

論文 / 著書情報  
Article / Book Information

題目(和文)	マルチウェイデータ解析のための特徴抽出および識別法
Title(English)	Feature Extraction and Classification Methods for Multi-way Data Analysis
著者(和文)	横田達也
Author(English)	Tatsuya Yokota
出典(和文)	学位:博士(工学), 学位授与機関:東京工業大学, 報告番号:甲第9466号, 授与年月日:2014年3月26日, 学位の種別:課程博士, 審査員:山下 幸彦,高田 潤一,阿部 直也,花岡 伸也,杉山 将, Andrzej Cichocki
Citation(English)	Degree:Doctor (Engineering), Conferring organization: Tokyo Institute of Technology, Report number:甲第9466号, Conferred date:2014/3/26, Degree Type:Course doctor, Examiner:,,,,,
学位種別(和文)	博士論文
Type(English)	Doctoral Thesis

Feature Extraction  
and Classification Methods  
for Multi-way Data Analysis

Tatsuya YOKOTA

Department of International Development Engineering  
Graduate School of Science and Engineering  
Tokyo Institute of Technology

February 2014

# Abstract

The world is filled with data which have various types such as numeric data, non-numeric data, continuous valued data, discrete valued data, scalar data, vector data, matrix data, and higher-dimensional array data. Continuous valued and multi-dimensional array data is denoted by the “multi-way data” which can be processed by the computers. Thus, the multi-way data includes a scalar, a matrix, and a higher-dimensional array of continuous values. Techniques of multi-way data analysis play important roles in wide research fields such as audio processing, image processing, bioinformatics, chemometrics, and brain science.

In this thesis, after the techniques of feature extraction and classification methods for multi-way data analysis are summarized; new criteria and solving algorithms are proposed for them, and their advantages are demonstrated by experiments. Matrix/tensor decomposition has been widely used for visualization, clustering, blind source separation, and dimensionality reduction of multi-way data. We often impose some constraints such as orthogonality, sparsity, and nonnegativity based on the characteristics of target data. Techniques for the constraints such as singular value decomposition (SVD), penalized matrix decomposition (PMD), and nonnegative matrix factorization (NMF) have been well studied.

Function approximation in NMF is a technique to represent nonnegative feature vectors by linear combinations of basis functions. For example, if we choose smooth basis functions, then smooth feature vectors can be obtained. Zdunek proposed this technique as the GRBF-NMF method which uses the Gaussian radial basis functions for function approximation. However, its algorithm, in which the QP optimization and the active-set algorithm run alternately and iteratively, is so slow that it can not be applied to large-scale problems. Then, I propose a new fast algorithm for the function approximation by NMF, and choose more effective basis functions for function approximation. Furthermore, I extend this method to the nonnegative Tucker decomposition and the nonnegative canonical polyadic (CP) decomposition.

The common and individual feature extraction is a very important concept for data analysis because real world data always have some common and individual features. For examples all human faces have two eyes, two ears, a nose, and a mouth (common feature); but, their positions and shapes are different as individual personalities (individual feature). However, existing techniques for it are not efficient very well, and their techniques have been developed only for matrices. Then, I propose a tensor based common and individual feature extraction method in CP and Tucker model, and impose orthogonality,

sparsity, and nonnegativity constraints to both models.

Classification techniques are used for various objectives such as brain computer interface (BCI), fingerprinting identification, hand-written character recognition, and face recognition. In recent years, the support vector machine (SVM) has been the most popular classifier; however, SVM does not always provide the best performance. The criterion of SVM consists of the hinge-loss minimization and a regularization. In this research, I propose two novel weighted regularization methods for a variety of applications. Furthermore, I apply these regularization methods to the conventional SVM.

The Fisher discriminant analysis (FDA) is another famous method for classification. However, FDA does give an optimal projection only for Gaussian distributions with equal covariance matrices. In other words, FDA is not optimal in the case of heteroscedastic Gaussian distributions. In this research, I propose a novel criterion for FDA including a correction term based on the Bhattacharyya distance which is closely related to classification rate. Furthermore, the Chernoff distance based criterion and its kernelized version are proposed as its extensions.

The final proposition of this thesis is a new criterion of classifier based on the maximum a posteriori (MAP) estimation for a binary problem without estimating the posterior probability, called the quadratically constrained maximum a posteriori (QCMAP) estimation. The QCMAP consists of the maximization of the expectation of a cost function, which is derived from the maximum a posteriori probability, and a quadratic constraint. This criterion is highly general since its forms include least squares regressions (LSRs) and a SVM. I propose several efficient classifiers from the criterion by selecting various weight functions.

Finally, I show the experimental results of all proposed methods to demonstrate their advantages for the feature extraction and the pattern classification. In the experiments, the proposed smooth NMF/NTF algorithms using function approximation are applied to reconstruction of two and three dimensional array data, blind source separation, part-based representation to compare with other NMF methods (unconstrained, sparse and smooth NMF algorithms); the proposed common and individual feature extraction algorithms based on Tucker and CP decomposition models are applied to reconstruction and multi-way blind source separation; the proposed regularized SVM, Chernoff based FDA, and QCMAP classifiers are applied to several toy well-known benchmarks and real-world brain computer interface (BCI) problems to compare with other classification methods.

# Acknowledgements

I would like to express my appreciation and gratitude to all those who supported me to complete my thesis. First, I thank to my advisers Prof. Yukihiro Yamashita and Prof. Andrzej Cichocki for giving me many advisement and supports. They have provided an environment conducive to leaning and research quality. Prof. Yamashita has introduced me to the field of interesting research on pattern recognition and machine learning from a mathematical point of view. Prof. Cichocki has introduced me to the field of interesting research on matrix/tensor based data analysis and common and individual feature extraction. Next, I also thank to Prof. Toru Wakahara and Dr. Hitoshi Sakano for giving me many advisement. Dr. Rafal Zdunek discussed with me in a research topic. Prof. Yoshikazu Washizawa supported me to create the environment of research. Ms. Wakako Honjo supported me for clerical work. I also thank to my colleagues of Yamashita Laboratory and Cichocki Laboratory for their encouragements and influence. Finally, I would like to acknowledge the support of my parents.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Background . . . . .	1
1.1.1	Motivation for tensor based feature extraction . . . . .	2
1.1.2	Motivation for pattern classification . . . . .	4
1.1.3	Social impacts . . . . .	4
1.2	Problems of existing methods and proposed methods . . . . .	5
1.2.1	FastGRBF . . . . .	5
1.2.2	LCPTD/LTD . . . . .	5
1.2.3	WRSVM . . . . .	5
1.2.4	CFDA/KCFDA . . . . .	6
1.2.5	QCMAP estimation . . . . .	6
1.3	Organization . . . . .	7
1.4	Notations . . . . .	8
<b>2</b>	<b>Existing methods and algorithms for Data Analysis</b>	<b>10</b>
2.1	Frequency Analysis . . . . .	10
2.2	Time-Frequency Analysis . . . . .	10
2.3	Dimensionality Reduction . . . . .	11
2.3.1	Principal Component Analysis . . . . .	12
2.3.2	Sparse Principal Component Analysis . . . . .	13
2.3.3	Example of Dimensionality Reduction . . . . .	15
2.4	Blind Source Separation . . . . .	17
2.4.1	Independent Component Analysis . . . . .	19
2.4.2	Kurtosis based ICA . . . . .	21
2.4.3	Neg-entropy based ICA . . . . .	23
<b>3</b>	<b>Matrix/Tensor based Feature Extraction</b>	<b>28</b>
3.1	Matrix decomposition model . . . . .	28
3.1.1	Singular Value Decomposition . . . . .	29
3.1.2	Penalized Matrix Decomposition . . . . .	32

3.1.3	Nonnegative Matrix Factorization . . . . .	34
3.1.4	Why is NMF important and well-researched? . . . . .	38
3.1.5	Combined methods of several constraints . . . . .	39
3.1.6	Function approximation based model . . . . .	43
3.2	Tensor algebra . . . . .	45
3.3	CP decomposition . . . . .	49
3.3.1	ALS algorithm . . . . .	51
3.3.2	Deflation based algorithm . . . . .	52
3.3.3	HALS algorithm . . . . .	52
3.3.4	Sparse/Nonnegative CP decomposition . . . . .	53
3.4	Tucker decomposition . . . . .	54
3.4.1	HOOI algorithm for orthogonal Tucker decomposition . . . . .	55
3.4.2	Sparse Tucker decomposition . . . . .	56
3.4.3	Nonnegative Tucker decomposition . . . . .	57
3.5	Common and Individual Feature Extraction . . . . .	58
3.5.1	JIVE . . . . .	59
3.5.2	Group NMF . . . . .	60
<b>4</b>	<b>Supervised Feature Extraction and Classification</b>	<b>62</b>
4.1	Pattern Recognition . . . . .	62
4.1.1	Heuristic Approach vs. Statistical Approach . . . . .	62
4.1.2	Supervised Learning . . . . .	63
4.2	Supervised Feature Extraction methods . . . . .	64
4.2.1	CSP filter . . . . .	65
4.2.2	CSSP filter . . . . .	67
4.2.3	SBCSP filter . . . . .	69
4.3	Classification methods . . . . .	70
4.3.1	Generative Model vs Discriminant Model . . . . .	71
4.3.2	Kernel Method . . . . .	73
4.3.3	Least Square Regression . . . . .	75
4.3.4	Fisher Linear Discriminant Analysis . . . . .	79
4.3.5	Support Vector Machine . . . . .	81
<b>5</b>	<b>Proposed Feature Extraction methods</b>	<b>86</b>
5.1	Improvements of NMF with function approximation . . . . .	86
5.1.1	Fast algorithm for GRBF-NMF . . . . .	86
5.1.2	Extensive Bases for the fastGRBF-NMF . . . . .	88
5.1.3	Reduction of the size of Basis Matrix . . . . .	90
5.1.4	Tensor Extension . . . . .	91

5.1.5	CP model . . . . .	92
5.2	Linked CP Tensor Decomposition . . . . .	94
5.2.1	LCPTD-HALS algorithm . . . . .	95
5.3	Linked Tucker Decomposition . . . . .	96
5.3.1	Orthogonal LTD method . . . . .	97
5.3.2	Sparse LTD . . . . .	98
5.3.3	Nonnegative LTD . . . . .	99
<b>6</b>	<b>Proposed Classifiers</b>	<b>103</b>
6.1	QCMAP . . . . .	103
6.1.1	QCMAP estimation . . . . .	103
6.1.2	Theoretical Interpretation of QCMAP estimation . . . . .	104
6.1.3	Derivation of Mathematical Programming Problem . . . . .	105
6.1.4	Relation to Hinge Loss Minimization . . . . .	106
6.1.5	Relation to the Tikhonov Regularization . . . . .	106
6.1.6	Uniqueness of QCMAP Solutions . . . . .	107
6.1.7	Solution Method for QCMAP and Approximation of Cost Function . . . . .	107
6.1.8	Equivalence of QCMAP to LSR and SVM . . . . .	111
6.1.9	Gaussian QCM Classifier . . . . .	116
6.1.10	Extended GQCM Classifiers . . . . .	118
6.2	Weighted regularization SVM (WRSVM) . . . . .	121
6.2.1	Basic Support Vector Machines . . . . .	121
6.2.2	Novel Weighted Regularization . . . . .	122
6.2.3	Analytical Calculation of Regularization Matrices . . . . .	123
6.2.4	Novel Classifiers . . . . .	123
6.3	Chernoff FDA . . . . .	124
6.3.1	Bhattacharyya and Chernoff Distances . . . . .	125
6.3.2	Chernoff FDA . . . . .	127
6.3.3	Kernel Extension of CFDA . . . . .	129
<b>7</b>	<b>Experiments</b>	<b>131</b>
7.1	fastGRBF based methods . . . . .	131
7.1.1	Low-rank approximation . . . . .	131
7.1.2	Blind source separation . . . . .	137
7.1.3	Block-based part representation . . . . .	140
7.2	LCPTD . . . . .	144
7.2.1	Low rank approximation . . . . .	144
7.2.2	Linked Blind Source Separation . . . . .	148
7.3	LTD . . . . .	149

7.3.1	Feature extraction for classification . . . . .	149
7.4	Classification . . . . .	150
7.4.1	One-tailed t-test . . . . .	150
7.4.2	WRSVM . . . . .	151
7.4.3	CFDA . . . . .	153
7.4.4	QCMAP . . . . .	157
<b>8</b>	<b>Conclusions</b>	<b>165</b>
8.1	Summary . . . . .	165
8.2	Open problems . . . . .	167
8.2.1	Feature Extraction . . . . .	167
8.2.2	Classification . . . . .	167

# List of Figures

1.1	Flowchart for visualization . . . . .	2
1.2	Flowchart for BCI . . . . .	3
1.3	Flow of chapters in this thesis . . . . .	6
2.1	Wavelet transform . . . . .	11
2.2	Concept of objective function of the PCA: black points stand for individual $\mathbf{x}_n$ , red line stands for the transformed space, and green arrows stand for the difference between the original signals $\mathbf{x}_n$ and projected signals $\mathbf{b}_n$ . The PCA minimize the sum of length of green arrows. . . . .	12
2.3	Results of dimensionality reduction . . . . .	16
2.4	Results of ERP averages with target stimuli and other stimuli . . . . .	17
2.5	Blind Source Separation . . . . .	18
2.6	Individual phases of ICA . . . . .	20
2.7	Kurtosis . . . . .	21
2.8	Example of kurtosis based ICA . . . . .	23
2.9	Neg-entropy based ICA (using $g_1$ ) . . . . .	27
2.10	Example of real image separation . . . . .	27
3.1	Comparison between only nonnegative factors and including negative factors . . . . .	38
3.2	Basis function $\phi_n$ . . . . .	44
3.3	Tensors for various orders . . . . .	46
3.4	Vectorization of 3d-tensor . . . . .	47
3.5	Matricization of each mode . . . . .	48
3.6	First mode tensor product . . . . .	49
3.7	All mode product . . . . .	50
3.8	Matricization of all mode product . . . . .	50
3.9	The CP decomposition . . . . .	51
3.10	The Tucker3 decomposition . . . . .	55
4.1	Basic concept of pattern classification . . . . .	63

4.2	Basic concept of supervised learning . . . . .	64
4.3	Standard feature extraction scheme . . . . .	65
4.4	Classification scheme . . . . .	65
4.5	Primary motor cortex . . . . .	68
4.6	Topographies of CSP filter: BCI competition III(IVa) . . . . .	68
4.7	2D-plot of CSP feature: classes 1 and 2 stand for right hand and foot imageries, respectively. Horizontal and vertical axes in CSP feature space are first and last components, respectively. . . . .	68
4.8	System Flowchart of SBCSP: for example frequency bands of individual filters are 4-8Hz, 8-12Hz, ... , 36-40Hz. . . . .	69
4.9	Example of training samples . . . . .	73
4.10	Squared error . . . . .	76
4.11	Examples of LSR and kernelized LSR . . . . .	78
4.12	Regularization . . . . .	78
4.13	4-fold cross validation . . . . .	79
4.14	Fisher's Linear Discriminant Analysis . . . . .	80
4.15	Example of which FDA does not work well . . . . .	81
4.16	Basic concept of optimal hyperplane classifier . . . . .	82
4.17	Hinge loss . . . . .	84
5.1	Decomposition of GRBF-NMF model . . . . .	87
5.2	Linked Tucker Decomposition (LTD): in special case, we obtain STD for $L_n = J_n$ and ITD for $L_n = 0$ . . . . .	97
6.1	The QCMAP cost function and the hinge loss function . . . . .	107
6.2	Two solution areas for the QCMAP problem: solid and dotted lines depict the contour lines of $J(\mathbf{w})$ and $\mathbf{w}^T H \mathbf{w} = 1$ , respectively . . . . .	108
6.3	Approximation Function $h(r_n t)$ for changing $t$ . . . . .	109
6.4	Difference of Gaussian function: Here, $\kappa = 0.9$ is fixed and $\nu$ varies. . . . .	120
6.5	Misclassification rate and Hellinger distance . . . . .	126
7.1	PSNR of NMF results on 10 faces for various ranks . . . . .	132
7.2	Values of $\mathbf{W}$ and $\mathbf{A}$ obtained by the fastGRBF-NMF and the GRBF-NMF . . . . .	133
7.3	PSNR of NMF results on 165 faces for various ranks . . . . .	134
7.4	PSNR of NTF results on 165 (11×15) faces for various ranks . . . . .	135
7.5	Denosing and face reconstruction experiments on CBCL faces database: noise level (PSNR) is 16.2 dB. . . . .	136
7.6	Experimental result shown by iso-surfaces at $R = 30$ . . . . .	138
7.7	Visualizations and experimental results on NMFLAB: 'X_5smooth' . . . . .	141

7.8	Visualizations and experimental results on far-infrared spectra . . . . .	142
7.9	Results of fastGRBF-block-NMF . . . . .	143
7.10	Denoising result . . . . .	145
7.11	Parts-based representation . . . . .	146
7.12	Face images corrupted by additive noise and the reconstructed images (PSNR= 15 dB, $J = 40$ ): 1st column: original images, 2nd column: noisy images, 3rd column: ICPTD model, 4th column: LCPTD ( $L_n = 35$ ), 5th column: LCPTD with sparse constraint ( $L_n = 35$ ), 6th column: LCPTD with nonnegative constraint ( $L_n = 35$ ), 7th column: SCPTD model. . . . .	146
7.13	PSNRs for various noisy data sets . . . . .	147
7.14	Linked Multi-block Tensor Factorization . . . . .	147
7.15	Experiments for LTD . . . . .	149
7.16	t-distribution and reject/accept region with significant level 5% . . . . .	150
7.17	Error rates with difference of covariance matrices . . . . .	154
7.18	Error rates with difference of prior probabilities . . . . .	155
7.19	Samples and projected space ( $N_2/N_1 = 1.5$ ) . . . . .	155
7.20	Distributions in projected spaces ( $N_2/N_1 = 1.5$ ) . . . . .	163
7.21	Score plot of UCI data sets . . . . .	164

# List of Tables

1.1	Notations . . . . .	8
2.1	Approximation functions for neg-entropy based ICA . . . . .	26
7.1	Parameter setting . . . . .	133
7.2	Computational times for various algorithms . . . . .	133
7.3	Individual parameters of fastGRBF-block-NMF . . . . .	143
7.4	UCI Data sets . . . . .	152
7.5	Experimental results . . . . .	152
7.6	Properties of UCI datasets . . . . .	156
7.7	Maximized values of Chernoff distance by each algorithm . . . . .	157
7.8	Number of iterations for main updates in each method . . . . .	157
7.9	Classification errors for UCI dataset . . . . .	158
7.10	Computational time of optimizers . . . . .	159
7.11	Proposed Classifiers tested in the experiment . . . . .	160
7.12	Experimental results of UCI data sets . . . . .	160
7.13	Experimental results of a BCI data set by the weight function selection . . . . .	162

# Chapter 1

## Introduction

### 1.1 Background

The world is filled with data which have various types such as numeric data, non-numeric data, continuous valued data, discrete valued data, scalar data, vector data, matrix data, and higher-dimensional array data. Continuous valued and multi-dimensional array data are denoted by the “multi-way data” which can be processed by the computers. Thus, the multi-way data include a scalar, a matrix, and a higher-dimensional array of continuous values.

Techniques of multi-way data analysis play important roles in science and engineering such as chemometric data analysis, economic data analysis, psychological data analysis, audio processing, image processing, text mining, spectrometry analysis, bioinformatics, and brain signal analysis. In a more mathematical point of view, the multi-way data stand for  $N$ -dimensional arrays including a scalar ( $N = 0$ ), a vector ( $N = 1$ ), a matrix ( $N = 2$ ), and a higher-dimensional array ( $N \geq 3$ ). We call also an  $N$ -way array the “ $N$ d-tensor”. For example, an image data is a 2d-tensor and a video data is a 3d-tensor. These data are obtained from various measurement devices such as a camera for images/videos, a microphone for audio signals, a scanner for fingerprintings, and a MRI (magnetic resonance imaging) device for 3d-images of a brain. When we treat such multi-way data, it is undesirable to unfold a higher-dimensional array to a vector because it destroys the data structure.

One of purposes of multi-way data analysis is to visualize to show extracted key information to the analyst, and main research objective is to provide analysis methods and algorithms which are accurate, robust, and computationally efficient. Since such general analysis methods can be applied to various types of data, the social impacts of data analysis are very wide.

There are many applications of data analysis because there are many kinds of needs for individual subjects such as the brain signal analysis, image/video data analysis, and financial data analysis. In recent years, the multi-way data analysis is closely-watched and critically-needed for the “BIG DATA” analysis in the world. Firstly, I discuss the brain signal analysis. Improvements of brain computer interface (BCI) systems contribute to the nursing care and the rehabilitation. If we come to be able to diagnose Alzheimer disease in an early phase, we can suppress the symptoms to the minimum.

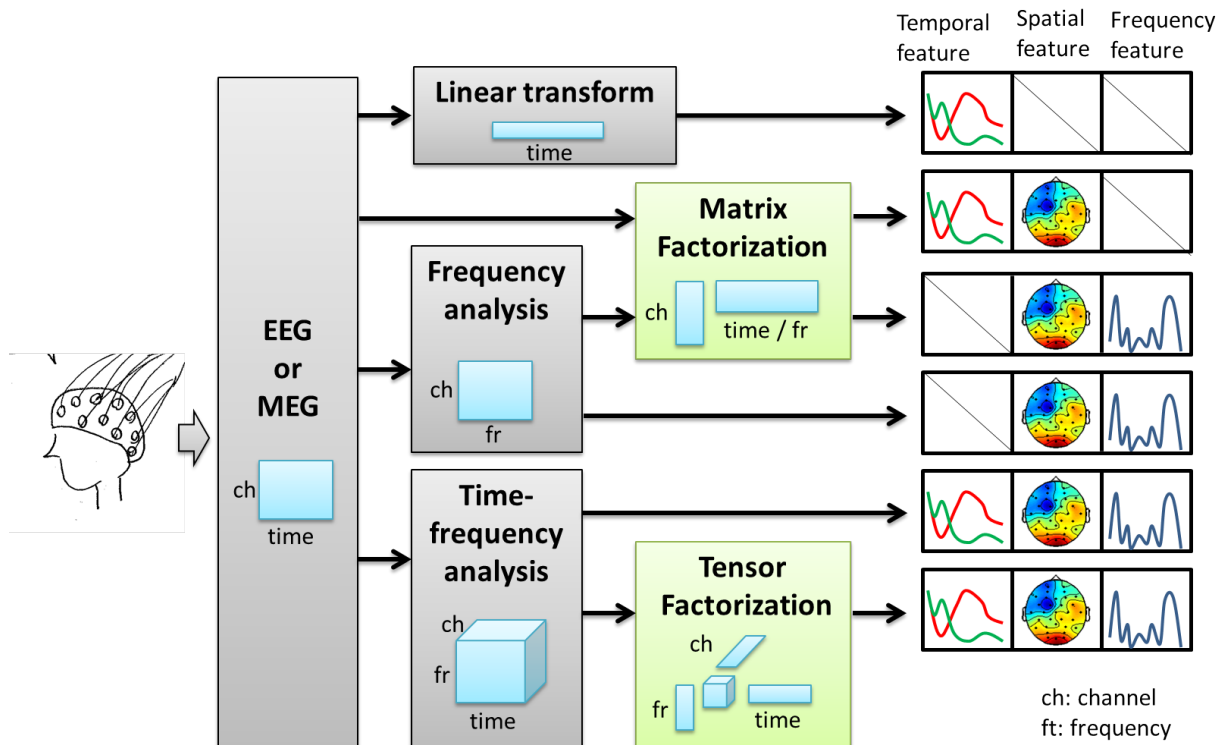


Figure 1.1: Flowchart for visualization

Improvements of brain fingerprinting systems contribute to the identification of the criminals. Secondly, I discuss the image analysis and recognition. Improvements of artificial intelligence for robots contribute to the industrial robots, the cleaning robots, the nursing-care robots, the rescue robots, and so on. Furthermore, improvements of artificial intelligence for security cameras contribute to the fire detection, the human tracking, the identification of criminals, and so on. Thirdly, I discuss the financial data analysis. If we can plan efficient economic strategies by data analysis, it contributes to the economic growth in the society. Precise prediction of stock prices will stabilize stock market. In other subjects, improvements of audio recognition contribute to the audio input systems of car navigation devices, mobile phones, computers, and so on; improvements of techniques of web data mining contribute to companies so that they can know customer's needs; improvements of weather prediction contribute to agriculture and preventing disaster; improvements of psychological data analysis contribute to curing mental disorders.

In this thesis, we focus on the tensor based feature extraction and classification methods.

### 1.1.1 Motivation for tensor based feature extraction

In this section, I introduce the motivation for tensor based feature extraction by focusing on time-series signals of EEG/MEG. Fig. 1.1 shows a flowchart for visualization of EEG/MEG signals. Temporal, spatial, and frequency features are often focused for this visualization. In this figure, several feature

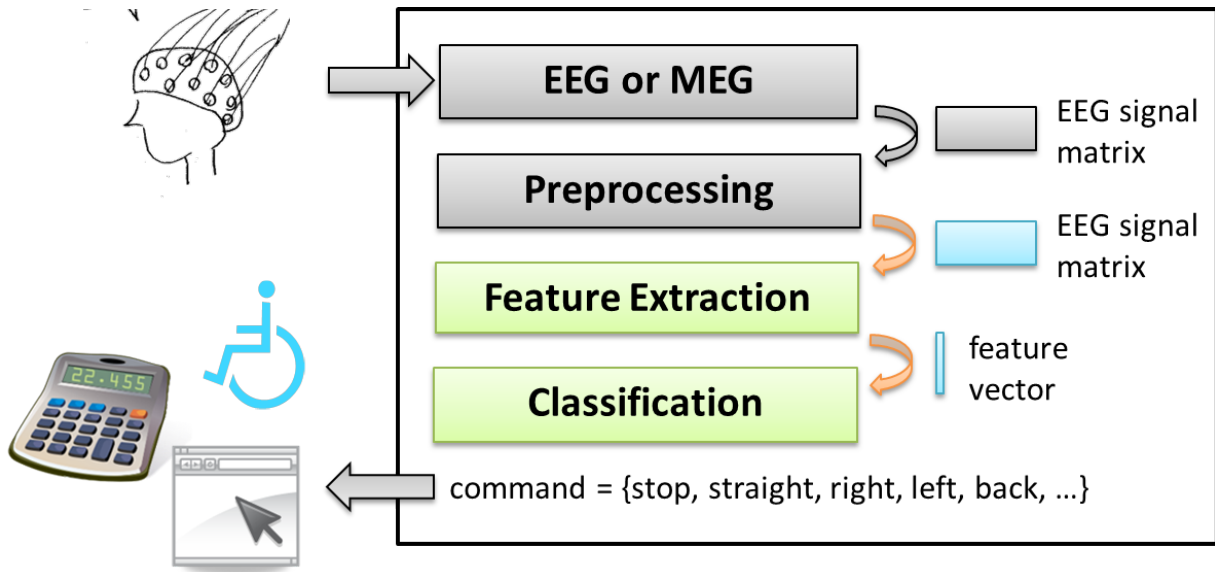


Figure 1.2: Flowchart for BCI

extraction methods are introduced. I explain frequency analysis, time-frequency analysis, and linear transform in Section 2.1, 2.2, 2.3 and 2.4. Matrix and tensor factorizations are explained in Chapter 3, and I discuss these techniques in this thesis, mainly.

From here, I mention about only EEG for simplicity; however almost all methods can be applied to also MEG. An EEG data is given as a  $(ch \times time)$ -matrix, where  $ch$  is the number of channels and  $time$  is the number of sample points in time domain. An EEG data stands for a set of sampled voltages of the multiple electrodes on a head while a subject carries out some tasks for a fixed time interval. In order to find some important factor from the EEG data, we need a large number of such samples. Thus the EEG dataset is given as a  $(ch \times time \times samples)$ -tensor, where I call a multi dimensional ( $N$ -way) array as a 'tensor' ( $N$ d-tensor). If we apply some time-frequency analysis technique to this tensor data, the 3d-tensor is transformed to a  $(ch \times time \times fr \times samples)$ -tensor (4d-tensor). The problem is how to process such a tensor data. One of the methods is to unfold a sample of tensor (e.g.,  $(ch \times time \times fr)$ -tensor) to a vector (e.g.,  $[(ch) \cdot (time) \cdot (fr)]$ -dimensional vector); however the unfolding destroys the data structure. Thus, tensor (including matrix) based feature extraction techniques are necessary for the brain signal processing, and also other multi-way data analysis such as video and fMRI image.

In general, matrix/tensor based feature extraction techniques have been widely studied to analyze various time-series signals, images, videos, physical spectra, and so on. The key-point of the matrix and tensor factorizations is to impose constraints such as orthogonality, sparsity, smoothness, and nonnegativity of factors. For example, EEGs are generally smooth, and frequency power spectrum, images, videos are always nonnegative. Therefore, the researches on algorithms and methods for matrix/tensor based feature extraction are very important. I introduce the details of matrix and tensor factorizations in Chapter 3.

### 1.1.2 Motivation for pattern classification

In this section, I introduce the motivation for pattern classification by focusing on EEG/MEG signal recognition as a subject. The pattern classification is a technique to assign some patterns into individual appropriate categories, automatically. For example, the EEG/MEG signal recognition is an important technique for BCI and the auto-diagnosis of the brain disease. BCI achieves to control a computer or a machine by using a brain signal without four limbs, for example typing text, inputting number, and controlling a wheel chair. These applications have possibilities of nursing care and welfare. Fig. 1.2 shows a flowchart for the EEG/MEG based BCI. Note that this flowchart is applied not only to BCI, but also to other objectives such as image recognition and text categorization. Preprocessing includes cutout, normalization, filtering, and so on. Feature extraction is to transform a EEG/MEG signal matrix (or other types of data) to a feature vector which includes important information for classification. Finally, classification is to assign a feature vector into a class which implies a command to control devices for BCI. In this thesis, we discuss the techniques of classification, mainly. Please note that the techniques of classification are not only for the brain signals, but also for other various types of patterns such as the handwritten characters, the fingerprintings, and the face images. Thus, the researches of classification are very important. I introduce the details of classification in Chapter 4.

### 1.1.3 Social impacts

The social impacts of multi-way data analysis are very wide and great. For example, if sufficient accuracy of pattern recognition which includes feature extraction and classification is achieved, the world would be significantly changed. Most of the works which persons have been done would be provided by machines, and persons have come to be dedicated to some more creative works. This technique has many possibilities such as the automated driving of car, the automated cleaning of building, the automated speech inputting, and the automated translation to different languages. Furthermore, if the auto-diagnosis of diseases becomes sufficiently accurate, people in regions where the number of medical doctors is not sufficient would be rescued by such systems.

On the other hand, there would be negative impacts caused by the world changes. Automation of works has a possibility to increase the jobless rate. For example, the automated driving system does away with the need to employ taxi drivers. There are risks that those systems are used for crimes such as the kidnapping, the smuggling, and the terror.

When the researchers publish their researches or new technologies, they should consider not only positive impacts but also negative impacts because their research might change the world significantly. Technology of nuclear energy provided new energy resources; however it threatened the world by the nuclear war. I hope that technologies complete the happiness of humankind.

## 1.2 Problems of existing methods and proposed methods

In this thesis, I discuss methods and algorithms for the matrix/tensor based feature extraction and the classification. The goal of this thesis is to contribute the progress of the feature extraction and the classification. In this section, I introduce existing methods of the matrix/tensor based feature extraction and classification, point out their problems, and propose new methods and algorithms to solve their problems, briefly.

### 1.2.1 FastGRBF

Function approximation by NMF is a technique to represent nonnegative feature vectors by linear combinations of basis functions. For example, if we choose smooth basis functions, then smooth feature vectors can be obtained. Zdunek proposed this technique as the GRBF-NMF method which uses the Gaussian radial basis functions for function approximation. However, its algorithm, in which the QP optimization and the active-set algorithm run alternately and iteratively, is so slow that it can not be applied to large-scale problems. Then, I propose a new fast algorithm for the function approximation by NMF, and choose more effective basis functions for function approximation. It is called the fastGRBF-NMF. Furthermore, I extend this method to the nonnegative Tucker decomposition (NTD) and the nonnegative canonical polyadic (CP) decomposition (NCPD), and they are called the fastGRBF-NTD and the fastGRBF-NCPD, respectively.

### 1.2.2 LCPTD/LTD

The common and individual feature analysis (CIFA) is a very important method for data analysis because real world data always have some common and individual features. For example all human faces have two eyes, two ears, a nose, and a mouse (common feature); but, their positions and shapes are different as individual personalities (individual feature). However, existing techniques for it are not very efficient, and their techniques have been developed only for matrix. Then, I propose a tensor based common and individual feature extraction method in CP and Tucker models, and impose orthogonality, sparsity, and nonnegativity constraints to both models. They are called the linked CP tensor decomposition (LCPTD) and the linked Tucker decomposition (LTD), respectively.

### 1.2.3 WRSVM

Support vector machine (SVM) has been the most popular classifier; however, SVM does not always provide the best performance. The criterion of SVM consists of the hinge-loss minimization and a regularization. In this research, I propose two novel weighted regularization methods for a variety of applications. Furthermore, I apply these regularization methods to the conventional SVM. It is called the SVM with weighted regularization (WRSVM)

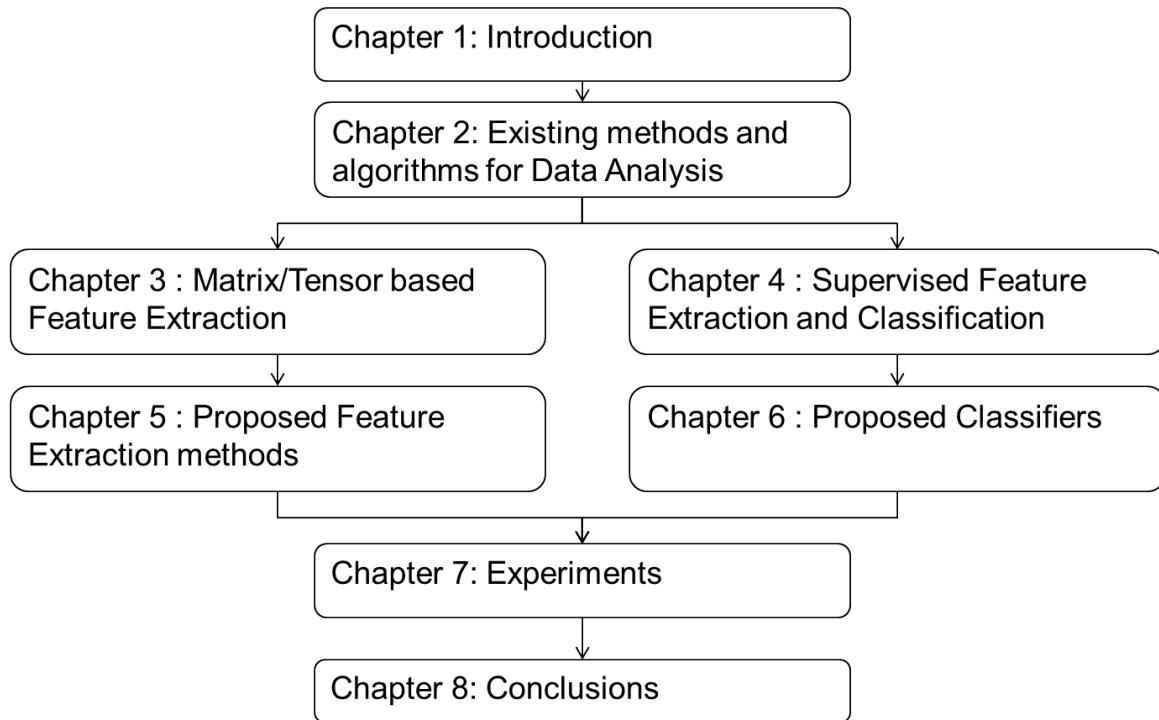


Figure 1.3: Flow of chapters in this thesis

#### 1.2.4 CFDA/KCFDA

Fisher discriminant analysis (FDA) is another famous method for classification. However, FDA does not give an optimal projection only for Gaussian distributions with equal covariance matrices. In other words, FDA is not optimal in the case of heteroscedastic Gaussian distributions. In this research, I propose a novel criterion for FDA including a correction term based on the Bhattacharyya distance which is closely related to classification rate. Furthermore, the Chernoff distance based criterion and its kernelized version are proposed as its extensions. They are called the Chernoff FDA (CFDA) and kernel CFDA (KCFDA), respectively.

#### 1.2.5 QCMAP estimation

Final proposition of this thesis is a new criterion of classifier based on the maximum a posteriori (MAP) estimation for a binary problem without estimating the posterior probability, called the quadratically constrained maximum a posteriori (QCMAP) estimation. The QCMAP consists of the maximization of the expectation of a cost function, which is derived from the maximum a posteriori probability and a quadratic constraint. This criterion is highly general since its forms include least squares regressions (LSRs) and a SVM. I propose several efficient classifiers from the criterion by selecting various weight functions.

### 1.3 Organization

This thesis is organized as follows. The flowchart of chapters is depicted in Fig. 1.3.

In Chapter 2, I introduce several existing methods for data analysis: the discrete Fourier transform (DFT), the discrete wavelet transform (DWT), the principal component analysis (PCA), the sparse principal component analysis (SPCA), and the independent component analysis (ICA).

In Chapter 3, I introduce matrix and tensor decomposition techniques. The most basic matrix decomposition is the singular value decomposition (SVD). SVD provides a full-rank complete decomposition which consists of orthonormal left and right singular matrices and a diagonal rectangular matrix. Truncated SVD (tSVD) is a low-rank approximation model of SVD. tSVD provides the minimum mean squared error between an original data and a decomposition model. The SVD and the tSVD features are obtained based on the orthogonality constraint; however, the orthogonality is not always efficient. We have to use also the sparsity, the smoothness and the nonnegativity constraints as the situation demands. The penalized matrix factorization (PMF) is a technique to impose the sparsity or the smoothness constraint to the matrix decomposition model, and the nonnegative matrix factorization (NMF) is a technique to impose the nonnegativity constraint to the matrix decomposition model. Furthermore, the sparse NMF, the smooth NMF methods are introduced. The tensor decomposition is an extended technique of the matrix decomposition to analyze multi-dimensional array (tensor) data. There are two models in the tensor decomposition: CP model and Tucker model. The CP model is a direct extension of the matrix decomposition, and the Tucker model is a more general extension. I introduce methods that impose the orthogonality, the sparsity and the nonnegativity constraints to both models. Furthermore, I introduce the common and individual feature analysis (CIFA), and several solution algorithms. The common and individual feature extraction is a very important technique for data analysis since its model is well fitted into real problems.

In Chapter 5, I propose new methods for the matrix and tensor based feature extractions including new fast algorithms for smooth nonnegative matrix and tensor factorizations, and tensor extensions of the CIFA.

In Chapter 4, I introduce a supervised-feature extraction and several basic classification methods: the common spatial pattern (CSP) filter, the kernel method, the least squares regression (LSR), the Fisher discriminant analysis (FDA), and the support vector machine (SVM).

In Chapter 6, I propose new classification methods such as the quadratically constrained maximum a posteriori estimation, SVM with weighted regularization, and the Chernoff FDA and its extension by the kernel method.

In Chapter 7, experiments are conducted to evaluate the proposed methods and their results are described. Finally, I conclude the research and provide future works in Chapter 8.

Table 1.1: Notations

$\mathbb{R}^N$	$N$ -dimensional Euclidean space
$\langle \mathbf{f}, \mathbf{g} \rangle$	inner product of two vectors $\mathbf{f}$ and $\mathbf{g}$
$\ \mathbf{f}\ $	Euclidean norm of $\mathbf{f}$
$f(i)$	$i$ -th component of $\mathbf{f}$
$\mathbf{I}_N$	$(N \times N)$ identity matrix
$\mathbf{A}^T$	transposition of a matrix $\mathbf{A}$
$\text{vec}(\mathbf{A})$	vectorization of $\mathbf{A}$
$[\mathbf{A}]_{ij}, \mathbf{A}(i, j), a_{ij}$	$(i, j)$ element of $\mathbf{A}$
$\text{tr}[\mathbf{A}]$	trace of $\mathbf{A}$
$\text{rank}(\mathbf{A})$	rank of $\mathbf{A}$

## 1.4 Notations

Table 1.1 describes a list of notations in this thesis. I denote individual variables by the following rules:

- Scalars are denoted by lowercase letters, e.g.,  $a \in \mathbb{R}$ .
- Vectors are denoted by boldface lowercase letters, e.g.,  $\mathbf{a} \in \mathbb{R}^I$ .
- Matrices are denoted by boldface capital letters, e.g.,  $\mathbf{A} \in \mathbb{R}^{I \times J}$ .
- Higher order tensors are denoted by underlined and boldface capital letter, e.g.,  $\underline{\mathbf{A}}$ .

Vectorization of a matrix is the transformation from the matrix to a vector. The vectorization of a matrix  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_J] \in \mathbb{R}^{I \times J}$  is given by

$$\text{vec}(\mathbf{A}) := \begin{pmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \vdots \\ \mathbf{a}_J \end{pmatrix} \in \mathbb{R}^{IJ}. \quad (1.1)$$

Next, I introduce four product rules of matrices. The Hadamard product (element-wise multiplication) is defined as

$$\mathbf{A} \otimes \mathbf{B} := \begin{pmatrix} a_{11}b_{11} & a_{12}b_{12} & \cdots \\ a_{21}b_{21} & a_{22}b_{22} & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix}. \quad (1.2)$$

In similar way, the element wise division is defined as

$$\mathbf{A} \oslash \mathbf{B} := \begin{pmatrix} \frac{a_{11}}{b_{11}} & \frac{a_{12}}{b_{12}} & \cdots \\ \frac{a_{21}}{b_{21}} & \frac{a_{22}}{b_{22}} & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix}. \quad (1.3)$$

Simply, we have

$$(\mathbf{A} \oplus \mathbf{B})_{ij} = a_{ij}b_{ij}, \quad (\mathbf{A} \oslash \mathbf{B})_{ij} = \frac{a_{ij}}{b_{ij}}. \quad (1.4)$$

Let  $\mathbf{A} \in \mathbb{R}^{I_a \times J_a}$  and  $\mathbf{B} \in \mathbb{R}^{I_b \times J_b}$ , the Kronecker product is defined by

$$\mathbf{A} \otimes \mathbf{B} := \begin{pmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix} \in \mathbb{R}^{I_a I_b \times J_a J_b}. \quad (1.5)$$

We have

$$(\mathbf{A} \otimes \mathbf{B})^T = \mathbf{A}^T \otimes \mathbf{B}^T, \quad (1.6)$$

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{A} \otimes \mathbf{B})^T = (\mathbf{A} \otimes \mathbf{B})(\mathbf{A}^T \otimes \mathbf{B}^T) \quad (1.7)$$

$$= \mathbf{A}\mathbf{A}^T \otimes \mathbf{B}\mathbf{B}^T. \quad (1.8)$$

The Khatri-Rao product is defined by

$$\mathbf{A} \odot \mathbf{B} := \begin{pmatrix} \mathbf{a}_1 \otimes \mathbf{b}_1 & \mathbf{a}_2 \otimes \mathbf{b}_2 & \cdots \end{pmatrix}, \quad (1.9)$$

$$= \begin{pmatrix} \text{vec}(\mathbf{b}_1 \mathbf{a}_1^T) & \text{vec}(\mathbf{b}_2 \mathbf{a}_2^T) & \cdots \end{pmatrix}. \quad (1.10)$$

We have

$$(\mathbf{A} \odot \mathbf{B})^T (\mathbf{A} \odot \mathbf{B}) = \begin{pmatrix} (\mathbf{a}_1^T \mathbf{a}_1)(\mathbf{b}_1^T \mathbf{b}_1) & (\mathbf{a}_1^T \mathbf{a}_2)(\mathbf{b}_1^T \mathbf{b}_2) & \cdots \\ (\mathbf{a}_2^T \mathbf{a}_1)(\mathbf{b}_2^T \mathbf{b}_1) & (\mathbf{a}_2^T \mathbf{a}_2)(\mathbf{b}_2^T \mathbf{b}_2) & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix} \quad (1.11)$$

$$= (\mathbf{A}^T \mathbf{A}) \otimes (\mathbf{B}^T \mathbf{B}). \quad (1.12)$$

## Chapter 2

# Existing methods and algorithms for Data Analysis

In this Chapter, I introduce existing methods and algorithms for multi-channel signal analysis including frequency analysis, time-frequency analysis, dimensionality reduction, and blind source separation.

### 2.1 Frequency Analysis

In this section, I introduce the Fourier transform as one of the techniques for frequency analysis. It is given by

$$F(\omega) = \int_{-\infty}^{\infty} x(t)e^{-i\omega t} dt. \quad (2.1)$$

In general,  $F(\omega)$  becomes a complex function, and we often focus on  $|F(\omega)|^2$  as a frequency power spectrum. Frequency is a very important factor for the brain signal analysis. Frequency bands of brain signals are called by various names. For example, a brain signal of which frequency components are in 8-13 Hz is called the alpha-wave. Furthermore, when we focus on the area of brain, the alpha-waves of the motor cortex, the temporal lobe, and the secondary somatosensory cortex are called the mu rhythm, the tau rhythm, and the sigma rhythm, respectively. Furthermore, the delta wave (0.1-4 Hz), the theta wave (4-7 Hz), the beta wave (12-30 Hz), and the gamma wave (25-100 Hz) are used in the brain signal analysis. In this way, it is very important to find out individual waves by frequency analysis.

### 2.2 Time-Frequency Analysis

Time-frequency analysis is a more important technique for the brain signal analysis than frequency analysis since we want to find out the temporal behavior for each frequency band. In this section, I introduce the wavelet transform (WT) as one of the techniques for time-frequency analysis. It is given

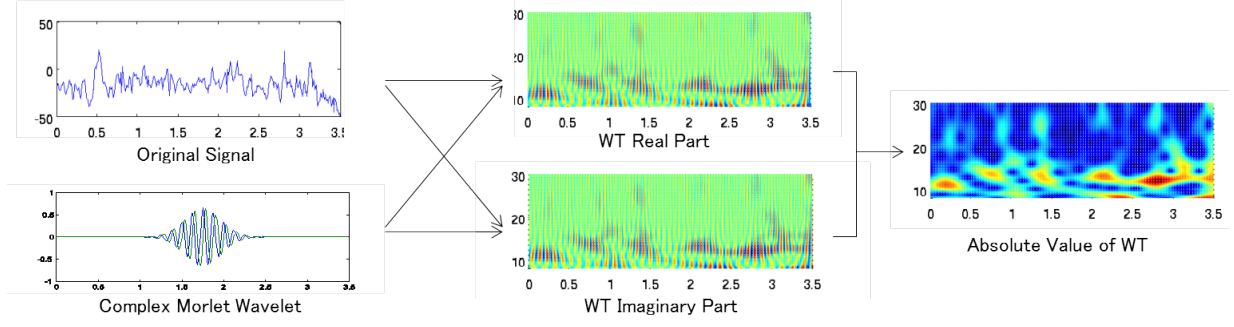


Figure 2.1: Wavelet transform

by

$$\mathbf{W}(a, b) = \frac{1}{\sqrt{a}} \int_{-\infty}^{\infty} x(t) \overline{\psi\left(\frac{t-b}{a}\right)} dt, \quad (2.2)$$

where  $x(t)$  is a signal,  $\psi(t)$  is a wavelet function.  $a$  and  $b$  stand for respectively variables of frequency domain and time domain. There are many kind of wavelets such as Haar wavelet, Meyer wavelet, Mexican Hat wavelet and Morlet wavelet. Fig. 2.1 shows an example of wavelet transform with the complex Morlet (CMOR) wavelet which is defined by

$$\psi_{f_b, f_c}(t) = \frac{1}{\sqrt{\pi f_b}} e^{i2\pi f_c t - (t^2/f_b)}, \quad (2.3)$$

where  $f_b$  and  $f_c$  are respectively a band width parameter and a frequency parameter. In this case,  $\mathbf{W}$  becomes complex matrix and we often focus on  $|\mathbf{W}|^2$  as a wavelet power spectrum matrix.

## 2.3 Dimensionality Reduction

When the dimension of data is very large, for example we analyze EEG signals which are recorded from 118 electrode channels to see an overall aspect of the brain without spatial information (e.g., P300), the dimensionality reduction techniques are quite useful. Furthermore, improvement of computational efficiency and noise reduction are also merits of the dimensionality reduction.

Its model is given by

$$\mathbf{b}_n = \mathbf{A}^T \mathbf{x}_n, \quad (2.4)$$

$$\mathbf{B} = \mathbf{A}^T \mathbf{X}, \quad (2.5)$$

where  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{M \times N}$  is a set of original signals,  $\mathbf{A} \in \mathbb{R}^{M \times R}$  is a transformation matrix, and  $\mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_N] \in \mathbb{R}^{R \times N}$  is a set of transformed signals. Basically, we set  $R$  to be smaller than  $M$  because we want to reduce the dimensionality of the signals and extract only the important components. The goal of this method is to obtain the transformation matrix  $\mathbf{A}$  which provides the meaningful feature  $\mathbf{b}_n$ .

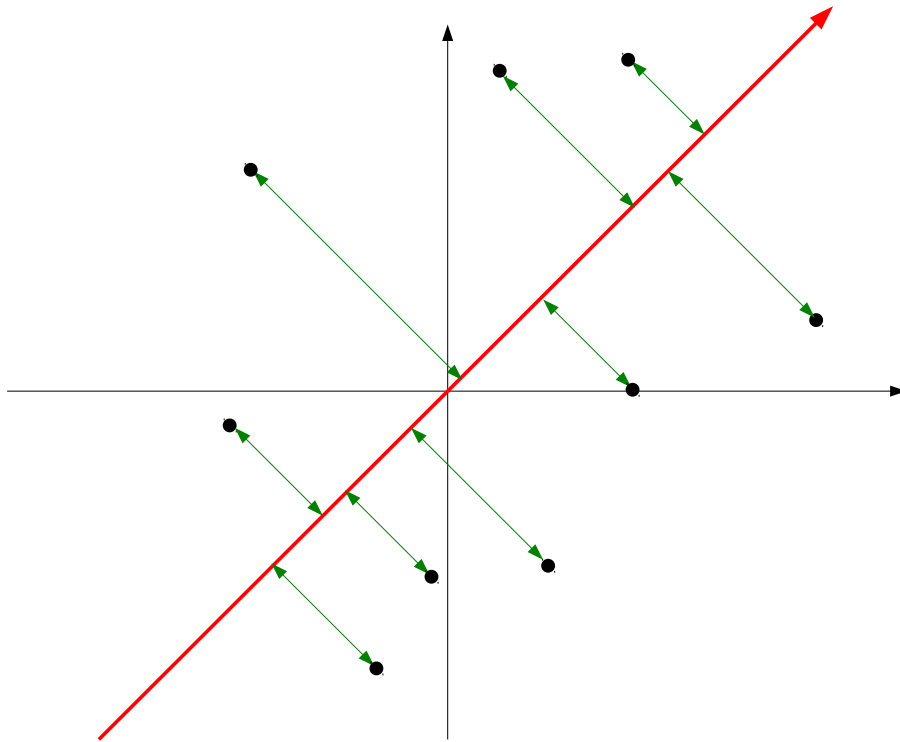


Figure 2.2: Concept of objective function of the PCA: black points stand for individual  $\mathbf{x}_n$ , red line stands for the transformed space, and green arrows stand for the difference between the original signals  $\mathbf{x}_n$  and projected signals  $\mathbf{b}_n$ . The PCA minimize the sum of length of green arrows.

### 2.3.1 Principal Component Analysis

Principal Component Analysis (PCA) [52] is a very typical method for feature extraction and dimensionality reduction. It is defined by

$$\begin{aligned} & \text{minimize} && \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{A}\mathbf{A}^T \mathbf{x}_n\|^2 = \|\mathbf{X} - \mathbf{A}\mathbf{A}^T \mathbf{X}\|_F^2 \\ & \text{subject to} && \mathbf{A}^T \mathbf{A} = \mathbf{I}_R, \end{aligned} \quad (2.6)$$

where  $\|\cdot\|_F$  stands for the Frobenius norm and  $\mathbf{I}_R \in \mathbb{R}^{R \times R}$  is an identity matrix. The objective function minimizes the Euclidean distance between the original signals  $\mathbf{x}_n$  and projected signals of  $\mathbf{b}_n$  (Fig. 2.2), and the constraint part imposes the orthogonality of the transformation matrix  $\mathbf{A}$ . This means that the PCA can be regarded as the minimization of the transformation loss.

On the other hand, the objective function can be transformed as

$$\begin{aligned}
\sum_{n=1}^N \|\mathbf{A}\mathbf{A}^T \mathbf{x}_n - \mathbf{x}_n\|^2 &= \|\mathbf{A}\mathbf{A}^T \mathbf{X} - \mathbf{X}\|_F^2 \\
&= \text{tr}[(\mathbf{A}\mathbf{A}^T \mathbf{X} - \mathbf{X})(\mathbf{A}\mathbf{A}^T \mathbf{X} - \mathbf{X})^T] \\
&= \text{tr}(\mathbf{A}\mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A}\mathbf{A}^T) - 2\text{tr}(\mathbf{A}\mathbf{A}^T \mathbf{X} \mathbf{X}^T) + \text{tr}(\mathbf{X} \mathbf{X}^T) \\
&= -\text{tr}(\mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A}) + \text{tr}(\mathbf{X} \mathbf{X}^T),
\end{aligned} \tag{2.7}$$

where  $\text{tr}(\cdot)$  describes ‘trace’. Since  $\text{tr}(\mathbf{X} \mathbf{X}^T)$  is a constant term, we can rewrite the PCA criterion as

$$\begin{aligned}
&\text{maximize} && \text{tr}(\mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A}), \\
&\text{subject to} && \mathbf{A}^T \mathbf{A} = \mathbf{I}_R.
\end{aligned} \tag{2.8}$$

Please note that this objective function is the same as  $\text{tr}(\mathbf{B}^T \mathbf{B})$ . Thus, the PCA can be also regarded as the maximization of the auto-correlations of the transformed signals  $\mathbf{b}_n$ .

### Solution of PCA

In this section, I provide a solution of the PCA criterion (2.8). First, its Lagrangian function is given by

$$L(\mathbf{A}, \boldsymbol{\Lambda}) = \text{tr}(\mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A}) - \text{tr}[(\mathbf{A}^T \mathbf{A} - \mathbf{I}_R)\boldsymbol{\Lambda}], \tag{2.9}$$

where  $\boldsymbol{\Lambda} \in \mathbb{R}^{R \times R}$  is a Lagrangian multiplier. Lagrange’s conditions are given by

$$\frac{\partial L}{\partial \mathbf{A}} = 2\mathbf{X} \mathbf{X}^T \mathbf{A} - 2\mathbf{A}\boldsymbol{\Lambda} = 0, \tag{2.10}$$

$$\frac{\partial L}{\partial \boldsymbol{\Lambda}} = \mathbf{A}^T \mathbf{A} - \mathbf{I}_R = 0. \tag{2.11}$$

The solution must satisfy the conditions and maximizes the PCA objective function. Let  $\lambda_1 \geq \dots \geq \lambda_M$  and  $\phi_1, \dots, \phi_M$  be respectively the eigenvalues and eigenvectors of  $(\mathbf{X} \mathbf{X}^T)$ , then the solution of the PCA criterion is given by

$$\mathbf{A}_{\text{PCA}} = [\phi_1, \dots, \phi_R]. \tag{2.12}$$

In other word, the solution is given as the eigenvectors of  $(\mathbf{X} \mathbf{X}^T)$  corresponding to the largest  $R$  eigenvalues.

### 2.3.2 Sparse Principal Component Analysis

Sparseness is an important property to analyze large scale data by sparse representations. Furthermore, sparseness often provides robust analysis for noisy data. From this reasons, sparse extensions of data analysis methods have been studied very well such as the sparse nonnegative matrix factorization (NMF) [48, 49, 57], the sparse canonical correlation analysis (CCA) [79, 105], and the sparse partial

least squares (SPLS) [15]. The sparse principal component analysis (SPCA) [53, 121, 105] can be regarded as the most basic method among them.

Jolliffe et al. proposed the SCoTLASS (Simplified Component Technique-LASSO) method [53] by adding the LASSO (least absolute shrinkage and selection operator) constraint [95] to the PCA. The SCoTLASS criterion is given by

$$\begin{aligned} & \text{maximize} && \text{tr}(\mathbf{a}_r^T \mathbf{X} \mathbf{X}^T \mathbf{a}_r), \\ & \text{subject to} && \mathbf{a}_r^T \mathbf{a}_r = 1, \mathbf{a}_h^T \mathbf{a}_r = 0 \text{ (for } r \geq 2 \text{ and } h < r), \\ & && \|\mathbf{a}_r\|_1 \leq t, \end{aligned} \quad (2.13)$$

where  $\|\mathbf{a}\|_1 := \sum_{m=1}^M |a_m|$  is a 11-norm, and  $1 \leq t \leq \sqrt{M}$  is a tuning parameter. The smaller  $t$ , the solution becomes more sparse. For example, when  $t = 1$ , we must have one nonzero  $a_m$ , and when  $t \geq \sqrt{M}$ , we get the PCA. Please note that the problem (2.13) have many local optima, thus it needs numerical optimization to estimate parameters. Jolliffe et al. proposed to use the projected gradient approach [14, 46] as a solution method for the SCoTLASS.

However, Zou et al. [121] noted that the SCoTLASS fails to achieve sufficient sparsity and there is a possibility that orthogonality and sparsity constraints of the PC vectors are not fully compatible, thus they proposed an alternate method for the sparse PCA by extending it to the elastic net regression [120]. The SPCA criterion is given by

$$\begin{aligned} & \text{minimize} && \|\mathbf{X} - \tilde{\mathbf{A}} \mathbf{A}^T \mathbf{X}\|_F^2 + \lambda \|\mathbf{A}\|_F^2 + \sum_{r=1}^R \lambda_1^{(r)} \|\mathbf{a}_r\|_1, \\ & \text{subject to} && \tilde{\mathbf{A}}^T \tilde{\mathbf{A}} = \mathbf{I}_R, \end{aligned} \quad (2.14)$$

where  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_R] \in \mathbb{R}^{M \times R}$  is a target sparse component,  $\tilde{\mathbf{A}} = [\tilde{\mathbf{a}}_1, \dots, \tilde{\mathbf{a}}_R] \in \mathbb{R}^{M \times R}$  is regarded as a relaxing principal component,  $\lambda$  and  $\lambda_1^{(r)}$  are nonnegative tuning parameters for penalty terms. Let consider  $\tilde{\mathbf{A}}_\perp \in \mathbb{R}^{M \times (M-R)}$  be an orthonormal matrix which satisfies that  $[\tilde{\mathbf{A}}, \tilde{\mathbf{A}}_\perp]$  is orthogonal, thus note that the first term of objective function (2.14) can be transformed as

$$\begin{aligned} \|\mathbf{X} - \tilde{\mathbf{A}} \mathbf{A}^T \mathbf{X}\|_F^2 &= \text{tr}(\mathbf{X} \mathbf{X}^T) - 2\text{tr}(\mathbf{X} \mathbf{X}^T \mathbf{A} \tilde{\mathbf{A}}^T) + \text{tr}(\tilde{\mathbf{A}} \mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A} \tilde{\mathbf{A}}^T) \\ &= \text{tr}(\tilde{\mathbf{A}}_\perp^T \mathbf{X} \mathbf{X}^T \tilde{\mathbf{A}}_\perp) + \text{tr}(\tilde{\mathbf{A}}^T \mathbf{X} \mathbf{X}^T \tilde{\mathbf{A}}) - 2\text{tr}(\tilde{\mathbf{A}}^T \mathbf{X} \mathbf{X}^T \mathbf{A}) + \text{tr}(\mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A}) \\ &= \text{tr}(\tilde{\mathbf{A}}_\perp^T \mathbf{X} \mathbf{X}^T \tilde{\mathbf{A}}_\perp) + \|\mathbf{X}^T \tilde{\mathbf{A}} - \mathbf{X}^T \mathbf{A}\|_F^2. \end{aligned} \quad (2.15)$$

Since the first term of (2.15) can be ignored for optimization, the first term of SPCA objective function can be replaced by the second term of (2.15).

### Solution algorithm for SPCA

The problem (2.14) can be solved by an alternating algorithm of two sub-optimizations. The first optimization updates  $\mathbf{A}$  with fixed  $\tilde{\mathbf{A}}$ , and the second optimization updates  $\tilde{\mathbf{A}}$  with fixed  $\mathbf{A}$ . Both

**Algorithm 1** SPCA algorithm

- 
- 1: **Input:**  $\mathbf{X}$ ,  $\lambda$ ,  $\lambda_1^{(r)}$  for  $r = 1, \dots, R$
  - 2: **Initialize:**  $\tilde{\mathbf{A}}$  as the PCA result
  - 3: **repeat**
  - 4:   Update  $\mathbf{a}_r$  to solve the problem (2.16) for all  $r$ ;
  - 5:   Update  $\tilde{\mathbf{A}}$  to solve the problem (2.19);
  - 6: **until** convergence
  - 7: Normalize  $\mathbf{A}$  by  $\mathbf{a}_r \leftarrow \mathbf{a}_r / \|\mathbf{a}_r\|$  for all  $r$ ;
  - 8: **Output:**  $\mathbf{A}$ ,  $\tilde{\mathbf{A}}$
- 

optimizations run alternately until the SPCA criterion is minimized. The first optimization problem is given by

$$\underset{\mathbf{a}_r}{\text{minimize}} \quad \|\mathbf{X}^T \tilde{\mathbf{a}}_r - \mathbf{X}^T \mathbf{a}_r\|^2 + \lambda \|\mathbf{a}_r\|^2 + \lambda_1^{(r)} \|\mathbf{a}_r\|_1. \quad (2.16)$$

Its problem is also called the elastic net problem [120], thus the solution is given by

$$\mathbf{a}_r(m) = \frac{\text{sgn}[\mathbf{v}_r(m)] \cdot \left[2|\mathbf{v}_r(m)| - \lambda_1^{(r)}\right]_+}{2(1 + \lambda)}, \quad (2.17)$$

where  $\mathbf{v}_r = \mathbf{X} \mathbf{X}^T \tilde{\mathbf{a}}_r$ ,  $[z]_+ := \max(z, \varepsilon)$  with a very small nonnegative parameter  $\varepsilon$  (e.g.,  $\varepsilon = 10^{-16}$ ), and  $\text{sgn}(\cdot)$  function is defined by

$$\text{sgn}(z) = \begin{cases} 1 & z > 0 \\ 0 & z = 0 \\ -1 & z < 0 \end{cases}. \quad (2.18)$$

The second optimization problem is given by

$$\begin{aligned} &\underset{\tilde{\mathbf{A}}}{\text{minimize}} \quad \|\mathbf{X} - \tilde{\mathbf{A}} \mathbf{A}^T \mathbf{X}\|_F^2, \\ &\text{subject to} \quad \tilde{\mathbf{A}}^T \tilde{\mathbf{A}} = \mathbf{I}_R. \end{aligned} \quad (2.19)$$

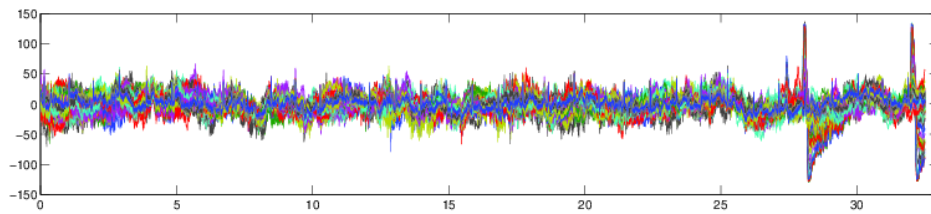
Its solution is given as

$$\tilde{\mathbf{A}} = \mathbf{U} \mathbf{V}^T, \quad (2.20)$$

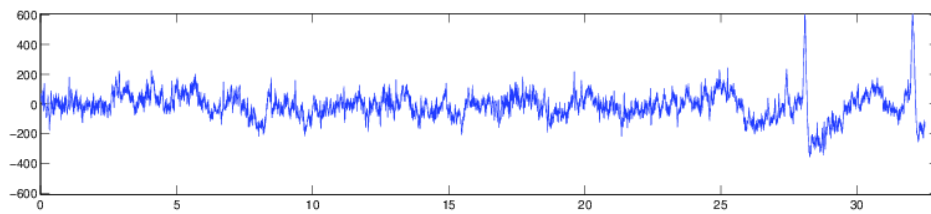
where  $\mathbf{U} \in \mathbb{R}^{M \times R}$  and  $\mathbf{V} \in \mathbb{R}^{R \times R}$  are respectively left and right factor matrices of singular value decomposition (SVD) of  $\mathbf{X} \mathbf{X}^T \mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^T$ . The details of SVD is described in Section 3.1.1. Finally, the SPCA algorithm can be summarized as Algorithm 1.

### 2.3.3 Example of Dimensionality Reduction

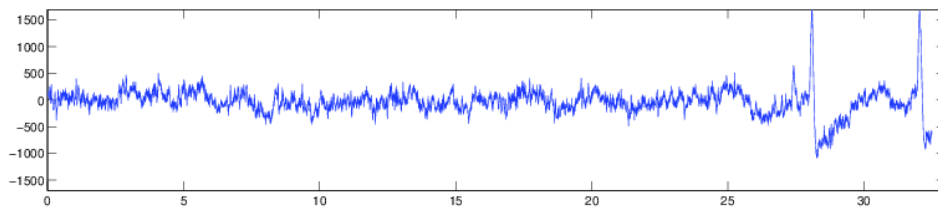
In this section, I introduce an example of application of dimensionality reduction. I used EEG signals from [8] of which sampling rate is 240 Hz and the number of channels is 64. This experiment was based



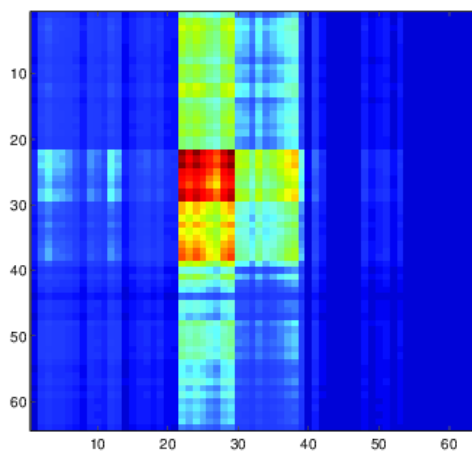
(a) 64 channels of EEG



(b) Embedded signal by the first PC component



(c) Embedded signal by the first SPC component



(d) Correlation matrix between PC (vertical) and SPC (horizontal) components

Figure 2.3: Results of dimensionality reduction

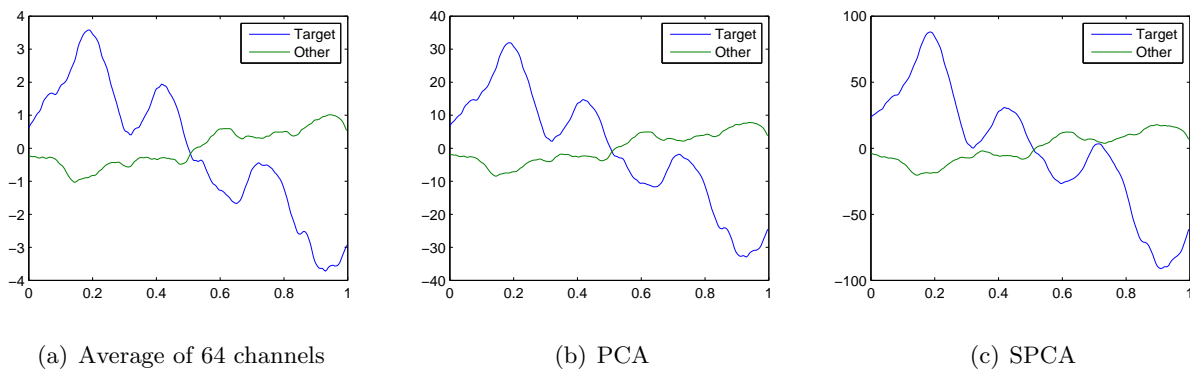


Figure 2.4: Results of ERP averages with target stimuli and other stimuli

on oddball paradigm, then I aimed to compare the event related potential (ERP) of target stimulus with ERP of other stimulus. Fig. 2.3 shows (a) the dataset of EEG signals, (b) reduced data by PCA, (c) reduced data by SPCA, and (d) correlation map between PC and SPC components. Their vertical and horizontal axis mean the signal amplitude and time [s], respectively. Seemingly, the results of PCA and SPCA are the almost same; however, we can see that only around 20-30 channels are important (correlated) from the correlation map. Fig. 2.4 shows ERPs of target and other stimuli compared by (a) average of 64 channels, (b) PCA, and (c) SPCA. Their vertical and horizontal axis mean the signal amplitude and time [s], respectively. In PCA and SPCA, the dimensionality was reduced from 64 into 1; however, results were very similar to original. Signal amplitudes of PCA/SPCA are significantly larger than original one because it is the first component including the largest eigenvalue for 64 eigenvalues. We can see that PCA and SPCA did not cut the important information.

## 2.4 Blind Source Separation

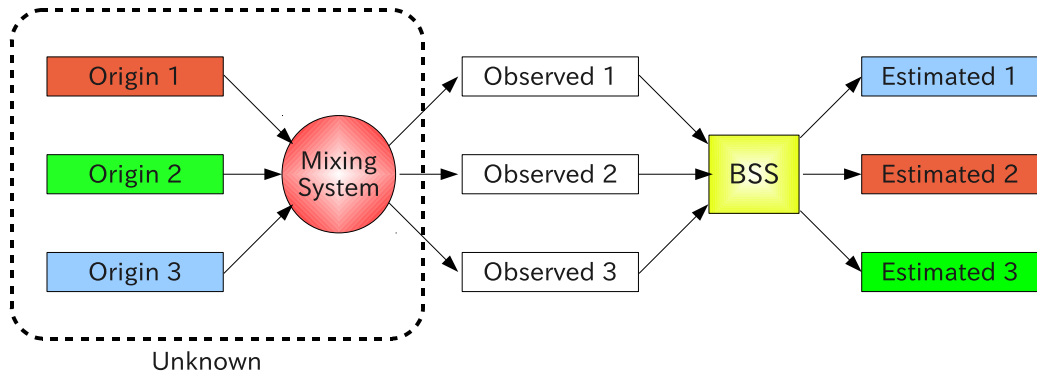
Blind Source Separation (BSS) is a method to estimate original signals from observed signals which consist of mixed original signals and noise. The cocktail party problem is a typical example of BSS. In case of the EEG signals, one electrode picks up signals from the surrounding numerous neurons, thus it can be considered that the regions covered by individual channels are overlapped each other (see Fig. 2.5).

Its model is given by

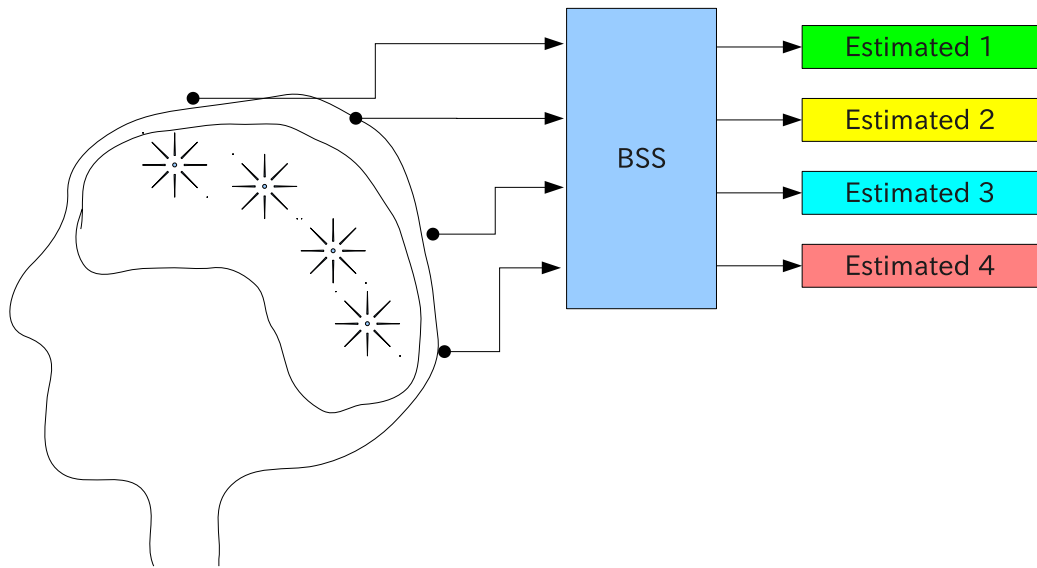
$$\mathbf{X} = \mathbf{B}\mathbf{S}, \quad (2.21)$$

where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_M]^T \in \mathbb{R}^{M \times N}$  and  $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_R]^T \in \mathbb{R}^{R \times N}$  are respectively a observed signal matrix and a source signal matrix, and  $\mathbf{B} \in \mathbb{R}^{M \times R}$  is a mixing matrix. When we regard  $\mathbf{x}_m$  and  $\mathbf{s}_r$  as functions  $x_m(n)$  and  $s_r(n)$  with respect to  $n$ , this model can be rewritten by

$$\mathbf{x}(n) = \mathbf{B}\mathbf{s}(n), \quad (2.22)$$



(a) General flow



(b) Example of brain

Figure 2.5: Blind Source Separation

where

$$\mathbf{x}(n) = \begin{pmatrix} x_1(n) \\ \vdots \\ x_M(n) \end{pmatrix}, \quad \mathbf{s}(n) = \begin{pmatrix} s_1(n) \\ \vdots \\ s_R(n) \end{pmatrix}. \quad (2.23)$$

The goal of BSS is to estimate  $\mathbf{B}$  and  $\mathbf{S}$  from  $\mathbf{X}$ . There are many methods for BSS such as the independent component analysis (ICA) [50], the sparse component analysis, and the nonnegative matrix factorization (NMF). In this section, I introduce ICA as the most typical approach for BSS, and the other methods are introduced in Chapter 3 since they are based on the matrix decomposition approach.

### 2.4.1 Independent Component Analysis

First, I introduce the hypothesis of ICA as follows:

Hypothesis of ICA

1.  $\{s_r\}$  are statistically independent of each other,

$$p(s_1, s_2, \dots, s_R) = p(s_1)p(s_2) \cdots p(s_R). \quad (2.24)$$

2.  $\{s_r\}$  follow non-Gaussian distributions.

If more than one of  $\{s_r\}$  follow the Gaussian distribution, then ICA is impossible.

3.  $\mathbf{B} \in \mathbb{R}^{M \times R}$  is a regular matrix (i.e.,  $M = R$ ).

Therefore, we can rewrite the model as

$$\mathbf{s}(n) = \mathbf{A}^T \mathbf{x}(n), \quad (2.25)$$

where  $\mathbf{A}^T = \mathbf{B}^{-1}$ . It is only necessary to estimate  $\mathbf{A}$ .

Before explaining the detail of ICA, it is helpful to show the overall procedures of ICA. ICA consists of two phases, mainly: the first phase is a whitening, and the second phase is a maximization of some non-Gaussianity measure. The reason why whitening is necessary is that the solution of  $\mathbf{A}$  can be limited to an orthogonal matrix after whitening. The non-Gaussianity measure is very important for ICA because it expresses the statistical independence of  $s_r$ . Figure 2.6 shows an example of original independent signals and observed signals with  $M = R = 2$ , and its whitening and ICA results. The blue lines of Fig. 2.6(d) show the directions of new axes, thus original sources are almost completely estimated excluding the scale and the order of them.

