

**GLOBAL-LOCAL PARTIAL LEAST SQUARES
DISCRIMINANT ANALYSIS AND ITS
EXTENSION IN REPRODUCING KERNEL
HILBERT SPACE**

AMINU MUHAMMAD

UNIVERSITI SAINS MALAYSIA

2021

**GLOBAL-LOCAL PARTIAL LEAST SQUARES
DISCRIMINANT ANALYSIS AND ITS
EXTENSION IN REPRODUCING KERNEL
HILBERT SPACE**

by

AMINU MUHAMMAD

**Thesis submitted in fulfilment of the requirements
for the degree of
Doctor of Philosophy**

April 2021

ACKNOWLEDGEMENT

I would like to express my deepest gratitude to you Dad, for the love, support and encouragement you gave me my entire life. I would also like to thank my supervisor, Assoc. Professor Dr. Noor Atinah Ahmad for her help and thoughtful advise. I am especially grateful to my co-supervisor Dr. Norhashidah Awang for her continued support and guidance during my graduate study.

TABLE OF CONTENTS

| | |
|---|--------------|
| ACKNOWLEDGEMENT | ii |
| TABLE OF CONTENTS | iii |
| LIST OF TABLES | vii |
| LIST OF FIGURES | ix |
| LIST OF ABBREVIATIONS | xii |
| LIST OF SYMBOLS | xiv |
| ABSTRAK | xvi |
| ABSTRACT | xviii |
| | |
| CHAPTER 1 INTRODUCTION | |
| 1.1 Overview | 1 |
| 1.2 Problem Statement | 2 |
| 1.3 Motivation of Study | 4 |
| 1.4 Research Objectives | 5 |
| 1.5 Research Contribution | 6 |
| 1.6 Scope of Study | 7 |
| 1.7 Thesis Organization | 7 |
| | |
| CHAPTER 2 SUBSPACE LEARNING | |
| 2.1 Classical Dimensionality Reduction Techniques | 9 |
| 2.1.1 Principal Component Analysis | 11 |
| 2.1.2 Linear Discriminant Analysis | 12 |
| 2.1.3 Partial Least Squares | 14 |

| | | |
|-------|---|----|
| 2.2 | Manifold Learning Techniques | 19 |
| 2.2.1 | Laplacian Eigenmap | 20 |
| 2.2.2 | Locally Linear Embedding | 22 |
| 2.2.3 | Locality Preserving Projections | 24 |
| 2.2.4 | Neighborhood Preserving Embedding | 26 |
| 2.2.5 | Summary | 28 |

CHAPTER 3 LOCALITY PRESERVING PLS-DA WITH APPLICATION IN FACE RECOGNITION

| | | |
|----------|--|----|
| 3.1 | Local Structure Modeling | 33 |
| 3.2 | Within Class Structure | 34 |
| 3.3 | Locality Preserving PLS-DA (LPPLS-DA) | 35 |
| 3.3.1 | The proposed method | 35 |
| 3.3.2 | Solution to the proposed LPPLS-DA method | 37 |
| 3.3.3 | Computational complexity | 39 |
| 3.4 | Experimental Results | 40 |
| 3.4.1 | Face representation using LPPLS-DA | 42 |
| 3.4.2 | Face recognition using LPPLS-DA | 43 |
| 3.4.2(a) | Small face databases | 43 |
| 3.4.2(b) | Moderate face databases | 47 |
| 3.4.2(c) | Large face database | 53 |
| 3.4.3 | Summary | 54 |

CHAPTER 4 DISCRIMINANT SUBSPACE LEARNING IN COMPLEX CHEMICAL DATA CLASSIFICATION AND DISCRIMINATION

| | | |
|-----|----------------------------|----|
| 4.1 | Experimental results | 58 |
|-----|----------------------------|----|

| | | |
|-------|--------------------------|----|
| 4.1.1 | Data visualization | 60 |
| 4.1.2 | Classification | 63 |
| 4.2 | Summary | 69 |

CHAPTER 5 KERNEL LOCALITY PRESERVING PLS-DA

| | | |
|-------|---|----|
| 5.1 | Derivation of LPPLS-DA in Reproducing Kernel Hilbert Space..... | 71 |
| 5.1.1 | Solution and computational analysis of KLPPLS-DA | 76 |
| 5.2 | Experiments..... | 78 |
| 5.2.1 | Datasets..... | 79 |
| 5.2.2 | Results | 81 |
| 5.2.3 | Summary | 86 |

CHAPTER 6 NEIGHBORHOOD PRESERVING PLS-DA FOR FACE RECOGNITION

| | | |
|----------|---|-----|
| 6.1 | Local geometric structure modeling | 89 |
| 6.2 | Global-Local PLS-DA..... | 92 |
| 6.2.1 | Neighborhood Preserving PLS-DA (NPPLS-DA)..... | 93 |
| 6.2.2 | Uncorrelated Neighborhood Preserving PLS-DA (UNPPLS-DA).. | 95 |
| 6.3 | Experimental Results | 98 |
| 6.3.1 | Face representation..... | 99 |
| 6.3.2 | Face recognition..... | 99 |
| 6.3.2(a) | Experiments on the Yale face database..... | 101 |
| 6.3.2(b) | Experiments on the ORL face database | 102 |
| 6.3.2(c) | Experiments on the Extended Yale B face database | 104 |
| 6.3.2(d) | Experiments on the CMU-PIE face database | 106 |
| 6.3.2(e) | Experiments on the AR face database..... | 108 |

| | | |
|----------|--|-----|
| 6.3.2(f) | Experiments on the Essex face database | 110 |
| 6.3.3 | Parameter selection for UNPPLS-DA | 112 |
| 6.3.4 | Summary | 113 |

CHAPTER 7 KERNEL NEIGHBORHOOD PRESERVING PLS-DA

| | | |
|-------|--|-----|
| 7.1 | Derivation of NPPLS-DA and UNPPLS-DA in Reproducing Kernel Hilbert Space | 117 |
| 7.1.1 | Kernel NPPLS-DA | 117 |
| 7.1.2 | Kernel UNPPLS-DA | 120 |
| 7.1.3 | Computational complexity | 127 |
| 7.2 | Experiments | 128 |
| 7.2.1 | Datasets | 128 |
| 7.2.2 | Compared algorithms | 130 |
| 7.2.3 | Results | 131 |
| 7.2.4 | Summary | 135 |

CHAPTER 8 CONCLUSION

| | |
|-------------------------|------------|
| REFERENCES | 143 |
|-------------------------|------------|

LIST OF PUBLICATIONS

LIST OF TABLES

| | | Page |
|-----------|---|------|
| Table 3.1 | Details of the different databases used in our experiments. | 42 |
| Table 3.2 | Best average recognition rate (in Percent), standard deviation and reduced dimensionality (in brackets) on Yale database over ten random splits..... | 43 |
| Table 3.3 | Best average recognition rate (in Percent), standard deviation and reduced dimensionality (in brackets) on ORL database over ten random splits..... | 46 |
| Table 3.4 | Best average recognition rate (in Percent), standard deviation and reduced dimensionality (in brackets) on Extended Yale database over ten random splits. | 48 |
| Table 3.5 | Best average recognition rate (in Percent), standard deviation and reduced dimensionality (in brackets) on AR database over ten random splits..... | 50 |
| Table 3.6 | Best average recognition rate (in Percent), standard deviation and reduced dimensionality (in brackets) on CMU PIE database over ten random splits. | 52 |
| Table 3.7 | Details of the different groups in the Essex face database. | 54 |
| Table 3.8 | Best average recognition rate (in Percent), standard deviation and reduced dimensionality (in brackets) on Essex face database over ten random splits. | 54 |
| Table 4.1 | Summary of datasets used in the experiments..... | 60 |
| Table 4.2 | Best average accuracy (in Percent), standard deviation and corresponding dimensionality (in brackets) obtained on the different datasets over ten random splits. For each dataset, 1/2 of the data samples are used for training and the remaining 1/2 are used for testing. | 67 |
| Table 4.3 | Best average accuracy (in Percent), standard deviation and corresponding dimensionality (in brackets) obtained on the different datasets over ten random splits. For each dataset, 2/3 of the data samples are used for training and the remaining 1/3 are used for testing. | 68 |

