted.

2.1.3 LDA Extensions to the SSS Problem

GSVD/LDA

The concept of GSVD was first introduced by Paige and Saunders (1981). Based on their idea, Ye et al. (2004) proposed a GSVD-based LDA extension. GSVD/LDA follows the Theorem 1 below:

Theorem 1 Let $\mathbf{e}_i = [1, \dots, 1]^T \in \mathcal{R}^{n_i \times 1}$,

$$\mathbf{H}_w = [\mathbf{X}_1 - \mathbf{m}_1 \mathbf{e}_1^T, \dots, \mathbf{X}_k - \mathbf{m}_k \mathbf{e}_k^T], \quad (2.7)$$

and

$$\mathbf{H}_b = [\sqrt{n_1}(\mathbf{m}_1 - \mathbf{m}), \dots, \sqrt{n_k}(\mathbf{m}_k - \mathbf{m})]. \quad (2.8)$$

Suppose that $\mathbf{H} = [\mathbf{H}_w \ \mathbf{H}_b]^T \in \mathcal{R}^{(k+n) \times d}$, then there exist orthogonal matrices $\mathbf{U} \in \mathcal{R}^{k \times k}$, $\mathbf{V} \in \mathcal{R}^{n \times n}$ and a nonsingular matrix $\mathbf{\Upsilon} \in \mathcal{R}^{n \times n}$, such that

$$\begin{bmatrix} \mathbf{U} & \mathbf{0} \\ \mathbf{0} & \mathbf{V} \end{bmatrix}^T \mathbf{H} \mathbf{\Psi} = \begin{bmatrix} \mathbf{\Sigma}_b & \mathbf{0} \\ \mathbf{\Sigma}_w & \mathbf{0} \end{bmatrix} \quad (2.9)$$

where

$$\mathbf{\Sigma}_b = \begin{bmatrix} \mathbf{I}_b & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_b & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0}_b \end{bmatrix} \quad \text{and} \quad \mathbf{\Sigma}_w = \begin{bmatrix} \mathbf{0}_w & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_w & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_w \end{bmatrix}$$

and $\mathbf{\Sigma}_b^T \mathbf{\Sigma}_b + \mathbf{\Sigma}_w^T \mathbf{\Sigma}_w = \mathbf{I}_d$.

In the Theorem 1, $\mathbf{I}_b \in \mathcal{R}^{t \times t}$ and $\mathbf{I}_w \in \mathcal{R}^{(t-s-r) \times (t-s-r)}$ are two identity matrices, in which $t = \text{rank}(\mathbf{H})$, $r = \text{rank}(\mathbf{H}) - \text{rank}(\mathbf{H}_w)$, and $s = \text{rank}(\mathbf{H}_b) + \text{rank}(\mathbf{H}_w) - \text{rank}(\mathbf{H})$. In addition, $\mathbf{D}_b = \text{diag}(\tau_{t+1}, \dots, \tau_{t+r}) \in \mathcal{R}^{r \times r}$ and $\mathbf{D}_w = \text{diag}(\nu_{t+1}, \dots, \nu_{t+s})$ are diagonal matrices with diagonal elements value positioning between 0 and 1.

From Theorem 1, as $\mathbf{S}_w = \mathbf{H}_w \mathbf{H}_w^T$ and $\mathbf{S}_b = \mathbf{H}_b \mathbf{H}_b^T$, we have

$$\Psi^T \mathbf{S}_b \Psi = \begin{bmatrix} \Sigma_b^T \Sigma_b & \\ & \mathbf{0} \end{bmatrix} \quad (2.10)$$

$$\Psi^T \mathbf{S}_w \Psi = \begin{bmatrix} \Sigma_w^T \Sigma_w & \\ & \mathbf{0} \end{bmatrix} \quad (2.11)$$

Accordingly, we can have

$$\nu_i^2 \mathbf{S}_b \tilde{\tau}_i = \nu_i^2 \mathbf{S}_w \tau_i, \quad (2.12)$$

where τ_i is the i -th column of the matrix \mathbf{Y} , $i = 1, \dots, d$.

Hence, the leftmost $k-1$ eigenvectors of \mathbf{Y} are obtained as the optimal transform matrix. The algorithm of GSVD/LDA is summarized in algorithm 2.

Algorithm 2 GSVD/LDA Algorithm

Input: Newly added training data matrix $\tilde{\mathbf{X}}$, and its class label set $\tilde{\mathbf{C}}$

Output: Optimal transform matrix ϕ , GSVD/LDA eigenspace model $\Omega = \{\mathbf{H}_w, \mathbf{H}_b, \mathbf{C}\}$.

- 1: Calculate matrix $\mathbf{H} = [\mathbf{H}_w \ \mathbf{H}_b]^T$;
- 2: Calculate the rank of \mathbf{H} as $t = \text{rank}(\mathbf{H})$;
- 3: Perform SVD on \mathbf{H} as $\mathbf{H} = \mathbf{P} \begin{bmatrix} \mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{\Pi}^T$;
- 4: Perform SVD on $\mathbf{P}(1:k, 1:t)$ as $\mathbf{P}(1:k, 1:t) = \mathbf{U} \mathbf{Y} \mathbf{\Xi}$;
- 5: Calculate submatrix $\mathbf{\Pi}_k$ as leading principal submatrix of $\mathbf{\Pi}$;
- 6: Calculate the rank of \mathbf{H}_b as $\tau = \text{rank}(\mathbf{H}_b)$;
- 7: Calculate matrix

$$\mathbf{T} = \mathbf{\Pi} \begin{bmatrix} \mathbf{R}^{-1} \mathbf{\Xi} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix};$$

- 8: Calculate optimal transform matrix $\phi = \mathbf{T}(:, 1:\tau)$.
-

LS/LDA

Ye (2007) investigated the relationship between LDA and the MLR problem and proposed LS/LDA. By adding a constraint to the optimization of the Fisher criterion (2.1), LDA can be integrated into the MLR schema, as

$$\begin{aligned} \max \operatorname{Tr}((\boldsymbol{\phi}^T \mathbf{S}_b \boldsymbol{\phi})(\boldsymbol{\phi}^T \mathbf{S}_t \boldsymbol{\phi})^+), \\ \text{s.t. } \boldsymbol{\phi} = (\mathbf{M}\mathbf{M}^T)^+ \mathbf{M}\mathbf{Y}, \end{aligned} \quad (2.13)$$

where \mathbf{M} is the centroid matrix of \mathbf{X} , and \mathbf{Y} is the indicator matrix.

From (2.13), we can compute \mathbf{Y} , as

$$\mathbf{Y}_{ji} = \begin{cases} \sqrt{\frac{n}{n_k} - \frac{n_k}{n}} & \text{if } \mathbf{X}_{\cdot i} \text{ belongs to class } k \\ -\sqrt{\frac{n_k}{n}} & \text{otherwise.} \end{cases} \quad (2.14)$$

Thus, the MLR solution $\boldsymbol{\phi}_{MLR}$ can be computed from (2.13) as

$$\boldsymbol{\phi}_{MLR} = (\mathbf{M}^+)^T \mathbf{Y} \quad (2.15)$$

Consequently, we have a 3-tuple model of LS/LDA as $\Omega = \{\mathbf{M}, \mathbf{Y}, \mathbf{C}\}$, and the pseudocode is shown in algorithm 3. In LS/LDA, the SSS problem is solved by transforming the LDA optimization problem into the MLR framework, since the singularity of \mathbf{S}_w or \mathbf{S}_t can be avoided in the MLR schema.

Algorithm 3 Least Square LDA Algorithm

Input: Newly added training data matrix $\widetilde{\mathbf{X}}$, and its class label set $\widetilde{\mathbf{C}}$

Output: Optimal transform matrix $\boldsymbol{\phi}$, LS/LDA component model $\Omega = \{\mathbf{M}, \mathbf{Y}, \mathbf{C}\}$.

- 1: Calculate centroid matrix \mathbf{M} ;
 - 2: Calculate pseudo inverse of centroid matrix \mathbf{M}^+ ;
 - 3: Calculate indicator matrix \mathbf{Y} by (2.14);
 - 4: Calculate optimal transform matrix $\boldsymbol{\phi}$ by (2.15);
-

QR/LDA

QR/LDA (Ye & Li, 2004; Ye et al., 2005) used QR-decomposition rather than SVD or GSVD to maximizing the Fisher criterion (2.1) by maximizing the separation

between different classes (i.e., maximize between-class scatter matrix \mathbf{S}_b), meanwhile minimizing the within-class distance (i.e., maximize within-class scatter matrix \mathbf{S}_w). Because QR/LDA is the basis of our proposed algorithm, the detailed review on this method is given in section 3.1.

For QR/LDA, to address SSS problem, the regularization method has been employed. Different from previous regularized LDAs, such as Regularized LDA (Lu et al., 2005) and Spectral Regression Discriminant Analysis (SRDA)(Cai, He & Han, 2008), the QR/LDA and sIncLDA applied the regularization step on the reduced within-class scatter matrix \mathbf{W} instead of \mathbf{S}_w . As discussed previously, the regularization method is the most efficient way to address SSS problem.

2.1.4 Summary

In many real-world applications, it is common for training data to be continuously presented over time rather than all being given in one batch in advance. Since the data is often given in random chunks, we cannot know in advance the next sample to process. In view of this, the LDAs discussed above suffer certain limitations. Particularly, the computational cost of batch LDAs, such as Fisher/LDA and GSVD/LDA, is extremely high when the dataset is large. Meanwhile, batch LDAs keep all data samples, incurring a high level of memory usage. This is obviously inefficient for those tasks needing online learning. Because of these limitations, it is necessary to conduct learning in an incremental way, where only information carried by newly presented data is accommodated for the existing model.

2.2 Incremental Linear Discriminant Analysis

This thesis focuses on incremental linear discriminant analysis. In this section, before studying IncLDA in depth, we first revisit the two important characteristics of incremental learning (IL), chunk and one-pass. Most existing IncLDAs will be reviewed and evaluated according to our proposed assessing criteria in Section 1.1.

2.2.1 Chunk Incremental Learning

IL operates numerous updates on a set of components (e.g., within-class scatter matrix \mathbf{S}_w) to renew the entire space/model. Sequential IL updates the learning

model for every data sample presented, causing many update calculations; As model updating may be computationally intensive, inaccuracies and inefficiencies may build up quickly (Hall, Marshall & Martin, 2000; Ozawa, Pang & Kasabov, 2006). Chunk IL is more likely to improve accuracy and efficiency, because it processes multiple samples at one time.

It is worth noting that chunk IL has recently gained more attention (Hao et al., 2004; Jiang, Song, Wu, Maurizio & Liang, 2006; Hall et al., 2000; Pang, Ozawa & Kasabov, 2004; Pang et al., 2005b; Pang, Ozawa & Kasabov, 2005a; Ozawa et al., 2006; J.-G. Wang, Sung & Yau, 2010). For incremental principle component analysis (IncPCA), Hall et al. (2000) enabled chunk IncPCA computing by merging original and newly created eigenspace; Ozawa et al. (2006) developed another version of chunk IncPCA by implementing both eigenspace rotation and augmentation. For incremental support vector machine (IncSVM), Karasuyama and Takeuchi (2010) extended a sequential (Cauwenberghs & Poggio, 2000) to chunk IncSVM algorithm by updating multiple Lagrange multipliers under the condition of keeping Karush-Kuhn-Tucker (KKT) balance.

2.2.2 One-pass Incremental Learning

When data are presented sequentially or in random chunks, despite incremental learning, a straightforward approach for learning is to collect the data and then conduct batch learning over the data so far collected. Obviously, this requires large memory and high computational cost, because the data presented so far would be stored in the memory until the learning process ends. Moreover, even if the learning is almost finished, adding the last single data sample, still requires repetition of the learning from the beginning (Pang et al., 2004, 2005b).

To cope with the problem, the one-pass concept has been widely applied to the incremental learning (Yen & Meesad, 1999; Loos, 1989; Pang et al., 2005b; Kidera, Ozawa & Abe, 2006; Pang, Ban, Kadobayashi & Kasabov, 2010). According to Loos (1989); Pang et al. (2005b), one-pass can be explained as, in the learning schema, only knowledge from a single presentation of the training data is acquired, learned, and retained, without keeping all data samples presented so far.

In the machine learning field, one-pass incremental learning has become more and more popular. For incremental PCA, Pang et al. (2004) proposed an one-pass PCA, which stores the eigenvectors of the new face data for updating the existing eigenspace

model rather than the raw face data. For incremental clustering, Ataa Allah, Grosky and Aboutajdine (2007) proposed an online single-pass (i.e., one-pass) method, which requires a single, sequential pass over the incoming data it attempts to cluster.

2.2.3 Incremental Linear Discriminant Analysis

A batch LDA can be formed as a component model Ω , and incremental LDA can be described as the update of Ω . In the following text, we will discuss several popular IncLDAs as the update of batch LDAs mentioned in section 2.1.

Pang's IncLDA

Pang et al. (2005b) proposed an incremental LDA based on the Fisher/LDA. As discussed above, Fisher/LDA can be formed as a 3-tuple model $\Omega = \{\mathbf{S}_w, \mathbf{S}_b, \mathbf{C}\}$, and thus can be updated in a sequential manner or chunk manner respectively,

$$\Omega' = \mathcal{F}_{sq}(\Omega, \tilde{\mathbf{x}}) = \{\mathbf{S}'_w, \mathbf{S}'_b, \mathbf{C}'\},$$

or

$$\Omega' = \mathcal{F}_{ck}(\Omega, \tilde{\mathbf{X}}) = \{\mathbf{S}'_w, \mathbf{S}'_b, \mathbf{C}'\}, \tilde{\mathbf{X}} = \{\tilde{\mathbf{x}}_i\}_{i=1}^s.$$

For Pang's sequential IncLDA, $\Omega' = \{\mathbf{S}'_w, \mathbf{S}'_b, \mathbf{C}'\}$ can be calculated following the steps below:

1. if $\tilde{\mathbf{x}}$ comes from an existing class, then $n' = n + 1$ and $\mathbf{C}' = \mathbf{C}$. The updated within-class scatter matrix \mathbf{S}'_w is calculated as,

$$\mathbf{S}'_w = \sum_{i=1, i \neq \rho}^k \Sigma_i + \Sigma'_\rho, \quad (2.16)$$

where

$$\Sigma_i = \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \mathbf{m}_i)(\mathbf{x}_{ij} - \mathbf{m}_i)^T, \quad (2.17)$$

and

$$\Sigma'_\rho = \Sigma_\rho + \frac{n_\rho}{n_\rho + 1} (\tilde{\mathbf{x}} - \mathbf{m}_\rho)(\tilde{\mathbf{x}} - \mathbf{m}_\rho)^T. \quad (2.18)$$

The updated between-class scatter matrix \mathbf{S}'_b can be calculated as,

$$\mathbf{S}'_b = \sum_{i=1, i \neq \rho}^k n_i (\mathbf{m}_i - \mathbf{m}') (\mathbf{m}_i - \mathbf{m}')^T + n'_\rho (\mathbf{m}'_\rho - \mathbf{m}') (\mathbf{m}'_\rho - \mathbf{m}')^T, \quad (2.19)$$

where $\mathbf{m}'_\rho = \frac{1}{n_\rho + 1} (n_\rho \mathbf{m}_\rho + \tilde{\mathbf{x}})$ and $n'_\rho = n_\rho + 1$.

2. if $\tilde{\mathbf{x}}$ belongs to a new class, then $n' = n + 1$ and $\mathbf{C}' = [\mathbf{C}, \tilde{\mathbf{c}}]$. The updated within-class scatter matrix does not change, namely $\mathbf{S}'_w = \mathbf{S}_w$.

The updated between-class scatter matrix \mathbf{S}'_b is calculated as,

$$\mathbf{S}'_b = \sum_{i=1}^k n_i (\mathbf{m}_i - \mathbf{m}') (\mathbf{m}_i - \mathbf{m}')^T + (\tilde{\mathbf{x}} - \mathbf{m}') (\tilde{\mathbf{x}} - \mathbf{m}')^T. \quad (2.20)$$

The pseudocode of the above algorithm is given as algorithm 4.