#### Whitening

White signals are defined as any  $\mathbf{z}$  which satisfy the following conditions

$$\mathbf{E}[\mathbf{z}] = \mathbf{0}, \quad \mathbf{E}[\mathbf{z}\mathbf{z}^T] = \mathbf{I}. \quad (2.26)$$

First, we modify  $\mathbf{X}$  as the zero average signal by

$$\mathbf{X}_0 = \mathbf{X} \left( \mathbf{I}_N - \frac{1}{N} \mathbf{1}_{N \times N} \right), \quad (2.27)$$

where  $\mathbf{1}_{N \times N}$  is an  $(N \times N)$ -matrix of which all elements are 1. We denote the zero average signal as  $\mathbf{x}_0(n)$ . Second, the whitening result is given by  $\mathbf{Z} = \mathbf{U}\mathbf{D}^{1/2}\mathbf{U}^T\mathbf{X}_0 \in \mathbb{R}^{R \times N}$  of eigenvalue decomposition of  $[\frac{1}{N}\mathbf{X}_0\mathbf{X}_0^T] = \mathbf{U}\mathbf{D}\mathbf{U}^T$ . The whitening results are also denoted as  $\mathbf{z}(n)$ , and they are also given by

$$\mathbf{z}(n) = \mathbf{U}\mathbf{D}^{-1/2}\mathbf{U}^T\mathbf{x}_0(n). \quad (2.28)$$

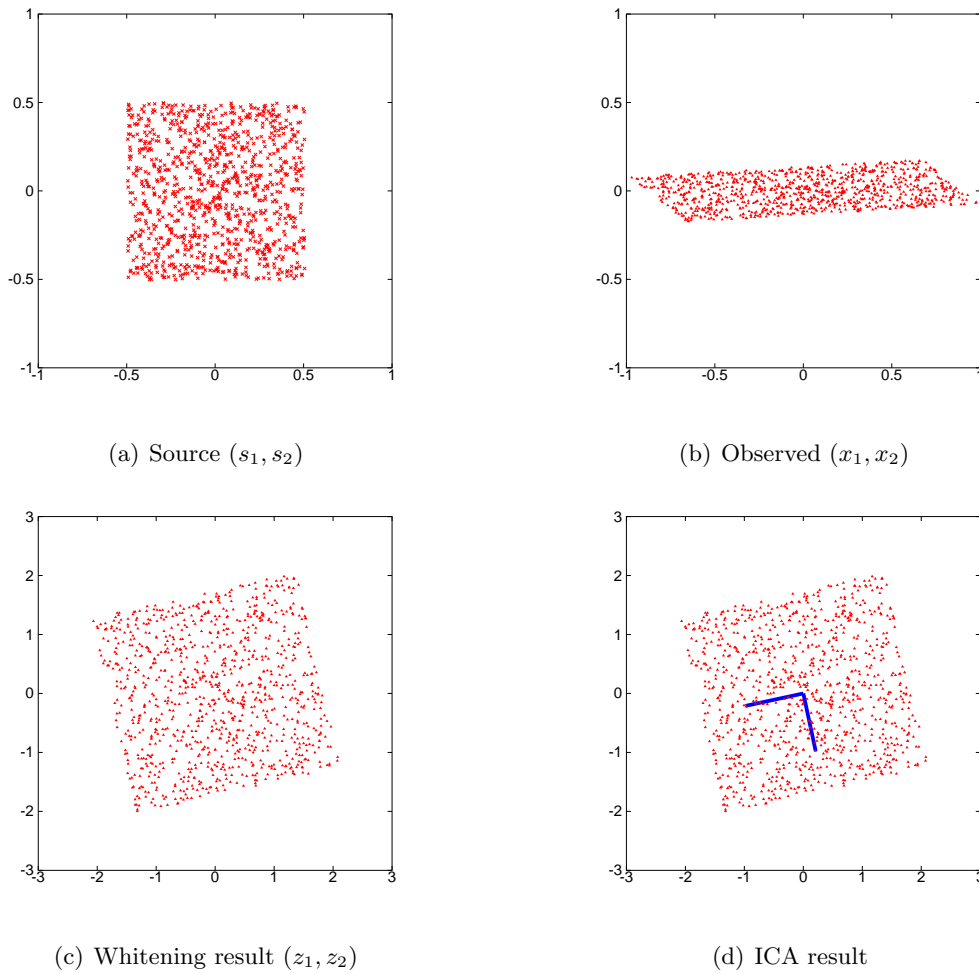


Figure 2.6: Individual phases of ICA

Thus, the model is given by

$$\mathbf{s}(n) = \mathbf{W}^T \mathbf{z}(n) = \mathbf{W}^T \mathbf{U} \mathbf{D}^{-1/2} \mathbf{U}^T \mathbf{x}_0(n) = \mathbf{A}^T \mathbf{x}_0(n), \quad (2.29)$$

where  $\mathbf{W} \in \mathbb{R}^{R \times R}$  is an orthogonal transform matrix. Therefore, it is only necessary to estimate  $\mathbf{W}$  after whitening.

### Non-Gaussianity

In order to separate the observed signals  $\mathbf{x}(n)$  into statistically independent sources  $\mathbf{s}(n)$ , some measure of statistical independence is necessary. Non-Gaussianity can be regarded as a measure of statistical independence according to the central limit theorem (CLT). From the CLT, we can know that a distribution which consists of several independent and identically distributed approaches a Gaussian distribution. Thus, if we assume that  $\{\mathbf{s}_r\}$  are identically distributed, then  $\mathbf{x}_m = \sum_{r=1}^R b_{mr} s_r$  are closer to Gaussian distributions than  $\{\mathbf{s}_r\}$ . In other words, the non-Gaussianity of  $\{\mathbf{s}_r\}$  are higher than  $\{\mathbf{x}_m\}$ .

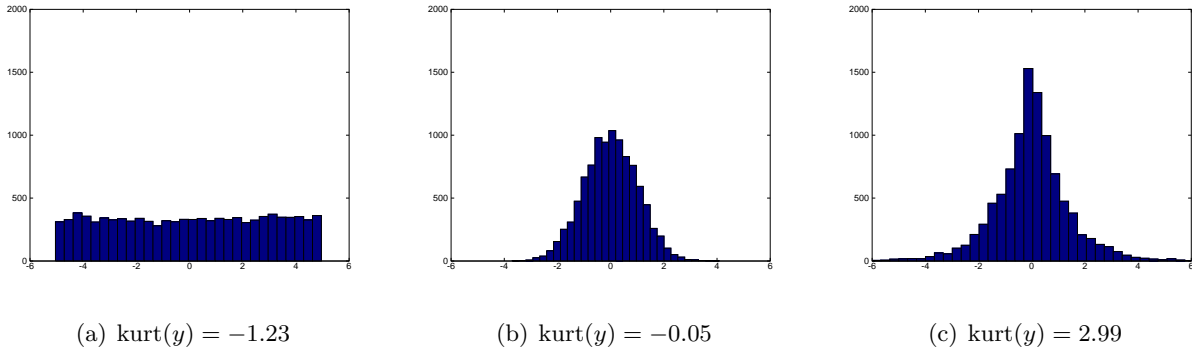


Figure 2.7: Kurtosis

Therefore, if some non-Gaussianity measure is given, we can estimate the parameter vector  $\mathbf{a}_r$  by maximizing the non-Gaussianity of  $\mathbf{a}_r^T \mathbf{x}(n)$ . The criterion is given by

$$\underset{\mathbf{a}_r}{\text{maximize}} \quad \text{Non-Gaussianity}(\mathbf{a}_r^T \mathbf{x}(n)). \quad (2.30)$$

However, we can estimate only the first parameter vector from this criterion. In practice, we can limit the parameter vectors in orthogonal vectors after whitening, thus the criterion is given by

$$\begin{aligned} \underset{\mathbf{W}}{\text{maximize}} \quad & \sum_{r=1}^R \text{Non-Gaussianity}(\mathbf{w}_r^T \mathbf{z}(n)), \\ \text{subject to} \quad & \mathbf{W}^T \mathbf{W} = \mathbf{I}_R. \end{aligned} \quad (2.31)$$

There are several measures of non-Gaussianity such as Kurtosis and Neg-entropy. I introduce the kurtosis based ICA in Section 2.4.2 and the neg-entropy based ICA in Section 2.4.3.

### 2.4.2 Kurtosis based ICA

In this section, I introduce an algorithm for ICA based on Kurtosis. When the average of  $y$  is zero, the Kurtosis can be given by

$$\text{kurt}(y) = \text{E}[y^4] - 3(\text{E}[y^2])^2. \quad (2.32)$$

Furthermore, we assume that  $y$  is white (i.e.  $\text{E}[y] = 0$ ,  $\text{E}[y^2] = 1$ ), then it is given by

$$\text{kurt}(y) = \text{E}[y^4] - 3. \quad (2.33)$$

When  $y$  is distributed normally,  $\text{kurt}(y)$  becomes zero. On the other hand,  $y$  follows a non-Gaussian distribution, the kurtosis becomes some positive or negative value. Thus, the kurtosis can be used as a non-Gaussianity measure for ICA. When the  $\text{kurt}(y)$  is positive, its distribution is called the super-gaussian, and when the  $\text{kurt}(y)$  is negative, its distribution is called the sub-gaussian. Figure 2.7 show examples of histograms and their kurtoses.

The kurtosis based ICA criterion is given by

$$\begin{aligned} & \text{maximize} && \sum_{r=1}^R |\text{kurt}(\mathbf{w}_r^T \mathbf{z}(n))|, \\ & \text{subject to} && \mathbf{W}^T \mathbf{W} = \mathbf{I}_R. \end{aligned} \quad (2.34)$$

I denote  $\mathbf{z}(n)$  as  $\mathbf{z}$  for simplicity, and firstly focus to estimate one parameter vector  $\mathbf{w}_r$  as

$$\text{maximize } |\text{kurt}(\mathbf{w}_r^T \mathbf{z})|, \text{ s.t. } \mathbf{w}_r^T \mathbf{w}_r = 1. \quad (2.35)$$

Differential of  $|\text{kurt}(\mathbf{w}_r^T \mathbf{z})|$  is given by

$$\begin{aligned} \frac{\partial |\text{kurt}(\mathbf{w}_r^T \mathbf{z})|}{\partial \mathbf{w}_r} &= \frac{\partial}{\partial \mathbf{w}_r} |\text{E}\{(\mathbf{w}_r^T \mathbf{z})^4\} - 3\text{E}\{(\mathbf{w}_r^T \mathbf{z})^2\}^2| \\ &= \frac{\partial}{\partial \mathbf{w}_r} |\text{E}\{(\mathbf{w}_r^T \mathbf{z})^4\} - 3\{|\mathbf{w}_r|^2\}^2| \quad (\text{because } \text{E}(\mathbf{z}\mathbf{z}^T) = \mathbf{I}) \\ &= 4\text{sign}[\text{kurt}(\mathbf{w}_r^T \mathbf{z})] [\text{E}\{\mathbf{z}(\mathbf{w}_r^T \mathbf{z})^3\} - 3\mathbf{w}_r \|\mathbf{w}_r\|^2] \end{aligned} \quad (2.36)$$

According to the gradient method, we can obtain following update rule:

$$\mathbf{w}_r \leftarrow \mathbf{w}_r + \Delta \mathbf{w}_r; \quad (2.37)$$

$$\mathbf{w}_r \leftarrow \frac{\mathbf{w}_r}{\|\mathbf{w}_r\|}; \quad (2.38)$$

where  $\Delta \mathbf{w}_r \propto \text{sign}[\text{kurt}(\mathbf{w}_r^T \mathbf{z})] [\text{E}\{\mathbf{z}(\mathbf{w}_r^T \mathbf{z})^3\} - 3\mathbf{w}_r]$ . Its algorithm converges when  $\mathbf{w}_r \propto \Delta \mathbf{w}_r$ , thus we can obtain another update rule as

$$\mathbf{w}_r \leftarrow \text{E}\{\mathbf{z}(\mathbf{w}_r^T \mathbf{z})^3\} - 3\mathbf{w}_r, \quad (2.39)$$

$$\mathbf{w}_r \leftarrow \frac{\mathbf{w}_r}{\|\mathbf{w}_r\|}. \quad (2.40)$$

It is well known as a fast fixed point algorithm for ICA.

Next, I introduce the algorithm to estimate orthogonal parameter vectors  $\mathbf{W}$ . This algorithm consists of two steps: the first step is to update  $\{\mathbf{w}_r\}$  individually, the second step is to project  $\mathbf{W}$  onto the orthogonal matrix. There are several orthogonalization methods, and they can be separated into the sequential approach and the symmetrical approach. The Gram-Schmidt orthogonalization is a very typical method for the sequential approach. It is given by

$$\mathbf{w}_r \leftarrow \mathbf{w}_r - \sum_{k=1}^{r-1} (\mathbf{w}_r^T \mathbf{w}_k) \mathbf{w}_k, \quad (2.41)$$

for  $r = 1, \dots, R$ . On the other hand, a symmetrical orthogonalization is given by

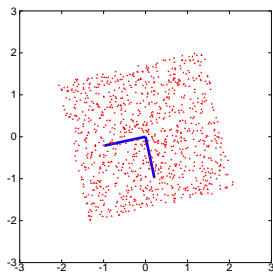
$$\mathbf{W} \leftarrow (\mathbf{W}\mathbf{W}^T)^{-1/2} \mathbf{W}. \quad (2.42)$$

When we want to treat individual vector fairly, the symmetrical approach is useful.

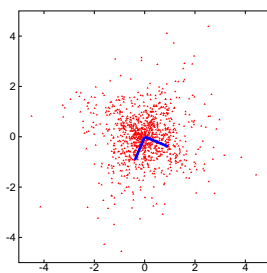
Finally, the kurtosis based ICA algorithm can be summarized as Algorithm 2.

**Algorithm 2** Kurtosis based ICA algorithm

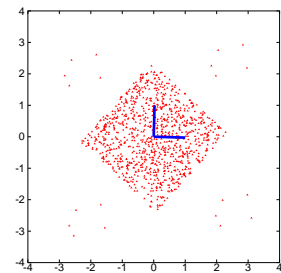
- 
- 1: **Input:**  $\mathbf{X} \in \mathbb{R}^{R \times N}$
  - 2: **Initialize:** orthogonal matrix  $\mathbf{W} \in \mathbb{R}^{R \times R}$
  - 3:  $\mathbf{X} \leftarrow \mathbf{X}(\mathbf{I}_N - \frac{1}{N}\mathbf{1}_{N \times N})$  ;
  - 4:  $\mathbf{Z} \leftarrow \mathbf{U}\mathbf{D}^{-1/2}\mathbf{U}^T\mathbf{X}$  with eigenvector  $\mathbf{U}$  and eigenvalue  $\mathbf{D}$  of  $\frac{1}{N}\mathbf{X}\mathbf{X}^T$  ;
  - 5: **repeat**
  - 6:    $\mathbf{w}_r \leftarrow \mathbb{E}\{\mathbf{z}(\mathbf{w}_r^T\mathbf{z})^3\} - 3\mathbf{w}_r$  for  $r = 1, \dots, R$  ;
  - 7:    $\mathbf{w}_r \leftarrow \frac{\mathbf{w}_r}{\|\mathbf{w}_r\|}$  for  $r = 1, \dots, R$  ;
  - 8:   Orthogonalize  $\mathbf{W}$  by (2.41) or (2.42) ;
  - 9: **until** convergence
  - 10:  $\mathbf{S} \leftarrow \mathbf{W}^T\mathbf{Z}$  ;
  - 11: **Output:**  $\mathbf{S}$
- 



(a) subgaussian



(b) supergaussian



(c) subgaussian with outliers (20:1000)

Figure 2.8: Example of kurtosis based ICA

**Example**

Figure 2.8 shows the results of the kurtosis based ICA. Subgaussian distribution is generated by the uniform distribution, and supergaussian distribution is generated by the Laplacian distribution. In both distributions, the kurtosis based ICA worked well; however, it did not work with outliers. Actually, the kurtosis is a very sensitive measure since it is the fourth order function. The rate of outliers in Figure 2.8(c) is only 2% ; however it is enough to show the invalidity of the criterion.

**2.4.3 Neg-entropy based ICA**

The neg-entropy is often used for ICA since it is a robust measure of non-Gaussianity for outliers. Strictly, the approximation of neg-entropy is often used because it is more robust for outliers.

The neg-entropy is defined by

$$J(y) = H(y_{\text{Gauss}}) - H(y), \quad (2.43)$$

where  $H(y) = -\int p_y(\eta) \log p_y(\eta) d\eta$  is a entropy and  $y_{\text{Gauss}}$  follows a Gaussian distribution with

$\mu = E[y]$  and  $\sigma^2 = E[(y - \mu)^2]$ . Thus, if  $y$  follows Gaussian distribution, then  $J(y) = 0$ , and if not, then  $J(y) > 0$ .

In this method, we approximate the neg-entropy by

$$J(y) \approx k_1(E\{G_1(y)\})^2 + k_2(E\{G_2(y)\} - E\{G_2(\nu)\})^2, \quad (2.44)$$

where  $G_1(\cdot)$  is a non-quadratically odd function,  $G_2(\cdot)$  is a non-quadratically even function,  $k_1$  and  $k_2$  are positive constants,  $\nu$  follows the standard Gaussian distribution, and we assume that  $y$  is white. This approximation is not strict, but it has the same properties of the neg-entropy. For example, this measure takes zero when  $y$  follows the Gaussian distribution, and positive when  $y$  does not follow the Gaussian distribution. Furthermore, this approximation can be robust measure by selecting appropriate functions for  $G_1(\cdot)$  and  $G_2(\cdot)$ . The following functions are well known as a useful selection:

$$G_1(y) = \frac{1}{a_1} \log \cosh a_1 y, \quad (2.45)$$

$$G_2(y) = -\exp(-y^2/2), \quad (2.46)$$

where  $1 \leq a_1 \leq 2$  (e.g.,  $a_1 = 1$  is often used).

If we use a single non-quadratic function, this approximation can be given by

$$J(y) \propto [E\{G(y)\} - E\{G(\nu)\}]^2. \quad (2.47)$$

This approximation has also the same properties of the neg-entropy, and it is simpler than the two functions version. Thus, we introduce the neg-entropy based ICA algorithm from this approximation. The objective function can be transformed as

$$\begin{aligned} J(y) &\propto [E\{G(y)\} - E\{G(\nu)\}]^2 \\ &\propto [E\{G(y)\}]^2 - 2E\{G(y)\}E\{G(\nu)\} + [E\{G(\nu)\}]^2. \end{aligned} \quad (2.48)$$

Furthermore, its differential function by  $\mathbf{w}$  can be given by

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{w}} &\propto 2E\{G(\mathbf{w}^T \mathbf{z})\}E\{\mathbf{z}g(\mathbf{w}^T \mathbf{z})\} - 2E\{G(\nu)\}E\{\mathbf{z}g(\mathbf{w}^T \mathbf{z})\} \\ &\propto [E\{G(\mathbf{w}^T \mathbf{z})\} - E\{G(\nu)\}]E\{\mathbf{z}g(\mathbf{w}^T \mathbf{z})\}, \end{aligned} \quad (2.49)$$

where  $g(\cdot)$  is the derivative of  $G(\cdot)$ . Thus, we can obtain the following updating rule based on a fixed point algorithm:

$$\mathbf{w} \leftarrow E\{\mathbf{z}g(\mathbf{w}^T \mathbf{z})\}, \quad (2.50)$$

$$\mathbf{w} \leftarrow \mathbf{w}/\|\mathbf{w}\|, \quad (2.51)$$

where we can ignore  $[E\{G(\mathbf{w}^T \mathbf{z})\} - E\{G(\nu)\}]$  since  $\mathbf{w}$  is normalized later. However, this algorithm do not converge well, because its approximation does not have a algebraic property that a true cumulant

**Algorithm 3** Neg-entropy based ICA algorithm

- 
- 1: **Input:**  $\mathbf{X} \in \mathbb{R}^{R \times N}$
  - 2: **Initialize:** orthogonal matrix  $\mathbf{W} \in \mathbb{R}^{R \times R}$
  - 3:  $\mathbf{X} \leftarrow \mathbf{X}(\mathbf{I}_N - \frac{1}{N}\mathbf{1}_{N \times N})$  ;
  - 4:  $\mathbf{Z} \leftarrow \mathbf{U}\mathbf{D}^{-1/2}\mathbf{U}^T\mathbf{X}$  with eigenvector  $\mathbf{U}$  and eigenvalue  $\mathbf{D}$  of  $\frac{1}{N}\mathbf{X}\mathbf{X}^T$  ;
  - 5: **repeat**
  - 6:    $\mathbf{w}_r \leftarrow \mathbf{E}\{\mathbf{z}g(\mathbf{w}_r^T\mathbf{z})\} - \mathbf{E}\{g'(\mathbf{w}_r^T\mathbf{z})\}\mathbf{w}_r$  for  $r = 1, \dots, R$  ;
  - 7:    $\mathbf{w}_r \leftarrow \frac{\mathbf{w}_r}{\|\mathbf{w}_r\|}$  for  $r = 1, \dots, R$  ;
  - 8:   Orthogonalize  $\mathbf{W}$  by (2.41) or (2.42) ;
  - 9: **until** convergence
  - 10:  $\mathbf{S} \leftarrow \mathbf{W}^T\mathbf{Z}$  ;
  - 11: **Output:**  $\mathbf{S}$
- 

has. Thus, we provide an improvement method for it by using a approximated newton's method. When we use (2.50) and (2.51), the convergence condition is given by

$$F(\mathbf{w}) = \mathbf{E}\{\mathbf{z}g(\mathbf{w}^T\mathbf{z})\} + \beta\mathbf{w} = \mathbf{0}, \quad (2.52)$$

where  $\beta$  is a parameter. We apply the approximated newton's method to the condition (2.52). Thus, the gradient of  $F(\mathbf{w})$  is given by

$$\frac{\partial F}{\partial \mathbf{w}} = \mathbf{E}\{\mathbf{z}\mathbf{z}^T g'(\mathbf{w}^T\mathbf{z})\} + \beta\mathbf{I}_R, \quad (2.53)$$

where  $g'(\cdot)$  is the derivative of  $g(\cdot)$ . Since  $\mathbf{z}$  is white, we apply the following approximation,

$$\mathbf{E}\{\mathbf{z}\mathbf{z}^T g'(\mathbf{w}^T\mathbf{z})\} \approx \mathbf{E}\{\mathbf{z}\mathbf{z}^T\}\mathbf{E}\{g'(\mathbf{w}^T\mathbf{z})\} = \mathbf{E}\{g'(\mathbf{w}^T\mathbf{z})\}\mathbf{I}_R. \quad (2.54)$$

Then,  $\frac{\partial F}{\partial \mathbf{w}}$  becomes a diagonal matrix and its inverse matrix can be obtained easily. We can derive the following approximated newton's method

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \frac{[\mathbf{E}\{\mathbf{z}g(\mathbf{w}^T\mathbf{z})\} + \beta\mathbf{w}]}{[\mathbf{E}\{g'(\mathbf{w}^T\mathbf{z})\} + \beta]} \propto \mathbf{w}[\mathbf{E}\{g'(\mathbf{w}^T\mathbf{z})\} + \beta] - [\mathbf{E}\{\mathbf{z}g(\mathbf{w}^T\mathbf{z})\} + \beta\mathbf{w}] \\ &= \mathbf{w}\mathbf{E}\{g'(\mathbf{w}^T\mathbf{z})\} - \mathbf{E}\{\mathbf{z}g(\mathbf{w}^T\mathbf{z})\} \end{aligned} \quad (2.55)$$

Finally, the fast fixed point update rule for neg-entropy based ICA is given by

$$\mathbf{w} \leftarrow \mathbf{E}\{\mathbf{z}g(\mathbf{w}^T\mathbf{z})\} - \mathbf{E}\{g'(\mathbf{w}^T\mathbf{z})\}\mathbf{w}; \quad (2.56)$$

$$\mathbf{w} \leftarrow \mathbf{w}/\|\mathbf{w}\|; \quad (2.57)$$

and its algorithm can be summarized as Algorithm 3.  $g(\cdot)$  and  $g'(\cdot)$  in the Table 2.1 are often used. Please note that  $(g_3, g'_3)$  is equivalent to Kurtosis based ICA.

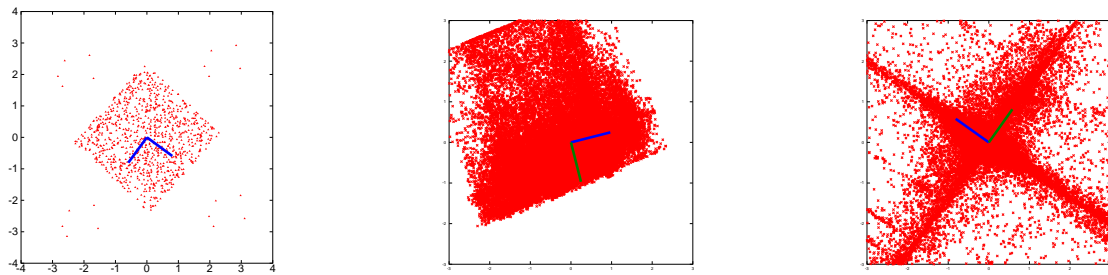
Table 2.1: Approximation functions for neg-entropy based ICA

label	$g(y)$	$g'(y)$
1	$g_1(y) = \tanh(a_1 y)$	$g'_1(y) = a_1(1 - \tanh^2(a_1 y))$
2	$g_2(y) = y \exp(-y^2/2)$	$g'_2(y) = (1 - y^2) \exp(-y^2/2)$
3	$g_3(y) = y^3$	$g'_3(y) = 3y^2$

**Example**

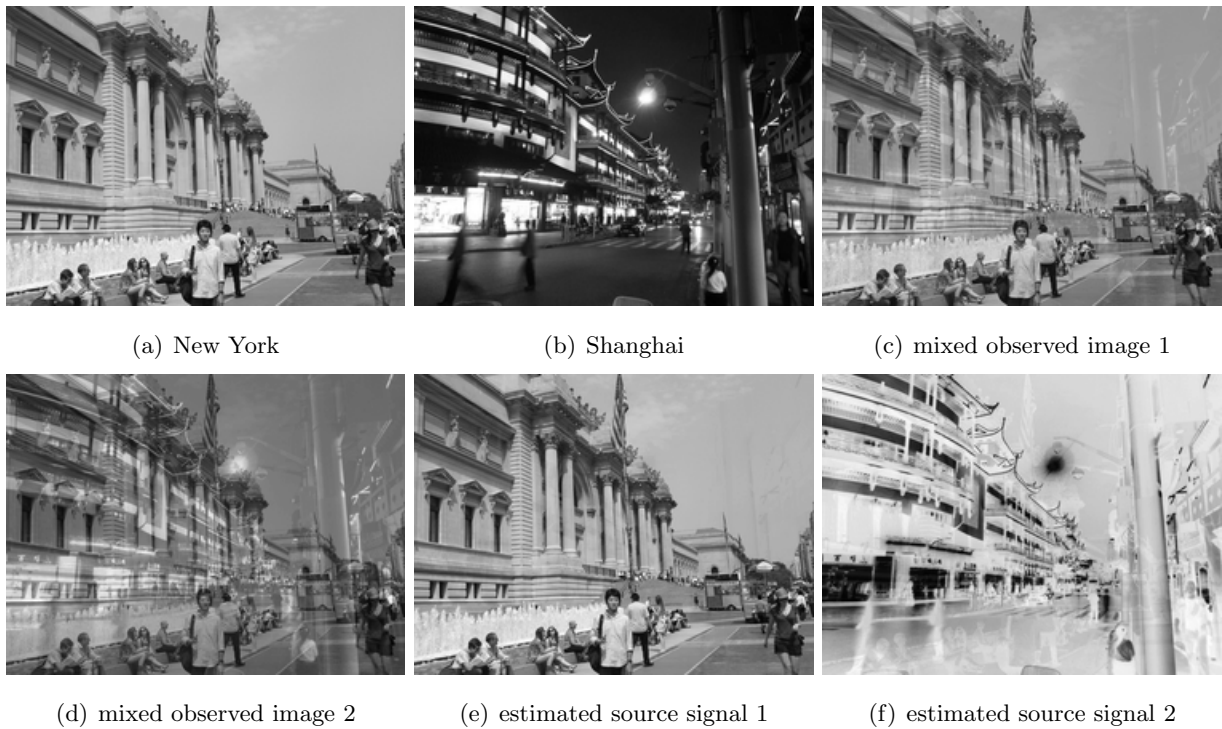
Figure 2.9 shows the results of the neg-entropy based ICA (using  $g_1$ ). We can see that the neg-entropy based ICA is robust for outliers and works well for real data.

Fig. 2.10 shows an example of (a,b) original images, (c,d) mixed observed images, and (e,f) estimated sources. These images were separated clearly; however, black and white are reversed in the estimated signal 2. In this way, ICA can not estimate either the correct or not for the opposite sign.



(a) subgaussian with outliers (20:1000)      (b) real data : subgaussian      (c) real data : supergaussian

Figure 2.9: Neg-entropy based ICA (using  $g_1$ )



(d) mixed observed image 2      (e) estimated source signal 1      (f) estimated source signal 2

Figure 2.10: Example of real image separation

## Chapter 3

# Matrix/Tensor based Feature Extraction

In this Chapter, I introduce matrix and tensor based feature extraction methods and point out their problems. Key-points of this chapter are decomposition models and imposed constraints. The decomposition models are separated into the matrix decomposition model and the tensor decomposition models. There are several constraints such as the orthogonality, the sparsity, the smoothness, and the nonnegativity. Properties of the method depend on the combination of the model and the constraint, and the models and the constraints are selected based on the situation demands.

### 3.1 Matrix decomposition model

In this section, I treat the feature extraction methods to analyze a matrix data. As an assumption, we have  $J$  observed signals  $\{\mathbf{y}_j \in \mathbb{R}^I\}_{j=1}^J$  which can be given as a matrix  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_J] \in \mathbb{R}^{I \times J}$ . The goal of feature extraction is to extract the latent important features included in the data. Matrix decomposition model is given by

$$\mathbf{Y} \simeq \mathbf{A}\mathbf{X}, \quad (3.1)$$

where  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_R] \in \mathbb{R}^{I \times R}$  and  $\mathbf{X} = [x_{rj}] \in \mathbb{R}^{R \times J}$ . This model can be also expressed as

$$\mathbf{y}_j \simeq \sum_{r=1}^R x_{rj} \mathbf{a}_r, \quad (3.2)$$

for  $j = 1, \dots, J$ . Eq. (3.2) means that the matrix decomposition approximates  $\mathbf{y}_j$  by a linear combination of bases  $[\mathbf{a}_r]$ . On the other hand, let  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_R] := \mathbf{X}^T \in \mathbb{R}^{J \times R}$ , then the model can be given by

$$\mathbf{Y} \simeq \mathbf{A}\mathbf{B}^T = \sum_{r=1}^R \mathbf{a}_r \mathbf{b}_r^T. \quad (3.3)$$

Thus, the matrix decomposition approximates  $\mathbf{Y}$  by a linear combination of one-rank matrices  $[\mathbf{a}_r \mathbf{b}_r^T]$ . The matrix decomposition model has a high-potential to apply to various real world problems such as clustering task, parts learning, denoising, and blind source separation.

When we try to extract the feature, it is important to make a criterion in concert with the target. For this reason, orthogonality, sparsity, smoothness, and nonnegativity are important factors to obtain such features. Actually, almost all methods consist of the minimization of some distance measure and some additional constraints. For example, a singular value decomposition (SVD) consists of the minimization of the Frobenius norm and the orthogonality constraint. In this way, we can say that the properties of methods are decided by the constraints.

### 3.1.1 Singular Value Decomposition

Singular value decomposition (SVD) is a basic matrix decomposition method. Its model is given by

$$\mathbf{Y} = \mathbf{U} \mathbf{D} \mathbf{V}^T, \quad (3.4)$$

where  $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_I] \in \mathbb{R}^{I \times I}$ ,  $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_J] \in \mathbb{R}^{J \times J}$  are unitary matrices, and  $\mathbf{D} = [d_{ij}] \in \mathbb{R}^{I \times J}$  is a diagonal matrix that its diagonal elements are nonnegative values. Each  $\mathbf{u}_r$ ,  $\mathbf{v}_r$ , and  $d_{rr}$  are called respectively a left-singular vector, a right-singular vector, and a singular value, thus they satisfy

$$\mathbf{u}_r^T \mathbf{Y} = \mathbf{v}_r^T d_{rr}, \quad (3.5)$$

$$d_{rr} \mathbf{u}_r = \mathbf{Y} \mathbf{v}_r, \quad (3.6)$$

for  $r = 1, \dots, \min(I, J)$ . Please note that the SVD model (3.4) is defined as an equation since this model can completely provide the original data matrix  $\mathbf{Y}$ . The criterion of SVD is given by

$$\begin{aligned} & \text{minimize } \|\mathbf{Y} - \mathbf{U} \mathbf{D} \mathbf{V}^T\|_F^2, \\ & \text{subject to } \mathbf{U}^T \mathbf{U} = \mathbf{U} \mathbf{U}^T = \mathbf{I}_I, \mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I}_J, \end{aligned} \quad (3.7)$$

where  $\mathbf{I}_I \in \mathbb{R}^{I \times I}$  and  $\mathbf{I}_J \in \mathbb{R}^{J \times J}$  are identity matrices. The objective function can be transformed as

$$\begin{aligned} \|\mathbf{Y} - \mathbf{U} \mathbf{D} \mathbf{V}^T\|_F^2 &= \text{tr}(\mathbf{Y}^T \mathbf{Y}) - 2\text{tr}(\mathbf{Y}^T \mathbf{U} \mathbf{D} \mathbf{V}^T) + \text{tr}(\mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{U} \mathbf{D} \mathbf{V}^T) \\ &= \text{tr}(\mathbf{Y}^T \mathbf{Y}) - 2\text{tr}(\mathbf{Y}^T \mathbf{U} \mathbf{D} \mathbf{V}^T) + \text{tr}(\mathbf{D}^T \mathbf{D}). \end{aligned} \quad (3.8)$$

Its Lagrange's condition with respect to  $\mathbf{D}$  is given by

$$\frac{\partial}{\partial \mathbf{D}} \|\mathbf{Y} - \mathbf{U} \mathbf{D} \mathbf{V}^T\|_F^2 = -2\mathbf{U}^T \mathbf{Y} \mathbf{V} + 2\mathbf{D} = \mathbf{0}. \quad (3.9)$$

Thus, the solution of  $\mathbf{D}$  can be given by

$$\mathbf{D} = \mathbf{U}^T \mathbf{Y} \mathbf{V}. \quad (3.10)$$

We can see that eq. (3.10) satisfies the conditions (3.5) and (3.6). By using (3.8) and (3.10), the Lagrange's function can be given by

$$\begin{aligned} L(\mathbf{U}, \mathbf{V}, \mathbf{\Lambda}_1, \mathbf{\Lambda}_2) = & -\text{tr}(\mathbf{U}^T \mathbf{Y} \mathbf{V} \mathbf{V}^T \mathbf{Y}^T \mathbf{U}) \\ & + \text{tr}[(\mathbf{U}^T \mathbf{U} - \mathbf{I}_I) \mathbf{\Lambda}_1] + \text{tr}[(\mathbf{V}^T \mathbf{V} - \mathbf{I}_J) \mathbf{\Lambda}_2], \end{aligned} \quad (3.11)$$

where  $\mathbf{\Lambda}_1$  and  $\mathbf{\Lambda}_2$  are Lagrange's coefficients. Its Lagrange's conditions with respect to  $\mathbf{U}$  and  $\mathbf{V}$  are given by

$$\frac{\partial L}{\partial \mathbf{U}} = -2(\mathbf{Y} \mathbf{Y}^T) \mathbf{U} + 2\mathbf{U} \mathbf{\Lambda}_1 = 0, \quad (3.12)$$

$$\frac{\partial L}{\partial \mathbf{V}} = -2(\mathbf{Y}^T \mathbf{Y}) \mathbf{V} + 2\mathbf{V} \mathbf{\Lambda}_2 = 0. \quad (3.13)$$

Eq. (3.12) is an eigenvalue problem, thus the solution of  $\mathbf{U}$  can be given as a set of eigenvectors of  $(\mathbf{Y} \mathbf{Y}^T)$ . Let  $[\phi_1, \phi_2, \dots, \phi_I]$  and  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_I$  be respectively eigenvectors and eigenvalues of  $(\mathbf{Y} \mathbf{Y}^T)$ , then  $\mathbf{U}$  is given by

$$\mathbf{U} = [\phi_1, \phi_2, \dots, \phi_I]. \quad (3.14)$$

In the same way, the solution of  $\mathbf{V}$  can be given as a set of eigenvectors of  $(\mathbf{Y}^T \mathbf{Y})$ . Let  $[\psi_1, \psi_2, \dots, \psi_J]$  and  $\eta_1 \geq \eta_2 \geq \dots \geq \eta_J$  be respectively eigenvectors and their corresponding eigenvalues of  $(\mathbf{Y}^T \mathbf{Y})$ , then  $\mathbf{V}$  is given by

$$\mathbf{V} = [\psi_1, \psi_2, \dots, \psi_J]. \quad (3.15)$$

Let  $R = \min(I, J)$ ,  $D_{rr} = \sqrt{\lambda_r} = \sqrt{\eta_r}$  for  $r = 1, \dots, R$ .

### Truncated SVD

Truncated singular value decomposition (tSVD) is a method to obtain a low-rank approximation of original data matrix  $\mathbf{Y}$ . Its model is given by

$$\mathbf{Y} \simeq \mathbf{U} \mathbf{D} \mathbf{V}^T, \quad (3.16)$$

where  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_R] \in \mathbb{R}^{I \times R}$ ,  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_R] \in \mathbb{R}^{J \times R}$  are orthonormal matrices, and  $\mathbf{D} \in \mathbb{R}^{R \times R}$  is a diagonal matrix which its diagonal elements are singular values. In general,  $R \leq \text{rank}(\mathbf{Y}) \leq \min(I, J)$ .

The criterion of tSVD is given by

$$\text{minimize } \|\mathbf{Y} - \mathbf{U} \mathbf{D} \mathbf{V}^T\|_F^2, \quad (3.17)$$

$$\text{subject to } \mathbf{U}^T \mathbf{U} = \mathbf{I}_R, \mathbf{V}^T \mathbf{V} = \mathbf{I}_R, \quad (3.18)$$

where  $\mathbf{I}_R \in \mathbb{R}^{R \times R}$  is an identity matrix. As its name suggests, the solution of tSVD is given by truncating the solution of SVD. Let  $\mathbf{U}_0 = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_I]$ ,  $\mathbf{V}_0 = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_J]$ , and  $\mathbf{D}_0 = [d_{ij}]$  be the

solution of SVD, then each factor matrices are given by

$$\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_R], \quad (3.19)$$

$$\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_R], \quad (3.20)$$

$$\mathbf{D} = [\mathbf{D}_0]_{1:R,1:R}. \quad (3.21)$$

The tSVD gives an  $R$ -rank approximation of  $\mathbf{Y}$  which minimizes the mean squared error.

### Deflation based algorithm for tSVD

The above mentioned algorithm for tSVD is consequential and wasteful since it must obtain a full-rank decomposition. In this section, I introduce a deflation based algorithm to obtain the solution of tSVD, directly. First we consider one-rank approximation criterion as follow:

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|\mathbf{Y} - d\mathbf{u}\mathbf{v}^T\|_F^2, \\ & \text{subject to} \quad \|\mathbf{u}\|^2 = 1, \|\mathbf{v}\|^2 = 1. \end{aligned} \quad (3.22)$$

The objective function of (3.22) can be transformed as

$$\frac{1}{2} \|\mathbf{Y} - d\mathbf{u}\mathbf{v}^T\|_F^2 = \frac{1}{2} \text{tr}(\mathbf{Y}\mathbf{Y}^T) - d\mathbf{u}^T\mathbf{Y}\mathbf{v} + \frac{1}{2}d^2. \quad (3.23)$$

Focusing on this objective function with respect to only  $\mathbf{u}$  and  $\mathbf{v}$ , a criterion (3.22) can be transformed to

$$\begin{aligned} & \text{maximize} \quad \mathbf{u}^T\mathbf{Y}\mathbf{v}, \\ & \text{subject to} \quad \|\mathbf{u}\|^2 = 1, \|\mathbf{v}\|^2 = 1. \end{aligned} \quad (3.24)$$

We solve this criterion by updating  $\mathbf{u}$  and  $\mathbf{v}$ , alternately and iteratively until convergence. Update rules for  $\mathbf{u}$  and  $\mathbf{v}$  are derived by the following criteria:

$$\mathbf{u} \leftarrow \underset{\mathbf{u}}{\text{argmax}} \quad \mathbf{u}^T\mathbf{Y}\mathbf{v}, \text{ s.t. } \|\mathbf{u}\| = 1, \quad (3.25)$$

$$\mathbf{v} \leftarrow \underset{\mathbf{v}}{\text{argmax}} \quad \mathbf{u}^T\mathbf{Y}\mathbf{v}, \text{ s.t. } \|\mathbf{v}\| = 1. \quad (3.26)$$

Thus, the update rule for  $\mathbf{u}$  is given by

$$\mathbf{u} \leftarrow \mathbf{Y}\mathbf{v}, \quad (3.27)$$

$$\mathbf{u} \leftarrow \mathbf{u}/\|\mathbf{u}\|. \quad (3.28)$$

In a similar way, the update rule for  $\mathbf{v}$  is given by

$$\mathbf{v} \leftarrow \mathbf{Y}^T\mathbf{u}, \quad (3.29)$$

$$\mathbf{v} \leftarrow \mathbf{v}/\|\mathbf{v}\|. \quad (3.30)$$

---

**Algorithm 4** One-rank tSVD by power method

---

```

1: Input:  $\mathbf{Y}$ 
2: Initialize:  $\mathbf{v}$  by randomly
3: repeat
4:    $\mathbf{u} \leftarrow \mathbf{Y}\mathbf{v}$ ;
5:    $\mathbf{u} \leftarrow \mathbf{u}/\|\mathbf{u}\|$ ;
6:    $\mathbf{v} \leftarrow \mathbf{Y}^T\mathbf{u}$ ;
7:    $\mathbf{v} \leftarrow \mathbf{v}/\|\mathbf{v}\|$ ;
8: until convergence
9:  $d \leftarrow \mathbf{u}^T\mathbf{Y}\mathbf{v}$ ;
10: Output:  $d$ ,  $\mathbf{u}$ , and  $\mathbf{v}$ 

```

---



---

**Algorithm 5**  $R$ -rank tSVD

---

```

1: Input:  $\mathbf{Y}$ ,  $R$ 
2: Initialize:  $\mathbf{Z}_0 := \mathbf{Y}$ 
3: for  $r = 1, \dots, R$  do
4:   Obtain  $(d_r, \mathbf{u}_r, \mathbf{v}_r)$  by Algorithm 4 inputting  $\mathbf{Z}_{r-1}$ ;
5:    $\mathbf{Z}_r \leftarrow \mathbf{Z}_{r-1} - d_r\mathbf{u}_r\mathbf{v}_r^T$ ;
6: end for
7: Output:  $[d_r, \mathbf{u}_r, \mathbf{v}_r]_{r=1}^R$ 

```

---

Its singular value is given by  $d = \mathbf{u}^T\mathbf{Y}\mathbf{v}$ . Finally, the solution algorithm for one-rank approximation of tSVD can be summarized as Algorithm 4. This iterative algorithm is often called the power method.

Next we consider  $R$ -rank approximation of tSVD. We can construct the algorithm for  $R$ -rank version of tSVD by using one-rank approximation of tSVD. Let  $d_r$ ,  $\mathbf{u}_r$ , and  $\mathbf{v}_r$  be the  $r$ -th singular value and left- and right-singular vectors, and define the  $r$ -tSVD model  $\widehat{\mathbf{Y}}_r = \sum_{k=1}^r d_k\mathbf{u}_k\mathbf{v}_k^T$  and the deflated matrix  $\mathbf{Z}_r := \mathbf{Y} - \widehat{\mathbf{Y}}_r$ , then we have

$$\begin{aligned} \text{rank}(\mathbf{Y}) &= \text{rank}(\widehat{\mathbf{Y}}_1) + \text{rank}(\mathbf{Z}_1) = \text{rank}(\widehat{\mathbf{Y}}_2) + \text{rank}(\mathbf{Z}_2) = \dots \\ &= \text{rank}(\widehat{\mathbf{Y}}_r) + \text{rank}(\mathbf{Z}_r) = \dots = \text{rank}(\widehat{\mathbf{Y}}_R) + \text{rank}(\mathbf{Z}_R), \end{aligned} \quad (3.31)$$

since  $\mathbf{Z}_r$  and  $\widehat{\mathbf{Y}}_r$  are orthogonal (i.e.,  $\mathbf{Z}_r^T\widehat{\mathbf{Y}}_r = \mathbf{0}$  and  $\mathbf{Z}_r\widehat{\mathbf{Y}}_r^T = \mathbf{0}$ ). Thus we can say that the  $r$ -th singular value and left- and right-singular vectors of  $\mathbf{Y}$  can be given as the first singular value and left- and right-singular vectors of  $\mathbf{Z}_{r-1}$ . Therefore, the algorithm for  $R$ -rank approximation of tSVD can be summarized as Algorithm 5.

### 3.1.2 Penalized Matrix Decomposition

Penalized matrix decomposition [105] can be regarded as a penalized tSVD. It was proposed to add several constraints into the tSVD criterion and solve the criterion based on a deflation based algorithm.

**Algorithm 6** One-rank PMD by power method

- 
- 1: **Input:**  $\mathbf{Y}$ ,  $P_1$ ,  $c_1$ ,  $P_2$  and  $c_2$
  - 2: **Initialize:**  $\mathbf{v}$  by randomly
  - 3: **repeat**
  - 4:    $\mathbf{u} \leftarrow \operatorname{argmax}_{\mathbf{u}} \mathbf{u}^T \mathbf{Y} \mathbf{v}$ , s.t.  $\|\mathbf{u}\| = 1$ ,  $P_1(\mathbf{u}) \leq c_1$ ;
  - 5:    $\mathbf{v} \leftarrow \operatorname{argmax}_{\mathbf{v}} \mathbf{u}^T \mathbf{Y} \mathbf{v}$ , s.t.  $\|\mathbf{v}\| = 1$ ,  $P_2(\mathbf{v}) \leq c_2$ ;
  - 6: **until** convergence
  - 7:  $d \leftarrow \mathbf{u}^T \mathbf{Y} \mathbf{v}$ ;
  - 8: **Output:**  $d$ ,  $\mathbf{u}$ , and  $\mathbf{v}$
- 

First we consider one-rank approximation of PMD as follow:

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|\mathbf{Y} - d\mathbf{u}\mathbf{v}^T\|_F^2, \\ & \text{subject to} \quad \|\mathbf{u}\|^2 = 1, \|\mathbf{v}\|^2 = 1, P_1(\mathbf{u}) \leq c_1, P_2(\mathbf{v}) \leq c_2, \end{aligned} \quad (3.32)$$

where  $P_1(\mathbf{u})$  and  $P_2(\mathbf{v})$  are convex penalty functions. In [105], two penalty functions are proposed, mainly: l1-norm penalty and fused lasso (FL) penalty. When we choose the l1-norm penalty, then PMD provides a sparse representation. When we choose the FL penalty, then PMD provides a sparse and smooth representation. Focusing on this objective function with respect to only  $\mathbf{u}$  and  $\mathbf{v}$ , a criterion (3.32) can be transformed to

$$\begin{aligned} & \text{maximize} \quad \mathbf{u}^T \mathbf{Y} \mathbf{v}, \\ & \text{subject to} \quad \|\mathbf{u}\|^2 = 1, \|\mathbf{v}\|^2 = 1, P_1(\mathbf{u}) \leq c_1, P_2(\mathbf{v}) \leq c_2. \end{aligned} \quad (3.33)$$

In a similar way of (3.25) and (3.26), the update rules for  $\mathbf{u}$  and  $\mathbf{v}$  are derived by the following criteria:

$$\mathbf{u} \leftarrow \operatorname{argmax}_{\mathbf{u}} \mathbf{u}^T \mathbf{Y} \mathbf{v}, \text{ s.t. } \|\mathbf{u}\| = 1, P_1(\mathbf{u}) \leq c_1, \quad (3.34)$$

$$\mathbf{v} \leftarrow \operatorname{argmax}_{\mathbf{v}} \mathbf{u}^T \mathbf{Y} \mathbf{v}, \text{ s.t. } \|\mathbf{v}\| = 1, P_2(\mathbf{v}) \leq c_2. \quad (3.35)$$

Therefore, the algorithms for one- and  $R$ -rank PMD are very similar to one- and  $R$ -rank tSVD, and they are summarized as Algorithms 6 and 7.

### l1-norm penalty

In this section, we consider to choose the l1-norm penalty function as  $P_1(\mathbf{u}) := \|\mathbf{u}\|_1 = \sum_{i=1}^I |u_i|$ . Let  $\mathbf{a} = \mathbf{Y} \mathbf{v}$ , the problem (3.34) can be rewritten by

$$\mathbf{u} \leftarrow \operatorname{argmax}_{\mathbf{u}} \mathbf{u}^T \mathbf{a}, \text{ s.t. } \|\mathbf{u}\| = 1, \|\mathbf{u}\|_1 \leq c_1. \quad (3.36)$$

If  $\left\| \frac{\mathbf{a}}{\|\mathbf{a}\|} \right\|_1 \leq c_1$ , this solution is given by  $\mathbf{u} = \frac{\mathbf{a}}{\|\mathbf{a}\|}$ ; otherwise this solution

$$\mathbf{u} = \frac{\mathbf{s}(\mathbf{a}, \Delta)}{\|\mathbf{s}(\mathbf{a}, \Delta)\|} \quad (3.37)$$

**Algorithm 7**  $R$ -rank PMD

- 
- 1: **Input:**  $\mathbf{Y}$ ,  $R$
  - 2: **Initialize:**  $\mathbf{Z}_0 := \mathbf{Y}$
  - 3: **for**  $r = 1, \dots, R$  **do**
  - 4:   Obtain  $(d_r, \mathbf{u}_r, \mathbf{v}_r)$  by Algorithm 6 inputting  $\mathbf{Z}_{r-1}$ ;
  - 5:    $\mathbf{Z}_r \leftarrow \mathbf{Z}_{r-1} - d_r \mathbf{u}_r \mathbf{v}_r^T$ ;
  - 6: **end for**
  - 7: **Output:**  $[d_r, \mathbf{u}_r, \mathbf{v}_r]_{r=1}^R$
- 

is obtained by a line-search algorithm for a scalar parameter  $\Delta$  so that it satisfy  $\|\mathbf{u}\|_1 = c_1$ , where  $\mathbf{s}(\mathbf{a}, \Delta)$  is calculated by  $s_i(\mathbf{a}, \Delta) = \text{sign}(a_i)[|a_i| - \Delta]_+$ . The function  $\mathbf{s}(\mathbf{a}, \Delta)$  is called the soft-thresholding operator. The PMD with l1-norm penalty can be regarded as a sparse representation of tSVD; however, sparse constraint and orthogonal constraint are incompatible, thus the PMD can not provide completely orthogonal factor matrices.

**Fused lasso penalty**

In this section, we consider to choose an FL penalty function as  $P_{FL}(\mathbf{u}) := \sum_{i=1}^I |u_i| + \lambda \sum_{i=2}^I |u_i - u_{i-1}|$ . In case of FL penalty, we solve a slightly modified Lagrange form, for simplicity, which is given by

$$\mathbf{u} \leftarrow \underset{\mathbf{u}}{\text{argmin}} \frac{1}{2} \|\mathbf{u} - \mathbf{a}\|_F^2 + \lambda_1 \|\mathbf{u}\|_1 + \lambda_2 \sum_{i=2}^I |u_i - u_{i-1}|, \quad (3.38)$$

where  $\lambda_1$  and  $\lambda_2$  are trade-off parameters. This problem can be solved by the fused lasso regression [36, 96, 47]. The PMD with the FL penalty can be regarded as a sparse and smooth representation of tSVD. Although it can not provide completely orthogonal factor matrices, the smooth representations are focused in wide research area. In brain signal processing, temporal or spatial smoothness is considered as an important feature.

**3.1.3 Nonnegative Matrix Factorization**

Nonnegative matrix factorization (NMF) is a technique to decompose a nonnegative matrix into two nonnegative matrices. The standard NMF model is given by

$$\mathbf{Y} \simeq \mathbf{A}\mathbf{X} \in \mathbb{R}_+^{I \times J}, \quad (3.39)$$

where  $\mathbf{A} \in \mathbb{R}_+^{I \times R}$ ,  $\mathbf{X} \in \mathbb{R}_+^{R \times J}$ .  $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_J]$  is an observation matrix consisting of  $J$  observations.  $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_R]$  is a feature matrix including feature vector bases, and  $R$  is the number of feature vectors.  $\mathbf{X}$  is a mixing matrix to reproduce the observation matrix  $\mathbf{Y}$  by  $\mathbf{y}_j \simeq \sum_{r=1}^R \mathbf{a}_r x_{rj}$ . The goal of NMF is to obtain  $\mathbf{A}$  and  $\mathbf{X}$  from the observation matrix  $\mathbf{Y}$  with a given parameter  $R$ . For example, we let  $R \leq J \ll I$  in BSS problems [18]. In this case, we want

to find  $R$  latent source signals from  $J$  mixed observations. The NMF gives  $\mathbf{A}$  as an estimator of the latent source signals. In case of extracting parts of face images [62],  $I$  and  $J$  are respectively the number of pixels and images, then the NMF represents each face image by the linear combination of  $R$  nonnegative parts. In case of clustering tasks [106],  $\mathbf{A}$  expresses the set of centroids of cluster, and  $\mathbf{X}$  represents the weight parameters of clusters. In this way, the NMF is a high-potential technique and it can be applied in a wide area of a real world data analysis.

NMF has a long history from 1971 in chemometrics under the name of “self modeling curve resolution” [61]. In early works, it had been researched under the name of “positive matrix factorization” [78] in the middle of 1990s. It becomes more widely known as “nonnegative matrix factorization” after the publication on the journal of Nature in 1999 [62].

In this section I introduce several basic algorithms for NMF. The standard NMF criterion can be given by

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2, \\ & \text{subject to} \quad \mathbf{A} \geq 0, \mathbf{X} \geq 0. \end{aligned} \quad (3.40)$$

Its criterion is very simple; however, this problem is not convex and its solution is not unique. Thus, many algorithms have been studied and the researches of solution algorithm of NMF are still going now. In this section, I introduce several standard algorithms to solve the NMF problem.

### ALS algorithm

In this section, I introduce the most standard algorithm for NMF which is called the alternating least squares (ALS) algorithm [26]. First we separate the NMF problem (3.40) into two sub-problems as follows:

$$\mathbf{A} \leftarrow \underset{\mathbf{A}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 \text{ s.t. } \mathbf{A} \geq 0, \quad (3.41)$$

$$\mathbf{X} \leftarrow \underset{\mathbf{X}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 \text{ s.t. } \mathbf{X} \geq 0. \quad (3.42)$$

Each problem is convex and have a unique solution. The procedure to solve the sub-problem for  $\mathbf{A}$  consists of only two steps as follow:

$$\mathbf{A} \leftarrow \underset{\mathbf{A}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2; \quad (3.43)$$

$$\mathbf{A} \leftarrow [\mathbf{A}]_+; \quad (3.44)$$

This procedure does not provide strict solution of (3.41); however, it provides an approximated solution of (3.41) and performs fast calculation. In a similar way, we can solve the sub-problem for  $\mathbf{X}$  as follow:

$$\mathbf{X} \leftarrow \underset{\mathbf{X}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2; \quad (3.45)$$

$$\mathbf{X} \leftarrow [\mathbf{X}]_+; \quad (3.46)$$

**Algorithm 8** ALS-NMF algorithm

- 
- 1: **Input:**  $\mathbf{Y}$ ,  $R$
  - 2: **Initialize:**  $\mathbf{X}$  randomly
  - 3: **repeat**
  - 4:    $\mathbf{A} \leftarrow [(\mathbf{X}\mathbf{X}^T)^{-1}\mathbf{X}\mathbf{Y}^T]^T$  ;
  - 5:    $\mathbf{A} \leftarrow [\mathbf{A}]_+$  ;
  - 6:    $\mathbf{X} \leftarrow (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{Y}$  ;
  - 7:    $\mathbf{X} \leftarrow [\mathbf{X}]_+$  ;
  - 8: **until** convergence
  - 9: **Output:**  $\mathbf{A}$ ,  $\mathbf{X}$
- 

The update procedures for  $\mathbf{A}$  and  $\mathbf{X}$  include standard least squares problems (3.43) and (3.45). Both problems can be solved by the same calculation method in principle. First the objective function can be transformed as

$$\begin{aligned}
L &:= \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 \\
&= \frac{1}{2} \|\mathbf{Y}^T - \mathbf{X}^T\mathbf{A}^T\|_F^2 \\
&= \frac{1}{2} \text{tr}[(\mathbf{Y} - \mathbf{A}\mathbf{X})^T(\mathbf{Y} - \mathbf{A}\mathbf{X})] \\
&= \frac{1}{2} \text{tr}(\mathbf{Y}^T\mathbf{Y}) - \text{tr}(\mathbf{X}^T\mathbf{A}^T\mathbf{Y}) + \frac{1}{2} \text{tr}(\mathbf{X}^T\mathbf{A}^T\mathbf{A}\mathbf{X}).
\end{aligned} \tag{3.47}$$

Thus the stationary points of  $\mathbf{X}$  can be founded by equating the gradient to 0:

$$\frac{\partial L}{\partial \mathbf{X}} = -\mathbf{A}^T\mathbf{Y} + \mathbf{A}^T\mathbf{A}\mathbf{X} = 0. \tag{3.48}$$

Thus, the solution is given by

$$\mathbf{X} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{Y}. \tag{3.49}$$

In a same way, the solution of  $\mathbf{A}$  is given by

$$\mathbf{A} = [(\mathbf{X}\mathbf{X}^T)^{-1}\mathbf{X}\mathbf{Y}^T]^T. \tag{3.50}$$

Finally, the ALS algorithm is summarized as Algorithm 8.

### Multiplicative update rule

The ALS-NMF algorithm includes an inverse problem and thresholding for nonnegativity which are not stable. In practice, ALS algorithm do not converge, sometimes. The multiplicative-NMF algorithm [62] is proposed to update them by multiplying each element of  $\mathbf{A}$  or  $\mathbf{X}$  by some nonnegative value. This update rule is appropriate for NMF since it does not need any thresholding. The multiplicative-NMF algorithm is summarized as Algorithm 9.

**Algorithm 9** Multiplicative-NMF algorithm

- 
- 1: **Input:**  $\mathbf{Y}$ ,  $R$
  - 2: **Initialize:**  $\mathbf{X}$  randomly
  - 3: **repeat**
  - 4:    $\mathbf{A} \leftarrow \mathbf{A} \circledast (\mathbf{Y}\mathbf{X}^T) \oslash (\mathbf{A}\mathbf{X}\mathbf{X}^T)$  ;
  - 5:    $\mathbf{X} \leftarrow \mathbf{X} \circledast (\mathbf{A}^T\mathbf{Y}) \oslash (\mathbf{A}^T\mathbf{A}\mathbf{X})$  ;
  - 6: **until** convergence
  - 7: **Output:**  $\mathbf{A}$ ,  $\mathbf{X}$
- 

**HALS algorithm**

Finally, I introduce the hierarchical alternating least squares (HALS) algorithm [24] for NMF. The HALS is related to ALS; however it does not include an inverse problem and is stable. The NMF model can be rewritten by

$$\mathbf{Y} \simeq \mathbf{A}\mathbf{B}^T = \sum_{r=1}^R \mathbf{a}_r \mathbf{b}_r^T, \quad (3.51)$$

where  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_R] \in \mathbb{R}_+^{J \times R}$  is equivalent to  $\mathbf{X}^T$ . In HALS algorithm, we consider a local-problem for  $\mathbf{a}_r$  and  $\mathbf{b}_r$  as follow:

$$\text{minimize } \|\mathbf{Y}_r - \mathbf{a}_r \mathbf{b}_r^T\|_F^2, \text{ s.t. } \mathbf{a}_r \geq 0, \mathbf{b}_r \geq 0, \|\mathbf{a}_r\| = 1, \quad (3.52)$$

where  $\mathbf{Y}_r := \mathbf{Y} - \sum_{k \neq r} \mathbf{a}_k \mathbf{b}_k^T$ . We iterate to solve this local-problem for  $r = 1, \dots, R$  until convergence. Thus, we need update rules for  $\mathbf{a}_r$  and  $\mathbf{b}_r$ . Since this local-problem (3.52) can be regarded as a nonnegativity constrained version of one-rank tSVD problem (3.22), then the update rules are given by

$$\mathbf{a}_r \leftarrow [\mathbf{Y}_r \mathbf{b}_r]_+; \quad (3.53)$$

$$\mathbf{a}_r \leftarrow \mathbf{a}_r / \|\mathbf{a}_r\|; \quad (3.54)$$

$$\mathbf{b}_r \leftarrow [\mathbf{Y}_r^T \mathbf{a}_r]_+; \quad (3.55)$$

Finally, the HALS-NMF algorithm is summarized as Algorithm 10. The HALS algorithm is similar to the deflation based algorithm for tSVD/PMD. The deflation based algorithm strictly solve one-rank approximation problem to repeat the updates inside the for-loop, one by one. On the other hand, the HALS algorithm has the repeat outside the for-loop. Thus, HALS can obtain more global solution than the deflation based algorithm. In fact, orthogonality is necessary to use deflation based algorithm properly. From this reason, the deflation based algorithm of PMD can not solve overall problem, but it solves just local-problem one by one using the matrix deflation technique. The HALS algorithm may have a possibility of improvement of PMD.

---

**Algorithm 10** HALS-NMF algorithm

---

- 1: **Input:**  $Y, R$
  - 2: **Initialize:**  $[\mathbf{a}_r, \mathbf{b}_r]_{r=1}^R$  randomly
  - 3:  $\mathbf{E} \leftarrow \mathbf{Y} - \sum_{r=1}^R \mathbf{a}_r \mathbf{b}_r^T$  ;
  - 4: **repeat**
  - 5:   **for**  $r = 1, \dots, R$  **do**
  - 6:      $\mathbf{Y}_r \leftarrow \mathbf{E} + \mathbf{a}_r \mathbf{b}_r^T$  ;
  - 7:      $\mathbf{a}_r \leftarrow [\mathbf{Y}_r \mathbf{b}_r]_+$  ;
  - 8:      $\mathbf{a}_r \leftarrow \mathbf{a}_r / \|\mathbf{a}_r\|$  ;
  - 9:      $\mathbf{b}_r \leftarrow [\mathbf{Y}_r^T \mathbf{a}_r]_+$  ;
  - 10:     $\mathbf{E} \leftarrow \mathbf{Y} - \mathbf{a}_r \mathbf{b}_r^T$  ;
  - 11:   **end for**
  - 12: **until** convergence
  - 13: **Output:**  $[\mathbf{a}_r, \mathbf{b}_r]_{r=1}^R$
- 

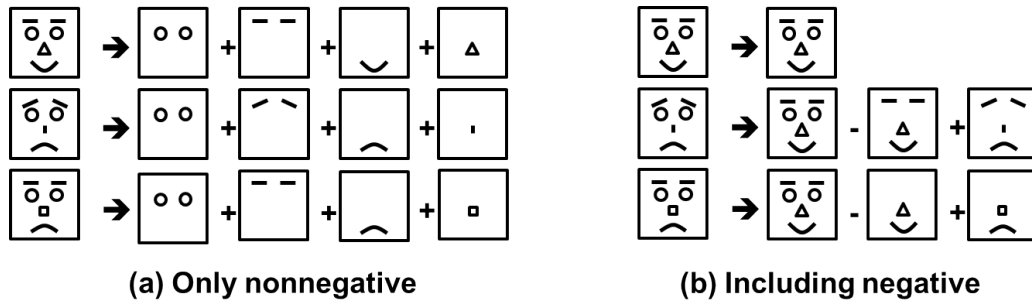


Figure 3.1: Comparison between only nonnegative factors and including negative factors

### 3.1.4 Why is NMF important and well-researched?

In this section, I discuss the importance of nonnegativity constraint. In the NMF, an observation is expressed by  $\mathbf{y}_j \simeq \sum_{r=1}^R x_{rj} \mathbf{a}_r$ , where  $x_{rj} \geq 0$  and  $\mathbf{a}_r \geq 0$ . We can not cancel out any wasted components because any subtraction is not allowed in the NMF. As a result, the NMF provides local independent features  $\{\mathbf{a}_r\}$ . Figure 3.1 shows two examples of the facial-parts representation. If the subtraction is allowed, these faces would be separated into meaningless parts (Fig. 3.1 (b)). On the other hand, NMF can separate faces into the meaningful facial parts (Fig. 3.1 (a)).

Furthermore, nonnegativity constraint is important to get nonnegative component because almost all components of nonnegative data should be nonnegative. For example, image, video, probability and stock value do not include any negative values, so their parts souled be also nonnegative.

However, an optimization criterion of the NMF (3.40) is not convex and its solution is not unique. Thus, many algorithms have been studied and the researches of solution algorithm of NMF are still going now. The criterion is often based on the minimization of Frobenius norm; however, the uses of

many divergence measures have been proposed such as KL divergence [62], Csiszar's divergence [23],  $\alpha$ - and  $\beta$ -divergences [25], Itakura-Saito divergence [34], and AB-divergence [20].

### 3.1.5 Combined methods of several constraints

In this section, I introduce several approaches to combine two or more constraints for matrix decomposition such as a combination of sparsity and nonnegativity.

#### Sparse NMF

A combination of sparsity and nonnegativity has been firstly proposed by Hoyer et al. [48] in 2002. After that, the algorithms for sparse NMF have been studied by many researchers [49, 94, 22, 55, 57, 40, 39]. In this section, I introduce two methods for the sparse NMF.

The first method is the regularized alternating least squares (RALS) algorithm [22]. The RALS algorithm prevents instability of the ALS algorithm to add some constraint terms for regularization and sparsity. These terms work to improve the uniqueness of NMF. The criterion of the RALS algorithm is given by

$$\begin{aligned} \text{minimize } D_F^{(\alpha)}(\mathbf{Y}||\mathbf{A}\mathbf{X}) &:= \frac{1}{2} \|\mathbf{W}^{-\frac{1}{2}}(\mathbf{Y} - \mathbf{A}\mathbf{X})\|_F^2 + \alpha_{A_s} \|\mathbf{A}\|_1 + \alpha_{X_s} \|\mathbf{X}\|_1 \\ &\quad + \frac{\alpha_{A_r}}{2} \|\mathbf{W}^{-\frac{1}{2}}\mathbf{A}\mathbf{L}_A\|_F^2 + \frac{\alpha_{X_r}}{2} \|\mathbf{L}_X\mathbf{X}\|_F^2, \\ \text{subject to } \mathbf{A} \geq 0, \mathbf{X} \geq 0, \end{aligned} \tag{3.56}$$

where  $\mathbf{W} \in \mathbb{R}^{I \times I}$  is a symmetric positive definite weighting matrix,  $\alpha_{A_s}$  and  $\alpha_{X_s}$  are trade-off parameters to tune sparsity level,  $\alpha_{A_r}$  and  $\alpha_{X_r}$  are trade-off parameters for regularization, and  $\mathbf{L}_A$  and  $\mathbf{L}_X$  are regularization matrices to enforce a certain application-dependent characteristics of the solution. The first term of the objective function can be transformed as

$$\|\mathbf{W}^{-\frac{1}{2}}(\mathbf{Y} - \mathbf{A}\mathbf{X})\|_F^2 = \text{tr}[(\mathbf{Y} - \mathbf{A}\mathbf{X})^T \mathbf{W}^{-1}(\mathbf{Y} - \mathbf{A}\mathbf{X})]. \tag{3.57}$$

If we set  $\mathbf{W}$  to be covariance matrix of  $\mathbf{Y}$ , then it can be regarded as the Mahalanobis distance. In the RALS algorithm, problem (3.56) is solved by the following update rules:

$$\mathbf{A} \leftarrow (\mathbf{Y}\mathbf{X}^T - \alpha_{A_s} \mathbf{W}\mathbf{S}_A + \alpha_{A_r} \mathbf{A}\mathbf{L}_A\mathbf{L}_A^T)(\mathbf{X}\mathbf{X}^T + \alpha_{A_r} \mathbf{L}_A\mathbf{L}_A^T)^\dagger, \tag{3.58}$$

$$\mathbf{X} \leftarrow (\mathbf{A}^T \mathbf{W}^{-1} \mathbf{A}^T + \alpha_{X_r} \mathbf{L}_X^T \mathbf{L}_X)^\dagger (\mathbf{A}^T \mathbf{W}^{-1} \mathbf{Y} - \alpha_{X_s} \mathbf{S}_X + \alpha_{X_r} \mathbf{L}_X^T \mathbf{L}_X \mathbf{X}), \tag{3.59}$$

where  $\mathbf{S}_A := \text{sign}(\mathbf{A})$  and  $\mathbf{S}_X := \text{sign}(\mathbf{X})$ .

The second method is the Kim&Park's sparse NMF (KP-NMF) algorithm [55, 57]. Since the KP-NMF algorithm was proposed to apply to clustering task, the left-factor matrix is regularized as centroids and the right-factor matrix is constrained to be sparse for the weighting parameter vectors.

Then the criterion is given by

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 + \eta \|\mathbf{A}\|_F^2 + \beta \sum_{j=1}^J \|\mathbf{x}_j\|_1 \\ & \text{subject to} \quad \mathbf{A} \geq 0, \mathbf{X} \geq 0. \end{aligned} \quad (3.60)$$

This problem can be separated into two sub-problems as follow:

$$\mathbf{X} \leftarrow \underset{\mathbf{X}}{\operatorname{argmin}} \left\| \begin{pmatrix} \mathbf{A} \\ \sqrt{\beta} \mathbf{1}_R^T \end{pmatrix} \mathbf{X} - \begin{pmatrix} \mathbf{Y} \\ \mathbf{0}_J^T \end{pmatrix} \right\|_F^2, \quad \text{s.t. } \mathbf{X} \geq 0, \quad (3.61)$$

$$\mathbf{A} \leftarrow \underset{\mathbf{A}}{\operatorname{argmin}} \left\| \begin{pmatrix} \mathbf{X}^T \\ \sqrt{\eta} \mathbf{I}_R^T \end{pmatrix} \mathbf{A}^T - \begin{pmatrix} \mathbf{Y}^T \\ \mathbf{0}_{R \times J}^T \end{pmatrix} \right\|_F^2, \quad \text{s.t. } \mathbf{A} \geq 0, \quad (3.62)$$

where  $\mathbf{1}_R$  is the  $R$ -dimensional vector of which all elements are 1,  $\mathbf{0}_J$  is the  $J$ -dimensional vector of which all elements are 0, and  $\mathbf{0}_{R \times J}$  is an  $(R \times J)$ -matrix of which all elements are 0. Both sub-problems can be characterized as the nonnegative least squares (NLS) problem. Solution methods for NLS problem are studied by many researches [10, 100, 56, 114]. The KP-NMF algorithm employs the active-set algorithm from [56].

### Smooth NMF

Temporal and spatial smoothness is focused in wide area of data analysis such as visual-audio data and brain signal processing [13, 116, 117, 32, 33]. Although real world data is often noisy, if we treat temporally or spatially smooth signals (e.g., natural image data, alpha brain wave, and financial data), smooth NMF [13, 116, 117, 32, 33, 115] is useful and effective. There are many methods for the smooth NMF and they can be separated into two approaches, mainly. The first one is to add some smooth constraint term into the NMF criterion. Chen et al. [13] proposed to add a temporal smoothness constraint and a spatial decorrelation constraint terms into the Frobenius norm and the Kullback-Leibler (KL) divergence based NMF for EEG analysis. Zdunek et al. [116, 117] proposed to add the Gibbs regularization term for smooth NMF. Drakakis et al. [32] proposed to add a sparseness constraint for mixing matrix and a smoothness constraint for feature matrix into the Frobenius norm and the KL divergence based NMF for analysis of financial data. Essid et al. [33] applied the KL divergence based smooth NMF to audiovisual document structuring. In this section, I introduce two methods for the smooth NMF.

The first method is the Gibbs regularized NMF [116, 117]. The Gibbs regularization is proposed for robustness to Gaussian noise, and its cost function can be expressed as

$$\Psi := \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 + 2\alpha U(\mathbf{X}), \quad (3.63)$$

where  $U(\mathbf{X})$  is a total energy function that measures the total roughness of  $\mathbf{X}$ , and  $\alpha$  is a regularization parameter. The stationary points of  $\Psi$  can be derived from the gradients of  $\Psi$  with respect to  $\mathbf{X}$  and

**A.** Its conditions are given by

$$\nabla_{\mathbf{X}}\Psi = 2\mathbf{A}^T(\mathbf{A}\mathbf{X} - \mathbf{Y}) + 2\alpha\nabla_{\mathbf{X}}U(\mathbf{X}) = \mathbf{0}, \quad (3.64)$$

$$\nabla_{\mathbf{A}}\Psi = (\mathbf{A}\mathbf{X} - \mathbf{Y})\mathbf{X}^T = \mathbf{0}. \quad (3.65)$$

From (3.64) and (3.65), we have

$$\frac{[\mathbf{A}^T\mathbf{Y} - \alpha\nabla_{\mathbf{X}}U(\mathbf{X})]_{rj}}{[\mathbf{A}^T\mathbf{A}\mathbf{X}]_{rj}} = 1, \quad \frac{[\mathbf{Y}\mathbf{X}^T]_{ir}}{[\mathbf{A}\mathbf{X}\mathbf{X}^T]_{ir}} = 1. \quad (3.66)$$

Using multiplicative update rules, we obtain the Gibbs regularized multiplicative NMF algorithm as follows:

$$\mathbf{X} \leftarrow \mathbf{X} \otimes [\mathbf{A}^T\mathbf{Y} - \alpha\nabla_{\mathbf{X}}U(\mathbf{X})]_+ \oslash [\mathbf{A}^T\mathbf{A}\mathbf{X}], \quad (3.67)$$

$$\mathbf{A} \leftarrow \mathbf{A} \otimes [\mathbf{Y}\mathbf{X}^T] \oslash [\mathbf{A}\mathbf{X}\mathbf{X}^T], \quad (3.68)$$

$$\mathbf{a}_r \leftarrow \mathbf{a}_r / \|\mathbf{a}_r\|_1, \quad (3.69)$$

Another important point is how to decide the function  $U(\mathbf{X})$ . Many candidates can be considered in [45, 42, 38, 6]. In [117], the green function is employed. Thus,  $U(\mathbf{X})$  is given by

$$U(\mathbf{X}) = \sum_{r=1}^R \sum_{j=1}^J \left[ \sum_{l \in S_j} w_{jl} \sigma \log \left\{ \cosh \left( \frac{x_{rj} - x_{rl}}{\sigma} \right) \right\} \right], \quad (3.70)$$

where  $w_{jl}$  is a weight parameter, and  $\sigma$  is a scaling parameter. Furthermore,  $\nabla_{\mathbf{X}}U(\mathbf{X})$  is given by

$$[\nabla_{\mathbf{X}}U(\mathbf{X})]_{rj} = \sum_{l \in S_j} w_{jl} \tanh \left( \frac{x_{rj} - x_{rl}}{\sigma} \right). \quad (3.71)$$

The second method is the Essid's smooth NMF [33]. Let  $\widehat{\mathbf{Y}} := \mathbf{A}\mathbf{X}$ , its criterion is given by

$$\text{minimize } D_{KL}(\mathbf{Y} \|\widehat{\mathbf{Y}}) + \beta S(\mathbf{X}), \text{ s.t. } \mathbf{A} \geq 0, \|\mathbf{a}_r\| = 1, \mathbf{X} \geq 0, \quad (3.72)$$

where  $\beta$  is a trade-off parameter for smoothness,  $D_{KL}(\mathbf{Y} \|\widehat{\mathbf{Y}})$  is the KL-divergence defined as

$$D_{KL}(\mathbf{Y} \|\widehat{\mathbf{Y}}) := \sum_{i=1}^I \sum_{j=1}^J \left( y_{ij} \log \frac{y_{ij}}{\widehat{y}_{ij}} - x_{ij} + y_{ij} \right), \quad (3.73)$$

and the smoothness constraint function  $S(\mathbf{X})$  is defined by

$$S(\mathbf{X}) := \frac{1}{2} \sum_{r=1}^R \sum_{j=2}^J (x_{rj} - x_{r(j-1)})^2. \quad (3.74)$$

We define  $\mathbf{L} \in \mathbb{R}^{J \times (J-1)}$  as

$$l_{mn} := \begin{cases} -1 & \text{if } m = n \\ 1 & \text{if } m - n = 1 \\ 0 & \text{otherwise} \end{cases}, \quad (3.75)$$

then  $S(\mathbf{X})$  can be expressed by

$$S(\mathbf{X}) = \frac{1}{2} \text{tr}(\mathbf{X}^T \mathbf{L}^T \mathbf{L} \mathbf{X}). \quad (3.76)$$

If  $J = 5$ ,  $\mathbf{L}$  can be given by

$$\mathbf{L} := \begin{pmatrix} -1 & 0 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 1 & -1 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & 0 & 1 \end{pmatrix}. \quad (3.77)$$

Update rule for  $\mathbf{X}$  can be given by

$$\mathbf{X} \leftarrow \frac{1}{2} \left[ \sqrt{\mathbf{X}_b \circledast \mathbf{X}_b + 4\mathbf{X}_a} - \mathbf{X}_b \right] \mathbf{D}_\beta^{-1}, \quad (3.78)$$

where

$$\mathbf{X}_a := \Psi \mathbf{D}_\beta, \quad (3.79)$$

$$\mathbf{X}_b := \mathbf{1}_{R \times J} + \mathbf{X} \mathbf{D}_\beta^{(2)}, \quad (3.80)$$

$$\Psi := \mathbf{X} \circledast \mathbf{A}^T (\mathbf{Y} \circledast \hat{\mathbf{Y}}), \quad (3.81)$$

$$\mathbf{D}_\beta := \text{diag}(\beta, 2\beta, \dots, 2\beta, \beta) \in \mathbb{R}^{J \times J}, \quad (3.82)$$

$$[\mathbf{D}_\beta^{(2)}]_{mn} := \begin{cases} \beta & \text{if } |m - n| = 1 \\ 0 & \text{otherwise} \end{cases} \in \mathbb{R}^{J \times J}. \quad (3.83)$$

If  $J = 5$ ,  $\mathbf{D}_\beta^{(2)}$  is given by

$$\mathbf{D}_\beta^{(2)} = \begin{pmatrix} 0 & \beta & 0 & 0 & 0 \\ \beta & 0 & \beta & 0 & 0 \\ 0 & \beta & 0 & \beta & 0 \\ 0 & 0 & \beta & 0 & \beta \\ 0 & 0 & 0 & \beta & 0 \end{pmatrix} \in \mathbb{R}^{J \times J}. \quad (3.84)$$

Update rule for  $\mathbf{A}$  can be given by

$$\mathbf{A} \leftarrow \frac{1}{2\beta} \left[ \sqrt{\mathbf{A}_b \circledast \mathbf{A}_b + 4\mathbf{A}_a} - \mathbf{A}_b \right] \mathbf{D}_s^{-1}, \quad (3.85)$$

$$\mathbf{a}_r \leftarrow \mathbf{a}_r / \|\mathbf{a}_r\|, \quad (3.86)$$

where

$$\mathbf{A}_a := \beta \Phi \mathbf{D}_s, \quad (3.87)$$

$$\mathbf{A}_b := \beta (\mathbf{1}_{I \times I} - \mathbf{I}_I) \mathbf{A} \mathbf{D}_s + \mathbf{1}_{I \times R} \mathbf{D}_x, \quad (3.88)$$

$$\Phi := \mathbf{A} \otimes \{(\mathbf{Y} \oslash \hat{\mathbf{Y}}) \mathbf{X}^T\}, \quad (3.89)$$

$$[\mathbf{D}_s]_{mn} := \begin{cases} [\mathbf{X}^T \mathbf{L}^T \mathbf{L} \mathbf{X}]_{mm} & \text{if } m = n \\ 0 & \text{otherwise} \end{cases} \in \mathbb{R}^{R \times R}, \quad (3.90)$$

$$[\mathbf{D}_x]_{mn} := \begin{cases} \sum_{j=1}^J x_{mj} & \text{if } m = n \\ 0 & \text{otherwise} \end{cases} \in \mathbb{R}^{R \times R}. \quad (3.91)$$

In summary, a standard approach of smooth NMF can be given by a combination of some divergence measure between  $\mathbf{Y}$  and  $\mathbf{A}\mathbf{X}$ , and some smooth constraint. In Section 3.1.6, I introduce a different approach for smooth NMF.