| | | |
|-----------|---|-----|
| Table 5.1 | Summary of datasets used in the experiments..... | 81 |
| Table 5.2 | Classification accuracy on the Brain dataset | 87 |
| Table 5.3 | Classification accuracy on the Colon dataset | 87 |
| Table 5.4 | Classification accuracy on the Leukemia dataset | 87 |
| Table 5.5 | Classification accuracy on the Lymphoma dataset..... | 88 |
| Table 5.6 | Classification accuracy on the Prostate dataset | 88 |
| Table 5.7 | Classification accuracy on the SRBCT dataset | 88 |
| Table 6.1 | Best average recognition accuracy (in Percent) on Yale database over ten random splits. | 102 |
| Table 6.2 | Best average recognition accuracy (in Percent) on ORL database over ten random splits. | 104 |
| Table 6.3 | Best average recognition accuracy (in Percent) on Extended Yale B database over ten random splits..... | 106 |
| Table 6.4 | Best average recognition accuracy (in Percent) on CMU PIE database over ten random splits. | 108 |
| Table 6.5 | Best average recognition accuracy (in Percent) on AR face database over ten random splits. | 110 |
| Table 6.6 | Best average recognition accuracy (in Percent) on Essex database over ten random splits. | 112 |
| Table 7.1 | Summary of datasets used in the experiments..... | 130 |
| Table 7.2 | Classification accuracy on the Coffee dataset..... | 136 |
| Table 7.3 | Classification accuracy on the Fruit dataset..... | 136 |
| Table 7.4 | Classification accuracy on the Meat dataset | 136 |
| Table 7.5 | Classification accuracy on the Oil dataset..... | 137 |

LIST OF FIGURES

| | | Page |
|-------------|---|-------------|
| Figure 1.1 | Cell nuclei graphs | 4 |
| Figure 3.1 | Experimental design | 41 |
| Figure 3.2 | First six Eigenfaces, Fisherfaces, LPPLS-DA and PLS-DA calculated on the Yale database. | 42 |
| Figure 3.3 | Sample face images from the Yale database. | 43 |
| Figure 3.4 | Recognition rate vs reduced dimensionality on Yale database (3, 5 and 7 Train). | 44 |
| Figure 3.5 | Sample face images from the ORL database. | 45 |
| Figure 3.6 | Recognition rate vs reduced dimensionality on ORL database (3, 5 and 7 Train). | 46 |
| Figure 3.7 | Sample face images from the Extended Yale database. | 47 |
| Figure 3.8 | Recognition rate vs reduced dimensionality on Extended Yale database (10, 30 and 50 Train). | 48 |
| Figure 3.9 | Sample face images from the AR database..... | 49 |
| Figure 3.10 | Recognition rate vs reduced dimensionality on AR database (3, 5 and 7 Train). | 50 |
| Figure 3.11 | Sample face images from the CMU PIE database. | 51 |
| Figure 3.12 | Recognition rate vs reduced dimensionality on CMU PIE database (3, 5 and 7 Train). | 52 |
| Figure 3.13 | Recognition rate vs reduced dimensionality on Essex face database (3, 5 and 7 Train). | 55 |
| Figure 4.1 | Visualizations of the Coffee and Pacific cod datasets in two-dimensional subspace: (a) PLS-DA on the Coffee data, (b) LPPLS-DA on the Coffee data, (c) PLS-DA on the Pacific Cod data, (d) LPPLS-DA on the Pacific Cod data..... | 61 |

| | | |
|------------|---|-----|
| Figure 4.2 | Visualizations of the Wood and Ink datasets in three-dimensional subspace: (a) PLS-DA on the Wood data, (b) LPPLS-DA on the Wood data, (c) PLS-DA on the Ink data, (d) LPPLS-DA on the Ink data. | 63 |
| Figure 4.3 | Classification accuracies in form of confusion matrices. (a) PLS-DA result on Coffee data, (b) LPPLS-DA result on Coffee data, (c) PLS-DA result on Pacific Cod data, (d) LPPLS-DA result on Pacific Cod data. | 66 |
| Figure 4.4 | Classification accuracies in form of confusion matrices. (a) PLS-DA result on Ink data, (b) LPPLS-DA result on Ink data, (c) PLS-DA result on Wood data, (d) LPPLS-DA result on Wood data. | 67 |
| Figure 4.5 | Average classification accuracy rates by a two-nearest neighbor classifier as a function of the reduced dimension. Here, half of the datasets are used as training sets and the remaining half as the test sets. | 68 |
| Figure 4.6 | Average classification accuracy rates by a two-nearest neighbor classifier as a function of the reduced dimension. Here, two-third of the datasets are used as training sets and the remaining one-third as the test sets. | 69 |
| Figure 5.1 | Average classification accuracies by a 2-NN classifier as a function of the reduced dimension. Here, two-third of the datasets are used as training sets and the remaining one-third as the test sets. | 82 |
| Figure 5.2 | Average classification accuracies by an SVM classifier as a function of the reduced dimension. Here, two-third of the datasets are used as training sets and the remaining one-third as the test sets. | 83 |
| Figure 6.1 | First six basis vectors (eigenvectors) calculated using the different methods on the Yale database. | 100 |
| Figure 6.2 | Recognition accuracy vs reduced dimensionality on Yale database (2, 4, 6 and 8 Train). | 103 |
| Figure 6.3 | Recognition accuracy vs reduced dimensionality on ORL database (2, 4, 6 and 8 Train). | 105 |
| Figure 6.4 | Recognition accuracy vs reduced dimensionality on Extended Yale B database (5, 10, 20 and 30 Train). | 107 |

| | | |
|------------|--|-----|
| Figure 6.5 | Recognition accuracy vs reduced dimensionality on CMU-PIE database (3, 5, 7 and 10 Train). | 109 |
| Figure 6.6 | Recognition accuracy vs reduced dimensionality on AR face database (3, 5, 7 and 10 Train). | 111 |
| Figure 6.7 | Recognition accuracy vs reduced dimensionality on Essex face database (3, 5, 7 and 10 Train). | 113 |
| Figure 6.8 | Recognition rates of UNPPLS-DA with respect to δ on the different databases. | 114 |
| Figure 7.1 | Average classification accuracies by a 2-NN classifier as a function of the reduced dimension. Here, two-third of the datasets are used as training sets and the remaining one-third as the test sets. | 133 |
| Figure 7.2 | Average classification accuracies by an SVM classifier as a function of the reduced dimension. Here, two-third of the datasets are used as training sets and the remaining one-third as the test sets. | 134 |