Algorithm 4 Pang's Sequential IncLDA Algorithm

Input: Newly added training data sample $\tilde{\mathbf{x}}$ and its class label $\tilde{\mathbf{c}}$, Fisher/LDA eigenspace model

$$\Omega = \{\mathbf{S}_w, \mathbf{S}_b, \mathbf{C}\}.$$

Output: Optimal transform matrix ϕ' , Pang's Sequential IncLDA eigenspace model $\Omega' =$

$$\{\mathbf{S}'_w, \mathbf{S}'_b, \mathbf{C}'\}.$$

- 1: **if** $\tilde{\mathbf{c}} \in \mathbf{C}$ **then**
 - 2: Calculate within-class scatter matrix \mathbf{S}'_w by (2.16);
 - 3: Calculate between-class scatter matrix \mathbf{S}'_b by (2.19);
 - 4: **else**
 - 5: Calculate $\mathbf{C}' = [\mathbf{C}, \tilde{\mathbf{c}}]$;
 - 6: Calculate within-class scatter matrix $\mathbf{S}'_w = \mathbf{S}_w$;
 - 7: Calculate between-class scatter matrix \mathbf{S}'_b by (2.20);
 - 8: **end if**
 - 9: Calculate matrix $\mathbf{Z}' = \mathbf{S}'_w^{-1} \mathbf{S}'_b$;
 - 10: Perform eigen-decomposition of \mathbf{Z}' : $\mathbf{Z}' = \phi' \mathbf{\Lambda} \phi'^T$.
-

Pang's chunk IncLDA requires the newly added samples to belong to one class, namely the class number of incoming data $\tilde{\mathbf{k}} = 1$. Hereby, $\Omega' = \{\mathbf{S}'_w, \mathbf{S}'_b, \mathbf{C}'\}$ is calculated in the following steps:

1. if $\widetilde{\mathbf{X}}_i$ comes from an existing class, then $n' = n + s$ and $\mathbf{C}' = \mathbf{C}$. The updated within-class scatter matrix \mathbf{S}'_w is calculated as,

$$\mathbf{S}'_w = \sum_{i=1}^k \boldsymbol{\Sigma}'_i, \quad (2.21)$$

where

$$\boldsymbol{\Sigma}'_i = \boldsymbol{\Sigma}_i + \frac{n_i s^2}{(n_i + s)^2} (\mathbf{D}_i) + \frac{n_i^2}{(n_i + s)^2} (\mathbf{E}_i) + \frac{s(s + 2n_i)}{(n_i + s)^2} (\mathbf{F}_i). \quad (2.22)$$

From (2.22), we have the second term as

$$\mathbf{D}_i = (\widetilde{\mathbf{m}}_i - \mathbf{m}_i)(\widetilde{\mathbf{m}}_i - \mathbf{m}_i)^T, \quad (2.23)$$

the third term as

$$\mathbf{E}_i = \sum_{j=1}^s (\widetilde{\mathbf{x}}_{ij} - \mathbf{m}_i)(\widetilde{\mathbf{x}}_{ij} - \mathbf{m}_i)^T, \quad (2.24)$$

and the fourth term as

$$\mathbf{F}_i = \sum_{j=1}^s (\widetilde{\mathbf{x}}_{ij} - \widetilde{\mathbf{m}}_i)(\widetilde{\mathbf{x}}_{ij} - \widetilde{\mathbf{m}}_i)^T, \quad (2.25)$$

The updated between-class scatter matrix \mathbf{S}'_b can be calculated as

$$\mathbf{S}'_b = \sum_{i=1}^k n'_i (\mathbf{m}'_i - \mathbf{m}')(\mathbf{m}'_i - \mathbf{m}')^T. \quad (2.26)$$

2. if $\widetilde{\mathbf{x}}$ belongs to a new class, then $n' = n + 1$ and $\mathbf{C}' = \mathbf{C} \cup \widetilde{\mathbf{C}}$. The updated within-class scatter matrix \mathbf{S}'_w is calculated as

$$\mathbf{S}'_w = \sum_{i=1}^k \boldsymbol{\Sigma}_i + \sum_{j=1}^s (\widetilde{\mathbf{x}}_j - \widetilde{\mathbf{m}})(\widetilde{\mathbf{x}}_j - \widetilde{\mathbf{m}})^T. \quad (2.27)$$

The updated between-class scatter matrix \mathbf{S}'_b is calculated as,

$$\mathbf{S}'_b = \sum_{i=1}^k n_i (\mathbf{m}_i - \mathbf{m}')(\mathbf{m}_i - \mathbf{m}')^T + (\widetilde{\mathbf{x}} - \mathbf{m}')(\widetilde{\mathbf{x}} - \mathbf{m}')^T. \quad (2.28)$$

The pseudocode of the above algorithm is given as algorithm 5.

Algorithm 5 Pang’s Chunk IncLDA Algorithm

Input: Newly added training data matrix $\widetilde{\mathbf{X}}_i$, $\widetilde{k} = 1$ and its class label set $\widetilde{\mathbf{C}}$, Fisher-LDA eigenspace model $\mathbf{\Omega} = \{\mathbf{S}_w, \mathbf{S}_b, \mathbf{C}\}$.

Output: Optimal transform matrix ϕ' , Pang’s Sequential IncLDA eigenspace model $\mathbf{\Omega}' = \{\mathbf{S}'_w, \mathbf{S}'_b, \mathbf{C}'\}$.

- 1: **if** $\widetilde{\mathbf{C}} \subset \mathbf{C}$ **then**
 - 2: Calculate within-class scatter matrix \mathbf{S}'_w by (2.21);
 - 3: Calculate between-class scatter matrix \mathbf{S}'_b by (2.26);
 - 4: **else**
 - 5: Calculate $\mathbf{C}' = \mathbf{C} \cup \widetilde{\mathbf{C}}$;
 - 6: Calculate within-class scatter matrix \mathbf{S}'_w by (2.27);
 - 7: Calculate between-class scatter matrix \mathbf{S}'_b by (2.28);
 - 8: **end if**
 - 9: Calculate matrix $\mathbf{Z}' = \mathbf{S}'_w^{-1} \mathbf{S}'_b$;
 - 10: Perform eigen-decomposition of \mathbf{Z}' : $\mathbf{Z}' = \phi' \mathbf{\Lambda} \phi'^T$.
-

GSVD IncLDA

GSVD/LDA models $\mathbf{\Omega} = \{\mathbf{H}_w, \mathbf{H}_b, \mathbf{C}\}$, according to our previous study, a direct incremental GSVD/LDA can update the $\mathbf{\Omega}$ as

$$\mathbf{\Omega}' = \mathcal{F}_{ck}(\mathbf{\Omega}, \widetilde{\mathbf{X}}) = \{\mathbf{H}'_w, \mathbf{H}'_b, \mathbf{C}'\}, \widetilde{\mathbf{X}} = \{\widetilde{\mathbf{x}}_i\}_{i=1}^s.$$

However, the computational cost of twice SVD on the large matrix \mathbf{H} is extremely high (H. Zhao & Yuen, 2008). To address the problem, it was necessary to use an efficient and effective SVD updating method, but it was found difficult to conduct SVD-updating on \mathbf{H} (H. Zhao & Yuen, 2008). Thus, H. Zhao and Yuen (2008) modified Ye’s GSVD/LDA by computing \mathbf{H}_t instead of \mathbf{H} , and derived an incremental GSVD/LDA (GSVD IncLDA) on the modified GSVD/LDA by updating singular vectors \mathbf{Q} and \mathbf{P} , respectively. Thus, the modified GSVD/LDA models $\mathbf{\Omega} = \{\mathbf{H}_b, \mathbf{Q}, \mathbf{P}, \mathbf{C}\}$.

For \mathbf{Q} in $\mathbf{\Omega}$, because it is the right singular vector of \mathbf{H} , and it is also the

eigenvector of the total scatter matrix \mathbf{S}_t , the method turns to update \mathbf{S}_t ,

$$\begin{aligned} \mathbf{S}'_t &= \mathbf{H}'_t \mathbf{H}'_t{}^T \\ &= \mathbf{H}_t \mathbf{H}_t{}^T + (\widetilde{\mathbf{X}} - \widetilde{\mathbf{m}} \widetilde{\mathbf{e}}^T)(\widetilde{\mathbf{X}} - \widetilde{\mathbf{m}} \widetilde{\mathbf{e}}^T)^T \\ &\quad + \frac{n\widetilde{n}}{n+\widetilde{n}}(\mathbf{m} - \widetilde{\mathbf{m}})(\mathbf{m} - \widetilde{\mathbf{m}})^T, \end{aligned} \quad (2.29)$$

Let

$$\mathbf{G} = \left[\mathbf{X} - \mathbf{m} \mathbf{e}^T, \widetilde{\mathbf{X}} - \widetilde{\mathbf{m}} \widetilde{\mathbf{e}}^T, \sqrt{\frac{n\widetilde{n}}{n+\widetilde{n}}}(\mathbf{m} - \widetilde{\mathbf{m}}) \right] \quad (2.30)$$

as we have

$$\widetilde{\mathbf{H}}_t = \left[\mathbf{X} - \mathbf{m} \mathbf{e}^T, \widetilde{\mathbf{X}} - \widetilde{\mathbf{m}} \widetilde{\mathbf{e}}^T \right], \quad (2.31)$$

then the updating of \mathbf{Q} can be obtained by calculating the SVD of \mathbf{G} , since $\mathbf{G} = \mathbf{H}'_t$, and

$$\mathbf{Q}'^T \mathbf{S}'_t \mathbf{Q} = \mathbf{Q}'^T \mathbf{H}'_t{}^T \mathbf{H}'_t \mathbf{Q} = \begin{bmatrix} \mathbf{R}'^2 & 0 \\ 0 & 0 \end{bmatrix} \quad (2.32)$$

When $\widetilde{\mathbf{X}}$ is inserted, the first k columns of \mathbf{G} are unchanged, thus the SVD of \mathbf{G} can take place using the SVD updating technique introduced by Zha and Simon (1997) as:

Given SVD on a $n \times d$ matrix \mathbf{A} as $\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{\Upsilon}$, the SVD of $[\mathbf{A}, \mathbf{B}]$ can be calculated as

$$[\mathbf{A}, \mathbf{B}] = ([\mathbf{U}_t, \mathbf{Q}] \widehat{\mathbf{U}}) \widehat{\mathbf{\Lambda}} \left(\begin{bmatrix} \mathbf{\Upsilon}_t & 0 \\ 0 & \mathbf{I} \end{bmatrix} \widehat{\mathbf{\Upsilon}} \right)^T \quad (2.33)$$

where $t = \text{rank}(\mathbf{A})$, \mathbf{U}_t and $\mathbf{\Upsilon}_t$ are constructed by the first t columns of \mathbf{U} and $\mathbf{\Upsilon}$, respectively. The \mathbf{Q} is from the QR decomposition $\mathbf{I} - \mathbf{U}_t \mathbf{U}_t^T = \mathbf{Q} \mathbf{L}$; The $\widehat{\mathbf{U}}$, $\widehat{\mathbf{\Lambda}}$, and $\widehat{\mathbf{\Upsilon}}$ comes from the following SVD

$$\begin{bmatrix} \mathbf{\Lambda}_t & \mathbf{U}_t^T \mathbf{B} \\ 0 & \mathbf{L} \end{bmatrix} = \widehat{\mathbf{U}} \widehat{\mathbf{\Lambda}} \widehat{\mathbf{\Upsilon}} \quad (2.34)$$

For updating \mathbf{P} in the modified GSVD/LDA,

$$\mathbf{P}(:, 1:t) = \mathbf{H} \mathbf{Q}(:, 1:t) \mathbf{R}^{-1}, \quad (2.35)$$

Algorithm 6 GSVD InclDA algorithm

Input: Newly added training data matrix $\widetilde{\mathbf{X}}$, and its class label set $\widetilde{\mathbf{C}}$, Modified GSVD/LDA model

$$\Omega = \{\mathbf{H}_b, \mathbf{Q}, \mathbf{P}, \mathbf{C}\}.$$

Output: Optimal transform matrix ϕ' , GSVD InclDA model $\Omega' = \{\mathbf{H}'_b, \mathbf{Q}', \mathbf{P}', \mathbf{C}'\}$.

- 1: Construct matrix $\widetilde{\mathbf{X}} = \{\widetilde{\mathbf{x}}_1, \dots, \widetilde{\mathbf{x}}_n\}$;
- 2: Calculate updated training data matrix $\mathbf{X}' = \begin{bmatrix} \mathbf{X} & \widetilde{\mathbf{X}} \end{bmatrix}$;
- 3: Calculate \mathbf{H}'_b on \mathbf{X}' ;
- 4: Calculate matrix \mathbf{G} by (2.31);
- 5: Calculate singular vectors \mathbf{Q}' and \mathbf{R}' by performing SVD updating (i.e., (2.33) and (2.34)) on \mathbf{G} ;
- 6: Calculate rank of \mathbf{H}'_b : $\delta \leftarrow \text{rank}(\mathbf{H}'_b)$;
- 7: Calculate matrix \mathbf{P}' by (2.37);
- 8: Conduct SVD on \mathbf{P}' : $\mathbf{P} = \mathbf{U}'\Sigma'\mathbf{W}'^T$;
- 9: Calculate matrix $\phi_o = \mathbf{Q}' \begin{bmatrix} \mathbf{R}'^{-1}\mathbf{W}' & 0 \\ 0 & \mathbf{I} \end{bmatrix}$;
- 10: Calculate transform matrix $\phi' = \phi_o(:, 1 : \delta)$.

According to Ye's GSVD/LDA (Ye et al., 2004), we have

$$\begin{aligned} \mathbf{P}(1 : k, 1 : t) &= \mathbf{H}(1 : k, :)\mathbf{Q}(:, 1 : t)\mathbf{R}^{-1} \\ &= \mathbf{H}_b\mathbf{Q}(:, 1 : t)\Xi(1 : t, 1 : t)^{-1}, \end{aligned} \quad (2.36)$$

by substituting $\mathbf{H} = \begin{bmatrix} \mathbf{H}_b^T \\ \mathbf{H}_w^T \end{bmatrix}$ to (2.36).

Therefore, \mathbf{P} can be updated as

$$\mathbf{P}'(:, 1 : t) = \mathbf{H}'_b\mathbf{Q}'(:, 1 : t)\Xi'(1 : t, 1 : t)^{-1} \quad (2.37)$$

In summary, algorithm 6 describes the pseudocode of the GSVD InclDA,

LS InclDA

L. Liu et al. (2009) proposed the incremental approach for LS/LDA. Given a 3-tuple LS/LDA model $\Omega = \{\mathbf{M}, \mathbf{Y}, \mathbf{C}\}$, the incremental LS/LDA (LS InclDA) updates the model as

$$\Omega' = \mathcal{F}_{sq}(\Omega, \widetilde{\mathbf{x}}) = \{\mathbf{M}', \mathbf{Y}', \mathbf{C}'\}.$$

Following (2.38), the centroid matrix \mathbf{M} is updated as:

$$\mathbf{M}' = \left[\mathbf{M} - \frac{1}{n+1}(\tilde{\mathbf{x}} - \mathbf{m}')\mathbf{e}^T \quad \frac{n}{n+1}(\tilde{\mathbf{x}} - \mathbf{m}') \right], \quad (2.38)$$

where $\mathbf{m}' = \mathbf{m} + \frac{1}{n+1}(\tilde{\mathbf{x}} - \mathbf{m})$.

The indicator matrix \mathbf{Y} can be updated as: For each new sample $\tilde{\mathbf{x}}$, a new indicator vector $\tilde{\mathbf{y}}$ is appended to \mathbf{Y} , and $\mathbf{Y} = \{y_1, \dots, y_k\}^T$ is defined as

$$\mathbf{Y}_{ji} = \begin{cases} \frac{1}{\sqrt{n_k}} & \text{if } \mathbf{X}_{.i} \text{ belongs to class } k \\ 0 & \text{otherwise.} \end{cases} \quad (2.39)$$

1. if $\tilde{c} \notin C$ then $k' = k + 1$,

$$\mathbf{Y}' = \begin{bmatrix} \mathbf{Y} & 0 \\ \mathbf{y}^T \end{bmatrix} \quad (2.40)$$

2. otherwise $\tilde{c} \in C$, then

$$\mathbf{Y}' = \begin{bmatrix} \mathbf{Y} \otimes_p \alpha_p \\ \mathbf{y}^T \end{bmatrix} \quad (2.41)$$

where

$$\alpha_p = \sqrt{\frac{n_p}{n_p + 1}} \quad (2.42)$$

the operator \otimes_p denotes multiplying the p -th column of \mathbf{Y} with α_p .

In general, the updating of \mathbf{Y} can be an unified form as

$$\mathbf{Y}' = \begin{bmatrix} \mathbf{Y} \bar{\otimes}_p \alpha_p \\ \mathbf{y}^T \end{bmatrix}, \quad (2.43)$$

in which $\bar{\otimes}_p$ is calculated if p is smaller than or equal to the width of \mathbf{Y} , the p -th column of \mathbf{Y} is multiplied with α_p ; otherwise, a new column with zero elements is appended to \mathbf{Y} .

Thus, an LS/LDA can be incremented by updating centroid data matrix \mathbf{M} and indicator matrix \mathbf{Y} , then recalculating the MLR by (2.15).

Towards a computational efficient IncLDA, L. Liu et al. (2009) extended the batch LS/LDA model as $\mathbf{\Omega}_{LSLDA} = \{\mathbf{M}, \mathbf{M}^+, \mathbf{Y}, \mathbf{C}\}$, and had (2.15) updated on $\tilde{\mathbf{x}}$ with the consideration of whether n is smaller than d or not.