### 3.1.6 Function approximation based model

In this section, I review the smooth NMF using function approximation proposed in [115]. Please note that we impose a smoothness constraint to  $\mathbf{A}$ . This method is to approximate a feature vector  $\mathbf{a}_r$  by the following model:

$$\mathbf{a}_r = \sum_{n=1}^N \phi_n w_{nr}, \quad (3.92)$$

where  $\{w_{nr}\}$  are real-valued coefficients, and  $\{\phi_n\}$  are smooth basis functions (e.g., the Gaussian function). Let  $\Phi = [\phi_1, \dots, \phi_N] \in \mathbb{R}_+^{I \times N}$  and  $\mathbf{W} = [w_{nr}] \in \mathbb{R}^{N \times R}$ , then we have the following model for NMF:

$$\mathbf{Y} \simeq \Phi \mathbf{W} \mathbf{X}, \text{ s.t. } \Phi \mathbf{W} \geq 0, \text{ and } \mathbf{X} \geq 0. \quad (3.93)$$

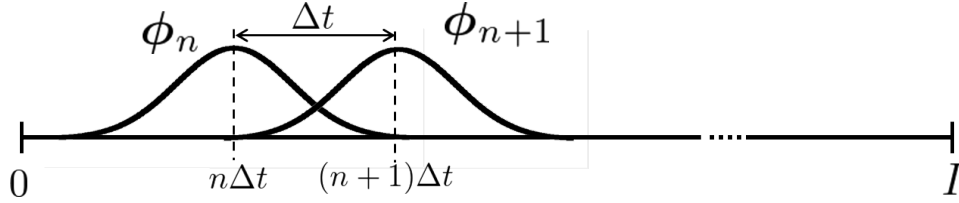
In this model, the feature matrix  $\mathbf{A}$  is approximated by  $\Phi \mathbf{W}$ . Instead the degree of freedom of estimation is decreased by function approximation  $\Phi \mathbf{W}$ , we can obtain a smooth estimator. When the observed data  $\mathbf{Y}$  includes some noise, this model can reduce affections of the noise. For optimization, we estimate two parameter matrices  $\mathbf{W}$  and  $\mathbf{X}$  since  $\Phi$  is fixed.

#### How to Choose $\Phi$

In [115] it was proposed to use the Gaussian radial basis functions (GRBF) for  $\Phi$  with a standard deviation  $\sigma$  as

$$\Phi(i, n) = \exp \left[ -\frac{(i - n\Delta t)^2}{2\sigma^2} \right], \quad (3.94)$$

where  $\Delta t$  is an interval which satisfies  $N = \lfloor (I - 1)/\Delta t \rfloor + 1$  (see Fig. 3.2). When  $\sigma$  is large, the degree of freedom of representation will decrease but it is expected that the NMF becomes robust for

Figure 3.2: Basis function  $\phi_n$ 

noisy data at the same time. On the other hand, when  $\sigma$  is small,  $\phi_n$  will reach orthogonal bases and the degree of freedom of representation will increase but it is expected that the NMF becomes weak for noisy data. Thus,  $\sigma$  can be regarded as a trade-off parameter.

### Algorithm

To estimate  $\mathbf{W}$  and  $\mathbf{X}$ , the most popular criterion minimizes the Frobenius norm as

$$\underset{\mathbf{W}, \mathbf{X}}{\text{minimize}} \frac{1}{2} \|\mathbf{Y} - \Phi \mathbf{W} \mathbf{X}\|_F^2, \text{ s.t. } \Phi \mathbf{W} \geq 0, \mathbf{X} \geq 0. \quad (3.95)$$

Since this is not convex, we separate this criterion into the following two sub-problems and solve them, alternately and iteratively:

$$\underset{\mathbf{X}}{\text{minimize}} \frac{1}{2} \|\mathbf{Y} - \Phi \mathbf{W} \mathbf{X}\|_F^2, \text{ s.t. } \mathbf{X} \geq 0, \quad (3.96)$$

$$\underset{\mathbf{W}}{\text{minimize}} \frac{1}{2} \|\mathbf{Y} - \Phi \mathbf{W} \mathbf{X}\|_F^2, \text{ s.t. } \Phi \mathbf{W} \geq 0. \quad (3.97)$$

To solve (3.96), we can apply basic NMF approach such as the right-side update rule of ALS [26]. The regularized fast combinatorial nonnegative least squares (FC-NNLS) algorithm [100] is used. This is based on the active-set algorithm. To solve (3.97), we transform the objective function to the following vectorized form (defined by (3.104)):

$$\begin{aligned} \frac{1}{2} \|\mathbf{Y} - \Phi \mathbf{W} \mathbf{X}\|_F^2 &= \frac{1}{2} \|\bar{\mathbf{y}} - (\mathbf{X}^T \otimes \Phi) \bar{\mathbf{w}}\|^2 \\ &= \frac{1}{2} \bar{\mathbf{y}}^T \bar{\mathbf{y}} - \bar{\mathbf{y}}^T (\mathbf{X}^T \otimes \Phi) \bar{\mathbf{w}} + \frac{1}{2} \bar{\mathbf{w}}^T (\mathbf{X} \mathbf{X}^T \otimes \Phi^T \Phi) \bar{\mathbf{w}}, \end{aligned} \quad (3.98)$$

where  $\bar{\mathbf{w}} = \text{vec}(\mathbf{W}) \in \mathbb{R}^{RN}$  and  $\bar{\mathbf{y}} = \text{vec}(\mathbf{Y}) \in \mathbb{R}^{IJ}$  are vectorized forms of the matrices  $\mathbf{W}$  and  $\mathbf{Y}$ , respectively, and  $\otimes$  stands for the Kronecker product. Finally, we solve the following quadratic programming (QP) problem:

$$\underset{\bar{\mathbf{w}}}{\text{minimize}} \frac{1}{2} \bar{\mathbf{w}}^T \mathbf{Q} \bar{\mathbf{w}} + \mathbf{c}^T \bar{\mathbf{w}}, \text{ s.t. } (\mathbf{I}_R \otimes \Phi) \bar{\mathbf{w}} \geq 0, \quad (3.99)$$

where  $\mathbf{Q} = \mathbf{X} \mathbf{X}^T \otimes \Phi^T \Phi \in \mathbb{R}^{RN \times RN}$ ,  $\mathbf{c} = -(\mathbf{X} \otimes \Phi^T) \bar{\mathbf{y}} \in \mathbb{R}^{RN}$ , and  $\mathbf{I}_R \in \mathbb{R}^{R \times R}$  is the identity matrix. This method is called the GRBF-NMF.

### Computational Issue

According to [115], the GRBF-NMF is robust in respect to noisy data. This optimization algorithm consists of the active set algorithm for  $\mathbf{X}$  and the QP optimization for  $\mathbf{W}$ . Each algorithm is excellent as a strict single optimization method; however, it is not so effective to use as a step of alternating optimization algorithm. The reason is that the dimension of parameter space of the QP optimization would be large (i.e., it is  $RN$ ) and the active-set method has to solve iteratively the least squares problem for evaluation of the optimality of current active-set. Thus, its computational cost becomes very high. The use of strict optimization methods is not indispensable for a step of iterative algorithm. A low cost step is better in many cases even though it does not give a strict solution but an approximated solution. Based on this philosophy, a new fast algorithm for the GRBF-NMF is proposed in Chapter 5.

## 3.2 Tensor algebra

We call a multi-dimensional array of numerical values as a tensor. An  $N$ d-tensor (i.e.,  $N$ -way array) has  $N$  indexes, for example,

- $N = 0 \rightarrow$  Scalar, e.g.,  $x$
- $N = 1 \rightarrow$  Vector, e.g.,  $x_i$
- $N = 2 \rightarrow$  Matrix, e.g.,  $x_{ij}$
- $N = 3 \rightarrow$  like a ‘‘Cube’’, e.g.,  $x_{ijk}$
- ...

Decompositions of higher-order tensors (i.e.,  $N$ -way arrays with  $N \geq 3$ ) have many applications [59] in psychometrics, chemometrics, signal processing, numerical linear algebra, computer vision, numerical analysis, data mining, neuroscience, graph analysis, and elsewhere. Higher-order (i.e.,  $N \geq 3$ ) tensors are denoted by underlined and boldface capital letters, e.g.,

$$\underline{\mathbf{Y}} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}, \quad (3.100)$$

$$\underline{\mathbf{G}} \in \mathbb{R}^{R_1 \times R_2 \times \cdots \times R_N}. \quad (3.101)$$

We denote these element-wise value by

$$y_{i_1 i_2 \cdots i_N}, \quad (3.102)$$

$$g_{r_1 r_2 \cdots r_N}. \quad (3.103)$$

Fig. 3.3 shows visualization of various orders of tensor data.

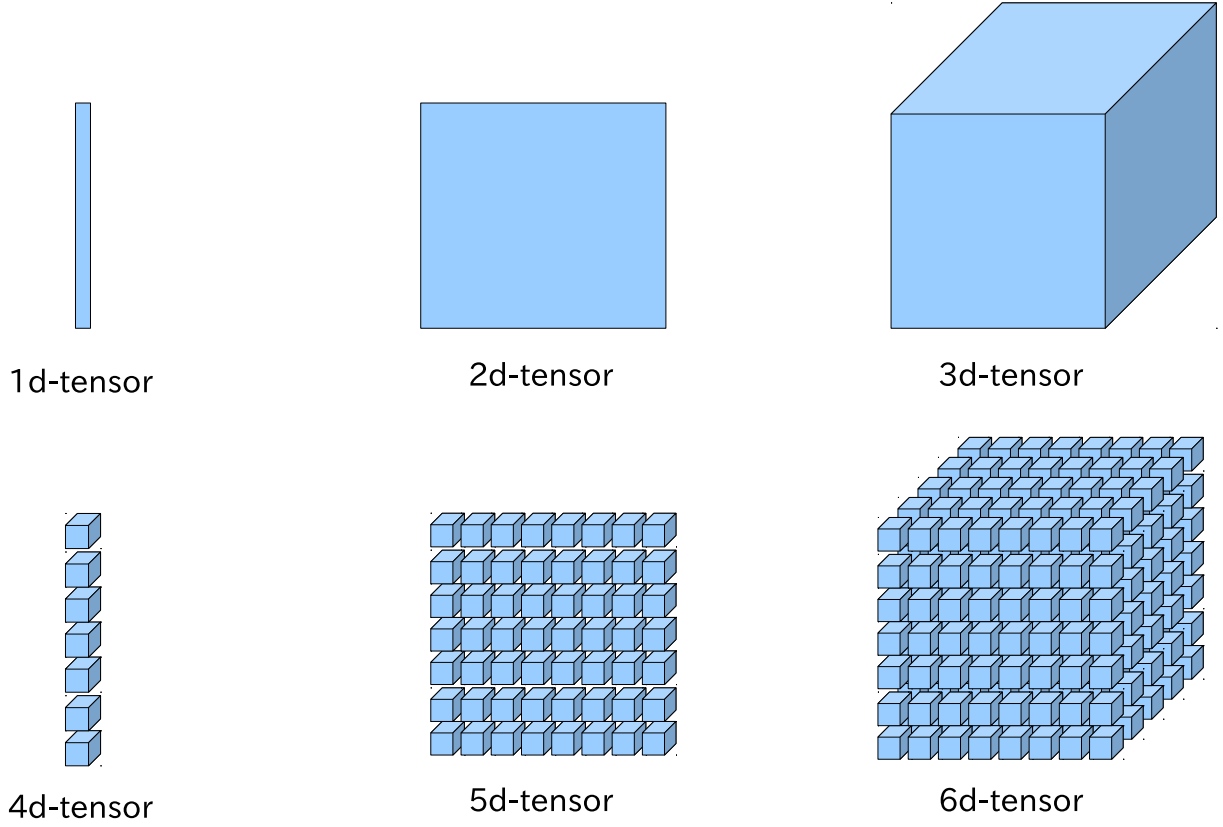


Figure 3.3: Tensors for various orders

**Vectorization** is the transformation from a tensor to a vector. To begin with, I introduce a vectorization of matrix. Vectorization of matrix  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_J] \in \mathbb{R}^{I \times J}$  is expressed by

$$\text{vec}(\mathbf{Y}) := \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_J \end{pmatrix} \in \mathbb{R}^{IJ}. \quad (3.104)$$

Fig. 3.4 shows a visualization of each-mode vectorization of a 3d-tensor. In regularly, the vectorization means the upper first-mode vectorization in Fig. 3.4. By using a notation of slice of tensor  $[\underline{\mathbf{Z}}_{::k}] \in \mathbb{R}^{I \times J}$  which is  $k$ -th slice of tensor  $\underline{\mathbf{Z}} \in \mathbb{R}^{I \times J \times K}$ , the vectorization of tensor  $\underline{\mathbf{Z}}$  is defined by

$$\text{vec}(\underline{\mathbf{Z}}) := \text{vec}([\underline{\mathbf{Z}}_{::1}, \underline{\mathbf{Z}}_{::2}, \dots, \underline{\mathbf{Z}}_{::K}]) \in \mathbb{R}^{IJK}. \quad (3.105)$$

We have

$$\|\mathbf{Y}\|_F^2 = \sum_{i,j} y_{ij}^2 = \text{tr}(\mathbf{Y}^T \mathbf{Y}) = \text{vec}(\mathbf{Y})^T \text{vec}(\mathbf{Y}) = \|\text{vec}(\mathbf{Y})\|^2, \quad (3.106)$$

$$\|\underline{\mathbf{Z}}\|_F^2 = \sum_{i,j,k} z_{ijk}^2 = \text{vec}(\underline{\mathbf{Z}})^T \text{vec}(\underline{\mathbf{Z}}) = \|\text{vec}(\underline{\mathbf{Z}})\|^2. \quad (3.107)$$

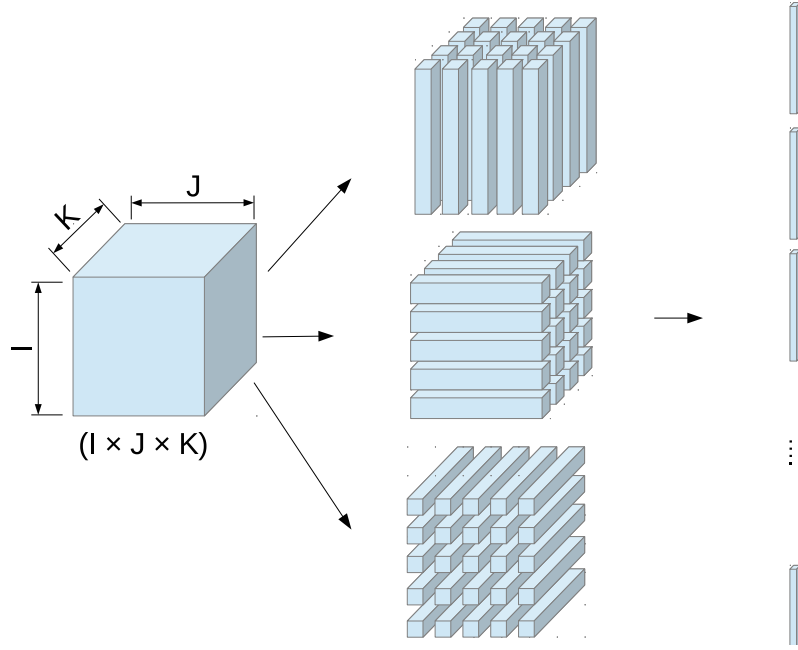


Figure 3.4: Vectorization of 3d-tensor

**Matricization** is transformation of a tensor into a matrix. The first-, second- and third-mode matricization of tensor  $\underline{\mathbf{Y}} \in \mathbb{R}^{I \times J \times K}$  can be expressed by  $\mathbf{Y}_{(1)} \in \mathbb{R}^{I \times JK}$ ,  $\mathbf{Y}_{(2)} \in \mathbb{R}^{J \times IK}$ , and  $\mathbf{Y}_{(3)} \in \mathbb{R}^{K \times IJ}$ , respectively. Fig. 3.5 shows a visualization of each-mode matricization of a 3d-tensor. Notation of a slice of tensor is often used, e.g.,  $[\underline{\mathbf{Y}}_{::k}] \in \mathbb{R}^{I \times J}$  is  $k$ -th slice of tensor  $\underline{\mathbf{Y}}$ . Using this notation, each-mode matricization can be defined as

$$\mathbf{Y}_{(1)} := [\text{vec}(\underline{\mathbf{Y}}_{1::}), \text{vec}(\underline{\mathbf{Y}}_{2::}), \dots, \text{vec}(\underline{\mathbf{Y}}_{I::})]^T \in \mathbb{R}^{I \times JK}, \quad (3.108)$$

$$\mathbf{Y}_{(2)} := [\text{vec}(\underline{\mathbf{Y}}_{:1:}), \text{vec}(\underline{\mathbf{Y}}_{:2:}), \dots, \text{vec}(\underline{\mathbf{Y}}_{:J:})]^T \in \mathbb{R}^{J \times KI}, \quad (3.109)$$

$$\mathbf{Y}_{(3)} := [\text{vec}(\underline{\mathbf{Y}}_{::1}), \text{vec}(\underline{\mathbf{Y}}_{::2}), \dots, \text{vec}(\underline{\mathbf{Y}}_{::K})]^T \in \mathbb{R}^{K \times IJ}. \quad (3.110)$$

In this way, three-order tensor has three ways of matricization. We have

$$\|\underline{\mathbf{Y}}\|_F^2 = \text{tr}(\mathbf{Y}_{(1)} \mathbf{Y}_{(1)}^T) = \|\mathbf{Y}_{(1)}\|_F^2 \quad (3.111)$$

$$= \text{tr}(\mathbf{Y}_{(2)} \mathbf{Y}_{(2)}^T) = \|\mathbf{Y}_{(2)}\|_F^2 \quad (3.112)$$

$$= \text{tr}(\mathbf{Y}_{(3)} \mathbf{Y}_{(3)}^T) = \|\mathbf{Y}_{(3)}\|_F^2. \quad (3.113)$$

**$n$ -mode product** of a tensor  $\underline{\mathbf{Z}} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$  and a matrix  $\mathbf{V} \in \mathbb{R}^{R \times I_n}$  is defined by

$$[\underline{\mathbf{Z}} \times_n \mathbf{V}]_{i_1 i_2 \dots i_{n-1} r i_{n+1} \dots i_N} := \sum_{i_n=1}^{I_n} (v_{r i_n}) (z_{i_1 i_2 \dots i_{n-1} i_n i_{n+1} \dots i_N}). \quad (3.114)$$

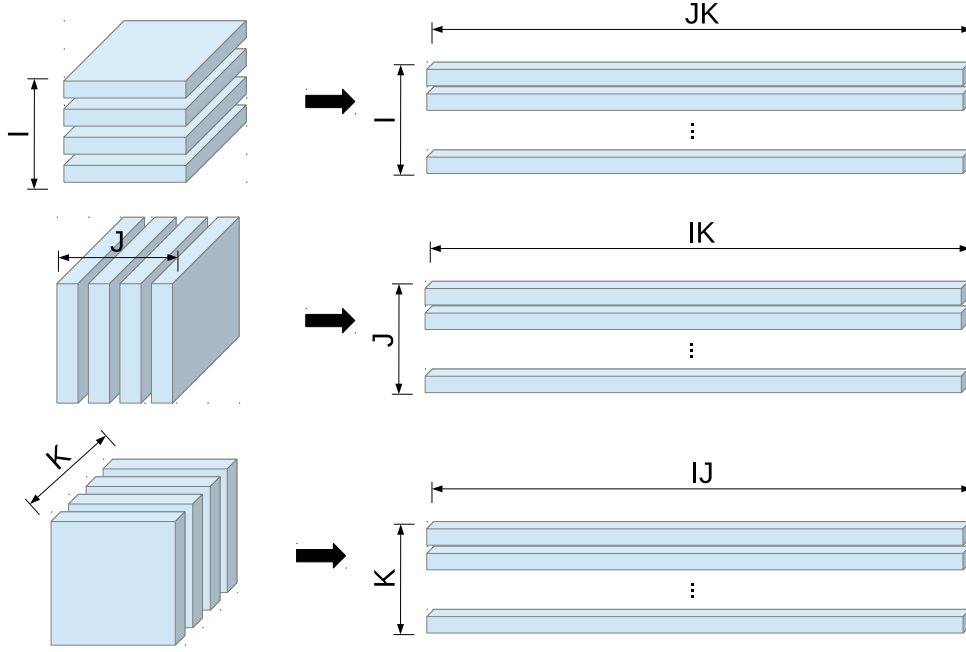


Figure 3.5: Matricization of each mode

For example, when we consider a 3d-tensor  $\underline{\mathbf{G}} \in \mathbb{R}^{R_1 \times R_2 \times R_3}$ , and the first-mode product of  $\underline{\mathbf{G}}$  and a matrix  $\mathbf{A} \in \mathbb{R}^{I \times R_1}$  is denoted by  $\underline{\mathbf{G}} \times_1 \mathbf{A} \in \mathbb{R}^{I \times R_2 \times R_3}$ . Its element-wise value is given by

$$[\underline{\mathbf{G}} \times_1 \mathbf{A}]_{ir_2r_3} = \sum_{r_1=1}^{R_1} a_{ir_1} g_{r_1r_2r_3}. \quad (3.115)$$

Fig. 3.6 shows a visualization of first-mode tensor product. In a matricized form, we have

$$[\underline{\mathbf{G}} \times_1 \mathbf{A}]_{(1)} = \mathbf{A}\mathbf{G}_{(1)}, \quad (3.116)$$

and

$$\|\underline{\mathbf{G}} \times_1 \mathbf{A}\|_F^2 = \text{tr}(\mathbf{A}\mathbf{G}_{(1)}\mathbf{G}_{(1)}^T\mathbf{A}^T). \quad (3.117)$$

**All mode product** of a tensor  $\underline{\mathbf{G}} \in \mathbb{R}^{R_1 \times R_2 \times R_3}$  and a set of matrices  $\{\mathbf{A} \in \mathbb{R}^{I \times R_1}, \mathbf{B} \in \mathbb{R}^{J \times R_2}, \mathbf{C} \in \mathbb{R}^{K \times R_3}\}$  is denoted by  $\underline{\mathbf{G}} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C} \in \mathbb{R}^{I \times J \times K}$ . We have

$$(\underline{\mathbf{G}} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C})_{(1)} = \mathbf{A}\mathbf{G}_{(1)}(\mathbf{C}^T \otimes \mathbf{B}^T), \quad (3.118)$$

$$(\underline{\mathbf{G}} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C})_{(2)} = \mathbf{B}\mathbf{G}_{(2)}(\mathbf{C}^T \otimes \mathbf{A}^T), \quad (3.119)$$

$$(\underline{\mathbf{G}} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C})_{(3)} = \mathbf{C}\mathbf{G}_{(3)}(\mathbf{B}^T \otimes \mathbf{A}^T). \quad (3.120)$$

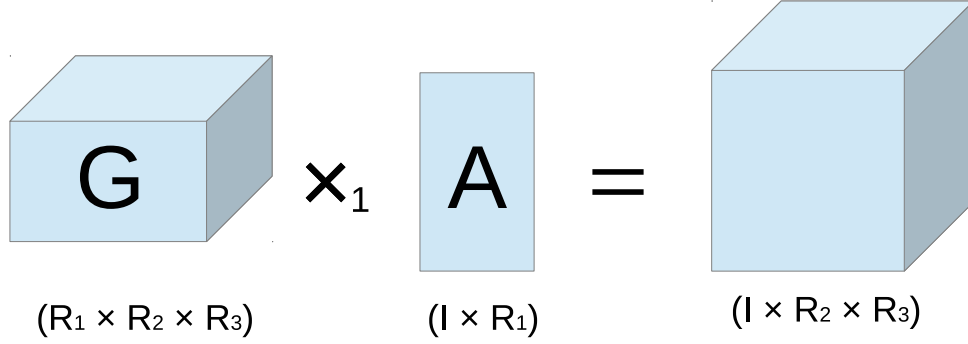


Figure 3.6: First mode tensor product

Figs. 3.7 and 3.8 show visualizations of all mode product in tensor and matricized forms. Moreover we have

$$\|\underline{\mathbf{G}} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C}\|_F^2 = \text{tr}(\mathbf{A}\mathbf{G}_{(1)}(\mathbf{C}^T\mathbf{C} \otimes \mathbf{B}^T\mathbf{B})\mathbf{G}_{(1)}^T\mathbf{A}^T), \quad (3.121)$$

$$= \text{tr}(\mathbf{B}\mathbf{G}_{(2)}(\mathbf{C}^T\mathbf{C} \otimes \mathbf{A}^T\mathbf{A})\mathbf{G}_{(2)}^T\mathbf{B}^T), \quad (3.122)$$

$$= \text{tr}(\mathbf{C}\mathbf{G}_{(3)}(\mathbf{B}^T\mathbf{B} \otimes \mathbf{A}^T\mathbf{A})\mathbf{G}_{(3)}^T\mathbf{C}^T). \quad (3.123)$$

If each matrices are orthonormal (i.e.,  $\mathbf{A}^T\mathbf{A} = \mathbf{I}_{R_1}$ ,  $\mathbf{B}^T\mathbf{B} = \mathbf{I}_{R_2}$ ,  $\mathbf{C}^T\mathbf{C} = \mathbf{I}_{R_3}$ ), thus we have

$$\|\underline{\mathbf{G}} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C}\|_F^2 = \|\underline{\mathbf{G}}\|_F^2. \quad (3.124)$$

### 3.3 CP decomposition

In this section, I introduce a practical tensor decomposition technique. First, we define a one-rank tensor that it can be written as the outer product of  $N$  vectors, i.e.,

$$\underline{\mathbf{U}} = \mathbf{u}^{(1)} \circ \mathbf{u}^{(2)} \circ \dots \circ \mathbf{u}^{(N)}. \quad (3.125)$$

The CP (canonical polyadic) decomposition factorizes a tensor into a sum of one-rank tensors. The CP model with  $N = 3$  (i.e., three way tensor decomposition) is given by

$$\underline{\mathbf{Y}} \simeq \sum_{r=1}^R \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket, \quad (3.126)$$

where  $\mathbf{A} := [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_R]$ ,  $\mathbf{B} := [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_R]$ , and  $\mathbf{C} := [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_R]$ . Fig. 3.9 shows a visualization of the CP model. The CP model is also called PARAFAC (parallel factor analysis) [43] or CANDECOMP (canonical decomposition) [11]. The CP model can be considered as the low rank approximation of tensor  $\underline{\mathbf{Y}}$ , when  $R$  is smaller than the rank of tensor  $\underline{\mathbf{Y}}$ . Thus, the CP decomposition can be considered as an extension of tSVD. We can denote the matricized forms of CP model as

$$\mathbf{Y}_{(1)} \simeq \mathbf{A}(\mathbf{C} \odot \mathbf{B})^T, \quad (3.127)$$

$$\mathbf{Y}_{(2)} \simeq \mathbf{B}(\mathbf{C} \odot \mathbf{A})^T, \quad (3.128)$$

$$\mathbf{Y}_{(3)} \simeq \mathbf{C}(\mathbf{B} \odot \mathbf{A})^T. \quad (3.129)$$

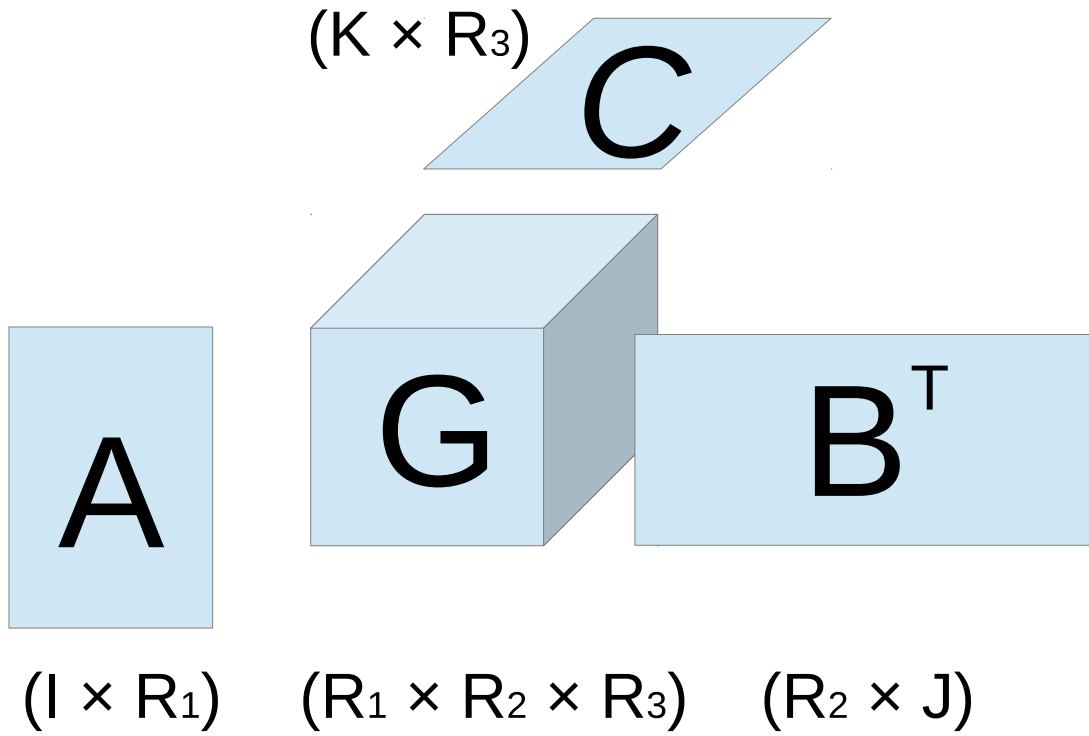


Figure 3.7: All mode product

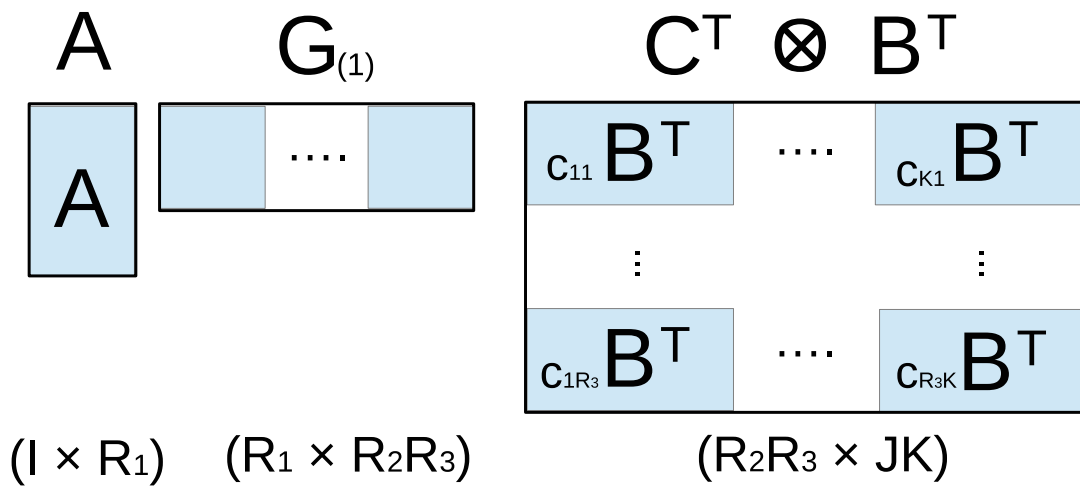


Figure 3.8: Matricization of all mode product

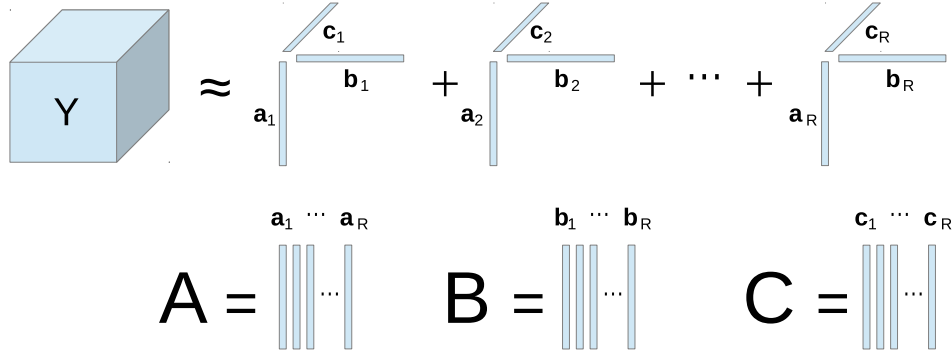


Figure 3.9: The CP decomposition

Criterion for the CP decomposition is given by

$$\text{minimize } \|\underline{\mathbf{Y}} - [\underline{\mathbf{A}}, \underline{\mathbf{B}}, \underline{\mathbf{C}}]\|_F^2. \quad (3.130)$$

When  $N \geq 3$ , the CP decomposition with  $R = \text{rank}(\underline{\mathbf{Y}})$  has an unique solution without any constraints unlike the SVD. But, when  $R < \text{rank}(\underline{\mathbf{Y}})$ , the CP decomposition can be considered as low-rank approximation of  $\underline{\mathbf{Y}}$ . It is considered that the best  $R$ -rank approximation does not exist, uniquely.

### 3.3.1 ALS algorithm

The ALS algorithm for CP decomposition is originally proposed in [43, 11]. In this method, we impose the normalization condition to factor matrices and add scalar parameters  $g_r$ , thus the criterion is given by

$$\begin{aligned} & \text{minimize } \left\| \underline{\mathbf{Y}} - \sum_{r=1}^R g_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r \right\|_F^2, \\ & \text{subject to } \|\mathbf{a}_r\| = \|\mathbf{b}_r\| = \|\mathbf{c}_r\| = 1. \end{aligned} \quad (3.131)$$

$g_r$  can be regarded as a singular value and  $\mathbf{a}_r$ ,  $\mathbf{b}_r$ , and  $\mathbf{c}_r$  can be regarded as singular vectors. Let  $\mathbf{D}_g := \text{diag}(g_1, \dots, g_R)$ , the cost function of (3.131) can be transformed as

$$\begin{aligned} L & := \|\mathbf{Y}_{(1)} - \mathbf{A} \mathbf{D}_g (\mathbf{C} \odot \mathbf{B})^T\|_F^2 \\ & = \text{tr}(\mathbf{Y}_{(1)} \mathbf{Y}_{(1)}^T) - 2 \text{tr}(\mathbf{A} \mathbf{D}_g (\mathbf{C} \odot \mathbf{B})^T \mathbf{Y}_{(1)}^T) + \text{tr}(\mathbf{A} \mathbf{D}_g (\mathbf{C}^T \mathbf{C} \otimes \mathbf{B}^T \mathbf{B}) \mathbf{D}_g \mathbf{A}^T). \end{aligned} \quad (3.132)$$

The stationary point of  $L$  with respect to  $\mathbf{A}$  is given by

$$\frac{\partial L}{\partial \mathbf{A}} = 2 \mathbf{A} \mathbf{D}_g (\mathbf{C}^T \mathbf{C} \otimes \mathbf{B}^T \mathbf{B}) \mathbf{D}_g - 2 \mathbf{Y}_{(1)} (\mathbf{C} \odot \mathbf{B}) \mathbf{D}_g = \mathbf{0}. \quad (3.133)$$

From (3.133), the update rules for  $\mathbf{A}$  can be given by

$$\mathbf{A} \leftarrow \mathbf{Y}_{(1)} (\mathbf{C} \odot \mathbf{B}) (\mathbf{C}^T \mathbf{C} \otimes \mathbf{B}^T \mathbf{B})^\dagger; \quad (3.134)$$

$$\mathbf{a}_r \leftarrow \mathbf{a}_r / \|\mathbf{a}_r\|; \quad (3.135)$$

for  $r = 1, \dots, R$ , where  $\dagger$  stands for the Moore-Penrose pseudo-inverse. Finally, ALS-CP algorithm is summarized as Algorithm 11.

**Algorithm 11** ALS-CP algorithm

---

```

1: Input:  $\underline{\mathbf{Y}}, R$ 
2: Initialize:  $\mathbf{B}, \mathbf{C}$  randomly
3: repeat
4:    $\mathbf{A} \leftarrow \mathbf{Y}_{(1)}(\mathbf{C} \odot \mathbf{B})(\mathbf{C}^T \mathbf{C} \otimes \mathbf{B}^T \mathbf{B})^\dagger$ ;
5:    $\mathbf{a}_r \leftarrow \mathbf{a}_r / \|\mathbf{a}_r\|$  for  $r = 1, \dots, R$ ;
6:    $\mathbf{B} \leftarrow \mathbf{Y}_{(2)}(\mathbf{C} \odot \mathbf{A})(\mathbf{C}^T \mathbf{C} \otimes \mathbf{A}^T \mathbf{A})^\dagger$ ;
7:    $\mathbf{b}_r \leftarrow \mathbf{b}_r / \|\mathbf{b}_r\|$  for  $r = 1, \dots, R$ ;
8:    $\mathbf{C} \leftarrow \mathbf{Y}_{(3)}(\mathbf{B} \odot \mathbf{A})(\mathbf{B}^T \mathbf{B} \otimes \mathbf{A}^T \mathbf{A})^\dagger$ ;
9:    $\mathbf{c}_r \leftarrow \mathbf{c}_r / \|\mathbf{c}_r\|$  for  $r = 1, \dots, R$ ;
10: until convergence
11:  $g_r \leftarrow \underline{\mathbf{Y}} \times_1 \mathbf{a}_r^T \times_2 \mathbf{b}_r^T \times_3 \mathbf{c}_r^T$  for  $r = 1, \dots, R$ ;
12: Output:  $g_r, \mathbf{A}, \mathbf{B}, \mathbf{C}$ 

```

---

**3.3.2 Deflation based algorithm**

In this section, I introduce a greedy deflation based algorithm for the CP decomposition. This is a straight-forward extension of the deflation based algorithm for tSVD. First, we consider the following one-rank approximation problem:

$$\begin{aligned}
& \text{minimize} \quad \|\underline{\mathbf{Y}} - g\mathbf{a} \circ \mathbf{b} \circ \mathbf{c}\|_F^2, \\
& \text{subject to} \quad \|\mathbf{a}\| = \|\mathbf{b}\| = \|\mathbf{c}\| = 1.
\end{aligned} \tag{3.136}$$

The cost function  $L$  can be transformed as

$$\begin{aligned}
L & := \|\mathbf{Y}_{(1)} - g\mathbf{a}(\mathbf{c} \otimes \mathbf{b})^T\|_F^2 \\
& = \text{tr}(\mathbf{Y}_{(1)}\mathbf{Y}_{(1)}^T) - 2\text{tr}(g\mathbf{a}(\mathbf{c} \otimes \mathbf{b})^T\mathbf{Y}_{(1)}^T) + g^2\|\mathbf{a}\|^2.
\end{aligned} \tag{3.137}$$

The stationary point of  $L$  with respect to  $\mathbf{a}_r$  is given by

$$\frac{\partial L}{\partial \mathbf{a}} = 2g^2\mathbf{a} - 2g\mathbf{Y}_{(1)}(\mathbf{c} \otimes \mathbf{b}) = \mathbf{0}. \tag{3.138}$$

From (3.138), the update rules for  $\mathbf{a}$  can be given by

$$\mathbf{a} \leftarrow \underline{\mathbf{Y}} \times_2 \mathbf{b}^T \times_3 \mathbf{c}^T; \tag{3.139}$$

$$\mathbf{a} \leftarrow \mathbf{a} / \|\mathbf{a}\|; \tag{3.140}$$

Finally, the power method for one-rank approximation and the deflation based algorithm for CP decomposition are summarized as Algorithms 12 and 13, respectively.

**3.3.3 HALS algorithm**

HALS algorithm for CP decomposition is also a straight-forward extension of the HALS-NMF algorithm. In fact, the HALS algorithm provides a better approximation than the greedy deflation based

**Algorithm 12** one-rank-CP decomposition by power method

---

```

1: Input:  $\underline{\mathbf{Y}}$ 
2: Initialize:  $\mathbf{b}, \mathbf{c}$  randomly
3: repeat
4:    $\mathbf{a} \leftarrow \underline{\mathbf{Y}} \times_2 \mathbf{b}^T \times_3 \mathbf{c}^T$  ;
5:    $\mathbf{a} \leftarrow \mathbf{a} / \|\mathbf{a}\|$  ;
6:    $\mathbf{b} \leftarrow \underline{\mathbf{Y}} \times_1 \mathbf{a}^T \times_3 \mathbf{c}^T$  ;
7:    $\mathbf{b} \leftarrow \mathbf{b} / \|\mathbf{b}\|$  ;
8:    $\mathbf{c} \leftarrow \underline{\mathbf{Y}} \times_1 \mathbf{a}^T \times_2 \mathbf{b}^T$  ;
9:    $\mathbf{c} \leftarrow \mathbf{c} / \|\mathbf{c}\|$  ;
10: until convergence
11:  $g \leftarrow \underline{\mathbf{Y}} \times_1 \mathbf{a}^T \times_2 \mathbf{b}^T \times_3 \mathbf{c}^T$  ;
12: Output:  $g, \mathbf{a}, \mathbf{b}, \mathbf{c}$ 

```

---

**Algorithm 13** Deflation-based-CP algorithm

---

```

1: Input:  $\underline{\mathbf{Y}}, R$ 
2: for  $r = 1, \dots, R$  do
3:   Obtain  $(g_r, \mathbf{a}_r, \mathbf{b}_r, \mathbf{c}_r)$  by the power method (Algorithm 12) inputting  $\underline{\mathbf{Y}}$  ;
4:    $\underline{\mathbf{Y}} \leftarrow \underline{\mathbf{Y}} - g_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$  ;
5: end for
6: Output:  $[g_r, \mathbf{a}_r, \mathbf{b}_r, \mathbf{c}_r]_{r=1}^R$ 

```

---

algorithm. Its criterion is given by

$$\begin{aligned}
& \text{minimize} && \|\underline{\mathbf{Y}}_r - g_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r\|_F^2, \\
& \text{subject to} && \|\mathbf{a}_r\| = \|\mathbf{b}_r\| = \|\mathbf{c}_r\| = 1,
\end{aligned} \tag{3.141}$$

where  $\underline{\mathbf{Y}}_r := \underline{\mathbf{Y}} - \sum_{k \neq r} g_k \mathbf{a}_k \circ \mathbf{b}_k \circ \mathbf{c}_k$ . Problems of (3.136) and (3.141) are basically equivalent. Thus the same update rules are given. Finally, we can get Algorithm 14.

### 3.3.4 Sparse/Nonnegative CP decomposition

Sparse CP decomposition can be available by imposing some sparsity constraint (e.g., l1-norm) into the CP criterion. In this case, techniques for the sparse PCA and the penalized matrix decomposition (PMD) are useful. For example, update rules (3.139) and (3.37) can be combined as

$$\mathbf{a}_r \leftarrow \frac{\mathbf{s}(\underline{\mathbf{Y}}_r \times_2 \mathbf{b}_r^T \times_3 \mathbf{c}_r^T, \Delta)}{\|\mathbf{s}(\underline{\mathbf{Y}}_r \times_2 \mathbf{b}_r^T \times_3 \mathbf{c}_r^T, \Delta)\|} \tag{3.142}$$

to obtain a sparse solution for the CP decomposition. The details are studied by [2].

In a similar way, a nonnegative version of the CP decomposition can be also available by imposing nonnegativity constraint into the CP criterion. In this case, techniques for the NMF are useful. Thus,

**Algorithm 14** HALS-CP algorithm

- 
- 1: **Input:**  $\underline{Y}$ ,  $R$
  - 2: **Initialize:**  $g_r$ ,  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$  randomly
  - 3:  $\underline{E} \leftarrow \underline{Y} - \sum_{r=1}^R g_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$  ;
  - 4: **repeat**
  - 5:    $\underline{Y}_r \leftarrow \underline{E} + g_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$  ;
  - 6:    $\mathbf{a}_r \leftarrow \underline{Y}_r \times_2 \mathbf{b}_r^T \times_3 \mathbf{c}_r^T$  ;
  - 7:    $\mathbf{a}_r \leftarrow \mathbf{a}_r / \|\mathbf{a}_r\|$  ;
  - 8:    $\mathbf{b}_r \leftarrow \underline{Y}_r \times_1 \mathbf{a}_r^T \times_3 \mathbf{c}_r^T$  ;
  - 9:    $\mathbf{b}_r \leftarrow \mathbf{b}_r / \|\mathbf{b}_r\|$  ;
  - 10:    $\mathbf{c}_r \leftarrow \underline{Y}_r \times_1 \mathbf{a}_r^T \times_2 \mathbf{b}_r^T$  ;
  - 11:    $\mathbf{c}_r \leftarrow \mathbf{c}_r / \|\mathbf{c}_r\|$  ;
  - 12:    $g_r \leftarrow \underline{Y}_r \times_1 \mathbf{a}_r^T \times_2 \mathbf{b}_r^T \times_3 \mathbf{c}_r^T$  ;
  - 13: **until** convergence
  - 14: **Output:**  $g$ ,  $\mathbf{a}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$
- 

the update rule for the nonnegative CP decomposition can be given by

$$\mathbf{a}_r \leftarrow [\underline{Y}_r \times_2 \mathbf{b}_r^T \times_3 \mathbf{c}_r^T]_+ . \quad (3.143)$$

Note that the greedy deflation based algorithm is not applicable for the nonnegative CP since the orthogonality of nonnegative factor matrices is too low. The HALS algorithm can be only applied to the nonnegative CP decomposition, directly. In the ALS algorithm for the nonnegative CP decomposition, the multiplicative update rule is stable. Since the matricized form of the CP criterion can be characterized as the standard NMF criterion, then the update rule can be given by

$$\mathbf{A} \leftarrow \mathbf{A} \circledast [\mathbf{Y}_{(1)}(\mathbf{C} \circ \mathbf{B})] \oslash [\mathbf{A}(\mathbf{C}^T \mathbf{C} \circledast \mathbf{B}^T \mathbf{B})]; \quad (3.144)$$

### 3.4 Tucker decomposition

The Tucker decomposition factorizes a tensor into some of factor matrices and a core tensor. The Tucker model with  $N = 3$  (i.e., Tucker3 model [98]) is given by

$$\underline{Y} \simeq \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \sum_{r_3=1}^{R_3} g_{r_1 r_2 r_3} \mathbf{a}_{r_1} \circ \mathbf{b}_{r_2} \circ \mathbf{c}_{r_3}, \quad (3.145)$$

$$= \underline{\mathbf{G}} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C}, \quad (3.146)$$

$$= \llbracket \underline{\mathbf{G}}; \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket. \quad (3.147)$$

Fig. 3.10 shows a visualization of Tucker3 model. The Tucker model can be considered as an generalized model of the CP decomposition since if a core tensor is diagonal tensor, then it is equivalent to the

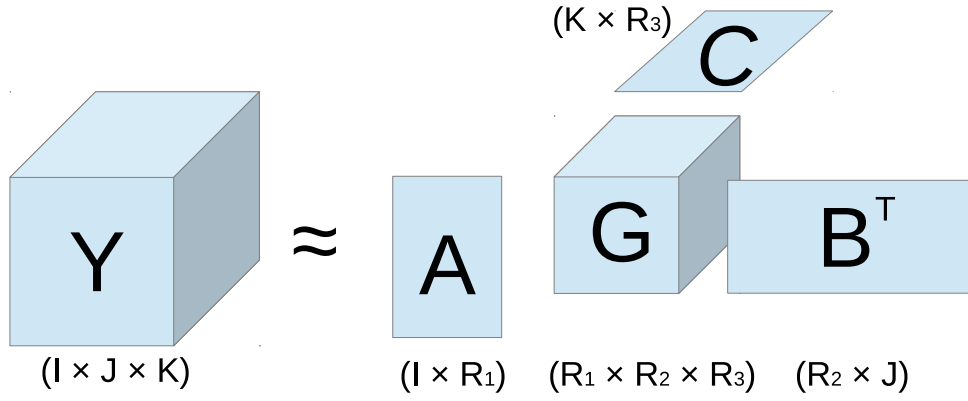


Figure 3.10: The Tucker3 decomposition

CP model (i.e., the CP model is a special case of the Tucker model). We can denote the matricized forms of the Tucker model as

$$\mathbf{Y}_{(1)} \simeq \mathbf{A}\mathbf{G}_{(1)}(\mathbf{C} \otimes \mathbf{B})^T, \quad (3.148)$$

$$\mathbf{Y}_{(2)} \simeq \mathbf{B}\mathbf{G}_{(2)}(\mathbf{C} \otimes \mathbf{A})^T, \quad (3.149)$$

$$\mathbf{Y}_{(3)} \simeq \mathbf{C}\mathbf{G}_{(3)}(\mathbf{B} \otimes \mathbf{A})^T. \quad (3.150)$$

Objective for the Tucker decomposition is given by

$$\text{minimize } \|\underline{\mathbf{Y}} - \llbracket \underline{\mathbf{G}}; \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket\|_F^2. \quad (3.151)$$

However, this solution is not unique and meaningless because the degree of freedom of the Tucker model is too high. In order to have a more meaningful decomposition, we need to impose appropriate constraints on factor matrices, e.g., orthogonality, sparsity, and nonnegativity.

### 3.4.1 HOOI algorithm for orthogonal Tucker decomposition

Degree of freedom of Tucker decomposition is too high to extract significant features since its solution is not unique. As the most basic approach, use of orthogonality constraint for each factor matrix has been introduced, thus the criterion for the three mode orthogonal Tucker decomposition can be given by

$$\begin{aligned} & \text{minimize } \|\underline{\mathbf{Y}} - \underline{\mathbf{G}} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C}\|_F^2, \\ & \text{subject to } \mathbf{A}^T \mathbf{A} = \mathbf{I}_{R_1}, \mathbf{B}^T \mathbf{B} = \mathbf{I}_{R_2}, \mathbf{C}^T \mathbf{C} = \mathbf{I}_{R_3}, \mathbf{G}_{(n)} \mathbf{G}_{(n)}^T = \mathbf{\Lambda}_{R_n}, \end{aligned} \quad (3.152)$$

for  $n = 1, 2, 3$ . Using vectorized form, the cost function can be transformed as

$$\begin{aligned} L & := \|\text{vec}(\underline{\mathbf{Y}}) - (\mathbf{C} \otimes \mathbf{B} \otimes \mathbf{A})\text{vec}(\underline{\mathbf{G}})\| \\ & = \|\text{vec}(\underline{\mathbf{Y}})\|_F^2 - 2\text{vec}(\underline{\mathbf{Y}})^T (\mathbf{C} \otimes \mathbf{B} \otimes \mathbf{A})\text{vec}(\underline{\mathbf{G}}) + \|\text{vec}(\underline{\mathbf{G}})\|_F^2. \end{aligned} \quad (3.153)$$

Thus the stationary point with respect to  $\text{vec}(\underline{\mathbf{G}})$  can be given by

$$\frac{\partial L}{\partial \text{vec}(\underline{\mathbf{G}})} = -2(\mathbf{C} \otimes \mathbf{B} \otimes \mathbf{A})^T \text{vec}(\underline{\mathbf{Y}}) + 2\text{vec}(\underline{\mathbf{G}}) = 0. \quad (3.154)$$

Finally, we can get the solution of  $\underline{\mathbf{G}}$  as

$$\underline{\mathbf{G}} = \underline{\mathbf{Y}} \times_1 \mathbf{A}^T \times_2 \mathbf{B}^T \times_3 \mathbf{C}^T. \quad (3.155)$$

By using (3.155), the cost function can be transformed as

$$\begin{aligned} L &= \|\underline{\mathbf{Y}}\|_F^2 - 2\text{vec}(\underline{\mathbf{Y}})^T (\mathbf{C} \otimes \mathbf{B} \otimes \mathbf{A}) \text{vec}(\underline{\mathbf{G}}) + \|\underline{\mathbf{G}}\|_F^2 \\ &= \|\underline{\mathbf{Y}}\|_F^2 - 2\|\underline{\mathbf{G}}\|_F^2 + \|\underline{\mathbf{G}}\|_F^2 \\ &= \|\underline{\mathbf{Y}}\|_F^2 - \|\underline{\mathbf{G}}\|_F^2 \\ &= \|\underline{\mathbf{Y}}\|_F^2 - \|\underline{\mathbf{Y}} \times_1 \mathbf{A}^T \times_2 \mathbf{B}^T \times_3 \mathbf{C}^T\|_F^2, \end{aligned} \quad (3.156)$$

thus the minimization of cost function is equivalent to the maximization of  $\|\underline{\mathbf{Y}} \times_1 \mathbf{A}^T \times_2 \mathbf{B}^T \times_3 \mathbf{C}^T\|_F^2$ .

This term can be transformed to

$$\begin{aligned} &\text{tr}[\mathbf{A}^T \mathbf{Y}_{(1)} (\mathbf{C} \mathbf{C}^T \otimes \mathbf{B} \mathbf{B}^T) \mathbf{Y}_{(1)}^T \mathbf{A}] \\ &= \text{tr}[\mathbf{B}^T \mathbf{Y}_{(2)} (\mathbf{C} \mathbf{C}^T \otimes \mathbf{A} \mathbf{A}^T) \mathbf{Y}_{(2)}^T \mathbf{B}] \\ &= \text{tr}[\mathbf{C}^T \mathbf{Y}_{(3)} (\mathbf{B} \mathbf{B}^T \otimes \mathbf{A} \mathbf{A}^T) \mathbf{Y}_{(3)}^T \mathbf{C}]. \end{aligned} \quad (3.157)$$

Thus, the criterion with respect to each factor matrix is equivalence to the PCA criterion. For example, solution of  $\mathbf{A}$  can be given as  $R_1$  eigenvectors of  $\mathbf{Y}_{(1)} (\mathbf{C} \mathbf{C}^T \otimes \mathbf{B} \mathbf{B}^T) \mathbf{Y}_{(1)}^T$  or  $R_1$  left-singular vectors of  $\mathbf{Y}_{(1)} (\mathbf{C} \otimes \mathbf{B})$ .

The first algorithm for the Tucker decomposition is proposed as ‘‘Tucker’s method I’’ in [99]. Today, this method is better-known as the ‘‘higher-order SVD’’ (HOSVD) from [28]. The HOSVD does not minimize the cost function, however it can be used as an initialization algorithm (see Algorithm 15). Thus, the first algorithm in a real sense could be proposed as ‘‘three-mode PCA’’ by using ALS approach in [60]. After that its  $N$ -way extension is proposed in [54]. Finally, its algorithm is well-known as the higher order orthogonal iteration (HOOI) algorithm from [29]. The HOOI algorithm can be summarized as Algorithm 15.

### 3.4.2 Sparse Tucker decomposition

Sparse Tucker decomposition is a method to impose some sparsity constraint into each factor matrix of the Tucker decomposition. As above mentioned, the criterion for each factor matrix can be characterized as PCA criterion. Thus, we can apply the techniques of the sparse PCA or the penalized matrix decomposition to the method. For example, let  $\underline{\mathbf{V}} := \underline{\mathbf{Y}} \times_2 \mathbf{B}^T \times_3 \mathbf{C}^T$ , the SPCA like criterion for  $\mathbf{A}$  of the Tucker decomposition can be given by

$$\begin{aligned} &\text{minimize} \quad \|\mathbf{V}_{(1)} - \tilde{\mathbf{A}} \mathbf{A}^T \mathbf{V}_{(1)}\| + \lambda \|\mathbf{A}\|_F^2 + \sum_{r_1=1}^{R_1} \|\mathbf{a}_{r_1}\|_1, \\ &\text{subject to} \quad \tilde{\mathbf{A}}^T \tilde{\mathbf{A}} = \mathbf{I}_{R_1}. \end{aligned} \quad (3.158)$$

**Algorithm 15** HOOI algorithm for orthogonal Tucker decomposition

- 
- 1: **Input:**  $\underline{\mathbf{Y}}, R_1, R_2, R_3$
  - 2: [Initialization by HOSVD algorithm]
  - 3:  $\mathbf{A} \leftarrow R_1$  left-singular vectors of  $\mathbf{Y}_{(1)}$  ;
  - 4:  $\mathbf{B} \leftarrow R_2$  left-singular vectors of  $\mathbf{Y}_{(2)}$  ;
  - 5:  $\mathbf{C} \leftarrow R_3$  left-singular vectors of  $\mathbf{Y}_{(3)}$  ;
  - 6:  $\underline{\mathbf{G}} \leftarrow \underline{\mathbf{Y}} \times_1 \mathbf{A}^T \times_2 \mathbf{B}^T \times_3 \mathbf{C}^T$  ;
  - 7: [Optimization step]
  - 8: **repeat**
  - 9:    $\mathbf{A} \leftarrow R_1$  left-singular vectors of  $\mathbf{Y}_{(1)}(\mathbf{C} \otimes \mathbf{B})$  ;
  - 10:    $\mathbf{B} \leftarrow R_2$  left-singular vectors of  $\mathbf{Y}_{(2)}(\mathbf{C} \otimes \mathbf{A})$  ;
  - 11:    $\mathbf{C} \leftarrow R_3$  left-singular vectors of  $\mathbf{Y}_{(3)}(\mathbf{B} \otimes \mathbf{A})$  ;
  - 12:    $\underline{\mathbf{G}} \leftarrow \underline{\mathbf{Y}} \times_1 \mathbf{A}^T \times_2 \mathbf{B}^T \times_3 \mathbf{C}^T$  ;
  - 13: **until** convergence
  - 14: **Output:**  $\underline{\mathbf{G}}, \mathbf{A}, \mathbf{B}, \mathbf{C}$
- 

Furthermore, the left factor matrix of PMD of  $\mathbf{V}_{(1)}$  can be also used to update  $\mathbf{A}$ . Note that the update rule (3.155) can not be used for the sparse tucker decomposition since it has no orthogonality. In this case, we must use

$$\underline{\mathbf{G}} = \underline{\mathbf{Y}} \times_1 \mathbf{A}^\dagger \times_2 \mathbf{B}^\dagger \times_3 \mathbf{C}^\dagger. \quad (3.159)$$

### 3.4.3 Nonnegative Tucker decomposition

In this section, we consider to impose the nonnegativity constraint to all factors of Tucker decomposition. In a similar way to CP decomposition, its update rules can be given by applying the multiplicative update rule for NMF to the matricized form of Tucker decomposition. Thus, let  $\underline{\mathbf{V}} := \underline{\mathbf{Y}} \times_2 \mathbf{B}^T \times_3 \mathbf{C}^T$ , the update rule for  $\mathbf{A}$  is given by

$$\mathbf{A} \leftarrow \mathbf{A} \circledast [\mathbf{Y}_{(1)} \mathbf{V}_{(1)}^T] \oslash [\mathbf{A} \mathbf{V}_{(1)} \mathbf{V}_{(1)}^T]; \quad (3.160)$$

The update rule for the core tensor  $\underline{\mathbf{G}}$  can be given by the vectorized form of Tucker decomposition, and expressed as

$$\underline{\mathbf{G}} \leftarrow \underline{\mathbf{G}} \circledast (\underline{\mathbf{Y}} \times_1 \mathbf{A}^T \times_2 \mathbf{B}^T \times_3 \mathbf{C}^T) \oslash (\underline{\mathbf{G}} \times_1 \mathbf{A} \mathbf{A}^T \times_2 \mathbf{B} \mathbf{B}^T \times_3 \mathbf{C} \mathbf{C}^T); \quad (3.161)$$

Finally, the algorithm for the nonnegative Tucker decomposition (NTD) [58] can be summarized as Algorithm 16. The algorithms for NTD are well studied in [83, 85, 84].

**Algorithm 16** ALS-NTD algorithm

- 
- 1: **Input:**  $\underline{Y}$ ,  $R_1$ ,  $R_2$ ,  $R_3$
  - 2: [Initialization step]
  - 3:  $\mathbf{A} \leftarrow$  left-factor matrix of NMF of  $\mathbf{Y}_{(1)}$  ;
  - 4:  $\mathbf{B} \leftarrow$  left-factor matrix of NMF of  $\mathbf{Y}_{(2)}$  ;
  - 5:  $\mathbf{C} \leftarrow$  left-factor matrix of NMF of  $\mathbf{Y}_{(3)}$  ;
  - 6:  $\underline{\mathbf{G}} \leftarrow \underline{\mathbf{G}} \circledast (\underline{\mathbf{Y}} \times_1 \mathbf{A}^T \times_2 \mathbf{B}^T \times_3 \mathbf{C}^T) \circledcirc (\underline{\mathbf{Y}} \times_1 \mathbf{A}^T \mathbf{A} \times_2 \mathbf{B}^T \mathbf{B} \times_3 \mathbf{C}^T \mathbf{C})$ ;
  - 7: [Optimization step]
  - 8: **repeat**
  - 9:    $\underline{\mathbf{V}} \leftarrow \underline{\mathbf{Y}} \times_2 \mathbf{B} \times_3 \mathbf{C}$  ;
  - 10:    $\mathbf{A} \leftarrow \mathbf{A} \circledast [\mathbf{Y}_{(1)} \mathbf{V}_{(1)}^T] \circledcirc [\mathbf{A} \mathbf{V}_{(1)} \mathbf{V}_{(1)}^T]$  ;
  - 11:    $\underline{\mathbf{V}} \leftarrow \underline{\mathbf{Y}} \times_1 \mathbf{A} \times_3 \mathbf{C}$  ;
  - 12:    $\mathbf{B} \leftarrow \mathbf{B} \circledast [\mathbf{Y}_{(2)} \mathbf{V}_{(2)}^T] \circledcirc [\mathbf{B} \mathbf{V}_{(2)} \mathbf{V}_{(2)}^T]$  ;
  - 13:    $\underline{\mathbf{V}} \leftarrow \underline{\mathbf{Y}} \times_1 \mathbf{A} \times_2 \mathbf{B}$  ;
  - 14:    $\mathbf{C} \leftarrow \mathbf{C} \circledast [\mathbf{Y}_{(3)} \mathbf{V}_{(3)}^T] \circledcirc [\mathbf{C} \mathbf{V}_{(3)} \mathbf{V}_{(3)}^T]$  ;
  - 15:    $\underline{\mathbf{G}} \leftarrow \underline{\mathbf{G}} \circledast (\underline{\mathbf{Y}} \times_1 \mathbf{A}^T \times_2 \mathbf{B}^T \times_3 \mathbf{C}^T) \circledcirc (\underline{\mathbf{G}} \times_1 \mathbf{A} \mathbf{A}^T \times_2 \mathbf{B} \mathbf{B}^T \times_3 \mathbf{C} \mathbf{C}^T)$ ;
  - 16: **until** convergence
  - 17: **Output:**  $\underline{\mathbf{G}}$ ,  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$
- 

### 3.5 Common and Individual Feature Extraction

In this section, I introduce common and individual feature extraction methods. When a group of data is given from multiple sources in a same situation, data have common features for all sources and individual features depending on individual sources. For example, when we record voices of the same sentence “this is a pen” from several persons. The audio data have a common feature as a content “this is a pen” and individual features depend on a personal differences which might be including intonations, accents, vocal sounds and so on. In this case, common feature is important for speech recognition, and individual feature is important for speaker recognition. Thus, it is very important to analyze common and individual features for group data. However, this research is quite primitive now.

**Algorithm 17** JIVE algorithm

- 
- 1: **Input:**  $\{\mathbf{Y}_k, R_k\}_{k=1}^K, R$
  - 2: **Initialize:**  $\mathbf{J} = [\mathbf{J}_1, \mathbf{J}_2, \dots, \mathbf{J}_K]$
  - 3: **repeat**
  - 4:    $\mathbf{H}_k \leftarrow \text{tSVD}(\mathbf{Y}_k - \mathbf{J}_k, R_k)$  for all  $k$ ;
  - 5:    $\mathbf{J} \leftarrow \text{tSVD}([\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_K] - [\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_K], R)$ ;
  - 6: **until** convergence
  - 7: **Output:**  $\mathbf{J}, \{\mathbf{H}_k\}_{k=1}^K$
- 

**3.5.1 JIVE**

The joint and individual variation explained (JIVE) method [66] can be regarded as the most basic model for common and individual feature analysis. JIVE model is given by

$$\mathbf{Y}_k \simeq \mathbf{J}_k + \mathbf{H}_k, \quad (3.162)$$

$$= \mathbf{A}_C \mathbf{B}_C^{(k)T} + \mathbf{A}_I^{(k)} \mathbf{B}_I^{(k)T}, \quad (3.163)$$

$$= [\mathbf{A}_C, \mathbf{A}_I^{(k)}][\mathbf{B}_C^{(k)}, \mathbf{B}_I^{(k)T}], \quad (3.164)$$

$$= \mathbf{A}^{(k)} \mathbf{B}^{(k)T}, \quad (3.165)$$

where  $\{\mathbf{Y}_k\}_{k=1}^K$  is a group of given data,  $\mathbf{J}_k$  and  $\mathbf{H}_k$  are joint and individual structure, respectively. Both structure matrices can be decomposed by

$$\mathbf{J}_k = \mathbf{A}_C \mathbf{B}_C^{(k)T}, \quad (3.166)$$

$$\mathbf{H}_k = \mathbf{A}_I^{(k)} \mathbf{B}_I^{(k)T}. \quad (3.167)$$

We put  $\mathbf{J} = [\mathbf{J}_1, \mathbf{J}_2, \dots, \mathbf{J}_K]$ , then  $\mathbf{J}$  can be factorized by

$$\mathbf{J} = \mathbf{A}_C [\mathbf{B}_C^{(1)T}, \mathbf{B}_C^{(2)T}, \dots, \mathbf{B}_C^{(K)T}]. \quad (3.168)$$

In the JIVE method, we consider an  $R$ -rank decomposition of  $\mathbf{J}$  and  $R_k$ -rank decompositions of  $\mathbf{H}_k$ . The cost function of this decomposition is defined by

$$L = \sum_{k=1}^K \|\mathbf{Y}_k - \mathbf{J}_k - \mathbf{H}_k\|_F^2. \quad (3.169)$$

In this method, we update iteratively and alternately  $\mathbf{J}$  and  $\{\mathbf{H}_k\}_{k=1}^K$  to minimize  $L$  until convergence. Thus, when  $\mathbf{J}$  is given, we find  $\{\mathbf{A}_I^{(k)}, \mathbf{B}_I^{(k)}\}_{k=1}^K$  to minimize  $\sum_{k=1}^K \|\mathbf{Y}_k - \mathbf{J}_k - \mathbf{A}_I^{(k)} \mathbf{B}_I^{(k)T}\|$ . This solution can be given by  $R_k$ -rank tSVD of  $(\mathbf{Y}_k - \mathbf{J}_k)$  for each  $\mathbf{H}_k$ . When  $\{\mathbf{H}_k\}_{k=1}^K$  is given, we find an optimal  $\mathbf{J}$  to minimize the cost function. We put  $\mathbf{Y} = [\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_K]$  and  $\mathbf{H} = [\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_K]$ , then the solution of  $\mathbf{J}$  can be given by the  $R$ -rank tSVD of  $(\mathbf{Y} - \mathbf{H})$ . Finally, the JIVE algorithm can be summarized as Algorithm 17.

### 3.5.2 Group NMF

The group NMF [63] is a technique of common and individual feature extraction for a group of nonnegative data. Its framework includes a nonnegative extension and a more generalized model of the JIVE method; however, it had been proposed before JIVE. The decomposition model of the group NMF is given by

$$\mathbf{Y}_k \simeq \mathbf{A}^{(k)} \mathbf{X}^{(k)} = [\mathbf{A}_C^{(k)} \mathbf{A}_I^{(k)}] \mathbf{X}^{(k)}, \quad (3.170)$$

where the basis matrix  $\mathbf{A}^{(k)} \in \mathbb{R}_+^{I \times R}$  is composed of two types of bases, where  $\mathbf{A}_C^{(k)} \in \mathbb{R}_+^{I \times R_C}$  consists of common bases for all sources, and  $\mathbf{A}_I^{(k)} \in \mathbb{R}_+^{I \times R_I}$  ( $R = R_C + R_I$ ) consists of bases which are depend on individual characteristics. When we impose a constraint  $\mathbf{A}_C^{(1)} = \mathbf{A}_C^{(2)} = \dots = \mathbf{A}_C^{(K)} =: \mathbf{A}_C$ , this model is equivalent to the model of JIVE. This model is called the ‘‘FFX-NMF’’. Furthermore, we consider a group of data has no individual factor (i.e.,  $R_I = 0$ ). In this case, its decomposition model can be given by

$$[\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_K] \simeq \mathbf{A}_C [\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(K)}]. \quad (3.171)$$

This model is called the ‘‘One-NMF’’.

#### Algorithm for One-NMF

One-NMF model is equivalent to the standard NMF model. We put  $\mathbf{Y} = [\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_K]$  and  $\mathbf{X} = [\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(K)}]$ , then it becomes  $\mathbf{Y} \simeq \mathbf{A}_C \mathbf{X}$  which can be solved by standard NMF algorithms such as multiplicative algorithm, ALS algorithm, and HALS algorithm. In addition, we can apply sparse NMF and smooth NMF algorithms to this model, directly.

#### Algorithm for FFX-NMF

Introducing  $\mathbf{X}^{(k)} = [\mathbf{X}_C^{(k)}, \mathbf{X}_I^{(k)}]$ , the objective function of FFX-NMF is given by

$$J_{FFX} = \sum_{k=1}^K \|\mathbf{Y}_k - \mathbf{A}_C \mathbf{X}_C^{(k)} - \mathbf{A}_I^{(k)} \mathbf{X}_I^{(k)}\|_F^2 + \gamma \left\{ K \|\mathbf{A}_C\|_F^2 + \sum_{k=1}^K \|\mathbf{A}_I^{(k)}\|_F^2 \right\}. \quad (3.172)$$

We find optimal  $\mathbf{A}_C$ ,  $\mathbf{A}_I^{(k)}$ , and  $\mathbf{X}^{(k)}$  to minimize  $J_{FFX}$ . Individual update rules are given by

$$\mathbf{A}_C \leftarrow \mathbf{A}_C \circledast \left( \sum_k \mathbf{Y}_k \mathbf{X}_C^{(k)T} \right) \oslash \left( \sum_k \mathbf{A}^{(k)} \mathbf{X}^{(k)} \mathbf{X}_C^{(k)T} + \gamma K \mathbf{A}_C \right), \quad (3.173)$$

$$\mathbf{A}_I^{(k)} \leftarrow \mathbf{A}_I^{(k)} \circledast \left( \mathbf{Y}_k \mathbf{X}_I^{(k)T} \right) \oslash \left( \mathbf{A}^{(k)} \mathbf{X}^{(k)} \mathbf{X}_I^{(k)T} + \gamma \mathbf{A}_I^{(k)} \right), \quad (3.174)$$

$$\mathbf{X}^{(k)} \leftarrow \mathbf{X}^{(k)} \circledast \left( \mathbf{A}^{(k)T} \mathbf{Y}_k \right) \oslash \left( \mathbf{A}^{(k)T} \mathbf{A}^{(k)} \mathbf{X}^{(k)} \right). \quad (3.175)$$

The solution can be obtained by the iteration to update individual parameter matrices until convergence.

### Algorithm for GNMF

One-NMF and FFX-NMF are special cases of group NMF (GNMF). In this section, I introduce an algorithm for GNMF. We consider the following objective function:

$$\begin{aligned}
J_{GNMF} = & \lambda \sum_{k=1}^K \|\mathbf{Y}_k - \mathbf{A}^{(k)} \mathbf{X}^{(k)}\|_F^2 + \gamma \sum_{k=1}^K \|\mathbf{A}^{(k)}\| \\
& + \frac{\alpha}{2} \sum_{i \neq j} \|\mathbf{A}_C^{(i)} - \mathbf{A}_C^{(j)}\|_F^2 - \frac{\beta}{2} \sum_{i \neq j} \|\mathbf{A}_I^{(i)} - \mathbf{A}_I^{(j)}\|_F^2,
\end{aligned} \tag{3.176}$$

where  $\lambda$ ,  $\gamma$ ,  $\alpha$ , and  $\beta$  are adjusting parameters. When  $\alpha = \beta = 0$  and  $\mathbf{A}_C = \mathbf{A}_C^{(k)}$  for all  $k$ , then its objective function is equivalent to  $J_{FFX}$ . The key point here is the third and the fourth terms in the objective function. The third term minimizes the squared error between individual  $\mathbf{A}_C^{(k)}$ , and the fourth term maximizes the squared difference between individual  $\mathbf{A}_I^{(k)}$ . We find optimal  $\mathbf{A}_C^{(k)}$ ,  $\mathbf{A}_I^{(k)}$ , and  $\mathbf{X}^{(k)}$  to minimize  $J_{GNMF}$ . Update rule of  $\mathbf{X}^{(k)}$  is the same as (3.175). Update rule of  $\mathbf{A}_C^{(k)}$  and  $\mathbf{A}_I^{(k)}$  are given by

$$\mathbf{A}_C^{(k)} \leftarrow \mathbf{A}_C^{(k)} \circledast \mathbf{H}_C, \tag{3.177}$$

$$\mathbf{A}_I^{(k)} \leftarrow \mathbf{A}_I^{(k)} \circledast \mathbf{H}_I, \tag{3.178}$$

where

$$\mathbf{H}_C = \left( \mathbf{Y}_k \mathbf{X}_C^{(k)T} + \frac{\alpha}{\lambda} \sum_{j \neq k} \mathbf{A}_C^{(j)} \right) \oslash \left( \mathbf{A}^{(k)} \mathbf{X}^{(k)} \mathbf{X}_C^{(k)T} + \frac{\alpha}{\lambda} (L-1) \mathbf{A}_C^{(k)} + \frac{\gamma}{\lambda} \mathbf{A}_C^{(k)} \right), \tag{3.179}$$

$$\mathbf{H}_I = \left( \mathbf{Y}_k \mathbf{X}_I^{(k)T} + \frac{\beta}{\lambda} (L-1) \mathbf{A}_I^{(k)} \right) \oslash \left( \mathbf{A}^{(k)} \mathbf{X}^{(k)} \mathbf{X}_I^{(k)T} + \frac{\beta}{\lambda} \sum_{j \neq k} \mathbf{A}_I^{(j)} + \frac{\gamma}{\lambda} \mathbf{A}_I^{(k)} \right). \tag{3.180}$$

The solution can be obtained by the iteration to update individual parameter matrices until convergence.