LIST OF ABBREVIATIONS

| | |
|------------|---|
| GLSP | Global Local Structure Preserving |
| KLPPLS-DA | Kernel Locality Preserving Partial Least Squares Discriminant Analysis |
| KNPPLS-DA | Kernel Neighborhood Preserving Partial Least Squares Discriminant Analysis |
| KNN | K-Nearest Neighbor |
| KUNPPLS-DA | Kernel Uncorrelated Neighborhood Preserving Partial Least Squares Discriminant Analysis |
| LE | Laplacian Eigenmap |
| LDA | Linear Discriminant Analysis |
| LLE | Locally Linear Embedding |
| LPP | Locality Preserving Projections |
| LPPLS-DA | Locality Preserving Partial Least Squares Discriminant Analysis |
| NIPALS | Nonlinear Iterative Partial Least Squares |
| NPE | Neighborhood Preserving Embedding |
| NPPLS-DA | Neighborhood Preserving Partial Least Squares Discriminant Analysis |
| PCA | Principal Component Analysis |
| PLS | Partial Least Squares |
| PLS-DA | Partial Least Squares Discriminant Analysis |

| | |
|-----------|---|
| RBF | Radial Basis Function |
| RKHS | Reproducing Kernel Hilbert Space |
| SVM | Support Vector Machines |
| UNPPLS-DA | Uncorrelated Neighborhood Preserving Partial Least Squares Discriminant Analysis |

LIST OF SYMBOLS

| | |
|--------------------------|--|
| n | The number of data points |
| m | The number of features |
| d | The number of reduced features |
| C | The number of classes |
| x_i | The i -th data point |
| X | The data matrix |
| \bar{X} | The centred data matrix |
| S_b | The between class scatter matrix |
| S_w | The within class scatter matrix |
| S_t | The total scatter matrix |
| S | The affinity matrix |
| L | The graph Laplacian matrix |
| W | The transformation matrix |
| \mathcal{H} | The feature space |
| K | The kernel matrix |
| \bar{K} | The centred kernel matrix |
| ϕ, ψ | The nonlinear mapping functions |
| Φ, Ψ | The data matrices in feature space |
| $\bar{\Phi}, \bar{\Psi}$ | The centred data matrices in feature space |

| | |
|----------------------|---|
| S_b^ϕ, S_b^ψ | The between class scatter matrices in feature space |
| S_w^ϕ, S_w^ψ | The within class scatter matrices in feature space |
| S_t^ϕ, S_t^ψ | The total class scatter matrices in feature space |

**ANALISIS DISKRIMINAN KUASA DUA TERKECIL SEPARA
GLOBAL-SETEMPAT DAN PELANJUTANNYA DALAM RUANG HILBERT
KERNEL PENGHASILAN SEMULA**

ABSTRAK

Pembelajaran subruang adalah satu pendekatan penting untuk mempelajari perwakilan dimensi rendah bagi suatu ruang dimensi tinggi. Apabila sampel data diwakili sebagai titik dalam ruang dimensi tinggi, pembelajaran dengan kedimensian tinggi menjadi mencabar kerana keberkesanan dan kecekapan algoritma pembelajaran turun dengan ketara apabila dimensi meningkat. Oleh itu, teknik pembelajaran subruang digunakan untuk mengurangkan kedimensian data sebelum menggunakan algoritma pembelajaran yang lain. Baru-baru ini, minat terhadap teknik pembelajaran subruang yang berdasarkan kerangka pemeliharaan struktur global dan tempatan (GLSP). Telah meningkat idea utama pendekatan GLSP adalah mencari transformasi data berdimensi tinggi kepada subruang berdimensi yang lebih rendah dengan maklumat struktur data global dan tempatan terpelihara dalam subruang berdimensi rendah. Tesis ini mempertimbangkan kes yang mana data daripada sampel dalam manifold dasar yang terbenam dalam ruang sekitar dimensi tinggi. Dua algoritma pembelajaran subruang baharu dipanggil 'locality preserving partial least squares discriminant analysis' (LPPLS-DA) dan 'neighborhood preserving partial least squares discriminant analysis' (NPPLS-DA) yang berdasarkan kerangka GLSP dicadangkan untuk pembelajaran subruang diskriminan. Tidak seperti analisis diskriminan kuasa dua terkecil separa konvensional (PLS-DA) yang bertujuan hanya untuk memelihara ruang data struktur Euclidan global ruang data, LPPLS-DA yang dicadangkan dan algoritma NPPLS-DA mencari pembenaman yang memelihara kedua-dua struktur global dan struktur manifold tempatan.

Hasilnya, kedua-dua LPPLS-DA dan NPPLS-DA mampu mengekstrak lebih banyak maklumat diskriminasi dalam data asal berbanding dengan PLS-DA dan sangat sesuai untuk pengurangan dimensi dan visualisasi kumpulan data yang kompleks. Selanjutnya, pelanjutan kernel LPPLS-DA dan NPPLS-DA dalam Ruang Hilbert (RKHS) kernel penghasilan semula dicadangkan untuk menangani situasi di mana wujud hubungan tak linear yang kuat di antara set-set data yang dicerap. Peningkatan prestasi algoritma yang dicadangkan berbanding dengan PLS-DA konvensional ditunjukkan melalui beberapa eksperimen. Ia telah menunjukkan bahawa LPPLS-DA dan NPPLS-DA sangat berkesan untuk analisis wajah (pengecaman dan perwakilan). Pelanjutan kernel kaedah-kaedah ini digunakan untuk tumor dan analisis data kimia dan ditunjukkan bahawa model dengan pelanjutan kernel mengatasi model linear apabila wujud hubungan tak linear kuat antara set data yang dicerap.

GLOBAL-LOCAL PARTIAL LEAST SQUARES DISCRIMINANT ANALYSIS AND ITS EXTENSION IN REPRODUCING KERNEL HILBERT SPACE

ABSTRACT

Subspace learning is an essential approach for learning a low dimensional representation of a high dimensional space. When data samples are represented as points in a high dimensional space, learning with the high dimensionality becomes challenging as the effectiveness and efficiency of the learning algorithms drops significantly as the dimensionality increases. Thus, subspace learning techniques are employed to reduce the dimensionality of the data prior to employing other learning algorithms. Recently, there has been a lot of interest in subspace learning techniques that are based on the global and local structure preserving (GLSP) framework. The main idea of the GLSP approach is to find a transformation of the high dimensional data into a lower dimensional subspace, where both the global and local structure information of the data are preserved in the lower dimensional subspace. This thesis consider the case where data is sampled from an underlying manifold embedded in a high dimensional ambient space. Two novel subspace learning algorithms called locality preserving partial least squares discriminant analysis (LPPLS-DA) and neighborhood preserving partial least squares discriminant analysis (NPPLS-DA) which are based on the GLSP framework are proposed for discriminant subspace learning. Unlike the conventional partial least squares discriminant analysis (PLS-DA) which aims at preserving only the global Euclidean structure of the data space, the proposed LPPLS-DA and NPPLS-DA algorithms find an embedding that preserves both the global and local manifold structure. As a result, both LPPLS-DA and NPPLS-DA can extract more discriminant information in the original data than PLS-DA and are well-suited for dimensionality reduction

and visualization of complex datasets. Furthermore, kernel extensions of LPPLS-DA and NPPLS-DA in reproducing kernel Hilbert space (RKHS) are proposed to handle situations where a strong nonlinear relation exist between the sets of observed data. Performance improvement of the proposed algorithms over the conventional PLS-DA is demonstrated through several experiments. It is shown that LPPPLS-DA and NPPLS-DA are very effective for face analysis (recognition and representation). Their kernel extensions are applied to tumor classification and chemical data analysis respectively and it was shown that the kernel extensions outperform their linear counterparts when a strong nonlinear relationship exist between the set of observed data.