In the case of $n < d$, assuming the rank of \mathbf{M} increases by 1 when a new sample

is added. \mathbf{M}^+ can be updated as

$$\mathbf{M}'^+ = \begin{bmatrix} \mathbf{M}^+ - \mathbf{M}^+(\tilde{\mathbf{x}} - \mathbf{m})\mathbf{h}^T - \frac{1}{n}\mathbf{e}\mathbf{h}^T \\ \mathbf{h}^T \end{bmatrix} \quad (2.44)$$

where

$$\mathbf{h} = \frac{(\tilde{\mathbf{x}} - \mathbf{m}) - \mathbf{M}\mathbf{M}^+(\tilde{\mathbf{x}} - \mathbf{m})}{(\tilde{\mathbf{x}} - \mathbf{m})^T(\tilde{\mathbf{x}} - \mathbf{m}) - (\tilde{\mathbf{x}} - \mathbf{m})^T\mathbf{M}\mathbf{M}^+(\tilde{\mathbf{x}} - \mathbf{m})} \quad (2.45)$$

In this way, (2.15) is updated for LS InclLDA by applying (2.44) and (2.43) to calculate \mathbf{M}'^+ and \mathbf{Y}' , respectively.

$$\begin{aligned} \phi' &= \mathbf{M}'^+\mathbf{Y}' \\ &= \left((\mathbf{M}^+)^T - \mathbf{h}(\tilde{\mathbf{x}} - \mathbf{m})^T(\mathbf{M}^+)^T - \frac{\mathbf{h}\mathbf{e}^T}{n} \right) (\mathbf{Y} \bar{\otimes}_p \alpha_p) + \mathbf{h}\mathbf{y}^T \\ &= \left(\phi - \mathbf{h}(\tilde{\mathbf{x}} - \mathbf{m})^T\phi - \frac{\mathbf{h}\mathbf{e}^T\mathbf{Y}}{n} \right) \bar{\otimes}_p \alpha_p + \mathbf{h}\mathbf{y}^T \end{aligned} \quad (2.46)$$

where \mathbf{h} is defined as (2.45), α_p as (2.42) and \mathbf{y} as (2.39).

When $n > d$, to avoid the computational cost of LS InclLDA growing with the number of samples, (2.15) can be interpolated as

$$\phi = \mathbf{T}^+\mathbf{M}\mathbf{Y}, \quad (2.47)$$

since $\mathbf{T} = \mathbf{M}\mathbf{M}^T$.

In this sense, an incremental LS/LDA model can be constructed by just updating \mathbf{T}^+ . To do that, \mathbf{T} is firstly updated as

$$\begin{aligned} \mathbf{T}' &= \mathbf{M}'\mathbf{M}'^T = \mathbf{T} + \frac{n}{n+1}(\tilde{\mathbf{x}} - \mathbf{m})(\tilde{\mathbf{x}} - \mathbf{m})^T, \\ &= \mathbf{T} + \boldsymbol{\mu}\boldsymbol{\mu}^T, \end{aligned} \quad (2.48)$$

where $\boldsymbol{\mu} = \sqrt{\frac{n}{n+1}}(\tilde{\mathbf{x}} - \mathbf{m})$. Then, \mathbf{T}^+ can be updated as

$$\mathbf{T}^+ = \begin{cases} \mathbf{T}^+ - \theta^{-1}\mathbf{s}\mathbf{s}^T & \text{if } t = 0 \\ \mathbf{T}^+ - \frac{\mathbf{s}\mathbf{r}^T + \mathbf{r}\mathbf{s}^T}{\mathbf{t}^T\mathbf{r}} + \frac{\theta\mathbf{r}\mathbf{r}^T}{(\mathbf{r}^T\mathbf{r})^2} & \text{otherwise,} \end{cases} \quad (2.49)$$

in which $\mathbf{s} = \mathbf{T}^+\boldsymbol{\mu}$, $\mathbf{r} = (\mathbf{I} - \mathbf{T}\mathbf{T}^+)\boldsymbol{\mu}$, and $\theta = 1 + \boldsymbol{\mu}^T\mathbf{T}^+\boldsymbol{\mu}$. The derivation of (2.49) can be found in L. Liu et al. (2009).

Applying (2.49) and (2.43) to calculate \mathbf{T}^+ and \mathbf{Y}' respectively, (2.47) can be

updated as

$$\begin{aligned}\phi' &= \mathbf{T}^+ \mathbf{M}' \mathbf{Y}', \\ &= \left(\mathbf{G} - \frac{\mathbf{T}^+ (\tilde{\mathbf{x}} - \mathbf{m}) e^T \mathbf{Y}}{n+1} \right) \bar{\otimes}_p \alpha_p + \frac{n \mathbf{T}^+ (\tilde{\mathbf{x}} - \mathbf{m}) \mathbf{y}}{n+1}\end{aligned}\quad (2.50)$$

where

$$\mathbf{G} = \begin{cases} \phi - \theta^{-1} s \mathbf{u}^T \phi & \text{if } \mathbf{r} = 0 \\ \phi - \frac{\mathbf{t} \mathbf{u}^T \phi}{\mathbf{r}^T \mathbf{r}} & \text{otherwise.} \end{cases}\quad (2.51)$$

In a summary, algorithm 7 presents the detailed computational steps for LS IncLDA by updating \mathbf{T}^+ . The algorithm integrates two updating algorithms proposed in L. Liu et al. (2009) for both cases of $n < d$ and $n > d$.

Sequential IDR/QR IncLDA

Ye et al. (2005) proposed the incremental version of QR/LDA, which is called sequential Incremental IDR/QR LDA. With the insertion of a new sample $\tilde{\mathbf{x}}$, the centroid matrix \mathbf{M} , \mathbf{H}_w , \mathbf{H}_b will change accordingly, where \mathbf{H}_w and \mathbf{H}_b can be obtained by $\mathbf{S}_w = \mathbf{H}_w \mathbf{H}_w^T$ and $\mathbf{S}_b = \mathbf{H}_b \mathbf{H}_b^T$, respectively. Thus, the reduced within-class scatter matrix \mathbf{W} and \mathbf{B} are changed as well.

According to Ye et al. (2005), sIncLDA utilizes three stages to accomplish the incremental learning task as follows:

1. Updating centroid matrix \mathbf{M} and its QR-decomposition.
2. Updating reduced within-class scatter matrix \mathbf{W} .
3. Updating reduced between-class scatter matrix \mathbf{B} .

Consequently, the sIncLDA can be formed as a 3-tuple model updating process from the original bath QR/LDA model Ω as

$$\Omega' = \mathcal{F}_{sq}(\Omega, \tilde{\mathbf{x}}).$$

The detailed review on this method is given in section 3.2.

2.2.4 Summary

As far as we know, the concept of chunk IncLDA is first discussed in Pang's IncLDA (Pang et al., 2005b), in which both sequential and chunk version of the incremental

Algorithm 7 LS InclLDA Algorithm

Input: Newly added training data sampe $\tilde{\mathbf{x}}$, and its class label \tilde{c} , LS/LDA eigenspace model

$$\Omega = \{M, M^+, Y, C\}.$$

Output: Optimal transform matrix ϕ' , LS InclLDA model $\Omega' = \{M', M'^+, Y', C'\}$.

- 1: **if** $\tilde{c} \in C$ **then**
 - 2: Calculate $n_i = n_i + 1$;
 - 3: **else**
 - 4: Set $n_i = 1$;
 - 5: Calculate $k = k + 1$;
 - 6: **end if**
 - 7: Calculate indicator matrix Y' by (2.40) and (2.41);
 - 8: Calculate general mean vector $\mathbf{m}' = \mathbf{m} + \frac{1}{n+1}(\tilde{\mathbf{x}} - \mathbf{m})$;
 - 9: Calculate n and d as the size of centroid matrix M ;
 - 10: **if** $d \leq n$ **then**
 - 11: Calculate matrix $T = MM^T$;
 - 12: **if** $\text{rank}(T) == d$ **then**
 - 13: Set $r = 0$;
 - 14: **else**
 - 15: Calculate $r = (I - TT^+)\mu$;
 - 16: Calculate total scatter matrix T' by (2.48);
 - 17: **end if**
 - 18: Calculate T^+ by (2.49);
 - 19: Calculate optimal transform matrix ϕ' by (2.50).
 - 20: **else**
 - 21: Calculate centroid matrix M' by (2.38);
 - 22: Calculate pseudo-inverse centroid matrix M'^+ by (2.44);
 - 23: Calculate optimal transform matrix ϕ' by (2.50).
 - 24: **end if**
-

Fisher LDA are developed and compared in terms of computation efficiency, memory usage, and eigenspace discriminability. The comparison results indicate that the running efficiency of chunk InclLDA outperforms sequential InclLDA, and improves continuously with the increase in chunk size. However, since this method cannot address the SSS problem (Song, Liu, Zhang & Yang, 2008), thus is not applicable to

any problem that has its number of samples less than dimensionality. To mitigate the problem, H. Zhao and Yuen (2008) proposed an alternative chunk IncLDA, called GSVD IncLDA. This method uses GSVD as its solution to the SSS problem. The disadvantage of the method is, it is not a one-pass incremental learning. In fact, it always retains in memory all data samples newly added at each incremental learning cycle. Apparently, this makes the algorithm memory-inefficient over large size datasets. Considering the drawbacks of GSVD IncLDA, L. Liu et al. (2009) developed another very different IncLDA, called LS IncLDA. This is superior to both IncLDAs above, because it simply updates a multivariate linear regression (Ye, 2007), instead of using eigen-decomposition as all previous methods. However, it is worth noting that the LS IncLDA is sequential-based IL. As a result, inefficiencies may build up over numerous learning cycles.

The sequential IDR/QR IncLDA, proposed by Ye et al. (2005), updates the original batch QR/LDA by using an efficient QR-updating method. The superiority of the sIncLDA as compared with other IncLDAs can be seen as follows:

1. it conducts incremental dimensionality reduction in $d - k$ scale with the least loss of discriminative information.
2. it updates a QR-decomposition of a $d \times k$ matrix and solves an eigen-decomposition of a $k \times k$ matrix. Thus, it is computationally efficient.
3. it is one-pass incremental learning, where IL is executed by a single presentation of new data.
4. SSS problem is well-handled in the algorithm by applying regularization on the reduced within-class scatter matrix \mathbf{W} .

Table 2.1 compares sIncLDA against all above IncLDAs on four IL algorithm assessing criteria.

2.3 Motivations for the Presented Research

As indicated in Table 2.1, we are aware that the drawback of sIncLDA is that it conducts only sequential IL. Also, we notice that it uses an approximation for updating the reduced within-class scatter matrix \mathbf{W} , which may be corrupted in practice for data with a large number of classes. In order to process chunk data, we intend

Algorithm	Sequential	Chunk	One-Pass	Solving SSS problems
Pang’s IncLDA (Pang et al., 2005b)	YES	YES	YES	NO
sIncLDA (Ye et al., 2005)	YES	NO	YES	YES
GSVD IncLDA (H. Zhao & Yuen, 2008)	YES	YES	NO	YES
LS IncLDA (L. Liu et al., 2009)	YES	NO	YES	YES

Table 2.1: Previous IncLDAs characterized in sequential, chunk, one-pass and SSS problem-solving properties.

to propose in this thesis a new type of chunk IncLDA on the basis of sIncLDA. The proposed cIncLDA not only enjoys the benefits of sIncLDA, but also significantly boosts its performance of the sIncLDA, because of 1) being capable of processing at one time multiple samples, so that its computational efficiency grows with the size of data chunk presented; it follows that the bigger the size of the dataset provided at one time, the faster the speed given by the proposed cIncLDA; and 2) the discriminability of eigenspace upgrades owing to the derivation of a more reliable and accurate update of the reduced within-class scatter matrix \mathbf{W} .

Chapter 3

Existing Sequential Incremental IDR/QR LDA

This chapter recaptures and analyzes the existing batch QR/LDA and sequential incremental IDR/QR LDA from an embedding point of view. Because our proposed cInclDA is on the basis of sInclDA, the technical details of sInclDA is important and therefore is presented in this chapter. We first recapitulate the principles of batch QR/LDA and sInclDA. Next, we summarize and discuss in depth the sInclDA on its limitations.

3.1 Batch QR/LDA

The batch QR/LDA also follows the Fisher criterion (1.1) in two stages:

The first stage maximizes the between-class distance by solving an optimization problem as follows:

$$\phi = \underset{\phi^T \phi = \mathbf{I}}{\operatorname{argmax}} \operatorname{Tr}(\phi^T \mathbf{S}_b \phi). \quad (3.1)$$

Ye et al. (2005) state that such a solution can be obtained by solving the eigenvalue problem on \mathbf{S}_b , or by applying QR-decomposition on the centroid matrix $\mathbf{M} = [\mathbf{m}_1, \dots, \mathbf{m}_k]$ which is the so-called Orthogonal Centroid Method (OCM)(Park, Jeon & Rosen, 2003). In batch QR/LDA, the latter solution was applied. Hereby, the optimization transformation ϕ can be obtained by $\phi = \mathbf{QZ}$ for any orthogonal matrix \mathbf{Z} , where $\mathbf{Q} \in \mathcal{R}^{d \times k}$ is orthonormal and $\mathbf{R} \in \mathcal{R}^{k \times k}$ is upper triangular. Both are obtained from the reduced QR-decomposition on \mathbf{M} , namely $\mathbf{M} = \mathbf{QR}$. The

Algorithm 8 QR/LDA Algorithm

Input: Training data matrix $\mathbf{X} \in \mathcal{R}^{d \times n}$ and its class label set \mathcal{C} .

Output: Optimal transform matrix ϕ , LDA eigenspace model $\Omega = \{\mathbf{M}, \mathbf{W}, \mathbf{B}, \mathbf{Q}, \mathbf{R}, \mathcal{C}\}$.

- 1: Calculate centroid matrix \mathbf{M} , within-class scatter matrix \mathbf{S}_w , between-class scatter matrix \mathbf{S}_B ;
- 2: Calculate QR-decomposition of \mathbf{M} by $\mathbf{M} = \mathbf{Q}\mathbf{R}$;
- 3: Calculate reduced within-class scatter matrix by $\mathbf{W} = \mathbf{Q}^T \mathbf{S}_w \mathbf{Q}$;
- 4: Calculate reduced between-class scatter matrix by $\mathbf{B} = \mathbf{Q}^T \mathbf{S}_B \mathbf{Q}$;
- 5: Calculate eigen-decomposition on $(\mathbf{W} + \mu \mathbf{I}_k)^{-1} \mathbf{B}$ to obtain \mathbf{Z} containing k eigenvectors φ_i with decreasing eigenvalues;
- 6: Calculate optimal transform matrix by $\phi = \mathbf{Q}\mathbf{Z}$.

detailed proof can be found in Ye et al. (2005). The second stage of batch QR/LDA seeks a transformation matrix ϕ such that $\phi = \mathbf{Q}\mathbf{Z}$ for some \mathbf{Z} , in order to minimize the within-classes distance. Since

$$\begin{aligned}\phi^T \mathbf{S}_b \phi &= \mathbf{Z}^T (\mathbf{Q}^T \mathbf{S}_b \mathbf{Q}) \mathbf{Z} \\ \phi^T \mathbf{S}_w \phi &= \mathbf{Z}^T (\mathbf{Q}^T \mathbf{S}_w \mathbf{Q}) \mathbf{Z},\end{aligned}$$

the original optimization problem on finding ϕ is equivalent to finding \mathbf{Z} . Here, two reduced scatter matrices $\mathbf{B} = \mathbf{Q}^T \mathbf{S}_b \mathbf{Q}$ and $\mathbf{W} = \mathbf{Q}^T \mathbf{S}_w \mathbf{Q}$ can be also obtained. Note that, \mathbf{W} and \mathbf{B} are much smaller in size than \mathbf{S}_w and \mathbf{S}_b . This greatly improves computational efficiency of the algorithm. In addition, Ye et al. (2005) give a solution to computing \mathbf{Z} by solving a small eigenvalue problem on $(\mathbf{W} + \mu \mathbf{I})^{-1} \mathbf{B}$ for $\mu > 0$. They also indicate that the solution is insensitive to the value of μ , where $\mu = 0.5$ was used in their experiments.

As discussed above, the batch QR/LDA can be written as a 3-tuple discriminant eigenspace model $\Omega = \{\mathbf{M}, \mathbf{W}, \mathbf{B}\}$.

The pseudocode is provided in algorithm 8.

3.2 Sequential Incremental IDR/QR LDA

Given batch QR/LDA eigenspace model on \mathbf{X} as $\Omega = \{\mathbf{M}, \mathbf{W}, \mathbf{B}\}$. When a new instance $\tilde{\mathbf{x}}$ is presented in the i th class, the problem of IncLDA can be described as

$$\Omega' = \mathcal{F}_{sq}(\Omega, \tilde{\mathbf{x}}) = \{\mathbf{M}', \mathbf{W}', \mathbf{B}'\}$$

1. if $\tilde{\mathbf{x}}$ belongs to an existing class (i.e., $i \leq k$), then $n'_i = n_i + 1$, $n' = n + 1$. The updating of the centroid matrix is calculated as

$$\mathbf{M}' = \mathbf{M} + f\mathbf{g}^T, \quad (3.2)$$

where $\mathbf{f} = \frac{\tilde{\mathbf{x}} - \mathbf{m}_i}{n'_i}$ and $\mathbf{g} = (0, \dots, 1, \dots, 0)$ where the 1 appears at the i th position. By implementing QR-updating in two stages: 1) Rank-One QR-updating; 2) QR-updating as the case of insertion with a new row, we have the updated QR-decomposition so that $\mathbf{Q}'\mathbf{R}' = \mathbf{M}'$.