## Chapter 4

# Supervised Feature Extraction and Classification methods for Pattern Recognition

In this Chapter, I introduce supervised signal processing methods for pattern recognition and point out their problems. Pattern recognition is a technique to recognize some patterns, automatically. For example, pattern recognition includes the face recognition, the handwritten character recognition, the fingerprinting identification, and the brain signal recognition.

### 4.1 Pattern Recognition

Flow of pattern recognition consists of pre-processing, feature extraction, and classification. The pre-processing includes cutout, normalization, filtering, and so on, and the feature extraction is to transform various types of data (tensor) to vector data. As an example, I introduce several supervised feature extraction techniques for the EEG based brain signal recognition in Sections 4.2. The classification is to assign an unknown pattern  $x$  into a category  $y$ . Fig 4.1 depicts the basic concept of pattern classification. Almost any data is allowed as a pattern  $x$  such as a discrete vector, a matrix of continuous values and a non-numeric data. The category  $y$  is a discrete value. I introduce several typical techniques of classification in Section 4.3. In this thesis, I focus on the classification methods, mainly.

#### 4.1.1 Heuristic Approach vs. Statistical Approach

There are two major approaches to find the rules for recognition, heuristic approach and statistical approach. In the heuristic approach, we find rules from a human subjective idea. For example, we know that apple and lemon can be classified by color, thus we implement a program so that the computer output 'apple' and 'lemon' if it is red and yellow, respectively. When computers did not

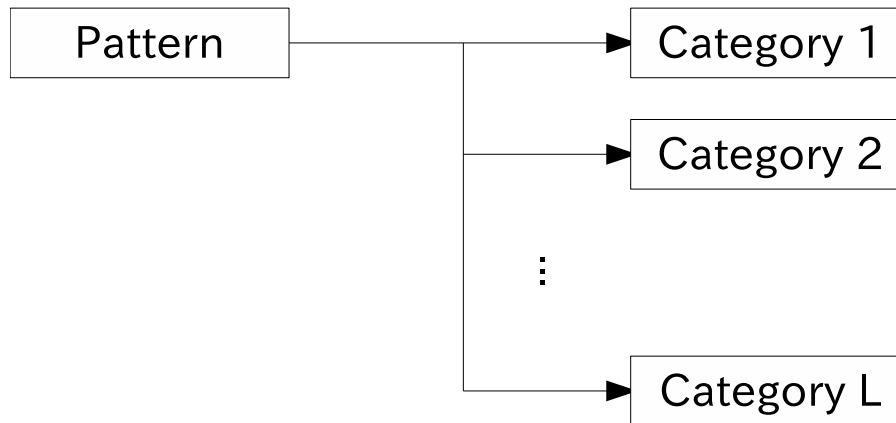


Figure 4.1: Basic concept of pattern classification

have a high computational capacity, the heuristic approach has been studied as a mainstream.

On the other hand, statistical approach finds the rules from observed pattern samples. In this case, pattern samples are very important since the classifier strongly depends on them. When we have enough number of samples, its classifier shows the excellent performance. With progress of computers, the statistical approach comes to be studied actively, and it became a mainstream. Statistical pattern recognition methods are based on supervised learning.

### 4.1.2 Supervised Learning

In this section, I explain the basic concept of supervised learning. We consider three components in supervised learning as follows [103].

- The generator of the data.
- The target operator (supervisor).
- The learning machine (classifier).

Fig 4.2 depicts the relation of three components. The generator outputs a random pattern vector  $\mathbf{x}$  that is independent and identically distributed (i.i.d.). The supervisor inputs  $\mathbf{x}$  and outputs its category  $y$ . But the supervisor can not realized in a computer. The learning machine has two functions. The first function inputs the training samples  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$  that are given from the generator and the supervisor. And a classifier is constructed to imitate or identify the supervisor by using the training samples. The second function inputs an unknown pattern  $\mathbf{x}$  and outputs its estimated category  $\hat{y}$ .

In the EEG based BCI system, we do not consider directly an EEG data as  $\mathbf{x}$  since the data size of an EEG matrix  $\mathbf{E}$  is very large which consists of many channels and many sampling points at every time. We consider the feature vector of an EEG matrix as  $\mathbf{x}$ . Thus, we first need a feature extraction method to transform EEG signal matrices to feature vectors. Fig 4.3 shows a standard feature extraction scheme for motor imagery based BCI from an EEG signal  $\mathbf{E}$  to a feature vector  $\mathbf{x}$ .

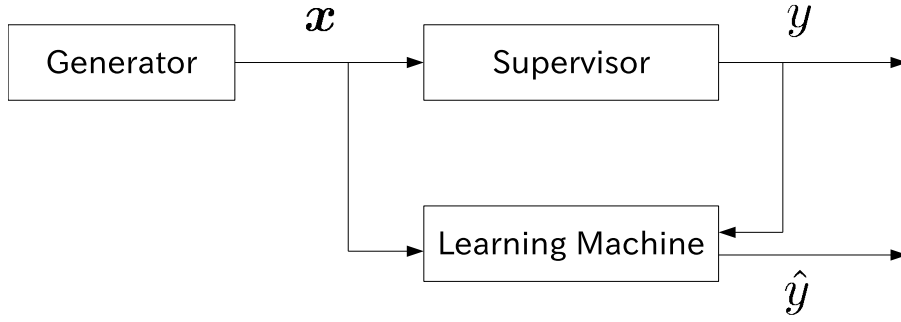


Figure 4.2: Basic concept of supervised learning

After the feature extraction, we shift to a classification procedure. Fig 4.4 shows a classification scheme. In order to classify a feature vector  $\mathbf{x}$  to a command  $y$ , a classification rule (i.e., classifier) is necessary. A classifier has to learn the classification rule from a training set  $[\mathbf{X}_{train}, \mathbf{y}_{train}]$ . Thus, a BCI system can be constructed.

## 4.2 Supervised Feature Extraction methods

In this section, I introduce a typical spatial filter and its extension methods for the EEG motor imagery classification. The proper spatial filter provides signals that are easily classified. The goal of this research is to design spatial filters that lead to better recognition rates.

EEG signals are formalized as

$$\{\mathbf{E}_n\}_{n=1}^N \in \mathbb{R}^{ch \times time}, \quad (4.1)$$

where  $ch$ ,  $time$ , and  $N$  are the number of channels, the number of sampling points of time, and the number of trials, respectively. To apply the EEG signals to classification, we have to transform the each EEG signal to a feature vector. Thus, some transformation from the EEG signal to the feature vector

$$\mathbf{E}_n \in \mathbb{R}^{ch \times time} \mapsto \mathbf{x}_n \in \mathbb{R}^d, \quad (4.2)$$

is necessary. Thus, we obtain a set of feature vectors  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{d \times N}$ .

The problem is how to extract important features for the best classification. The key-points here are

- EEG signals are very noisy (noise reduction),
- There exists the best frequency band for classification (frequency band selection),
- There exists important channels and un-important channels for classification (channel selection or spatial filters).

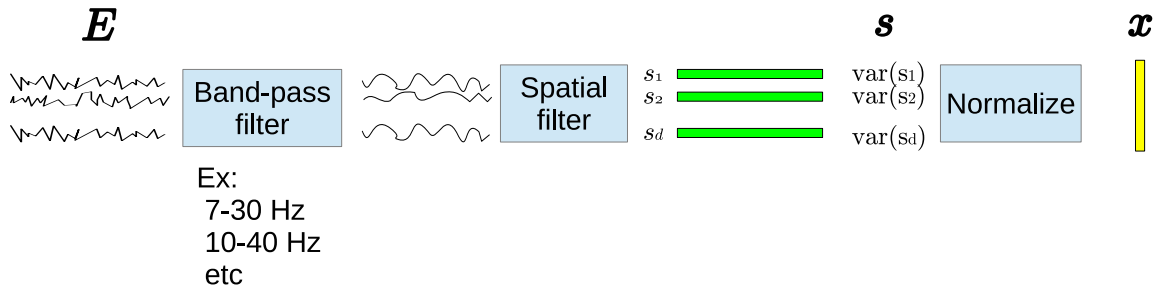


Figure 4.3: Standard feature extraction scheme

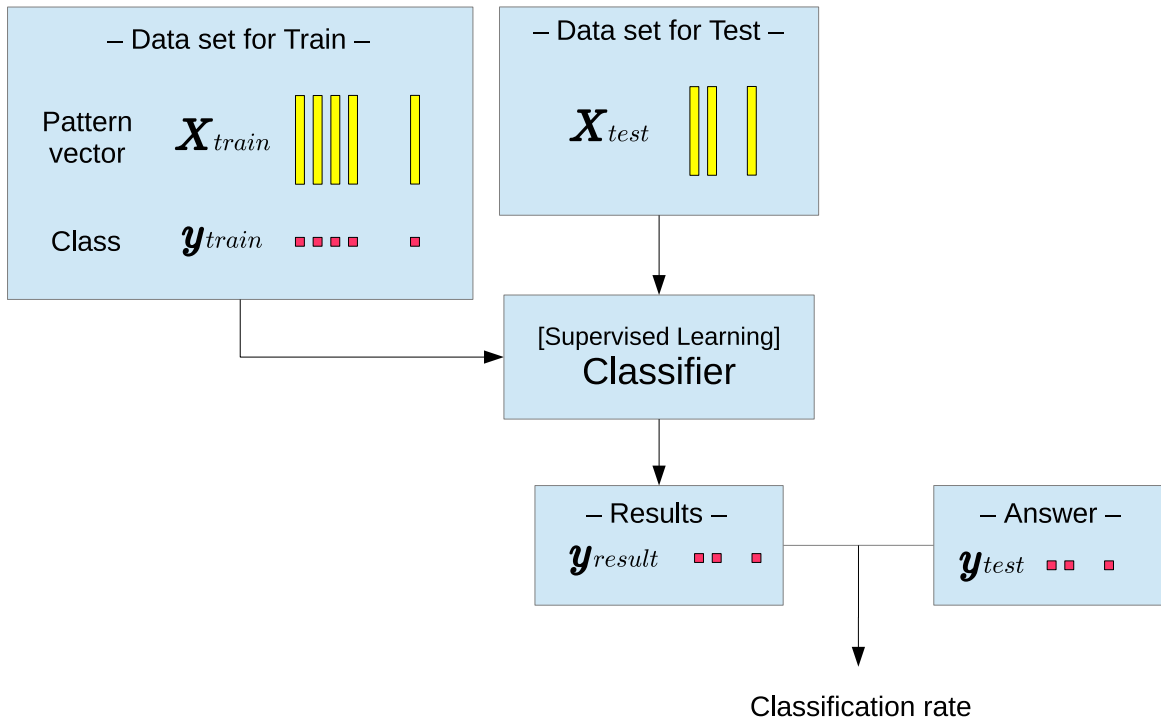


Figure 4.4: Classification scheme

### 4.2.1 CSP filter

The “Common Spatial Pattern” (CSP) filter [75] is a typical spatial filter for motor imagery classification by using many channels of EEG signals. We denote the CSP filter by

$$\mathbf{S} = \mathbf{W}^T \mathbf{E} \text{ or } s(t) = \mathbf{W}^T \mathbf{e}(t), \quad (4.3)$$

where  $\mathbf{W} \in \mathbb{R}^{d \times ch}$  is the spatial filter matrix,  $\mathbf{S} \in \mathbb{R}^{d \times time}$  is the filtered signal matrix.

The criterion of CSP is given by

$$\text{maximize } \text{tr} \mathbf{W}^T \boldsymbol{\Sigma}_1 \mathbf{W}, \quad (4.4)$$

$$\text{subject to } \mathbf{W}^T (\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2) \mathbf{W} = \mathbf{I}, \quad (4.5)$$

where

$$\mathbf{\Sigma}_1 = \text{Exp}_{\mathbf{E}_n \in \{\text{class 1}\}} \frac{\mathbf{E}_n \mathbf{E}_n^T}{\text{tr} \mathbf{E}_n \mathbf{E}_n^T}, \quad (4.6)$$

$$\mathbf{\Sigma}_2 = \text{Exp}_{\mathbf{E}_n \in \{\text{class 2}\}} \frac{\mathbf{E}_n \mathbf{E}_n^T}{\text{tr} \mathbf{E}_n \mathbf{E}_n^T}. \quad (4.7)$$

This problem can be solved by the generalized eigenvalue problem. However, we can also solve it by solving two standard eigenvalue problems.

First we decompose  $\mathbf{\Sigma}_1 + \mathbf{\Sigma}_2$  as

$$\mathbf{\Sigma}_1 + \mathbf{\Sigma}_2 = \mathbf{U} \mathbf{D} \mathbf{U}^T, \quad (4.8)$$

where  $\mathbf{U}$  is a set of eigenvectors, and  $\mathbf{D}$  is a diagonal matrix of eigenvalues.

Next, let  $\mathbf{P} := \sqrt{\mathbf{D}^{-1}} \mathbf{U}^T$ , and calculate

$$\widehat{\mathbf{\Sigma}}_1 = \mathbf{P} \mathbf{\Sigma}_1 \mathbf{P}^T, \quad (4.9)$$

$$\widehat{\mathbf{\Sigma}}_2 = \mathbf{P} \mathbf{\Sigma}_2 \mathbf{P}^T. \quad (4.10)$$

Please note that we have  $\widehat{\mathbf{\Sigma}}_1 + \widehat{\mathbf{\Sigma}}_2 = \mathbf{I}$ , here. Thus, any orthonormal matrices  $\mathbf{V}$  satisfy  $\mathbf{V}^T (\widehat{\mathbf{\Sigma}}_1 + \widehat{\mathbf{\Sigma}}_2) \mathbf{V} = \mathbf{I}$ .

Finally, we decompose  $\widehat{\mathbf{\Sigma}}_1$  as

$$\widehat{\mathbf{\Sigma}}_1 = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T, \quad (4.11)$$

where  $\mathbf{V}$  is a set of eigenvectors, and  $\mathbf{\Lambda}$  is a diagonal matrix of eigenvalues.

A set of CSP filters is obtained as

$$\mathbf{W} = \mathbf{P}^T \mathbf{V}. \quad (4.12)$$

We have

$$\mathbf{W}^T \mathbf{\Sigma}_1 \mathbf{W} = \mathbf{\Lambda} = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_{ch} \end{pmatrix}, \quad (4.13)$$

$$\mathbf{W}^T \mathbf{\Sigma}_2 \mathbf{W} = \mathbf{I} - \mathbf{\Lambda} = \begin{pmatrix} 1 - \lambda_1 & & \\ & \ddots & \\ & & 1 - \lambda_{ch} \end{pmatrix}, \quad (4.14)$$

where  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{ch}$ . Therefore, first CSP filter  $\mathbf{w}_1$  provides maximum variance for class 1, and last CSP filter  $\mathbf{w}_{ch}$  provides the maximum variance for class 2.

We select first and last  $m$  filters and construct a matrix  $\mathbf{W}$  as

$$\mathbf{W}_{csp} = \begin{pmatrix} \mathbf{w}_1 & \cdots & \mathbf{w}_m & \mathbf{w}_{ch-m+1} & \cdots & \mathbf{w}_{ch} \end{pmatrix} \in \mathbb{R}^{2m \times ch}, \quad (4.15)$$

and the filtered signal matrix  $\mathbf{s}(t)$  is given by

$$\mathbf{s}(t) = \mathbf{W}_{csp}^T \mathbf{e}(t) = \begin{pmatrix} s_1(t) & \cdots & s_d(t) \end{pmatrix}^T, \quad (4.16)$$

where  $d = 2m$ .

Finally, the feature vector  $\mathbf{x} = (x_1, x_2, \dots, x_d)^T$  is calculated by

$$x_i = \log \left( \frac{\text{var}[s_i(t)]}{\sum_{i=1}^d \text{var}[s_i(t)]} \right). \quad (4.17)$$

### Example

As an example, I show feature vectors for the motor imagery EEG dataset [8]. The motor imagery EEG is closely related to primary motor cortex in brain (see Fig 4.5). Fig. 4.6 shows topographs of CSP weight at right hand and foot imagery. We can see that CSP filter provides optimal weightings so that left area is weighted for right hand imagery and central area is weighted for foot imagery. Right hand imagery pattern and foot imagery pattern of EEG signals can be easily separated by using CSP filter unlike PCA (see Fig. 4.7).

### 4.2.2 CSSP filter

The Common Spatio-Spectral Pattern (CSSP) filter is an extension of the CSP filter [64]. The CSSP can be regarded as a CSP method with the time delay embedding.

The algorithm is similar to that of the standard CSP. In CSP, we consider the following transformation

$$\mathbf{S} = \mathbf{W}^T \mathbf{E} \quad \text{or} \quad \mathbf{s}(t) = \mathbf{W}^T \mathbf{e}(t). \quad (4.18)$$

However, the CSSP's transform is given by

$$\mathbf{S} = \mathbf{W}^T \mathbf{E} + \mathbf{W}_\tau^T \mathbf{E}_\tau = \widehat{\mathbf{W}}^T \begin{pmatrix} \mathbf{E} \\ \mathbf{E}_\tau \end{pmatrix}, \quad (4.19)$$

$$\text{or} \quad \mathbf{s}(t) = \mathbf{W}^T \mathbf{e}(t) + \mathbf{W}_\tau^T \mathbf{e}(t + \tau) = \widehat{\mathbf{W}}^T \begin{pmatrix} \mathbf{e}(t) \\ \mathbf{e}(t + \tau) \end{pmatrix}, \quad (4.20)$$

where  $\mathbf{E}_\tau$  is a  $\tau$ -time delayed signal matrix of  $\mathbf{E}$ , and  $\widehat{\mathbf{W}}^T = [\mathbf{W}^T, \mathbf{W}_\tau^T]$  is a CSSP matrix. It can be also regarded that the number of channels increases to double. The difference between CSP and CSSP is only this point, and then we can apply this method easily in a same way as the CSP algorithm. However,  $\tau$  is a hyper-parameter.

The key-point here is this method can be interpreted into a spatial and a spectral filters. Therefore let  $\widehat{\mathbf{w}}$  denote the  $i$ -th column of the CSSP matrix  $\widehat{\mathbf{W}}$ , then the projected signal  $s(t)$  can be expressed

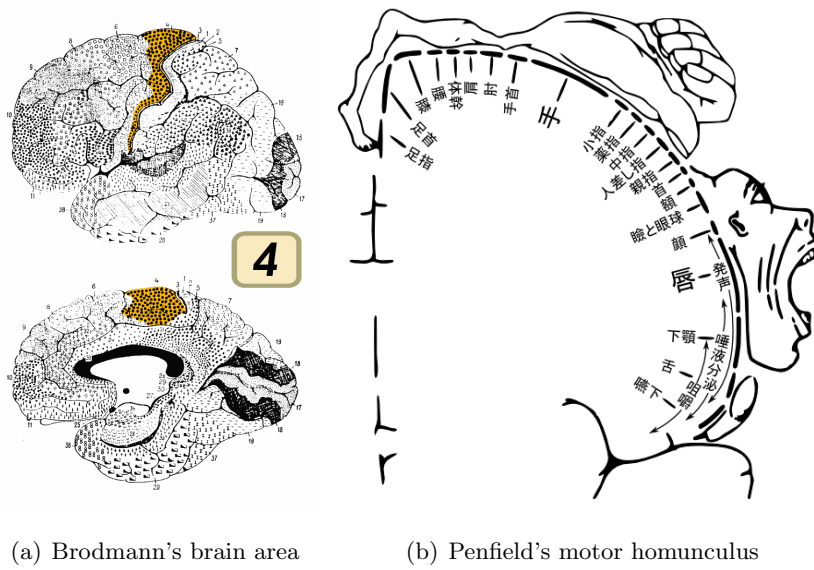


Figure 4.5: Primary motor cortex

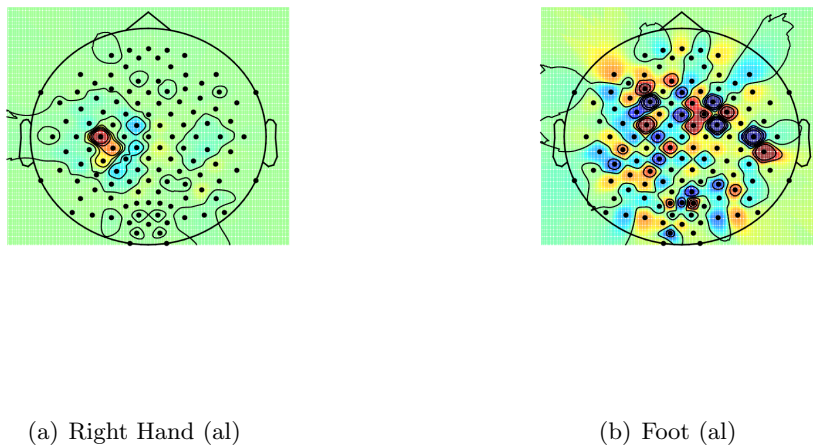


Figure 4.6: Topographies of CSP filter: BCI competition III(IVa)

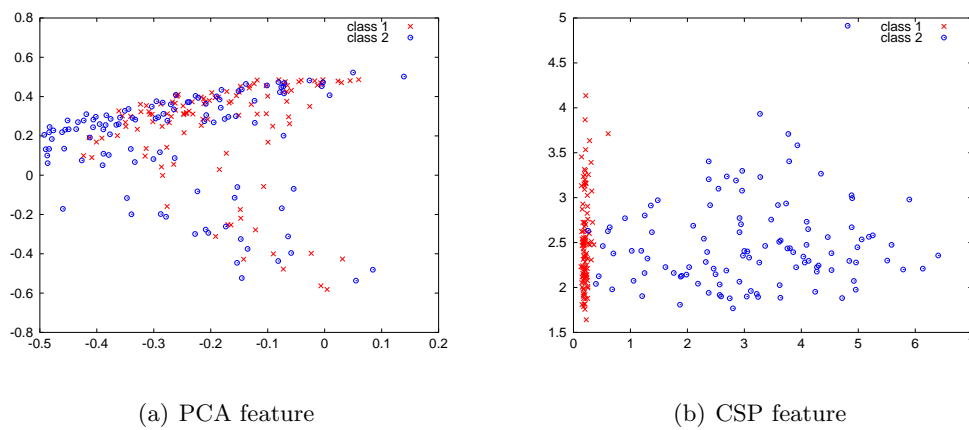


Figure 4.7: 2D-plot of CSP feature: classes 1 and 2 stand for right hand and foot imageries, respectively. Horizontal and vertical axes in CSP feature space are first and last components, respectively.

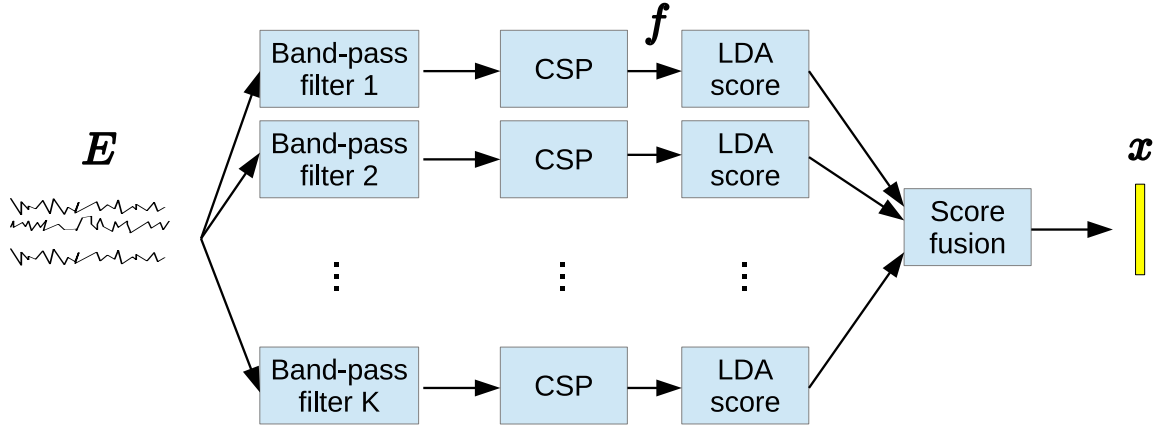


Figure 4.8: System Flowchart of SBCSP: for example frequency bands of individual filters are 4-8Hz, 8-12Hz, ... , 36-40Hz.

as

$$s(t) = \mathbf{w}^T \mathbf{e}(t) + \mathbf{w}_\tau^T \mathbf{e}(t + \tau) \quad (4.21)$$

$$= \sum_{j=1}^{ch} w_j e_j(t) + (w_\tau)_j e_j(t + \tau) \quad (4.22)$$

$$= \sum_{j=1}^{ch} \gamma_j \left( \frac{w_j}{\gamma_j} e_j(t) + \frac{(w_\tau)_j}{\gamma_j} e_j(t + \tau) \right), \quad (4.23)$$

where

- $\gamma_j$  can be regarded as a pure spatial filter,
- $\left[ \frac{w_j}{\gamma_j}, 0, \dots, 0, \frac{(w_\tau)_j}{\gamma_j} \right]$  can be regarded as a finite impulse response (FIR) filter for each channel.

Since the CSSP filter includes a temporal filter, it can be considered that the CSSP becomes robust for noise or outliers.

### 4.2.3 SBCSP filter

Here, I introduce an alternative method based on Sub-band CSP (SBCSP) and score fusion [77]. Fig. 4.8 shows a flowchart of feature extraction by using the SBCSP filter. How to decide a frequency band of band-pass filter is a difficult problem of CSP/CSSP filter; however, the SBSCP can automatically decide the frequency band of the band-pass filter. First, we make candidates of frequency bands, for example 4-8 Hz, 8-12 Hz, 12-16 Hz, ..., 36-40 Hz. Next we use all candidates of frequency bands for band-pass filter, and obtain CSP features from each filtered signals. Thus, we have individual CSP feature data matrices  $\{\mathbf{F}^{(1)}, \mathbf{F}^{(2)}, \dots, \mathbf{F}^{(K)}\}$ . Next, the process of LDA score is given by follow:

- Calculate the LDA projection vector  $\hat{\mathbf{w}}^{(k)}$  for each data matrices  $\mathbf{F}^{(k)}$ .

- LDA score vector is given by

$$\mathbf{s}_n = \begin{pmatrix} \hat{\mathbf{w}}^{(1)T} \mathbf{f}_n^{(1)} \\ \hat{\mathbf{w}}^{(2)T} \mathbf{f}_n^{(2)} \\ \vdots \\ \hat{\mathbf{w}}^{(K)T} \mathbf{f}_n^{(K)} \end{pmatrix}. \quad (4.24)$$

There are two approaches of score fusion method.

- Recursive Band Elimination (RBE)
- Meta-Classifier (MC)

In this way, we do not need to select frequency bands in the SBCSP by hand, thus it can be regarded as an automatically selection of the frequency band.

### Recursive Band Elimination

In RBE, first we set a RBE-order that is an integer number  $r \in [1, K]$ . The first data set is  $\mathbf{X} = [\mathbf{s}_1, \dots, \mathbf{s}_N] \in \mathbb{R}^{K \times N}$ . In this method, we remove  $(K - r)$  rows from  $\mathbf{X}$  by support vector machine (SVM) feature selection (SVM is explained in Section 4.3.5). The algorithm is as follow

1. Train the parameter  $\mathbf{w}$  in linear model  $\langle \mathbf{w}, \mathbf{x} \rangle$  by SVM
2. Remove the row of  $\mathbf{X}$  with the smallest  $w_k^2$
3. If number of rows of  $\mathbf{X}$  becomes  $r$ , this algorithm is finished, else back to 1.

The survival data  $\mathbf{X}$  can be used as a training data matrix.

### Meta-Classifier

In MC, we use a Bayesian classifier for score fusion. We assume that each band score  $s_k \in \{\text{class 1}\}$  and  $s_k \in \{\text{class 2}\}$  are distributed normally. And we can estimate individual parameters of normal distributions (i.e.,  $\{\mu_1^{(k)}, \sigma_1^{(k)}\}$  and  $\{\mu_2^{(k)}, \sigma_2^{(k)}\}$ ). We calculate the training data matrix by

$$x_{kn} = \log \left\{ \frac{p(s_{kn} | \mu_1^{(k)}, \sigma_1^{(k)})}{p(s_{kn} | \mu_2^{(k)}, \sigma_2^{(k)})} \right\}, \quad (4.25)$$

where  $x_{kn}$  is the  $(k, n)$ -element of  $\mathbf{X}$ .

## 4.3 Classification methods

In this section, I introduce practical classification methods.

### 4.3.1 Generative Model vs Discriminant Model

For classification, it is very important to decide the model of classifier. Models of classifier can be divided into either generative model or discriminant model.

Generative model is based on probability such that we can classify any patterns  $\mathbf{x}$  into the category  $\hat{y}$  based on a posterior probability  $P(y|\mathbf{x})$  as

$$\hat{y} = \operatorname{argmax}_y P(y|\mathbf{x}). \quad (4.26)$$

It is called maximum a *posteriori* (MAP) estimation. In terms of recognition rate, the MAP estimation is the best estimator. For these reasons, it is very famous and extensively studied. A key point of MAP estimation is how to estimate the posterior probabilities accurately.

The discriminant model is based on some discriminant functions  $D_y(\mathbf{x})$  as

$$D_y : \mathbb{R}^d \rightarrow \mathbb{R} \quad (4.27)$$

$$\hat{y} = \operatorname{argmax}_y D_y(\mathbf{x}). \quad (4.28)$$

If we assume that  $D_y(\mathbf{x}) := P(y|\mathbf{x})$ , we can see that the discriminant model includes the generative model. However, their approaches are different. Generally, we define some non-negative-valued function  $R(\mathbf{D}(\mathbf{x}), y)$  named error function, where  $\mathbf{D}(\mathbf{x})$  is a vector-valued function of which element is  $D_y(\mathbf{x})$ . The discriminant functions are estimated by minimizing the expectation of error function as follows:

$$\operatorname{minimize} \sum_y \int_{\mathcal{D}} P(\mathbf{x}, y) R(\mathbf{D}(\mathbf{x}), y) d\mathbf{x}. \quad (4.29)$$

When we have training samples, it can be approximated by replacing ensemble mean by sample mean as

$$\sum_y \int_{\mathcal{D}} P(\mathbf{x}, y) R(\mathbf{D}(\mathbf{x}), y) d\mathbf{x} \approx \frac{1}{N} \sum_{n=1}^N R(\mathbf{D}(\mathbf{x}_n), y_n). \quad (4.30)$$

Then we can obtain the discriminant functions by minimizing this formula if the definition of error function and the model of discriminant functions are given.

### Generative models

In this section, I introduce the generative model based classification methods. In general, its approach based on Bayes theorem:

$$P(y|\mathbf{x}) = \frac{P(y)P(\mathbf{x}|y)}{P(\mathbf{x})}. \quad (4.31)$$

If we estimate each probabilities of right side, then any  $\mathbf{x}$  can be classified by the MAP rule. Since  $P(y|\mathbf{x}) \propto P(y)P(\mathbf{x}|y)$ , we can assign any  $\mathbf{x}$  into  $\hat{y}$  by  $\operatorname{argmax}_y P(y)f_y(\mathbf{x}|\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$ .  $P(y)$  can be estimated as the appearance ratio of individual  $y$ . Thus, the key-point here is to estimate  $P(\mathbf{x}|y)$  by

some generative model function. There are three approaches to estimate the generative model: the parametric model, the non-parametric model, and the semi-parametric model.

The parametric model is generally given by  $f_y(\mathbf{x}; \Theta_y)$ , where  $\Theta_y$  is a set of parameters. For example, when we assume  $f_y$  follows the Gaussian distribution, then it is given by

$$f_y(\mathbf{x}|\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y) = \frac{1}{(\sqrt{2\pi})^d \sqrt{|\boldsymbol{\Sigma}_y|}} \exp \left[ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_y)^T \boldsymbol{\Sigma}_y^{-1} (\mathbf{x} - \boldsymbol{\mu}_y) \right], \quad (4.32)$$

where  $\Theta_y = \{\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y\}$ . In this case, we only need to estimate a mean parameter vector  $\boldsymbol{\mu}_y$  and a covariance parameter matrix  $\boldsymbol{\Sigma}_y$  for each  $y$ . Maximum likelihood estimation is often used to estimate these parameters.

The non-parametric model does not have any parameter to estimate; but it has some hyper-parameters. Kernel density estimation is often used for non-parametric estimation.

The semi-parametric model can be considered as an intermediate model between parametric and non-parametric models. Gaussian mixture model is often used, and it is often obtained by EM-algorithm.

### Discriminant models

There are several discriminant models. Most basic model is a linear model which is given by

$$D(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + \theta, \quad (4.33)$$

where  $\mathbf{w}$  and  $\theta$  are the model parameters, and  $\langle \cdot, \cdot \rangle$  stands for the inner-product. It can be expanded by using basis functions  $\boldsymbol{\phi}(\mathbf{x})$  as

$$D(\mathbf{x}) = \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}) \rangle + \theta. \quad (4.34)$$

Note that, the dimensions of parameter vectors  $\mathbf{w}$  in Eq.(4.33) and in Eq.(4.34) are different. The linear model is often used for theoretical study of classifier. For example, perceptron learning [88], optimal hyperplane classifier (OHC) [103] and Fisher discriminant analysis (FDA) [72] are using the linear model.

Next I introduce the quadratic model as

$$D(\mathbf{x}) = \langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle + \langle \mathbf{b}, \mathbf{x} \rangle + \theta, \quad (4.35)$$

where  $\mathbf{A}$ ,  $\mathbf{b}$  and  $\theta$  are parameters. Its typical example is given by the class-featuring information compression (CLAFIC) method [104].

In general, if the order of the model or the number of parameters is higher, we can construct a more complex classifier. However, if it is too high, the classifier over-fits to the training samples and its generalization capability is decreased. Model selection and regularization techniques are very important for classification to prevent the over-fitting. Especially, the cross validation method is very famous and useful for model selection. I explain its details in section 4.3.3.

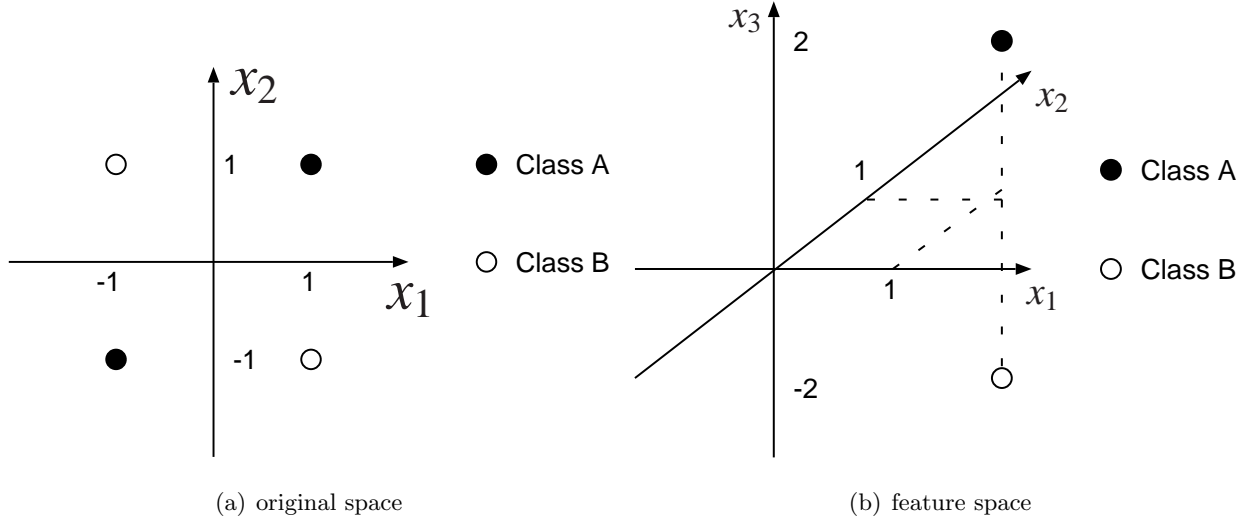


Figure 4.9: Example of training samples

### 4.3.2 Kernel Method

There are many cases such that training samples can not be separated linearly. In such cases, it is very efficient to map pattern vectors to a higher dimensional space. For example, we consider a mapping  $\phi$  as follows:

$$\begin{aligned} \phi : \mathbb{R}^2 &\rightarrow \mathbb{R}^3, \\ (t_1, t_2)^T &\mapsto (t_1^2, t_2^2, 2t_1t_2)^T, \end{aligned} \tag{4.36}$$

where  $(\cdot)^T$  denotes the transposition. This mapping is based on correlations of each element. We assume that there are four training samples  $\{(\mathbf{x}_n, y_n)\}_{n=1}^4$  with categories  $y \in \{A, B\}$  as follow:

$$\mathbf{x}_1 = (1, 1)^T, \quad y_1 = A, \tag{4.37}$$

$$\mathbf{x}_2 = (-1, -1)^T, \quad y_2 = A, \tag{4.38}$$

$$\mathbf{x}_3 = (1, -1)^T, \quad y_3 = B, \tag{4.39}$$

$$\mathbf{x}_4 = (-1, 1)^T, \quad y_4 = B. \tag{4.40}$$

Fig 4.9(a) depicts the training samples in a 2-dimensional original Euclid space. In this case, we can not separate the training samples into 2 categories linearly. However, if we consider the mapping (4.36), it can be separated linearly. Feature vectors as mapped samples are given as

$$\phi(\mathbf{x}_1) = (1, 1, 2)^T, \quad y_1 = A, \tag{4.41}$$

$$\phi(\mathbf{x}_2) = (1, 1, 2)^T, \quad y_2 = A, \tag{4.42}$$

$$\phi(\mathbf{x}_3) = (1, 1, -2)^T, \quad y_3 = B, \tag{4.43}$$

$$\phi(\mathbf{x}_4) = (1, 1, -2)^T, \quad y_4 = B. \tag{4.44}$$

Fig 4.9(b) depicts the training samples in the 3-dimensional feature space. In this space, we can separate the training samples linearly.

It often happens that the performance of classifier is improved by mapping data into the high dimensional feature space. In this way, classification capability can be improved by using a high-dimensional feature mapping  $\phi$ ; however,  $\phi(\mathbf{x})$  has a very large dimension and its computational cost is also very high. Fortunately, we can solve this problem by using the kernel method.

In the kernel method, instead of calculating the high-dimensional feature vector, we calculate the inner-product between vectors in feature space by using a kernel function:

$$k(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle. \quad (4.45)$$

In such case, we have to show that a mapping  $\phi$  that satisfies Eq.(4.45) exists. If the kernel function  $k(\mathbf{x}, \mathbf{y})$  satisfies Mercer's theorem, the existence of  $\phi$  is guaranteed. Mercer's theorem is described as follows:

**Theorem 1** (Mercer [103]). *Let  $\mathcal{X}$  be a pattern space,  $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ , and  $\mathcal{L}_2(\mathcal{X})$  be a set of all quadratically integrable functions, to guarantee that a continuous symmetric function  $k(\mathbf{x}, \mathbf{y})$  in  $\mathcal{L}_2(\mathcal{X})$  has an expansion*

$$k(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{\infty} a_i \phi_i(\mathbf{x}) \phi_i(\mathbf{y}) \quad (4.46)$$

with positive coefficients  $a_i > 0$  (i.e.,  $k(\mathbf{x}, \mathbf{y})$  describes an inner product in some feature space), it is necessary and sufficient that the condition

$$\int_{\mathcal{X}} \int_{\mathcal{X}} k(\mathbf{x}, \mathbf{y}) f(\mathbf{x}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} \geq 0 \quad (4.47)$$

be valid for all  $f \in \mathcal{L}_2(\mathcal{X})$   $\square$

Now, we define the kernel discriminant function by using sample patterns  $\{\mathbf{x}_n\}_{n=1}^N$  as

$$D(\mathbf{x}) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x}). \quad (4.48)$$

This discriminant function model is called the kernel model. We can handle classification problems in a high dimensional feature space without the calculation of mapping by using the kernel model. This technique is called "kernel trick". For examples, the following kernel functions satisfy Mercer's theorem.

$$\begin{aligned} k(\mathbf{x}, \mathbf{y}) &= \langle \mathbf{x}, \mathbf{y} \rangle && : \text{Linear kernel} \\ k(\mathbf{x}, \mathbf{y}) &= (\langle \mathbf{x}, \mathbf{y} \rangle + 1)^d && : \text{Polynomial kernel} \\ k(\mathbf{x}, \mathbf{y}) &= \exp(-\gamma \|\mathbf{x} - \mathbf{y}\|^2) && : \text{Gaussian kernel} \end{aligned} \quad (4.49)$$

Since  $k(\mathbf{x}_n, \mathbf{x}) = \langle \phi(\mathbf{x}_n), \phi(\mathbf{x}) \rangle$ , the discriminant function is given by

$$D(\mathbf{x}) = \sum_{n=1}^N \alpha_n \langle \phi(\mathbf{x}_n), \phi(\mathbf{x}) \rangle \quad (4.50)$$

$$= \left\langle \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n), \phi(\mathbf{x}) \right\rangle. \quad (4.51)$$

Then, this classifier in the feature space is expressed as

$$D(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle, \quad (4.52)$$

$$\mathbf{w} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n). \quad (4.53)$$

From Eq.(4.53), the parameter vector  $\mathbf{w}$  is confined in a subspace  $[\phi(\mathbf{x}_n)]_{n=1}^N$  that is spanned by feature vectors  $\{\phi(\mathbf{x}_n)\}_{n=1}^N$ . It is a basic property of kernel models that  $\|\mathbf{w}\|^2$  is expressed as

$$\|\mathbf{w}\|^2 = \left\langle \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n), \sum_{m=1}^N \alpha_m \phi(\mathbf{x}_m) \right\rangle \quad (4.54)$$

$$= \sum_{n=1}^N \alpha_n \sum_{m=1}^N \alpha_m k(\mathbf{x}_n, \mathbf{x}_m) \quad (4.55)$$

$$= \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha}, \quad (4.56)$$

where

$$\boldsymbol{\alpha} := (\alpha_1 \quad \alpha_2 \quad \cdots \quad \alpha_N)^T, \quad (4.57)$$

$$\mathbf{K} := \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}. \quad (4.58)$$

### 4.3.3 Least Square Regression

In this section, I explain the least square regression (LSR). LSR is not a classification method but a regression method for general purposes, however we use LSR as a classifier here. And we also assume that the discriminant function is given by

$$D(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle. \quad (4.59)$$

Furthermore, the training samples are given by  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ , and we consider a binary classification problem as  $y \in \{+1, -1\}$ .

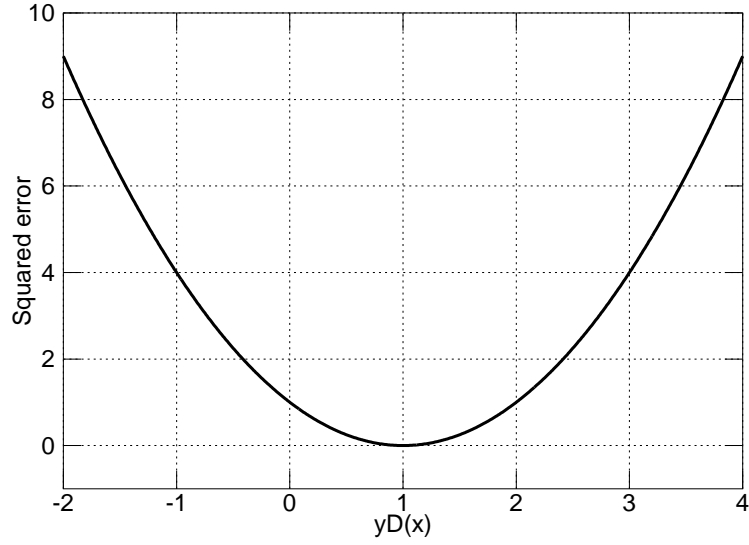


Figure 4.10: Squared error

The basic idea of LSR is to minimize the expectation of the squared error  $(y - \langle \mathbf{w}, \mathbf{x} \rangle)^2$ . Since  $y^2 = 1$ , the error function is given by

$$R(D(\mathbf{x}), y) = (y - \langle \mathbf{w}, \mathbf{x} \rangle)^2 \quad (4.60)$$

$$\begin{aligned} &= y^2 - 2y \langle \mathbf{w}, \mathbf{x} \rangle + \langle \mathbf{w}, \mathbf{x} \rangle^2 \\ &= 1 - 2y \langle \mathbf{w}, \mathbf{x} \rangle + y^2 \langle \mathbf{w}, \mathbf{x} \rangle^2 \\ &= (1 - y \langle \mathbf{w}, \mathbf{x} \rangle)^2. \end{aligned} \quad (4.61)$$

Fig 4.10 depicts the squared error function based on Eq.(4.61). Note that the squared error is a convex function. The expectation of the squared error is given by

$$E[(y - \langle \mathbf{w}, \mathbf{x} \rangle)] = \sum_y \int_{\mathcal{D}} P(\mathbf{x}, y) (y - \langle \mathbf{w}, \mathbf{x} \rangle)^2 d\mathbf{x}. \quad (4.62)$$

It can be approximated by replacing ensemble mean by sample mean as

$$E[(y - \langle \mathbf{w}, \mathbf{x} \rangle)] \approx \frac{1}{N} \sum_{n=1}^N (y_n - \langle \mathbf{w}, \mathbf{x}_n \rangle)^2 =: J. \quad (4.63)$$

It can be transformed as

$$J = \frac{1}{N} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \quad (4.64)$$

$$= \frac{1}{N} (\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}), \quad (4.65)$$

where

$$\mathbf{y} := \begin{pmatrix} y_1 & y_2 & \cdots & y_N \end{pmatrix}^T, \quad (4.66)$$

$$\mathbf{X} := \begin{pmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_N \end{pmatrix}^T. \quad (4.67)$$

The goal of LSR is to calculate the optimal  $\mathbf{w}$  that minimizes Eq.(4.65). Then, its solution  $\hat{\mathbf{w}}_{\text{LSR}}$  can be calculated analytically as follow:

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{1}{N}(-2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \mathbf{w}) = \mathbf{0}, \quad (4.68)$$

$$\hat{\mathbf{w}}_{\text{LSR}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (4.69)$$

We can classify an unlearned pattern  $\mathbf{x}$  based on the following rule:

$$\hat{y} = \begin{cases} +1 & \text{if } \langle \hat{\mathbf{w}}_{\text{LSR}}, \mathbf{x} \rangle > 0 \\ -1 & \text{if } \langle \hat{\mathbf{w}}_{\text{LSR}}, \mathbf{x} \rangle < 0 \end{cases}. \quad (4.70)$$

### Kernelized LSR

In this section, I explain LSR with the kernel model. We can adapt the kernel model to LSR easily by replacing  $\mathbf{X}\mathbf{w}$  by  $\mathbf{K}\boldsymbol{\alpha}$ . Then,  $J$  is given by

$$J = \frac{1}{N}(\mathbf{y} - \mathbf{K}\boldsymbol{\alpha})^T(\mathbf{y} - \mathbf{K}\boldsymbol{\alpha}) \quad (4.71)$$

$$= \frac{1}{N}(\mathbf{y}^T \mathbf{y} - 2\boldsymbol{\alpha}^T \mathbf{K} \mathbf{y} + \boldsymbol{\alpha}^T \mathbf{K}^2 \boldsymbol{\alpha}). \quad (4.72)$$

Its solution  $\hat{\boldsymbol{\alpha}}_{\text{LSR}}$  can be calculated analytically as follow:

$$\frac{\partial J}{\partial \boldsymbol{\alpha}} = \frac{1}{N}(-2\mathbf{K} \mathbf{y} + 2\mathbf{K}^2 \boldsymbol{\alpha}) = \mathbf{0}, \quad (4.73)$$

$$\hat{\boldsymbol{\alpha}}_{\text{LSR}} = \mathbf{K}^{-1} \mathbf{y}. \quad (4.74)$$

### Regularization

Fig 4.11 is depicts the discriminant functions trained by LSR and kernelized LSR. In LSR, model function is too simple to fit to training samples. On the other hand, function with kernel model is too free, and fits to training samples too precisely. Its problem is called ‘‘over-fitting’’. Regularization can prevent the over-fitting.

Adding  $\lambda \|\mathbf{w}\|^2$  to  $J$  is often used for regularization, where  $\lambda > 0$  is called the regularization parameter. Since  $\|\mathbf{w}\|^2 = \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha}$  and  $\frac{1}{N}$  can be abbreviated,  $J$  is given by

$$J = (\mathbf{y} - \mathbf{K}\boldsymbol{\alpha})^T(\mathbf{y} - \mathbf{K}\boldsymbol{\alpha}) + \lambda \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} \quad (4.75)$$

$$= (\mathbf{y}^T \mathbf{y} - 2\boldsymbol{\alpha}^T \mathbf{K} \mathbf{y} + \boldsymbol{\alpha}^T \mathbf{K}(\mathbf{K} + \lambda \mathbf{I})\boldsymbol{\alpha}), \quad (4.76)$$

where  $\mathbf{I}$  is the identity matrix. Then its solution is given by

$$\frac{\partial J}{\partial \boldsymbol{\alpha}} = \frac{1}{N}(-2\mathbf{K} \mathbf{y} + 2\mathbf{K}(\mathbf{K} + \lambda \mathbf{I})\boldsymbol{\alpha}) = \mathbf{0}, \quad (4.77)$$

$$\hat{\boldsymbol{\alpha}}_{\text{LSR}} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}. \quad (4.78)$$

Fig 4.12 depicts the discriminant function trained by LSR with this regularization method. We can see that the over-fitting is prevented.

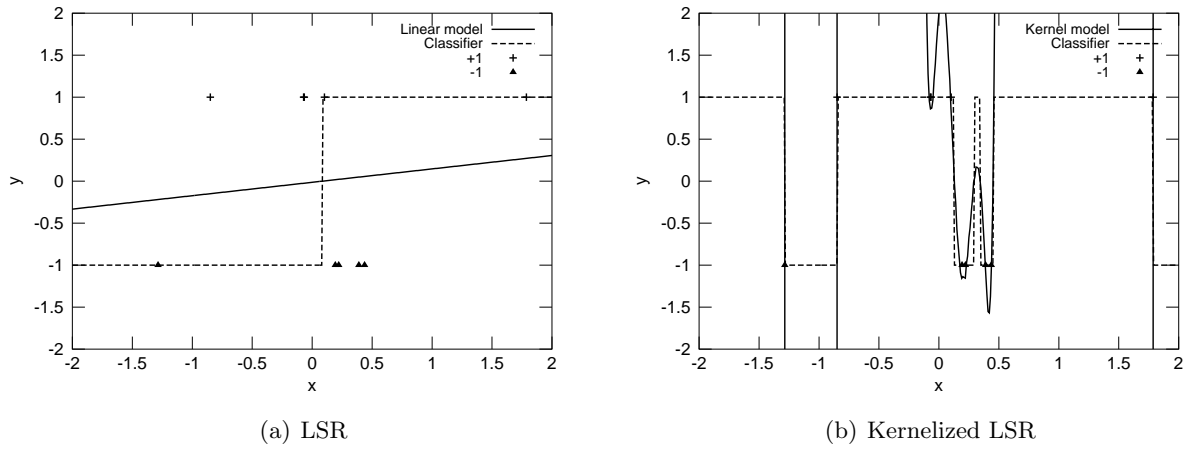


Figure 4.11: Examples of LSR and kernelized LSR

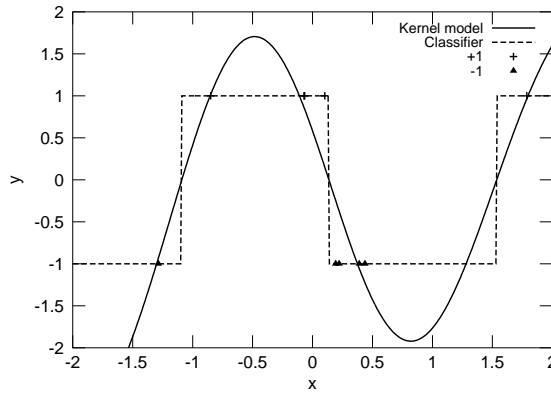


Figure 4.12: Regularization

### Model Selection

Many kernel functions have kernel parameters such as  $d$  of the polynomial kernel and  $\gamma$  of the Gaussian kernel. Such a parameter is called the model parameter too. If the model parameter is chosen improperly, it degrades the classifiers. The task of the model is to estimate the optimal model parameter such that the best classifier is obtained. Here, I introduce the cross validation method [7] that is very famous and effective.

For the cross validation, we separate training samples into a training set and a test set, and to select the optimal parameter based on its accuracy. For example, we divide samples into 4 parts and there are 4 combinations to separate them into 3 parts of training set and 1 part of test set. Fig 4.13 depicts the 4-fold cross validation method. This method is generalized as the  $k$ -fold cross validation. Especially, when we assume that the number of samples is  $N$ , a method that separates them into  $(N - 1)$  samples in training set and 1 samples in test set is called the leave-one-out cross validation.

These methods are useful to select not only the model parameter, but also the regularization pa-

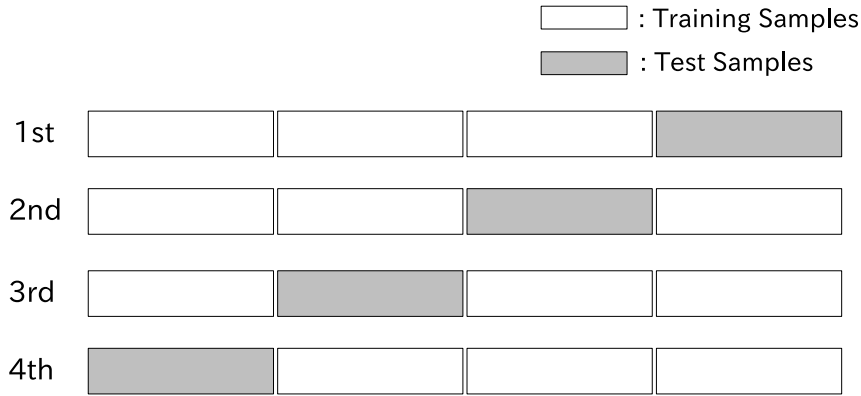


Figure 4.13: 4-fold cross validation

parameter. But the computational complexity becomes much higher.

#### 4.3.4 Fisher Linear Discriminant Analysis

I introduce the Fisher's linear discriminant analysis (FDA). The FDA is a very famous binary classification method. It is based on mean vectors and covariance matrices of patterns for individual classes. We consider the linear model as  $z = \mathbf{w}^T \mathbf{x}$ . The FDA gives an optimal projection  $\mathbf{w}$  so that the distribution of  $z$  is easily discriminated. Fig. 4.14 shows the concept of FDA, and there are two projection vectors  $\mathbf{w}$  and  $\mathbf{w}^*$ . We can see that  $\mathbf{w}^*$  is clearly easier to discriminate patterns than  $\mathbf{w}$ . This concept can be provided by the following FDA criterion

$$\text{maximize } J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1 + s_2}, \quad (4.79)$$

where

- $m_1$  and  $m_2$  denote averages for  $z_n \in \{ \text{class 1} \}$  and  $z_n \in \{ \text{class 2} \}$ , respectively,
- $s_1$  and  $s_2$  denote variances for  $z_n \in \{ \text{class 1} \}$  and  $z_n \in \{ \text{class 2} \}$ , respectively.

We have

$$\begin{aligned} (m_1 - m_2)^2 &= (\mathbf{w}^T \boldsymbol{\mu}_1 - \mathbf{w}^T \boldsymbol{\mu}_2) (\mathbf{w}^T \boldsymbol{\mu}_1 - \mathbf{w}^T \boldsymbol{\mu}_2)^T \\ &= \mathbf{w}^T (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{w} = \mathbf{w}^T \mathbf{S}_B \mathbf{w}, \end{aligned} \quad (4.80)$$

and

$$\begin{aligned} s_1 + s_2 &= \mathbf{w}^T \boldsymbol{\Sigma}_1 \mathbf{w} + \mathbf{w}^T \boldsymbol{\Sigma}_2 \mathbf{w} \\ &= \mathbf{w}^T (\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2) \mathbf{w} = \mathbf{w}^T \mathbf{S}_W \mathbf{w}, \end{aligned} \quad (4.81)$$

where

- $\boldsymbol{\mu}_1$  and  $\boldsymbol{\mu}_2$  denote mean vectors for  $\mathbf{x}_n \in \{ \text{class 1} \}$  and  $\mathbf{x}_n \in \{ \text{class 2} \}$ , respectively,

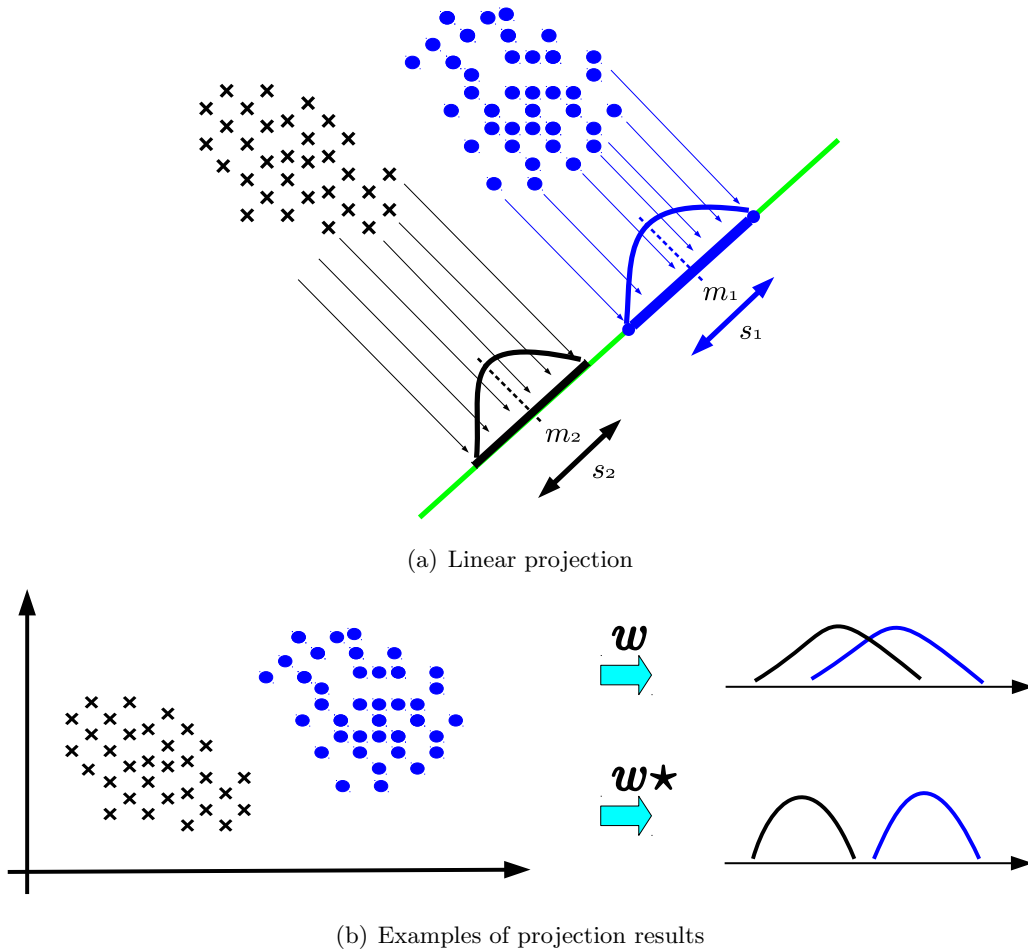


Figure 4.14: Fisher's Linear Discriminant Analysis

- $\Sigma_1$  and  $\Sigma_2$  denote covariance matrices for  $\mathbf{x}_n \in \{ \text{class 1} \}$  and  $\mathbf{x}_n \in \{ \text{class 2} \}$ , respectively.

The cost function  $J(\mathbf{w})$  in original space can be written as

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}, \tag{4.82}$$

and then this solution is given by

$$\hat{\mathbf{w}} \propto \mathbf{S}_W^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2). \tag{4.83}$$

Finally, we choose an optimal threshold  $z_0$ , then we can classify any  $\mathbf{x}$  by

$$\hat{\mathbf{w}}^T \mathbf{x} \geq z_0 \rightarrow \mathbf{x} \in \{ \text{class 1} \}, \tag{4.84}$$

$$\hat{\mathbf{w}}^T \mathbf{x} < z_0 \rightarrow \mathbf{x} \in \{ \text{class 2} \}. \tag{4.85}$$

For example,  $z_0 = (m_1 + m_2)/2$  could be usable.

### Problem of FDA

In this section, I explain that FDA does not provide the best optimal hyper-plane in heteroscedastic Gaussian case. We denote the probability density functions of  $y$  which belongs to class 1 and 2 by

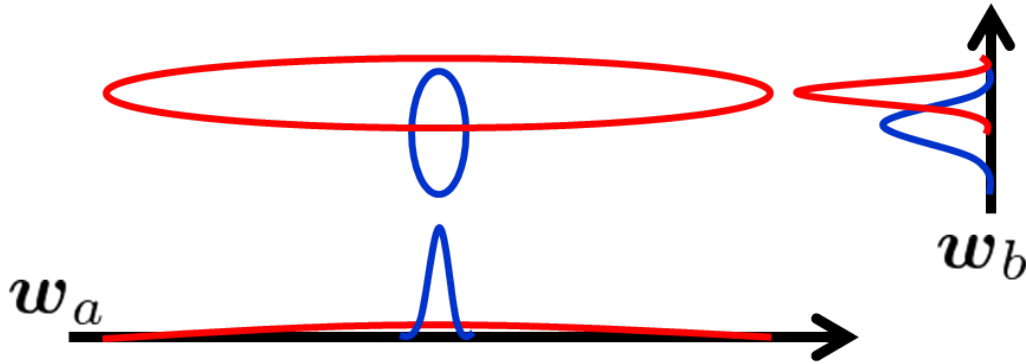


Figure 4.15: Example of which FDA does not work well

$f_1(y)$  and  $f_2(y)$ , and the prior probabilities of class 1 and 2 by  $p_1$  and  $p_2$ , respectively. We assume  $f_1(y)$  and  $f_2(y)$  follows the Gaussian distributions. Then, the Bayes decision rule for the minimum error is given by

$$\mathbf{x} \text{ is assigned to } \begin{cases} \Omega_1 & \text{if } p_1 f_1(\langle \mathbf{w}, \mathbf{x} \rangle) > p_2 f_2(\langle \mathbf{w}, \mathbf{x} \rangle) \\ \Omega_2 & \text{if } p_1 f_1(\langle \mathbf{w}, \mathbf{x} \rangle) < p_2 f_2(\langle \mathbf{w}, \mathbf{x} \rangle) \end{cases} \quad (4.86)$$

Fig. 4.15 shows a simple example of which FDA does not work well. Red and blue ellipses show individual Gaussian distributions in this figure. We consider two projection vectors  $\mathbf{w}_a$  and  $\mathbf{w}_b$ . The FDA evaluation function of  $\mathbf{w}_a$  becomes  $J(\mathbf{w}_a) = 0$  which is the lowest value; then, FDA has to select  $\mathbf{w}_b$ , preferentially. However, the classification rate by  $\mathbf{w}_a$  is better than  $\mathbf{w}_b$  under the classification rule (4.86). In this way, the improvement of evaluation function is necessary for the heteroscedastic Gaussian case.

### 4.3.5 Support Vector Machine

Support vector machine (SVM) is a typical kernel classifier and it is an extension of the optimal hyperplane classifier (OHC) by the kernel method. [102, 101] OHC and SVM are trained by quadratic programming. In this section, I explain the OHC at first. After that, I explain the SVM as an extension of OHC.

#### Optimal Hyperplane classifier

We assume a binary classification problem as  $y \in \{+1, -1\}$ . The discriminant function of OHC is given by

$$D(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + \theta. \quad (4.87)$$

Fig. 4.16 depicts an example of OHC in 2-dimensional space. The boundary of two categories is given by a  $(d - 1)$ -dimensional hyperplane described as  $D(\mathbf{x}) = 0$ , where  $d$  is the dimension of  $\mathbf{x}$ . In this

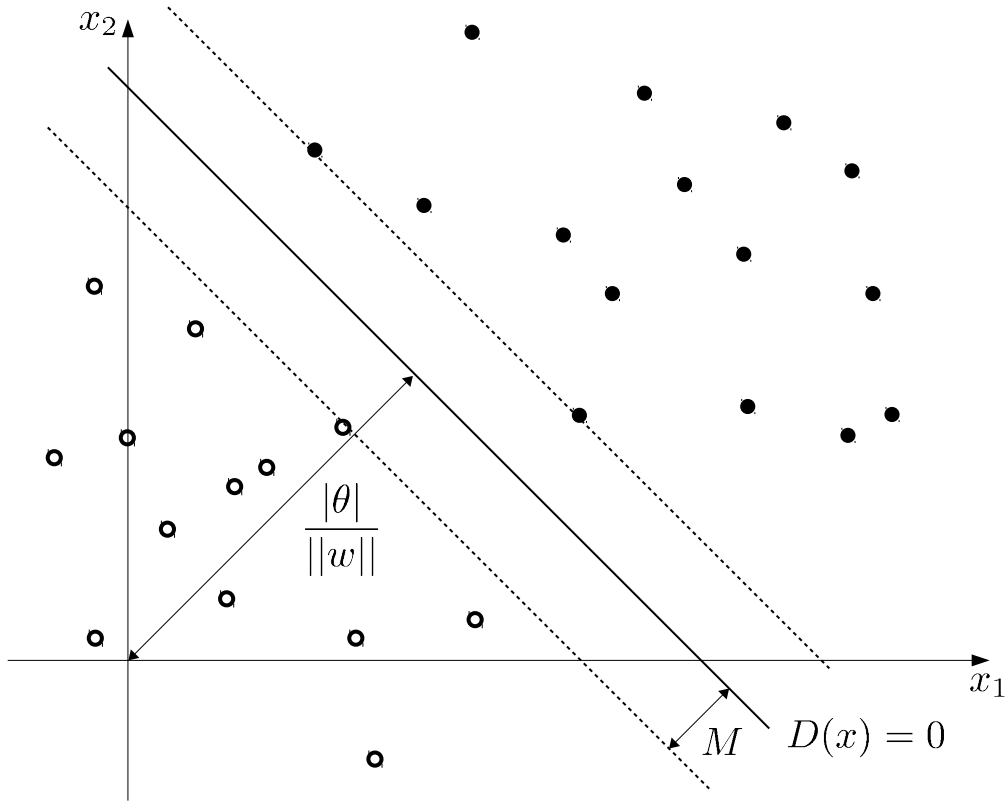


Figure 4.16: Basic concept of optimal hyperplane classifier

figure,  $M$  stands for “margin” that is the minimum distance between the boundary and its nearest samples. Basic concept of OHC is to maximize “margin” that is described by “ $M$ ” in Fig. 4.16. The margin is given by

$$M = \min_{n=1, \dots, N} \frac{y_n D(\mathbf{x}_n)}{\|\mathbf{w}\|}. \quad (4.88)$$

Furthermore, we normalize the numerator of Eq. (4.88) as

$$\min_{n=1, \dots, N} y_n D(\mathbf{x}_n) = 1, \quad (4.89)$$

then the margin is given by

$$M = \frac{1}{\|\mathbf{w}\|}. \quad (4.90)$$

From Eqs.(4.89) and (4.90), OHC is trained by the following quadratic programming problem,

$$\text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 \quad (4.91)$$

$$\text{subject to } y_n (\langle \mathbf{w}, \mathbf{x}_n \rangle + \theta) \geq 1, \quad n = 1, \dots, N. \quad (4.92)$$

We can solve it directly, however it can be simplified by transforming to a dual problem [68]. We define Lagrange's function  $L$  as

$$L(\mathbf{w}, \theta, \boldsymbol{\gamma}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N \gamma_n (y_n (\langle \mathbf{w}, \mathbf{x}_n \rangle + \theta) - 1), \quad (4.93)$$

where  $\gamma_n \geq 0$  are Lagrange coefficients. Lagrange primal and dual problems are expressed as

$$P : \min_{\mathbf{w}, \theta} \sup_{\boldsymbol{\gamma} \geq 0} L(\mathbf{w}, \theta, \boldsymbol{\gamma}), \quad (4.94)$$

$$D : \max_{\boldsymbol{\gamma} \geq 0} \inf_{\mathbf{w}, \theta} L(\mathbf{w}, \theta, \boldsymbol{\gamma}), \quad (4.95)$$

respectively. The two problems are equivalent. From the Karuch-Kuhn-Tucker (KKT) conditions we have

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \mathbf{X}^T \mathbf{Y} \boldsymbol{\gamma} = \mathbf{0}, \quad \frac{\partial L}{\partial \theta} = \boldsymbol{\gamma}^T \mathbf{y} = 0, \quad (4.96)$$

$$\gamma_n (y_n (\langle \mathbf{w}, \mathbf{x}_n \rangle + \theta) - 1) = 0, \quad (4.97)$$

$$y_n (\langle \mathbf{w}, \mathbf{x}_n \rangle + \theta) - 1 \geq 0, \quad \gamma_n \geq 0 \quad n = 1, \dots, N, \quad (4.98)$$

where

$$Y := \text{diag}(\mathbf{y}). \quad (4.99)$$

Then, the dual problem is given by

$$\text{minimize} \quad -\frac{1}{2} \boldsymbol{\gamma}^T \mathbf{Y} \mathbf{X} \mathbf{X}^T \mathbf{Y} \boldsymbol{\gamma} + \mathbf{1}^T \boldsymbol{\gamma} \quad (4.100)$$

$$\text{subject to} \quad \boldsymbol{\gamma}^T \mathbf{y} = 0, \quad \gamma_n \geq 0, \quad n = 1, \dots, N, \quad (4.101)$$

where  $\mathbf{1}$  is an  $N$ -dimensional vector of which all elements are one. Its solution  $\hat{\boldsymbol{\gamma}}$  is almost always sparse. Then, let  $S$  be the set of indexes of non-zero  $\hat{\gamma}_n$  as

$$S = \{n | \gamma_n \neq 0, \quad n = 1, \dots, N\}. \quad (4.102)$$

The solution of the original problem is described as

$$\hat{\mathbf{w}} = \sum_{n=1}^N y_n \hat{\gamma}_n \mathbf{x}_n = \sum_{s \in S} y_s \hat{\gamma}_s \mathbf{x}_s, \quad (4.103)$$

$$\hat{\theta} = \frac{1}{|S|} \sum_{s \in S} (\hat{\mathbf{w}}^T \mathbf{x}_s - y_s). \quad (4.104)$$

Eqs. (4.103) and (4.104) imply the optimal parameters depend only on the support vectors.

This method is also called the hard margin method. In this method, training samples that can not be separated linearly are not allowed. It is a serious issue, so that the soft margin method was proposed.

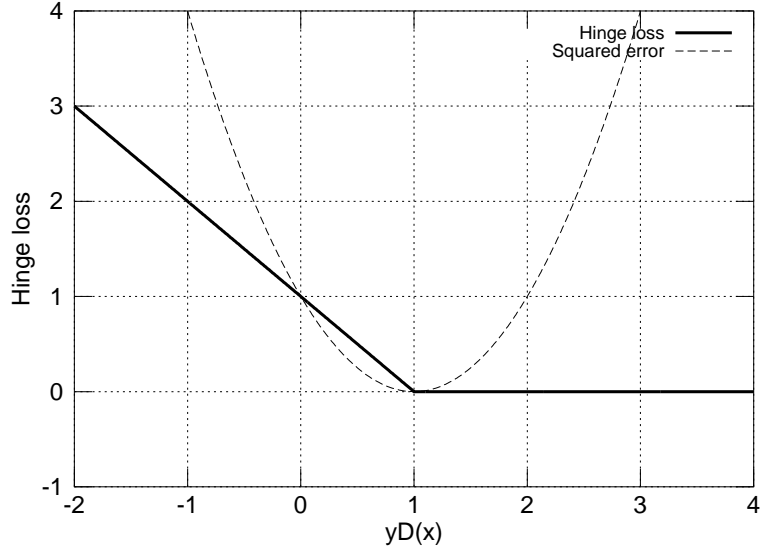


Figure 4.17: Hinge loss

### Soft margin

For the soft margin, we introduce slack variables  $\xi_n \geq 0$  to relent the constraints Eq. (4.92) as follow:

$$y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + \theta) \geq 1 - \xi_n \quad n = 1, \dots, N. \quad (4.105)$$

The variable  $\xi_n$  expresses how much condition (4.92) is not satisfied.  $\xi_n$  are called the hinge loss and it should be minimized as well as  $\|\mathbf{w}\|$ .  $\xi_n$  is also given by

$$\xi_n = \max(1 - y_n D(\mathbf{x}_n), 0). \quad (4.106)$$

Fig. 4.17 depicts the comparison between the hinge loss function and the squared loss function. The hinge loss is a piecewise linear function, and note that the hinge loss is a convex function. Since the hinge loss is evaluated linearly, then it is more robust for outlier than the squared error.

The OHC with soft margin is trained by the following quadratic programming:

$$\text{minimize} \quad \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \xi_i \quad (4.107)$$

$$\text{subject to} \quad y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + \theta) \geq 1 - \xi_n \quad n = 1, \dots, N, \quad (4.108)$$

where  $\lambda$  is a constant parameter to adjust weights between the hinge loss term and the margin because both terms have to be minimized. The margin term  $\|\mathbf{w}\|^2$  and  $\lambda$  are regarded as a regularization term and a regularization parameter.

Its dual problem is given by

$$\text{minimize} \quad -\frac{1}{2\lambda} \boldsymbol{\gamma}^T \mathbf{Y} \mathbf{X} \mathbf{X}^T \mathbf{Y} \boldsymbol{\gamma} + \mathbf{1}^T \boldsymbol{\gamma} \quad (4.109)$$

$$\text{subject to} \quad 0 \leq \gamma_n \leq 1, \quad \mathbf{y}^T \boldsymbol{\gamma} = 0, \quad n = 1, \dots, N, \quad (4.110)$$

where  $\boldsymbol{\gamma}$  is a dual parameter vector. Then, the solution is calculated by

$$\hat{\boldsymbol{w}} = \frac{1}{\lambda} \sum_{s \in S} y_s \hat{\gamma}_s \boldsymbol{x}_s, \quad (4.111)$$

$$\hat{\theta} = \frac{1}{|S|} \sum_{s \in S} (\hat{\boldsymbol{w}}^T \boldsymbol{x}_s - y_s). \quad (4.112)$$

### Support Vector Machine

Support Vector Machine (SVM) is a kernelized version of OHC. That is the OHC defined in high dimensional feature space.

Since  $D(\boldsymbol{x}) = \sum_{n=1}^N \alpha_n k(\boldsymbol{x}_n, \boldsymbol{x}) + \theta$  and  $\|\boldsymbol{w}\| = \boldsymbol{\alpha}^T \boldsymbol{K} \boldsymbol{\alpha}$ , SVM is trained by the following quadratic programming:

$$\text{minimize} \quad \frac{\lambda}{2} \boldsymbol{\alpha}^T \boldsymbol{K} \boldsymbol{\alpha} + \sum_{i=1}^n \xi_i \quad (4.113)$$

$$\text{subject to} \quad y_n \left( \sum_{m=1}^N \alpha_m k(\boldsymbol{x}_m, \boldsymbol{x}_n) + \theta \right) \geq 1 - \xi_n \quad n = 1, \dots, N. \quad (4.114)$$

And its dual problem is given by

$$\text{minimize} \quad -\frac{1}{2\lambda} \boldsymbol{\gamma}^T \boldsymbol{K} \boldsymbol{\gamma} + \mathbf{1}^T \boldsymbol{\gamma} \quad (4.115)$$

$$\text{subject to} \quad 0 \leq \gamma_n \leq 1, \quad \boldsymbol{y}^T \boldsymbol{\gamma} = 0, \quad n = 1, \dots, N, \quad (4.116)$$

where  $\boldsymbol{\gamma}$  is a dual parameter vector. Its solution is calculated by

$$\hat{\boldsymbol{\alpha}} = \frac{1}{\lambda} \boldsymbol{Y} \hat{\boldsymbol{\gamma}}, \quad (4.117)$$

$$\hat{\theta} = \frac{1}{|S|} \sum_{s \in S} \left( \sum_{n=1}^N \alpha_n k(\boldsymbol{x}_n, \boldsymbol{x}_s) - y_s \right). \quad (4.118)$$

## Chapter 5

# Proposed Feature Extraction methods

In this Chapter, I propose several new methods for matrix and tensor based feature extractions. This chapter consists of three parts, mainly. The first part is a proposition of new methods and algorithms for smooth nonnegative matrix and tensor factorizations. I propose a new fast algorithm for the GRBF-NMF, several extension methods, and its tensor versions with the CP and the Tucker models. The second part is a proposition of new model and algorithms for the common and individual feature analysis (CIFA) for the multi-block CP decomposition. I call this method as the “linked CP tensor decomposition” (LCPTD). Furthermore, I propose the sparse LCPTD and the nonnegative LCPTD. The third part is a proposition of new model and algorithms for the common and individual feature analysis (CIFA) for the multi-block Tucker decomposition. I call this method as the “linked Tucker decomposition” (LTD). Furthermore, I propose the orthogonal LTD, the sparse LTD, and the nonnegative LTD.