CHAPTER 1

INTRODUCTION

1.1 Overview

In many research fields such as machine learning, computer vision, bioinformatics and pattern recognition, data points are represented as points in a high dimensional space. Researchers in such areas encounter difficulties working with such high dimensional data sets. The effectiveness and efficiency of many learning algorithms, such as clustering and classification algorithms, drop rapidly as the dimensionality increases (Guo and Dyer, 2005; Souza et al., 2016). A lot of techniques have been proposed in the past to reduce the dimensionality of the data by either selecting the most representative features from the original ones (feature selection) or by creating new features as linear combinations of the original features (feature extraction). These techniques include principal component analysis (PCA) (Yi et al., 2017; Zhao et al., 2019b), partial least squares (PLS) (Boulesteix and Strimmer, 2007; Rosipal and Krämer, 2005) and linear discriminant analysis (LDA) (Belhumeur et al., 1997). PCA is an unsupervised dimension reduction technique which captures most of the variance of a data, while LDA is a supervised dimension reduction technique which aim at discriminating the different classes in a data.

PLS is a statistical method that models the linear relationship between sets of observed variables X and Y by means of latent variables (components). The method was first developed by Herman Wold (Wold, 1966) and since then it gained wide acceptance in fields such as chemometrics, bioinformatics, social sciences, medicine etc.

The ability of PLS in handling high dimensionality and collinearity problems in spectral data makes it a powerful and standard tool for the analysis of chemical data in chemometrics (Aliakbarzadeh et al., 2016; Bai et al., 2017; Borràs et al., 2014; Høbro et al., 2010; Kemsley, 1996). Although PLS was not designed for discrimination and classification tasks, it has been successfully applied to these problems and its performance was outstanding (Barker and Rayens, 2003; Huang et al., 2005). PLS for discrimination or better known as partial least squares discriminant analysis (PLS-DA) was shown to have a statistical relationship with LDA and it was further suggested that PLS should be used instead of PCA when discrimination is the goal and dimension reduction is needed (Barker and Rayens, 2003). PLS-DA combines feature extraction and discriminant analysis into one algorithm and is well applicable for high dimensional data sets. Theoretically, PLS-DA finds a transformation of the high dimensional data into a lower dimensional subspace in which data samples of different classes are mapped far apart. The transformation is readily computed using the nonlinear iterative partial least squares (NIPALS) algorithm (Wold, 1966).

1.2 Problem Statement

Although PLS-DA provide a principle way of dealing with high dimensional data, the method does not automatically lead to extraction of relevant features. Many studies (Brereton and Lloyd, 2014; Goodhue et al., 2012; Gromski et al., 2015; Mendez et al., 2020) have pointed out that when classification is the goal and dimension reduction is needed, PLS-DA should not be preferred over other traditional methods as it has no significant advantages over them. Some recent studies (Brereton and Lloyd, 2014; Lee et al., 2018b; Pomerantsev and Rodionova, 2018) also indicate the need to refine the

PLS-DA modeling practice strategies, especially in complex data sets such as multi-class, colossal and imbalanced data sets.

Another major drawback of the PLS-DA method is its lack of ability to preserve the local structure of data. PLS-DA sees only the global Euclidean structure of data. It fails to preserve the local structure of data point if the data points lie on a nonlinear manifold hidden in the high dimensional Euclidean space. Fortunately, several techniques that can effectively preserve the local structure of data point have been proposed (He et al., 2005a; He and Niyogi, 2004; Roweis and Saul, 2000; Shikkenawis and Mitra, 2016). However, when treating multi-clustered data as in the case of appearance-based face recognition and cancer classification, there is a need to treat simultaneously, both the global clustering structure as well as the local clustering structure (Cai, 2017; Liu et al., 2013).

In digital pathology image analysis, spatial arrangement of nuclei in histopathological images has been shown to be able to predict patient outcomes (Lu et al., 2021; Nguyen et al., 2014; Zhou et al., 2019). Cell graphs have been proposed to model the relationship between different cell nuclei and the tissue micro-environment using graph features. The graph can be constructed via global approaches such as Voronoi or Delaunay triangulation methods (Basavanhally et al., 2009) or via local approach such as the FLock method (Lu et al., 2021). After cell graph construction, features related to edge length and node density are then extracted to predict disease outcome. Figure 1.1 show graphs constructed using both the global and local approaches. The PLS-DA method is design to capture only the global information uncover from the global graph approaches, important information involving local spatial interaction may be left un-

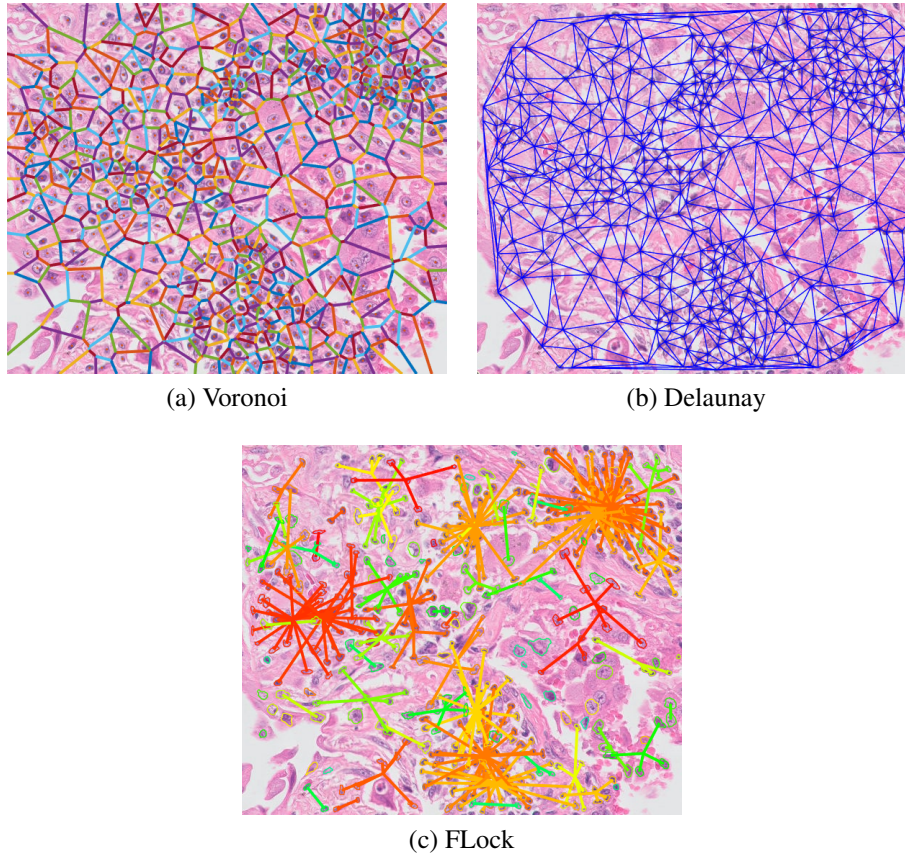


Figure 1.1: Cell nuclei graphs

exploited. Since both the global and local features are useful in the context of cancer grading, a method that capture both the global and local information is highly desirable. Also, capturing both the global and local information will lead to a better approach for modeling the tissue micro-environment.