The updating of the reduced within-class scatter matrix \mathbf{W} is

$$\mathbf{W}' = \mathbf{W} + (\boldsymbol{\alpha} - \boldsymbol{\beta})(\boldsymbol{\alpha} - \boldsymbol{\beta})^T + n_i\boldsymbol{\beta}\boldsymbol{\beta}^T, \quad (3.3)$$

where $\boldsymbol{\alpha} = \mathbf{Q}'^T(\tilde{\mathbf{x}} - \mathbf{m}_i)$ and $\boldsymbol{\beta} = \mathbf{Q}'^T(\mathbf{m}'_i - \mathbf{m}_i)$.

The updating of the reduced between-class scatter matrix \mathbf{B} is

$$\mathbf{B}' = (\mathbf{R}'\mathbf{D} - (\frac{1}{n'}\mathbf{R}' \cdot r) \cdot h^T)(\mathbf{R}'\mathbf{D} - (\frac{1}{n'}\mathbf{R}' \cdot r) \cdot h^T)^T, \quad (3.4)$$

where $\mathbf{D} = \text{diag}(\sqrt{n'_1}, \dots, \sqrt{n'_k})$, $h = [\sqrt{n'_1}, \dots, \sqrt{n'_k}]^T$, and $r = (n'_1, \dots, n'_k)$

2. if $\tilde{\mathbf{x}}$ belongs to a new class (i.e., $i > k$), then for $i = 1, \dots, k$, $n'_i = n_i$, $n'_{k+1} = 1$, and $n' = n + 1$. The updated centroid matrix is

$$\mathbf{M}' = [\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_k, \tilde{\mathbf{x}}] = [\mathbf{M}, \tilde{\mathbf{x}}]. \quad (3.5)$$

Then, the updated reduced within-class scatter matrix is

$$\mathbf{W}' = \begin{pmatrix} \mathbf{W} & 0 \\ 0 & 0 \end{pmatrix} \quad (3.6)$$

The reduced between-class scatter matrix can be updated similarly to (3.4) as

$$\mathbf{B}' = (\mathbf{R}'\mathbf{D} - (\frac{1}{n'}\mathbf{R}' \cdot r) \cdot h^T)(\mathbf{R}'\mathbf{D} - (\frac{1}{n'}\mathbf{R}' \cdot r) \cdot h^T)^T, \quad (3.7)$$

with $\mathbf{D} = \text{diag}(\sqrt{n'_1}, \dots, \sqrt{n'_{k+1}})$, $h = [\sqrt{n'_1}, \dots, \sqrt{n'_{k+1}}]^T$, and $r = (n'_1, \dots, n'_{k+1})$

The steps of sInclDA above are summarized in algorithm 9.

Algorithm 9 Sequential IDR/QR IncLDA Algorithm

Input: Newly added training data sample $\tilde{\mathbf{x}}$ and its class label \tilde{c} , LDA/sIncLDA eigenspace model

$$\Omega = \{\mathbf{M}, \mathbf{W}, \mathbf{B}, \mathbf{Q}, \mathbf{R}, \mathbf{C}\}.$$

Output: Optimal transform matrix ϕ' , sIncLDA eigenspace model $\Omega' = \{\mathbf{M}', \mathbf{W}', \mathbf{B}', \mathbf{Q}', \mathbf{R}', \mathbf{C}'\}$.

1: **if** $\tilde{c} \in \mathbf{C}$ **then**

2: Calculate centroid matrix \mathbf{M}' by (3.2);

3: Calculate reduced between-class matrix \mathbf{B}' by (3.4);

4: Calculate reduced within-class matrix \mathbf{W}' by (3.3);

5: **else**

6: Calculate centroid matrix \mathbf{M}' by(3.5);

7: Calculate reduced between-class matrix \mathbf{B}' by (3.7);

8: Calculate reduced within-class matrix \mathbf{W}' by (3.6);

9: **end if**

10: Calculate k' eigenvectors $\varphi_{i'}$ of $(\mathbf{W}' + \mu\mathbf{I}_{k'})^{-1}\mathbf{B}'$ with decreasing eigenvalues;

11: Calculate optimal transform matrix $\phi' = \mathbf{Q}'\mathbf{Z}'$.

3.3 Discussion

It is worth noting that when data is presented in a chunk manner, the sIncLDA performs learning inefficiently for a single sample at a time. This is because a large number of iterative updates are carried out in sIncLDA for just one sample, and this might be discriminatively redundant (i.e., gives no contribution to the existing discriminant model). Furthermore, the class separability of the sIncLDA deducts for datasets with a large class number. This is because when the \mathbf{W} is being updated, discriminative information loss occurs, 1) when the newly presented sample belongs to an existing class, the sIncLDA has the assumption of $\mathbf{Q}' \approx \mathbf{Q}$ in (3.3); and 2) when the newly presented sample comes from a new class, the sIncLDA has another assumption of $\mathbf{W}' \approx \mathbf{W}$ in (3.6). However, these two assumptions often lose generality in practice, especially when the class number is large.

Chapter 4

Proposed Chunk Incremental IDR/QR LDA

In this chapter, we propose our new chunk incremental IDR/QR LDA. Given an eigenspace model $\Omega = \{\mathbf{M}, \mathbf{W}, \mathbf{B}\}$, and a chunk of samples $\widetilde{\mathbf{X}}$ in \widetilde{k} classes, regardless of existing or entirely new categories. Simply put, the set of total class labels is updated as $\mathbf{C}' = \mathbf{C} \cup (\widetilde{\mathbf{C}} - \mathbf{C} \cap \widetilde{\mathbf{C}})$ and $k' = |\mathbf{C}'|$. The number of total training samples is renewed as $n' = n + \widetilde{n}$. Then, the problem of the proposed IncLDA can be described as

$$\Omega' = \mathcal{F}_{ck}(\Omega, \widetilde{\mathbf{X}}) = \{\mathbf{M}', \mathbf{W}', \mathbf{B}'\}.$$

\mathcal{F}_c updates Ω on $\widetilde{\mathbf{X}}$ in three steps: (1) Updating centroid matrix \mathbf{M} and its QR decomposition; (2) Updating the reduced within-class scatter matrix \mathbf{W} ; and (3) Updating the reduced between-class scatter matrix \mathbf{B} . Note that differing from Sequential IDR/QR IncLDA, Chunk IDR/QR IncLDA takes care of newly added data from existing and new classes simultaneously.

4.1 Updating the Centroid Matrix and its QR-decomposition

To update $\mathbf{M} \in \mathcal{R}^{d \times k}$ to $\mathbf{M}' \in \mathcal{R}^{d \times k'}$, the dimensionality of \mathbf{M} is augmented from k to k' for accommodating new classes in the centroid matrix. In doing so, we insert zero vectors $\mathbf{0}$ at i th column of \mathbf{M} for $i = k + 1, \dots, k'$ and $\widetilde{\mathbf{M}}$ for $i = \widetilde{k} + 1, \dots, k'$. Thus, $\mathbf{M} \in \mathcal{R}^{d \times k}$ is augmented to $\widehat{\mathbf{M}} \in \mathcal{R}^{d \times k'}$ upon the presence of $\widetilde{\mathbf{X}}$. Here, ‘ $\widehat{}$ ’

identifies the result of a matrix augmentation, which will be applied throughout the paper.

Applying matrix augmentation to $\mathbf{M} = \mathbf{QR}$, we have $\hat{\mathbf{M}} = \hat{\mathbf{Q}}\hat{\mathbf{R}}$. $\hat{\mathbf{Q}}\hat{\mathbf{R}}$ can be obtained by recomputing a QR decomposition on $\hat{\mathbf{M}}$, or by a QR-updating on \mathbf{QR} . However, both approaches are computationally expensive due to the high complexity of matrix decomposition or multiplication. For efficient $\hat{\mathbf{Q}}\hat{\mathbf{R}}$ augmentation, we propose the rule of matrix augmentation for QR decomposition as Proposition 1, which involves only manipulations of matrix partition or combination.

Proposition 1 *Given $\mathbf{A} = (\mathbf{A}_1, \mathbf{A}_2) = \mathbf{QR} \in \mathcal{R}^{d \times k}$ with $\mathbf{A}_1 \in \mathcal{R}^{d \times (i-1)}$, and $\mathbf{A}_2 \in \mathcal{R}^{d \times (k-i)}$. If a zero vector $\mathbf{0}$ is inserted as $\hat{\mathbf{A}} = (\mathbf{A}_1, \mathbf{0}, \mathbf{A}_2)$, then the QR decomposition $\hat{\mathbf{Q}}\hat{\mathbf{R}}$ of $\hat{\mathbf{A}}$ can be computed as*

$$\hat{\mathbf{Q}}\hat{\mathbf{R}} = (\mathbf{Q}_1, \mathbf{0}, -\mathbf{Q}_2) \begin{pmatrix} \mathbf{R}_1 & \mathbf{0} & \mathbf{R}_2 \\ \mathbf{0}^T & 1 & \mathbf{0}^T \\ \mathbf{R}_3 & \mathbf{0} & -\mathbf{R}_4 \end{pmatrix}$$

where $\mathbf{Q}_1 \in \mathcal{R}^{d \times (i-1)}$, $\mathbf{Q}_2 \in \mathcal{R}^{d \times (k-i)}$, $\mathbf{R}_1 \in \mathcal{R}^{(i-1) \times (i-1)}$, $\mathbf{R}_2 \in \mathcal{R}^{(i-1) \times (k-i)}$, $\mathbf{R}_3 \in \mathcal{R}^{(k-i) \times (i-1)}$, and $\mathbf{R}_4 \in \mathcal{R}^{(k-i) \times (k-i)}$. Note that, $\mathbf{Q} = (\mathbf{Q}_1, \mathbf{Q}_2)$, $\mathbf{R} = \begin{pmatrix} \mathbf{R}_1 & \mathbf{R}_2 \\ \mathbf{R}_3 & \mathbf{R}_4 \end{pmatrix}$, and \mathbf{R}_3 is a zero matrix. The detailed proof is given in Appendix A.

According to the definition of centroid matrix, we reformulate $\mathbf{M}' = [\frac{\mathbf{m}_1 n_1 + \tilde{\mathbf{m}}_1 \tilde{n}_1}{n'_1}, \dots, \frac{\mathbf{m}_{k'} n_{k'} + \tilde{\mathbf{m}}_{k'} \tilde{n}_{k'}}{n'_{k'}}]$ as,

$$\mathbf{M}' = \hat{\mathbf{M}} + \Delta. \quad (4.1)$$

Then, we have $\Delta = [\frac{\tilde{n}_1}{n'_1}(\tilde{\mathbf{m}}_1 - \mathbf{m}_1), \dots, \frac{\tilde{n}_{k'}}{n'_{k'}}(\tilde{\mathbf{m}}_{k'} - \mathbf{m}_{k'})] \in \mathcal{R}^{d \times k'}$. Note that Δ here identifies the residue information of newly presented against existing data.

To calculate $\mathbf{Q}'\mathbf{R}'$, we apply the thin Singular Value Decomposition on Δ , namely $\Delta = \mathbf{U}\Sigma\mathbf{V}^T$ and obtain a summation of rank one matrices products (Golub & Van Loan, 1996) as,

$$\Delta = \sum_{j=1}^{\eta} \Sigma(j, j) \mathbf{U}(1 : d, j) \mathbf{V}^T(j, 1 : k') = \sum_{j=1}^{\eta} \mathbf{u}_j \mathbf{v}_j^T \quad (4.2)$$

where $\eta = \text{rank}(\Delta)$, $\mathbf{u}_j = \Sigma(j, j) \mathbf{U}(1 : d, j)$ and $\mathbf{v}_j = \mathbf{V}(j, 1 : k')$.

Substituting (4.2) into (4.1), we can calculate $\mathbf{Q}'\mathbf{R}'$ by an iterative QR-updating procedure as,

$$\begin{aligned}
\mathbf{Q}'\mathbf{R}' &= \hat{\mathbf{Q}}^{(\eta+1)} \hat{\mathbf{R}}^{(\eta+1)} \\
&= \hat{\mathbf{Q}}^{(\eta)} \hat{\mathbf{R}}^{(\eta)} + \mathbf{u}_\eta \mathbf{v}_\eta^T \\
&\dots \\
&= (\hat{\mathbf{Q}}^{(j)} \hat{\mathbf{R}}^{(j)} + \mathbf{u}_j \mathbf{v}_j^T) + \sum_{j=j+1}^\eta \mathbf{u}_j \mathbf{v}_j^T \\
&\dots \\
&= (\hat{\mathbf{Q}}^{(2)} \hat{\mathbf{R}}^{(2)} + \mathbf{u}_2 \mathbf{v}_2^T) + \sum_{j=3}^\eta \mathbf{u}_j \mathbf{v}_j^T \\
&= (\hat{\mathbf{Q}}^{(1)} \hat{\mathbf{R}}^{(1)} + \mathbf{u}_1 \mathbf{v}_1^T) + \sum_{j=2}^\eta \mathbf{u}_j \mathbf{v}_j^T \\
&= \hat{\mathbf{Q}} \hat{\mathbf{R}} + \sum_{j=1}^\eta \mathbf{u}^j \mathbf{v}_j^T.
\end{aligned} \tag{4.3}$$

Here, we define the recursive function \mathcal{P} such that $\hat{\mathbf{Q}}^{(j+1)} \hat{\mathbf{R}}^{(j+1)} = \mathcal{P}(\hat{\mathbf{Q}}^{(j)} \hat{\mathbf{R}}^{(j)})$. Thus, for each iteration of (4.3),

$$\begin{aligned}
\mathcal{P}(\hat{\mathbf{Q}}^{(j)} \hat{\mathbf{R}}^{(j)}) &= \hat{\mathbf{Q}}^{(j)} \hat{\mathbf{R}}^{(j)} + \mathbf{u}_j \mathbf{v}_j^T \\
&= \hat{\mathbf{Q}}^{(j)} (\hat{\mathbf{R}}^{(j)} + \hat{\mathbf{Q}}^{(j)T} \mathbf{u}_j \mathbf{v}_j^T) + (\mathbf{I} - \hat{\mathbf{Q}}^{(j)} \hat{\mathbf{Q}}^{(j)T}) \mathbf{u}_j \mathbf{v}_j^T \\
&= \hat{\mathbf{Q}}^{(j)} (\hat{\mathbf{R}}^{(j)} + \mathbf{w}^{(j)} \mathbf{v}_j^T) + \mathbf{f}^{(j)} \mathbf{v}_j^T,
\end{aligned} \tag{4.4}$$

where $j = [1, \dots, \eta]$, and $\mathbf{w}^{(j)} = \hat{\mathbf{Q}}^{(j)T} \mathbf{u}_j$ and $\mathbf{f}^{(j)} = (\mathbf{I} - \hat{\mathbf{Q}}^{(j)} \hat{\mathbf{Q}}^{(j)T}) \mathbf{u}_j$. Now, we can calculate $\mathbf{Q}'\mathbf{R}'$ by applying the two-stage QR-updating method on (4.4) as in Ye et al. (2005).

4.2 Updating the Reduced Within-Class Scatter Matrix

To update the reduced within-class scatter matrix \mathbf{W} , we augment, similar to the updating of \mathbf{M} , the k -dimension \mathbf{W} to k' -dimension $\hat{\mathbf{W}}$. Again to avoid matrix multiplication, we simplify the augmentation of \mathbf{W} following the rule of Proposition 2:

Proposition 2 *Given $\mathbf{W} = \mathbf{Q}^T \mathbf{H}_w \mathbf{H}_w^T \mathbf{Q} \in \mathcal{R}^{k \times k}$, if a zero vector $\mathbf{0}$ is inserted as the i -th column into $\mathbf{Q} \in \mathcal{R}^{d \times k}$, such that $\hat{\mathbf{Q}} = (\mathbf{Q}_1, \mathbf{0}, -\mathbf{Q}_2) \in \mathcal{R}^{d \times (k+1)}$ where*

$\mathbf{Q}_1 \in \mathcal{R}^{d \times (i-1)}$ and $\mathbf{Q}_2 \in \mathcal{R}^{d \times (k-i)}$, then we have

$$\begin{aligned} \hat{\mathbf{W}} &= \hat{\mathbf{Q}}^T \mathbf{H}_w \mathbf{H}_w^T \hat{\mathbf{Q}} \\ &= \begin{pmatrix} \mathbf{W}_1 & \mathbf{0} & -\mathbf{W}_2 \\ \mathbf{0}^T & 0 & \mathbf{0}^T \\ -\mathbf{W}_3 & \mathbf{0} & \mathbf{W}_4 \end{pmatrix}, \end{aligned}$$

where $\mathbf{W}_1 \in \mathcal{R}^{(i-1) \times (i-1)}$, $\mathbf{W}_2 \in \mathcal{R}^{(i-1) \times (k-i)}$, $\mathbf{W}_3 \in \mathcal{R}^{(k-i) \times (i-1)}$, $\mathbf{W}_4 \in \mathcal{R}^{(k-i) \times (k-i)}$, and $\mathbf{W} = \begin{pmatrix} \mathbf{W}_1 & \mathbf{W}_2 \\ \mathbf{W}_3 & \mathbf{W}_4 \end{pmatrix}$. The detailed proof is given in Appendix B.