### 5.1 Improvements of NMF with function approximation

In this section, I propose an improved algorithm for function approximation based smooth NMF (GRBF-NMF) introduced in Section 3.1.6. Furthermore, I propose some extended methods and tensor versions.

#### 5.1.1 Fast algorithm for GRBF-NMF

First, I modify the GRBF-NMF optimization problem (3.95) as

$$\begin{aligned} & \underset{\mathbf{W}, \mathbf{X}}{\text{minimize}} \quad \|\mathbf{Y} - \Phi \mathbf{W} \mathbf{X}\|_F^2, \\ & \text{s.t.} \quad \mathbf{W} \geq 0, \mathbf{X} \geq 0, \|\mathbf{x}_r\|^2 = 1 \text{ for } r = 1, \dots, R, \end{aligned} \quad (5.1)$$

where  $\mathbf{x}_r^T$  is the  $r$ -th row vector of  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_R]^T$ . In this problem, the constraint  $\Phi \mathbf{W} \geq 0$  is replaced by  $\mathbf{W} \geq 0$ . Since  $\Phi \geq 0$ ,  $\mathbf{W} \geq 0$  is a sufficient condition for  $\Phi \mathbf{W} \geq 0$ . Under this constraint, the flexibility of model will be decreased; however it is not necessarily bad. In many regression models,

$$\begin{array}{l}
\text{a)} \quad \boxed{\Phi} \quad \boxed{W} \quad \boxed{X} = \boxed{\Phi} \left\| \begin{array}{c} \overline{x_1^T} \\ w_1 \end{array} \right. + \dots + \boxed{\Phi} \left\| \begin{array}{c} \overline{x_R^T} \\ w_R \end{array} \right. \\
\quad \quad \quad (I \times N) \quad (N \times R) \quad (R \times J) \\
\text{b)} \quad \boxed{W} = \left\| \begin{array}{c} w_1 \\ \vdots \\ w_R \end{array} \right. \\
\quad \quad \quad (N \times R) \\
\text{c)} \quad \boxed{X} = \left\| \begin{array}{c} \overline{x_1^T} \\ \vdots \\ \overline{x_R^T} \end{array} \right. \\
\quad \quad \quad (R \times J)
\end{array}$$

Figure 5.1: Decomposition of GRBF-NMF model

a highly flexible model without any regularization has a problem of over-fitting. When the data is noisy, such a highly flexible model often represents not only the main feature but also noise and outliers. In order to prevent the over-fitting, one of the solutions is to use a simpler model, and the proposed model is simpler than the original model since I do not allow negative value for  $W$ . Thus, the proposed model can be considered as more robust for noise and outliers than the original model. In addition, I impose a constraint  $\|\mathbf{x}_r\|^2 = 1$  which does not change the flexibility. This constraint becomes a key-point for the new fast algorithm.

In order to obtain the solution of (5.1), I consider to separate (5.1) into sub-problems based on the HALS method [24]. Since the GRBF-NMF model can be decomposed as  $\Phi W X = \Phi w_1 x_1^T + \Phi w_2 x_2^T + \dots + \Phi w_R x_R^T$  (See Fig 5.1), the problem (5.1) can be separated as the following sub-problems:

$$\underset{\mathbf{x}_r}{\text{minimize}} \quad \|\mathbf{Y}_r - \Phi \mathbf{w}_r \mathbf{x}_r^T\|_F^2, \quad \text{s.t. } \mathbf{x}_r \geq 0, \quad \|\mathbf{x}_r\|^2 = 1. \quad (5.2)$$

$$\underset{\mathbf{w}_r}{\text{minimize}} \quad \|\mathbf{Y}_r - \Phi \mathbf{w}_r \mathbf{x}_r^T\|_F^2, \quad \text{s.t. } \mathbf{w}_r \geq 0, \quad (5.3)$$

where  $\mathbf{Y}_r := \mathbf{Y} - \sum_{k \neq r} \Phi \mathbf{w}_k \mathbf{x}_k^T$ .

First I consider the problem (5.2). Since  $\Phi \mathbf{w}_r$  is currently fixed and it is equivalent to a simple least squares problem with a nonnegativity constraint. Thus, its problem has been already studied in many papers such as [19, 24, 26]. Then, the update rule is given by

$$\mathbf{x}_r \leftarrow [\mathbf{Y}_r^T \Phi \mathbf{w}_r]_+, \quad (5.4)$$

$$\mathbf{x}_r \leftarrow \mathbf{x}_r / \|\mathbf{x}_r\|, \quad (5.5)$$

where  $[x]_+ := \max(x, \varepsilon)$ , and  $\varepsilon$  is a very small positive value (e.g.,  $\varepsilon = 10^{-16}$ ).

Next, I consider the problem (5.3). The objective function can be transformed as

$$\begin{aligned}
& \|\mathbf{Y}_r - \Phi \mathbf{w}_r \mathbf{x}_r^T\|_F^2 \\
& = \text{tr}(\mathbf{Y}_r^T \mathbf{Y}_r) - 2\text{tr}(\mathbf{Y}_r^T \Phi \mathbf{w}_r \mathbf{x}_r^T) + \text{tr}(\mathbf{w}_r^T \Phi^T \Phi \mathbf{w}_r)
\end{aligned} \quad (5.6)$$

since  $\mathbf{x}_r^T \mathbf{x}_r = 1$ . By computing the gradient of the objective function, the condition for a solution is

given by

$$\frac{\partial}{\partial \mathbf{w}_r} \|\mathbf{Y}_r - \Phi \mathbf{w}_r \mathbf{x}_r^T\|_F^2 = 2\Phi^T \Phi \mathbf{w}_r - 2\Phi^T \mathbf{Y}_r \mathbf{x}_r = 0. \quad (5.7)$$

This direct solution will be  $\mathbf{w}_r = [(\Phi^T \Phi)^{-1} \Phi^T \mathbf{Y}_r \mathbf{x}_r]_+$ ; however, the inverse problem is unstable and high-cost. Thus, I do not employ this update rule, but propose a more effective iterative update rule. Note that the following problem is equivalent to (5.3) via the same differential since  $\|\mathbf{x}_r\| = 1$ :

$$\underset{\mathbf{w}_r}{\text{minimize}} \|\mathbf{Y}_r \mathbf{x}_r - \Phi \mathbf{w}_r\|^2, \text{ s.t. } \mathbf{w}_r \geq 0. \quad (5.8)$$

This problem can be assigned to NNLS (nonnegative least squares). Finally, I employ the following multiplicative update rule:

$$\mathbf{w}_r \leftarrow [\mathbf{w}_r \circledast (\Phi^T \mathbf{Y}_r \mathbf{x}_r) \oslash (\Phi^T \Phi \mathbf{w}_r)]_+, \quad (5.9)$$

where  $\circledast$  and  $\oslash$  are respectively element-wise multiplication and element-wise division.

The algorithm is summarized in Algorithm 18. I call this method the fastGRBF-NMF. The computational cost for each iteration of the fastGRBF-NMF is very low because it consists of only four operations and a thresholding without any complex optimization procedure. If we assume that  $N = I$  and  $I > J > R$ , the computational complexity of one update for  $\mathbf{W}$  in the fastGRBF-NMF algorithm is  $\mathcal{O}(I^2 R)$ . On the other hand, if we assume the QP optimization in the GRBF-NMF algorithm includes the Cholesky decomposition, its computational complexity is  $\mathcal{O}(I^3 R^3)$  because the number of parameters are  $IR$ .

The computational efficiency is often considered as a more important factor than the strictness of solution, especially for a large-scale problem. Since the NMF solution is not unique, the most important viewpoint is a balance between computational efficiency and a meaningful solution. In Section 7.1, we can confirm that the proposed method provides not only the improvement of computational efficiency but also meaningful results.

### 5.1.2 Extensive Bases for the fastGRBF-NMF

In this section, I discuss how to choose  $\Phi$  more efficiently. I propose to use various bases  $\phi$  to make  $\Phi$ , for example multiple bases with various  $\sigma$  can be mixed. Thus, I propose to construct  $\Phi$  by

$$\Phi = [\Phi_{\sigma_1}, \Phi_{\sigma_2}, \dots, \Phi_{\sigma_U}] \in \mathbb{R}^{I \times N}, \quad (5.10)$$

where  $U$  is the number of  $\sigma_u$  and  $\Phi_{\sigma_u}$  stands for the basis matrix with a standard deviation  $\sigma_u$ . For example, if we set

$$\sigma_1 = \sigma, \sigma_2 = 2\sigma, \sigma_3 = 4\sigma, \dots, \sigma_U = 2^{U-1}\sigma, \quad (5.11)$$

$$\Delta t_1 = \delta t, \Delta t_2 = 2\delta t, \Delta t_3 = 4\delta t, \dots, \Delta t_U = 2^{U-1}\delta t, \quad (5.12)$$

then  $N < 2I/\delta t$  holds for any  $U$ . Let us put the horizontal size of  $\Phi_{\sigma_1}$  as  $N_0$ , then the horizontal size of  $\Phi$  is roughly bounded as  $N = N_0 + \frac{1}{2}N_0 + \frac{1}{4}N_0 + \dots + \frac{1}{2^{U-1}N_0} < 2N_0$  for any  $U$ . Thus, the

**Algorithm 18** FastGRBF-HALS algorithm

---

```

1: Input:  $\mathbf{Y}$ ,  $R$ , and  $\Phi$ 
2: Initialize:  $\mathbf{W}$  and  $\mathbf{X}$ 
3:  $\mathbf{E} = \mathbf{Y} - \Phi \mathbf{W} \mathbf{X}$  ;
4: repeat
5:   for  $r = 1, \dots, R$  do
6:      $\mathbf{Y}_r \leftarrow \mathbf{E} + (\Phi \mathbf{w}_r) \mathbf{x}_r^T$  ;
7:      $\mathbf{x}_r \leftarrow [\mathbf{Y}_r^T (\Phi \mathbf{w}_r)]_+$  ;
8:      $\mathbf{x}_r \leftarrow \mathbf{x}_r / \|\mathbf{x}_r\|$  ;
9:      $\mathbf{w}_r \leftarrow [\mathbf{w}_r \circledast \{\Phi^T (\mathbf{Y}_r \mathbf{x}_r)\} \oslash \{\Phi^T (\Phi \mathbf{w}_r)\}]_+$  ;
10:     $\mathbf{E} \leftarrow \mathbf{Y}_r - (\Phi \mathbf{w}_r) \mathbf{x}_r^T$ 
11:   end for
12: until  $\|\mathbf{E}\|_F^2$  converges
13: Output:  $\mathbf{W}$  and  $\mathbf{X}$ 

```

---

increase of complexity (which is double at the maximum) is not a large problem since the proposed algorithm is well improved to be faster by the HALS based method. Furthermore, I propose to add an additional DC component, then  $\Phi$  is given by

$$\Phi = [\Phi_{\sigma_1}, \Phi_{\sigma_2}, \dots, \Phi_{\sigma_U}, \mathbf{1}]. \quad (5.13)$$

Finally, we have  $N = \sum_{u=1}^U [\lfloor (I-1)/(2^{u-1}\delta t) \rfloor + 1] + 1$ . By this extension we can express the data with various resolutions.

**2D-basis**

I propose another kind of the Gaussian basis function for smooth image signals. I assume that the observed signal  $\mathbf{y}_j \in \mathbb{R}^I$  is an unfolded vector of an  $(I_1 \times I_2)$ -matrix  $\mathbf{Y}_j \in \mathbb{R}^{I_1 \times I_2}$ , where  $I = I_1 I_2$ . In this case, if we use  $\phi_n$  in (3.94), it represents only vertical smoothness. However, natural images usually have vertical and horizontal smoothnesses. Thus I consider the following 2D-basis:

$$\Phi_n^{(2)}(i_1, i_2) := \exp \left[ \frac{1}{2\sigma^2} \left\| \begin{pmatrix} i_1 \\ i_2 \end{pmatrix} - \begin{pmatrix} n_1 \\ n_2 \end{pmatrix} \right\|^2 \right] \in \mathbb{R}^{I_1 \times I_2}, \quad (5.14)$$

where

$$n_1 = \{(n\Delta t - 1) \bmod I_1\} + 1, \quad (5.15)$$

$$n_2 = \lfloor (n\Delta t - 1)/I_1 \rfloor + 1. \quad (5.16)$$

Thus, we have  $n\Delta t = (n_2 - 1)I_1 + n_1$ ,  $1 \leq n_1 \leq I_1$ , and  $1 \leq n_2 \leq I_2$ . In practice, I unfold the matrix  $\Phi_n^{(2)}$  into a vector as  $\phi_n^{(2)} = \text{vec}(\Phi_n^{(2)}) \in \mathbb{R}^I$ , and construct the basis matrix  $\Phi^{(2)}$  as

$$\Phi^{(2)} = [\phi_1^{(2)}, \phi_2^{(2)}, \dots, \phi_N^{(2)}] \in \mathbb{R}^{I \times N}. \quad (5.17)$$

Its multi- $\sigma$  version can be also applied in the same way as

$$\Phi = [\Phi_{\sigma_1}^{(2)}, \Phi_{\sigma_2}^{(2)}, \dots, \Phi_{\sigma_U}^{(2)}, \mathbf{1}]. \quad (5.18)$$

### 5.1.3 Reduction of the size of Basis Matrix

The function approximation method occupies a lot of memory when the size of basis matrix  $\Phi$  is very large. It will be  $8I^2$  bytes in double-precision floating-point with  $U = 1$  and  $\delta t = 1$ . For example, when we factorize several  $(256 \times 256)$ -images by NMF, each image must be often unfolded to a  $65536$ -dimensional vector. Then the size of basis matrix  $\Phi$  becomes  $(65536 \times 65536)$  which needs about  $34.36$  Gbytes. This is a critical issue of the GRBF and the fastGRBF methods. In this section, I propose a method to reduce the size of basis matrix  $\Phi$ . We fold a observed signal  $\mathbf{y}_j$  to matrix  $\mathbf{Y}^{(j)}$ , for example a  $1000$ -dimensional-vector is folded to a  $(100 \times 10)$ -matrix. Next, we replace the observed matrix by  $\mathbf{Y}' = [\mathbf{Y}^{(1)}, \dots, \mathbf{Y}^{(J)}]$ . In the previous example we fold the  $1000$ -dimensional-vectors to  $(100 \times 10)$ -matrices, then we can reduce the matrix size of  $\Phi$  by almost  $99\%$  of the conventional method.

For image data, it is effective to split an image into several blocks and vectorize these blocks. I assume the image  $\mathbf{Y}_j$  is given by

$$\mathbf{Y}_j = \begin{pmatrix} \mathbf{D}_{11}^{(j)} & \dots & \mathbf{D}_{1q}^{(j)} \\ \vdots & \ddots & \vdots \\ \mathbf{D}_{p1}^{(j)} & \dots & \mathbf{D}_{pq}^{(j)} \end{pmatrix} \in \mathbb{R}^{I_1 \times I_2}. \quad (5.19)$$

Then its folded data is given by  $\mathbf{Y}^{(j)} = [\text{vec}(\mathbf{D}_{11}^{(j)}), \text{vec}(\mathbf{D}_{12}^{(j)}), \dots, \text{vec}(\mathbf{D}_{pq}^{(j)})] \in \mathbb{R}^{I_1 I_2 / (pq) \times pq}$ . In this case, 2D-basis can be applicable and useful. I call this method the fastGRBF-block-NMF.

Please note that it is related to the block transform by 2D discrete cosine transform (2D-DCT) [1] for image compression because it often splits an image into  $(8 \times 8)$ -blocks and factorizes each block by  $64$  2D-cosine bases. If we apply the 2D-cosine bases to  $\Phi$ , it becomes similar to the block transform. However, it is different from the block transform because we estimate constrained (e.g., low-rank approximation) and the most correlated  $R$  bases modeled with the 2D-cosine bases by a function approximation.

A problem is how to fold the large-scale vector when the dimension is a prime number. In this case, there are two approaches: to reduce the dimension or to expand the dimension. Although the reduction approach has a possibility to lose some information, its computational cost will be decreased. On the other hand, the expand approach does not lose any information; however, computational cost and wasted information will be added. I propose a repeat type and a symmetry type as the methods of the expansion approach. Let us consider to expand  $\mathbf{y} \in \mathbb{R}^N$  to  $\mathbf{z} \in \mathbb{R}^{N+M}$ . When  $\mathbf{y}$  has a cyclic feature, it might be proper to use the repeat type which is given as the expansion from  $\mathbf{y} = [y_1, y_2, \dots, y_N]^T$  to  $\mathbf{z} = [y_1, y_2, \dots, y_N, y_1, y_2, \dots, y_M]^T$ . The symmetry type is given as the expansion from  $\mathbf{y} = [y_1, y_2, \dots, y_N]^T$  to  $\mathbf{z} = [y_1, y_2, \dots, y_N, y_N, y_{N-1}, \dots, y_{N-M+1}]^T$ , and it could be well worked to prevent the influence of additional information for any  $\mathbf{y}$ .

### 5.1.4 Tensor Extension

Nonnegative tensor decompositions found applications in positron emission tomography (PET), electroencephalography (EEG), spectroscopy, chemometrics and environmental science [26, 9, 82]. There are two models in tensor decompositions, mainly.

The first model is Tucker model [98]. The three-mode version of the Tucker (Tucker-3) model can be expressed as

$$\begin{aligned} \underline{\mathbf{Y}} &\simeq \underline{\mathbf{G}} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C} \\ &= \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \sum_{r_3=1}^{R_3} g_{r_1, r_2, r_3} \mathbf{a}_{r_1} \circ \mathbf{b}_{r_2} \circ \mathbf{c}_{r_3}, \end{aligned} \quad (5.20)$$

where  $\underline{\mathbf{Y}} \in \mathbb{R}_+^{I \times J \times K}$  is a observed tensor signal,  $\underline{\mathbf{G}} \in \mathbb{R}_+^{R_1 \times R_2 \times R_3}$  is a core tensor,  $\mathbf{A} \in \mathbb{R}_+^{I \times R_1}$ ,  $\mathbf{B} \in \mathbb{R}_+^{J \times R_2}$ ,  $\mathbf{C} \in \mathbb{R}_+^{K \times R_3}$  are factor matrices,  $\times_k$  is the  $k$ -th way tensor and matrix product, and  $\circ$  means outer product. Its nonnegativity constrained decomposition is called the nonnegative Tucker decomposition (NTD). The Tucker model can be rewritten by the following matricized and vectorized forms as

$$\mathbf{Y}_{(1)} \simeq \mathbf{A} \mathbf{G}_{(1)} (\mathbf{C}^T \otimes \mathbf{B}^T), \quad (5.21)$$

$$\mathbf{Y}_{(2)} \simeq \mathbf{B} \mathbf{G}_{(2)} (\mathbf{C}^T \otimes \mathbf{A}^T), \quad (5.22)$$

$$\mathbf{Y}_{(3)} \simeq \mathbf{C} \mathbf{G}_{(3)} (\mathbf{B}^T \otimes \mathbf{A}^T), \quad (5.23)$$

$$\bar{\mathbf{y}} \simeq (\mathbf{C} \otimes \mathbf{B} \otimes \mathbf{A}) \bar{\mathbf{g}}, \quad (5.24)$$

where  $\mathbf{Y}_{(1)} \in \mathbb{R}_+^{I \times JK}$  and  $\mathbf{G}_{(1)} \in \mathbb{R}_+^{R_1 \times R_2 R_3}$  are the first-way matricized form of tensor  $\underline{\mathbf{Y}}$  and  $\underline{\mathbf{G}}$ , respectively, others are their second-way and third-way versions. If we put  $\Xi := \mathbf{G}_{(1)} (\mathbf{C}^T \otimes \mathbf{B}^T)$ , it can be regarded as NMF  $\mathbf{Y}_{(1)} \simeq \mathbf{A} \Xi$ . Thus,  $\mathbf{A}$  can be updated by NMF based update rules. In similar way,  $\mathbf{B}$ ,  $\mathbf{C}$ , and  $\underline{\mathbf{G}}$  can be also updated by NMF based update rules (e.g., by applying ALS).

The second model is the Canonical Polyadic (CP) model which is also called PARAFAC [43] or CANDECOMP [11]. The CP model is a special case of the Tucker model that  $R_1 = R_2 = R_3 = R$  and  $\underline{\mathbf{G}} = \underline{\mathbf{A}} \in \mathbb{R}_+^{R \times R \times R}$  is a diagonal tensor with entries  $g_r$  on the main diagonal as

$$\underline{\mathbf{Y}} \simeq \sum_{r=1}^R g_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r. \quad (5.25)$$

The CP model can be regarded as a straightforward tensor extension of NMF since NMF can be rewritten by  $\mathbf{Y} \simeq \sum_{r=1}^R g_r \mathbf{a}_r \circ \mathbf{b}_r$ . Thus, the CP model gives an  $R$ -rank approximation of the observed tensor signal.

In this section, I consider how to extend the fastGRBF-NMF into tensor decompositions. For example, if we assume  $\mathbf{a}_r$ ,  $\mathbf{b}_r$ , and/or  $\mathbf{c}_r$  are smooth, then these are approximated by  $\Phi_a \mathbf{w}_r$ ,  $\Phi_b \mathbf{v}_r$ , and/or  $\Phi_c \mathbf{h}_r$ , respectively. The update rule for  $\mathbf{w}_r$  can be applied to  $\mathbf{v}_r$  and  $\mathbf{h}_r$ , hence I provide an update rule for only  $\mathbf{w}_r$  in Sections 5.1.4 and 5.1.5.

### Tucker model

In this section, I extend the fastGRBF-NMF into the nonnegative Tucker decomposition (NTD) to obtain  $\mathbf{A}$  by the function approximation of  $\Phi\mathbf{W}$ . The criterion is given by

$$\begin{aligned} & \underset{\underline{\mathbf{G}}, \mathbf{W}, \mathbf{B}, \mathbf{C}}{\text{minimize}} \quad \|\underline{\mathbf{Y}} - \underline{\mathbf{G}} \times_1 \Phi\mathbf{W} \times_2 \mathbf{B} \times_3 \mathbf{C}\|_F^2, \\ & \text{s.t.} \quad \underline{\mathbf{G}} \geq 0, \mathbf{W} \geq 0, \mathbf{B} \geq 0, \mathbf{C} \geq 0. \end{aligned} \quad (5.26)$$

I call this method the fastGRBF-NTD. The key point of this problem is how to update  $\mathbf{W}$ . This is because the other factors ( $\mathbf{B}, \mathbf{C}, \underline{\mathbf{G}}$ ) can be updated by the ALS (alternating least square) based algorithm for the nonnegative Tucker decomposition [84, 85, 83, 58]. Focusing on only  $\mathbf{W}$  and using a matricized form, the problem can be rewritten as

$$\underset{\mathbf{W}}{\text{minimize}} \quad \|\mathbf{Y}_{(1)} - \Phi\mathbf{W}\mathbf{G}_{(1)}(\mathbf{C}^T \otimes \mathbf{B}^T)\|_F^2, \quad \text{s.t.} \quad \mathbf{W} \geq 0. \quad (5.27)$$

Putting  $\mathbf{Z}_r := \mathbf{Y}_{(1)} - \sum_{k \neq r} \Phi\mathbf{w}_k \xi_k^T$  and  $\Xi = [\xi_1, \dots, \xi_{R_1}]^T := \mathbf{G}_{(1)}(\mathbf{C}^T \otimes \mathbf{B}^T) \in \mathbb{R}^{R_1 \times JK}$ , I consider the problem for  $\mathbf{w}_r$  as

$$\underset{\mathbf{w}_r}{\text{minimize}} \quad \|\mathbf{Z}_r - \Phi\mathbf{w}_r \xi_r^T\|_F^2, \quad \text{s.t.} \quad \mathbf{w}_r \geq 0. \quad (5.28)$$

(5.28) is essentially equivalent to (5.3). Since the condition of  $\|\xi_r\| = 1$  must be satisfied, then an update rule for  $\mathbf{w}_r$  is given by

$$\xi_r \leftarrow [\mathbf{Z}_r^T \Phi\mathbf{w}_r]_+; \quad (5.29)$$

$$\xi_r \leftarrow \xi_r / \|\xi_r\|; \quad (5.30)$$

$$\mathbf{w}_r \leftarrow [\mathbf{w}_r \otimes (\Phi^T \mathbf{Z}_r \xi_r) \odot (\Phi^T \Phi\mathbf{w}_r)]_+; \quad (5.31)$$

Note that the update rules (5.29) and (5.30) are important steps to update  $\mathbf{w}_r$  accurately.

The algorithm for fastGRBF-NTD can be summarized as Algorithm 19. The 5-13th lines are update procedures for  $\mathbf{W}$ . If we apply the function approximation to factor matrices  $\mathbf{B}$  and/or  $\mathbf{C}$ , the 14th and/or 15th lines must be modified in a similar way to  $\mathbf{W}$ .

### 5.1.5 CP model

In this section, I extend the fastGRBF-NMF into the nonnegative CP decomposition (NCPD) to obtain  $\mathbf{A}$  by function approximation  $\Phi\mathbf{W}$ . The criterion is given by

$$\begin{aligned} & \underset{g_r, \mathbf{w}_r, \mathbf{b}_r, \mathbf{c}_r}{\text{minimize}} \quad \|\underline{\mathbf{Y}} - \sum_{r=1}^R g_r(\Phi\mathbf{w}_r) \circ \mathbf{b}_r \circ \mathbf{c}_r\|_F^2, \\ & \text{s.t.} \quad g_r \geq 0, \mathbf{w}_r \geq 0, \mathbf{b}_r \geq 0, \mathbf{c}_r \geq 0, \\ & \quad \|\Phi\mathbf{w}_r\| = \|\mathbf{b}_r\| = \|\mathbf{c}_r\| = 1, \end{aligned} \quad (5.32)$$

**Algorithm 19** FastGRBF-NTD algorithm

---

```

1: Input:  $\underline{Y}$ ,  $R_1$ ,  $R_2$ ,  $R_3$ , and  $\Phi$ 
2: Initialize:  $\underline{W}$  by randomly, and  $\underline{B}$ ,  $\underline{C}$ ,  $\underline{G}$  by some initialization method for NTD
3:  $\underline{E} \leftarrow \underline{Y} - \underline{G} \times_1 \Phi \underline{W} \times_2 \underline{B} \times_3 \underline{C}$ ;
4: repeat
5:    $\Xi \leftarrow \mathbf{G}_{(1)}(C^T \otimes B^T)$ ;
6:   for  $r = 1, \dots, R_1$  do
7:      $\mathbf{Z}_r \leftarrow \mathbf{E}_{(1)} + (\Phi \mathbf{w}_r) \xi_r^T$ ;
8:      $\xi_r \leftarrow [\mathbf{Z}_r^T (\Phi \mathbf{w}_r)]_+$ ;
9:      $\xi_r \leftarrow \xi_r / \|\xi_r\|$ ;
10:     $\mathbf{w}_r \leftarrow [\mathbf{w}_r \otimes \{\Phi^T(\mathbf{Z}_r \xi_r)\} \otimes \{\Phi^T(\Phi \mathbf{w}_r)\}]_+$ ;
11:     $\mathbf{E}_{(1)} \leftarrow \mathbf{Z}_r - (\Phi \mathbf{w}_r) \xi_r^T$ ;
12:  end for
13:   $\underline{A} \leftarrow \Phi \underline{W}$ ;
14:   $\underline{B} \leftarrow \underline{B} \otimes (\mathbf{Y}_{(2)}(\underline{A} \otimes \underline{C}) \mathbf{G}_{(2)}^T) \otimes (\underline{B} \mathbf{G}_{(2)}(\underline{A}^T \underline{A} \otimes \underline{C}^T \underline{C}) \mathbf{G}_{(2)}^T)$ ;
15:   $\underline{C} \leftarrow \underline{C} \otimes (\mathbf{Y}_{(3)}(\underline{A} \otimes \underline{B}) \mathbf{G}_{(3)}^T) \otimes (\underline{C} \mathbf{G}_{(3)}(\underline{A}^T \underline{A} \otimes \underline{B}^T \underline{B}) \mathbf{G}_{(3)}^T)$ ;
16:   $\underline{G} \leftarrow \underline{G} \otimes (\underline{Y} \times_1 \underline{A}^T \times_2 \underline{B}^T \times_3 \underline{C}^T) \otimes (\underline{G} \times_1 \underline{A} \underline{A}^T \times_2 \underline{B} \underline{B}^T \times_3 \underline{C} \underline{C}^T)$ ;
17:   $\underline{E} \leftarrow \underline{Y} - \underline{G} \times_1 \Phi \underline{W} \times_2 \underline{B} \times_3 \underline{C}$ ;
18: until  $\|\underline{E}\|_F^2$  converges
19: Output:  $\underline{W}$ ,  $\underline{B}$ ,  $\underline{C}$ , and  $\underline{G}$ 

```

---

for  $r = 1, 2, \dots, R$ . I call this method the fastGRBF-NCPD. We apply the HALS based algorithm to the problem (5.32). Putting  $\underline{Y}_r := \underline{Y} - \sum_{k \neq r} g_r(\Phi \mathbf{w}_k) \circ \mathbf{b}_k \circ \mathbf{c}_k$  and using the first-way matricized form, the sub-problem for  $\mathbf{w}_r$  can be rewritten by

$$\begin{aligned}
& \underset{\mathbf{w}_r}{\text{minimize}} \|\mathbf{Y}_{r(1)} - g_r(\Phi \mathbf{w}_r)(\mathbf{c}_r^T \otimes \mathbf{b}_r^T)\|_F^2, \\
& \text{s.t. } \|\Phi \mathbf{w}_r\| = 1, \mathbf{w}_r \geq 0.
\end{aligned} \tag{5.33}$$

In the similar relation between (5.3) and (5.8), problem (5.33) can be transformed to

$$\begin{aligned}
& \underset{\mathbf{w}_r}{\text{minimize}} \|\mathbf{Y}_{r(1)}(\mathbf{c}_r \otimes \mathbf{b}_r) - g_r \Phi \mathbf{w}_r\|_F^2, \\
& \text{s.t. } \|\Phi \mathbf{w}_r\| = 1, \mathbf{w}_r \geq 0,
\end{aligned} \tag{5.34}$$

where we have  $(\mathbf{c}_r^T \otimes \mathbf{b}_r^T)(\mathbf{c}_r \otimes \mathbf{b}_r) = (\mathbf{c}_r^T \mathbf{c}_r \otimes \mathbf{b}_r^T \mathbf{b}_r) = 1$ . Since there is a constraint  $\|\Phi \mathbf{w}_r\| = 1$  in (5.34), its solution is not dependent on  $g_r$ . Thus, the update rule for  $\mathbf{w}_r$  is given by

$$\mathbf{w}_r \leftarrow [\mathbf{w}_r \otimes \{\Phi^T \mathbf{Y}_{r(1)}(\mathbf{c}_r \otimes \mathbf{b}_r)\} \otimes \{\Phi^T \Phi \mathbf{w}_r\}]_+, \tag{5.35}$$

$$\mathbf{w}_r \leftarrow \mathbf{w}_r / \|\Phi \mathbf{w}_r\|. \tag{5.36}$$

**Algorithm 20** FastGRBF-NCPD algorithm

---

```

1: Input:  $\underline{Y}$ ,  $R$ , and  $\Phi$ 
2: Initialize:  $\underline{W}$  by randomly, and  $\underline{B}$ ,  $\underline{C}$ ,  $\underline{G}$  by some initialization method for nonnegative CP
   decomposition
3:  $\underline{E} \leftarrow \underline{Y} - \underline{G} \times_1 \Phi \underline{W} \times_2 \underline{B} \times_3 \underline{C}$  ;
4: repeat
5:   for  $r = 1, \dots, R$  do
6:      $\underline{Y}_r \leftarrow \underline{E} + g_r(\Phi \underline{w}_r) \circ \underline{b}_r \circ \underline{c}_r$  ;
7:      $\underline{w}_r \leftarrow [\underline{w}_r \circledast \{\Phi^T(\underline{Y}_r \times_2 \underline{b}_r^T \times_3 \underline{c}_r^T)\} \circledcirc \{\Phi^T \Phi \underline{w}_r\}]_+$ ;
8:      $\underline{w}_r \leftarrow \underline{w}_r / \|\Phi \underline{w}_r\|$  ;
9:      $\underline{b}_r \leftarrow [\underline{Y}_r \times_1 (\Phi \underline{w}_r)^T \times_3 \underline{c}_r^T]_+$  ;
10:     $\underline{b}_r \leftarrow \underline{b}_r / \|\underline{b}_r\|$  ;
11:     $\underline{c}_r \leftarrow [\underline{Y}_r \times_1 (\Phi \underline{w}_r)^T \times_2 \underline{b}_r^T]_+$  ;
12:     $\underline{c}_r \leftarrow \underline{c}_r / \|\underline{c}_r\|$  ;
13:     $g_r \leftarrow [\underline{Y}_r \times_1 (\Phi \underline{w}_r)^T \times_2 \underline{b}_r^T \times_3 \underline{c}_r^T]_+$  ;
14:     $\underline{E} \leftarrow \underline{Y}_r - g_r(\Phi \underline{w}_r) \circ \underline{b}_r \circ \underline{c}_r$  ;
15:   end for
16: until  $\|\underline{E}\|_F^2$  converge
17: Output:  $\underline{W}$ ,  $\underline{B}$ ,  $\underline{C}$ , and  $\underline{G}$ 

```

---

The fastGRBF-NCPD method can be summarized as Algorithm 20. The 7-8th lines are update procedures for  $\underline{w}_r$ . If we apply the function approximation to also  $\underline{b}_r$  and/or  $\underline{c}_r$ , the 9-10th and/or 11-12th lines must be modified in a similar way to  $\underline{w}_r$ .

## 5.2 Linked CP Tensor Decomposition

In Section 3.5, I introduced the matrix based common and individual feature extraction methods; however, their methods can not be applied to multiway-array data. In this section, I extend this common and individual feature extraction model to the CP tensor decomposition model, called the linked CP tensor decomposition (LCPTD) [107]. Furthermore, I impose the sparsity and nonnegativity constraints to this model.

LCPTD model is given by

$$\underline{\mathbf{Z}}^{(s)} \approx \widehat{\underline{\mathbf{Z}}}^{(s)} = \llbracket \underline{\mathbf{G}}^{(s)}; \underline{\mathbf{U}}^{(1,s)}, \dots, \underline{\mathbf{U}}^{(N,s)} \rrbracket = \sum_{j=1}^J g_j^{(s)} \mathbf{u}_j^{(1,s)} \circ \dots \circ \mathbf{u}_j^{(N,s)}, \quad (5.37)$$

where each factor matrix  $\underline{\mathbf{U}}^{(n,s)} = [\underline{\mathbf{U}}_C^{(n)}, \underline{\mathbf{U}}_I^{(n,s)}] \in \mathbb{R}^{I_n \times J}$  is composed of two set of bases:  $\underline{\mathbf{U}}_C^{(n)} \in \mathbb{R}^{I_n \times L_n}$  (with  $0 \leq L_n \leq J$ ), which is a common factor matrix for all blocks and corresponds to the same or maximally correlated components, and  $\underline{\mathbf{U}}_I^{(n,s)} \in \mathbb{R}^{I_n \times J - L_n}$ , which corresponds to different

individual characteristics.

The LCPTD can be considered as a generalized model of the simultaneous decomposition. When we let  $L_n = J$ , its decomposition is equivalent to the simultaneous common factor decomposition [82]. On the other hand, when  $L_n = 0$ , its decomposition of each subject is equivalent to the standard tensor decomposition. Then the LTD is an intermediate decomposition between the simultaneous and the ordinary tensor decompositions.

### 5.2.1 LCPTD-HALS algorithm

In this section, I introduce a new HALS algorithm for LCPTD. Optimization criterion for LCPTD is given by

$$\text{minimize } \sum_{s=1}^S \left\| \underline{\mathbf{Z}}^{(s)} - \sum_{j=1}^J g_j^{(s)} \mathbf{u}_j^{(1,s)} \circ \dots \circ \mathbf{u}_j^{(N,s)} \right\|_F^2, \quad (5.38)$$

$$\text{subject to } \mathbf{u}_j^{(n,1)} = \dots = \mathbf{u}_j^{(n,S)} \text{ for } j \leq L_n, \|\mathbf{u}_j^{(n,s)}\| = 1, \quad (5.39)$$

for all  $n$ ,  $s$ , and  $j$ . Furthermore, I add  $\|\mathbf{u}_j^{(n,s)}\|_1 < v$  or  $u_{ji}^{(n,s)} \geq 0$ ,  $g_j^{(s)} \geq 0 \quad \forall i, j, n, s$  into (5.39) if we want to obtain sparse or nonnegative components.

The hierarchical ALS (HALS) algorithm was first proposed for the NMF and NTF in [24]. The HALS algorithm was applied to the CP model and it achieved a high performance in [19]. In this algorithm, I consider  $J$  local-problems and solve them sequentially and iteratively instead of solving (5.38) and (5.39), directly. Let  $\underline{\mathbf{Y}}_j^{(s)} := \underline{\mathbf{Z}}^{(s)} - \sum_{i \neq j} g_i^{(s)} \mathbf{u}_i^{(1,s)} \circ \dots \circ \mathbf{u}_i^{(N,s)}$ , the  $j$ -th local problem is given by

$$\text{minimize } \sum_{s=1}^S \|\underline{\mathbf{Y}}_j^{(s)} - g_j^{(s)} \mathbf{u}_j^{(1,s)} \circ \dots \circ \mathbf{u}_j^{(N,s)}\|_F^2, \quad (5.40)$$

$$\text{subject to } \mathbf{u}_j^{(n,1)} = \dots = \mathbf{u}_j^{(n,S)} \text{ if } j \leq L_n, \|\mathbf{u}_j^{(n,s)}\| = 1, \quad (5.41)$$

for all  $n$  and  $s$ . The LTD-HALS algorithm can be summarized as Algorithm 21; note that it does not require matrix inversion and is solved by only simple calculations.

If we want to obtain sparse components, I implement the following updates after (5.42) in Algorithm 21:

$$\mathbf{u}_j^{(n,s)} \leftarrow \text{sign}(\mathbf{u}_j^{(n,s)}) \otimes [\text{abs}(\mathbf{u}_j^{(n,s)}) - \xi_n \mathbf{1}]_+ \text{ for all } s; \quad (5.44)$$

where  $\xi_n$  is a positive parameter deciding on their sparsity.

If we want to obtain nonnegative components, I implement the following updates after (5.42) and (5.43) in Algorithm 21:

$$\mathbf{u}_j^{(n,s)} \leftarrow [\mathbf{u}_j^{(n,s)}]_+ \text{ for all } s, \quad (5.45)$$

$$g_j^{(s)} \leftarrow [g_j^{(s)}]_+ \text{ for all } s. \quad (5.46)$$

**Algorithm 21** LTD-HALS algorithm

---

**Input:**  $\{\underline{\mathbf{Z}}^{(s)}\}_{s=1}^S$ ,  $J$ , and  $\{L_n\}_{n=1}^N$   
**Initialize:**  $\{\mathbf{g}^{(s)}, \{\mathbf{U}^{(n,s)}\}_{n=1}^N\}_{s=1}^S$ .  
 $\underline{\mathbf{E}}^{(s)} = \underline{\mathbf{Z}}^{(s)} - \sum_{j=1}^J g_j^{(s)} \mathbf{u}_j^{(1,s)} \circ \dots \circ \mathbf{u}_j^{(N,s)}$  for all  $s$ ;

**repeat**

**for**  $j = 1, \dots, J$  **do**

$\underline{\mathbf{Y}}_j^{(s)} = \underline{\mathbf{E}}^{(s)} + g_j^{(s)} \mathbf{u}_j^{(1,s)} \circ \dots \circ \mathbf{u}_j^{(N,s)}$  for all  $s$ ;

**for**  $n = 1, \dots, N$  **do**

      Updating  $\mathbf{u}_j^{(n,s)}$ :

$$\begin{aligned} \mathbf{u}_j^{(n,s)} \leftarrow & g_j^{(s)} \underline{\mathbf{Y}}_j^{(s)} \times_1 \mathbf{u}_j^{(1,s)} \dots \times_{n-1} \mathbf{u}_j^{(n-1,s)} \\ & \times_{n+1} \mathbf{u}_j^{(n+1,s)} \dots \times_N \mathbf{u}_j^{(N,s)} \text{ for all } s; \end{aligned} \quad (5.42)$$

**if**  $j \leq L_n$ ,  $\mathbf{t} \leftarrow \sum_{s=1}^S \mathbf{u}_j^{(n,s)}$ ;  $\mathbf{u}_j^{(n,s)} \leftarrow \mathbf{t}$  for all  $s$ ; **end if**

    Normalize  $\mathbf{u}_j^{(n,s)} \leftarrow \mathbf{u}_j^{(n,s)} / \|\mathbf{u}_j^{(n,s)}\|$  for all  $s$ ;

**end for**

  Update  $g_j^{(s)}$ :

$$g_j^{(s)} \leftarrow \underline{\mathbf{Y}}_j^{(s)} \times_1 \mathbf{u}_j^{(1,s)} \dots \times_N \mathbf{u}_j^{(N,s)} \text{ for all } s; \quad (5.43)$$

$\underline{\mathbf{E}}^{(s)} = \underline{\mathbf{Y}}_j^{(s)} - g_j^{(s)} \mathbf{u}_j^{(1,s)} \circ \dots \circ \mathbf{u}_j^{(N,s)}$  for all  $s$ ;

**end for**

**until**  $\sum_{s=1}^S \|\underline{\mathbf{E}}^{(s)}\|_F^2$  converge

**Output:**  $\{\mathbf{g}^{(s)}, \{\mathbf{U}^{(n,s)}\}_{n=1}^N\}_{s=1}^S$

---

### 5.3 Linked Tucker Decomposition

In this section, I extend the LCPTD model to the Tucker model, called the linked Tucker decomposition (LTD). Since Tucker model is very general and does not have a unique solution, non-constraint is not suitable. I consider orthogonality, sparsity, and nonnegativity constraints for LTD model.

The LTD model is given by

$$\begin{aligned} \underline{\mathbf{Z}}^{(s)} \approx \widehat{\underline{\mathbf{Z}}}^{(s)} &= \underline{\mathbf{G}}^{(s)} \times \{\mathbf{U}^{(s)}\} \\ &= \underline{\mathbf{G}}^{(s)} \times_1 \mathbf{U}^{(1,s)} \dots \times_N \mathbf{U}^{(N,s)}, \end{aligned} \quad (5.47)$$

where each factor matrix  $\mathbf{U}^{(n,s)} = [\mathbf{U}_C^{(n)}, \mathbf{U}_I^{(n,s)}] \in \mathbb{R}^{I_n \times J_n}$  is composed of two set of bases:  $\mathbf{U}_C^{(n)} \in \mathbb{R}^{I_n \times L_n}$  (with  $0 \leq L_n \leq J_n$ ), which is common factor matrices for all blocks and corresponds to the same or maximally correlated components, and  $\mathbf{U}_I^{(n,s)} \in \mathbb{R}^{I_n \times R_n}$  (i.e.  $L_n + R_n = J_n$ ), which corresponds to different individual characteristics. Fig. 5.2 depicts a three-modes case of LTD model. This model was first briefly introduced as an idea for the linked multiway blind source separation in

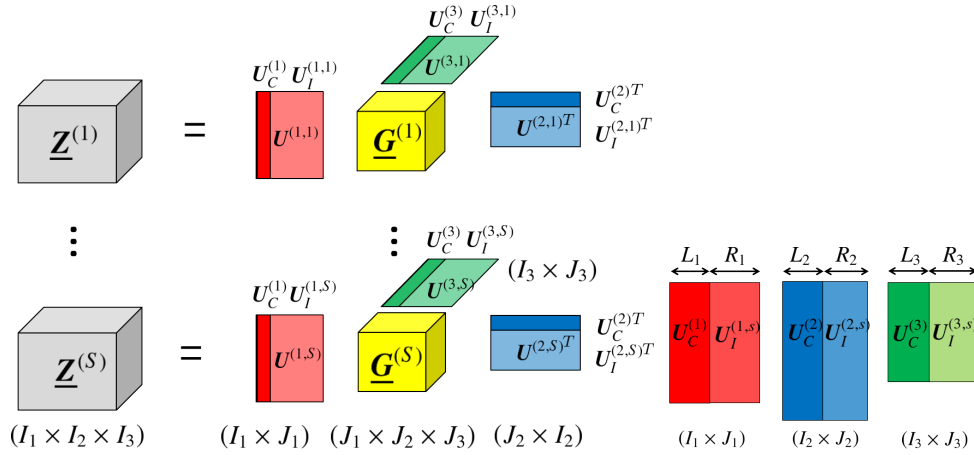


Figure 5.2: Linked Tucker Decomposition (LTD): in special case, we obtain STD for  $L_n = J_n$  and ITD for  $L_n = 0$ .

[17], but its algorithms and applications have not been proposed yet. In this paper, I propose a new approach and its application.

The LTD can be considered as a generalized model of the STD. When we put  $L_n = J_n$ , its decomposition is equivalent to the simultaneous common factor decomposition (i.e. the STD model) [82]. On the other hand, when  $L_n = 0$ , its decomposition of each subject is equivalent to the standard tensor decomposition (i.e. the ITD model). Then the LTD is an intermediate decomposition between the simultaneous and the individual tensor/Tucker decompositions.

### 5.3.1 Orthogonal LTD method

In this section, I propose a novel method called the “linked multi-way principal component analysis” (LMPCA) and associated algorithm called the LTD-ALS algorithm. The LMPCA can be considered as a prototype of all linked group tensor analysis methods, just as PCA is a prototype of single block data analysis methods. We could derive many extensions for examples the linked multiway sparse component analysis (LMSCA), the linked nonnegative Tucker decomposition (LNTD), the linked multiway independent component analysis (LMICA) and so on. The ALS algorithm was first proposed for the three-way PCA in [60]. The LTD-ALS can be derived in a similar way. Optimization criterion for the LMPCA is given by

$$\begin{aligned} \text{minimize} \quad & C(\Theta) := \sum_{s=1}^S \|\underline{\mathbf{Z}}^{(s)} - \underline{\mathbf{G}}^{(s)} \times \{\mathbf{U}^{(s)}\}\|_F^2, \\ \text{subject to} \quad & \mathbf{U}^{(n,s)T} \mathbf{U}^{(n,s)} = \mathbf{I}_{J_n} \text{ for all } n \text{ and } s, \end{aligned} \quad (5.48)$$

where  $C(\Theta)$  is the cost function and  $\Theta$  represents a set of parameters (i.e.,  $\{\{\mathbf{U}^{(n,s)}\}_{n=1}^N, \underline{\mathbf{G}}^{(s)}\}_{s=1}^S\}$ ). Substituting one of the KKT conditions  $\underline{\mathbf{G}}^{(s)} = \underline{\mathbf{Z}}^{(s)} \times \{\mathbf{U}^{(s)T}\}$  into (5.48) and transforming this, then

---

**Algorithm 22** Local algorithm to update  $n$ -mode factor matrices for orthogonal LTD
 

---

```

1: Input:  $\{\mathbf{Y}_n^{(s)}, \mathbf{U}_I^{(n,s)}\}_{s=1}^S, \mathbf{U}_C^{(n)}, L_n,$  and  $R_n$ 
2: repeat
3:    $\mathbf{Y}_n \leftarrow [\mathbf{Y}_n^{(1)} - \mathbf{U}_I^{(n,1)} \mathbf{U}_I^{(n,1)T} \mathbf{Y}_n^{(1)}, \mathbf{Y}_n^{(2)} - \mathbf{U}_I^{(n,2)} \mathbf{U}_I^{(n,2)T} \mathbf{Y}_n^{(2)},$ 
4:      $\dots, \mathbf{Y}_n^{(S)} - \mathbf{U}_I^{(n,S)} \mathbf{U}_I^{(n,S)T} \mathbf{Y}_n^{(S)}]$ ;
5:    $\mathbf{U}_C^{(n)} \leftarrow$  left singular matrix of  $\text{tSVD}(\mathbf{Y}_n, L_n)$ ;
6:   for  $s = 1, \dots, S$  do
7:      $\mathbf{X} \leftarrow \mathbf{Y}_n^{(s)} - \mathbf{U}_C^{(n)} \mathbf{U}_C^{(n)T} \mathbf{Y}_n^{(s)}$ ;
8:      $\mathbf{U}_I^{(n,s)} \leftarrow$  left singular matrix of  $\text{tSVD}(\mathbf{X}, R_n)$ ;
9:   end for
10: until  $C(\Theta)$  converge
11: Output:  $\mathbf{U}_C^{(n)}, \{\mathbf{U}_I^{(n,s)}\}_{s=1}^S$ 

```

---

we can obtain a simplified optimization problem with respect to  $n$ -mode factor matrices as

$$\begin{aligned}
 & \text{minimize} \quad \sum_{s=1}^S \|\mathbf{Y}_n^{(s)} - \mathbf{U}^{(n,s)} \mathbf{U}^{(n,s)T} \mathbf{Y}_n^{(s)}\|_F^2, \\
 & \text{subject to} \quad \mathbf{U}^{(n,s)T} \mathbf{U}^{(n,s)} = \mathbf{I}_{J_n} \text{ for all } s,
 \end{aligned} \tag{5.49}$$

where  $\mathbf{Y}_n^{(s)} := [\underline{\mathbf{Z}}^{(s)} \times_{-n} \{\mathbf{U}^{(s)T}\}]_{(n)} \in \mathbb{R}^{I_n \times \prod_{i \neq n} J_i}$ . Problem (5.49) can be considered as the PCA problem with respect to  $\mathbf{U}^{(n,s)}$ . Please note that the solutions of PCA and the left factor matrix of truncated SVD (tSVD) are the same. Thus I consider to factorize  $\mathbf{Y}_n^{(s)}$  by the common and the individual left factor matrices in a similar way to JIVE algorithm. To update  $\mathbf{U}_C^{(n)}$ , I run the tSVD as

$$\begin{aligned}
 & [\mathbf{Y}_n^{(1)} - \mathbf{U}_I^{(n,1)} \mathbf{U}_I^{(n,1)T} \mathbf{Y}_n^{(1)}, \mathbf{Y}_n^{(2)} - \mathbf{U}_I^{(n,2)} \mathbf{U}_I^{(n,2)T} \mathbf{Y}_n^{(2)}, \\
 & \quad \dots, \mathbf{Y}_n^{(S)} - \mathbf{U}_I^{(n,S)} \mathbf{U}_I^{(n,S)T} \mathbf{Y}_n^{(S)}] \simeq \mathbf{U}_C \mathbf{D}_C \mathbf{V}_C^T.
 \end{aligned} \tag{5.50}$$

To update  $\mathbf{U}_I^{(n,s)}$ , I run the tSVD one by one, individually, which is given by

$$[\mathbf{Y}_n^{(s)} - \mathbf{U}_C \mathbf{U}_C^T \mathbf{Y}_n^{(s)}] \simeq \mathbf{U}_I^{(n,s)} \mathbf{D}_I^{(s)} \mathbf{V}_I^{(s)T}. \tag{5.51}$$

In order to obtain optimal parameters, I iterate the common and the individual steps alternately until convergence. Thus, the algorithm for local problem (5.49) can be summarized as Algorithm 22.

In the orthogonal LTD algorithm, I solve the local problems (5.49), sequentially and iteratively. To solve (5.49) I consider two steps: the first and the second steps are for common and individual factors, respectively. The orthogonal LTD algorithm can be summarized as Algorithm 23.

### 5.3.2 Sparse LTD

In this section, I consider to impose the sparse constraint to each factor matrices  $\mathbf{U}^{(n,s)}$ . Sparsity is a very important property for some special sparse data analysis such as text-mining data. The one of

**Algorithm 23** Orthogonal LTD algorithm

---

```

1: Input:  $\{\underline{\mathbf{Z}}^{(s)}\}_{s=1}^S$  and  $\{L_n, R_n\}_{n=1}^N$ 
2: Initialize: orthonormal matrices  $\{\mathbf{U}_C^{(n)}, \{\mathbf{U}_I^{(n,s)}\}_{s=1}^S\}_{n=1}^N$ 
3: repeat
4:   for  $n = 1, \dots, N$  do
5:      $\mathbf{Y}_n^{(s)} \leftarrow [\underline{\mathbf{Z}}^{(s)} \times_{-n} \{\mathbf{U}^{(s)T}\}]_{(n)}$  for all  $s$ ;
6:     Update  $\mathbf{U}_C^{(n)}$  and  $\{\mathbf{U}_I^{(n,s)}\}_{s=1}^S$  by Algorithm 22;
7:   end for
8:    $\underline{\mathbf{G}}^{(s)} = \underline{\mathbf{Z}}^{(s)} \times \{\mathbf{U}^{(s)T}\}$  for all  $s$ ;
9: until  $C(\Theta)$  converge
10: Output:  $\{\underline{\mathbf{G}}^{(s)}, \{\mathbf{U}^{(n,s)}\}_{n=1}^N\}_{s=1}^S$ 

```

---

most popular approaches to give a sparse representation is the l1-norm regularization. For example, there are the sparse PCA [121] and the penalized matrix decomposition (PMD) [105] as the methods of sparse analysis. Since the PMD can be regarded as a sparse extension of tSVD, we can apply the PMD method to LTD. The local problem with respect to  $n$ -mode factor matrices can be given by

$$\begin{aligned}
& \text{minimize} && \sum_{s=1}^S \|\mathbf{Y}_n^{(s)} - \mathbf{U}^{(n,s)} \mathbf{U}^{(n,s)\dagger} \mathbf{Y}_n^{(s)}\|_F^2, \\
& \text{subject to} && \|\mathbf{u}_{j_n}^{(n,s)}\|_1 \leq c_n \text{ for all } s,
\end{aligned} \tag{5.52}$$

where  $\dagger$  describes the Moore-Penrose's pseudo inverse (i.e.,  $\mathbf{U}^{(n,s)\dagger} = (\mathbf{U}^{(n,s)T} \mathbf{U}^{(n,s)})^{-1} \mathbf{U}^{(n,s)T}$ ),  $\|\cdot\|_1$  is the l1-norm, and  $c_n$  is a trade-off parameter for l1-norm regularization. I set this parameter in  $1 < c_n < \sqrt{J_n}$ . This solution can be obtained by a similar way to Algorithm 22. The differences are that the transposition  $T$  becomes the pseudo inverse  $\dagger$  in the matrix deflation formula, and tSVD becomes PMD for update. Thus the local and global algorithm for sparse LTD can be summarized as Algorithms 24 and 25.

### 5.3.3 Nonnegative LTD

In this section, I impose the nonnegative constraint to the LTD method. I call this method the nonnegative linked Tucker decomposition (NLTD). In order to give an algorithm for NLTD, techniques of nonnegative Tucker decomposition [58, 83, 84, 85] and group NMF [63] can be used. In the same manner as that the algorithm for NMF is quite different from the algorithm for tSVD, the algorithm for NLTD is also different from the OLTD algorithm. Optimization criterion is given by

$$\begin{aligned}
& \text{minimize} && C(\Theta) := \sum_{s=1}^S \|\underline{\mathbf{Z}}^{(s)} - \underline{\mathbf{G}}^{(s)} \times \{\mathbf{U}^{(s)}\}\|_F^2, \\
& \text{subject to} && \mathbf{U}^{(n,s)} \geq 0, \underline{\mathbf{G}}^{(s)} \geq 0 \text{ for all } n \text{ and } s.
\end{aligned} \tag{5.53}$$

**Algorithm 24** Local algorithm to update  $n$ -mode factor matrices for sparse LTD

---

```

1: Input:  $\{\mathbf{Y}_n^{(s)}, \mathbf{U}_I^{(n,s)}\}_{s=1}^S$ ,  $\mathbf{U}_C^{(n)}$ ,  $L_n$ ,  $R_n$ , and  $c_n$ 
2: repeat
3:    $\mathbf{Y}_n \leftarrow [\mathbf{Y}_n^{(1)} - \mathbf{U}_I^{(n,1)} \mathbf{U}_I^{(n,1)\dagger} \mathbf{Y}_n^{(1)}, \mathbf{Y}_n^{(2)} - \mathbf{U}_I^{(n,2)} \mathbf{U}_I^{(n,2)\dagger} \mathbf{Y}_n^{(2)},$ 
4:      $\dots, \mathbf{Y}_n^{(S)} - \mathbf{U}_I^{(n,S)} \mathbf{U}_I^{(n,S)\dagger} \mathbf{Y}_n^{(S)}]$ ;
5:    $\mathbf{U}_C^{(n)} \leftarrow$  left factor matrix of  $\text{PMD}(\mathbf{Y}_n, L_n, c_n)$ ;
6:   for  $s = 1, \dots, S$  do
7:      $\mathbf{X} \leftarrow \mathbf{Y}_n^{(s)} - \mathbf{U}_C^{(n)} \mathbf{U}_C^{(n)\dagger} \mathbf{Y}_n^{(s)}$ ;
8:      $\mathbf{U}_I^{(n,s)} \leftarrow$  left factor matrix of  $\text{PMD}(\mathbf{X}, R_n, c_n)$ ;
9:   end for
10: until  $C(\Theta)$  converge
11: Output:  $\mathbf{U}_C^{(n)}, \{\mathbf{U}_I^{(n,s)}\}_{s=1}^S$ 

```

---

**Algorithm 25** Sparse LTD algorithm

---

```

1: Input:  $\{\mathbf{Z}^{(s)}\}_{s=1}^S$  and  $\{L_n, R_n, c_n\}_{n=1}^N$ 
2: Initialize: orthonormal matrices  $\{\mathbf{U}_C^{(n)}, \{\mathbf{U}_I^{(n,s)}\}_{s=1}^S\}_{n=1}^N$ 
3: repeat
4:   for  $n = 1, \dots, N$  do
5:      $\mathbf{Y}_n^{(s)} \leftarrow [\mathbf{Z}^{(s)} \times_{-n} \{\mathbf{U}^{(s)T}\}]_{(n)}$  for all  $s$ ;
6:     Update  $\mathbf{U}_C^{(n)}$  and  $\{\mathbf{U}_I^{(n,s)}\}_{s=1}^S$  by Algorithm 24;
7:   end for
8:    $\mathbf{G}^{(s)} = \mathbf{Z}^{(s)} \times \{\mathbf{U}^{(s)\dagger}\}$  for all  $s$ ;
9: until  $C(\Theta)$  converge
10: Output:  $\{\mathbf{G}^{(s)}, \{\mathbf{U}^{(n,s)}\}_{n=1}^N\}_{s=1}^S$ 

```

---

Please note that nonnegative constraints are imposed to not only factor matrices  $\mathbf{U}^{(n,s)}$  but also core tensors  $\mathbf{G}^{(s)}$ . Thus, we can not obtain  $\mathbf{G}^{(s)}$ , analytically, unlike OLTD and SLTD. From this reason, the algorithm for NLTD becomes quite complex. Let us use the following notations

$$\mathbf{U}_{(-n)}^{(s)} := \mathbf{U}^{(N,s)} \otimes \dots \otimes \mathbf{U}^{(n+1,s)} \otimes \mathbf{U}^{(n-1,s)} \otimes \dots \otimes \mathbf{U}^{(1,s)}, \quad (5.54)$$

$$\mathbf{Y}_{(n)}^{(s)} := \mathbf{G}_{(n)}^{(s)} \mathbf{U}_{(-n)}^{(s)T}, \quad (5.55)$$

$$\mathbf{Y}_{C(n)}^{(s)} := \text{first } L_n \text{ row vectors of } \mathbf{Y}_{(n)}^{(s)}, \quad (5.56)$$

$$\mathbf{Y}_{I(n)}^{(s)} := \text{last } R_n \text{ row vectors of } \mathbf{Y}_{(n)}^{(s)}, \quad (5.57)$$

where  $\otimes$  describes Kronecker product, then the cost function can be rewritten as

$$\begin{aligned}
C(\Theta) &= \sum_{s=1}^S \|\mathbf{Z}_{(n)}^{(s)} - \mathbf{U}^{(n,s)} \mathbf{G}_{(n)}^{(s)} \mathbf{U}_{(-n)}^{(s)T}\|_F^2 \\
&= \sum_{s=1}^S \|\mathbf{Z}_{(n)}^{(s)} - \mathbf{U}^{(n,s)} \mathbf{Y}_{(n)}^{(s)}\|_F^2 \\
&= \sum_{s=1}^S \|\mathbf{Z}_{(n)}^{(s)} - \mathbf{U}_C^{(n,s)} \mathbf{Y}_{C(n)}^{(s)} - \mathbf{U}_I^{(n,s)} \mathbf{Y}_{I(n)}^{(s)}\|_F^2.
\end{aligned} \tag{5.58}$$

Therefore, the update rules of common and individual factor matrices are given by

$$\mathbf{U}_C^{(n)} \leftarrow \mathbf{U}_C^{(n)} \otimes \left[ \sum_{s=1}^S \mathbf{Z}_{(n)}^{(s)} \mathbf{Y}_{C(n)}^{(s)T} \right] \oslash \left[ \sum_{s=1}^S \mathbf{U}^{(n,s)} \mathbf{Y}_{(n)}^{(s)} \mathbf{Y}_{C(n)}^{(s)T} \right], \tag{5.59}$$

$$\mathbf{U}_I^{(n,s)} \leftarrow \mathbf{U}_I^{(n,s)} \otimes \left[ \mathbf{Z}_{(n)}^{(s)} \mathbf{Y}_{I(n)}^{(s)T} \right] \oslash \left[ \mathbf{U}^{(n,s)} \mathbf{Y}_{(n)}^{(s)} \mathbf{Y}_{I(n)}^{(s)T} \right] \tag{5.60}$$

for all  $s$ , where  $\otimes$  and  $\oslash$  are element-wise multiplication and division, respectively. Finally, a technique in NTD [58] can be applied to the update rule of core tensors  $\underline{\mathbf{G}}^{(s)}$  as follow:

$$\underline{\mathbf{G}}^{(s)} \leftarrow \underline{\mathbf{G}}^{(s)} \otimes \underline{\mathbf{Z}}^{(s)} \times \{\mathbf{U}^{(s)T}\} \oslash \underline{\mathbf{G}}^{(s)} \times \{\mathbf{U}^{(s)T} \mathbf{U}^{(s)}\} \tag{5.61}$$

for all  $s$ . Furthermore, we have to be careful for the rotations of updates. In this algorithm, I separate the updates of all common factor matrices and the updates of all individual factor matrices, completely. In the common step, I update all common factor matrices and individual core tensors while individual factor matrices are fixed. In the individual step, I update all individual factor matrices and individual core tensors while common factor matrices are fixed. I iterate both steps until convergence.

Thus, the algorithms for nonnegative LTD can be summarized as Algorithm 26. In this algorithm, I iterate  $K$  times to update by each step. When  $K$  is very small, convergence becomes very slow; however, when  $K$  is too large, calculation time of one cycle becomes very long. For example, I set  $K = 50$  or  $100$ .

---

**Algorithm 26** Nonnegative LTD algorithm

---

```

1: Input:  $\{\underline{Z}^{(s)}\}_{s=1}^S$ ,  $\{L_n, R_n\}_{n=1}^N$ , and  $K$ 
2: Initialize: nonnegative matrices  $\{U_C^{(n)}\}$ ,  $\{U_I^{(n,s)}\}_{s=1}^S\}_{n=1}^N$  and core tensors  $\{\underline{G}^{(s)}\}_{s=1}^S$ 
3: repeat
4:   % common step
5:   for  $k = 1, \dots, K$  do
6:     for  $n = 1, \dots, N$  do
7:       Update  $U_C^{(n)}$  by calculation (5.59) ;
8:     end for
9:     for  $s = 1, \dots, S$  do
10:      Update  $\underline{G}^{(s)}$  by calculation (5.61);
11:    end for
12:  end for
13:  % individual step
14:  for  $s = 1, \dots, S$  do
15:    for  $k = 1, \dots, K$  do
16:      for  $n = 1, \dots, N$  do
17:        Update  $\{U_I^{(n,s)}\}_{s=1}^S$  by calculation (5.60) ;
18:      end for
19:      Update  $\underline{G}^{(s)}$  by calculation (5.61);
20:    end for
21:  end for
22: until  $C(\Theta)$  converge
23: Output:  $\{\underline{G}^{(s)}\}$ ,  $\{U^{(n,s)}\}_{n=1}^N\}_{s=1}^S$ 

```

---

## Chapter 6

# Proposed Classifiers

In this Chapter, I propose new three classification methods: the QCMAP estimation in Section 6.1, SVM with weighted regularization in Section 6.2, and the Chernoff FDA in Section 6.3.

### 6.1 QCMAP

#### 6.1.1 QCMAP estimation

Let  $Q(\mathbf{x})$  be a weight function that satisfies

$$\int_{\mathcal{D}} Q(\mathbf{x})d\mathbf{x} = 1, \quad Q(\mathbf{z}) > 0 \text{ for all } \mathbf{z} \in \mathcal{D}, \quad (6.1)$$

where  $\mathcal{D}$  is the data domain (e.g.,  $\mathbb{R}^d$ ).

The *quadratically constrained MAP* (QCMAP) classifier [113, 109, 111] is then defined as follows:

$$\text{maximize } \sum_{y \in \{+1, -1\}} \int_{\mathcal{D}} P(\mathbf{x}, y) \min(yW'(\mathbf{x}), 1)d\mathbf{x}, \quad (6.2)$$

$$\text{subject to } \int_{\mathcal{D}} Q(\mathbf{x})|W'(\mathbf{x})|^2d\mathbf{x} \leq 1. \quad (6.3)$$

Then the solution of the QCMAP problem is given by

$$W'(\mathbf{x}) \begin{cases} \geq 1 & \text{if } P(+1|\mathbf{x}) > P(-1|\mathbf{x}) = 0 \\ = 1 & \text{if } P(+1|\mathbf{x}) > P(-1|\mathbf{x}) \neq 0 \\ = \zeta(x) & \text{if } P(+1|\mathbf{x}) = P(-1|\mathbf{x}) \\ = -1 & \text{if } P(-1|\mathbf{x}) > P(+1|\mathbf{x}) \neq 0 \\ \leq -1 & \text{if } P(-1|\mathbf{x}) > P(+1|\mathbf{x}) = 0 \end{cases}, \quad (6.4)$$

where  $\zeta(x)$  is an arbitrary number ( $-1 \leq \zeta(x) \leq 1$ ). In addition, if the measure of  $\{x|P(+1|x) = P(-1|x)\}$  is zero (a condition that is satisfied in almost all classification problems), then  $W'(\mathbf{x})$  is given by

$$W'(\mathbf{x}) = \begin{cases} 1 & \text{if } P(+1|\mathbf{x}) > P(-1|\mathbf{x}) \\ -1 & \text{if } P(+1|\mathbf{x}) < P(-1|\mathbf{x}) \end{cases}. \quad (6.5)$$

In both cases,  $W'(\mathbf{x})$  yields the same results as the MAP classifier.

Furthermore, we can interpret that the objective and constraint functions (6.2) and (6.3) play fitting and regularization roles, respectively. Classification criteria typically have a tradeoff parameter between fitting and regularization functions; however our criterion does not need such a tradeoff parameter. This lack of tradeoff is an interesting and important point in using our classifier.

### 6.1.2 Theoretical Interpretation of QCMAPE estimation

As previously stated, the QCMAPE estimation yields the same results as the MAP classifier. In this section, I discuss a theoretical interpretation of the QCMAPE classifier.

Consider the following problem:

$$\text{maximize } \sum_{y \in \{+1, -1\}} \int_{\mathcal{D}} P(\mathbf{x}, y) y W'(\mathbf{x}) d\mathbf{x}, \quad (6.6)$$

$$\text{subject to } \int_{\mathcal{D}} P(\mathbf{x}) |W'(\mathbf{x})|^2 d\mathbf{x} \leq 1. \quad (6.7)$$

Here, the  $\min(\cdot)$  function is removed from the QCMAPE estimation and  $P(\mathbf{x})$  is substituted for  $Q(\mathbf{x})$ . This can be transformed to

$$\max \int_{\mathcal{D}} P(\mathbf{x}) [W'(\mathbf{x}) \{P(+1|\mathbf{x}) - P(-1|\mathbf{x})\}] d\mathbf{x}, \quad (6.8)$$

$$\text{s.t. } \int_{\mathcal{D}} P(\mathbf{x}) |W'(\mathbf{x})|^2 d\mathbf{x} \leq 1. \quad (6.9)$$

From the Cauchy-Schwarz inequality, the solution is given by

$$W'(\mathbf{x}) \propto P(+1|\mathbf{x}) - P(-1|\mathbf{x}). \quad (6.10)$$

This is a direct estimation of a difference of posterior probabilities up to a constant factor, so  $W'(\mathbf{x})$  yields the same result as the MAP classifier.

If  $Q(\mathbf{x}) \neq P(\mathbf{x})$ , the problem can be written as

$$\max \int_{\mathcal{D}} Q(\mathbf{x}) \left[ W'(\mathbf{x}) \frac{P(\mathbf{x})}{Q(\mathbf{x})} \{P(+1|\mathbf{x}) - P(-1|\mathbf{x})\} \right] d\mathbf{x}, \quad (6.11)$$

$$\text{s.t. } \int_{\mathcal{D}} Q(\mathbf{x}) |W'(\mathbf{x})|^2 d\mathbf{x} \leq 1. \quad (6.12)$$

In this case, the solution is given by

$$W'(\mathbf{x}) \propto \frac{P(\mathbf{x})}{Q(\mathbf{x})} \{P(+1|\mathbf{x}) - P(-1|\mathbf{x})\}, \quad (6.13)$$

and  $W'(\mathbf{x})$  too yields the same result as the MAP classifier. I call this the direct difference estimation of the posterior probabilities (DDEP).

Standard MAP estimation provides  $P(+1|\mathbf{x})$  and  $P(-1|\mathbf{x})$ , individually. After estimation U compare these values, but their difference is sufficient for classification. A direct estimation of the difference (i.e., DDEP) may improve the accuracy of estimation as compared with the difference between individual estimations (the standard MAP estimation).

If  $yW'(\mathbf{x}) > 1$ , then it could be clearly discriminative. Because the important region is that where  $yW'(\mathbf{x})$  is small, it is not smart to include large values of  $yW'(\mathbf{x})$  in the evaluation, as large values of  $yW'(\mathbf{x})$  affect the evaluation of the important region. In the QCMAP estimation, I clip the larger parts of  $yW'(\mathbf{x})$  using the cost function  $\min(yW'(\mathbf{x}), 1)$ . Therefore we can consider the QCMAP estimation as a clipped DDEP.

### 6.1.3 Derivation of Mathematical Programming Problem

I define a discriminant model  $D(\mathbf{x})$  as an approximation of  $W'(\mathbf{x})$  that can be expressed as a linear combination of the components of the basis functions  $\phi_i(\mathbf{x})$  and their parameters  $w_i$ :

$$W'(\mathbf{x}) \approx D(\mathbf{x}) := \sum_{i=1}^M w_i \phi_i(\mathbf{x}) = \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}) \rangle, \quad (6.14)$$

where  $\mathbf{w}$  and  $\boldsymbol{\phi}(\mathbf{x})$  are a parameter vector and a basis function vector, respectively. Also assume that  $\{\phi_i(\mathbf{x})\}_{i=1}^M$  is a set of linear independent functions.

Now, I describe QCMAP learning from training samples  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ . Since the QCMAP method includes integration with  $P(\mathbf{x}, y)$  and  $Q(\mathbf{x})$ , finding the optimal value of estimates is nontrivial. The objective part of QCMAP (6.2) can be considered as the maximization of the expected value of a cost function  $\min(yW'(\mathbf{x}), 1)$  with respect to  $(\mathbf{x}, y)$ . Since expectation can be approximated by a sample mean, the method can be simplified through an approximation by replacing the expectation with the sample mean and by substituting Eq. (6.14) into the QCMAP estimation such that

$$\begin{aligned} \sum_{y \in \{+1, -1\}} \int_{\mathcal{D}} P(\mathbf{x}, y) \min(yD(\mathbf{x}), 1) d\mathbf{x} &= E_{\mathbf{x}, y}[\min(y \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}) \rangle, 1)] \\ &\simeq \frac{1}{N} \sum_{n=1}^N \min(y_n \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}_n) \rangle, 1). \end{aligned} \quad (6.15)$$

The left part of Eq. (6.3) can be written as

$$\begin{aligned} \int_{\mathcal{D}} Q(\mathbf{x}) |D(\mathbf{x})|^2 d\mathbf{x} &= \int_{\mathcal{D}} Q(\mathbf{x}) \left| \sum_{i=1}^M w_i \phi_i(\mathbf{x}) \right|^2 d\mathbf{x} \\ &= \sum_{i=1}^M \sum_{j=1}^M w_i w_j \int_{\mathcal{D}} Q(\mathbf{x}) \phi_i(\mathbf{x}) \phi_j(\mathbf{x}) d\mathbf{x} = \mathbf{w}^T \mathbf{H} \mathbf{w}, \end{aligned} \quad (6.16)$$

where the  $(i, j)$  element of matrix  $\mathbf{H}$  is given by

$$\mathbf{H}(i, j) := \int_{\mathcal{D}} Q(\mathbf{x}) \phi_i(\mathbf{x}) \phi_j(\mathbf{x}) d\mathbf{x}. \quad (6.17)$$

We must calculate the integral in (6.17) to obtain  $\mathbf{H}$ . One way to do this is to apply a numerical integration method such as the Monte Carlo method [71], but this approach is computationally expensive. In Sec. 6.1.9 I provide an alternative approach that allows analytical calculation in a special

case. Finally, the QCMAP estimation criterion can be summarized as

$$\text{maximize } \sum_{n=1}^N \min(y_n \langle \mathbf{w}, \phi(\mathbf{x}_n) \rangle, 1), \quad (6.18)$$

$$\text{subject to } \mathbf{w}^T \mathbf{H} \mathbf{w} \leq 1. \quad (6.19)$$

Here, the objective function is a convex piecewise linear function because of the presence of the min function. If  $\mathbf{w}^T \mathbf{H} \mathbf{w} = 0$ , we have  $D(\mathbf{x}) = 0$  for all  $\mathbf{x}$  in  $\mathcal{D}$  because  $Q(\mathbf{x}) > 0$  for all  $\mathbf{x} \in \mathcal{D}$  and we have  $\mathbf{w} = 0$  or  $\phi(\mathbf{x})$  are linearly dependent. Therefore  $\mathbf{H}$  is a positive definite matrix and the constraint function is convex. QCMAP problems are convex and can be solved (trained) by the primal-dual interior point method [122] described in Section 6.1.7.

#### 6.1.4 Relation to Hinge Loss Minimization

Hinge loss is defined as

$$J(\mathbf{w}) = \sum_{n=1}^N \max(1 - y_n \langle \mathbf{w}, \phi(\mathbf{x}_n) \rangle, 0). \quad (6.20)$$

In fact, maximization of the QCMAP cost function and minimization of the hinge loss function are equivalent, since we have

$$\max(1 - r_n, 0) = 1 - \min(r_n, 1), \quad (6.21)$$

where  $r_n := y_n \langle \mathbf{w}, \phi(\mathbf{x}_n) \rangle$ . Fig. 6.1 shows a single term of the QCMAP cost function and the hinge loss function. The vertical axis expresses the evaluation value and horizontal axis expresses the value of  $r_n$ . The QCMAP estimation maximizes the evaluation value, so  $r_n$  seeks to be at least 1. On the other hand, the hinge loss minimization expects the evaluation value to be as small as possible, so again  $r_n$  seeks to be at least 1. The results of both approaches are therefore equal.