1.3 Motivation of Study

Since global and local structures are both important for many pattern classification problems, a method that can reduce the dimensionality of a data set while preserving both of its global and local structures is thus highly desirable. The most current trends in feature extraction methods in general are seeing more and more approaches embracing the global and local structure preserving (GLSP) framework and are reporting very

promising results (Abeo et al., 2019; Lee, 2018; Song and Shi, 2018; Wan et al., 2018; Yao et al., 2018; Zhao et al., 2019a; Zhao and Jia, 2018). Their approaches are similar which involves developing an objective function that includes both the global and local structure features of a data set. The optimal direction for projection is thus the optimal solution of the objective. The resulting optimization problem is also equivalent to a generalized eigenvalue problem and can be solved using existing high performance computational methods for eigenvalue problems.

1.4 Research Objectives

The overall aim of this study is to develop an effective subspace learning (feature extraction) technique. Specifically, this study aims to enhance the overall performance of PLS-DA in feature extraction for complex high dimensional datasets. The objectives of this study include:

1. To show the effectiveness of global local PLS-DA method in discrimination and classification of face and chemical datasets.
2. To propose variants of the PLS-DA method which are more effective for discrimination and classification problems.
3. To propose new nonlinear subspace learning techniques based on the improved techniques to handle situations where strong nonlinear relationship exist between sets of observed data.
4. To show the effectiveness of the nonlinear extensions of the global local PLS-DA method in discrimination and classification of chemical and gene expression datasets.

1.5 Research Contribution

In an effort to make PLS-DA more effective for feature extraction and discrimination, we propose modifications to PLS-DA termed locality preserving PLS-DA (LPPLS-DA) and neighborhood preserving PLS-DA (NPPLS-DA), which are based on the GLSP framework. The local geometric structure of data is integrated into PLS-DA, which was almost never considered in the literature. We employed two different criteria to model the local geometric structure and obtain two feature extraction algorithms. Two efficient algorithms are developed to solve the resulting optimization problems, and their computational complexities are carefully discussed. The new algorithms are interesting in a number of perspectives.

1. The proposed LPPLS-DA and NPPLS-DA algorithms are fundamentally based on discriminant and spectral graph analysis. They are designed to solve different criteria from the PLS-DA method.
2. LPPLS-DA shares some similar properties with NPPLS-DA. Both techniques aim to discover both the global and the local geometric structure of data. However, their objective functions are entirely different.
3. Both LPPLS-DA and NPPLS-DA construct graphs over labeled data points to uncover the intrinsic discriminant structure in the data.
4. Both LPPLS-DA and NPPLS-DA are linear techniques which makes them suitable for practical applications. They may be conducted in the original space or in the reproducing kernel Hilbert space (RKHS) into which data points are mapped. This approach gives rise to nonlinear variants of LPPLS-DA and NPPLS-DA

called kernel LPPLS-DA (KLPPLS-DA) and kernel NPPLS-DA (KNPPLS-DA) respectively.

1.6 Scope of Study

This study focuses on improving the performance of PLS-DA in dealing with complex high dimensional datasets. Several approaches are proposed to refine the PLS-DA modeling practice strategies, especially in complex datasets such as multi-class datasets, imbalanced datasets, highly nonlinear datasets and manifold learning. All of the proposed approaches need to construct a graph to capture the local geometric structure of the data. Since, PLS-DA is already a supervised dimension reduction technique, we used the knowledge of the class labels while constructing the graph in the newly proposed variants of PLS-DA. This way, the discriminating ability of the PLS-DA technique can be enhanced to a higher extent.

The newly proposed approaches are applied to a large number of complex high dimensional datasets including face datasets, chemical datasets and biomedical datasets. Note that, all the datasets used in this study are publicly available as benchmark datasets for evaluating the performance of newly proposed machine learning techniques.

1.7 Thesis Organization

The organization of this thesis is as follows. Chapter 2 provides a detailed review of some of the most popular subspace learning techniques. Specifically, we give a detailed review of classical and manifold based subspace learning techniques. Chapter 3 intro-

duces the newly proposed locality preserving PLS-DA (LPPLS-DA) algorithm, plus a detailed computational analysis of LPPLS-DA. Extensive experimental results were also carried out on face databases to demonstrate the effectiveness of the LPPLS-DA method. The LPPLS-DA algorithm is also applied to complex chemical datasets and the experimental results are presented in Chapter 4. In chapter 5, a nonlinear extension of the LPPLS-DA algorithm in reproducing kernel Hilbert space (RKHS) is introduced to handle situations in which data are highly nonlinear. A detailed computational analysis of the nonlinear version of LPPLS-DA as well as extensive experimental results on gene expression datasets are also presented in Chapter 5. In Chapter 6, we introduce two new algorithms called neighborhood preserving PLS-DA (NPPLS-DA) and uncorrelated neighborhood preserving PLS-DA (UNPPLS-DA) for discriminant feature extraction. The extensive experimental results on face databases are also presented in Chapter 6. The kernel extensions of NPPLS-DA and UNPPLS-DA in reproducing kernel Hilbert space (RKHS), and the computational analysis of the algorithms are presented in Chapter 7. The extensive experimental results on some spectra datasets are also presented in Chapter 7. Finally, we provide some concluding remarks and suggestions for future research direction in Chapter 8.

CHAPTER 2

SUBSPACE LEARNING

Subspace learning is a framework applicable in research areas such as machine learning and pattern recognition where data samples are represented as points in high-dimensional spaces. Learning in high dimensional spaces becomes challenging because the performance of learning algorithms drops drastically as the number of dimensions increases. This phenomenon is known as “*curse of dimensionality*”. Thus, subspace learning techniques are first employed to reduce the dimensionality of the data before other learning techniques are applied. Subspace learning techniques consist of classical dimensionality reduction techniques and manifold learning techniques (Li and Allinson, 2009). In this chapter, a detailed review of some of the most popular subspace learning techniques is provided.

2.1 Classical Dimensionality Reduction Techniques

With the recent advances in computer technologies, there has been an explosion in the amount of data generated, stored and analyze. Most of this data are high dimensional in nature, ranging from several hundreds to thousands. Clustering or classification of such high dimensional datasets is almost infeasible. Thus, classical dimensionality reduction algorithms are used to map the high dimensional data into a lower dimensional subspace prior to the application of the conventional clustering or classification algorithms. The most popular algorithms for this purpose are principal component analysis (PCA) (Bartenhagen et al., 2010; Ma and Dai, 2011; Ma and Kosorok, 2009;

Yeung and Ruzzo, 2001), linear discriminant analysis (LDA) (Brereton, 2009; Cai et al., 2007; Hastie et al., 1995) and partial least squares discriminant analysis (PLS-DA) (Boulesteix and Strimmer, 2007; Nguyen and Rocke, 2002a,0; Pérez-Enciso and Tenenhaus, 2003; Tan et al., 2004). PCA is an unsupervised dimension reduction algorithm which aims at maximizing the variance of the new representations in the lower dimensional subspace. While LDA and PLS-DA are supervised dimensionality reduction algorithms which attempts to maximize between class covariance in the projected space. In addition to maximizing the between class covariance, LDA also attempts to minimize within class covariance in the projected space. These algorithms have proved successful for dimension reduction in many fields of research (Bahreini et al., 2019; Lee et al., 2018a; Li et al., 2020; Mas et al., 2020; Sitnikova et al., 2020; Xie et al., 2019; Yang et al., 2018; Zhao et al., 2018). The classical PCA, LDA and PLS-DA algorithms are linear dimensionality reduction algorithms and their performance can be restrictive when handling highly nonlinear datasets (Cao et al., 2011; Mika et al., 1998). To overcome this limitation, nonlinear extensions of PCA (Bartenhagen et al., 2010; Liu et al., 2005; Schölkopf et al., 1997), LDA (Baudat and Anouar, 2000; Cai et al., 2011) and PLS-DA (Song et al., 2018; Srinivasan et al., 2013; Štruc and Pavešić, 2009) through “kernel trick” have been proposed. The main idea of the kernel based techniques is to map the data into a feature space using a nonlinear mapping function. For a properly chosen nonlinear mapping function, an inner product can be defined in the feature space by a kernel function without defining the nonlinear mapping explicitly. The nonlinear extensions of the PCA and PLS-DA algorithms usually outperforms the linear PCA and PLS-DA algorithms when the data is highly nonlinear. In what follows, we give a brief review of the classical PCA, LDA and PLS-DA algorithms.