Based on $\hat{\mathbf{W}}$, we can compute the reduced within-class scatter matrix \mathbf{W}' by the following proposition 3.

Proposition 3 *Given an augmented reduced within-class scatter matrix $\hat{\mathbf{W}}$ and a set of new training samples $\tilde{\mathbf{X}}$, the reduced within-class scatter matrix can be updated as*

$$\mathbf{W}' = \bar{\mathbf{W}} + \tilde{\mathbf{W}} + \mathbf{E} \quad (4.5)$$

where $\bar{\mathbf{W}} = \mathbf{Q}'^T \mathbf{H}_w \mathbf{H}_w^T \mathbf{Q}'$, $\tilde{\mathbf{W}} = \mathbf{Q}'^T \tilde{\mathbf{H}}_w \tilde{\mathbf{H}}_w^T \mathbf{Q}'$, $\mathbf{E} = \sum_{i=1}^{k'} (\tilde{n}_i (\boldsymbol{\alpha}' - \boldsymbol{\beta}') (\boldsymbol{\alpha}' - \boldsymbol{\beta}')^T + n_i \boldsymbol{\beta}' \boldsymbol{\beta}'^T)$, $\boldsymbol{\alpha}' = \mathbf{Q}'^T (\tilde{\mathbf{m}}_i - \mathbf{m}_i)$, and $\boldsymbol{\beta}' = \mathbf{Q}'^T (\mathbf{m}'_i - \mathbf{m}_i)$. The detailed proof is given in Appendix C.

Consider $\bar{\mathbf{W}} = \mathbf{Q}'^T \mathbf{H}_w \mathbf{H}_w^T \mathbf{Q}'$ and $\mathbf{Q}^{(j+1)} = \mathcal{P}(\mathbf{Q}^{(j)})$, there should exist an iterative function \mathcal{P}_W for \mathbf{W}

$$\mathbf{W}^{(j+1)} = \mathcal{P}_W(\mathbf{W}^{(j)}), j = 1, \dots, \eta \quad (4.6)$$

such that $\bar{\mathbf{W}} = \mathbf{W}^{(\eta+1)}$ when $j = \eta$.

Following (4.6), we have

$$\begin{aligned} \mathbf{W}^{(j+1)} &= \hat{\mathbf{Q}}^{(j+1)T} \mathbf{H}_w \mathbf{H}_w^T \hat{\mathbf{Q}}^{(j+1)} \\ &= [\mathcal{P}(\hat{\mathbf{Q}}^{(j)T})] \mathbf{H}_w \mathbf{H}_w^T [\mathcal{P}(\hat{\mathbf{Q}}^{(j)})] \\ &\approx \mathbf{G}^{(j)T} \cdot \mathbf{J}^{(j)T} \cdot \hat{\mathbf{Q}}^{(j)T} \mathbf{H}_w \mathbf{H}_w^T \hat{\mathbf{Q}}^{(j)} \cdot \mathbf{J}^{(j)} \cdot \mathbf{G}^{(j)} \\ &= \mathbf{G}^{(j)T} \cdot \mathbf{J}^{(j)T} \cdot \hat{\mathbf{W}}^{(j)} \cdot \mathbf{J}^{(j)} \cdot \mathbf{G}^{(j)}, \end{aligned} \quad (4.7)$$

where $\mathbf{G}^{(j)} = [\mathbf{G}_{1,2}^{(j)}, \dots, \mathbf{G}_{k'-1,k'}^{(j)}]$ and $\mathbf{J}^{(j)} = [\mathbf{J}_{k'-1,k'}^{(j)}, \dots, \mathbf{J}_{1,2}^{(j)}]$ are two sets of Givens rotations obtained from the first QR-updating stage; “.” denotes that applying the Givens rotations on $\hat{\mathbf{W}}^{(j)}$ one by one. Note that the approximation in (4.7) assumes that $\hat{\mathbf{W}}^{(j)}$ is not varied in the second QR-updating stage.

4.3 Updating the Reduced Between-Class Scatter Matrix

We use the same process for updating \mathbf{B} as in (Ye et al., 2005), which can be written as

$$\mathbf{B}' = \left(\mathbf{R}'\mathbf{D} - \left(\frac{1}{n'}\mathbf{R}'\mathbf{r}\right)\mathbf{h}^T \right) \left(\mathbf{R}'\mathbf{D} - \left(\frac{1}{n'}\mathbf{R}'\mathbf{r}\right)\mathbf{h}^T \right)^T, \quad (4.8)$$

where $\mathbf{D} = \text{diag}(\sqrt{n'_1}, \dots, \sqrt{n'_{k'}})$, $\mathbf{r} = (n'_1, \dots, n'_{k'})^T$, and $\mathbf{h} = [\sqrt{n'_1}, \dots, \sqrt{n'_{k'}}]^T$.

4.4 The Pseudocode of the Proposed Algorithm

The pseudocode of the proposed algorithm is presented in algorithm 10.

4.5 Time Complexity Analysis

This section provides the time complexity analysis of the proposed cInclDA with a comparison to the sInclDA. The term ‘flam’, which denotes an addition and a multiplication, is used for presenting operation counts (G.W., 1998). Consider that both InclDAs calculate an update of the 3-tuple model $\{\mathbf{M}, \mathbf{W}, \mathbf{B}, \}$. We measure the time complexity of each variable update for every single new sample presented for incremental learning. Here, for each variable updating, we only consider the maximum time cost regardless of whatever class the incoming sample belongs to.

For sInclDA, the maximum time cost occurs when the new sample belongs to an existing class, since this involves more complicated calculations on \mathbf{W} and \mathbf{QR} than when a new sample comes from a new class. In this case, the updating of \mathbf{M} requires dk flam, namely $O(dk)$. According to Daniel, Gragg, Kaufman and Stewart (1976), total operation counts for the two-stage QR-updating are $9dk + \frac{9}{2}k^2$ flam, so its time complexity is written as $O(dk + k^2)$. With regard to the scatter matrices, it requires $dk + 2k^2 + k$ flam for updating \mathbf{W} and $\frac{3}{2}k^3 + 2k^2$ flam for \mathbf{B} . Hence, the

Algorithm 10 Chunk IDR/QR IncLDA Algorithm

Input: Newly added training data matrix $\tilde{\mathbf{X}}$ and its class label set $\tilde{\mathbf{C}}$, LDA/cIncLDA eigenspace model $\Omega = \{\mathbf{M}, \mathbf{W}, \mathbf{B}, \mathbf{Q}, \mathbf{R}, \mathbf{C}\}$.

Output: Optimal transform matrix ϕ' , cIncLDA eigenspace model $\Omega' = \{\mathbf{M}', \mathbf{W}', \mathbf{B}', \mathbf{Q}', \mathbf{R}', \mathbf{C}'\}$.

- 1: Calculate current class label set $\mathbf{C}' = \mathbf{C} \cup (\tilde{\mathbf{C}} - \mathbf{C} \cap \tilde{\mathbf{C}})$, $\tilde{k} = \|\tilde{\mathbf{C}}\|$, and $k' = \|\mathbf{C}'\|$;
- 2: Calculate centroid matrix $\tilde{\mathbf{M}} = [\tilde{\mathbf{m}}_1, \dots, \tilde{\mathbf{m}}_{\tilde{k}}]$, and within-class scatter matrix $\tilde{\mathbf{S}}_{\mathbf{w}} = \sum_{i=1}^{\tilde{k}} \sum_{j=1}^{\tilde{n}_i} (\tilde{\mathbf{x}}_{ij} - \tilde{\mathbf{m}}_i)(\tilde{\mathbf{x}}_{ij} - \tilde{\mathbf{m}}_i)^T$;
- 3: **if** $k' > \tilde{k}$ **then**
- 4: Conduct matrix augmentation on $\tilde{\mathbf{M}}$ by inserting $k' - \tilde{k}$ zero vectors;
- 5: **end if**
- 6: **if** $k' > k$ **then**
- 7: Conduct matrix augmentation on $\mathbf{M}, \mathbf{Q}, \mathbf{R}, \mathbf{W}$ by Proposition 1 and Proposition 2;
- 8: **end if**
- 9: Calculate centroid matrix \mathbf{M}' and residue matrix Δ by (4.1);
- 10: Calculate the rank $\eta = \text{rank}(\Delta)$;
- 11: Perform SVD of Δ as, $\Delta = \mathbf{U}\Sigma\mathbf{V}^T$;
- 12: **for** $j = 1$ to η **do**
- 13: Calculate vector $\mathbf{w} = \mathbf{Q}^T \Sigma(j, j) \mathbf{U}(:, j)$, and $\mathbf{v} = \mathbf{Q}^T \mathbf{V}(j, :)$;
- 14: Calculate $\mathbf{Q}'\mathbf{R}'$ by updating QR-decomposition as (4.3) and (4.4);
- 15: Calculate matrix $\tilde{\mathbf{W}}$ by (4.7);
- 16: **end for**
- 17: Calculate \mathbf{E} and $\tilde{\mathbf{W}}$ by Proposition 3;
- 18: Calculate reduced within-class scatter matrix \mathbf{W}' by Proposition 3;
- 19: Calculate reduced between-class scatter matrix \mathbf{B}' by (4.8);
- 20: Calculate k' eigenvectors φ'_i of $(\mathbf{W}' + \mu \mathbf{I}_{k'})^{-1} \mathbf{B}'$ with decreasing eigenvalues;
- 21: Calculate optimal transform matrix $\phi' = \mathbf{Q}'\mathbf{Z}'$, where $\mathbf{Z}' = [\varphi'_1, \dots, \varphi'_{k'}]$.

time complexity for updating \mathbf{W} and \mathbf{B} is simplified as $O(dk + k^2)$ and $O(k^3 + k^2)$, respectively.

In the proposed cIncLDA, a matrix augmentation operation is employed before all updates. We have $\mathbf{M} \in \mathcal{R}^{d \times k}$, $\tilde{\mathbf{M}} \in \mathcal{R}^{d \times \tilde{k}}$, $\mathbf{Q} \in \mathcal{R}^{d \times k}$, $\mathbf{R} \in \mathcal{R}^{k \times k}$ and $\mathbf{W} \in \mathcal{R}^{k \times k}$ augmented from k or \tilde{k} to k' , so the total time complexity for augmentation is $O(dk' + k'^2)$. Updating $\hat{\mathbf{M}}$ requires dk' flam whose time complexity is $O(dk')$. For

Main step	Sequential IDR/QR IncLDA	Chunk IDR/QR IncLDA
Matrices Augmentation	NO	$O(dk' + k'^2)$
Update \mathbf{M}	$O(dk)$	$O(dk')$
Update \mathbf{Q} and \mathbf{R}	$O(dk + k^2)$	$O(dk'^2 + k'^3 + dk' + k'^2)$
Update \mathbf{W}	$O(dk + k^2)$	$O(k'^3 + dk' + k'^2)$
Update \mathbf{B}	$O(k^3 + k^2)$	$O(k'^3 + k'^2)$
General	$O(k^3 + dk)$	$O(dk'^2 + k'^3 + dk')$

Table 4.1: Comparison of Time Complexity for incremental learning on a single sample, sIncLDA vs cIncLDA

updating the $\hat{\mathbf{Q}}\hat{\mathbf{R}}$, the SVD on $\mathbf{\Delta}$ requires $4dk'^2 + 8k'^3$ flam according to G.W. (1998), and the following two-stage QR-updating on $\hat{\mathbf{Q}}\hat{\mathbf{R}}$ is a η -iteration process, which requires $9\eta dk' + \frac{9}{2}\eta k'^2$ flam. Adding the above two steps, it gives the total time complexity $O(dk'^2 + k'^3 + dk' + k'^2)$ for updating $\hat{\mathbf{Q}}\hat{\mathbf{R}}$. As discussed in Section 4.2, the update of the reduced within-class matrix $\hat{\mathbf{W}}$ follows Proposition 3, in which the operation counts for computing $\bar{\mathbf{W}}$ by (4.6) and (4.7) are $6\eta k'^2 + k'^3 + dk' + 6\eta k' + k'^2$ flam, and those for computing $\tilde{\mathbf{W}}$ and \mathbf{E} are $d'k' + k'^2$ and k'^3 flam. Hereby for updating $\hat{\mathbf{W}}$, the total operation counts are $2k'^3 + 2dk' + 6\eta k'^2 + 2k'^2 + 6\eta k'$ flam, namely $O(k'^3 + dk' + k'^2)$. For updating the reduced between-class scatter matrix \mathbf{B} , we can calculate the time complexity similarly to sIncLDA as $O(k'^3 + k'^2)$ by (4.8).

The time complexity discussed above is summarized in Table 4.1. As seen in Table 4.1, for incremental learning on a single sample, the general complexity of cIncLDA $O(dk'^2 + k'^3 + dk')$ is higher than the sIncLDA $O(k^3 + dk)$ since $d \gg k$ and $k' \geq k$.

However, for incremental learning on a chunk of s samples, sIncLDA can only process the data iteratively in s cycles. Thus, it gives $O(s(k^3 + dk))$. In contrast, the proposed cIncLDA processes whole chunk s samples at one time, and the actual time cost by Proposition 3 is determined by η , the rank of $\mathbf{\Delta}$. It follows that s samples are compressed informatively into η samples in the proposed cIncLDA. Thus, the time complexity for cIncLDA is $O(\eta(dk'^2 + k'^3 + dk'))$.

Matching the running efficiency of the cIncLDA with sIncLDA, $s(k^3 + dk) = \eta(dk'^2 + k'^3 + dk')$, since $k \leq k'$, we have the minimal chunk size s as,

$$s = \frac{\eta(dk'^2 + k'^3 + dk')}{k^3 + dk} \geq \frac{\eta k'^2 d}{k^3 + dk} + \eta. \quad (4.9)$$

Here, s is greater than η as proved below, and it follows that the cIncLDA matches

the sIncLDA efficiency when the chunk size $s > \eta$. Because s often grows much faster than the η for cIncLDA, the general efficiency of cIncLDA (in terms of the total time cost to complete incremental learning of a whole dataset) is expected to surpass the sIncLDA with the increase of the chunk size s .

Proof: Consider $k \ll d$, without loss of generality, we enlarge the denominator of (4.9) by replacing the k^3 item with dk^2 , thus we have

$$s \geq \frac{\eta k'^2 d}{dk^2 + dk} + \eta = \frac{\eta k'^2}{k(k+1)} + \eta > \eta \quad (4.10)$$

Chapter 5

Experiment Evaluation

In this chapter, we examine the accuracy and efficiency of the proposed cInclLDA by comparing it with sInclLDA and batch QR/LDA. We have specifically examined its equality to the ground truth eigenspace, its class separability of the embedding eigenspace, and its execution time. All the experiments were conducted on a Intel 2.26GHZ Core I5 PC with 4GB Ram.

5.1 Data Description

As we have discussed previously, the performance of sInclLDA may suffer when data with a large number of classes are presented. Thus, we need to choose datasets with diverse class numbers to examine the performance of both InclLDAs. In addition, to our best, we can find these five benchmarked datasets, which are widely used in other related work (T.-K. Kim, Kim, Hwang & Kittler, 2005; T.-K. Kim et al., 2011; Pang, Kim & Bang, 2003; M.-S. Kim, Kim & Lee, 2003; Pang et al., 2004). Therefore, the following five benchmarked facial recognition datasets and one combined face dataset have been selected for the following performance and efficiency evaluation experiments. Those have different class numbers from 40 to 1010.

1. The AT&T (ORL) face dataset ¹ consists of 400 face images from 40 persons (10 images per person). Each image has the size of 92×112 . The images were taken at the same background but at different times, with varied facial expressions (e.g., open/closed eyes and smiling/non-smiling) or facial details (e.g., glasses/no-glasses).

¹<http://www.cl.cam.ac.uk/research/dtg/attarchive/facedatabase.html>.

2. The Altkom dataset ² (T.-K. Kim et al., 2005, 2011) is composed of 1200 face images from 80 persons (15 images per person) with different poses. The size of image is normalized to 46×56 .
3. The MPEG-7 dataset has 1355 face images from 271 persons (five images per person) (Pang et al., 2003; M.-S. Kim et al., 2003; Pang et al., 2004). Each image has the size of 23×28 pixels. The dataset is collected from AR(Purdue), AT&T, Yale, UMIST, University of Berne, and some of them were obtained from MPEG-7 news Videos.
4. The XM2VTS dataset ³ (T.-K. Kim et al., 2005, 2011) contains 2950 face images from five recordings of 295 persons in which each recording has two head shots taken at two different time (i.e., 10 images per person). Each image has the same size of 46×56 .
5. The extended version 1 MPEG dataset ⁴ (T.-K. Kim et al., 2005, 2011) is a combination of several public face sets (e.g., AR and ORL). It collects 3175 images from 635 persons (five images per person) showing illumination and view variations. The size of image is the same, 46×56 .
6. The combined face dataset includes 5050 images from 1010 persons (five images per person), in which we selected 400 images from the Altkom dataset (80 persons), 1475 images from XM2VTS dataset (295 persons), and 3175 images from the version 1 MPEG dataset (635 persons). The size of image is kept the same, 46×56 .