#### 6.1.5 Relation to the Tikhonov Regularization

The constraint (6.19) on the QCMAP problem can be considered as a special case of the generalized Tikhonov regularization (TR) [97]. As described above, the matrix  $\mathbf{H}$  is a positive definite matrix, and so can be decomposed as  $\mathbf{H} = \mathbf{\Gamma}^T \mathbf{\Gamma}$ , where  $\mathbf{\Gamma}$  is called the Tikhonov matrix. Moreover, the TR constraint can be written as  $\|\mathbf{\Gamma} \mathbf{w}\|^2 \leq \sigma$ , where  $\sigma$  is a regularization (tradeoff) parameter. The key point of TR is choosing the Tikhonov matrix  $\mathbf{\Gamma}$  and  $\sigma$ , but this selection is nontrivial due to the large number of degrees of freedom.

On the other hand, the key point of the QCMAP approach is to select the weight function  $Q(\mathbf{x})$ . Unfortunately, unlike  $\mathbf{\Gamma}$  we cannot arbitrarily select  $Q(\mathbf{x})$ . For some  $\mathbf{\Gamma}$ , there may not exist a  $Q(\mathbf{x})$  that satisfies condition (6.1). However, we can select  $Q(\mathbf{x})$  instead of  $\mathbf{\Gamma}$  in the TR method because the existence of a  $\mathbf{\Gamma}$  is ensured for any  $Q(\mathbf{x})$  satisfying Eq. (6.1).

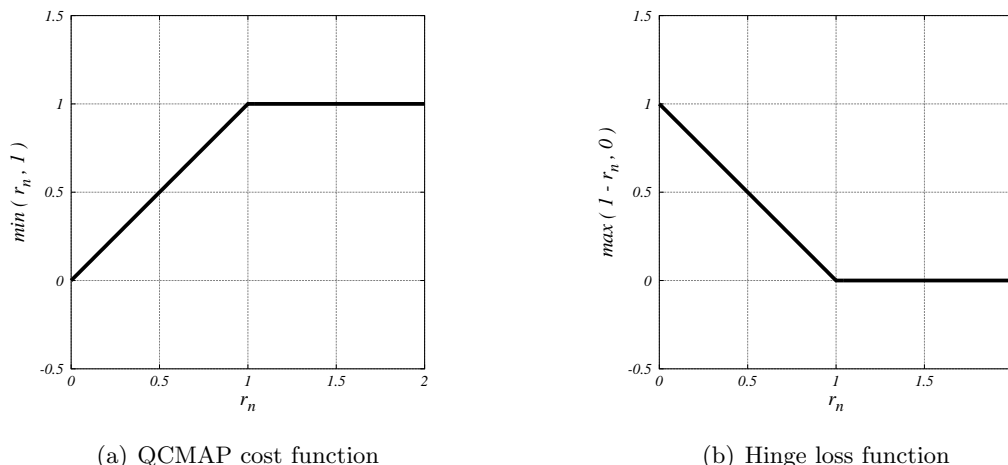


Figure 6.1: The QCMAP cost function and the hinge loss function

We can therefore say that the quadratic constraint of QCMAP is a special case of the TR method. The contribution of our method is that if  $\mathbf{H}$  is obtained by a  $Q(\mathbf{x})$  that satisfies condition (6.1), and if we let  $\sigma = 1$ , then it is sufficient to provide a MAP-based classification.

### 6.1.6 Uniqueness of QCMAP Solutions

I next examine whether a QCMAP problem has a unique solution. In fact, we cannot strictly say that QCMAP problems always have a unique solution. However, the QCMAP problems have unique solutions in many cases.

Fig. 6.2 illustrates two potential solutions to a QCMAP problem. If the hinge loss function cannot be optimized to zero because  $\mathbf{H}$  is a positive definite matrix then the QCMAP problem has a unique solution (Fig. 6.2(a)). On the other hand, if the hinge loss function can be easily optimized to zero, then the solution to the QCMAP problem permits some range within  $\{\mathbf{w} \mid J(\mathbf{w}) = 0 \cap \mathbf{w}^T \mathbf{H} \mathbf{w} \leq 1\}$  (Fig. 6.2(b)). An empirical risk (including hinge loss) of zero implies over-fitting and that the degree of freedom of the model is too high, or that the problem is overly simplistic. In such cases, we might have to reduce the degree of freedom of the model to have a unique solution.

I propose using an approximation function as the hinge loss function in Sec. 6.1.7. In this case, the empirical risk cannot be optimized to zero, and a unique solution is always given.

### 6.1.7 Solution Method for QCMAP and Approximation of Cost Function

To solve the QCMAP problem, I use the primal-dual interior point algorithm [122]. Let  $\mathbf{w} \in \mathbb{R}^M$  be a primal parameter vector and  $f(\mathbf{w})$  be a convex function, and consider the following problem:

$$\text{minimize } f(\mathbf{w}), \tag{6.22}$$

$$\text{subject to } \mathbf{w}^T \mathbf{H} \mathbf{w} \leq 1, \tag{6.23}$$

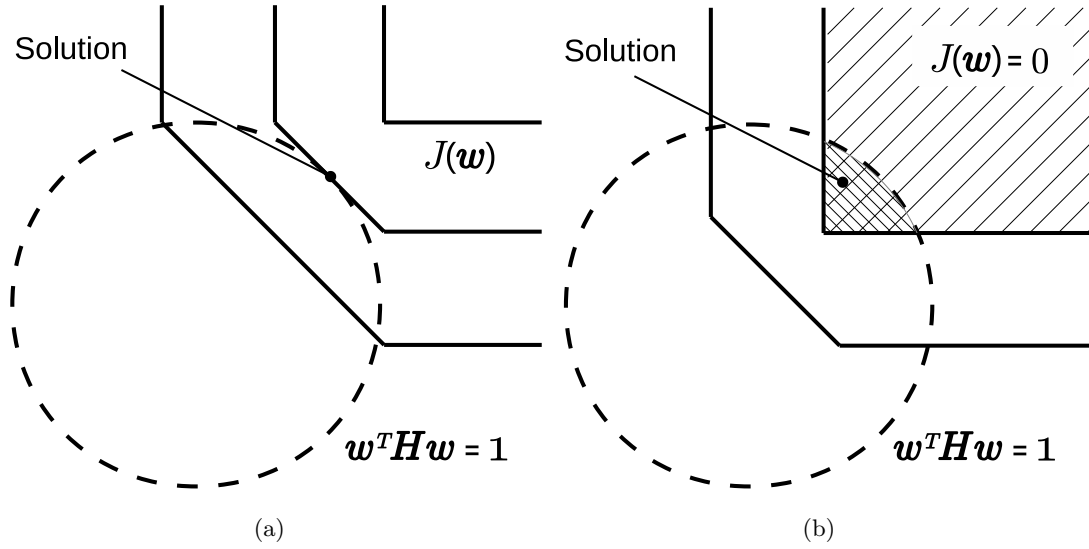


Figure 6.2: Two solution areas for the QCMAP problem: solid and dotted lines depict the contour lines of  $J(\mathbf{w})$  and  $\mathbf{w}^T \mathbf{H} \mathbf{w} = 1$ , respectively

where  $\mathbf{H}$  is an  $(M, M)$ -positive semidefinite matrix. A solution to this problem requires that the function  $f(\mathbf{w})$  is differentiable, yet (6.18) and (6.20) are not differentiable. Therefore, I define  $f(\mathbf{w})$  as an approximation function of (6.20) given by

$$J(\mathbf{w}) \approx f(\mathbf{w}) := \sum_{n=1}^N h(r_n | t), \quad (6.24)$$

where

$$h(r_n | t) := \frac{1}{t} \log(1 + e^{t(1-r_n)}), \quad (6.25)$$

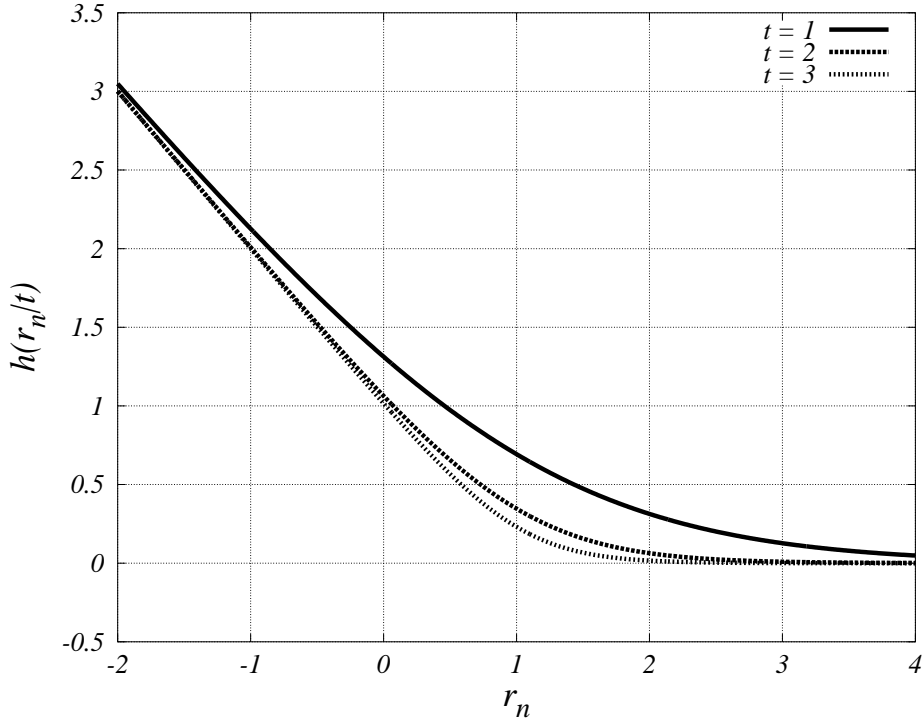
and  $t$  is a parameter [87]. As  $t \rightarrow \infty$ , the objective function  $f(\mathbf{w})$  approaches the hinge loss function. If  $r_n \geq 1$  (i.e.,  $t(1-r_n)$  is zero or negative), then  $\lim_{t \rightarrow \infty} h(r_n | t) = 0$ . If  $r_n < 1$  (i.e.,  $t(1-r_n)$  is positive) and  $t$  is very large, then  $e^{t(1-r_n)} \gg 1$  and  $\log(1 + e^{t(1-r_n)}) \approx \log(e^{t(1-r_n)}) = t(1-r_n)$ . Therefore, we can derive

$$\begin{aligned} \lim_{t \rightarrow \infty} h(r_n | t) &= \lim_{t \rightarrow \infty} \frac{1}{t} \log(1 + e^{t(1-r_n)}) \\ &= \begin{cases} 0 & r_n \geq 1 \\ 1 - r_n & r_n < 1 \end{cases}. \end{aligned} \quad (6.26)$$

This is equivalent to the hinge loss, i.e.,  $\max(1 - r_n, 0)$ .

Smaller values of  $t$  yield smoother approximations of the hinge loss. For example, Fig. 6.3 shows  $h(r_n | t)$  for  $t = 1, 2$ , and 3. Smoothness of the objective function is an important factor for smooth convergence in optimization methods. However, smoothness and sharpness are tradeoffs: sharpness yields a better approximation of the hinge loss but complicates the optimization procedure.

In this experiment, I fixed  $t = 2$  (see Sec. 7.4.4) because empirical investigation showed that this value led to good, smooth optimizations. Methods of selecting the value of  $t$  remains an open problem.

Figure 6.3: Approximation Function  $h(r_n|t)$  for changing  $t$ 

### Primal-Dual Interior Point Algorithm for QCMAP

First, I consider the Lagrange function of this problem:

$$L(\mathbf{w}, z) = f(\mathbf{w}) - z(1 - \mathbf{w}^T \mathbf{H} \mathbf{w}), \quad (6.27)$$

where  $z$  is a dual parameter. From the Karush-Kuhn-Tucker (KKT) conditions, we have

$$\frac{\partial f}{\partial \mathbf{w}} + 2z\mathbf{H}\mathbf{w} = \mathbf{0}, \quad (6.28)$$

$$z(1 - \mathbf{w}^T \mathbf{H} \mathbf{w}) = 0, \quad (6.29)$$

$$z \geq 0, \quad 1 - \mathbf{w}^T \mathbf{H} \mathbf{w} \geq 0. \quad (6.30)$$

A pair  $(\mathbf{w}, z)$  that satisfies condition (6.30) is called a “primal-dual interior point.” Furthermore, if  $(\mathbf{w}, z)$  satisfies Eqs. (6.28) and (6.29), then the pair is the optimal solution of the primal dual problem of QCMAP.

Thus, I define the residual function as

$$\mathbf{r}(\mathbf{w}, z) := \begin{pmatrix} \frac{\partial f}{\partial \mathbf{w}} + 2z\mathbf{H}\mathbf{w} \\ z(1 - \mathbf{w}^T \mathbf{H} \mathbf{w}) \end{pmatrix}, \quad (6.31)$$

and the goal of the solution algorithm is to optimize  $(\mathbf{w}, z)$  such that

$$\|\mathbf{r}(\mathbf{w}, z)\| \rightarrow 0, \quad (6.32)$$

**Algorithm 27** Primal-dual interior point (PDIP) algorithm

---

Input: training samples,  $\mathbf{H}$ ,  $T$ ,  $I_{\max}$  and  $\alpha_{\min}$ ;  
Initialize:  $\mathbf{w}_0 = \mathbf{0}$ ,  $z_0 = 1$ ,  $\alpha = 1$  and  $k = 0$ ;  
**while**  $\frac{\|\mathbf{r}(\mathbf{w}_k, z_k)\|}{M+1} > T$  &  $k < I_{\max}$  &  $\alpha > \alpha_{\min}$  **do**  
    Update  $(\mathbf{w}_{k+1}, z_{k+1})$  by (6.36), (6.37), and Algorithm 28;  
     $k = k + 1$ ;  
**end while**  
Output:  $\mathbf{w}_k$ ;

---

while satisfying Eq. (6.30). In the actual algorithm, I use a tolerance  $T$  (e.g.,  $T = 10^{-5}$ ) as a threshold for convergence of the residual function. Specifically, the convergence condition is given by

$$\frac{\|\mathbf{r}(\mathbf{w}, z)\|}{M+1} < T. \quad (6.33)$$

The basic concept of the primal-dual interior point algorithm is as follows: First, I set the values of an initial pair  $(\mathbf{w}_0, z_0)$  such that the pair is an interior point, for example,  $\mathbf{w}_0 = \mathbf{0}$  and  $z_0 = 1$ . Next, I iteratively update  $(\mathbf{w}_k, z_k)$  through

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \Delta \mathbf{w}, \quad (6.34)$$

$$z_{k+1} = z_k - \alpha \Delta z. \quad (6.35)$$

Here,  $(\Delta \mathbf{w}, \Delta z)$  is the direction for the update and  $\alpha$  is the step size. Once the residual function  $\mathbf{r}(\mathbf{w}_k, z_k)$  converges, the algorithm is completed. An outline of the primal-dual interior point (PDIP) algorithm is described in Algorithm 27.

In this regard, I set the maximum number of iterations  $I_{\max}$  and the minimum step size  $\alpha_{\min}$  (e.g.,  $I_{\max} = 500$  and  $\alpha_{\min} = 10^{-6}$ ).

The key point in executing this algorithm is correctly determining  $(\Delta \mathbf{w}, \Delta z)$  and  $\alpha$ . Many methods and algorithms have been proposed and developed for doing so. In this paper, I introduce an interior path-following method to decide  $(\Delta \mathbf{w}, \Delta z)$  and a back-tracking search algorithm to decide  $\alpha$ .

### Path-Following Method

The interior path-following method provides a  $(\Delta \mathbf{w}, \Delta z)$  that approaches the optimal solution following the analytic center path [41]. In this method,  $(\Delta \mathbf{w}, \Delta z)$  is given by

$$\Delta \mathbf{w} = \left( A + \frac{z_k}{c} \mathbf{b} \mathbf{b}^T \right)^{-1} \left( \mathbf{g} + \frac{\mu}{c} \mathbf{b} \right), \quad (6.36)$$

$$\Delta z = z_k - \frac{\mu}{c} + \frac{z_k}{c} \mathbf{b}^T \Delta \mathbf{w}_k, \quad (6.37)$$

**Algorithm 28** Back tracking search algorithm

---

Input:  $0 < \alpha_{\min} < \alpha_{\max}$ ,  $0 < \beta < 1$ ,  $(\mathbf{w}, z)$  and  $(\Delta\mathbf{w}, \Delta z)$ ;
Initialize:  $\alpha = \alpha_{\max}$ ; $(\mathbf{w}', z') = (\mathbf{w}, z) - \alpha(\Delta\mathbf{w}, \Delta z)$ ;**repeat** $\alpha \leftarrow \beta\alpha$ ; $(\mathbf{w}', z') = (\mathbf{w}, z) - \alpha(\Delta\mathbf{w}, \Delta z)$ ;**until**  $(\mathbf{w}', z')$  satisfies conditions (6.30) and (6.42), or  $\alpha < \alpha_{\min}$ Output:  $\alpha$ ,  $(\mathbf{w}', z')$ ;

where

$$A = \nabla^2 f(\mathbf{w}_k) + 2z_k \mathbf{H}, \quad \mathbf{b} = 2\mathbf{H}\mathbf{w}_k, \quad (6.38)$$

$$c = 1 - \mathbf{w}_k^T \mathbf{H}\mathbf{w}_k, \quad \mathbf{g} = \frac{\partial}{\partial \mathbf{w}} f(\mathbf{w}_k), \quad (6.39)$$

$$\mu = \gamma z_k (1 - \mathbf{w}_k^T \mathbf{H}\mathbf{w}_k), \quad (6.40)$$

$$\gamma = \min \left( \sqrt{\frac{\|\mathbf{r}(\mathbf{w}_k, z_k)\|}{M+1}}, 0.5 \right), \quad (6.41)$$

and  $\mu$  is a barrier parameter to prevent  $(\mathbf{w}, z)$  from falling outside the region of interior points.

**Backtracking Search Algorithm**

In the backtracking search algorithm, I iteratively decrease the value of  $\alpha_{\min} \leq \alpha \leq \alpha_{\max}$  until it satisfies condition (6.30) and

$$\|\mathbf{r}_\mu(\mathbf{w}', z')\| \leq (1 - \alpha(1 - \eta))\|\mathbf{r}_\mu(\mathbf{w}, z)\|, \quad (6.42)$$

where

$$\mathbf{r}_\mu(\mathbf{w}, z) := \begin{pmatrix} \frac{\partial f}{\partial \mathbf{w}} + 2z\mathbf{H}\mathbf{w} \\ z(1 - \mathbf{w}^T \mathbf{H}\mathbf{w}) - \mu \end{pmatrix}, \quad (6.43)$$

$$\eta = \min \left( \frac{\|\mathbf{r}_\mu(\mathbf{w}, z)\|}{M+1}, 0.25 \right). \quad (6.44)$$

Algorithm 28 gives pseudocode for this search method, where values for  $\alpha_{\min}$ ,  $\alpha_{\max}$ , and  $\beta$  are selected a priori, for example,  $\alpha_{\min} = 10^{-6}$ ,  $\alpha_{\max} = 0.995$ ,  $\beta = 0.75$ .

**6.1.8 Equivalence of QCMAF to LSR and SVM**

In this section, I briefly discuss the equivalence of the QCMAF estimation with LSR and the linear SVM. The QCMAF estimation has a weight function  $Q(\mathbf{x})$ , and an LSR or the linear SVM can be derived from the QCMAF estimation by selecting appropriate functions for  $Q(\mathbf{x})$ .

### Equivalence to LSR

I first present the equivalence to LSR. The solution of an LSR problem,  $\hat{\mathbf{w}}_{\text{LSR}}$ , is given by

$$\hat{\mathbf{w}}_{\text{LSR}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}, \quad (6.45)$$

where

$$\Phi := \begin{pmatrix} \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_N) \end{pmatrix}^T, \quad (6.46)$$

$$\mathbf{y} := \begin{pmatrix} y_1 & y_2 & \cdots & y_N \end{pmatrix}^T. \quad (6.47)$$

If we choose  $Q_A(\mathbf{x}) = P(\mathbf{x})$  and the expectation is approximated by the sample mean, then the constraint matrix is given by

$$\begin{aligned} \mathbf{H}_A(i, j) &= \int_{\mathcal{D}} P(\mathbf{x}) \phi(\mathbf{x}) \phi(\mathbf{x})^T d\mathbf{x} = \mathbb{E}_{\mathbf{x}}[\phi(\mathbf{x}) \phi(\mathbf{x})^T] \\ &\simeq \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T = \frac{1}{N} \Phi^T \Phi. \end{aligned} \quad (6.48)$$

The constraint is then written in the form

$$\mathbf{w}^T \Phi^T \Phi \mathbf{w} \leq N. \quad (6.49)$$

Therefore, the QCMAP problem with  $Q_A(\mathbf{x})$  is given by

$$\text{maximize} \quad \sum_{n=1}^N \min(y_n \langle \mathbf{w}, \phi(\mathbf{x}_n) \rangle, 1), \quad (6.50)$$

$$\text{subject to} \quad \mathbf{w}^T \Phi^T \Phi \mathbf{w} \leq N. \quad (6.51)$$

**Theorem 2.** *The solution of a QCMAP problem with  $Q_A(\mathbf{x})$  is equivalent to  $\hat{\mathbf{w}}_{\text{LSR}}$ .*

### Proof of Theorem 2

I prove that  $\hat{\mathbf{w}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$  is the solution of the following optimization problem:

$$\text{maximize} \quad \sum_{n=1}^N \min(y_n \langle \mathbf{w}, \phi(\mathbf{x}_n) \rangle, 1), \quad (6.52)$$

$$\text{subject to} \quad \mathbf{w}^T \Phi^T \Phi \mathbf{w} \leq N, \quad (6.53)$$

where

$$\Phi := \begin{pmatrix} \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_N) \end{pmatrix}^T, \quad (6.54)$$

$$\mathbf{y} := \begin{pmatrix} y_1 & y_2 & \cdots & y_N \end{pmatrix}^T. \quad (6.55)$$

*Proof.* First, this objective can be rewritten using slack variables  $\xi_n$ :

$$\text{minimize } \sum_{n=1}^N \xi_n, \quad (6.56)$$

$$\text{subject to } y_n \langle \mathbf{w}, \phi(\mathbf{x}_n) \rangle \geq 1 - \xi_n, \quad (6.57)$$

$$\xi_n \geq 0, \quad n = 1, \dots, N, \quad (6.58)$$

$$\mathbf{w}^T \Phi^T \Phi \mathbf{w} \leq 1. \quad (6.59)$$

Next I derive its dual problem in terms of a Lagrange function with KKT conditions.

The corresponding Lagrangian of this problem is then

$$\begin{aligned} L(\boldsymbol{\xi}, \mathbf{w}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \eta) &= \sum_{n=1}^N \xi_n - \sum_{n=1}^N \beta_n \xi_n \\ &\quad - \sum_{n=1}^N \gamma_n (\xi_n - 1 + y_n \mathbf{w}^T \phi(\mathbf{x}_n)) - \frac{\eta}{2} (N - \mathbf{w}^T \Phi^T \Phi \mathbf{w}), \end{aligned} \quad (6.60)$$

where  $\beta_n$ ,  $\gamma_n$ , and  $\eta$  are Lagrange coefficients. From the KKT conditions we have

$$\frac{\partial L}{\partial \xi_n} = 1 - \beta_n - \gamma_n = 0, \quad (6.61)$$

$$\frac{\partial L}{\partial w_i} = \sum_{n=1}^N \phi_i(\mathbf{x}_n) (\eta \mathbf{w}^T \phi(\mathbf{x}_n) - y_n \gamma_n) = 0, \quad (6.62)$$

$$\beta_n \xi_n = 0, \quad \gamma_n (\xi_n - 1 + y_n \mathbf{w}^T \phi(\mathbf{x}_n)) = 0, \quad (6.63)$$

$$\frac{\eta}{2} (N - \mathbf{w}^T \Phi^T \Phi \mathbf{w}) = 0, \quad (6.64)$$

$$\beta_n \geq 0, \quad \gamma_n \geq 0, \quad \eta \geq 0, \quad n = 1, \dots, N. \quad (6.65)$$

Thus, the dual problem can be expressed as

$$\text{maximize } F(\boldsymbol{\gamma}, \eta) = -\frac{1}{2\eta} \sum_{n=1}^N \gamma_n^2 + \sum_{n=1}^N \gamma_n - \frac{\eta}{2} N,$$

$$\text{subject to } 0 \leq \gamma_n \leq 1, \quad \eta \geq 0, \quad n = 1, \dots, N.$$

The solution of this problem is given by

$$\frac{\partial F}{\partial \gamma_n} = -\frac{1}{\eta} \gamma_n + 1 = 0 \quad \Rightarrow \quad \gamma_n = \eta. \quad (6.66)$$

The optimal parameter vector  $\hat{\mathbf{w}}$  from Eq. (6.62) is then

$$\begin{aligned} \sum_{n=1}^N \phi(\mathbf{x}_n) (\eta \mathbf{w}^T \phi(\mathbf{x}_n) - y_n \eta) &= \mathbf{0} \quad (\text{since } \gamma_n = \eta) \\ \eta \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \mathbf{w} &= \eta \sum_{n=1}^N \phi(\mathbf{x}_n) y_n \\ \Phi^T \Phi \mathbf{w} &= \Phi^T \mathbf{y} \\ \Rightarrow \hat{\mathbf{w}} &= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}. \end{aligned} \quad (6.67)$$

□

### Equivalence to SVM

Here, I show the relation between the QCMAP estimation and the linear SVM. To describe the key points first, the QCMAPs with uniform weight function are equivalent to a special case of SVM in a linear discriminant model  $D(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle$ . Similarly, if we use some orthonormal basis functions and a uniform weight function for the QCMAP, it is equivalent to a special case of SVM. Unfortunately, the QCMAP is not equivalent to SVM in general, non-linear models. For example, we cannot claim equivalence to SVM when we use a Gaussian kernel model. However, we can show that the QCMAP is similar to the SVM when a weighting function reaches a uniform function.

**Theorem 3.** *Assume that  $\|\mathbf{x}\| = 1$  and the data domain  $\mathcal{D}$  is equal to the unit  $(d-1)$ -dimensional hyperspherical surface  $S^{d-1}$  (i.e.,  $S^{d-1} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| = 1\}$ ). Choose  $Q_B(\mathbf{x}) = 1/S_{d-1}$ , where  $S_{d-1}$  is the surface area of  $S^{d-1}$ . Then condition (6.1) and*

$$\int_{\mathcal{D}} Q_B(\mathbf{x}) |\langle \mathbf{w}, \mathbf{x} \rangle|^2 d\mathbf{x} = \frac{1}{d} \|\mathbf{w}\|^2 \quad (6.68)$$

are satisfied.

### Proof of Theorem 3

I prove that if we assume  $\mathbf{x} \in S^{d-1}$ ,  $Q_B(\mathbf{x}) := \frac{1}{S_{d-1}}$ , where  $S^{d-1}$  presents the unit  $(d-1)$ -dimensional hyper-spherical surface (i.e.,  $S^{d-1} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| = 1\}$ ), and  $S_{d-1}$  is the surface area of  $S^{d-1}$ , then it is satisfied that

$$\int_{S^{d-1}} Q_B(\mathbf{x}) d\mathbf{x} = 1, \quad Q_B(\mathbf{z}) > 0 \text{ for all } \mathbf{z} \in \mathcal{D}, \quad (6.69)$$

$$\int_{S^{d-1}} Q_B(\mathbf{x}) |\langle \mathbf{w}, \mathbf{x} \rangle|^2 d\mathbf{x} = \frac{1}{d} \|\mathbf{w}\|^2. \quad (6.70)$$

*Proof.* First, we can derive

$$\int_{S^{d-1}} Q_B(\mathbf{x}) d\mathbf{x} = \frac{1}{S_{d-1}} \int_{S^{d-1}} dV_{d-1}(\mathbf{x}) = 1, \quad (6.71)$$

where  $dV_{d-1}(\mathbf{x})$  is the volume element of  $S^{d-1}$  and is given by

$$\begin{aligned} dV_{M-1} &= \sin^{M-2}(\psi_1) \sin^{M-3}(\psi_2) \\ &\quad \cdots \sin(\psi_{M-2}) d\psi_1 d\psi_2 \cdots d\psi_{M-1}. \end{aligned} \quad (6.72)$$

Then,  $Q_B(\mathbf{x})$  satisfies eq. (6.69).

Furthermore, we have

$$\int_{S^{d-1}} \frac{1}{S_{d-1}} |\langle \mathbf{w}, \mathbf{x} \rangle|^2 d\mathbf{x} = \frac{\|\mathbf{w}\|^2}{S_{d-1}} \int_{S^{d-1}} \cos^2 \theta dV_{d-1}(\mathbf{x}), \quad (6.73)$$

where  $\theta$  is the angle between  $\mathbf{w}$  and  $\mathbf{x}$ . We can assume  $\psi_1 = \theta$  because the integral is isometrically defined. Thus, Eq. (6.73) is rewritten as

$$\begin{aligned} & \frac{\|\mathbf{w}\|^2}{S_{d-1}} \int_0^\pi \cdots \int_0^\pi \int_0^{2\pi} \cos^2 \theta dV_{d-1} \\ &= \frac{\|\mathbf{w}\|^2}{S_{d-1}} \int_0^\pi \cdots \int_0^\pi \int_0^{2\pi} \cos^2 \theta \sin^{d-2}(\theta) d\theta dV_{d-2} \\ &= \|\mathbf{w}\|^2 \frac{S_{d-2}}{S_{d-1}} \int_0^\pi \cos^2 \theta \sin^{d-2}(\theta) d\theta = \frac{1}{d} \|\mathbf{w}\|^2. \end{aligned} \quad (6.74)$$

□

From Theorem 3, the constraint of QCMAP with  $Q_B(\mathbf{x})$  is given by  $\|\mathbf{w}\|^2 \leq d$ , and the constraint matrix becomes  $\mathbf{H}_B = \frac{1}{d} \mathbf{I}_d$ , where  $\mathbf{I}_d \in \mathbb{R}^{d \times d}$  is an identity matrix. In this case, the criterion of QCMAP estimation with  $Q_B(\mathbf{x})$  is given by

$$\text{minimize } J(\mathbf{w}), \quad (6.75)$$

$$\text{subject to } \|\mathbf{w}\|^2 \leq d. \quad (6.76)$$

On the other hand, the criterion of SVM estimation is given by

$$\text{minimize } J(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2, \quad (6.77)$$

where  $\lambda$  is a regularization parameter. The differences between the two criteria are only whether  $\|\mathbf{w}\|^2$  is in the objective or the constraint part of the criterion, and whether a regularization parameter exists. Because the QCMAP estimation does not have a regularization parameter, we are unable to state that the two criteria are strictly equivalent. However, the resulting effects of the two criteria are equivalent: they both minimize  $J(\mathbf{x})$  and regularize  $\|\mathbf{w}\|^2$ .

Here, I prove that QCMAP and SVM estimators are equivalent under the condition  $\|\mathbf{w}\|^2 = d$ . First, we have that the solution of QCMAP is always included in an area where  $\|\mathbf{w}\|^2 = d$  is satisfied. Next, I assume that  $\|\mathbf{w}\|^2 = d$  is satisfied in the SVM estimation. Since  $\frac{\lambda}{2} \|\mathbf{w}\|^2$  is now constant, it can be disregarded. In this case, the criterion of SVM estimation can be expressed by minimization of  $J(\mathbf{w})$ , s.t.  $\|\mathbf{w}\|^2 = d$ . Therefore, the QCMAP estimation with  $Q_B(\mathbf{x})$  is realized by an SVM when its solution satisfies  $\|\mathbf{w}\|^2 = d$ . If we choose  $\lambda$  such that  $\|\mathbf{w}\|^2 = d$  in SVM, then SVM and QCMAP estimators are equivalent. Specifically, an SVM with  $\|\mathbf{w}\| = d$  can be regarded as a MAP-based classifier.

Next, I consider a set of orthonormal basis functions  $\{\phi_i(\mathbf{x})\}_{i=1}^M$  as

$$\int_{\mathcal{D}} \phi_i(\mathbf{x}) \phi_j(\mathbf{x}) d\mathbf{x} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}. \quad (6.78)$$

If we assume that  $\Omega$  is the volume of  $\mathcal{D}$ , which is bounded, and uniform weight function  $Q(\mathbf{x}) = \frac{1}{\Omega}$  exists, then we have

$$\frac{1}{\Omega} \sum_{i=1}^M \sum_{j=1}^M w_i w_j \int_{\mathcal{D}} \phi_i(\mathbf{x}) \phi_j(\mathbf{x}) d\mathbf{x} = \frac{1}{\Omega} \|\mathbf{w}\|^2. \quad (6.79)$$

Similar to the previous description, this is equivalent to a special case of SVM.

Next, I consider a Gaussian kernel model as

$$D(\mathbf{x}) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x}), \quad (6.80)$$

where  $k(\mathbf{x}, \mathbf{y}) = \exp(-\gamma\|\mathbf{x} - \mathbf{y}\|^2)$  and  $\alpha_n$  are parameters. In the kernel method, the criterion of SVM is given by

$$\text{minimize } J(\boldsymbol{\alpha}) + \frac{\lambda}{2} \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha}, \quad (6.81)$$

where  $\mathbf{K}(i, j) = k(\mathbf{x}_i, \mathbf{x}_j)$ . Assume that a uniform weight function  $Q(\mathbf{x}) = \frac{1}{\Omega}$  can be defined, then the constraint matrix is given by

$$H(i, j) = \frac{1}{\Omega} \int_{\mathcal{D}} k(\mathbf{x}_i, \mathbf{x}) k(\mathbf{x}_j, \mathbf{x}) d\mathbf{x} \quad (6.82)$$

$$= \frac{1}{\Omega} \left( \frac{\pi}{2\gamma} \right)^{\frac{d}{2}} \sqrt{k(\mathbf{x}_i, \mathbf{x}_j)}. \quad (6.83)$$

In SVM with a Gaussian kernel, the regularization term depends on  $k(\mathbf{x}_i, \mathbf{x}_j)$ . On the other hand, the constraint of QCMAP depends on  $\sqrt{k(\mathbf{x}_i, \mathbf{x}_j)}$ . The two approaches are different, but they are similar. This is indicating that a QCMAP with Gaussian kernel model and the uniform weight function is similar to the kernel SVM.

### 6.1.9 Gaussian QCM Classifier

I next introduce the Gaussian QCM (GQCM) classifier. The key points are that we choose a pdf with a normal distribution for  $Q(\mathbf{x})$ , then fit a Gaussian kernel model to  $D(\mathbf{x})$ . Many applications using the Gaussian kernel model are studied and developed for various research topics [76, 81]. As a result, we can analytically calculate the integral to obtain  $\mathbf{H}$  in (6.3).

#### Construction of GQCM Classifiers

In this section I consider a concrete definition of the discriminant function, given by

$$D(\mathbf{x}) := \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x}) + \beta, \quad (6.84)$$

where  $k(\mathbf{x}, \mathbf{y}) = \exp(-\gamma\|\mathbf{x} - \mathbf{y}\|^2)$  and,  $\alpha_n$  and  $\beta$  are parameters. We can regard this definition (6.84) as the inner product between a parameter vector and a basis function vector, redefining it as

$$\mathbf{w} := \left( \alpha_1 \quad \cdots \quad \alpha_N \quad \beta \right)^T, \quad (6.85)$$

$$\boldsymbol{\phi}(\mathbf{x}) := \left( k(\mathbf{x}_1, \mathbf{x}) \quad \cdots \quad k(\mathbf{x}_N, \mathbf{x}) \quad 1 \right)^T. \quad (6.86)$$

Therefore, we have

$$\langle \mathbf{w}, \phi(\mathbf{x}) \rangle = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x}) + \beta. \quad (6.87)$$

Moreover, if we use the following weight function

$$\begin{aligned} Q(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) &:= N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= \frac{1}{(\sqrt{2\pi})^d \sqrt{|\boldsymbol{\Sigma}|}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}, \end{aligned} \quad (6.88)$$

then the constraint matrix  $\mathbf{H}$  is given by

$$\mathbf{H}(i, j) := \int_{\mathcal{D}} \phi_i(\mathbf{x}) \phi_j(\mathbf{x}) N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}. \quad (6.89)$$

QCMAP estimation using Eq. (6.89) thus provides the GQCM classifiers. To calculate  $\mathbf{H}$ , we must compute the integral given in Eq. (6.89). However, if Gaussian kernels are employed as basis functions and  $Q(\mathbf{x})$  is still a normal distribution in the original space, we can analytically calculate  $\mathbf{H}$ .

### Analytical Calculation of Constraint Matrix

Since the basis functions are either Gaussian kernels or 1,  $\mathbf{H}$  is given by

$$\mathbf{H} = \begin{pmatrix} \mathbf{U} & \mathbf{u} \\ \mathbf{u}^T & u \end{pmatrix}, \quad (6.90)$$

where,  $\mathbf{U}$ ,  $\mathbf{u}$ , and  $u$  are an  $(N \times N)$ -matrix, an  $N$ -dimensional vector, and a value, respectively:

$$\mathbf{U}(i, j) = \int_{\mathcal{D}} N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) k(\mathbf{x}_i, \mathbf{x}) k(\mathbf{x}_j, \mathbf{x}) d\mathbf{x}, \quad (6.91)$$

$$\mathbf{u}(i) = \int_{\mathcal{D}} N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) k(\mathbf{x}_i, \mathbf{x}) d\mathbf{x}, \quad (6.92)$$

$$u = \int_{\mathcal{D}} N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} = 1. \quad (6.93)$$

By using the general equation of a Gaussian integral [92, 16],

$$\int e^{-\frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}} d\mathbf{x} = \sqrt{\frac{(2\pi)^d}{|\mathbf{A}|}} e^{\frac{1}{2}\mathbf{b}^T \mathbf{A}^{-1} \mathbf{b}}, \quad (6.94)$$

we can analytically calculate Eqs. (6.91) and (6.92) as follows. From Eq. (6.91)

$$\mathbf{U}(i, j) = \frac{1}{\sqrt{|4\gamma\boldsymbol{\Sigma} + \mathbf{I}_d|}} \exp\left(\frac{1}{2}\mathbf{b}_{ij}^T \mathbf{A}^{-1} \mathbf{b}_{ij} + c_{ij}\right), \quad (6.95)$$

where

$$\mathbf{A} = 4\gamma\mathbf{I}_d + \boldsymbol{\Sigma}^{-1}, \quad (6.96)$$

$$\mathbf{b}_{ij} = 2\gamma(\mathbf{x}_i + \mathbf{x}_j) + \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}, \quad (6.97)$$

$$c_{ij} = -\gamma(\|\mathbf{x}_i\|^2 + \|\mathbf{x}_j\|^2) - \frac{1}{2}\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}. \quad (6.98)$$

The  $N$ -dimensional vector  $\mathbf{u}$  is given by

$$\mathbf{u}(i) = \frac{1}{\sqrt{|2\gamma\mathbf{\Sigma} + \mathbf{I}_d|}} \exp\left(\frac{1}{2}\mathbf{b}'_i{}^T \mathbf{A}'^{-1} \mathbf{b}'_i + c'_i\right), \quad (6.99)$$

where

$$\mathbf{A}' = 2\gamma\mathbf{I}_d + \mathbf{\Sigma}^{-1}, \quad (6.100)$$

$$\mathbf{b}'_i = 2\gamma\mathbf{x}_i + \mathbf{\Sigma}^{-1}\boldsymbol{\mu}, \quad (6.101)$$

$$c'_i = -\gamma\|\mathbf{x}_i\|^2 - \frac{1}{2}\boldsymbol{\mu}^T \mathbf{\Sigma}^{-1}\boldsymbol{\mu}. \quad (6.102)$$

### 6.1.10 Extended GQCM Classifiers

The weight functions used thus far are a pdf of  $\mathbf{x}$ , a uniform distribution, and a normal distribution. However, we do not know what constitutes a good weight function. To ascertain the ideal weight function, experiments involving various new weight functions are necessary.

#### Mixture of Gaussians Model

First, I propose to use the mixture of Gaussians model for  $Q(\mathbf{x})$  given by

$$Q(\mathbf{x}|\rho_1, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \dots, \rho_L, \boldsymbol{\mu}_L, \boldsymbol{\Sigma}_L) := \sum_{l=1}^L \rho_l N(\mathbf{x}|\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l), \quad (6.103)$$

where  $\rho_l$  are weight parameters that satisfy

$$0 \leq \rho_l \leq 1, \quad \sum_{l=1}^L \rho_l = 1. \quad (6.104)$$

From the equality in (6.89),  $\mathbf{H}$  can be calculated analytically, and

$$\mathbf{H}(i, j) = \int_{\mathcal{D}} \phi_i(\mathbf{x})\phi_j(\mathbf{x}) \sum_{l=1}^L \rho_l N(\mathbf{x}|\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l) d\mathbf{x} \quad (6.105)$$

$$= \sum_{l=1}^L \rho_l \int_{\mathcal{D}} \phi_i(\mathbf{x})\phi_j(\mathbf{x}) N(\mathbf{x}|\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l) d\mathbf{x} \quad (6.106)$$

$$= \sum_{l=1}^L \rho_l \mathbf{H}_l(i, j), \quad (6.107)$$

where

$$\mathbf{H}_l(i, j) = \int_{\mathcal{D}} \phi_i(\mathbf{x})\phi_j(\mathbf{x}) N(\mathbf{x}|\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l) d\mathbf{x}. \quad (6.108)$$

QCMAP estimation with (6.107) therefore provides a novel classifier, which I term the ‘‘mixture of Gaussians QCMAP’’ (MOGQCM) classifier. The problem in using this method is then how to select values for parameters  $L$ ,  $\rho_l$ ,  $\boldsymbol{\mu}_l$ , and  $\boldsymbol{\Sigma}_l$ .

I propose to choose parameters  $L$ ,  $\rho_l$ ,  $\boldsymbol{\mu}_l$ , and  $\boldsymbol{\Sigma}_l$  based on the strategy of the kernel density estimation. Using the kernel density function as the weight function gives

$$Q(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N N(\mathbf{x}|\mathbf{x}_n, s\mathbf{I}_d). \quad (6.109)$$

In other words, I choose  $L := N$ ,  $\rho_l := \frac{1}{N}$ ,  $\boldsymbol{\mu}_l = \mathbf{x}_l$ , and  $\boldsymbol{\Sigma}_l = s\mathbf{I}_d$ , for  $l = 1, \dots, N$ .

This classifier is similar to an LSR for very small  $s$ , since we can prove that if  $s \rightarrow 0$ , then

$$\mathbf{H}_{\text{kdw}}(i, j) \rightarrow \frac{1}{N} \sum_{n=1}^N \phi_i(\mathbf{x}_n)\phi_j(\mathbf{x}_n). \quad (6.110)$$

Conversely, this classifier is similar to an SVM for very large  $s$ , because  $Q(\mathbf{x})$  then approximates a uniform distribution.

I call this the ‘‘KDEQCM’’ classifier.

### Difference of Gaussians Function

In general, a mixture of Gaussian distributions must satisfy condition (6.104). However, we can break condition (6.104) when condition (6.1) would still hold. As a result, I define the difference of Gaussians (DOG) function as

$$Q(\mathbf{x}) := \left(1 + \frac{\kappa}{\nu^d - 1}\right) N(\mathbf{x}|\boldsymbol{\mu}, \nu^2\boldsymbol{\Sigma}) - \frac{\kappa}{\nu^d - 1} N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad (6.111)$$

where  $0 < \kappa < 1$  and  $\nu > 1$ . In general, the DOG function is often used as a window function for an edge extraction filter [27]; however, here I use the DOG function in a very different way (as a weight function).

If we assume that  $\boldsymbol{\Sigma}$  is a diagonal matrix, then this weight function always satisfies (6.1). Fig. 6.4 shows examples of such weight functions, and we see that if  $\nu$  is increased, the weight function has a smoother and wider distribution. It is effective to set the value of  $\kappa$  to be near 1 (e.g.,  $\kappa = 0.9$ ), and if  $\kappa$  is small, the function becomes similar to a Gaussian distribution. The associated constraint matrix  $\mathbf{H}$  is given by

$$\mathbf{H}(i, j) = \left(1 + \frac{\kappa}{\nu^d - 1}\right) \mathbf{H}_1(i, j) - \frac{\kappa}{\nu^d - 1} \mathbf{H}_2(i, j), \quad (6.112)$$

where,

$$\mathbf{H}_1(i, j) = \int_{\mathcal{D}} \phi_i(\mathbf{x})\phi_j(\mathbf{x})N(\mathbf{x}|\boldsymbol{\mu}, \nu^2\boldsymbol{\Sigma})d\mathbf{x}, \quad (6.113)$$

$$\mathbf{H}_2(i, j) = \int_{\mathcal{D}} \phi_i(\mathbf{x})\phi_j(\mathbf{x})N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})d\mathbf{x}. \quad (6.114)$$

Hence, MAP estimation with (6.112) also provides a novel classifier, which I term the ‘‘difference of Gaussians QCMAP’’ (DOGQCM) classifier.

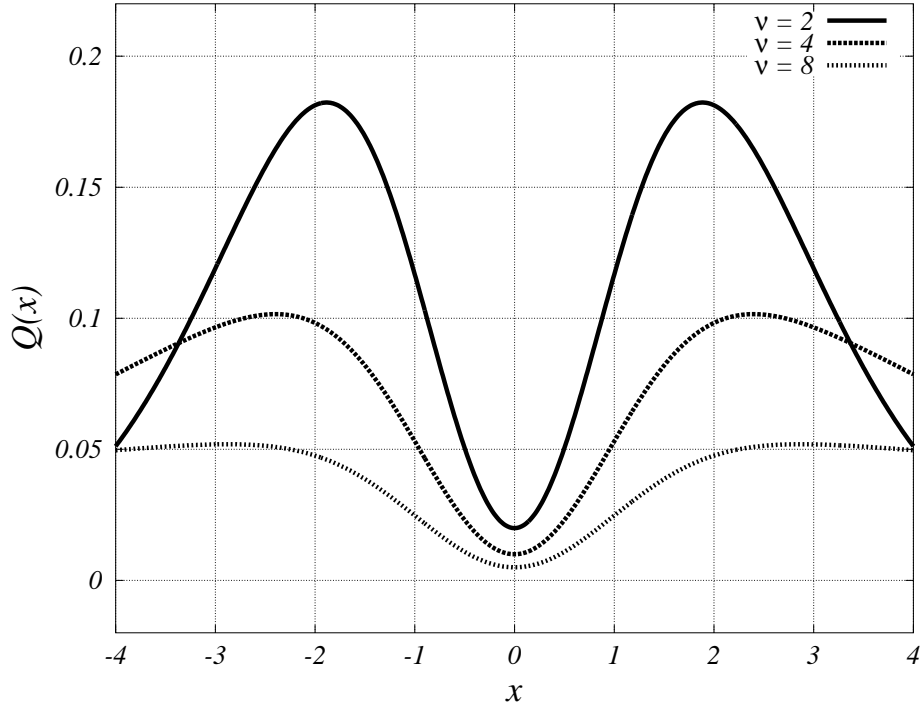


Figure 6.4: Difference of Gaussian function: Here,  $\kappa = 0.9$  is fixed and  $\nu$  varies.

Regularization is not so important in the central area of distribution if the samples are normally distributed because we have a sufficient number of samples. Conversely, regularization is very important in the peripheral area of the distribution because we have only a small number of samples. Therefore, the DOG weight function is effective in such a case.

### Weight Function Combinations

I have introduced and proposed a number of weight functions, and calculated their constraint matrices. Thus far we have treated these weight functions individually; however, they can also be used in combination. For example, if we consider two constraints  $\mathbf{w}^T \mathbf{H}_1 \mathbf{w} \leq 1$  and  $\mathbf{w}^T \mathbf{H}_2 \mathbf{w} \leq 1$ , we can make a new constraint by combining  $\mathbf{H}_1$  and  $\mathbf{H}_2$ :

$$\mathbf{w}^T \mathbf{H} \mathbf{w} = p \mathbf{w}^T \mathbf{H}_1 \mathbf{w} + (1 - p) \mathbf{w}^T \mathbf{H}_2 \mathbf{w} \leq 1, \quad 0 \leq p \leq 1. \quad (6.115)$$

In general, when we consider  $L$  constraints:

$$\mathbf{w}^T \mathbf{H}_l \mathbf{w} \leq 1, \quad l = 1, \dots, L, \quad (6.116)$$

I unite the constraints as

$$\sum_{l=1}^L p_l \mathbf{w}^T \mathbf{H}_l \mathbf{w} \leq 1 \quad \text{and} \quad \sum_{l=1}^L p_l = 1. \quad (6.117)$$

I therefore propose a combination DOGQCM + KDEQCM with the strategies combined equally (i.e.,  $p = 0.5$ ).

## 6.2 Weighted regularization SVM (WRSVM)

In this section, I propose a novel regularization criterion for robust classifiers [110]. The criterion can produce many types of regularization terms by selecting an appropriate weighting function. L2 regularization terms, which are used for support vector machines (SVMs), can be realized by this criterion when the norm of patterns is normalized. In this regard, I propose two novel regularization terms based on the new criterion for a variety of applications. Furthermore, I propose new classifiers by applying these regularization terms to conventional SVMs.

Let a weighting function  $Q(\mathbf{x})$  satisfy

$$Q(\mathbf{z}) > 0 \quad (6.118)$$

for all  $\mathbf{z} \in \mathcal{D}$ , where  $\mathcal{D}$  is our data domain. The new regularization criterion is given by

$$R := \int_{\mathcal{D}} Q(\mathbf{x}) |\langle \mathbf{w}, \mathbf{x} \rangle|^2 d\mathbf{x}. \quad (6.119)$$

This regularization term can be rewritten as

$$\int_{\mathcal{D}} Q(\mathbf{x}) |\langle \mathbf{w}, \mathbf{x} \rangle|^2 d\mathbf{x} = \mathbf{w}^T H \mathbf{w}. \quad (6.120)$$

where  $H(i, j) := \int_{\mathcal{D}} Q(\mathbf{x}) x_i x_j d\mathbf{x}$  and  $H(i, j)$  denotes element  $(i, j)$  of the regularization matrix  $H$ . Note that  $H$  becomes a positive definite matrix from condition (6.118).

Combining our regularization approach with the hinge loss function, I propose a classification criterion whereby

$$\text{minimize } \mathbf{w}^T H \mathbf{w} + c \sum_{n=1}^N \xi_n, \quad (6.121)$$

$$\text{subject to } y_n (\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq 1 - \xi_n, \quad (6.122)$$

$$\xi_n \geq 0, \quad n = 1, \dots, N, \quad (6.123)$$

where  $\xi_n$  are slack variables. The proposed criterion can produce not only a basic SVM, but also various new classifiers. I demonstrate this relation in the following sections 6.2.1 and 6.2.2. In this regard, I refer to the proposed classifier as ‘‘Weighted Regularization SVM’’ (WR-SVM).

### 6.2.1 Basic Support Vector Machines

In this section, I demonstrate that our regularization criterion includes the basic regularization term  $\|\mathbf{w}\|^2$ . In other words, WR-SVM includes basic SVM.

Let us assume that  $\|\mathbf{x}\| = 1$ , and  $\{x_i, x_j\} (i \neq j)$  are orthogonal. The following assumption holds in the Gaussian kernel model:

$$\|\phi(\mathbf{x})\|^2 = k(\mathbf{x}, \mathbf{x}) = \exp(-\gamma\|\mathbf{x} - \mathbf{x}\|^2) = 1, \quad (6.124)$$

where  $k(\mathbf{x}, \mathbf{y}) = \exp(-\gamma\|\mathbf{x} - \mathbf{y}\|^2)$  is the Gaussian kernel function. We let  $Q(\mathbf{x})$  be a constant function:  $Q_1(\mathbf{x}) := \frac{1}{S}$ , where  $S$  is the volume of  $\mathcal{D}$ . Then, the constraint matrix is given by

$$H(i, j) = \frac{1}{S} \int_{\mathcal{D}} x_i x_j d\mathbf{x} \propto \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}. \quad (6.125)$$

We can see that the regularization matrix is defined as  $H_1 := I_M$ , as well as that it is equivalent to  $\|\mathbf{w}\|^2$ . Thus, we can regard our regularization method as an extension of the basic regularization term. Also, we can infer that the weighted regularization becomes basic if  $Q(\mathbf{x})$  is constant (i.e., no weight).

### 6.2.2 Novel Weighted Regularization

Next, I search for an appropriate weighting function. There are two approaches to this purpose. One is to make  $Q(\mathbf{x})$  large in a mixed area of categories. Therefore, I define  $Q_2(\mathbf{x})$  as a normal distribution:

$$Q_2(\mathbf{x}) := N(\mathbf{x}|\boldsymbol{\mu}, s\Sigma) = \frac{1}{(\sqrt{2\pi})^d \sqrt{|s\Sigma|}} \exp\left(-\frac{1}{2s}(\mathbf{x} - \boldsymbol{\mu})\Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right), \quad (6.126)$$

where

$$\boldsymbol{\mu} := \frac{\bar{\mathbf{x}}_{+1} + \bar{\mathbf{x}}_{-1}}{2}, \quad \Sigma(i, j) := \begin{cases} \frac{1}{N-1} \sum_{n=1}^N (\boldsymbol{\mu}(i) - \mathbf{x}_n(i))^2 & i = j \\ 0 & i \neq j \end{cases}, \quad (6.127)$$

and  $\bar{\mathbf{x}}_{+1}$  and  $\bar{\mathbf{x}}_{-1}$  denote the mean vectors of patterns labeled by  $+1$  and  $-1$ , respectively. The classifier becomes robust if patterns of different categories are mixed in the central area of the pattern distribution. Furthermore, if we let the parameter  $s$  is sufficiently large, this function is similar to a uniform function. Hence, its classifier becomes to the standard SVM.

The other approach is to make  $Q(\mathbf{x})$  small in dense areas and large in sparse areas. Then, this classifier becomes robust for outliers. Thus, I define a weighting function as the difference between two types of normal distributions:

$$Q_3(\mathbf{x}) := \left(1 + \frac{\kappa}{\nu^d - 1}\right) N(\mathbf{x}|\boldsymbol{\mu}, \nu^2\Sigma) - \frac{\kappa}{\nu^d - 1} N(\mathbf{x}|\boldsymbol{\mu}, \Sigma), \quad (6.128)$$

where  $0 < \kappa < 1$  and  $\nu > 1$ . If we assume that  $\Sigma$  is a diagonal matrix, then this weighting function always satisfies Eq. (6.118). If  $\nu$  increases, the weighting function becomes smoother and wider. Essentially,  $\kappa$  should be near 1 (e.g.,  $\kappa = 0.9$ ), and if  $\kappa$  is small,  $Q_3(\mathbf{x})$  is similar to  $Q_2(\mathbf{x})$ .

The calculation of these regularization matrices includes integration; however, if we use the Gaussian kernel as a basis function, then we can calculate  $H$  analytically since  $Q(\mathbf{x})$  consists of the Gaussian functions. I present the details of this approach in Section 6.2.3.

### 6.2.3 Analytical Calculation of Regularization Matrices

I define the regularization matrices  $H_2$  and  $H_3$  as

$$H_t(i, j) = \int_{\mathcal{D}} Q_t(\mathbf{x}) k(\mathbf{x}_i, \mathbf{x}) k(\mathbf{x}_j, \mathbf{x}) d\mathbf{x}, \quad t = 2, 3. \quad (6.129)$$

Note that it is only necessary to integrate the product of the normal distribution and two Gaussian kernel functions analytically. Then, I consider only the following integration:

$$U(i, j) = \int_{\mathcal{D}} N(\mathbf{x}|\boldsymbol{\mu}, \Sigma) k(\mathbf{x}_i, \mathbf{x}) k(\mathbf{x}_j, \mathbf{x}) d\mathbf{x}. \quad (6.130)$$

Using the general formula for a Gaussian integral;

$$\int e^{-\frac{1}{2}\mathbf{x}^T A \mathbf{x} + \mathbf{b}^T \mathbf{x}} d\mathbf{x} = \sqrt{\frac{(2\pi)^M}{|A|}} e^{\frac{1}{2}\mathbf{b}^T A^{-1} \mathbf{b}}, \quad (6.131)$$

we can calculate Eqs. (6.130) analytically as follows.

$$U(i, j) = \frac{1}{\sqrt{|4\gamma\Sigma + I_N|}} \exp\left(\frac{1}{2}\mathbf{b}_{ij}^T A^{-1} \mathbf{b}_{ij} + C_{ij}\right), \quad (6.132)$$

$$A = 4\gamma I_N + \Sigma^{-1}, \quad (6.133)$$

$$\mathbf{b}_{ij} = 2\gamma(\mathbf{x}_i + \mathbf{x}_j) + \Sigma^{-1}\boldsymbol{\mu}, \quad (6.134)$$

$$C_{ij} = -\gamma(\|\mathbf{x}_i\|^2 + \|\mathbf{x}_j\|^2) - \frac{1}{2}\boldsymbol{\mu}^T \Sigma^{-1} \boldsymbol{\mu}. \quad (6.135)$$

$H_2$  and  $H_3$  can also be calculated in a similar manner.

In practice, the regularization matrix  $H_t$  is normalized by  $(NH_t)/\text{tr}(H_t)$  so that the constant factor of  $U(i, j)$  is independent of multiplication factor.

### 6.2.4 Novel Classifiers

In this section, I propose novel classifiers by using weighted regularization terms.

I assume that the discriminant function is given by

$$D(\mathbf{x}|\boldsymbol{\alpha}, b) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x}) + b. \quad (6.136)$$

Then, the training problem is given by

$$\text{minimize } \frac{1}{2}\boldsymbol{\alpha}^T H_t \boldsymbol{\alpha} + c \sum_{n=1}^N \xi_n, \quad (6.137)$$

$$\text{subject to } y_n \left( \sum_{i=1}^N \alpha_i k(\mathbf{x}_i, \mathbf{x}_n) + b \right) \geq 1 - \xi_n, \quad (6.138)$$

$$\xi_n \geq 0, \quad n = 1, \dots, N. \quad (6.139)$$

I solve this problem by two steps. First, I solve its dual problem:

$$\text{maximize } -\frac{1}{2}\boldsymbol{\beta}^T Y K H_t^{-1} K Y \boldsymbol{\beta} + \sum_{n=1}^N \beta_n, \quad (6.140)$$

$$\text{subject to } 0 \leq \beta_n \leq c, \quad \sum_{n=1}^N \beta_n y_n = 0, \quad n = 1, \dots, N, \quad (6.141)$$

where  $Y := \text{diag}(\mathbf{y})$ ,  $\boldsymbol{\beta}$  is a dual parameter vector, and its solution  $\hat{\boldsymbol{\beta}}$  can be obtained by quadratic programming [73]. For the problem, a number of quadratic programming solvers have been developed thus far, such as LOQO [5]. Second, the estimated parameters  $\hat{\boldsymbol{\alpha}}$  and  $\hat{b}$  are given by

$$\hat{\boldsymbol{\alpha}} = H_t^{-1} K Y \hat{\boldsymbol{\beta}}, \quad \hat{b} = \frac{\mathbf{1}^T \mathbf{y} - \mathbf{k}^T \hat{\boldsymbol{\alpha}}}{N}. \quad (6.142)$$

Substituting  $H_2$  or  $H_3$  into Eq.(6.140), we can construct two novel classifiers. I denote these classifiers as “SVM<sub>G</sub>” and “SVM<sub>D</sub>”, respectively (based on the initials of Gaussian and Difference).

### 6.3 Chernoff FDA

Fisher discriminant analysis (FDA) is a popular supervised feature extractor. FDA has many applications in machine learning and pattern recognition [7, 37]. However, FDA does not give an optimal projection when two covariance matrices do not have the same shape even if we assume each category follows a Gaussian distribution (i.e., heteroscedastic case).

In order to handle a heteroscedastic case, it is a key-point how to evaluate a distance between two heteroscedastic Gaussian distributions, and various methods have been studied [30, 93, 118, 67, 89, 70, 69]. Decell et al. [30] proposed to evaluate the distance of distributions by the Kullback-Leibler (KL) divergence, and Tao et al. [93] generalized the KL divergence based FDA by using Bregman divergence and generalized average. Bhattacharyya and Chernoff distances have been also used to evaluate the distance of distributions for a heteroscedastic case in many studies [118, 67, 89, 70]. Zhang et al. [118] proposed to maximize Fisher and Bhattacharyya criteria for binary- and multi-class problems via a Fukunaga-Koontz transform. Loog et al. [67] proposed to use an approximated criterion for the Chernoff distance and Rueda et al. [89] developed a gradient-based algorithm to maximize the Chernoff distance. The Matusita separability measure [70] and a sufficient statistic [69] also have been used.

On the other hand, there are several approaches to extend the objects of the FDA to non-Gaussian distributions. These approaches include non-linear methods using a quadratic polynomial model [91] and a kernel trick [3, 72, 12], Gaussian mixture model based methods [44, 119], and an extended FDA based on auto-correlation matrices [90].

In this section, I propose a new algorithm to maximize the Bhattacharyya distance as an extension of FDA [112, 108]. The Bhattacharyya distance provides a new criterion which consists of the conventional FDA criterion and an additional term. This additional term can be regarded as a correction term for the heteroscedastic Gaussian case since the Bhattacharyya distance is closely related

to misclassification rate via the Hellinger distance. Moreover, the Bhattacharyya distance is a special case of the Chernoff distance. Hence, I propose to use the Chernoff distance as an extension of the Bhattacharyya distance. I call these methods by ‘‘Bhattacharyya FDA’’ (BFDA) and ‘‘Chernoff FDA’’ (CFDA). BFDA is applicable to heteroscedastic Gauss distributions, and CFDA is robust against unbalanced training sample size. Furthermore, I propose their kernelized versions and conduct experiments to evaluate their abilities as binary classifiers. These kernelized classifiers work well even when given samples follow non-Gaussian distributions.

Our approach is related to the prior works [67, 89] because they also tried to maximize the Chernoff distance. However, the former did not provide the strict solution for maximizing the Chernoff distance. Although the latter proposed a strict solution algorithm for maximizing the Chernoff distance, the gradient-based algorithm was not stable. Furthermore, both works did not discuss a hyper-parameter selection for the Chernoff distance carefully just by substituting a prior probability for the hyper-parameter. The contributions and advantages of our research are to propose a new stable algorithm for maximizing the Chernoff distance, to propose several methods of hyper-parameter estimation, and to provide their kernelized versions as applied to the non-Gaussian case.

### 6.3.1 Bhattacharyya and Chernoff Distances

In this section, I show the theoretical relationships between misclassification rate and the Bhattacharyya distance, when the class distributions are heteroscedastic Gaussian.

First, I describe the definition of misclassification rate. Suppose we classify  $\mathbf{x}$  according to the rule of (4.86) and the prior probabilities of class 1 and 2 are the same (i.e.,  $p_1 = p_2 = 0.5$ ), the misclassification rate is defined by

$$\frac{1}{2} \int_{y|f_1(y) < f_2(y)} f_1(y) dy + \frac{1}{2} \int_{y|f_1(y) > f_2(y)} f_2(y) dy \quad (6.143)$$

The minimization of (6.143) is closely related to the maximization of the Hellinger distance [65] that is a measure of distance between two probability density functions:

$$\begin{aligned} h(f_1, f_2) &:= \frac{1}{2} \int_{-\infty}^{\infty} \left( \sqrt{f_1(y)} - \sqrt{f_2(y)} \right)^2 dy \\ &= 1 - \int_{-\infty}^{\infty} \sqrt{f_1(y)f_2(y)} dy. \end{aligned} \quad (6.144)$$

Figure 6.5 illustrates the conceptual sketch of relationship between misclassification rate and the Hellinger distance. We have

$$\left( \sqrt{f_1(y)} - \sqrt{f_2(y)} \right)^2 \simeq \begin{cases} f_1(y) & f_1(y) \gg f_2(y) \\ f_2(y) & f_2(y) \ll f_1(y) \end{cases}. \quad (6.145)$$

Although the maximization of the Hellinger distance and the minimization of misclassification rate are not equivalent, we can know from (6.145) that the former can be used as an approximation of the latter.

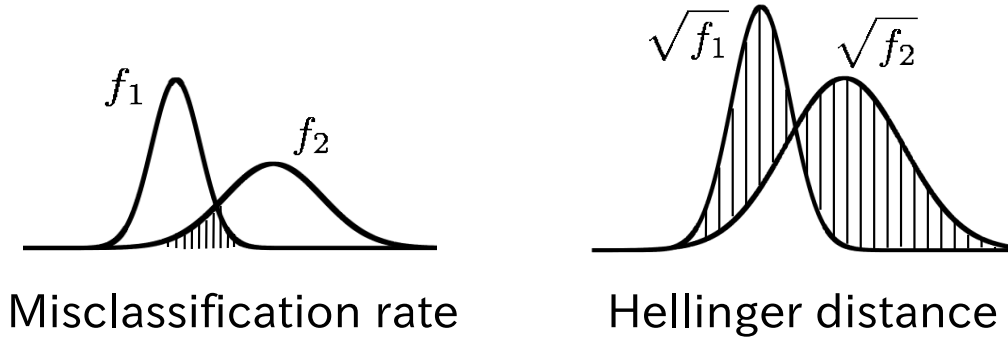


Figure 6.5: Misclassification rate and Hellinger distance

The second term of (6.144) is called the Bhattacharyya coefficient. The Bhattacharyya distance is also defined by

$$B_h(f_1, f_2) := -\log \int_{-\infty}^{\infty} \sqrt{f_1(y)f_2(y)} dy. \quad (6.146)$$

The maximization of the Bhattacharyya distance and that of the Hellinger distance are clearly equivalent. Therefore, it is effective to use the Bhattacharyya distance as a criterion of supervised feature extraction for classification.

In addition, I introduce the Chernoff distance by

$$C_h(f_1, f_2) := -\log \int_{-\infty}^{\infty} f_1(y)^{\rho_1} f_2(y)^{\rho_2} dy, \quad (6.147)$$

where  $\rho_1$  and  $\rho_2$  are positive weighting parameters, and  $\rho_1 + \rho_2 = 1$ . The Chernoff distance is an extension of the Bhattacharyya distance since the Chernoff distance with  $\rho_1 = \rho_2 = 0.5$  equals to the Bhattacharyya distance. It should be noted that when  $\rho_1 < \rho_2$  the weight of  $f_1(y)$  is larger than that of  $f_2(y)$ ; the smaller value of  $\rho$  yields the higher weight. This fact can be reconfirmed by experiments. Loog et al. suggested that  $\rho_1$  and  $\rho_2$  are estimated by prior probabilities of class 2 and class 1, respectively (i.e.  $\rho_1 \leftarrow p_2$  and  $\rho_2 \leftarrow p_1$ ) [67].

Next, I define the probability density function  $f_c(y)$  for category  $c$  by

$$f_c(y) = \frac{1}{\sqrt{2\pi\sigma_c^2}} \exp\left[-\frac{(y - m_c)^2}{2\sigma_c^2}\right] \quad (c = 1, 2). \quad (6.148)$$

Substituting (6.148) into (6.146), I can obtain the following explicit form:

$$\begin{aligned} B_h(f_1, f_2) &= \frac{(m_1 - m_2)^2}{4(\sigma_1^2 + \sigma_2^2)} - \frac{1}{2} \log \frac{2\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2} \\ &= \frac{\langle \mathbf{w}, \mathbf{S}_B \mathbf{w} \rangle}{4 \langle \mathbf{w}, \mathbf{S}_W \mathbf{w} \rangle} - \frac{1}{2} \log \frac{2\sqrt{\langle \mathbf{w}, \mathbf{\Sigma}_1 \mathbf{w} \rangle} \sqrt{\langle \mathbf{w}, \mathbf{\Sigma}_2 \mathbf{w} \rangle}}{\langle \mathbf{w}, \mathbf{S}_W \mathbf{w} \rangle}. \end{aligned} \quad (6.149)$$

The first term of (6.149) is equivalent to the FDA criterion. If  $\mathbf{\Sigma}_1 = \beta \mathbf{\Sigma}_2$  for any positive  $\beta$ , the second

term of (6.149) is given by

$$\begin{aligned} & \log \frac{2\sqrt{\langle \mathbf{w}, \beta \boldsymbol{\Sigma}_2 \mathbf{w} \rangle} \sqrt{\langle \mathbf{w}, \boldsymbol{\Sigma}_2 \mathbf{w} \rangle}}{\langle \mathbf{w}, (\beta \boldsymbol{\Sigma}_2 + \boldsymbol{\Sigma}_2) \mathbf{w} \rangle} \\ &= \log \frac{2\sqrt{\beta} \langle \mathbf{w}, \boldsymbol{\Sigma}_2 \mathbf{w} \rangle}{(1 + \beta) \langle \mathbf{w}, \boldsymbol{\Sigma}_2 \mathbf{w} \rangle} = \log \frac{2\sqrt{\beta}}{(1 + \beta)}. \end{aligned} \quad (6.150)$$

Because this is constant, the criterion does not depend on the second term in such a case. Since most of classification problems will violate the assumption, I must say that the FDA criterion is not enough for general-purpose classification.

Substituting (6.148) into (6.147), the Chernoff version is obtained as

$$\begin{aligned} C_h(f_1, f_2) &= \frac{\rho_1 \rho_2 (m_1 - m_2)^2}{2(\rho_2 \sigma_1^2 + \rho_1 \sigma_2^2)} - \frac{1}{2} \log \frac{\sigma_1^{2\rho_2} \sigma_2^{2\rho_1}}{\rho_2 \sigma_1^2 + \rho_1 \sigma_2^2} \\ &= \frac{\rho_1 \rho_2 \langle \mathbf{w}, \mathbf{S}_B \mathbf{w} \rangle}{2 \langle \mathbf{w}, \mathbf{S}_C \mathbf{w} \rangle} - \frac{1}{2} \log \frac{\langle \mathbf{w}, \boldsymbol{\Sigma}_1 \mathbf{w} \rangle^{\rho_2} \langle \mathbf{w}, \boldsymbol{\Sigma}_2 \mathbf{w} \rangle^{\rho_1}}{\langle \mathbf{w}, \mathbf{S}_C \mathbf{w} \rangle}, \end{aligned} \quad (6.151)$$

where  $\mathbf{S}_C := \rho_2 \boldsymbol{\Sigma}_1 + \rho_1 \boldsymbol{\Sigma}_2$ . Note that (6.151) is equal to (6.149) at  $\rho_1 = \rho_2 = 0.5$ .

### 6.3.2 Chernoff FDA

According to the discussion in Section 6.3.1, I propose to maximize (6.149) and (6.151) as new criteria. The second term will play an important role in many actual classification problems as a correction term of FDA. Since the Chernoff distance includes the Bhattacharyya distance, I discuss only CFDA. As an important property, (6.151) is invariant with respect to the scale of  $\mathbf{w}$ . Then, I can introduce a constraint that  $\langle \mathbf{w}, \mathbf{S}_C \mathbf{w} \rangle = 1$ , and a new criterion for CFDA is given by

$$\begin{aligned} \max_{\mathbf{w}} C_h(f_1, f_2) &= \rho_1 \rho_2 \langle \mathbf{w}, \mathbf{S}_B \mathbf{w} \rangle - \log \langle \mathbf{w}, \boldsymbol{\Sigma}_1 \mathbf{w} \rangle^{\rho_2} \\ &\quad - \log \langle \mathbf{w}, \boldsymbol{\Sigma}_2 \mathbf{w} \rangle^{\rho_1}, \\ \text{s.t. } &\langle \mathbf{w}, \mathbf{S}_C \mathbf{w} \rangle = 1. \end{aligned} \quad (6.152)$$

When  $\rho_1 = \rho_2 = 0.5$ , (6.152) is the criterion of BFDA.

### Optimization Algorithm

In this section, I explain an optimization algorithm for the CFDA criterion. From the CFDA criterion (6.152), its Lagrangian is given by

$$\begin{aligned} L(\mathbf{w}) &= \rho_1 \rho_2 \langle \mathbf{w}, \mathbf{S}_B \mathbf{w} \rangle - \rho_2 \log \langle \mathbf{w}, \boldsymbol{\Sigma}_1 \mathbf{w} \rangle \\ &\quad - \rho_1 \log \langle \mathbf{w}, \boldsymbol{\Sigma}_2 \mathbf{w} \rangle - \lambda (\langle \mathbf{w}, \mathbf{S}_C \mathbf{w} \rangle - 1), \end{aligned} \quad (6.153)$$

where  $\lambda$  is a Lagrangian multiplier. Then, Lagrange's conditions are given by

$$\left[ \mathbf{S}_B - \frac{\boldsymbol{\Sigma}_1}{\rho_1 \langle \mathbf{w}, \boldsymbol{\Sigma}_1 \mathbf{w} \rangle} - \frac{\boldsymbol{\Sigma}_2}{\rho_2 \langle \mathbf{w}, \boldsymbol{\Sigma}_2 \mathbf{w} \rangle} \right] \mathbf{w} = \lambda \mathbf{S}_C \mathbf{w}, \quad (6.154)$$

$$\langle \mathbf{w}, \mathbf{S}_C \mathbf{w} \rangle = 1. \quad (6.155)$$

Since it is difficult to solve (6.154) directly, I propose an updating and iterating approach. For the  $k$ -th iteration, I solve

$$\left[ \mathbf{S}_B - \frac{\mathbf{\Sigma}_1}{\rho_1 \langle \mathbf{w}_k, \mathbf{\Sigma}_1 \mathbf{w}_k \rangle} - \frac{\mathbf{\Sigma}_2}{\rho_2 \langle \mathbf{w}_k, \mathbf{\Sigma}_2 \mathbf{w}_k \rangle} \right] \mathbf{v}_i = \lambda_i \mathbf{S}_C \mathbf{v}_i \quad (6.156)$$

as a generalized eigenvalue problem. Then, I obtain a set of eigenvectors  $\{\mathbf{v}_i\}_{i=1}^d$ . Next I normalize individual eigenvectors by

$$\mathbf{v}_i \leftarrow \frac{\mathbf{v}_i}{\sqrt{\langle \mathbf{v}_i, \mathbf{S}_C \mathbf{v}_i \rangle}}, \text{ for } i = 1, \dots, d. \quad (6.157)$$

Finally, I update

$$\mathbf{w}_{k+1} \leftarrow \operatorname{argmax}_{\mathbf{v}_i} C_h(f_1, f_2) \text{ with } \mathbf{v}_i, \quad (6.158)$$

$$k \leftarrow k + 1. \quad (6.159)$$

In this way,  $\mathbf{w}_{k+1}$  always satisfies equations (6.156) and (6.155), and maximizes the Chernoff distance. When  $\|\mathbf{w}_{k+1} - \mathbf{w}_k\| < \varepsilon$ ,  $\mathbf{w}_{k+1}$  satisfies both Lagrange's conditions (6.154) and (6.155), the algorithm is finished and output  $\mathbf{w}_{k+1}$  as a solution of CFDA criterion, where  $\varepsilon$  is a very small positive tolerance parameter (e.g.,  $\varepsilon = 10^{-8}$ ). I summarize this algorithm in Algorithm 29.

Sometimes, numerical calculation of the generalized eigenvalue problem provides a strange solution so that the solution includes a complex number. In this case, it is effective to add some regularization term. For example,  $\mathbf{\Sigma}_c \leftarrow \mathbf{\Sigma}_c + \eta \mathbf{I}$  where  $\mathbf{I}$  is the identity matrix, and  $\eta$  is a regularization parameter.

---

#### Algorithm 29 CFDA-algorithm

---

**Input:**  $\mathbf{\Sigma}_1, \mathbf{\Sigma}_2, \mathbf{S}_B, \mathbf{S}_C, \rho_1, \rho_2$  and tolerance  $\varepsilon > 0$ .

**Initialize:**  $k = 0$ , and  $\mathbf{w}_1$  so that it satisfies (6.155).

**repeat**

$k \leftarrow k + 1$ ,

Obtain the set of eigenvectors  $\{\mathbf{v}_i\}_{i=1}^d$  by (6.156),

Normalize  $\{\mathbf{v}_i\}_{i=1}^d$  by (6.157),

Update  $\mathbf{w}_{k+1} \leftarrow \operatorname{argmax}_{\mathbf{v}_i} C_h(f_1, f_2)$  with  $\mathbf{v}_i$

**until**  $\|\mathbf{w}_k - \mathbf{w}_{k+1}\|^2 < \varepsilon$ .

**Output:**  $\mathbf{w}_k$ .

---

### Hyper-parameter Estimation

In this section, I discuss how to estimate  $\rho_1$  and  $\rho_2$  of CFDA. There are two parameters, although I need to estimate only one parameter  $\rho_1$  since  $\rho_2$  is decided as  $1 - \rho_1$ . I propose two methods to estimate these hyper-parameters.

In the first method, I estimate the parameters based on prior probabilities. Thus,  $\rho_1 \leftarrow p_2$  and  $\rho_2 \leftarrow p_1$ . I denote this method by ‘CFDA- $p$ ’. From the view point of weighting parameter, it is expected as good estimators.