2.1.1 Principal Component Analysis

PCA is a well-known feature extraction technique in machine learning and pattern recognition. PCA seeks directions on which the data points are distributed with maximum variance. Given a set of n data points $\mathbf{x}_1, \dots, \mathbf{x}_n \in R^m$. Let the vector $\mathbf{z} = \{z_1, \dots, z_n\}$ represent the m -dimensional data points such that $z_i = \mathbf{w}^T \mathbf{x}_i$ is the one-dimensional map (representation) of \mathbf{x}_i ($i = 1, 2, \dots, n$), and $\mathbf{w} \in R^m$ denotes the transformation vector. The variance of the data points in the one-dimensional space can be calculated as follows (Jolliffe, 1986):

$$\frac{1}{n} \sum_{i=1}^n (z_i - \bar{z}_i)^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - \bar{z}_i)^2 \quad (2.1)$$

where $\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i$. Equation (2.1) can be reduced to

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (z_i - \bar{z}_i)^2 &= \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - \frac{1}{n} \sum_{i=1}^n z_i)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - \frac{1}{n} \sum_{i=1}^n \mathbf{w}^T \mathbf{x}_i)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - \mathbf{w}^T \bar{\mathbf{x}})^2 \\ &= \mathbf{w}^T \left(\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^T (\mathbf{x}_i - \bar{\mathbf{x}}) \right) \mathbf{w} \\ &= \mathbf{w}^T \mathbf{S} \mathbf{w} \end{aligned} \quad (2.2)$$

where $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^T (\mathbf{x}_i - \bar{\mathbf{x}})$ denotes the variance of the data points in the original space. Consequently, the objective function of PCA is given as follows:

$$\max_{\mathbf{w}^T \mathbf{w} = 1} \mathbf{w}^T \mathbf{S} \mathbf{w} \quad (2.3)$$

Using Lagrange multiplier method, the objective function (2.3) can be converted into an eigenproblem. Let

$$L(\mathbf{w}, \lambda) = \mathbf{w}^T \mathbf{S} \mathbf{w} - \lambda (\mathbf{w}^T \mathbf{w} - 1) \quad (2.4)$$

where λ is the Lagrange multiplier. Differentiating $L(\mathbf{w}, \lambda)$ and setting the result to zero, one can get

$$\frac{\partial L(\mathbf{w}, \lambda)}{\partial \mathbf{w}} = \mathbf{S} \mathbf{w} - \lambda \mathbf{w} = 0 \quad (2.5)$$

or,

$$\mathbf{S} \mathbf{w} = \lambda \mathbf{w} \quad (2.6)$$

where \mathbf{w} is the eigenvector of \mathbf{S} and λ is the corresponding eigenvalue.

2.1.2 Linear Discriminant Analysis

LDA aims at finding a lower dimensional subspace in which data samples from the same class remain close to each other while data samples from different class are mapped far apart. Given a set of n data points $\mathbf{x}_1, \dots, \mathbf{x}_n \in R^m$ belonging to C different classes. The LDA method solve the following objective function (Cai et al., 2007):

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}, \quad (2.7)$$

$$\mathbf{S}_b = \sum_{c=1}^C n_c (\boldsymbol{\mu}^{(c)} - \boldsymbol{\mu})(\boldsymbol{\mu}^{(c)} - \boldsymbol{\mu})^T, \quad (2.8)$$

$$\mathbf{S}_w = \sum_{c=1}^C \sum_{j=1}^{n_c} (\mathbf{x}_j^{(c)} - \boldsymbol{\mu}^{(c)})(\mathbf{x}_j^{(c)} - \boldsymbol{\mu}^{(c)})^T, \quad (2.9)$$

where $\boldsymbol{\mu}$ denotes the global centroid, $\boldsymbol{\mu}^{(c)}$ denotes the centroid of the c th class, n_c denotes the size of data samples in the c th class and $\mathbf{x}_j^{(c)}$ denotes the j th sample in the c th class. The matrices $\mathbf{S}_b \in R^{m \times m}$ and $\mathbf{S}_w \in R^{m \times m}$ are called the between-class and the within-class matrices, respectively. Let

$$L(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} \quad (2.10)$$

To optimize the objective function (2.7), LDA requires the derivative of $L(\mathbf{w})$ and sets it to zero (Fukunaga, 2013):

$$\begin{aligned} \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} &= \frac{(\mathbf{w}^T \mathbf{S}_w \mathbf{w})(2\mathbf{S}_b \mathbf{w}) - (\mathbf{w}^T \mathbf{S}_b \mathbf{w})(2\mathbf{S}_w \mathbf{w})}{(\mathbf{w}^T \mathbf{S}_w \mathbf{w})^2} \\ &= \frac{\mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} - \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_w \mathbf{w})^2} \mathbf{S}_w \mathbf{w} \\ &= 0 \end{aligned} \quad (2.11)$$

Equation (2.11) can be reduced to

$$\mathbf{S}_b \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w}. \quad (2.12)$$

where $\lambda = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$. Thus, the optimal \mathbf{w} can be computed as the generalized eigenvector of \mathbf{S}_b and \mathbf{S}_w corresponding to the eigenvalue λ .

2.1.3 Partial Least Squares

PLS is a well known method for modeling the linear relationship between two sets of observed variables. The method is widely used as a feature extraction method to deal with undersampled and multi-collinearity issues usually encountered in high dimensional data (Ahmad et al., 2006; Jia et al., 2016). Given two set of observed variables $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times m}$ and $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n]^T \in \mathbb{R}^{n \times N}$. PLS decomposes the zero mean data matrices $\bar{\mathbf{X}}$ and $\bar{\mathbf{Y}}$ into the following form (Rosipal and Krämer, 2005):

$$\begin{aligned}\bar{\mathbf{X}} &= \mathbf{T}\mathbf{P}^T + \mathbf{E} \\ \bar{\mathbf{Y}} &= \mathbf{U}\mathbf{Q}^T + \mathbf{F}\end{aligned}\tag{2.13}$$

where the $n \times d$ matrices \mathbf{T} and \mathbf{U} represents the score matrices of the d extracted components, \mathbf{P} and \mathbf{Q} are the $m \times d$ and $N \times d$ loading matrices of \mathbf{X} and \mathbf{Y} respectively, and the $n \times m$ matrix \mathbf{E} and the $n \times N$ matrix \mathbf{F} correspond to residual matrices of \mathbf{X} and \mathbf{Y} respectively. The PLS method is based on the nonlinear iterative partial least squares (NIPALS) algorithm (Geladi and Kowalski, 1986; Wold, 1975) which find weight vectors \mathbf{w} and \mathbf{c} such that

$$[\text{cov}(\mathbf{t}, \mathbf{u})]^2 = \max_{\mathbf{w}^T \mathbf{w} = \mathbf{c}^T \mathbf{c} = 1} [\text{cov}(\bar{\mathbf{X}}\mathbf{w}, \bar{\mathbf{Y}}\mathbf{c})]^2\tag{2.14}$$

where $\text{cov}()$ denotes the sample covariance between variables. The outline of the NIPALS algorithm can be summarized as follows:

Step 1: Randomly initializes \mathbf{u} , usually \mathbf{u} is set to be one of the columns of $\bar{\mathbf{Y}}$

Step 2: Compute the $\bar{\mathbf{X}}$ weights \mathbf{w} :

$$\mathbf{w} = \bar{\mathbf{X}}^T \mathbf{u} / \mathbf{u}^T \mathbf{u} \quad (2.15)$$

\mathbf{w} can be normalized, i.e. $\|\mathbf{w}_1\| = 1$.