5.2 Experimental setup

To conduct incremental learning on each dataset, we first construct an initial discriminant eigenspace using 10% of the face images, in which at least two classes of data are guaranteed to be included according to the definition of classic LDA. The remaining 90% training instances are divided equally into nine chunks/subsets and presented to the learner sequentially. As a reference, we also apply batch QR/LDA

²<http://www.iis.ee.ic.ac.uk/~tkkim/code.htm>.

³<http://www.iis.ee.ic.ac.uk/~tkkim/code.htm>.

⁴<http://www.iis.ee.ic.ac.uk/~tkkim/code.htm>.

to perform the same incremental learning task. As the incremental learning proceeds, the learner is provided with chunks of new data in a sequential manner until all subsets are exhausted.

5.3 Results Evaluation

5.3.1 Similarity to Ground Truth Eigenspace

To examine the similarity between the discriminant eigenspace from the cInclLDA and that from the ground truth QR/LDA, we set the chunk size s as 1, and use the cInclLDA for sequential incremental learning. We experiment cInclLDA, sInclLDA and the batch LDA on the UCI Iris dataset (Frank & Asuncion, 2010) which has three classes and 150 samples in a 4-dimension space. We calculate the inner products between the discriminant eigenvectors from the cInclLDA and the counterpart of the batch QR/LDA at every learning stage, and record the converging procedure of the cInclLDA and sInclLDA, respectively.

As a result, Fig. 5.1 reveals the difference between two InclLDAs on their converging procedure to batch QR/LDA. As seen, both InclLDAs converge to the Batch QR/LDA as the incremental learning proceeds, demonstrating the cInclLDA equivalent to the ground truth on Iris data. The proposed cInclLDA has not only fewer but also far smaller fluctuations than the sInclLDA on all three axes. Specifically on the d_1 axis, the cInclLDA is found converging to the ground truth over 100 stages earlier than the sInclLDA; and on the other two axes performing similarly to the sInclLDA. This suggests that with reference to the ground truth, the proposed cInclLDA has deterministically less discriminative information lost for incremental QR/LDA learning, and thus is expected to have better stage learning performance (in terms of stage class separability evaluation) than the sInclLDA even if it is being used just in sequential learning mode.

5.3.2 Class Separability

To evaluate the class separability of the proposed cInclLDA, we project the data presented so far to the current LDA eigenspace and then classify the data using a k -Nearest Neighbor (k NN) classifier. The dimensions of the embedding spaces of both reported methods are set to the number of classes, and the classification accuracy is

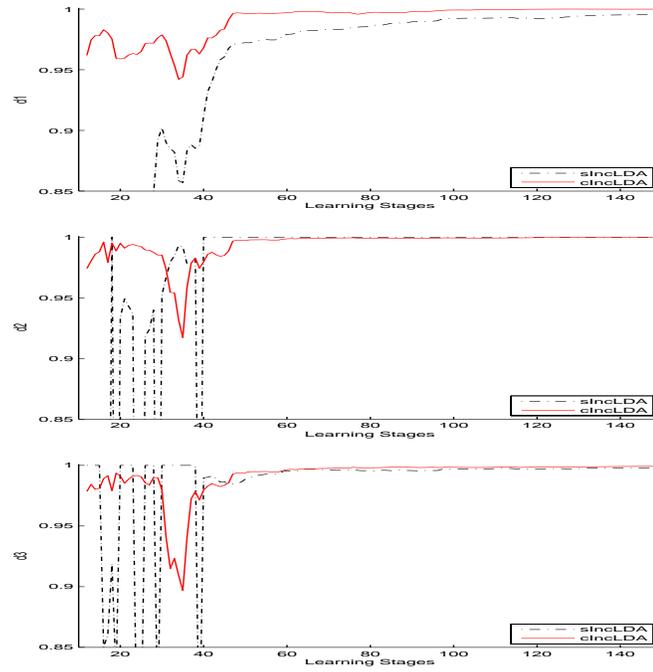
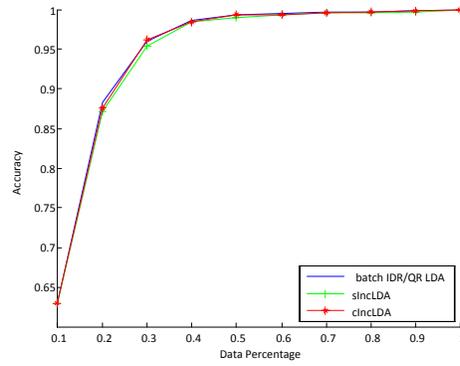


Figure 5.1: inner products of eigenvectors, cIncLDA vs. sIncLDA on their converging procedure to batch QR/LDA

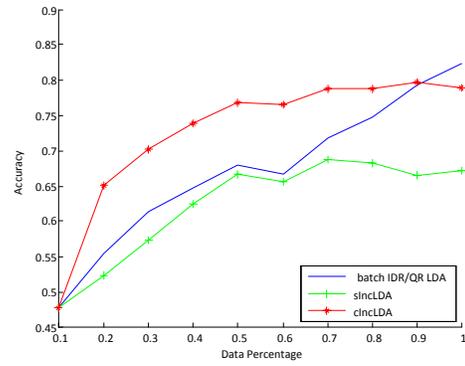
measured by leave-one-out (LOO) cross-validation.

Fig 5.2 shows the experimental results, where the proposed cIncLDA is compared to sIncLDA and batch QR/LDA, at every incremental learning stage, which is identified as the percentage of data presented so far. As seen in the figure, the cIncLDA resembles the trend of sIncLDA for all six datasets; however the performance of the cIncLDA apparently approaches the ground truth batch QR/LDA better than the sIncLDA. Such superiority arises at the beginning stage and amplifies quickly in a few cycles of learning. For dataset Altkom and XM2V., it grows to even outperform the batch QR/LDA. Most importantly, it turns to be stable at later stages to the end of incremental learning.

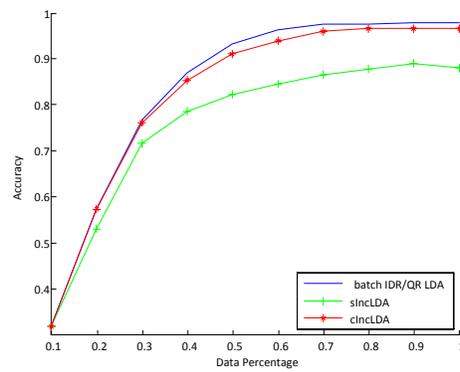
We further measure the statistical difference from the proposed cIncLDA to the existing sIncLDA on class separability by calculating the average stage classification accuracy difference (i.e., stage Acc. diff.) and the final stage classification accuracy difference (i.e., final Acc. diff.). Table 5.1 gives the comparison for all six case studies, in which the quantitative properties of each dataset are detailed, and the class separability differences are calculated as the sIncLDA minus cIncLDA on classification rate.



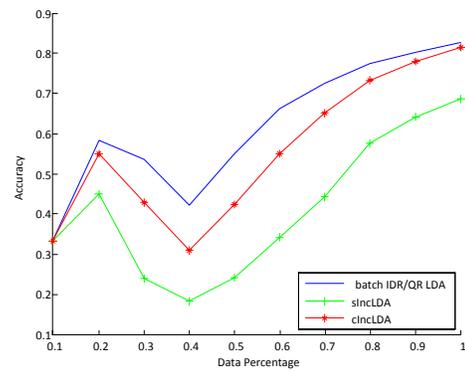
(a) AT&T



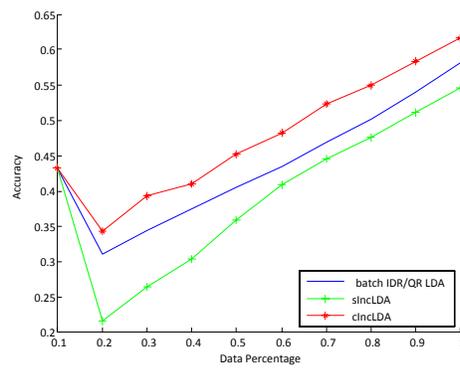
(b) Altkom



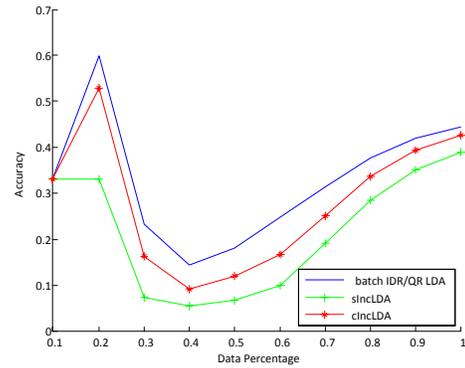
(c) MPEG-7



(d) MPEG



(e) XM2V



(f) Combined

Figure 5.2: The comparison of class separability between the proposed cIncLDA and sIncLDA, with reference to the batch QR/LDA

As seen from the table, both type differences are positive for all six datasets. This indicates that the proposed cInclDA retains more discriminant information than the sInclDAs at every incremental learning stage. The stage difference is seen to be smaller than the final difference for most cases except the At&T dataset. This suggests that the superiority of the cInclDA to sInclDA grows at every incremental learning stage, ensuring a reliable and significant superiority on the learning of the whole data.

For the AT&T dataset with a small number of classes ($k=40$), only subtle differences are exhibited between two InclDAs, whereas for datasets with larger number of classes, such as the MPEG and XM2VTS data, the differences are more evident. However, the superiority is not determined by merely the class number, but also the number of samples and the dimensionality, despite the class number playing an important role.

Dataset	# of samples	# of dimensions	# of classes	stage Acc. diff.	final Acc. diff.
AT&T(ORL)	400	$92 \times 112 = 10304$	40	+0.0018	+0.0005
Altkom	1220	$46 \times 56 = 2576$	80	+0.1037	+0.1172
MPEG-7	1355	$23 \times 28 = 644$	271	+0.1851	+0.2590
XM2VTS	2950	$46 \times 56 = 2576$	295	+0.0825	+0.0707
MPEG	3175	$46 \times 56 = 2576$	635	+0.1439	+0.1272
Combined	5050	$46 \times 56 = 2576$	1010	+0.0636	+0.0382

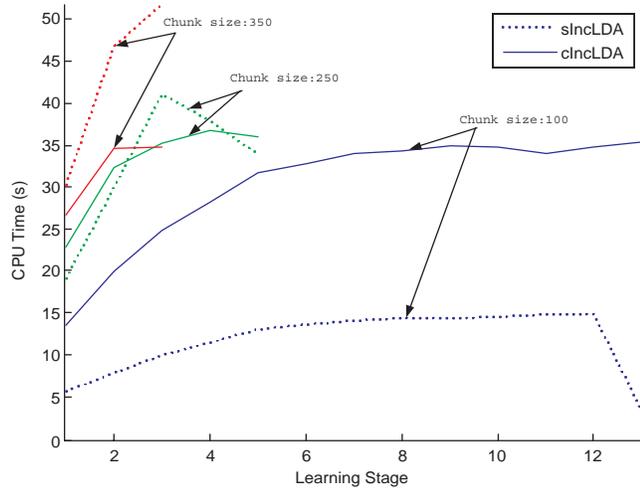
Table 5.1: The stage accuracy differences and final accuracy differences of the incremental learning for six datasets, sInclDA versus cInclDA

5.3.3 Computational Efficiency

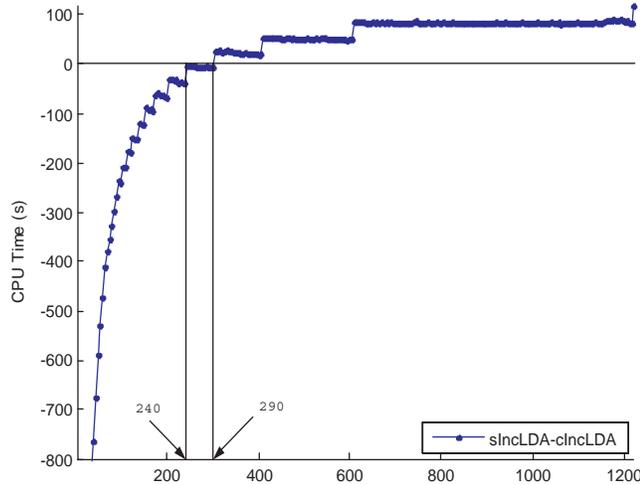
In this experiment, cInclDA is compared with sInclDA on CPU time cost. To observe the relationship between the execution time and chunk size (size of the added subsets), we use different chunk sizes ranging from 5 to 1220 with a step of five in the experiment, and record the differences between cInclDA and sInclDA at every incremental learning stage.

Fig. 5.3 (a) gives a comparison of the time cost at each stage between cInclDA and sInclDA for chunk sizes at 100, 250, and 350, respectively. It is seen that when the chunk size is 100, the stage learning time cost of cInclDA is less than that of sInclDA, and the difference is as large as 20 seconds on average. However at 250, the efficiency of cInclDA is competitive with sInclDA, and its efficiency surpasses the sInclDA entirely, when the chunk size is just 100 bigger.

We further measure the total time cost by cInclLDA to compute a discriminant model for all 1,355 samples under different chunk sizes. To better clarify the comparison between cInclLDA and sInclLDA, we show the sInclLDA time minus the cInclLDA time in Fig. 5.3 (b). As is seen, the efficiency of cInclLDA grows nicely with the increase of the chunk size. When the chunk size reaches 240, cInclLDA matches sInclLDA in general. The superiority of cInclLDA occurs and remarks when the incremental sizes are larger than 290. The difference scales up to over 100 seconds, when the chunk size increases to 1220.



(a)



(b)

Figure 5.3: Comparison on time cost with different incremental sizes, cInclLDA vs. sInclLDA.

To discover how the running efficiency of the proposed cInclDA is maximized by a large chunk size, we conduct incremental learning additionally in two settings: the first uses half, and the second uses all the remaining dataset as one chunk data input to InclDAs. We measure the time costs of each incremental learning for all the above six datasets, respectively. Table 5.2 presents all the time costs from the proposed cInclDA with a comparison to those of sInclDA, where the time saving is calculated in both CPU seconds and the percentage to the corresponding sInclDA cost.

As seen in Table 5.2, the time cost of the sInclDA learning varies little on chunk size choices. In contrast, when the chunk size is doubled, the proposed cInclDA reduces the time cost to half for all 6 datasets regardless of their differences on sample size (i.e., n), dimensionality (i.e., d), and the number of classes (i.e., k). Moreover, the proposed cInclDA saves time against sInclDA consistently in all cases. The minimum saving is on the MPEG-7 dataset, where the saved 84.46 CPU seconds is 54% of the sInclDA time cost. The maximum time saving (in terms of percentage to sInclDA) is on the Altkom dataset, where 473.16 seconds are equivalent to 92.67% of the sInclDA time cost.

Dataset (#sample/#dimension/#class)	Chunk Size	CPU Time(s)		Time saving(s)
		sInclDA	cInclDA	
AT&T(ORL)	180	2279.30	473.40	1805.9 / 79.23%
(400/10304/40)	360	2278.20	242.29	2035.9 / 89.37%
Altkom	540	509.50	70.14	439.36 / 86.23%
(1200/2576/80)	1080	510.58	37.42	473.16 / 92.67%
MPEG-7	610	153.68	69.22	84.46 / 54.96%
(1355/644/271)	1220	152.74	35.76	116.99 / 76.59%
XM2VTS	1327	2893.30	462.09	2431.2 / 84.03%
(2950/2576/295)	2655	2893.20	235.16	2658.1 / 91.87%
MPEG	1428	6194.80	1807.10	4387.7 / 70.83%
(3175/2576/635)	2857	6195.60	913.40	5282.2 / 85.26%
Combined	2272	17419.00	4500.80	12918.00 / 74.16%
(5050/2576/1010)	4545	17417.20	2291.70	15125.00 / 86.84%

Table 5.2: The execution time of incremental learning on six datasets, sInclDA versus cInclDA, with two different chunk sizes.

5.4 Summary

The experimental result indicates that our algorithm achieves an accuracy level that is competitive to batch QR/LDA and is consistently higher than sInclDA. The proposed algorithm also produces the optimal transform matrix that is very similar to those obtained by batch QR/LDA. The efficiency of our algorithm can match and surpass sInclDA as the chunk size increases for multiple instances processing. In view of this, we believe that our proposed cInclDA is much better than sInclDA, and will be applied to more real-world tasks.

Chapter 6

Conclusions and Directions for Future Research

This chapter provides a brief retrospect of the content presented in the thesis. We summarize the work reported in each chapter and identify the contributions. In addition, several directions for future research will be highlighted.