In the second method, I estimate the parameters based on the Chernoff distance by

$$(\rho_1^*, \rho_2^*) = \underset{\rho_1, \rho_2}{\operatorname{argmax}} C_h(f_1, f_2), \text{ s.t. } \rho_1 + \rho_2 = 1. \quad (6.160)$$

However, this optimization problem is not convex. Therefore, I use a simple one dimensional search. First I consider a set of candidates of parameters. The number of candidates can be selected to few or much. Note that too few or much number of candidates are not so effective in this method. I propose to use  $P = \{0.1, 0.2, \dots, 0.8, 0.9\}$  as the set of candidates. Thus, I estimate the weighting parameters from these candidates by

$$(\rho_1^*, \rho_2^*) = \underset{\rho_1, \rho_2 \in P}{\operatorname{argmax}} C_h(f_1, f_2), \text{ s.t. } \rho_1 + \rho_2 = 1. \quad (6.161)$$

This solution can be obtained by nine iterations of CFDA-algorithm. I denote this method by ‘CFDA- $m$ ’. I summarize the CFDA- $m$ -algorithm in Algorithm 30.

---

**Algorithm 30** CFDA- $m$ -algorithm
 

---

**Input:**  $\Sigma_1, \Sigma_2, \mathbf{S}_B, \mathbf{S}_C$ , and tolerance  $\varepsilon > 0$ .

**for**  $k = 1, 2, \dots, 9$  **do**

$$\rho_1 \leftarrow \frac{k}{10},$$

$$\rho_2 \leftarrow 1 - \rho_1,$$

$\mathbf{w}_k$  is obtained by CFDA-algorithm with  $(\rho_1, \rho_2)$ ,

$$J_k \leftarrow C_h(f_1, f_2) \text{ with } (\mathbf{w}_k, \rho_1, \rho_2).$$

**end for**

$$k^* \leftarrow \operatorname{argmax}_k J_k.$$

**Output:**  $\mathbf{w}_{k^*}$ .

---

### 6.3.3 Kernel Extension of CFDA

In this section, I kernelize the CFDA criterion. Kernel method is to map from an original space to a high dimensional feature space by  $\phi(\mathbf{x})$ , and perform some linear analysis method in the high dimensional feature space. In this regard, I do not calculate  $\phi(\mathbf{x})$  directly, but I need to calculate only an inner product  $k(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$ . In this paper, I use the Gaussian RBF kernel function as  $k(\mathbf{x}, \mathbf{y}) := \exp(-\gamma \|\mathbf{x} - \mathbf{y}\|^2)$ . In a similar way of the kernel FDA [72, 3], CFDA can be also extended to a kernel classifier.

In this method,  $\mathbf{w}$  is given by

$$\mathbf{w} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n) = \mathbf{\Phi} \boldsymbol{\alpha}, \quad (6.162)$$

where  $\mathbf{\Phi} := [\phi(\mathbf{x}_1) \cdots \phi(\mathbf{x}_N)]$ ,  $\boldsymbol{\alpha} := [\alpha_1, \dots, \alpha_N]^T$ , and  $N$  is the number of samples. Then, I define

$$\mathbf{k}(\mathbf{x}) := \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}) \\ \vdots \\ k(\mathbf{x}_N, \mathbf{x}) \end{pmatrix} = \mathbf{\Phi}^T \phi(\mathbf{x}), \quad (6.163)$$

$$\mathbf{k}_c := \frac{1}{N_c} \sum_{\mathbf{x}_n \in \Omega_c} \mathbf{k}(\mathbf{x}_n), \quad (6.164)$$

$$\mathbf{R}_c := \frac{1}{N_c} \sum_{\mathbf{x} \in \Omega_c} (\mathbf{k}(\mathbf{x}_n) - \mathbf{k}_c)(\mathbf{k}(\mathbf{x}_n) - \mathbf{k}_c)^T, \quad (6.165)$$

where  $N_c$  is the number of samples belonging to class  $c$ . Projected one-dimensional vector is given by

$$y = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x}) = \boldsymbol{\alpha}^T \mathbf{k}(\mathbf{x}). \quad (6.166)$$

Average and variance in the projected space are given by

$$m_c = \boldsymbol{\alpha}^T \mathbf{k}_c, \quad (6.167)$$

$$\sigma_c^2 = \boldsymbol{\alpha}^T \mathbf{R}_c \boldsymbol{\alpha}. \quad (6.168)$$

Finally, the criterion of the ‘kernel CFDA’ (KCFDA) is given by

$$\begin{aligned} \max \quad & \rho_1 \rho_2 \langle \boldsymbol{\alpha}, \mathbf{K}_B \boldsymbol{\alpha} \rangle - \rho_2 \log \langle \boldsymbol{\alpha}, \mathbf{R}_1 \boldsymbol{\alpha} \rangle - \rho_1 \log \langle \boldsymbol{\alpha}, \mathbf{R}_2 \boldsymbol{\alpha} \rangle, \\ \text{s.t.} \quad & \langle \mathbf{w}, \mathbf{K}_C \mathbf{w} \rangle = 1, \end{aligned} \quad (6.169)$$

where

$$\mathbf{K}_B = (\mathbf{k}_1 - \mathbf{k}_2)(\mathbf{k}_1 - \mathbf{k}_2)^T, \quad (6.170)$$

$$\mathbf{K}_C = \rho_2 \mathbf{R}_1 + \rho_1 \mathbf{R}_2. \quad (6.171)$$

When  $\rho_1 = \rho_2 = 0.5$ , I call it the ‘kernel BFDA’ (KBFDA). In this regard, I do not calculate the feature mapping  $\phi(\mathbf{x})$  directly, but I need to calculate only the kernel function.

The numerical calculation of KCFDA is given by just replacing  $\mathbf{x}$  by  $\mathbf{k}(\mathbf{x})$ , and applying the CFDA algorithm to obtain  $\boldsymbol{\alpha}$ . I also propose KCFDA- $p$  and KCFDA-m by applying the hyper-parameter estimation methods.

# Chapter 7

## Experiments

In this Chapter, I show the experimental results of new feature extraction and classification methods. To evaluate proposed feature extraction methods, I conduct experiments for the low-rank approximation, the face reconstruction, the blind source separation (BSS), the denoising, the parts-based representation, and the compression in Sections 7.1, 7.2, and 7.3. To evaluate proposed classification methods, I conduct experiments of classification with UCI datasets and BCI dataset in Section 7.4.

### 7.1 fastGRBF based methods

#### 7.1.1 Low-rank approximation

##### The Yale Face Database

In this experiment, I examine the robustness for the salt-and-pepper noise and computational times of our algorithms. First, I use 10 face images of a subject in the Yale Face Database [4] for denoising experiment. Its images are down-sampled and unfolded to 99 dimensional vectors. Then, an image is converted to a  $(99 \times 10)$ -matrix. The reason why one image is so small is that the original GRBF algorithm is slow even if I try such a small matrix factorization, and I conduct this experiment to compare the proposed fastGRBF-NMF with the original GRBF. Salt-and-pepper noise was added to the dataset so that PSNRs of datasets are 10, 15, and 20 dB, where PSNR (Peak Signal to Noise Ratio) of  $\mathbf{Y}, \mathbf{Y}' \in \mathbb{R}^{I \times J}$  is defined by

$$\text{PSNR}(\mathbf{Y}, \mathbf{Y}') := 10 \log_{10} \left[ \frac{(255^2)IJ}{\|\mathbf{Y}' - \mathbf{Y}\|_F^2} \right]. \quad (7.1)$$

For tensors  $\underline{\mathbf{Y}}, \underline{\mathbf{Y}}' \in \mathbb{R}^{I \times J \times K}$ , it is also defined by

$$\text{PSNR}(\underline{\mathbf{Y}}, \underline{\mathbf{Y}}') := 10 \log_{10} \left[ \frac{(255^2)IJK}{\|\underline{\mathbf{Y}}' - \underline{\mathbf{Y}}\|_F^2} \right]. \quad (7.2)$$

I tried to run the HALS-NMF [24], the multiplicative NMF [62], the GRBF-NMF, fastGRBF-NMF, the fastGRBF-NMF-2Dbasis for all PSNRs. Table 7.1 shows the parameter settings for the GRBF

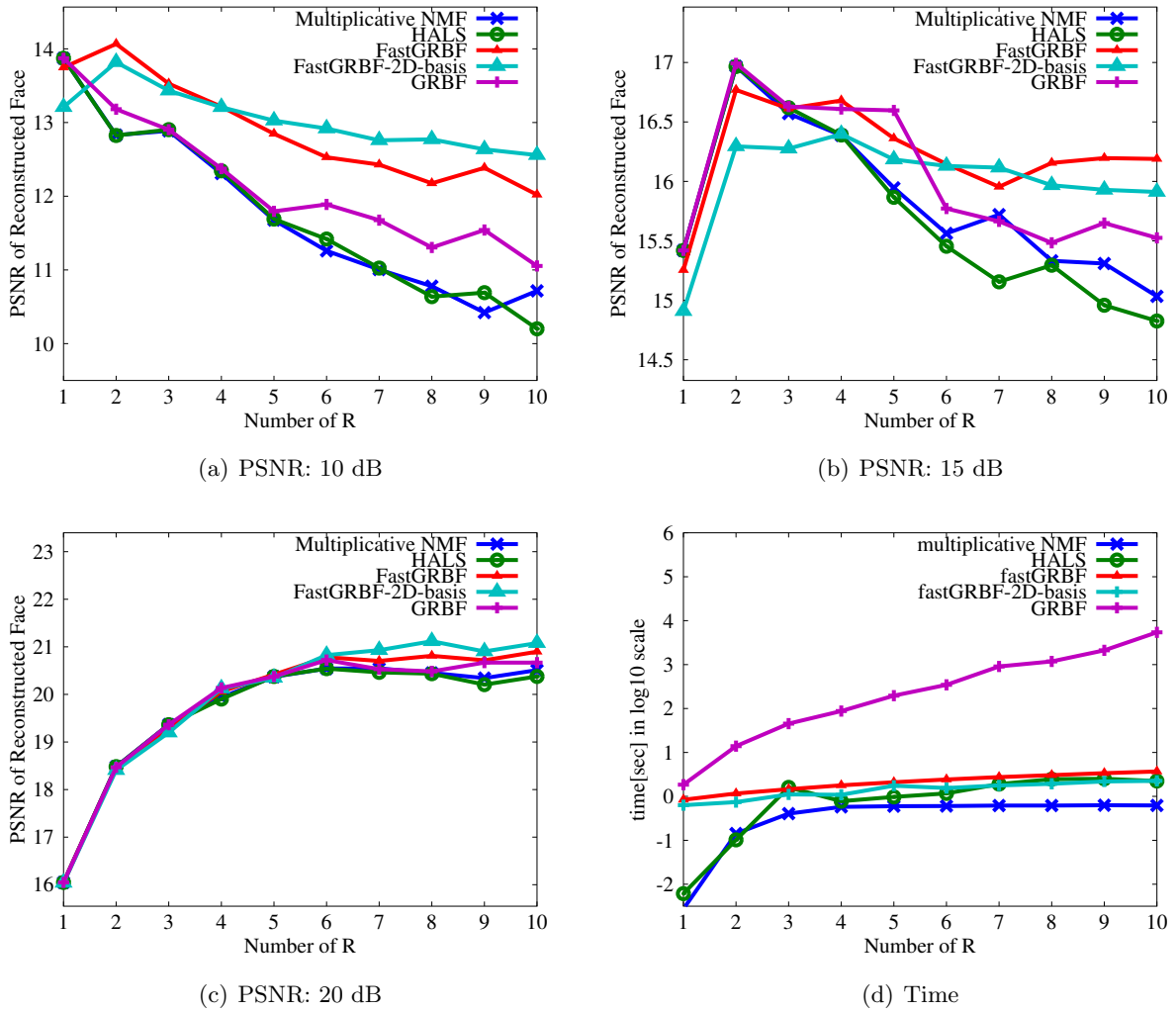


Figure 7.1: PSNR of NMF results on 10 faces for various ranks

and the fastGRBF methods. The degree of freedom of the function approximation in the GRBF is higher than in the fastGRBF since the constraints  $\Phi \mathbf{W} \geq 0$  for the GRBF is weaker than  $\mathbf{W} \geq 0$  for the fastGRBF. Therefore, I set the standard deviation parameters of the GRBF to a little larger value than those of the fastGRBF.

Figure 7.1 shows PSNRs and log 10-scale computational times for various rank decompositions. The HALS and the multiplicative algorithms decreased the PSNRs for large  $R$ ; however, the GRBF and the fastGRBF algorithms kept PSNRs for large  $R$ . The fastGRBF performed almost always the best PSNRs in case of 10 dB. The constraint of  $\mathbf{W} \geq 0$  would work well for noisy data used in the fastGRBF based methods. In Figure 7.1(d), the computational time of the GRBF exponentially increased for large  $R$ . On the other hand, the fastGRBF spent a much less to obtain its solution. Table 7.2 describes the average computational times of one decomposition by four algorithms. The fastGRBF is about 300 times faster than the GRBF.

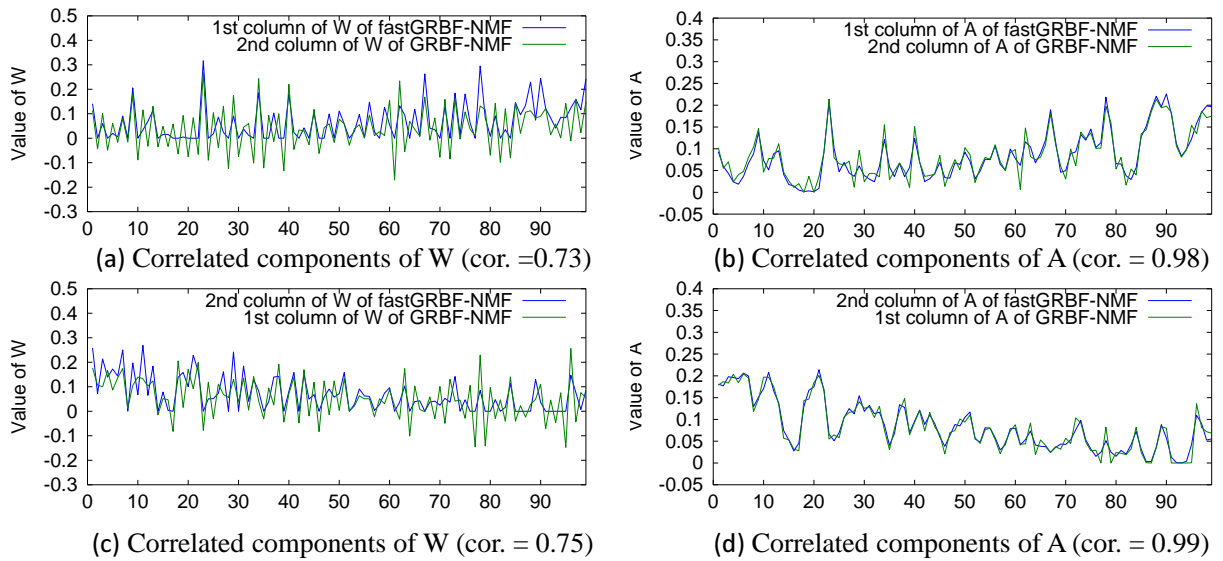
Next, I try to compare the values of  $\mathbf{W}$  and  $\mathbf{A}$  between by the proposed fastGRBF-NMF and by

Table 7.1: Parameter setting

PSNR	GRBF	FastGRBF
10	$\Delta t = 1, \sigma = 1.2$	$U = 4, \delta t = 1, \sigma = 1.0$
15	$\Delta t = 1, \sigma = 0.8$	$U = 4, \delta t = 1, \sigma = 0.7$
20	$\Delta t = 1, \sigma = 0.5$	$U = 4, \delta t = 1, \sigma = 0.5$

Table 7.2: Computational times for various algorithms

Algorithms	Multiplicative	HALS	GRBF	FastGRBF
Ave. of times [s]	0.499	0.759	720.0	2.23

Figure 7.2: Values of  $\mathbf{W}$  and  $\mathbf{A}$  obtained by the fastGRBF-NMF and the GRBF-NMF

the original GRBF-NMF. I set the parameters as follows:  $U = 1, \delta t = 1, \sigma = 0.8$ , and  $R = 2$  for the fastGRBF-NMF,  $\Delta t = 1, \sigma = 0.8$ , and  $R = 2$  for the GRBF-NMF. Fig. 7.2 shows the values of  $\mathbf{W}$  and  $\mathbf{A}$  obtained by the fastGRBF-NMF and the GRBF-NMF. In spite of the difference between the constraints  $\mathbf{W} \geq 0$  and  $\Phi \mathbf{W} \geq 0$ , the results of  $\mathbf{A}$  are highly correlated. Furthermore, the fastGRBF-NMF provides smoother components than the original GRBF-NMF. Many spikes in the result by the GRBF-NMF are smoothed by the fastGRBF-NMF method. Thus, it can be considered that the modification from  $\Phi \mathbf{W} \geq 0$  to  $\mathbf{W} \geq 0$  gives smoother features.

Next, I used also the Yale Faces Database including 165 face images of 15 subjects for experiments on NMF and NTF. One face image can be expressed as an 858 dimensional vector. The data matrix is an  $(858 \times 165)$ -matrix  $\mathbf{Y}$ . The data tensor is an  $(858 \times 11 \times 15)$ -tensor  $\underline{\mathbf{Y}}$ . The salt-and-pepper noise was added to the dataset so that PSNRs of datasets are 10, 15, and 20 dB. I run NMF and NTF with the HALS, the Multiplicative, the fastGRBF, the fastGRBF-2Dbasis, the ALS-NTD [58], the HALS-NCPD [24], the fastGRBF-NTD, the fastGRBF-NCPD, the fastGRBF-NTD-2Dbasis, and the

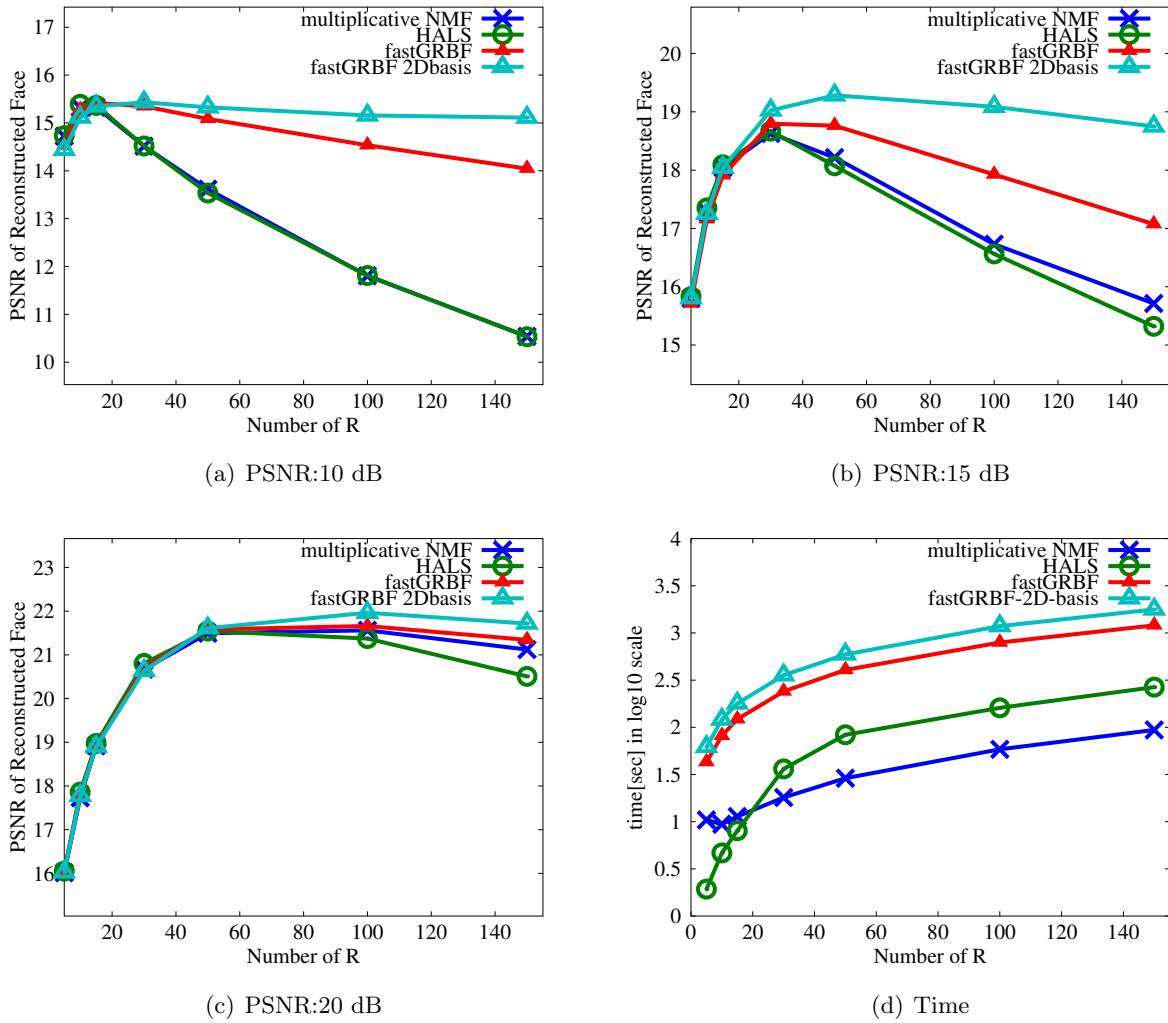


Figure 7.3: PSNR of NMF results on 165 faces for various ranks

fastGRBF-NCPD-2Dbasis algorithms for all PSNRs. The rank of NMF and CP models was changed from 5 to 150.  $R_1$  was changed from 5 to 150, and  $R_2 = R_3 = 10$  were fixed in Tucker model. I set  $U = 4$ ,  $\delta t = 1$ ,  $\sigma = 2.0$  for PSNR=10 dB,  $\sigma = 1.0$  for PSNR=15 dB, and  $\sigma = 0.5$  for PSNR=20 dB in all fastGRBF based methods.

Figures 7.3 and 7.4 show the PSNRs and the log10-scale computational times for various rank decompositions. The fastGRBF-NMF, the fastGRBF-NTD, and the fastGRBF-NCPD were almost better than traditional algorithms. Furthermore, the fastGRBF-NMF-2Dbasis, the fastGRBF-NTD-2Dbasis, and the fastGRBF-NCPD-2Dbasis were almost better than the fastGRBF-NMF, the fastGRBF-NTD, and the fastGRBF-NCPD algorithms. The computational times of the proposed methods were not shorter than those of the traditional algorithms since the sizes of parameter spaces for the proposed methods are larger. However, they are promising since we can obtain better performances.

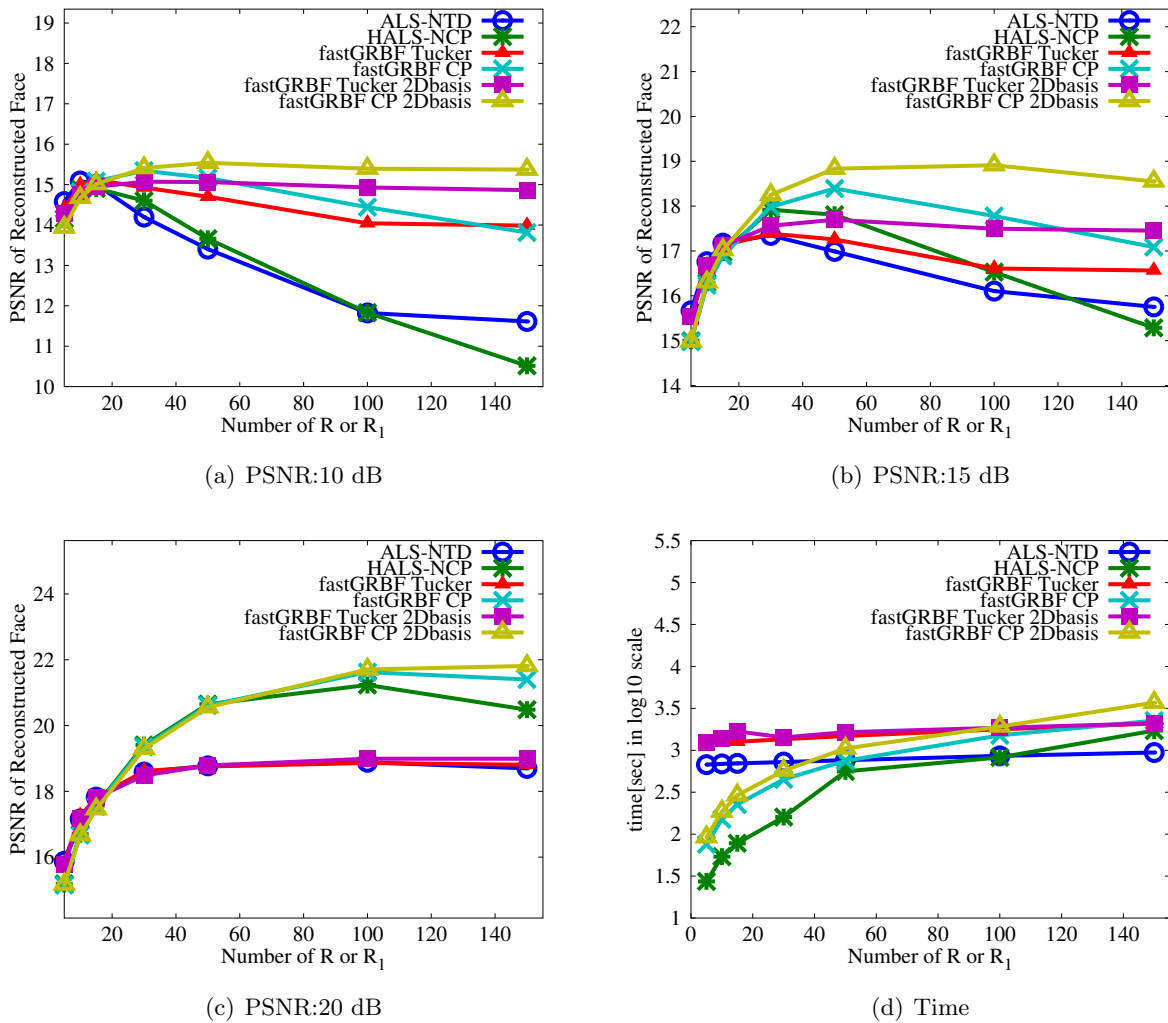
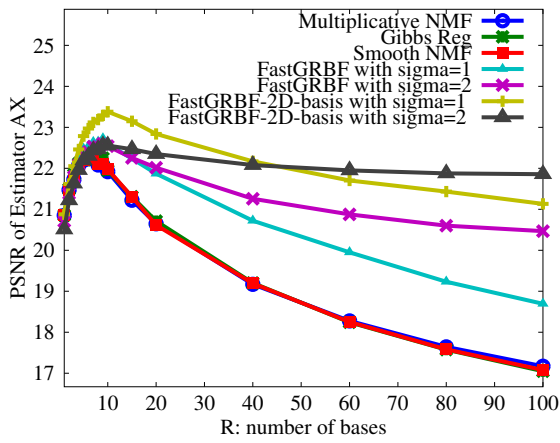


Figure 7.4: PSNR of NTF results on 165 (11x15) faces for various ranks

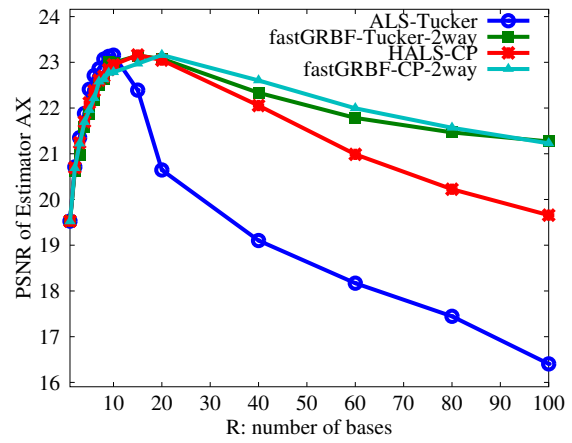
### CBCL faces

In this experiment, I examined the robustness of the proposed methods for Gaussian noise and applied to the two-way function approximation for the tensor decomposition using 100 images of the CBCL face dataset<sup>2</sup>. A face image is a 361-dimensional vector or a (19 x 19)-matrix, and its data matrix is a (361 x 100)-matrix  $\mathbf{Y}$  or a (19 x 19 x 100)-tensor  $\underline{\mathbf{Y}}$ . The Gaussian noise (PSNR=16.2 dB, Fig. 7.5(c)) is added to all images in the dataset. I try to run the NMF with multiplicative updates, the Gibbs Regularized NMF [116, 117], the Essid's Smooth NMF [33], the fastGRBF and the fastGRBF-2Dbasis with  $\sigma = 1.0$  and  $\sigma = 2.0$  for decompositions of various ranks. To tensor data, I applied the function approximation by two ways as  $\mathbf{A} = \Phi_a \mathbf{W}$  and  $\mathbf{B} = \Phi_b \mathbf{V}$ . I denote the methods by fastGRBF-NTD-2way and fastGRBF-NCPD-2way, respectively. The size of the core tensor is changed from

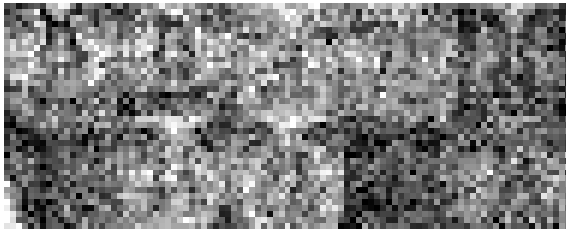
<sup>2</sup>CBCL Face Database #1, MIT Center For Biological and Computation Learning, <http://www.ai.mit.edu/projects/cbcl>



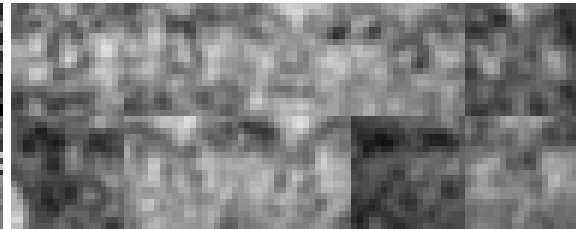
(a) PSNRs of NMF



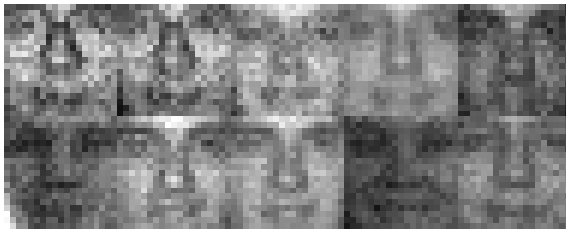
(b) PSNRs of NTF



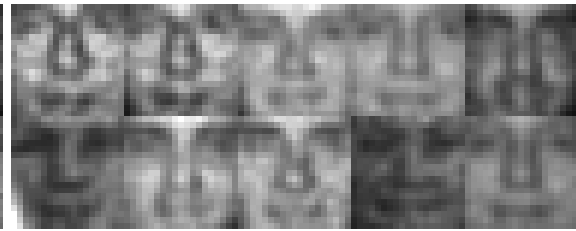
(c) CBCL faces : 60 of 100 images (16.2 dB)



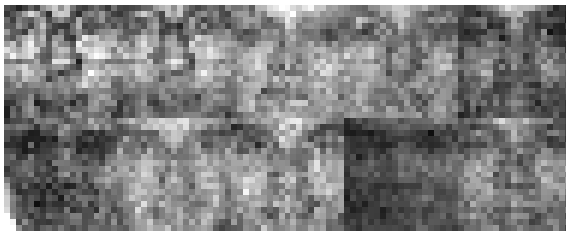
(d) 3x3 moving average filter (22.8 dB)



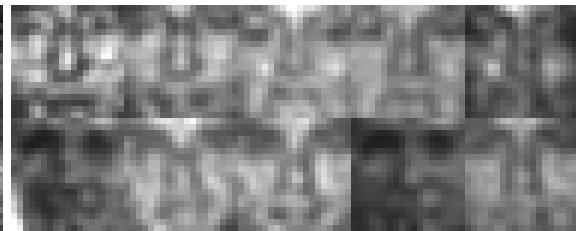
(e) Multiplicative NMF (21.9 dB)



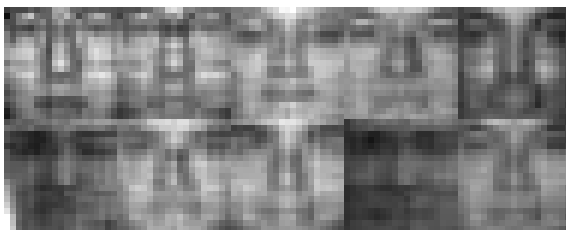
(f) fastGRBF-2Dbasis (23.4 dB)



(g) ALS-NTD (20.7 dB)



(h) fastGRBF-NTD-2way (23.2 dB)



(i) HALS-NCPD (23.2 dB)



(j) fastGRBF-NCPD-2way (23.3 dB)

Figure 7.5: Denoising and face reconstruction experiments on CBCL faces database: noise level (PSNR) is 16.2 dB.

$R_1 = R_2 = R_3 = R = 1$  to  $R_1 = R_2 = R_3 = R = 100$ . I set  $U = 4$ ,  $\delta t = 1$  for the fastGRBF-based methods.

Figure 7.5 shows the results of this experiment. We can see that the proposed method provided better results than the traditional and state-of-the-art methods. The fastGRBFs with  $\sigma = 2.0$  kept higher PSNRs than with  $\sigma = 1.0$  for large  $R$ . The two-way function approximation worked well in NTF. Figure 7.5 (d) shows the result by the moving average filter with 3x3 window for reference. Figures 7.5 (e)-(j) show the reconstructed results by the individual methods with  $R = 10$  for NMF, and with  $R = 20$  for NTF. We can see that the fastGRBF based methods provided smoother images than the traditional methods because the proposed methods are constructed by smooth bases.

In practical cases, we can not measure PSNR between a reconstructed image and an original image since the original image is unknown. Thus, it is difficult to estimate optimal  $R$ . However, we can expect that the fastGRBF-based approach gives better results than traditional approaches in a very wide range of  $R$  if we assume the images are smooth. This fact is remarkable for practical use in signal processing.

### 3D-Tensor Data Reconstruction

In this experiment, I applied the fastGRBF-NTD/NCPD-3way to the reconstruction of a noisy 3D-tensor data (Fig. 7.6(c)). The size of the tensor is  $(50 \times 50 \times 50)$ , and its original data (Fig. 7.6(b)) is generated by a mixture of several 3D-Gaussian functions. Figure 7.6(a) shows the results of SIR (signal to inference ratio) of each estimator for various  $R$ . The fastGRBF-NCPD and the fastGRBF-NTD kept high SIRs for large  $R$  and Figures 7.6(d)-(g) show the robustness of the two methods for noisy data.

#### 7.1.2 Blind source separation

In this experiment, I applied the fastGRBF to the BSS problem by using artificial and real-world datasets. The generative model is given by

$$\mathbf{Y} = [\mathbf{S}\mathbf{X}_0 + \mathbf{E}_0]_+, \quad (7.3)$$

where  $\mathbf{S} \in \mathbb{R}_+^{I \times R}$  is an original source signal matrix,  $\mathbf{X}_0 \in \mathbb{R}_+^{R \times J}$  is a mixing matrix,  $\mathbf{E}_0 \in \mathbb{R}^{I \times J}$  is a Gaussian noise matrix, and  $\mathbf{Y} \in \mathbb{R}_+^{I \times J}$  is an observed signal matrix. The SNR (Signal to Noise Ratio) is defined by

$$\text{SNR} := 10 \log_{10} \left[ \frac{\|\mathbf{S}\mathbf{X}_0\|_F^2}{\|\mathbf{S}\mathbf{X}_0 - \mathbf{Y}\|_F^2} \right]. \quad (7.4)$$

I use the SIR (Signal to Inference Ratio) measure to evaluate the results  $\mathbf{A}\mathbf{X}$  obtained by NMF. The SIR is defined by

$$\text{SIR} := 10 \log_{10} \left[ \frac{\|\mathbf{S}\mathbf{X}_0\|_F^2}{\|\mathbf{S}\mathbf{X}_0 - \mathbf{A}\mathbf{X}\|_F^2} \right]. \quad (7.5)$$

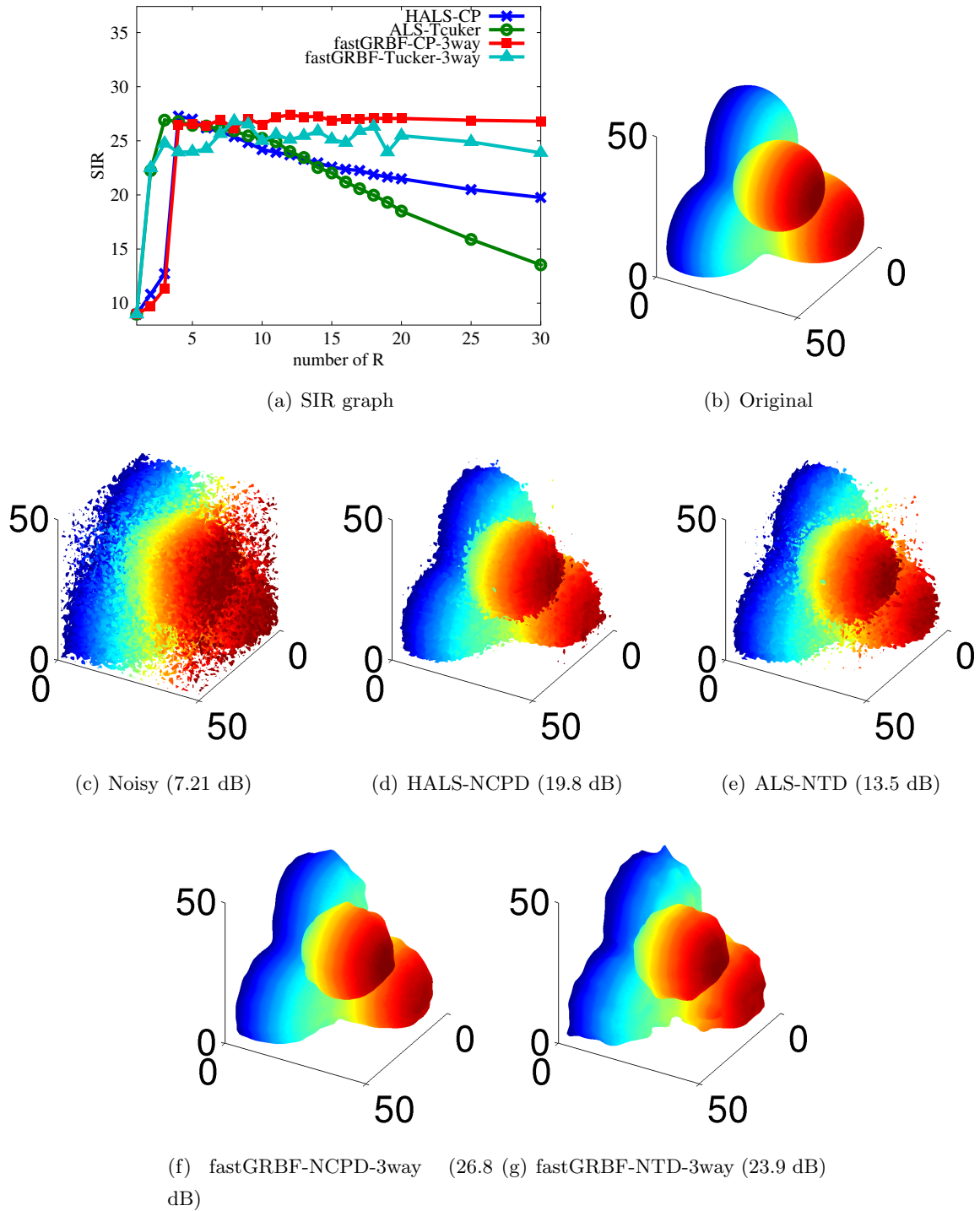


Figure 7.6: Experimental result shown by iso-surfaces at  $R = 30$

By comparing the SIR with the SNR, I can obtain information with respect to the reduction of the noise level. Furthermore, I evaluated the estimated source  $\mathbf{A} = \Phi \mathbf{W}$  by  $\text{SIR}_M$  which is calculated by Algorithm 31. First, the normalization of each signal is necessary because the NMF problem does not have a unique solution. Next, the SIR combination matrix is calculated by  $\mathbf{M}(r_1, r_2) = \text{SIR}(\mathbf{s}_{r_1}, \mathbf{a}_{r_2})$ . For example, assume that  $\mathbf{M}$  is given as

$$\mathbf{M} = \begin{pmatrix} 3.02 & 0.74 & 4.13 \\ 5.37 & 2.38 & 3.30 \\ 3.51 & 3.54 & 5.58 \end{pmatrix}. \quad (7.6)$$

The maximum element of this matrix is  $\mathbf{M}(3, 3) = 5.58$ . Then I set  $-\infty$  on the third row and the third column of  $\mathbf{M}$ :

$$\mathbf{M} = \begin{pmatrix} 3.02 & 0.74 & -\infty \\ 5.37 & 2.38 & -\infty \\ -\infty & -\infty & -\infty \end{pmatrix}. \quad (7.7)$$

In similar way, the next maximum element is  $\mathbf{M}(2, 1) = 5.37$ . Then I set  $-\infty$  on the second row and the first column of  $\mathbf{M}$ :

$$\mathbf{M} = \begin{pmatrix} -\infty & 0.74 & -\infty \\ -\infty & -\infty & -\infty \\ -\infty & -\infty & -\infty \end{pmatrix}. \quad (7.8)$$

Finally,  $\text{SIR}_M$  is given by  $(5.58 + 5.37 + 0.74)/3 = 3.90$ . Actually the maximum combination is given by  $(5.37 + 4.13 + 3.54)/3 = 4.35$ , however we need to calculate  $R!$  combinations to obtain this solution.  $\text{SIR}_M$  gives a useful approximated solution of the combinatorial maximization problem of the mean of SIR for  $\mathbf{A}$ .

First, I use the 'X.5smooth' benchmark from NMFLAB [21]. Figs. 7.7(a) and 7.7(b) show the original sources  $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_R] \in \mathbb{R}_+^{I \times R}$  and the observed signals  $\mathbf{Y}$  where  $R = 5$ ,  $I = 200$ , and  $J = 20$ . Each element of mixing matrix  $\mathbf{X}_0$  is generated by the absolute value of the Gaussian distribution, randomly. I try to run the traditional and state-of-the-art NMF methods of the multiplicative NMF [62], the Kim & Park's Sparse NMF [57], the Regularized ALS Sparse NMF [22], the Gibbs Regularized NMF [116, 117], the Essid's Smooth NMF [33], and the fastGRBF-NMF to separate blind source signals for various noise levels. The parameters of the fastGRBF were used as  $U = 4$ ,  $\delta t = 1$ , and  $\sigma = 1.0$ . Figs. 7.7(c) and 7.7(d) show the result of Fig. 7.7(b) by the multiplicative NMF and the fastGRBF, respectively. Figs. 7.7(e) and 7.7(f) show the result of SIR and  $\text{SIR}_M$  of all the methods. We can see that the fastGRBF is the robustest to noise in both measures, significantly.

Next, I applied the fastGRBF-NMF to BSS of real world data, far-infrared (FIR) spectra from the database of California Institute of Technology<sup>1</sup> in the same conditions, excluding  $\sigma = 4.0$ . Fig. 7.8(a) shows FIR spectra of 'Phosgenite', 'Rhodocrosite', and 'Cancrinite'. Fig. 7.8(b) shows observed mixed

<sup>1</sup>[http://minerals.gps.caltech.edu/FILES/Infrared\\_Far/Index.html](http://minerals.gps.caltech.edu/FILES/Infrared_Far/Index.html)

---

**Algorithm 31** Calculation algorithm of  $\text{SIR}_M$ 

---

- 1: **Input:**  $\mathbf{S}, \mathbf{A} \in \mathbb{R}^{I \times R}$
  - 2: **Initialize:**  $v = 0$
  - 3:  $\mathbf{s}_r \leftarrow \mathbf{s}_r / \|\mathbf{s}_r\|$  for  $r = 1, \dots, R$  ;
  - 4:  $\mathbf{a}_r \leftarrow \mathbf{a}_r / \|\mathbf{a}_r\|$  for  $r = 1, \dots, R$  ;
  - 5: Calculate  $\mathbf{M}(r_1, r_2) = 10 \log_{10} \left[ \frac{\|\mathbf{s}_{r_1}\|^2}{\|\mathbf{s}_{r_1} - \mathbf{a}_{r_2}\|^2} \right]$  for all  $r_1, r_2$ ;
  - 6: **for**  $r = 1, \dots, R$  **do**
  - 7:    $[r_1^*, r_2^*] \leftarrow \operatorname{argmax}_{r_1, r_2} \mathbf{M}(r_1, r_2)$  ;
  - 8:    $v \leftarrow v + \mathbf{M}(r_1^*, r_2^*) / R$  ;
  - 9:    $\mathbf{M}(r_1^*, :) \leftarrow -\infty$  ;
  - 10:    $\mathbf{M}(:, r_2^*) \leftarrow -\infty$  ;
  - 11: **end for**
  - 12: **Output:**  $v$
- 

signals with additive Gaussian noise. Fig. 7.8(c) and Fig. 7.8(d) show the results of estimated sources obtained by the multiplicative NMF and the fastGRBF, respectively. The fastGRBF gave smooth representations of signals, and SIR and  $\text{SIR}_M$  were better than those by the multiplicative NMF. Fig. 7.8(e) and Fig. 7.8(f) show SIR and  $\text{SIR}_M$  for various noise levels with various state-of-the-art methods. We can see that the fastGRBF almost always achieved the best performance in comparison with all the methods and for all noise levels.

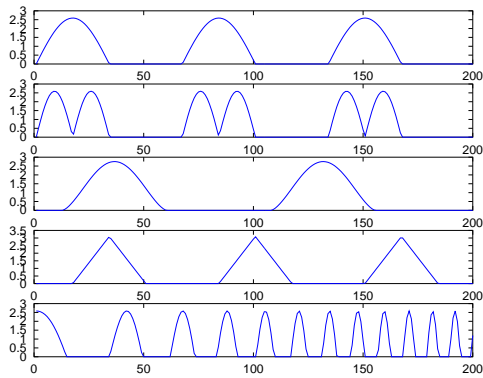
### 7.1.3 Block-based part representation

#### Block-wise Image Factorization

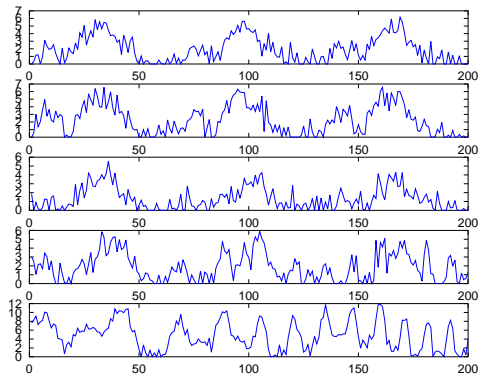
In this section, I apply the fastGRBF-block-NMF to the decomposition of several noise-free ( $64 \times 64$ )-image data. I try to run the fastGRBF methods with the parameters settings shown in Table 7.3. Please note that the parameter in the first line implies a block for an image so that the fastGRBF-NMF and the fastGRBF-block-NMF are equivalent. Other parameters are given as  $U = 1$ ,  $\delta t = 1$ , and  $\sigma = 1$ . Fig. 7.9 shows the results of this experiment. We can see that the methods with smaller block size gave high PSNRs in a short time, however they did not give a smooth representation. On the other hand, the methods with a larger block size gave smooth representations, however they are very slow. Thus, the block size can be regarded as a trade-off parameter. Furthermore, the block distortion occurs in this method along with the DCT.

#### Large-scale image analysis

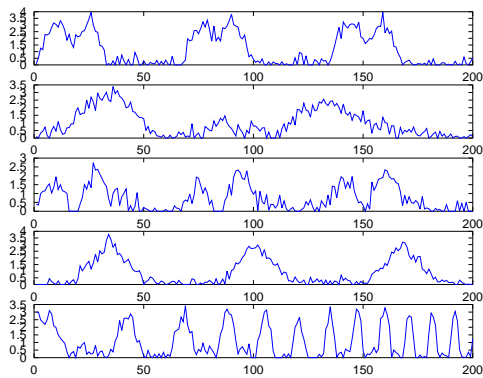
Next I apply the fastGRBF-block-NMF to large-scale image analysis: denoising and local-parts analysis. A noisy ( $3456 \times 4608$ )-image of which noise level is 10 dB is used. Image can be transformed as ( $1024 \times 15552$ )-nonnegative matrix data by the unfolding of individual ( $32 \times 32$ )-blocks. The



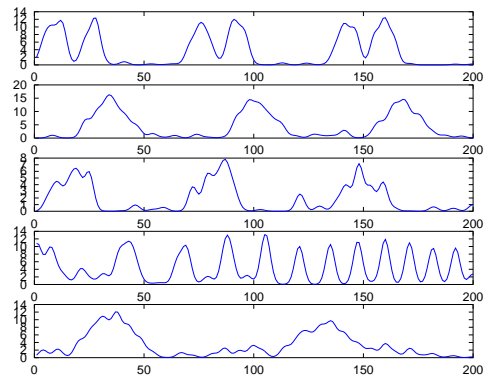
(a) Source signals  $\mathcal{S}$



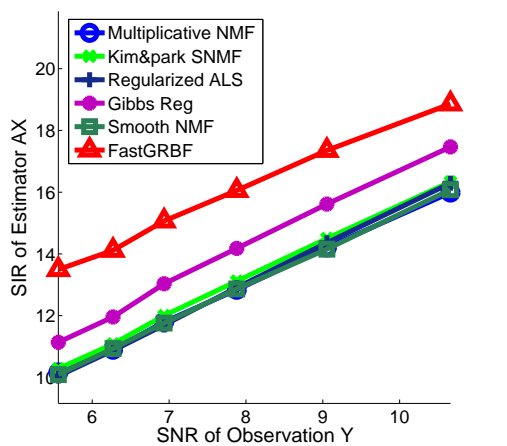
(b) First 5 of 20 observed signals (SNR=12.4 dB)



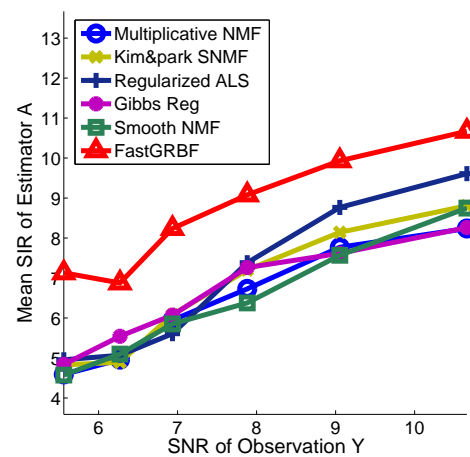
(c) Estimated signals  $\mathbf{A}$  (Multiplicative NMF, SIR=17.6 dB,  $\text{SIR}_M = 9.3$  dB)



(d) Estimated signals  $\mathbf{A}$  (FastGRBF, SIR=19.6 dB,  $\text{SIR}_M=13.3$  dB)

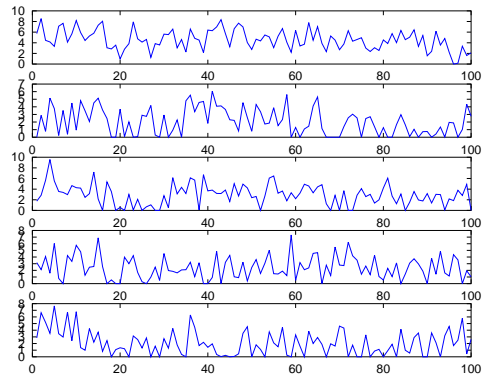
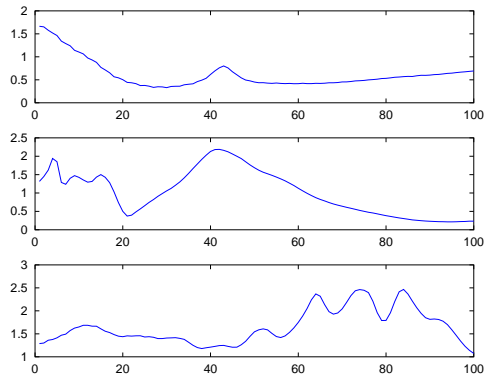


(e) Evaluation with SIR of  $\mathbf{AX}$

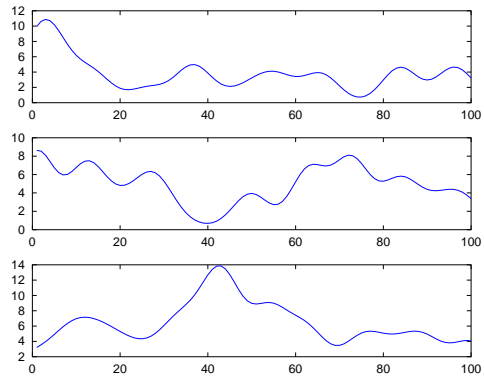
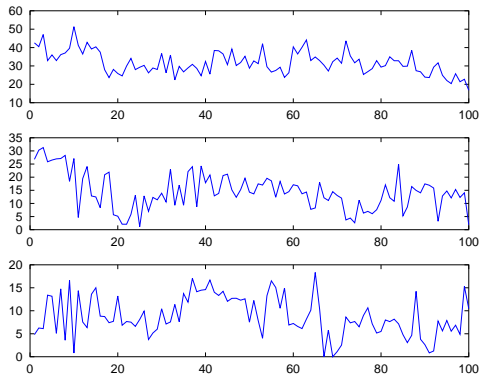


(f) Evaluation with  $\text{SIR}_M$  of  $\mathbf{A}$

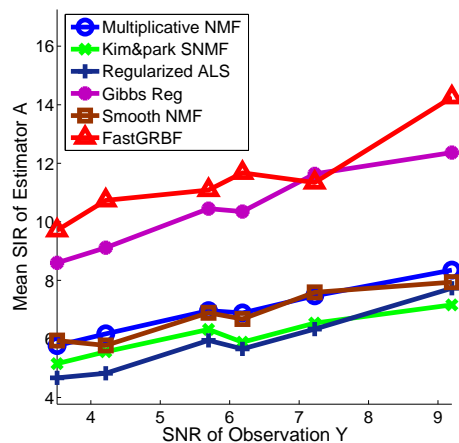
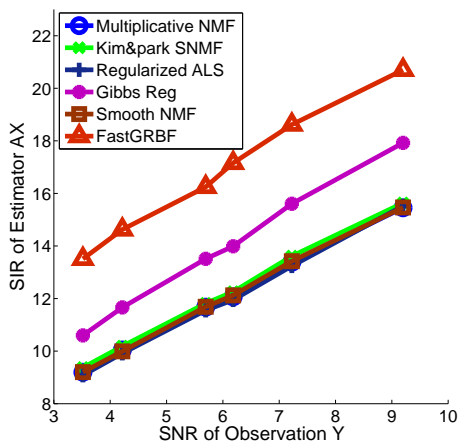
Figure 7.7: Visualizations and experimental results on NMFLAB: 'X\_5smooth'



(a) Source signals  $\mathbf{S}$ : Phosgenite (top), Rhodocrosite (middle), Cancrinite (bottom) (b) First 5 of 20 observed signals (SNR=4.4 dB)



(c) Estimated signals  $\mathbf{A}$  (Gibbs reg. NMF, SIR= 11.5 dB,  $SIR_M = 8.7$  dB) (d) Estimated signals  $\mathbf{A}$  (FastGRBF, SIR= 15.1 dB,  $SIR_M=10.7$  dB)



(e) Evaluation with SIR of  $\mathbf{AX}$

(f) Evaluation with  $SIR_M$  of  $\mathbf{A}$

Figure 7.8: Visualizations and experimental results on far-infrared spectra

Table 7.3: Individual parameters of fastGRBF-block-NMF

block size	# of blocks	# of $I$	# of $R$	# of $J$	# of param.
$64 \times 64$	$1 \times 1$	4096	4	4	16640
$32 \times 32$	$2 \times 2$	1024	4	16	4160
$16 \times 16$	$4 \times 4$	256	4	64	1280
$8 \times 8$	$8 \times 8$	64	4	256	1280
$4 \times 4$	$16 \times 16$	16	4	1024	4160
$2 \times 2$	$32 \times 32$	4	4	4096	16400

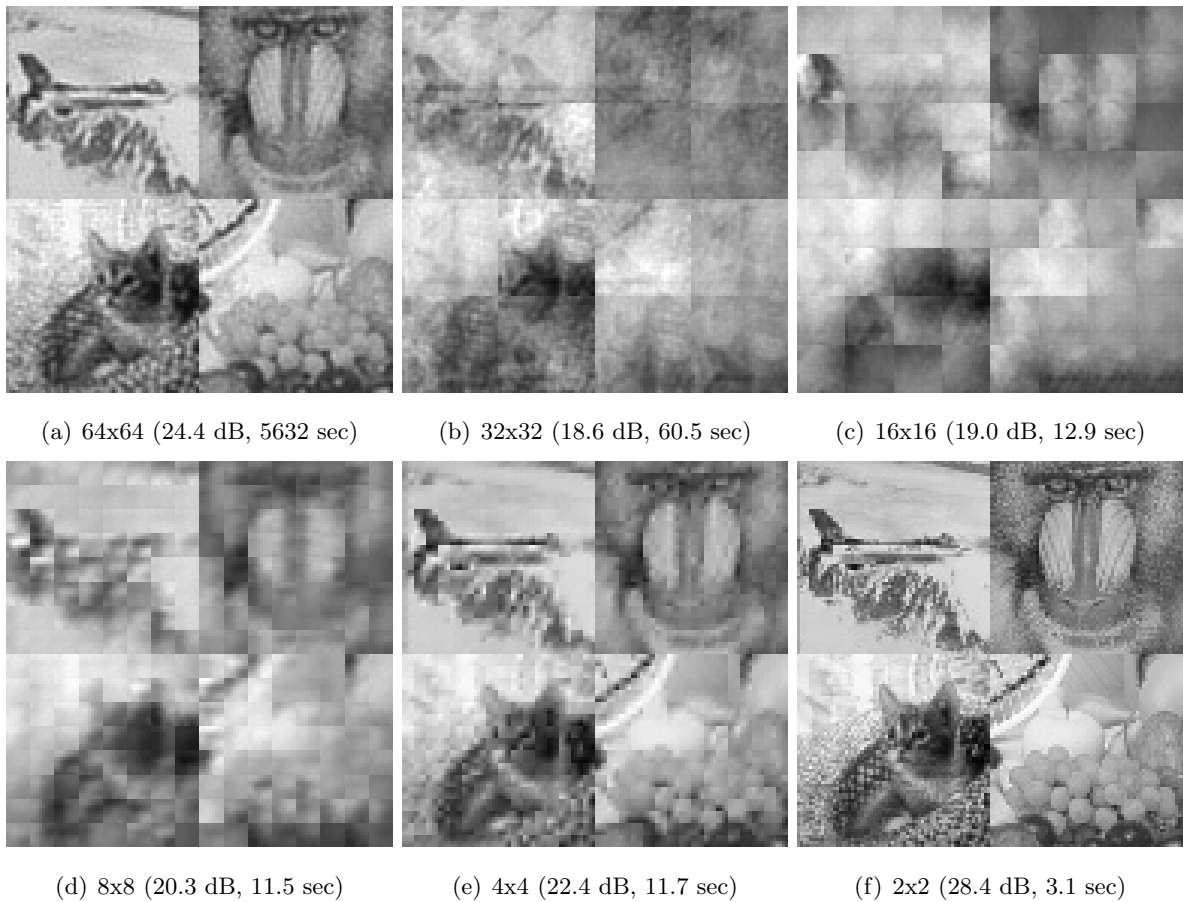


Figure 7.9: Results of fastGRBF-block-NMF

fastGRBF-NMF-2Dbasis is applied to this data, and finally its reconstructed matrix is transformed to a  $(3456 \times 4608)$ -image.

Fig. 7.10 shows the reconstruction result. I tried to compare it with results by the moving average filter, the low-pass (butter worth) filter, and the wavelet shrinkage denoising [31]. Individual top and bottom images are zoomed  $(801 \times 1001)$ - and  $(201 \times 401)$ -images, respectively. Wavelet shrinkage denoising seems to reduce the many noise visually; however, image is blurred and small noise is still remained. In the proposed method, the block distortion can be recognized, clearly. According to the value of PSNR, the low-pass filter is the best method and the stripe pattern of the roof was not destroyed. In this way, the denoising result of the proposed method is not the best; nonetheless, the proposed method provided the second best PSNR, and it can be improved by a deblocking filter.

Next I apply the fastGRBF-block-NMF to the local parts analysis of a noisy large-scale image which is the same as the above one. The  $(1024 \times 15552)$ -matrix data is generated by the same unfolding method as above one. I extract  $R = 20$  parts by the fastGRBF-block-NMF method. Fig. 7.11 shows the result of local parts (reverse image). I tried to compare it with the standard multiplicative NMF [62], the non-smooth NMF (nsNMF) [80], the Chen's smooth NMF [13] and the Gibbs regularized smooth NMF [116, 117]. Only (a) is learned from the original noise free image, and the others (b-f) are learned only from the noisy image (10 dB). We can see that almost all parts includes a lot of noise excluding the nsNMF and the fastGRBF-NMF of which parts are clear. SIRs between the result of the multiplicative NMF from the noise free image and the results of the other methods are (b) 14.36 dB, (c) 16.65 dB, (d) 13.68 dB, (e) 14.42 dB, and (f) 16.99 dB, then the fastGRBF provides the best result among these methods. In fact, a little noise is still remained in the parts of nsNMF. As is well-known, the sparseness is efficiently worked for parts-based representation. The fastGRBF does not imposed any sparse constraint but smoothness by the function approximation, it provide the best result. Thus, the fastGRBF-NMF is promising to combine the smoothness with the sparseness for parts-based representation.

## 7.2 LCPTD

### 7.2.1 Low rank approximation

In this part, the LCPTD was applied to face reconstruction problems and the performances were compared with other models. The Yale face database consists of 165 gray-scale images of 15 individuals. There are 11 images per subject with different facial expressions or configurations. In this experiment, I used 15 full-face images; I took an image from each subject. The size of an image is  $215 \times 171$  pixels, then we considered that  $I_1 = 215$ ,  $I_2 = 171$ , and  $S = 15$ . Furthermore, I prepared salt-and-pepper noised data sets:  $\text{SNR} \in \{5, 10, 20\}$  [dB].

I applied our new LCPTD model, the sparse LCPTD, and the nonnegative LCPTD with various numbers of common components for the noise free and the noisy data sets; the number of bases was fixed as  $J = 40$ , and the number of common components was changed as  $L_n \in \{0, 5, 10, \dots, 40\}$ . I

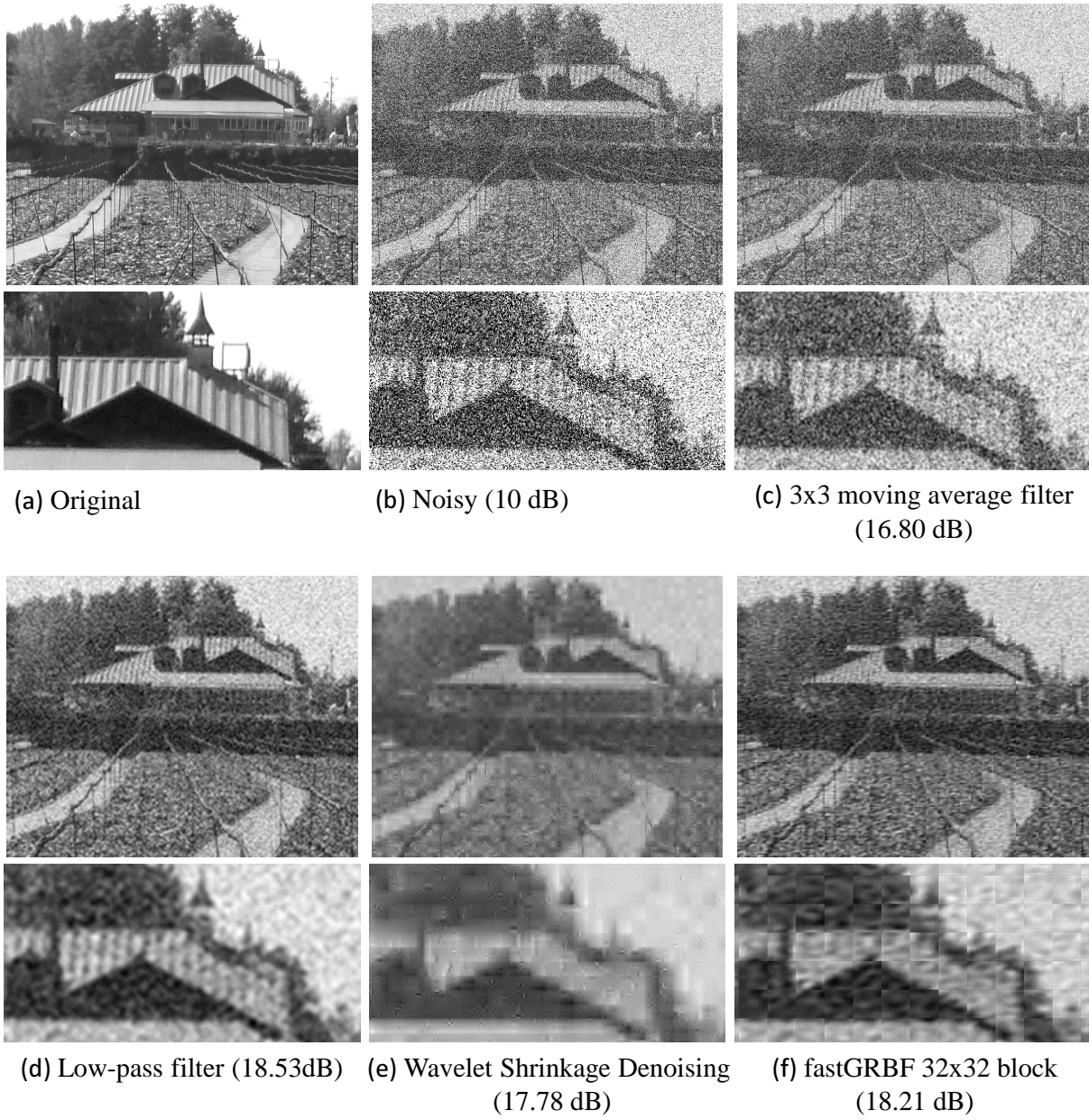


Figure 7.10: Denoising result

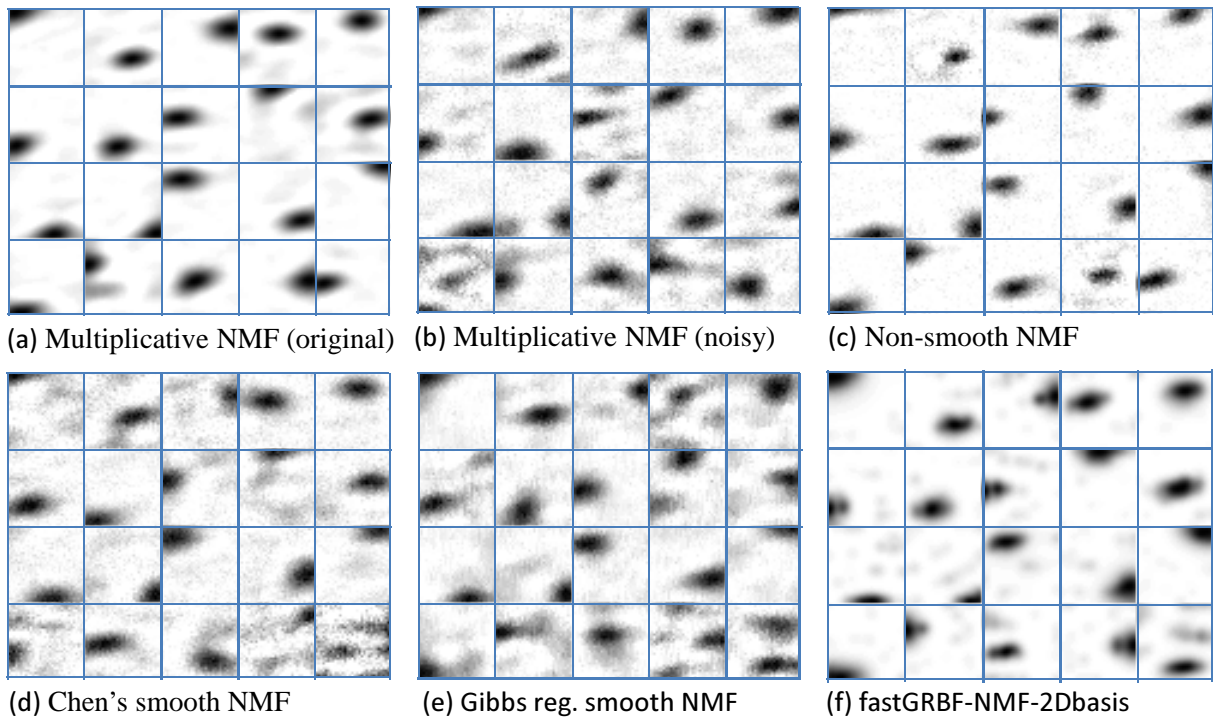


Figure 7.11: Parts-based representation

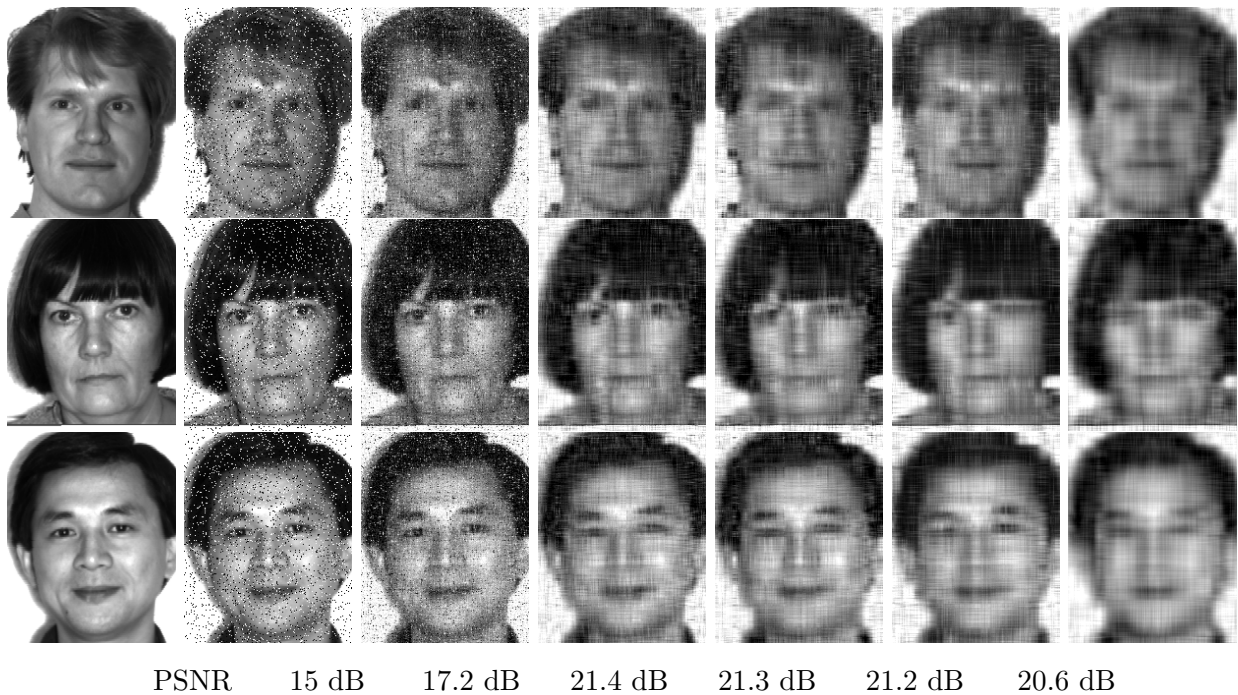


Figure 7.12: Face images corrupted by additive noise and the reconstructed images (PSNR= 15 dB,  $J = 40$ ): 1st column: original images, 2nd column: noisy images, 3rd column: ICPTD model, 4th column: LCPTD ( $L_n = 35$ ), 5th column: LCPTD with sparse constraint ( $L_n = 35$ ), 6th column: LCPTD with nonnegative constraint ( $L_n = 35$ ), 7th column: SCPTD model.

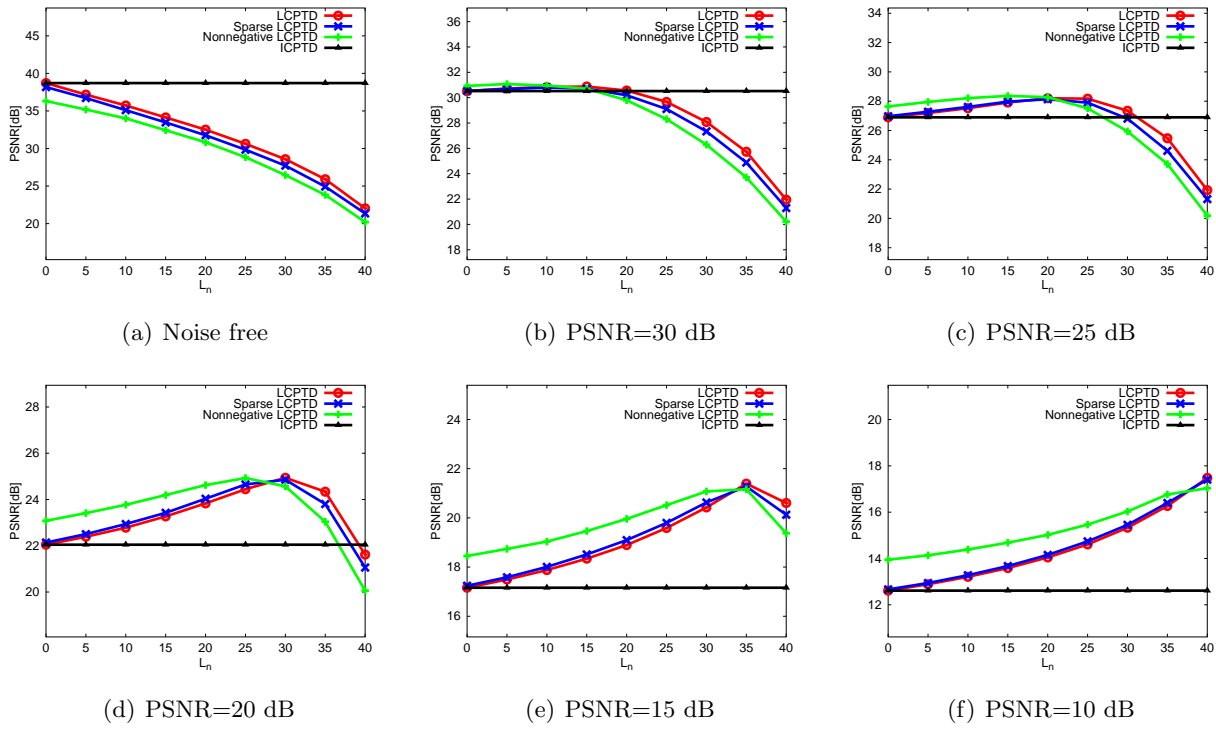


Figure 7.13: PSNRs for various noisy data sets

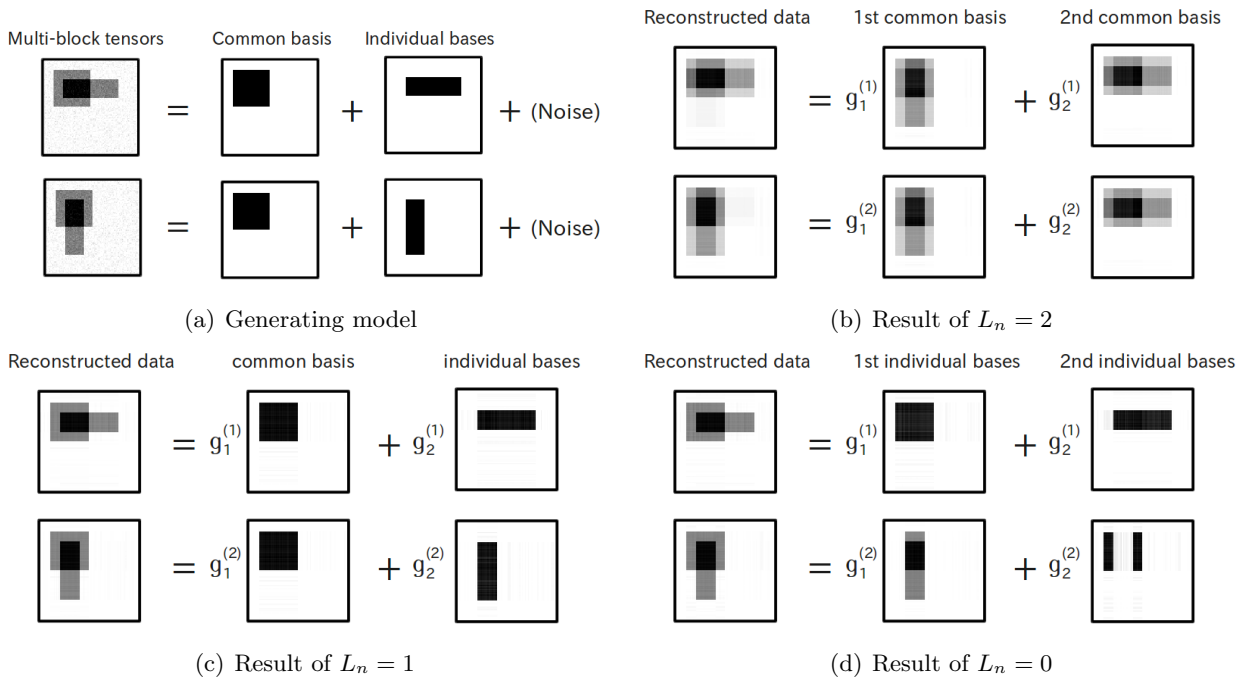


Figure 7.14: Linked Multi-block Tensor Factorization

computed the PSNRs between the original faces and the reconstructed faces.