Step 3: Compute the $\bar{\mathbf{X}}$ scores \mathbf{t} :

$$\mathbf{t} = \bar{\mathbf{X}} \mathbf{w} \quad (2.16)$$

Step 4: Compute the $\bar{\mathbf{Y}}$ weights \mathbf{c} :

$$\mathbf{c} = \bar{\mathbf{Y}}^T \mathbf{t} / \mathbf{t}^T \mathbf{t} \quad (2.17)$$

Step 5: Compute an updated set of $\bar{\mathbf{Y}}$ scores \mathbf{u} :

$$\mathbf{u} = \bar{\mathbf{Y}} \mathbf{c} \quad (2.18)$$

Step 6: Test for convergence on the change in \mathbf{t} , if \mathbf{t} converges proceeds to the next step, else return to Step 2.

Step 7: Deflate the data matrices $\bar{\mathbf{X}}$ and $\bar{\mathbf{Y}}$:

$$\begin{aligned} \bar{\mathbf{X}} &= \bar{\mathbf{X}} - \mathbf{t}(\mathbf{t}^T \bar{\mathbf{X}}) / \mathbf{t}^T \mathbf{t} \\ \bar{\mathbf{Y}} &= \bar{\mathbf{Y}} - \mathbf{t}(\mathbf{t}^T \bar{\mathbf{Y}}) / \mathbf{t}^T \mathbf{t} \end{aligned} \quad (2.19)$$

The PLS method has been used in discrimination problems (i.e., separating distinct data samples) and classification problems (i.e., assigning new data samples to predefined groups) (Bai et al., 2017; Nguyen and Rocke, 2002b). In this case, the input data

matrix \mathbf{Y} is replaced by a dummy (class membership) matrix containing class information and the procedure is called partial least squares discriminant analysis (PLS-DA). The objective function of PLS-DA is as follows:

$$\mathbf{w}^* = \max_{\mathbf{w}^T \mathbf{w} = 1} [\text{cov}(\bar{\mathbf{X}} \mathbf{w}, \mathbf{Y})]^2 \quad (2.20)$$

where \mathbf{Y} denotes the class membership matrix defined as:

$$\mathbf{Y} = \begin{pmatrix} 1_{n_1} & 0_{n_1} & \dots & 0_{n_1} \\ 0_{n_2} & 1_{n_2} & \dots & 0_{n_2} \\ \vdots & \vdots & \ddots & \vdots \\ 0_{n_C} & 0_{n_C} & \dots & 1_{n_C} \end{pmatrix} \quad (2.21)$$

where n_i (for $i = 1, 2, \dots, C$) represents the number of samples in the i -th class, $\sum_{i=1}^C n_i = n$ (total number of samples), and 0_{n_i} and 1_{n_i} are $n_i \times 1$ vectors of zeros and ones respectively. For example, if the data set contains two classes, then the matrix \mathbf{Y} is designed as a single-column vector with entries of 1 for all samples in the first class and 0 for

samples in the second class, i.e.

$$\mathbf{Y} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix},$$

Further, if the data have three classes, then the \mathbf{Y} matrix is encoded with three columns as follows,

$$\mathbf{Y} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \end{pmatrix},$$

One may choose to centre the class membership matrix \mathbf{Y} to have zero mean. Since

$$\begin{aligned} [\text{cov}(\bar{\mathbf{X}}\mathbf{w}, \mathbf{Y})]^2 &= \frac{1}{(n-1)^2} (\mathbf{Y}^T \bar{\mathbf{X}}\mathbf{w})^T (\mathbf{Y}^T \bar{\mathbf{X}}\mathbf{w}) \\ &= \frac{1}{(n-1)^2} \mathbf{w}^T \mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X} \mathbf{w}, \end{aligned} \quad (2.22)$$

the objective function (2.20) can be rewritten in the following equivalent form:

$$\max_{\mathbf{w}^T \mathbf{w} = 1} \mathbf{w}^T \bar{\mathbf{X}}^T \mathbf{Y} \mathbf{Y}^T \bar{\mathbf{X}} \mathbf{w} \quad (2.23)$$

Using Lagrange multiplier method, the objective function (2.23) can be reduced to an eigenproblem of the form:

$$\bar{\mathbf{X}}^T \mathbf{Y} \mathbf{Y}^T \bar{\mathbf{X}} \mathbf{w} = \lambda \mathbf{w} \quad (2.24)$$

Thus, the optimal weight (projection) vector \mathbf{w} in (2.23) can be obtained as the eigenvector of $\bar{\mathbf{X}}^T \mathbf{Y} \mathbf{Y}^T \bar{\mathbf{X}}$ corresponding to the eigenvalue λ in (2.24). It was shown that the eigenstructure (2.24) is basically that of a slightly altered version of the between class scatter matrix in LDA (Aminu and Ahmad, 2019; Barker and Rayens, 2003). Therefore, what PLS-DA does is basically maximizing between class separation.

2.2 Manifold Learning Techniques

The classical PCA, LDA and PLS-DA methods aim at preserving the global Euclidean structure of the data space; if data points happen to reside on nonlinear submanifold embedded in the high dimensional ambient space, this can pose a problem for these methods. Thus, several manifold based dimension reduction algorithms have been proposed to discover the local manifold structure. These algorithms include Laplacian eigenmaps (LE) (Belkin and Niyogi, 2002,0), locally linear embedding (LLE) (Roweis and Saul, 2000), neighborhood preserving embedding (NPE) (He et al., 2005a) and locality preserving projections (LPP) (He and Niyogi, 2004). These methods are designed to determine a subspace where local structure of data points are well preserved, but how well they are able to capture the global structure of a dataset is still not well understood. In what follows, we give a brief review of the LE, LLE, LPP and NPE algorithms.