6.1 Conclusions

In this thesis, we first reviewed major incremental learning methods including LDAs and IncLDAs. Based on the review, we proposed a new framework for assessing IncLDAs. Further, we intensively analyzed current popular IncLDAs in terms of a component model, and evaluated those IncLDAs based on the framework. In particular, we analyzed existing sequential IDR/QR IncLDA from an embedding point of view, and indicated its limitations, such as inefficiency of processing chunk data, etc.

To resolve the limitations of sIncLDA, we proposed a new Chunk IDR/QR IncLDA algorithm that is capable of processing a whole chunk of data for every incremental learning cycle. Following this, we explicitly compared the time complexity of cIncLDA and sIncLDA. Our analysis showed that the computational efficiency of cIncLDA is capable of matching and surpassing that of sIncLDA as the chunk increases.

Experimental results showed that performance (i.e., accuracy) is improved in our solution, especially when processing data with large number of classes. In addition,

compared with sInclDA, the proposed cInclDA also produces the transform matrix much similar to the one obtained by batch QR/LDA. Moreover, the experimental study on computational efficiencies agrees well with theoretical time complexity analysis. As evidenced in the thesis, the running efficiency of the proposed cInclDA can be maximized by large chunks.

The main contributions of this paper are summarized as below:

1. A new evaluation criteria framework for InclDAs has been proposed. It allows us to assess InclDAs from a new point of view, not just by comparisons on performance, efficiency and so forth.
2. Two efficient matrix augmentation methods have been proposed. They allow us to easily accommodate information of new classes contained in chunk data through expanding the dimensions of matrices \mathbf{Q} and \mathbf{W} .
3. A new method for incremental updating of \mathbf{W} , more accurate than the existing approach, has been successfully developed in this thesis. This is achieved by relaxing a key assumption on \mathbf{Q} and \mathbf{W} that is vital to the existing algorithm but may be significantly violated as the number of classes increases.
4. Unlike Pang’s InclDAs (Pang et al., 2005b) where only the chunk data in one class can be acquired each time, our method is capable of absorbing information obtained from newly added samples in whatever classes.

6.2 Directions for Future Research

As we noted in the thesis, the reduced between-class scatter matrix \mathbf{B} in the existing sInclDA and the proposed cInclDA are not incrementally updated but completely recomputed. As a matter of fact, whenever new data are presented for the incremental learning procedure, there should exist an incremental learning rule such as $\mathbf{B}' = \mathbf{B} + \mathbf{\Sigma}$. One important future challenge would be to develop such an incremental mechanism for calculating $\mathbf{\Sigma}$ and subsequently updating \mathbf{B} . Once realized, we believe the efficiency of cInclDA could be improved further. As analyzed in section 4.5, the time complexity for recomputing \mathbf{B} is $O(k^3)$. It is straightforward to see that if the reduced between-class scatter matrix can be updated through $\mathbf{B}' = \mathbf{B} + \mathbf{\Sigma}$, which only requires $O(k^2)$, significant reduction of computation cost can be achieved.

Our time complexity analysis also showed that when $s > \eta$, the running efficiency of cIncLDA can match and even exceed that of sIncLDA. This has been mathematically proved. However, it remains as a challenge to find a general formula for computing the exact match point and surpass point.

Finally, in the thesis, only face recognition application has been touched. It is also interesting to investigate other potential applications of our proposed cIncLDA, and design experiments on other real-world datasets to see the generability and stability of our work.

Appendix A

Proof of Proposition 1

Given $\mathbf{A} = (\mathbf{A}_1, \mathbf{A}_2) = \mathbf{Q}(\mathbf{R}_1, \mathbf{R}_2) = \mathbf{Q}\mathbf{R} \in \mathcal{R}^{d \times k}$ with $\mathbf{A}_1 \in \mathcal{R}^{d \times (i-1)}$, and $\mathbf{A}_2 \in \mathcal{R}^{d \times (k-i)}$. According to Daniel et al. (1976), if a zero vector $\mathbf{0}$ is inserted at the i -th column of \mathbf{A} , then we have the augmented matrix $\hat{\mathbf{A}}$ as,

$$\hat{\mathbf{A}} = (\mathbf{A}_1, \mathbf{0}, \mathbf{A}_2) = (\mathbf{Q}, \mathbf{0}) \begin{pmatrix} \mathbf{R}_1 & \mathbf{0} & \mathbf{R}_2 \\ \mathbf{0}^T & 1 & \mathbf{0}^T \end{pmatrix} = \Psi \Xi,$$

where $\Psi = (\mathbf{Q}, \mathbf{0})$ and $\Xi = \begin{pmatrix} \mathbf{R}_1 & \mathbf{0} & \mathbf{R}_2 \\ \mathbf{0}^T & 1 & \mathbf{0}^T \end{pmatrix}$.

For QR-updating, the re-orthogonalization step stated by Daniel et al. (1976) can be avoided, because $\mathbf{0}$ is orthogonal to any vectors. We choose directly the Givens rotations $\mathbf{G} = [\mathbf{G}_{i,i+1}, \dots, \mathbf{G}_{k-1,k}, \mathbf{G}_{k,k+1}]$ so that

$$\mathbf{G}^T \cdot \begin{pmatrix} \mathbf{0} \\ 1 \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ 1 \\ \mathbf{0} \end{pmatrix} \in \mathcal{R}^{k \times 1} \quad (\text{A.1})$$

where 1 positions at the i -th row, and “.” denotes the application of individual Givens rotation $\mathbf{G}_{i,i+1}$ on $\begin{pmatrix} \mathbf{0} \\ 1 \end{pmatrix}$ one by one. As a result, we obtain $\mathbf{G}_{i,i+1} = \dots = \mathbf{G}_{k-1,k} = \mathbf{G}_{k,k+1} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$.

Consequently, we can apply \mathbf{G}^T on Ξ to compute $\hat{\mathbf{R}}$ as,

$$\begin{aligned}
\mathbf{G}^T \cdot \Xi &= \mathbf{G}_{i,i+1}^T \cdots \mathbf{G}_{k,k+1}^T \begin{pmatrix} \mathbf{R}_1 & \mathbf{0} & \mathbf{R}_2 \\ \mathbf{0}^T & 1 & \mathbf{0}^T \end{pmatrix} \\
&= \mathbf{G}_{i,i+1}^T \cdots \mathbf{G}_{k-1,k}^T \begin{pmatrix} \mathbf{R}_1^{(1)} & \mathbf{0} & \mathbf{R}_2^{(1)} \\ \mathbf{0}^T & 1 & \mathbf{0}^T \\ \mathbf{R}_3^{(1)} & \mathbf{0} & -\mathbf{R}_4^{(1)} \end{pmatrix} \\
&\dots \\
&= \mathbf{G}_{i,i+1}^T \begin{pmatrix} \mathbf{R}_1^{(k-i)} & \mathbf{0} & \mathbf{R}_2^{(k-i)} \\ \mathbf{0}^T & 1 & \mathbf{0}^T \\ \mathbf{R}_3^{(k-i)} & \mathbf{0} & -\mathbf{R}_4^{(k-i)} \end{pmatrix} \\
&= \begin{pmatrix} \mathbf{R}_1^{(k+1-i)} & \mathbf{0} & \mathbf{R}_2^{(k+1-i)} \\ \mathbf{0}^T & 1 & \mathbf{0}^T \\ \mathbf{R}_3^{(k+1-i)} & \mathbf{0} & -\mathbf{R}_4^{(k+1-i)} \end{pmatrix} \\
&= \hat{\mathbf{R}}.
\end{aligned} \tag{A.2}$$

Here $\mathbf{R}_1^{(k+1-i)} \in \mathcal{R}^{(i-1) \times (i-1)}$, $\mathbf{R}_2^{(k+1-i)} \in \mathcal{R}^{(i-1) \times (k-i)}$, $\mathbf{R}_3^{(k+1-i)} \in \mathcal{R}^{(k-i) \times (i-1)}$, $\mathbf{R}_4^{(k+1-i)} \in \mathcal{R}^{(k-i) \times (k-i)}$, and $\begin{pmatrix} \mathbf{R}_1^{(k+1-i)} & \mathbf{R}_2^{(k+1-i)} \\ \mathbf{R}_3^{(k+1-i)} & \mathbf{R}_4^{(k+1-i)} \end{pmatrix} = \mathbf{R}$ with its submatrix $\mathbf{R}_3^{(k+1-i)}$ being a zero matrix. Note that (A.2) performs an iterative calculation, because an individual Givens rotation $\mathbf{G}_{i,i+1}$ effects only two rows of Ξ (Golub & Van Loan, 1996).

Similarly, we apply \mathbf{G} on Ψ to compute $\hat{\mathbf{Q}}$ as,

$$\begin{aligned}
\Psi \cdot \mathbf{G} &= (\mathbf{Q}, \mathbf{0}) \mathbf{G}_{k,k+1} \cdots \mathbf{G}_{i,i+1} \\
&= (\mathbf{Q}_1^{(1)}, \mathbf{0}, -\mathbf{Q}_2^{(1)}) \mathbf{G}_{k-1,k} \cdots \mathbf{G}_{i,i+1} \\
&\dots \\
&= (\mathbf{Q}_1^{(k-i)}, \mathbf{0}, -\mathbf{Q}_2^{(k-i)}) \mathbf{G}_{i,i+1} \\
&= (\mathbf{Q}_1^{(k+1-i)}, \mathbf{0}, -\mathbf{Q}_2^{(k+1-i)}) \\
&= \hat{\mathbf{Q}}.
\end{aligned} \tag{A.3}$$

Here $\mathbf{Q}_1^{(k+1-i)} \in \mathcal{R}^{d \times (i-1)}$, $\mathbf{Q}_2^{(k+1-i)} \in \mathcal{R}^{d \times (k-i)}$, and $(\mathbf{Q}_1^{(k+1-i)}, \mathbf{Q}_2^{(k+1-i)}) = \mathbf{Q}$. Note that, only two columns of Ψ are effected by an individual Givens rotation $\mathbf{G}_{i,i+1}$ (Golub & Van Loan, 1996).

Hereby, we have $\hat{\mathbf{A}} = \hat{\mathbf{Q}}\hat{\mathbf{R}}$ as required.

Appendix B

Proof of Proposition 2

Given $\mathbf{W} = \mathbf{Q}^T \mathbf{H}_w \mathbf{H}_w^T \mathbf{Q} \in \mathcal{R}^{k \times k}$, if a zero vector $\mathbf{0}$ is inserted as the i -th column into $\mathbf{Q} \in \mathcal{R}^{d \times k}$, so that $\hat{\mathbf{Q}} = (\mathbf{Q}_1, \mathbf{0}, -\mathbf{Q}_2) \in \mathcal{R}^{(k+1) \times (k+1)}$ where $\mathbf{Q}_1 \in \mathcal{R}^{d \times (i-1)}$ and $\mathbf{Q}_2 \in \mathcal{R}^{d \times (k-i)}$. According to (A.3), we derive the augmented reduced within-class matrix as follows,

$$\begin{aligned}
 \hat{\mathbf{W}} &= \hat{\mathbf{Q}}^T \mathbf{H}_w \mathbf{H}_w^T \hat{\mathbf{Q}} \\
 &= (\mathbf{G}^T \cdot \Psi^T) \mathbf{H}_w \mathbf{H}_w^T (\Psi \cdot \mathbf{G}) \\
 &= \mathbf{G}^T \cdot (\mathbf{Q}, \mathbf{0})^T \mathbf{H}_w \mathbf{H}_w^T (\mathbf{Q}, \mathbf{0}) \cdot \mathbf{G} \\
 &= \mathbf{G}^T \cdot \begin{pmatrix} \mathbf{W} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \cdot \mathbf{G} \\
 &= \mathbf{G}_{i,i+1}^T \cdots \mathbf{G}_{k,k+1}^T \begin{pmatrix} \mathbf{W} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{G}_{k,k+1} \cdots \mathbf{G}_{i,i+1} \\
 &= \mathbf{G}_{i,i+1}^T \cdots \mathbf{G}_{k-1,k}^T \begin{pmatrix} \mathbf{W}_1^{(1,0)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}^T & 0 & 0 \\ -\mathbf{W}_3^{(1,0)} & -\mathbf{W}_4^{(1,0)} & \mathbf{0} \end{pmatrix} \mathbf{G}_{k,k+1} \cdots \mathbf{G}_{i,i+1} \\
 &\dots \\
 &= \begin{pmatrix} \mathbf{W}_1^{(k+1-i,0)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}^T & \mathbf{0}^T & 0 \\ -\mathbf{W}_3^{(k+1-i,0)} & -\mathbf{W}_4^{(k+1-i,0)} & \mathbf{0} \end{pmatrix} \mathbf{G}_{k,k+1} \cdots \mathbf{G}_{i,i+1} \\
 &\dots \\
 &= \begin{pmatrix} \mathbf{W}_1^{(k+1-i,k+1-i)} & \mathbf{0} & -\mathbf{W}_2^{(k+1-i,k+1-i)} \\ \mathbf{0}^T & 0 & \mathbf{0}^T \\ -\mathbf{W}_3^{(k+1-i,k+1-i)} & \mathbf{0} & \mathbf{W}_4^{(k+1-i,k+1-i)} \end{pmatrix},
 \end{aligned}$$

where $\mathbf{W}_1^{(k+1-i, k+1-i)} \in \mathcal{R}^{(i-1) \times (i-1)}$, $\mathbf{W}_2^{(k+1-i, k+1-i)} \in \mathcal{R}^{(i-1) \times (k-i)}$, $\mathbf{W}_3^{(k+1-i, k+1-i)} \in \mathcal{R}^{(k-i) \times (i-1)}$, $\mathbf{W}_4^{(k+1-i, k+1-i)} \in \mathcal{R}^{(k-i) \times (k-i)}$, and $\begin{pmatrix} \mathbf{W}_1^{(k+1-i, k+1-i)} & \mathbf{W}_2^{(k+1-i, k+1-i)} \\ \mathbf{W}_3^{(k+1-i, k+1-i)} & \mathbf{W}_4^{(k+1-i, k+1-i)} \end{pmatrix} = \mathbf{W}$.

Appendix C

Proof of Proposition 3

As $W = Q^T H_w H_w^T Q$ (Ye et al., 2005), its update for incremental learning is calculated as,

$$\begin{aligned} W' &= Q'^T H'_w H'^T_w Q' \\ &= Q'^T \sum_{i=1}^{k'} (H'_{w_i} H'^T_{w_i}) Q' . \end{aligned} \quad (\text{C.1})$$

As we know, $H'_w = [X'_1 - m'_1 e'^T_1, \dots, X'_{k'} - m'_{k'} e'^T_{k'}] \in \mathcal{R}^{d \times k'}$ and $e'_i = \begin{bmatrix} e_i \\ \tilde{e}_i \end{bmatrix} = (1, \dots, 1)^T \in \mathcal{R}^{n'_i}$, H'_{w_i} in (C.1) gives,

$$\begin{aligned} H'_{w_i} &= X'_i - m'_i e'^T_i \\ &= [X_i, \tilde{X}_i] - m'_i e'^T_i \\ &= [X_i - m_i e_i^T, \tilde{X}_i - \tilde{m}_i \tilde{e}_i^T] - m'_i e'^T_i + m_i e_i^T \\ &= [H_{w_i}, \tilde{X}_i - \tilde{m}_i \tilde{e}_i^T + \tilde{m}_i \tilde{e}_i^T - m_i \tilde{e}_i^T] - (m'_i - m_i) e'^T_i \\ &= [H_{w_i}, \tilde{H}_{w_i} + (\tilde{m}_i - m_i) \tilde{e}_i^T] - (m'_i - m_i) e'^T_i \\ &= [H_{w_i}, \tilde{H}_{w_i} + \alpha \tilde{e}_i^T] - \beta e'^T_i, \end{aligned} \quad (\text{C.2})$$

where $\alpha = \tilde{m}_i - m_i$ and $\beta = m'_i - m_i$. Hereby, $H'_{w_i} H'^T_{w_i}$ in (C.1) is computed as

$$\begin{aligned}
\mathbf{H}'_{w_i} \mathbf{H}'_{w_i T} &= ([\mathbf{H}_{w_i}, \widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T] - \beta \mathbf{e}'_i^T) ([\mathbf{H}_{w_i}, \widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T] - \beta \mathbf{e}'_i^T)^T \\
&= ([\mathbf{H}_{w_i}, \widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T] - \beta \mathbf{e}'_i^T) \left(\begin{bmatrix} \mathbf{H}_{w_i}^T \\ (\widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T)^T \end{bmatrix} - \mathbf{e}'_i \beta^T \right) \\
&= [\mathbf{H}_{w_i}, \widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T] \begin{bmatrix} \mathbf{H}_{w_i}^T \\ (\widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T)^T \end{bmatrix} - \beta \mathbf{e}'_i^T \begin{bmatrix} \mathbf{H}_{w_i}^T \\ (\widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T)^T \end{bmatrix} \\
&\quad - [\mathbf{H}_{w_i}, \widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T] \mathbf{e}'_i \beta^T + \beta \mathbf{e}'_i^T \mathbf{e}'_i \beta^T \\
&= \mathbf{H}_{w_i} \mathbf{H}_{w_i}^T + (\widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T) (\widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T)^T - \beta [\mathbf{e}'_i^T, \widetilde{\mathbf{e}}_i^T] \begin{bmatrix} \mathbf{H}_{w_i}^T \\ (\widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T)^T \end{bmatrix} \\
&\quad - [\mathbf{H}_{w_i}, \widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T] \begin{bmatrix} \mathbf{e}_i \\ \widetilde{\mathbf{e}}_i \end{bmatrix} \beta^T + n'_i \beta \beta^T \\
&= \mathbf{H}_{w_i} \mathbf{H}_{w_i}^T + (\widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T) (\widetilde{\mathbf{H}}_{w_i} + \widetilde{\mathbf{e}}_i \alpha^T) - \beta (\mathbf{e}'_i^T \mathbf{H}_{w_i}^T + \widetilde{\mathbf{e}}_i^T (\widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T)^T) \\
&\quad - (\mathbf{H}_{w_i} \mathbf{e}_i + (\widetilde{\mathbf{H}}_{w_i} + \alpha \widetilde{\mathbf{e}}_i^T) \widetilde{\mathbf{e}}_i) \beta^T + n'_i \beta \beta^T \\
&= \mathbf{H}_{w_i} \mathbf{H}_{w_i}^T + \widetilde{\mathbf{H}}_{w_i} \widetilde{\mathbf{H}}_{w_i}^T + \alpha \widetilde{\mathbf{e}}_i \widetilde{\mathbf{H}}_{w_i}^T + \widetilde{n}_i \alpha \alpha^T + \widetilde{\mathbf{H}}_{w_i} \widetilde{\mathbf{e}}_i \alpha^T \\
&\quad - \beta (\mathbf{e}'_i^T \mathbf{H}_{w_i}^T + \widetilde{\mathbf{e}}_i^T \widetilde{\mathbf{H}}_{w_i}^T + \widetilde{\mathbf{e}}_i^T \widetilde{\mathbf{e}}_i \alpha^T) - (\mathbf{H}_{w_i} \mathbf{e}_i + \widetilde{\mathbf{H}}_{w_i} \widetilde{\mathbf{e}}_i + \alpha \widetilde{\mathbf{e}}_i^T \widetilde{\mathbf{e}}_i) \beta^T + n'_i \beta \beta^T \\
&= \mathbf{H}_{w_i} \mathbf{H}_{w_i}^T + \widetilde{\mathbf{H}}_{w_i} \widetilde{\mathbf{H}}_{w_i}^T + \widetilde{n}_i (\alpha - \beta) (\alpha - \beta)^T + n_i \beta \beta^T,
\end{aligned} \tag{C.3}$$

since $\mathbf{H}_{w_i} \mathbf{e}_i = \widetilde{\mathbf{H}}_{w_i} \widetilde{\mathbf{e}}_i = 0$, $\widetilde{\mathbf{e}}_i^T \widetilde{\mathbf{e}}_i = \widetilde{n}_i$ and $\mathbf{e}'_i^T \mathbf{e}'_i = n'_i$. Substituting $\mathbf{H}'_{w_i} \mathbf{H}'_{w_i T}$ in $\mathbf{H}'_w \mathbf{H}'_w T = \sum_{i=1}^{k'} \mathbf{H}'_{w_i} \mathbf{H}'_{w_i T}$, we have

$$\mathbf{H}'_w \mathbf{H}'_w T = \mathbf{H}_w \mathbf{H}_w T + \widetilde{\mathbf{H}}_w \widetilde{\mathbf{H}}_w T + \sum_{i=1}^{k'} (\widetilde{n}_i (\alpha - \beta) (\alpha - \beta)^T + n_i \beta \beta^T). \tag{C.4}$$

Further substituting (C.4) into (C.1), we obtain the update of reduced within-class scatter matrix as,

$$\begin{aligned}
\mathbf{W}' &= \mathbf{Q}'^T \mathbf{H}'_w \mathbf{H}'_w T \mathbf{Q}' \\
&= \mathbf{Q}'^T \mathbf{H}_w \mathbf{H}_w T \mathbf{Q}' + \mathbf{Q}'^T \widetilde{\mathbf{H}}_w \widetilde{\mathbf{H}}_w T \mathbf{Q}' \\
&\quad + \mathbf{Q}'^T (\sum_{i=1}^{k'} (\widetilde{n}_i (\alpha - \beta) (\alpha - \beta)^T + n_i \beta \beta^T)) \mathbf{Q}' \\
&= \bar{\mathbf{W}} + \widetilde{\mathbf{W}} + (\sum_{i=1}^{k'} (\widetilde{n}_i (\alpha' - \beta') (\alpha' - \beta')^T + n_i \beta' \beta'^T)) \\
&= \bar{\mathbf{W}} + \widetilde{\mathbf{W}} + \mathbf{E},
\end{aligned} \tag{C.5}$$

where $\bar{\mathbf{W}} = \mathbf{Q}'^T \mathbf{H}_w \mathbf{H}_w T \mathbf{Q}'$, $\widetilde{\mathbf{W}} = \mathbf{Q}'^T \widetilde{\mathbf{H}}_w \widetilde{\mathbf{H}}_w T \mathbf{Q}'$, $\mathbf{E} = \sum_{i=1}^{k'} (\widetilde{n}_i (\alpha' - \beta') (\alpha' - \beta')^T + n_i \beta' \beta'^T)$, $\alpha' = \mathbf{Q}'^T \alpha$ and $\beta' = \mathbf{Q}'^T \beta$.

References

- Ataa Allah, F., Grosky, W. & Aboutajdine, D. (2007). On-line single-pass clustering based on diffusion maps. In Z. Kedad, N. Lammari, E. Mtais, F. Meziane & Y. Rezgui (Eds.), *Natural language processing and information systems* (Vol. 4592, p. 107-118). Springer Berlin / Heidelberg.
- Baldi, P. & Hatfield, G. (2002). *Dna microarrays and gene expression: from experiments to data analysis and modeling*. Cambridge University Press.
- Belhumeur, P. N., Hespanha, J. P. & Kriegman, D. J. (1997). Eigenfaces vs. fisherfaces: recognition using class specific linear projection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(7), 711-720.
- Bishop, C. (2006). *Pattern recognition and machine learning*. Springer.
- Cai, D., He, X. & Han, J. (2008). Srda: An efficient algorithm for large-scale discriminant analysis. *IEEE Transactions on Knowledge and Data Engineering*, 20, 1-12.
- Cauwenberghs, G. & Poggio, T. (2000). *Incremental and decremental support vector machine learning*.
- Chan, C.-H., Kittler, J. & Messer, K. (2007). Multi-scale local binary pattern histograms for face recognition. In S.-W. Lee & S. Li (Eds.), *Advances in biometrics* (Vol. 4642, p. 809-818). Springer Berlin Heidelberg.
- Chen, L.-F., Liao, H.-Y. M., Ko, M.-T., Lin, J.-C. & Yu, G.-J. (2000). A new lda-based face recognition system which can solve the small sample size problem. *Pattern Recognition*, 33(10), 1713-1726.
- Chu, D. & Thye, G. S. (2010). A new and fast implementation for null space based linear discriminant analysis. *Pattern Recognition*, 43(4), 1373-1379.
- Daniel, J. W., Gragg, W. B., Kaufman, L. & Stewart, G. W. (1976). Reorthogonalization and stable algorithms for updating the gram-schmidt qr factorization. *Mathematics of Computation*, 30(136), 772-795.
- Duda, R., Hart, P. & Stork, D. (2001). *Pattern classification*. Wiley.
- Dudoit, S., Fridlyand, J. & Speed, T. P. (2002). Comparison of discrimination methods for the classification of tumors using gene expression data. *Journal of the American Statistical Association*, 97(457), 77-87.
- Fisher, R. A. (1936). The use of multiple measurements in taxonomic problems. *Annals of Human Genetics*, 7(2), 179-188.
- Frakes, W. & Baeza-Yates, R. (1992). *Information retrieval: data structures &*

- algorithms*. Prentice Hall.
- Frank, A. & Asuncion, A. (2010). *UCI machine learning repository*. Available from <http://archive.ics.uci.edu/ml>
- Friedman, J. H. (1989). Regularized discriminant analysis. *Journal of the American Statistical Association*, 84(405), 165-175.
- Fukunaga, K. (1990). *Introduction to statistical pattern recognition*. Academic Press.
- Giraud-Carrier, C. (2000). A note on the utility of incremental learning. *AI Communications*, 13(4), 215-223.
- Golub, G. H. & Van Loan, C. F. (1996). *Matrix computations* (third ed.). Baltimore, M.D.: Johns Hopkins University Press.
- G.W., S. (1998). Baise decompositions. In *Matrix algorithms* (Vol. I). SIAM.
- Hall, P., Marshall, D. & Martin, R. (2000). Merging and splitting eigenspace models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(9), 1042-1049.
- Hao, Z., Yu, S., Yang, X., Zhao, F., Hu, R. & Liang, Y. (2004). Online ls-svm learning for classification problems based on incremental chunk. In F.-L. Yin, J. Wang & C. Guo (Eds.), *Advances in neural networks c issn 2004* (Vol. 3173, p. 558-564). Springer Berlin / Heidelberg.
- Hong, Z.-Q. & Yang, J.-Y. (1991). Optimal discriminant plane for a small number of samples and design method of classifier on the plane. *Pattern Recognition*, 24(4), 317-324.
- Howland, P., Jeon, M. & Park, H. (2003). Structure preserving dimension reduction for clustered text data based on the generalized singular value decomposition. *SIAM Journal on Matrix Analysis and Applications*, 25(1), 165-179.
- Huang, R., Liu, Q., Lu, H. & Ma, S. (2002). *Solving the small sample size problem of lda* (Vol. 3).
- Jain, A. K. & Dubes, R. C. (1988). *Algorithms for clustering data*. Upper Saddle River, NJ, USA: Prentice-Hall, Inc.
- Jiang, J., Song, C., Wu, C., Maurizio, M. & Liang, Y. (2006). Support vector machine regression algorithm based on chunking incremental learning. In V. Alexandrov, G. van Albada, P. Sloot & J. Dongarra (Eds.), *Computational science c iccs 2006* (Vol. 3991, p. 547-554). Springer Berlin / Heidelberg.
- Jin, Z., Yang, J., Hu, Z. & Lou, Z. (2001). Face recognition based on the uncorrelated discriminant transformation. *Pattern Recognition*, 34(7), 1405-1416.

- Karasuyama, M. & Takeuchi, I. (2010). Multiple incremental decremental learning of support vector machines. *Trans. Neur. Netw.*, *21*(7), 1048-1059.
- Kidera, T., Ozawa, S. & Abe, S. (2006). *An incremental learning algorithm of ensemble classifier systems*.
- Kim, M.-S., Kim, D. & Lee, S.-Y. (2003). Face recognition using the embedded hmm with second-order block-specific observations. *Pattern Recognition*, *36*(11), 2723-2735.
- Kim, T.-K., Kim, H., Hwang, W. & Kittler, J. (2005). Component-based lda face description for image retrieval and mpeg-7 standardisation. *Image Vision Comput.*, *23*(7), 631-642.
- Kim, T.-K., Stenger, B., Kittler, J. & Cipolla, R. (2011). Incremental linear discriminant analysis using sufficient spanning sets and its applications. *International Journal of Computer Vision*, *91*(2), 216-232.
- Kowalski, G. (1997). *Information retrieval systems: theory and implementation*. Kluwer Academic Publishers.
- Liu, C. & Wechsler, H. (2002). Gabor feature based classification using the enhanced fisher linear discriminant model for face recognition. *IEEE Transactions on Image Processing*, *11*(4), 467-476.
- Liu, L., Jiang, Y. & Zhou, Z. (2009). *Least square incremental linear discriminant analysis* (Vol. 0).
- Loos, H. G. (1989, 18-22 Jun 1989). *Parity madeline: a neural net with complete boolean repertoire capable of one-pass learning*.
- Lu, J., Plataniotis, K. N. & Venetsanopoulos, A. N. (2003a). Face recognition using kernel direct discriminant analysis algorithms. *IEEE Transactions on Neural Networks*, *14*(1), 117-126.
- Lu, J., Plataniotis, K. N. & Venetsanopoulos, A. N. (2003b). Face recognition using lda-based algorithms. *IEEE Transactions on Neural Networks*, *14*(1), 195-200.
- Lu, J., Plataniotis, K. N. & Venetsanopoulos, A. N. (2005). Regularization studies of linear discriminant analysis in small sample size scenarios with application to face recognition. *Pattern Recognition Letters*, *26*(2), 181-191.
- McLachlan, G. (2004). *Discriminant analysis and statistical pattern recognition*. Wiley-Interscience.
- Ozawa, S., Pang, S. & Kasabov, N. (2006). *An incremental principal component analysis for chunk data*.

- Paige, C. C. & Saunders, M. A. (1981). Towards a generalized singular value decomposition. *SIAM Journal on Numerical Analysis*, 18(3), 398-405.
- Pang, S., Ban, T., Kadobayashi, Y. & Kasabov, N. (2010, 18-23 July 2010). *Incremental and decremental lda learning with applications*.
- Pang, S., Kim, D. & Bang, S. Y. (2003). Membership authentication in the dynamic group by face classification using svm ensemble. *Pattern Recognition Letters*, 24(1-3), 215-225.
- Pang, S., Ozawa, S. & Kasabov, N. (2004). One-pass incremental membership authentication by face classification. In D. Zhang & A. Jain (Eds.), *Biometric authentication* (Vol. 3072, p. 1-44). Springer Berlin / Heidelberg.
- Pang, S., Ozawa, S. & Kasabov, N. (2005a). Chunk incremental lda computing on data streams. In J. Wang, X.-F. Liao & Z. Yi (Eds.), *Advances in neural networks* (Vol. 3497, p. 829-831). Springer Berlin / Heidelberg.
- Pang, S., Ozawa, S. & Kasabov, N. (2005b). Incremental linear discriminant analysis for classification of data streams. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 35(5), 905-914.
- Park, H., Jeon, M. & Rosen, J. B. (2003). Lower dimensional representation of text data based on centroids and least squares. *BIT Numerical Mathematics*, 43(2), 427-448.
- Polikar, R., Upda, L., Upda, S. S. & Honavar, V. (2001). Learn++: an incremental learning algorithm for supervised neural networks. *IEEE Transactions on Systems, Man, and Cybernetics, Part C: Applications and Reviews*, 31(4), 497-508.
- Raudys, S. J. & Jain, A. K. (1990). *Small sample size effects in statistical pattern recognition: recommendations for practitioners and open problems* (Vol. i).
- Song, F., Liu, H., Zhang, D. & Yang, J. (2008). A highly scalable incremental facial feature extraction method. *Neurocomputing*, 71(10-12), 1883-1888.
- Swets, D. L. & Weng, J. J. (1996). Using discriminant eigenfeatures for image retrieval. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 18(8), 831-836.
- Tian, Q., Barbero, M., Gu, Z. & Lee, S. (1986). *Image classification by the foley-sammon transform*.
- Wang, J.-G., Sung, E. & Yau, W.-Y. (2010). Incremental two-dimensional linear discriminant analysis with applications to face recognition. *Journal of Network*

- and Computer Applications*, 33(3), 314-322.
- Wang, X., Han, T. X. & Yan, S. (2009). *An hog-lbp human detector with partial occlusion handling*.
- Ye, J. (2004). Using uncorrelated discriminant analysis for tissue classification with gene expression data. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 1, 181-190.
- Ye, J. (2005). Characterization of a family of algorithms for generalized discriminant analysis on undersampled problems. *J. Mach. Learn. Res.*, 6, 483-502.
- Ye, J. (2007). *Least squares linear discriminant analysis*. ACM.
- Ye, J., Janardan, R., Park, C. H. & Park, H. (2004). An optimization criterion for generalized discriminant analysis on undersampled problems. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(8), 982-994.
- Ye, J. & Li, Q. (2004). Lda/qr: an efficient and effective dimension reduction algorithm and its theoretical foundation. *Pattern Recognition*, 37(4), 851-854.
- Ye, J., Li, Q., Xiong, H., Park, H., Janardan, R. & Kumar, V. (2005). Idr/qr: An incremental dimension reduction algorithm via qr decomposition. *IEEE Transactions on Knowledge and Data Engineering*, 17(9), 1208-1222.
- Yen, O. & Meesad, P. (1999). *Pattern classification by an incremental learning fuzzy neural network* (Vol. 5).
- Yu, H. & Yang, H. (2001). *A direct lda algorithm for high-dimensional data - with application to face recognition*.
- Zha, H. & Simon, H. (1997). *On updating problems in latent semantic indexing* (Tech. Rep.). ((United States) [Lawrence Berkeley National Lab., CA] [Pennsylvania State Univ., University Park, PA . Dept. of Computer Science and Engineering] Other Information: PBD: Nov 1997)
- Zhao, H. & Yuen, P. C. (2008). Incremental linear discriminant analysis for face recognition. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 38(1), 210-221.
- Zhao, W., Chellappa, R. & Phillips, P. (1999). *Subspace linear discriminant analysis for face recognition*. Computer Vision Laboratory, Center for Automation Research, University of Maryland. ((College Park, Md.))