2.2.1 Laplacian Eigenmap

Laplacian eigenmaps (LE) (Belkin and Niyogi, 2002) is a local nonlinear subspace learning technique based on spectral graph theory. The method attempts to preserve the local geometrical structure of data after dimension reduction. Specifically, LE seek an embedding such that nearby points on the manifold are mapped close to each other in the low dimensional subspace. Suppose $\mathbf{x}_1, \dots, \mathbf{x}_n$ denotes the set of n data points sampled from an underlying manifold \mathcal{M} embedded in a high dimensional ambient space R^m . LE first construct a graph G with the data points \mathbf{x}_i ($i = 1, \dots, n$) as nodes and a weight matrix \mathbf{S} that assign weights to the edges between the nodes. The weight matrix \mathbf{S} can be computed as follows:

$$S_{ij} = \begin{cases} 1; & \text{if } \mathbf{x}_i \in N_p(\mathbf{x}_j) \text{ or } \mathbf{x}_j \in N_p(\mathbf{x}_i) \\ 0; & \text{otherwise.} \end{cases} \quad (2.25)$$

where $N_p(\mathbf{x}_i)$ denotes the set of p nearest neighbors of \mathbf{x}_i . Let $\mathbf{z} = (z_1, \dots, z_n)^T$ denotes the map determined by LE where z_i is the one-dimensional map of \mathbf{x}_i ($i = 1, \dots, n$). LE realize the optimal map by solving the following minimization problem:

$$\min_{\mathbf{z}} \sum_{i,j=1}^n (z_i - z_j)^2 S_{ij} \quad (2.26)$$

The minimization problem (2.26) can be reduced to

$$\begin{aligned}
\min_{\mathbf{z}} \sum_{i,j=1}^n (z_i - z_j)^2 S_{ij} &= \min_{\mathbf{z}} \sum_{i=1}^n (z_i^2 + z_j^2 - 2z_i z_j) S_{ij} \\
&= \min_{\mathbf{z}} \left(\sum_{i=1}^n z_i^2 D_{ii} + \sum_{j=1}^n z_j^2 D_{jj} - 2 \sum_{i,j=1}^n z_i z_j S_{ij} \right) \\
&= 2 \min_{\mathbf{z}} \left(\sum_{i=1}^n z_i^2 D_{ii} - \sum_{i=1}^n z_i^2 S_{ij} \right) \\
&= 2 \min_{\mathbf{z}} (\mathbf{z}^T \mathbf{D} \mathbf{z} - \mathbf{z}^T \mathbf{S} \mathbf{z}) \\
&= 2 \min_{\mathbf{z}} \mathbf{z}^T (\mathbf{D} - \mathbf{S}) \mathbf{z} \\
&= \min_{\mathbf{z}} 2 \mathbf{z}^T \mathbf{L} \mathbf{z} \tag{2.27}
\end{aligned}$$

where $\mathbf{L} = \mathbf{D} - \mathbf{S}$ is the graph Laplacian (Chung and Graham, 1997) and \mathbf{D} is a diagonal matrix whose entries are column (or row) sum of \mathbf{S} , $D_{ii} = \sum_j S_{ij}$. In order to remove an arbitrary scaling factor in the embedding, the following constraint is impose:

$$\mathbf{z}^T \mathbf{D} \mathbf{z} = 1$$

Finally, the minimization problem (2.26) reduces to

$$\arg \min_{\mathbf{z}^T \mathbf{D} \mathbf{z} = 1} \mathbf{z}^T \mathbf{L} \mathbf{z} \tag{2.28}$$

Minimizing the objective function (2.28) is an attempt to ensure that if \mathbf{x}_i and \mathbf{x}_j are close on the manifold, then their low dimensional representations z_i and z_j are close as well.

Solving the minimization problem (2.28) is equivalent to finding the eigenvector

corresponding to the smallest eigenvalue of the following generalized eigen-problem

$$\mathbf{L}\mathbf{z} = \lambda\mathbf{D}\mathbf{z} \quad (2.29)$$

Equivalently, the embedding \mathbf{z} can be obtained as the eigenvector corresponding to the largest eigenvalue of the following generalized eigen-problem

$$\mathbf{S}\mathbf{z} = \lambda\mathbf{D}\mathbf{z} \quad (2.30)$$

2.2.2 Locally Linear Embedding

Locally linear embedding (LLE) (Roweis and Saul, 2000) is another nonlinear subspace learning technique which is also based on spectral graph theory. The main idea in LLE is that each data point resides on a local linear patch of a manifold and can be reconstructed by a linear combination of its nearest neighbors. Suppose $\mathbf{x}_1, \dots, \mathbf{x}_n$ are the set of n data points sampled from an underlying manifold \mathcal{M} embedded in a high dimensional ambient space R^m . The reconstruction errors are then measured by the cost function:

$$\varepsilon(\mathbf{S}) = \sum_{i=1}^n \left\| \mathbf{x}_i - \sum_{j=1}^n S_{ij}\mathbf{x}_j \right\|^2. \quad (2.31)$$

The weights S_{ij} characterized the contribution of the data point \mathbf{x}_j to the reconstruction of \mathbf{x}_i . To compute the weights S_{ij} , the cost function (2.31) is minimized subject to two constraints: 1) the rows of the weight matrix sum to one, i.e., $\sum_j S_{ij} = 1$, 2) each data point \mathbf{x}_i is reconstructed only from its nearest neighbors, enforcing $S_{ij} = 0$ if \mathbf{x}_i and \mathbf{x}_j are not neighbors. Let $\mathbf{z} = (z_1, \dots, z_n)^T$ denotes the map of the original data points to a line where z_i represents \mathbf{x}_i ($i = 1, \dots, n$). The LLE algorithm determines a neighbor-

hood preserving mapping by minimizing the following embedding cost function:

$$\varphi(\mathbf{z}) = \sum_{i=1}^n (z_i - \sum_{j=1}^n S_{ij}z_j)^2. \quad (2.32)$$

Let

$$p_i = z_i - \sum_{j=1}^n S_{ij}z_j \quad (i = 1, \dots, n)$$

which can be written in vector form as

$$\begin{aligned} \mathbf{p} &= \mathbf{z} - \mathbf{S}\mathbf{z} \\ &= (\mathbf{I} - \mathbf{S})\mathbf{z} \end{aligned}$$

The embedding cost function (2.32) can be reduced to

$$\begin{aligned} \varphi(\mathbf{z}) &= \sum_{i=1}^n (z_i - \sum_{j=1}^n S_{ij}z_j)^2. \\ &= \sum_{i=1}^n (p_i)^2 \\ &= \mathbf{p}^T \mathbf{p} \\ &= \mathbf{z}^T (\mathbf{I} - \mathbf{S})^T (\mathbf{I} - \mathbf{S}) \mathbf{z} \end{aligned} \quad (2.33)$$

Thus, the embedding \mathbf{z} that minimizes the cost function (2.32) is given by the eigenvector corresponding to the smallest eigenvalue of the following eigen-problem:

$$(\mathbf{I} - \mathbf{S})^T (\mathbf{I} - \mathbf{S}) \mathbf{z} = \lambda \mathbf{z} \quad (2.34)$$

2.2.3 Locality Preserving Projections

Locality preserving projections (LPP) (He and Niyogi, 2004; He et al., 2005b) is basically a linear approximation of the nonlinear LE technique. Similar to the LE technique, LPP seek an embedding that preserves the local geometrical structure of the original data. Given a set of n data points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{M}$ and \mathcal{M} is a nonlinear manifold embedded in R^m , LPP models the local geometrical structure of the data by an adjacency graph G with the data points \mathbf{x}_i ($i = 1, \dots, n$) as nodes and a weight matrix \mathbf{S} that assign weights to the edges between the nodes. The weight matrix \mathbf{S} can be define as follows:

$$S_{ij} = \begin{cases} \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{t}\right); & \text{if } \mathbf{x}_i \in N_p(\mathbf{x}_j) \text{ or } \mathbf{x}_j \in N_p(\mathbf{x}_i) \\ 0; & \text{otherwise.} \end{cases} \quad (2.35)$$

where $N_p(\mathbf{x}_i)$ denotes the set of p nearest neighbors of \mathbf{x}_i . Let z_i denotes a one dimensional map of \mathbf{x}_i ($i = 1, \dots, n$). LPP minimizes the following objective function:

$$\sum_{i,j=1}^n (z_i - z_j)^2 S_{ij}. \quad (2.36)$$