Fig. 7.12 depicts the results of face reconstruction. We can see that the ICPTD model couldn't reduce noise well, and the SCPTD model was robust with respect to noise but reconstructed faces were too fuzzy (distorted). However, the LCPTD based methods gave the appropriate and intermediate reconstructions comparing to the two models.

Fig. 7.13 depicts the graphs of PSNRs with various number of common components. We can see that if the noise level is larger, the maximum points of PSNR is obtained by a larger  $R$ . In noise free data set Fig. 7.13 (a), the maximum PSNR was obtained at  $L_n = 0$  for all methods; so the ICPTD model is the best in this case. On the other hand, the maximum PSNRs were obtained by the LCPTD based methods in Fig. 7.13 (b,c,d,e). It is also interesting that the nonnegative LCPTD kept a high PSNR for smaller number of common components in comparison with the other methods in high noise level (see Fig. 7.13 (d,e,f)). Thus, we can say that the nonnegative constraint is useful for this problem.

In general, because real data often includes some noise, the proposed method could be very useful and practical for the real tensor data analysis. The higher noise level requires a large number of common components, but the multitude of common components is often only harmful for fitting. Then, we have to select the best parameter of  $L_n$  and its selection method is an open problem. We may be able to select  $L_n$  depending on PSNR if it is known as prior information.

### 7.2.2 Linked Blind Source Separation

In this part, I applied the LCPTD to the blind source separation of a simple data for the linked multi-block tensor factorization. I generated two block data tensors consisting of a one common basis factor and two individual basis factors with noise (see Fig. 7.14(a)). And I decompose them by our LCPTD model with nonnegative constraints for various number of common bases  $L_n \in \{2, 1, 0\}$  for  $n = 1, 2$ . Figs. 7.14 (b,c,d) depict the results of this experiment. It is obvious that the result of  $L_n = 2$  couldn't represent the original data tensors since its degree of freedom of model is not sufficient. On the other hand, the result of  $L_n = 0$  could represent the original data tensors, but each basis is not matched, completely. The result of  $L_n = 1$  could represent not only the original data tensors, but also each basis; besides, the additive noise were reduced.

We can see that the LCPTD model can be very useful assuming that some components are common in generating model. The blind source separation can separate into two original sources from two observed signals. It is very interesting that the LCPTD can separate three bases (i.e., a common and two individual bases) from only two observed tensors. We should note that the selection of  $L_n$  could be a very important deciding factor to obtain proper decomposition for the LCPTD method from this experiment.

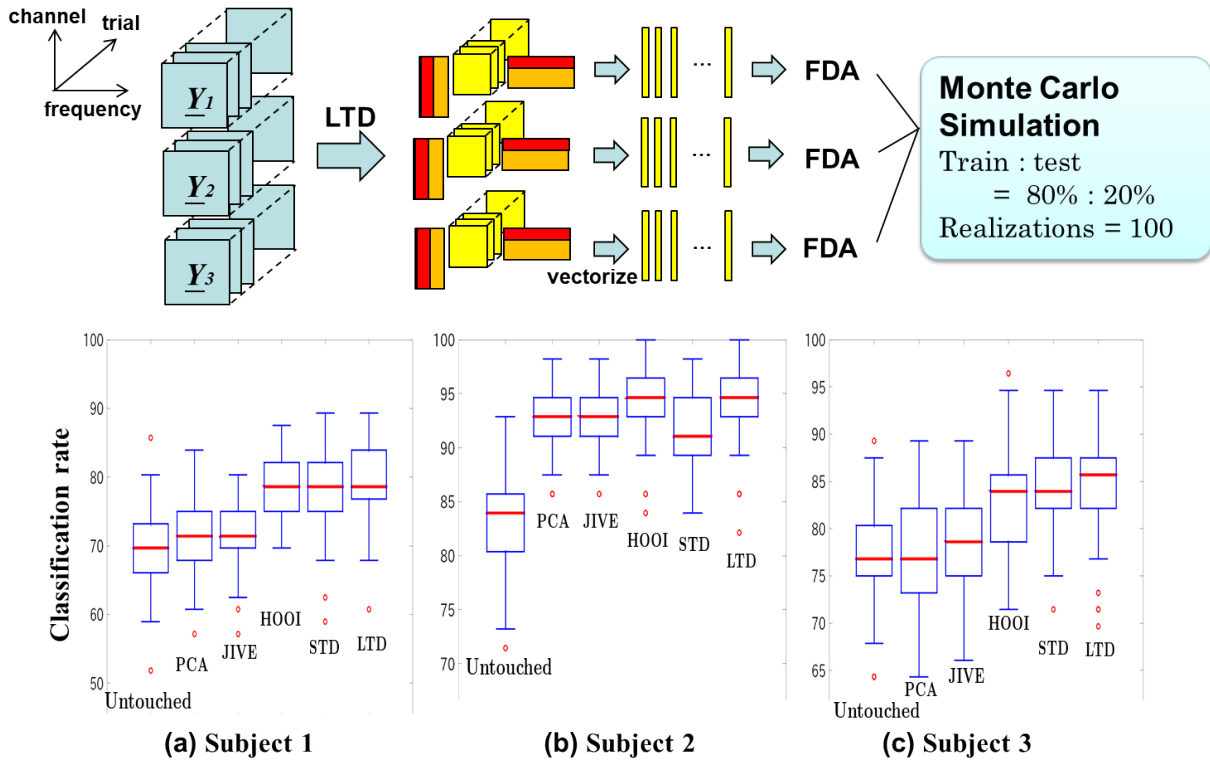


Figure 7.15: Experiments for LTD

## 7.3 LTD

### 7.3.1 Feature extraction for classification

In this experiment, I applied the LTD to the feature extraction for the motor imagery BCI. Figure 7.15 shows a flow of this experiment. At first, I used three frequency power tensors which were taken from three persons. Each tensor is constructed by 21 channels, 25 frequency bands, and 280 trials. I decomposed these three tensors by the orthogonal LTD with  $L_1 = L_2 = R_1 = R_2 = 4$  without the third-mode factorization. As the result of LTD, we can obtain three core tensors. After that I vectorize each trial matrix, and try to classify by the FDA. In this regard, I iterate to calculate classification rates for 100 times. In each iteration, I randomly separate the dataset into two parts for train (80%) and test (20%), respectively. I compared the LTD with the untouched data, the PCA, the JIVE, the HOOI, and the STD. Figure 7.15 (a,b,c) show the results of statistics of classification rates by individual feature extraction methods. We can see that the LTD outperform the other feature extraction methods.

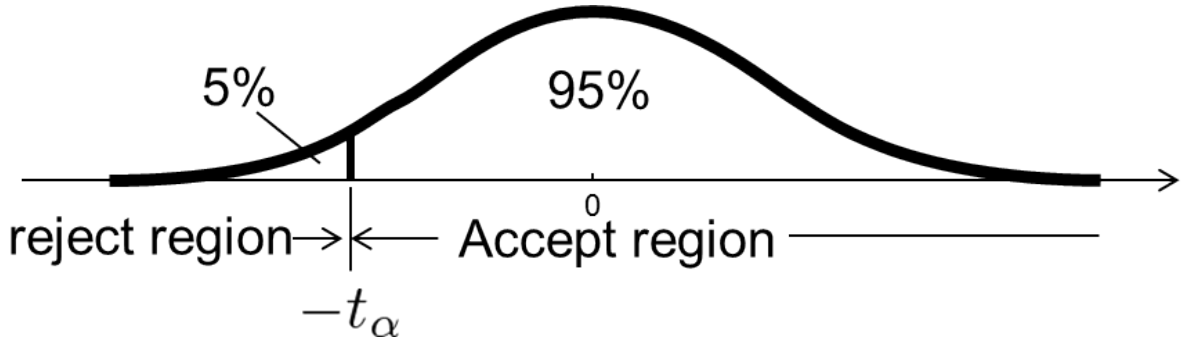


Figure 7.16: t-distribution and reject/accept region with significant level 5%

## 7.4 Classification

### 7.4.1 One-tailed t-test

Before showing the experimental results, I introduce the one-tailed t-test that a method to compare a result of proposed method with a result of existing method. The one-tailed t-test plays important roles in statistical data analysis. The classifiers are evaluated by misclassification rates. The lower misclassification rate, the better. The population means of the misclassification rates with some dataset by the proposed method and the existing method are denoted as  $\mu_1$  and  $\mu_2$ , respectively. We test that  $\mu_1$  is significantly smaller than  $\mu_2$ . Thus, the null hypothesis  $\mathcal{H}_0$  is given by  $\mu_1 = \mu_2$ , and the alternative hypothesis  $\mathcal{H}_1$  is given by  $\mu_1 < \mu_2$ . If the null hypothesis is rejected and the alternative hypothesis is accepted with some significant level  $\alpha = \alpha_0$ , then we can say that  $\mu_1$  is significantly smaller than  $\mu_2$  with significant level  $\alpha_0$ .

In order to test  $\mathcal{H}_0$  and  $\mathcal{H}_1$ , samples of misclassification rates by individual methods, that are denoted by  $\{x_1^{(1)}, x_1^{(2)}, \dots, x_1^{(N_1)}\}$  and  $\{x_2^{(1)}, x_2^{(2)}, \dots, x_2^{(N_2)}\}$ , are necessary, where  $N_1$  and  $N_2$  are numbers of samples. Sample means  $\{\bar{x}_1, \bar{x}_2\}$  and sample variances  $\{s_1^2, s_2^2\}$  are given by

$$\bar{x}_i = \frac{1}{N_i} \sum_{n=1}^{N_i} x_i^{(n)}, \quad (7.9)$$

$$s_i^2 = \frac{1}{N_i - 1} \sum_{n=1}^{N_i} (x_i - \bar{x}_i)^2, \quad (7.10)$$

for  $i = 1, 2$ . In the one-tailed t-test, we consider the value of  $t$  which is given by

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{s \sqrt{1/N_1 + 1/N_2}}, \quad (7.11)$$

where

$$s^2 = \frac{(N_1 - 1)s_1^2 + (N_2 - 1)s_2^2}{N_1 + N_2 - 2}, \quad (7.12)$$

and assume that  $t$  follows the t-distribution with the degree of freedom is  $N_1 + N_2 - 2$ . Fig. 7.16 shows the image of t-distribution and its reject/accept region with significant level 5%. In this figure,  $-t_\alpha$

is a boundary between reject and accept regions. Assuming that  $\mathcal{H}_0$  is true, then  $t$  is given by

$$t_0 = \frac{(\bar{x}_1 - \bar{x}_2)}{s\sqrt{1/N_1 + 1/N_2}}. \quad (7.13)$$

Under this assumption, if  $t_0 < -t_\alpha$ , then  $\mathcal{H}_0$  is rejected, and if  $t_0 \geq -t_\alpha$ , then  $\mathcal{H}_0$  is accepted. Therefore, we can test that the proposed classifier is significantly better than the existing classifier by the one-tailed t-test with some significant level.

## 7.4.2 WRSVM

### UCI database

In this experiment, I used thirteen UCI data sets for binary problems to compare the two novel classifiers SVM<sub>G</sub> and SVM<sub>D</sub> with SVM and L1-norm regularized SVM (L1-SVM).

These data sets are summarized in Table 7.4, which lists the data set name, the respective numbers of training samples, test samples, realizations, and dimensions.

Several hyper parameters must be optimized, namely, the kernel parameter  $\gamma$ , the adjusting parameter  $c$ , and the weighting parameters  $\lambda$  of  $Q_2(\mathbf{x})$  and  $\nu$  of  $Q_3(\mathbf{x})$ , but  $\rho = 0.9$  is fixed. These parameters are optimized on the first five realizations of each data set. The best values of each parameter are obtained by using each realization. Finally, the median of the five values is selected. After that, the classifiers are trained and tested for all of the remaining realizations (i.e., 95 or 15 realizations) by using the same parameters.

Table 7.5 contains the results of this experiment. The values in the classifier name column show “average  $\pm$  standard deviation” of the error rates for all of the remaining realizations, and the minimum values among all classifiers are marked with bold font. The values in the columns for  $\lambda$  and  $\nu$  show the value selected through model selection for each data set. The signs in columns L2 and L1 show the results of a significance test (t-test with  $\alpha = 5\%$ ) for the differences between the SVM<sub>G</sub>/SVM<sub>D</sub> and the SVM/L1-SVM, respectively. “+” indicates that the error obtained with the novel classifier is significantly smaller, while “-” indicates that this error is significantly larger. The penultimate line for “Mean %”, is computed by using the average values for all data sets as follows. First, I normalize the error rates by taking

$$\left\{ \frac{(\text{particular value})}{(\text{minimum value})} - 1 \right\} \times 100[\%] \quad (7.14)$$

for each data set. Next, the “average” values are computed for each classifier. This evaluation method is taken from [86]. The last line shows the average of the p-value between “particular” and “minimum” (i.e., the minimum p-value is 50 %).

SVM<sub>G</sub> provides the best results for two data sets. Compared to the SVM, the SVM<sub>G</sub> is significantly better for four data sets and significantly worse for five data sets. Compared to the L1-SVM, the SVM<sub>G</sub> is significantly better for five data sets and significantly worse for three data sets.

Table 7.4: UCI Data sets

No.	Name	# of training samples	# of test samples	# of realizations	Dimensions
1	Banana	400	4900	100	2
2	Breast Cancer	200	77	100	9
3	Diabetes	468	300	100	8
4	Flare-Solar	666	400	100	9
5	German	700	300	100	20
6	Heart	170	100	100	13
7	Image	1300	1010	20	18
8	Ringnorm	400	7000	100	20
9	Splice	1000	2175	20	60
10	Thyroid	140	75	100	5
11	Titanic	150	2051	100	3
12	Twonorm	400	700	100	20
13	Waveform	400	4600	100	21

Table 7.5: Experimental results

	$\lambda$	SVM <sub>G</sub>	L2	L1	$\nu$	SVM <sub>D</sub>	L2	L1	SVM	L1-SVM
Banana	1	10.6 ± 0.4	+		8	<b>10.4 ± 0.4</b>	+		11.5 ± 0.7	10.5 ± 0.4
B.Cancer	.01	26.3 ± 5.2			4	26.0 ± 4.4			26.0 ± 4.7	<b>25.4 ± 4.5</b>
Diabetes	100	24.1 ± 2.0	-	-	8	24.0 ± 2.0	-	-	23.5 ± 1.7	<b>23.4 ± 1.7</b>
F.Solar	10	36.7 ± 5.2	-	-	2	<b>32.4 ± 1.8</b>		+	<b>32.4 ± 1.8</b>	32.9 ± 2.7
German	10	24.0 ± 2.4			16	24.7 ± 2.3			<b>23.6 ± 2.1</b>	24.0 ± 2.3
Heart	10	<b>15.3 ± 3.2</b>			2	15.6 ± 3.2			16.0 ± 3.3	15.4 ± 3.4
Image	100	4.3 ± 0.9	-		16	4.1 ± 0.8	-	+	<b>3.0 ± 0.6</b>	4.8 ± 1.3
Ringnorm	100	<b>1.5 ± 0.1</b>	+	+	8	<b>1.5 ± 0.1</b>	+	+	1.7 ± 0.1	1.6 ± 0.1
Splice	10	11.3 ± 0.7	-	+	8	11.0 ± 0.5		+	<b>10.9 ± 0.7</b>	12.4 ± 0.9
Thyroid	1	7.3 ± 2.9	-	-	2	<b>4.1 ± 2.0</b>	+	+	4.8 ± 2.2	5.4 ± 2.4
Titanic	.01	22.4 ± 1.0		+	2	22.5 ± 0.5		+	<b>22.4 ± 1.0</b>	23.0 ± 2.1
Twonorm	10	2.4 ± 0.1	+	+	4	<b>2.3 ± 0.1</b>	+	+	3.0 ± 0.2	2.7 ± 0.2
Waveform	1	9.7 ± 0.4	+	+	2	<b>9.6 ± 0.5</b>	+	+	9.9 ± 0.4	10.1 ± 0.5
Mean %		12.0				<b>4.2</b>			6.1	10.7
P-value %		87.2				<b>71.5</b>			79.2	87.5

Furthermore, the SVM<sub>D</sub> provides the best results for six data sets. Compared to the SVM, the SVM<sub>D</sub> is significantly better for five data sets and significantly worse for two data sets. Compared to the L1-SVM, the SVM<sub>D</sub> is significantly better for eight data sets and significantly worse for one data sets.

According to the results for both “mean” and “p-value”, the SVM<sub>D</sub> classifier is the best among the four classifiers considered.

### 7.4.3 CFDA

#### Toy Simulation I

In this simulation, I show the classification error with respect to the difference of two covariance matrices. I assume that there are two normal distributions in the original feature space characterized by

$$\boldsymbol{\mu}_1 = \begin{pmatrix} -1 \\ -1 \end{pmatrix}, \quad \boldsymbol{\Sigma}_1 = \begin{pmatrix} 4 & 0 \\ 0 & 4 \end{pmatrix},$$

$$\boldsymbol{\mu}_2 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \boldsymbol{\Sigma}_2 = \begin{pmatrix} 1 & 0 \\ 0 & q^2 \end{pmatrix}.$$

The two prior probabilities are equal (i.e.,  $p_1 = p_2 = 0.5$ ), and I let the weighting parameters as  $\rho_1 = \rho_2 = 0.5$ . I evaluate classification errors of the FDA, the Loog's approximated Chernoff FDA (aCFDA) [67], the BFDA, and the CFDA-m by changing the value of  $q$ . Figure 7.17 shows the result of this simulation. When  $q = 1$ , all criteria gave the same result since the two covariance matrices have the same shape. The BFDA and the CFDA-m gave almost the same results. When the value of  $q$  increase, the difference of errors between the FDA and the BFDA/CFDA-m increase. Note that classification errors of the BFDA and the CFDA-m were dramatically decrease from  $q = 8$ . The BFDA and the CFDA-m are very effective when covariance matrices are very different. The aCFDA error was larger than the FDA from  $q = 5$  to  $q = 7$  although it almost reached to the BFDA/CFDA-m from  $q = 8$ .

#### Toy Simulation II

In this simulation, I show the classification error with respect to the difference of prior probabilities of two classes. I assume that there are two classes of samples in a two dimensional space. The number of samples of class 1 is fixed as  $N_1 = 100$ , and that of class 2 is changed as  $N_2 \in \{100, 150, 200, 300\}$ . I assumed the samples of each class are distributed normally, and generated samples by Matlab/Octave function 'normrnd'. Mean vectors and covariance matrices are calculated from samples.

Fig. 7.18 shows the result of the classification errors with respect to the difference of prior probabilities. The FDA and the aCFDA gave the same result at  $N_2/N_1 = 1$ . The aCFDA and the CFDA- $p$  gave similar results. When  $N_2/N_1$  increases, the aCFDA and the CFDA- $p$  almost reached the CFDA-m. Note that the CFDA-m always gave the best results.

Fig. 7.19 shows the samples and the projected space of the FDA, the aCFDA, the CFDA- $p$ , and the CFDA-m at  $N_2 = 150$ . Class 1 and 2 are denoted by 'x' and 'o', respectively. The obtained projected one-dimensional spaces were different among the methods. They can be ranked in the order of the CFDA-m, the CFDA- $p$ , the aCFDA, and the FDA.

Fig. 7.20 shows two distributions in the projected one-dimensional space of individual methods at  $N_2 = 150$ . They are approximations by normal distribution with the sample mean and the sample

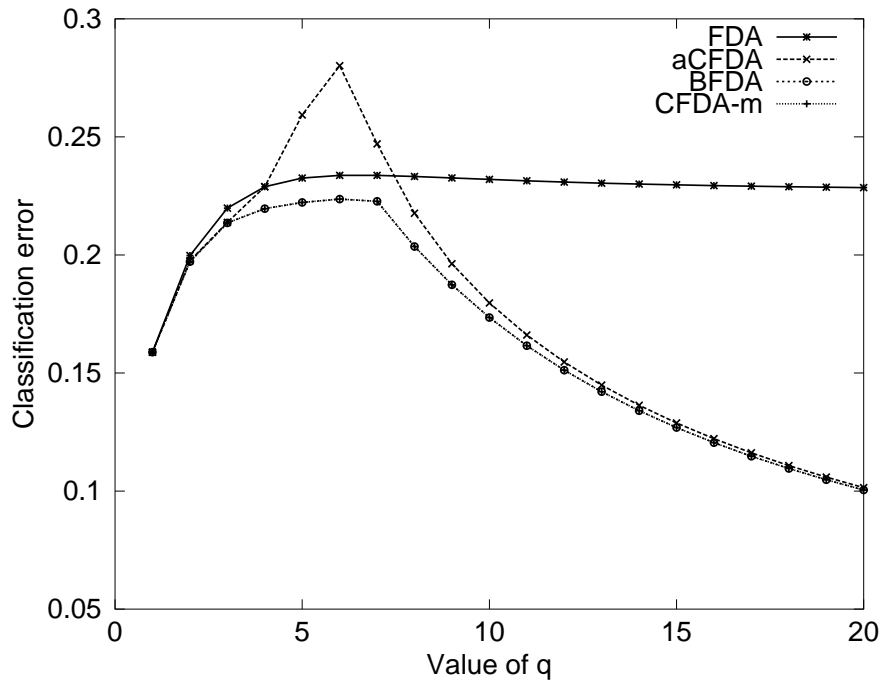


Figure 7.17: Error rates with difference of covariance matrices

variance. Results of the kernel based methods are calculated by (6.167) and (6.168) with  $\gamma = 10^{-1}$ . Note that one distribution is very narrow and the another distribution is wide in the Chernoff based methods. These can be ranked in the order of the CFDA-m, the CFDA- $p$ , and the aCFDA. The KCFDA- $p$  worked well in the feature space. Note that the KCFDA-m gave an impressive result that the variance of one distribution extremely approached to zero.

### UCI database

In this section, I show the results for actual classification datasets. I used nine datasets from the UCI machine learning repository [35]. Table 7.6 describes their properties. Since most of datasets are not for binary classification, I separated them into two classes or focused on two classes. Two datasets of “MNIST” and “Poker hand” have too many instances to apply the kernel method. Therefore, I trained a kernel classifier by using randomly selected 500 instances for training of their datasets. I calculate the 5-fold cross validation (5-CV) errors for all datasets and all methods. In kernel methods, I selected the best kernel parameter  $\gamma$  by hand tuning from candidates of  $\{10^{-10}, 10^{-9}, \dots, 10^2, 10^3\}$ .

Table 7.7 describes the averages of maximized values of the Chernoff distance obtained by the CFDA-algorithm and the Rueda’s gradient-based algorithm [89]. The Rueda’s algorithm can not maximize the Chernoff distance, sometimes. On the other hand, our algorithm provided larger or equal value of the Chernoff distance to the Rueda’s algorithm for all datasets.

Table 7.8 describes the number of iterations for main updates in these methods. The values show

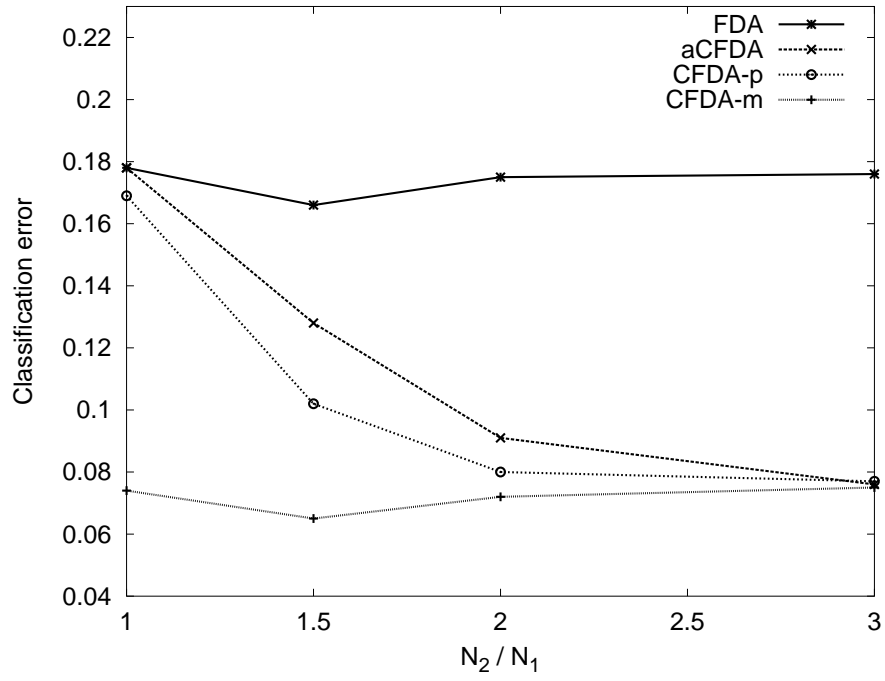


Figure 7.18: Error rates with difference of prior probabilities

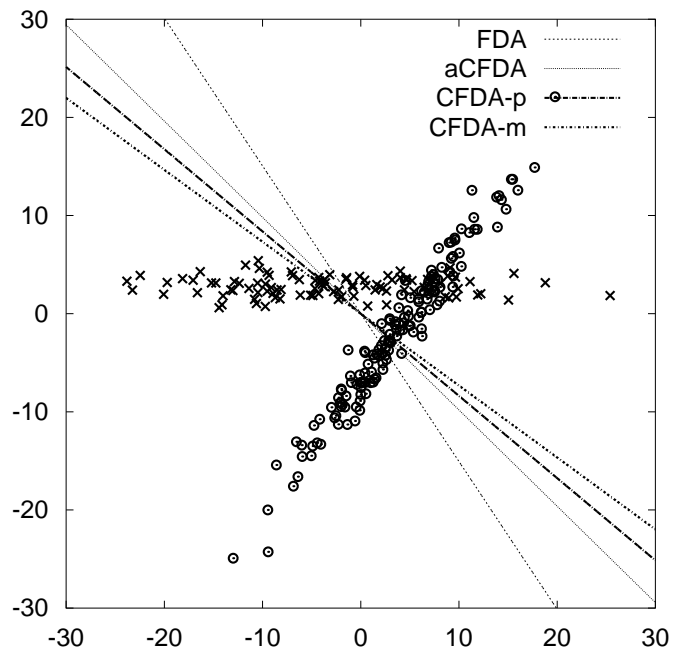


Figure 7.19: Samples and projected space ( $N_2/N_1 = 1.5$ )

Table 7.6: Properties of UCI datasets

label	name	attributes	instance	class	binalize	$p_1$	$p_2$
1	Breast cancer wisconsin	9	569	2	1 vs 2	0.65	0.35
2	Pima indians diabetes	8	768	2	1 vs 2	0.35	0.65
3	SPECTF heart	44	267	2	1 vs 2	0.79	0.21
4	Thyroid gland	5	215	2	1 vs 2	0.30	0.70
5	Iris plants	4	150	3	1,3 vs 2	0.67	0.33
6	Wine	13	178	3	1,2 vs 3	0.73	0.27
7	Glass identification	9	163	7	1,2 vs 3,4	0.90	0.10
8	MNIST digit-recognition	256	1000	10	4 vs 9	0.50	0.50
9	Poker hand	10	25010	10	1 vs 2	0.46	0.54

“average  $\pm$  standard deviation” of number of iterations. Almost all results of the CFDA-algorithm converged by several iterations. The numbers of iterations of the CFDA-m-algorithm are about nine times of those of the CFDA-algorithm, respectively. These results show that the convergence of the CFDA-algorithm does not depend on the number of attributes and instances. However, we can see more iterations are necessary for the kernel extensions. The Rueda’s gradient-based algorithm [89] needed several iterations to converge for almost all data set, however it could not maximize the Chernoff distance for several datasets signed by ‘\*’.

Table 7.9 describes the results of misclassification rates for nine UCI datasets. I described “average  $\pm$  standard deviation” of 20 realizations of classification errors, which are calculated by using two-third instances for training and one-third instances for test, and the average of estimated values of  $\rho_1$  in the CFDA- $p$ , the CFDA-m, the KCFDA- $p$ , and the KCFDA-m. The Rueda’s method is not included since it provides the same result of the CFDA- $p$  when it reached to the maxima of Chernoff distance. The minimum errors are signed by ‘o’ and significantly better results than existing methods are described by bold font. They were calculated by 5 % significant test. In strictly, it is the one-tailed t-test.

Proposed methods provided minimum error results for five datasets. In linear methods, the proposed methods provided significant improvements for “SPECTF”, “Iris”, and “Poker hand”. In kernel methods, the proposed methods provided significant improvements for “SPECTF”, “Glass”, and “Poker hand”. On the other hand, the linear Bhattacharyya/Chernoff FDA provided clearly wrong results for “Thyroid” and “Glass”. These results might be derived from over-fittings or non-Gaussianity. The linear aCFDA provided a significant result for “Thyroid” because of proper fit to heteroscedastic Gaussians. For “Glass”, extending a linear discriminant to a non-linear discriminant provided dramatically improvements, and this dataset is made from a complex combination of heteroscedastic and non-Gaussian distributions.

The difference among the BFDA, the CFDA- $p$  and the CFDA-m is only in values of the hyper-parameters  $\rho_1$  and  $\rho_2$ . The estimation of the best value of  $\rho_1$  is still a open problem for the Chernoff approaches.

Table 7.7: Maximized values of Chernoff distance by each algorithm

label	1	2	3	4	5	6	7	8	9
CFDA- $p$	5.71	0.437	4.50	3.09	1.02	7.33	0.851	15.0	0.0342
Gradient	5.71	0.436	1.37	3.09	1.00	0.905	0.340	15.0	0.00424

Table 7.8: Number of iterations for main updates in each method

label	Gradient	BFDA	CFDA- $p$	CFDA-m	KBFDA	KCFDA- $p$	KCFDA-m
1	$4.0 \pm 0.0$	$2.4 \pm 0.5$	$2.4 \pm 0.6$	$23.3 \pm 2.1$	$5.5 \pm 2.0$	$6.7 \pm 2.5$	$58.4 \pm 17.8$
2	$2.6 \pm 0.7$	$2.8 \pm 0.9$	$3.1 \pm 1.0$	$27.6 \pm 2.9$	$11.8 \pm 4.4$	$12.9 \pm 3.8$	$89.5 \pm 25.8$
3	* $111.8 \pm 102.6$	$4.4 \pm 0.6$	$4.5 \pm 0.5$	$41.2 \pm 4.9$	$6.2 \pm 0.6$	$6.5 \pm 0.6$	$56.5 \pm 2.8$
4	$17.7 \pm 2.7$	$6.0 \pm 0.6$	$5.8 \pm 0.5$	$52.2 \pm 4.0$	$5.7 \pm 2.5$	$6.5 \pm 2.5$	$57.7 \pm 18.4$
5	$21.1 \pm 5.3$	$5.9 \pm 0.8$	$6.7 \pm 1.5$	$54.7 \pm 8.3$	$9.1 \pm 3.8$	$10.7 \pm 3.7$	$87.2 \pm 35.4$
6	* $656.5 \pm 450.4$	$2.9 \pm 0.4$	$2.9 \pm 0.5$	$28.0 \pm 1.5$	$4.4 \pm 0.6$	$8.0 \pm 1.3$	$46.0 \pm 2.9$
7	* $20.5 \pm 16.0$	$5.1 \pm 0.6$	$5.1 \pm 0.6$	$47.0 \pm 5.6$	$5.5 \pm 0.9$	$5.2 \pm 0.7$	$55.5 \pm 4.2$
8	$2.0 \pm 0.0$	$5.9 \pm 0.9$	$6.2 \pm 0.9$	$58.2 \pm 5.4$	$4.0 \pm 0.5$	$4.0 \pm 0.5$	$47.6 \pm 4.8$
9	* $2.4 \pm 1.3$	$3.5 \pm 0.5$	$3.5 \pm 0.5$	$31.6 \pm 3.0$	$7.0 \pm 1.7$	$6.9 \pm 1.2$	$60.6 \pm 11.2$

(\*: could not reach to the maxima of Chernoff distance)

#### 7.4.4 QCMAP

##### UCI database: Optimization

In this section, I discuss the computational cost of the QCMAP estimation in comparison with other methods. I implemented an optimizer based on the primal-dual interior point (PDIP) algorithm for the QCMAP problem in the programming language C. The PDIP algorithm was explained in Section 6.1.7.

I used the PR\_LOQO optimizer for the SVM problem for comparison. The PR\_LOQO is an optimizer for convex quadratic programming, created by A. Smola in 1997 and implemented in the programming language C. It is based on the PDIP algorithm including a path-following method.

In both implementations, individual calculation procedures (e.g., matrix and vector multiplication, and Cholesky decomposition) are almost equivalent. Therefore, the two optimizers operate under nearly equal conditions.

For reference, I also tested the IBM ILOG CPLEX Optimizer [51] for the SVM problem.

I tested their optimizers for each problem using the UCI data sets [35] (see Table 7.4), and measured computational times [sec] for each realization when running on a computer with an Intel Core i7 3.33 GHz CPU and 12 GBytes of RAM. Table 7.10 gives average times  $\pm$  standard deviations for all realizations.

Table 7.9: Classification errors for UCI dataset

label	FDA	aCFDA	BFDA	CFDA- $p$	$\rho_1$ ave.	CFDA-m	$\rho_1$ ave.
1	2.6 ± 0.9	2.4 ± 0.8	2.6 ± 0.9	2.4 ± 0.9	0.35	2.4 ± 0.9	0.32
2	23.7 ± 1.9	23.5 ± 2.0	23.6 ± 2.1	23.5 ± 1.8	0.65	23.5 ± 2.1	0.50
3	31.7 ± 3.9	34.8 ± 3.7	○ 19.4 ± 3.5	○19.4 ± 3.5	0.21	○19.4 ± 3.5	0.90
4	13.1 ± 4.1	<b>8.4 ± 3.0</b>	18.2 ± 4.3	18.2 ± 4.3	0.69	18.2 ± 4.3	0.80
5	26.6 ± 5.0	27.4 ± 4.6	<b>6.6 ± 2.7</b>	<b>6.6 ± 2.7</b>	0.33	<b>6.6 ± 2.7</b>	0.80
6	1.0 ± 1.4	1.9 ± 1.9	1.0 ± 1.4	1.9 ± 1.9	0.28	○ 0.9 ± 1.2	0.70
7	<b>17.2 ± 5.6</b>	21.3 ± 7.3	27.0 ± 14.9	27.0 ± 14.9	0.11	27.0 ± 14.9	0.83
8	3.8 ± 1.2	3.8 ± 1.2	3.8 ± 1.2	3.8 ± 1.2	0.50	3.9 ± 1.4	0.55
9	50.4 ± 0.9	50.3 ± 1.0	○ <b>41.1 ± 0.4</b>	○ <b>41.1 ± 0.4</b>	0.54	○ <b>41.1 ± 0.4</b>	0.50

label	KFDA	Kernel aCFDA	KBFDA	KCFDA- $p$	$\rho_1$ ave.	KCFDA-m	$\rho_1$ ave.
1	○ 2.2 ± 0.7	○ 2.2 ± 0.7	○ 2.2 ± 0.8	○ 2.2 ± 0.7	0.35	○ 2.2 ± 0.8	0.51
2	23.4 ± 1.6	○ 22.9 ± 1.8	23.5 ± 1.8	23.4 ± 2.0	0.65	23.7 ± 1.9	0.52
3	25.7 ± 4.2	28.5 ± 3.4	○ <b>19.4 ± 3.1</b>	○ <b>19.4 ± 3.1</b>	0.21	○ <b>19.4 ± 3.1</b>	0.90
4	6.2 ± 2.4	○ 5.3 ± 2.0	6.2 ± 2.5	5.5 ± 2.0	0.69	6.5 ± 2.3	0.73
5	○ 2.8 ± 2.2	3.1 ± 2.6	3.1 ± 2.4	3.3 ± 2.5	0.33	3.4 ± 2.4	0.66
6	3.7 ± 2.6	3.6 ± 3.1	6.9 ± 3.6	3.7 ± 3.3	0.28	6.6 ± 3.6	0.58
7	18.3 ± 3.7	26.8 ± 5.2	<b>9.5 ± 3.0</b>	20.2 ± 5.8	0.11	○ <b>9.4 ± 3.3</b>	0.89
8	0.7 ± 0.4	0.7 ± 0.3	0.7 ± 0.4	○ 0.6 ± 0.4	0.50	0.8 ± 0.4	0.51
9	44.6 ± 1.0	44.7 ± 1.2	<b>41.9 ± 0.8</b>	<b>42.2 ± 0.9</b>	0.54	<b>41.8 ± 1.0</b>	0.56

### UCI database: Classification

I tested four classifiers provided by individual and combined novel weight function strategies. These classifiers and their parameter settings are listed in Table 7.11, where  $\bar{\mathbf{x}}$  and  $\hat{\Sigma}$  are given by

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n, \quad (7.15)$$

$$\hat{\Sigma} = \begin{cases} \frac{1}{N-1} \sum_{n=1}^N (\boldsymbol{\mu}(i) - \mathbf{x}_n(i))^2 & i = j \\ 0 & i \neq j \end{cases}. \quad (7.16)$$

I also have to select the parameter  $s$  in the KDEQCM method, and if we want to make the constraints strong, then  $s$  must be small. However, if  $s$  is too small in KDEQCM, then the classifier is almost the same as LSR and its results are meaningless. Therefore, I selected  $s = 0.01$  in KDEQCM. Nevertheless, we cannot guarantee these parameter settings are optimal.

Conversely, the kernel parameter  $\gamma$  is optimized as follows:

- Repeat the five-fold cross validation five times using different realizations, and record their optimal parameters:  $\{\gamma^{(1)}, \dots, \gamma^{(5)}\}$ .
- Select a median value from their values:

$$\gamma^* = \text{Median} \left[ \gamma^{(1)}, \dots, \gamma^{(5)} \right]. \quad (7.17)$$

Table 7.10: Computational time of optimizers

Optimizer	PR_LOQO in C for SVM	CPLEX in C for SVM	PDIP in C for QCMAp
Banana	5.93 ± 0.05	0.26 ± 0.10	5.09 ± 1.22
Breast Cancer	0.76 ± 0.01	0.08 ± 0.01	1.13 ± 0.30
Diabetes	9.24 ± 0.06	0.37 ± 0.03	18.60 ± 3.33
Flare-Solar	27.29 ± 0.16	0.09 ± 0.00	93.13 ± 45.12
German	30.45 ± 0.15	0.89 ± 0.02	211.20 ± 39.16
Heart	0.53 ± 0.03	0.04 ± 0.00	0.65 ± 0.29
Image	207.28 ± 1.17	5.86 ± 0.19	346.48 ± 63.91
Ringnorm	6.44 ± 0.03	0.36 ± 0.01	5.67 ± 1.11
Splice	86.51 ± 0.36	2.36 ± 0.04	318.77 ± 87.80
Thyroid	0.33 ± 0.02	0.03 ± 0.00	0.22 ± 0.06
Titanic	0.35 ± 0.01	0.02 ± 0.00	0.30 ± 0.20
Twonorm	5.55 ± 0.03	0.35 ± 0.02	6.83 ± 3.14
Waveform	6.08 ± 0.04	0.30 ± 0.07	7.15 ± 1.44

This cross-validation procedure is taken from [86, 74] and following that study. I selected the kernel parameters from among:

$$\{10^{-6.0}, 10^{-5.75}, 10^{-5.50}, \dots, 10^{0.50}, 10^{0.75}, 10^{1.00}\}.$$

For QCMAp estimation, we have only one hyper-parameter for model selection. This is an important property of the QCMAp method since computation time of cross-validation increases exponentially with the number of hyper-parameters.

The classifiers were trained by using the PDIP algorithm described in Section 6.1.7. In this regard, I set the approximation parameter  $t = 2$ . Although this value does not result in highly accurate approximations of  $\min(r_n, 1)$ , it maintains smooth progression of the algorithm.

The k-nearest neighbor (kNN) method, the regularized AdaBoost ( $AB_R$ ) method and SVM were used for comparison with the proposed classifiers. In the kNN method, I select the optimal value of  $k$  from  $\{1, 3, 5, \dots, 61\}$ . Results of  $AB_R$  and SVM are referenced from [86, 74]. The RBF Gaussian kernel was used in SVM. Regularization parameters of the  $AB_R$  and the tradeoff (regularization) parameter and kernel parameter of SVM were selected in the same way (the five-fold cross-validation method). Table 7.12 shows the results of this experiment.

The values listed in this table show the mean value  $\pm$  standard deviation of the error rates for all realizations, and the minimum error rate for each data set is denoted by a bold font. Symbols ‘●’, ‘○’, and ‘\*’ indicate that the error of QCMAp is significantly smaller than the errors of kNN,  $AB_R$ , and SVM, respectively, by one tailed t-test with significance level of 5%.

Table 7.11: Proposed Classifiers tested in the experiment

Number	Classifier(s)	Parameter Setting (First Classifier)	Parameter Setting (Second Classifier)	Mixing Rate
1	GQCM	$\mu = \bar{x}, \Sigma = \hat{\Sigma}$		
2	DOGQCM	$\mu = \bar{x}, \Sigma = \hat{\Sigma}, \nu = 4, \kappa = 0.9$		
3	KDEQCM	$\mu_n = \mathbf{x}_n, s = 0.01, n = 1, \dots, N$		
4	DOGQCM + KDEQCM	$\mu = \bar{x}, \Sigma = \hat{\Sigma}, \nu = 4, \kappa = 0.9$	$\mu_n = \mathbf{x}_n, s = 0.01, n = 1, \dots, N$	$p = 0.5$

Table 7.12: Experimental results of UCI data sets

	kNN	AB <sub>R</sub>	SVM	GQCM	DOGQCM	KDEQCM	DOG+KDE
Banana	11.3 ± 0.6	10.9 ± 0.4	11.5 ± 0.7	●○* <b>10.5 ± 0.4</b>	●○* <b>10.5 ± 0.4</b>	●○* <b>10.5 ± 0.5</b>	●○* <b>10.5 ± 0.4</b>
Breast Cancer	<b>25.3 ± 4.0</b>	26.5 ± 4.5	26.0 ± 4.5	○ 26.0 ± 4.4	○ 25.4 ± 4.3	26.0 ± 4.5	25.5 ± 4.5
Diabetes	25.4 ± 1.8	23.8 ± 1.8	23.5 ± 1.7	●○ <b>23.2 ± 2.0</b>	●○ <b>23.2 ± 1.7</b>	●○ <b>23.2 ± 2.0</b>	●○ <b>23.2 ± 1.8</b>
Flare-Solar	34.8 ± 1.9	34.2 ± 2.2	<b>32.4 ± 1.8</b>	●○ 33.3 ± 1.8	●○ 33.3 ± 1.7	●○ 33.3 ± 1.8	●○ 33.3 ± 1.7
German	25.2 ± 2.3	24.3 ± 2.1	<b>23.6 ± 2.1</b>	●○ 23.7 ± 2.2	●○ 23.7 ± 2.2	●○ 23.7 ± 2.2	● 23.8 ± 2.3
Heart	15.7 ± 3.4	16.5 ± 3.5	16.0 ± 3.3	○ <b>15.6 ± 3.2</b>	○ <b>15.6 ± 3.3</b>	15.7 ± 3.4	15.7 ± 3.4
Image	3.4 ± 0.5	2.7 ± 0.6	3.0 ± 0.6	3.6 ± 0.6	● 2.7 ± 0.6	3.1 ± 0.7	●* <b>2.6 ± 0.5</b>
Ringnorm	35.0 ± 1.4	<b>1.6 ± 0.1</b>	1.7 ± 0.1	● 1.9 ± 0.1	● 1.8 ± 0.1	● 4.6 ± 0.8	● 2.0 ± 0.4
Splice	26.1 ± 2.2	<b>9.5 ± 0.7</b>	10.9 ± 0.7	● 10.8 ± 0.6	● 10.7 ± 0.7	● 10.9 ± 0.7	● 10.8 ± 0.8
Thyroid	4.4 ± 2.2	4.6 ± 2.2	4.8 ± 2.2	* 4.1 ± 2.7	●○* 3.8 ± 2.1	6.9 ± 3.0	●○* <b>3.7 ± 2.1</b>
Titanic	22.8 ± 1.1	22.6 ± 1.2	22.4 ± 1.0	●○ 22.2 ± 1.0	●○ 22.2 ± 1.0	●○* <b>21.9 ± 1.0</b>	●○ 22.3 ± 1.1
Twonorm	2.6 ± 0.2	2.7 ± 0.2	3.0 ± 0.2	●○* 2.4 ± 0.1	●○* <b>2.3 ± 0.1</b>	●○* 2.4 ± 0.1	●○* <b>2.3 ± 0.1</b>
Waveform	10.7 ± 1.0	<b>9.8 ± 0.8</b>	9.9 ± 0.4	● 10.2 ± 0.6	● <b>9.8 ± 0.5</b>	● 10.1 ± 0.5	● 9.9 ± 0.5

- : Superior to kNN solution by t-test with significance level 5%.
- : Superior to AB<sub>R</sub> solution by t-test with significance level 5%.
- \*: Superior to SVM solution by t-test with significance level 5%.

### Weight function selection

In this section, I try to select a good weight function for a data set by using the results of previous experiments. From the results, I found that similar weight functions performed in similar accuracy, DOGQCM was superior to GQCM for all data sets, and DOG+KDE was partially superior to DOGQCM. Therefore, the number of candidates of weighting functions can be reduced. In this experiment, I employed DOGQCM and DOG+KDE as the candidates, and propose a method to select either of the weight functions. DOG+KDE is suitable for an easy data set and DOGQCM is suitable for a difficult data set because DOG+KDE would be more similar to LSR than DOGQCM. Fig. 7.21 shows the 2d-plot of NN and LDA error rates for 13 UCI data sets. NN error rates are calculated as the rate of samples of which the nearest neighbor belongs to the other class. LDA error rates are calculated by using the same samples for training and testing. From the previous results, DOG+KDE was almost always better for ‘Image’, ‘Thyroid’, and ‘Twonorm’, and DOG+KDE was almost always better for the other data sets. Then I decided the boundary line between the two weight functions by LDA. When another data set is given, we can select either of the weight functions by calculating NN and LDA error rates and separating it with the boundary line. Since the computational costs of NN and LDA are much smaller than kernelized methods, we can select a weight function effectively. I call it the weight function selection QCMAP (WS-QCM).

### Application to Brain Signal

I conducted an experiment by using the challenging brain computer interface (BCI) data set [8]. This is for motor imagery classification to assign EEG signals into either ‘foot’ imagery or ‘right hand’ imagery. It consists of 280 EEG signals for 5 subjects. I extracted features from EEG signals into six dimensional feature vectors by the common spacial pattern (CSP) filter [75]. And I calculated generalization errors by 5-fold cross validation (CV) for each subject. Table 7.13 describes the result of this experiment. I compared the proposed method to SVM with two types of hyper parameter selection approaches. SVM has two hyper parameters: a kernel parameter  $\gamma$  and a tradeoff parameter  $\lambda$ . I selected an optimal combination from grid-like candidates by the 5-fold CV. The number of candidates of  $\gamma$  and  $\lambda$  are 29 and 13, then the number of all combinations is 377 in this method. The first approach is to optimize two parameters from all 377 combinations of grid-like candidates by the 5-fold CV. I denoted this method as “SVM-1” in the table. The other approach consists of two phases. In the first phase, I obtain a combination of values from rough grid-like candidates by the 5-fold CV. In the second phase, I optimize two parameters from fine grid-like candidates around the combination by the 5-fold CV. I denoted this method as “SVM-2” in the table. From the result, SVM-2 (two phase CV) was better than SVM-1 (all search CV) for three subjects. The fact implies that the CV based hyperparameter selection is very difficult.

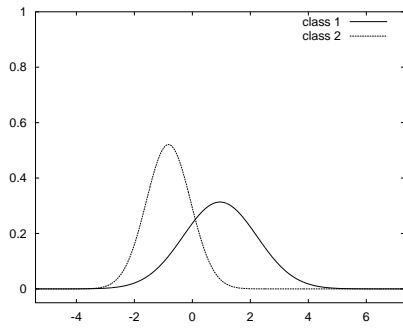
Furthermore, I applied the  $SVM_G$ , the  $SVM_D$ , and the KCFDA-m to this experiment. Hyperparameters of the  $SVM_G$  and the  $SVM_D$  were optimized by the two phase CV. Kernel parameter  $\gamma$  of the KCFDA-m was optimized by the 5-fold CV.

QCMAP approach has one hyper parameter as a kernel parameter  $\gamma$  and two candidates of weight functions. I first selected either of the weight functions by the boundary line in the NN and LDA error rate space. After that, I optimized its kernel parameter  $\gamma$  by the 5-fold CV. I denoted this method as “WS-QCM” in the table. The column of “DOG:D&K” shows numbers of times that DOGQCM and DOG+KDE were selected, respectively. The columns of “only DOG” and “only D&K” show the error rates of individual weight functions. From the results, weight function selection worked well for three out of five subjects. The data of Subjects 2 and 4 are easy for classification, and DOG+KDE provided better performances. The data of Subjects 3 is difficult, and DOGQCM provided a better performance. The data of Subjects 1 and 5 are intermediate, and we could not select the proper weight function. However, the differences between error rates of DOGQCM and DOG+KDE are small for the two subjects. The experimental result shows the proposed selection method can be applied to practical problems. Another notable result is that QCMAP approaches provided better performances than SVM for almost all subjects. Furthermore, I have to recall that one hyper parameter is more easily optimized than two hyper parameters, and we can use QCMAP more easily than SVM.

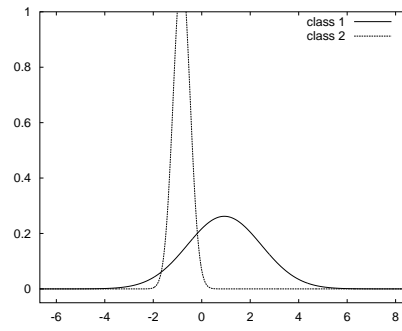
Table 7.13: Experimental results of a BCI data set by the weight function selection

	SVM-1(all search)	SVM-2(two phase)	SVM <sub>G</sub>	SVM <sub>D</sub>	KCFDA-m
Subject 1	28.2 ± 5.8	28.6 ± 5.9	23.6 ± 6.1	24.3 ± 7.4	23.9 ± 4.5
Subject 2	2.1 ± 1.5	2.9 ± 2.0	2.1 ± 1.5	1.4 ± 1.5	1.8 ± 1.3
Subject 3	37.1 ± 10.8	35.7 ± 7.7	36.4 ± 11.1	38.6 ± 6.0	47.9 ± 8.8
Subject 4	8.6 ± 2.0	8.2 ± 3.0	8.2 ± 3.0	8.6 ± 2.9	<b>6.4 ± 3.0</b>
Subject 5	15.0 ± 4.8	13.6 ± 3.9	15.7 ± 5.0	15.7 ± 2.0	13.2 ± 6.3

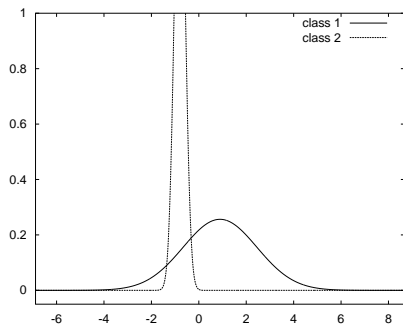
WS-QCM	DOG:D&K	only DOG	only D&K
23.9 ± 4.8	5 : 0	23.9 ± 4.8	<b>23.2 ± 4.6</b>
<b>1.1 ± 1.6</b>	0 : 5	1.4 ± 1.5	<b>1.1 ± 1.6</b>
<b>33.6 ± 6.5</b>	5 : 0	<b>33.6 ± 6.5</b>	35.7 ± 6.1
7.9 ± 2.0	0 : 5	9.3 ± 2.9	7.9 ± 2.0
13.6 ± 5.0	1 : 4	<b>12.9 ± 5.8</b>	13.6 ± 5.0



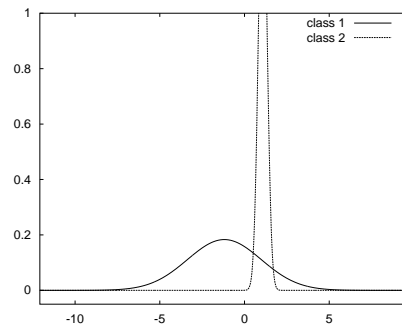
(a) FDA



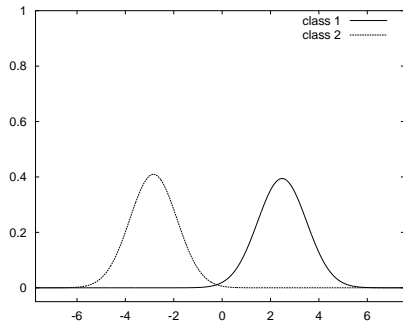
(b) aCFDA



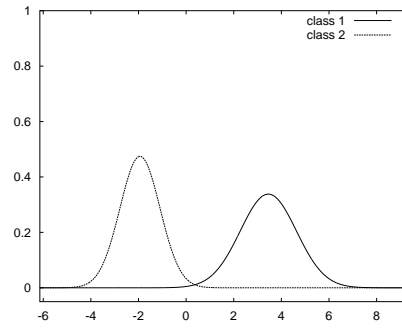
(c) CFDA- $p$



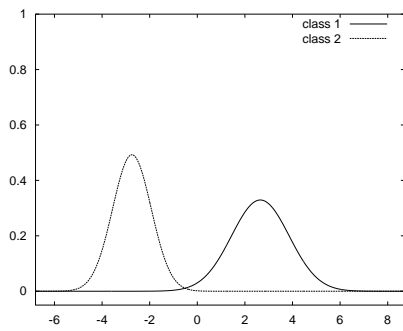
(d) CFDA- $m$



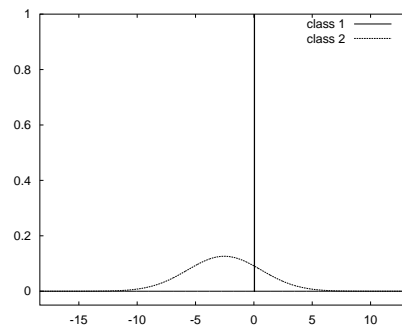
(e) KFDA



(f) Kernel aCFDA



(g) KCFDA- $p$



(h) KCFDA- $m$

Figure 7.20: Distributions in projected spaces ( $N_2/N_1 = 1.5$ )

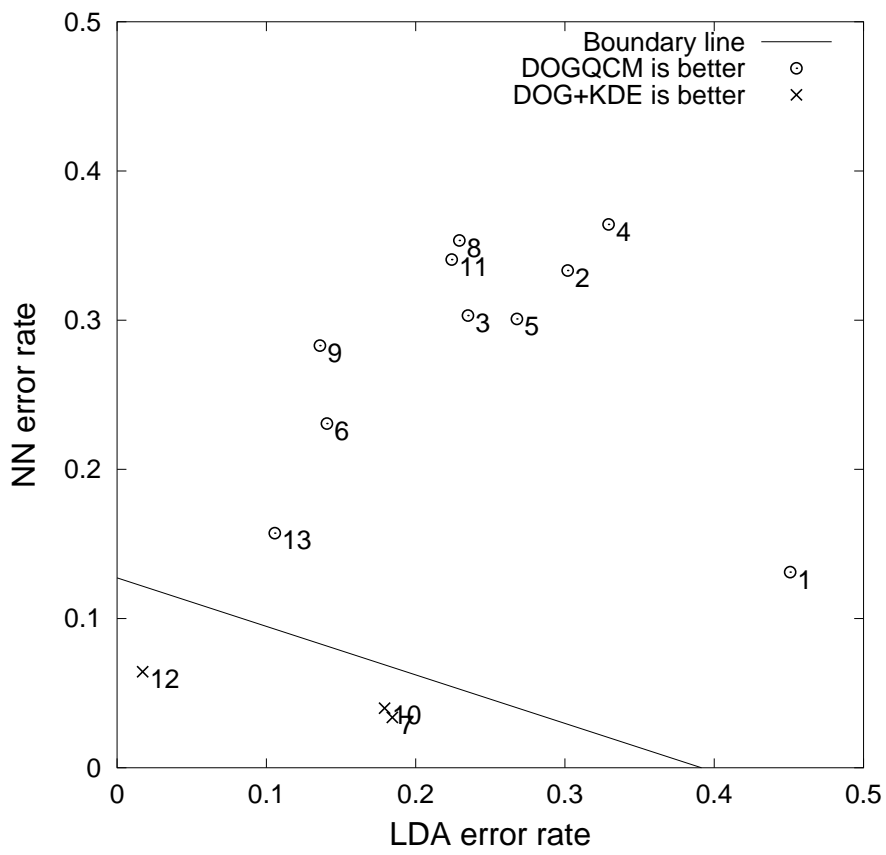


Figure 7.21: Score plot of UCI data sets

# Chapter 8

## Conclusions

### 8.1 Summary

In this thesis, I have explained the matrix and tensor based feature extraction and the classification methods, and proposed several new methods for them.

In Chapter 2, I have introduced several existing methods for data analysis: the discrete Fourier transform (DFT), the discrete wavelet transform (DWT), the principal component analysis (PCA), the sparse principal component analysis (SPCA), and the independent component analysis (ICA).

In Chapter 3, I have introduced matrix and tensor decomposition techniques. The most basic matrix decomposition is the singular value decomposition (SVD). The SVD provides a full-rank complete decomposition which consists of orthonormal left and right singular matrices and a diagonal rectangular matrix. The truncated SVD (tSVD) is a low-rank approximation model of SVD. The tSVD provides the minimum mean squared error between an original data and a decomposition model. The SVD and the tSVD features are obtained based on orthogonality constraint; however, orthogonality is not always efficient. I have to use also sparsity, smoothness and nonnegativity constraints as the situation demands. The penalized matrix factorization (PMF) is a technique to impose a sparsity or a smoothness constraint to the matrix decomposition model, and the nonnegative matrix factorization (NMF) is a technique to impose a nonnegativity constraint to the matrix decomposition model. Furthermore, the sparse NMF, the smooth NMF methods have been introduced.

The tensor decomposition is an extension technique of the matrix decomposition to analyze multi-dimensional array (tensor) data. There are two models in the tensor decomposition: CP model and Tucker model. The CP model is a direct extension of the matrix decomposition, and the Tucker model is more general extension. I have introduced to impose the orthogonality, the sparsity and the nonnegativity constraints to both models.

Furthermore, I have introduced the common and individual feature extraction, and several solution algorithms. The common and individual feature extraction is a very important technique for data analysis; however, this is technically complex.

In Chapter 4, I have introduced a supervised-feature extraction and several basic classification meth-

ods: the common spatial pattern (CSP) filter, the kernel method, the least squares regression (LSR), the Fisher discriminant analysis (FDA), and the support vector machine (SVM).

In Chapter 5, I have proposed the following new algorithms for matrix and tensor based feature extraction:

- Fast algorithm and extension methods for Smooth NMF with function approximation
  - HALS based fast improvement algorithm
  - Extensive basis function
  - Reduction method of the size of basis matrix
  - CP tensor extension
  - Tucker tensor extension
- Linked (common and individual) CP tensor decomposition (LCPTD)
  - Nonconstrained LCPTD
  - Sparse LCPTD
  - Nonnegative LCPTD
- linked (common and individual) Tucker decomposition (LTD)
  - Orthogonal LTD
  - Sparse LTD
  - Nonnegative LTD

In Chapter 6, I have proposed the following new classification methods:

- QCMAP classifiers
  - QCMAP classifier with the Gaussian weight function (GQCM)
  - QCMAP classifier with the difference of the Gaussian weight functions (DoGQCM)
  - QCMAP classifier with the kernel density estimated weight function (KDEQCM)
  - QCMAP classifier with the DoG and the KDE weight functions (Dog+KDE-QCM)
- Support vector machine with weighted regularizations (WRSVM)
  - WRSVM with the Gaussian weight function ( $SVM_G$ )
  - WRSVM with the difference of the Gaussian weight functions ( $SVM_D$ )
- Chernoff Fisher discriminant analysis (CFDA) classifiers and their kernel extensions
  - Bhattacharyya FDA (BFDA)

- Chernoff FDA with prior based parameter (CFDA-p)
- Chernoff FDA with error minimization (CFDA-m)
- their kernel extensions (KBFDA, KCFDA-p, KCFDA-m)

In Chapter 7, I have shown the experimental results of new feature extraction and classification methods. To evaluate proposed feature extraction methods, I conducted experiments for the low-rank approximation, the face reconstruction, the blind source separation (BSS), the denoising, the parts-based representation, and the compression. I showed the advantages of a new fast algorithm for smooth NMF/NTF and new linked tensor decomposition models. To evaluate proposed classification methods, I conducted experiments of classification with UCI datasets and BCI dataset. I showed the advantages of the proposed classifiers.

## 8.2 Open problems

In this section, I mention about open problems for individual methods.

### 8.2.1 Feature Extraction

In almost all cases, matrix/tensor decomposition is used as a low-rank approximation; however, the best rank is usually unknown. For example, it is difficult to know the number of original sources in BSS problem. It should be careful to decide the amount of reduction of the information of data in dimensionality reduction. Another problem is how to apply our new methods to big data. In order to process the big data, it is necessary to reduce computation time and online algorithm is also important. Furthermore, researches to apply the proposed method to various types of brain data is necessary.

### 8.2.2 Classification

Hyper parameter selection methods can be considered as an open problem in classification methods. In recent research, most of efficient classifiers are based on kernel method such as SVM. Furthermore, most of robust approaches are based on regularization techniques. Therefore, most of recent classifiers have at least two hyper parameters: a kernel parameter and a regularization parameter. For two hyper parameters, their parameter domain is a 2-dimensional hyper plane. In this case, there are too many candidates of parameter combinations to select. WRSVM has this problem. CFDA based classifiers have only the kernel parameter; however, they still have a possibility of over-fitting. Furthermore, I have proposed only one-dimensional feature model, then the research for multi-dimensional feature model is necessary. The QCMAP classifiers have only a kernel parameter and it is robust for outliers. How to select kernel parameter is an open problem. It can be considered to apply the multiple kernel learning (MKL) methods to the QCMAP estimation. This will provide parameter-less efficient classifiers. Reduction of computational time (especially, learning time) is sometimes considered less serious; however, it is also very important from the view point of online learning. Difficult and complex

criteria and optimization algorithms are not efficient for online learning. It is necessary to apply the proposed methods to various types of brain data (especially multi-categorical data).

# Bibliography

- [1] T. Ahmed, T. Natarajan, and K.R. Rao. Discrete cosine transform. *Computers, IEEE Transactions on*, 100(1):90–93, 1974.
- [2] G.I. Allen. Regularized tensor factorizations and higher-order principal components analysis. *arXiv preprint arXiv:1202.2476*, 2012.
- [3] G. Baudat and F. Anouar. Generalized discriminant analysis using a kernel approach. *Neural Comput.*, 12(10):2385–2404, October 2000.
- [4] P. Belhumeur, J. Hespanha, and D. Kriegman. Eigenfaces vs. fisherfaces: Recognition using class specific linear projection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 711–720, 1997.
- [5] H.Y. Benson and R.J. Vanderbei. Solving problems with semidefinite and related constraints using interior-point methods for nonlinear programming, 2002.
- [6] J. Besag. Towards bayesian image analysis\*. *Journal of Applied statistics*, 20(5-6):107–119, 1993.
- [7] C.M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [8] B. Blankertz, K.R. Muller, D. Krusienski, G. Schalk, J.R. Wolpaw, A. Schlogl, G. Pfurtscheller, J.R. Millan, M. Schröder, and N. Birbaumer. The BCI competition III: Validating alternative approaches to actual BCI problems. *IEEE Trans. Neural Systems and Rehabilitation Engineering*, 14:153–159, 2006.
- [9] R. Bro. Multi-way analysis in the food industry - models, algorithms, and applications. Technical report, MRI, EPG and EMA, " Proc ICSLP 2000, 1998.
- [10] R. Bro and S.D. Jong. A fast non-negativity-constrained least squares algorithm. *Journal of chemometrics*, 11(5):393–401, 1997.
- [11] J. Carroll and J.J. Chang. Analysis of individual differences in multidimensional scaling via an n-way generalization of " eckart-young " decomposition. *Psychometrika*, 35:283–319, 1970.

- [12] Y.H. Chao, H.M. Wang, and R.C. Chang. GMM-based Bhattacharyya kernel fisher discriminant analysis for speaker recognition. In *Proc. ICASSP*, volume 3, pages 3–6, 2005.
- [13] Z. Chen, A. Cichocki, and T.M. Rutkowski. Constrained non-negative matrix factorization method for EEG analysis in early detection of alzheimer disease. In *Acoustics, Speech and Signal Processing, 2006. ICASSP 2006 Proceedings. 2006 IEEE International Conference on*, volume 5, pages 893–896. IEEE, 2006.
- [14] M.T. Chu and N.T. Trendafilov. The orthogonally constrained regression revisited. *Journal of Computational and Graphical Statistics*, 10(4):pp. 746–771, 2001.
- [15] H. Chun and S. Keleos. Sparse partial least squares regression for simultaneous dimension reduction and variable selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(1):3–25, 2010.
- [16] R.V. Churchill. *Complex Variables and Applications*. McGraw-Hill, New York, 1960.
- [17] A. Cichocki. Tensor decompositions: New concepts for brain data analysis? *Journal of Control, Measurement, and System Integration (SICE)*, 7:507–517, 2011.
- [18] A. Cichocki, S. Amari, et al. *Adaptive blind signal and image processing*. John Wiley Chichester, 2002.
- [19] A. Cichocki and A.H. Phan. Fast local algorithms for large scale nonnegative matrix and tensor factorizations. *IEICE transactions on fundamentals of electronics, communications and computer sciences*, 92(3):708–721, 2009.
- [20] A. Cichocki, S. Cruces, and S. Amari. Generalized alpha-beta divergences and their application to robust nonnegative matrix factorization. *Entropy*, 13(1):134–170, 2011.
- [21] A. Cichocki and R. Zdunek. NMFLAB-NTFLAB for signal and image processing. *Technical Report, Laboratory for Advanced Brain Processing, BSI, RIKEN*, 2006.
- [22] A. Cichocki and R. Zdunek. Regularized alternating least squares algorithms for non-negative matrix/tensor factorization. In *Advances in Neural Networks-ISNN 2007*, pages 793–802. Springer, 2007.
- [23] A. Cichocki, R. Zdunek, and S. Amari. Csiszar 's divergences for non-negative matrix factorization: Family of new algorithms. In *Independent Component Analysis and Blind Signal Separation*, pages 32–39. Springer, 2006.
- [24] A. Cichocki, R. Zdunek, and S. Amari. Hierarchical ALS algorithms for nonnegative matrix and 3d tensor factorization. In *Independent Component Analysis and Signal Separation*, pages 169–176. Springer, 2007.

- [25] A. Cichocki, R. Zdunek, S. Choi, R. Plemmons, and S. Amari. Non-negative tensor factorization using alpha and beta divergences. In *Acoustics, Speech and Signal Processing, 2007. ICASSP 2007. IEEE International Conference on*, volume 3, pages III–1393. IEEE, 2007.
- [26] A. Cichocki, R. Zdunek, A.H. Phan, and S. Amari. *Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source Separation*. Wiley Publishing, 2009.
- [27] J.L. Crowley and A.C. Parker. A representation for shape based on peaks and ridges in the difference of low-pass transform. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, PAMI-6(2):156–170, march 1984.
- [28] L.D. Lathauwer, B.D. Moor, and J. Vandewalle. A multilinear singular value decomposition. *SIAM journal on Matrix Analysis and Applications*, 21(4):1253–1278, 2000.
- [29] L.D. Lathauwer, B.D. Moor, and J. Vandewalle. On the best rank-1 and rank- $(r_1, r_2, \dots, r_n)$  approximation of higher-order tensors. *SIAM Journal on Matrix Analysis and Applications*, 21(4):1324–1342, 2000.
- [30] H.P. Decell Jr. and S.M. Mayekar. Feature combinations and the divergence criterion. *Computers & Mathematics with Applications*, 3(1):71–76, 1977.
- [31] D.L. Donoho and I.M. Johnstone. Adapting to unknown smoothness via wavelet shrinkage. *Journal of the american statistical association*, 90(432):1200–1224, 1995.
- [32] K. Drakakis, S. Rickard, R.D. Frein, and A. Cichocki. Analysis of financial data using non-negative matrix factorization, 2008.
- [33] S. Essid and C. Fevotte. Smooth nonnegative matrix factorization for unsupervised audiovisual document structuring. *Multimedia, IEEE Transactions on*, 15(2):415–425, 2013.
- [34] C. Fevotte, N. Bertin, and J.L. Durrieu. Nonnegative matrix factorization with the itakura-saito divergence: With application to music analysis. *Neural computation*, 21(3):793–830, 2009.
- [35] A. Frank and A. Asuncion. UCI machine learning repository, 2010.
- [36] J. Friedman, T. Hastie, H. Höfling, and R. Tibshirani. Pathwise coordinate optimization. *The Annals of Applied Statistics*, 1(2):302–332, 2007.
- [37] K. Fukunaga. *Introduction to Statistical Pattern Recognition*. John Wiley & Sons, second edition, 1990.
- [38] D. Geman and G. Reynolds. Constrained restoration and the recovery of discontinuities. *IEEE Transactions on pattern analysis and machine intelligence*, 14(3):367–383, 1992.

- [39] N. Gillis. Sparse and unique nonnegative matrix factorization through data preprocessing. *arXiv preprint arXiv:1204.2436*, 2012.
- [40] N. Gillis and F. Glineur. Using under approximations for sparse nonnegative matrix factorization. *Pattern Recognition*, 43(4):1676–1687, 2010.
- [41] C.C. Gonzaga. Path-following methods for linear programming. *SIAM Review*, 34(2):167–224, 1992.
- [42] P.J. Green. Bayesian reconstructions from emission tomography data using a modified em algorithm. *Medical Imaging, IEEE Transactions on*, 9(1):84–93, 1990.
- [43] R.A. Harshman. Foundations of the PARAFAC procedure: Model and conditions for an 'explanatory' multi-mode factor analysis. *UCLA Working Papers in phonetics*, 16:1–84, 1970.
- [44] T. Hastie and R. Tibshirani. Discriminant analysis by Gaussian mixtures. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):pp. 155–176, 1996.
- [45] T. Hebert and R. Leahy. A generalized EM algorithm for 3-D bayesian reconstruction from poisson data using gibbs priors. *Medical Imaging, IEEE Transactions on*, 8(2):194–202, 1989.
- [46] U. Helmke and J.B. Moore. *Optimization and dynamical systems*. Springer, 1994.
- [47] H. Hoefling. A path algorithm for the fused lasso signal approximator. *Journal of Computational and Graphical Statistics*, 19(4):984–1006, 2010.
- [48] P.O. Hoyer. Non-negative sparse coding. In *Neural Networks for Signal Processing, 2002. Proceedings of the 2002 12th IEEE Workshop on*, pages 557–565. IEEE, 2002.
- [49] P.O. Hoyer. Non-negative matrix factorization with sparseness constraints. *The Journal of Machine Learning Research*, 5:1457–1469, 2004.
- [50] A. Hyvärinen, J. Karhunen, and E. Oja. *Independent Component Analysis*. Wiley, 2001.
- [51] IBM ILOG CPLEX Optimizer. <http://www.ibm.com/software/integration/optimization/cplex-optimizer/>.
- [52] I.T. Jolliffe. *Principal component analysis*, volume 487. Springer-Verlag New York, 1986.
- [53] I.T. Jolliffe, N.T. Trendafilov, and M. Uddin. A modified principal component technique based on the LASSO. *Journal of Computational and Graphical Statistics*, 12(3):531–547, 2003.
- [54] A. Kapteyn, H. Neudecker, and T. Wansbeek. An approach to n-mode components analysis. *Psychometrika*, 51(2):269–275, 1986.
- [55] H. Kim and H. Park. Sparse non-negative matrix factorizations via alternating non-negativity-constrained least squares for microarray data analysis. *Bioinformatics*, 23(12):1495–1502, 2007.

- [56] H. Kim and H. Park. Nonnegative matrix factorization based on alternating nonnegativity constrained least squares and active set method. *SIAM Journal on Matrix Analysis and Applications*, 30(2):713–730, 2008.
- [57] J. Kim and H. Park. Sparse nonnegative matrix factorization for clustering. *Technical report, Georgia Institute of Technology*, 2008.
- [58] Y.D. Kim and S. Choi. Nonnegative Tucker decomposition. In *Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on*, pages 1–8. IEEE, 2007.
- [59] T.G. Kolda and B.W. Bader. Tensor decompositions and applications. *SIAM REVIEW*, 51(3):455–500, 2009.
- [60] P.M. Kroonenberg and J.D. Leeuw. Principal component analysis of three-mode data by means of alternating least squares algorithms. *Psychometrika*, 45(1):69–97, 1980.
- [61] W.H. Lawton and E.A. Sylvestre. Self modeling curve resolution. *Technometrics*, 13(3):617–633, 1971.
- [62] D.D. Lee and H.S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):789, 1999.
- [63] H. Lee and S. Choi. Group nonnegative matrix factorization for eeg classification. *Journal of Machine Learning Research - Proceedings Track*, 5:320–327, 2009.
- [64] S. Lemm, B. Blankertz, G. Curio, and K.R. Muller. Spatio-spectral filters for improving the classification of single trial EEG. *Biomedical Engineering, IEEE Transactions on*, 52(9):1541–1548, sept. 2005.
- [65] F. Liese and K.J. Miescke. *Statistical Decesion Theory: Estimation, Testing, and Selection*. Statistics. Springer, 2008.
- [66] E.F. Lock, K.A. Hoadley, J.S. Marron, and A.B. Nobel. Joint and individual variation explained (JIVE) for integrated analysis of multiple datatypes. *The Annals of Applied Statistics*, 7(1):523–542, 2013.
- [67] M. Loog and R.P.W. Duin. Linear dimensionality reduction via a heteroscedastic extension of LDA: The chernoff criterion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26:732–739, 2004.
- [68] D.G. Luenberger. *Optimization by vector space methods*. John Wiley and sons, 1969.
- [69] M.S. Mahanta, A.S. Aghaei, K.N. Plataniotis, and S. Pasupathy. Heteroscedastic linear feature extraction based on sufficiency conditions. *Pattern Recognition*, 45(2):821 – 830, 2012.

- [70] M.S. Mahanta and K.N. Plataniotis. A heteroscedastic extension of LDA based on multi-class matusita affinity. In *Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on*, pages 1921–1924, 2012.
- [71] N. Metropolis and S. Ulam. The monte carlo method. *Journal of the American Statistical Association*, 44(247):335–341, 1949.
- [72] S. Mika, G. Ratsch, J. Weston, B. Scholkopf, and K.R. Muller. Fisher discriminant analysis with kernels. In *Neural Networks for Signal Processing IX, 1999. Proceedings of the 1999 IEEE Signal Processing Society Workshop*, pages 41 –48, August 1999.
- [73] V. Moraru. An algorithm for solving quadratic programming problems. *Computer Science Journal of Moldova*, 5(2):223–235, 1997.
- [74] K.R. Muller, S. Mika, G. Ratsch, K. Tsuda, and B. Scholkopf. An introduction to kernel-based learning algorithms. *Neural Networks, IEEE Transactions on*, 12(2):181–201, 2001.
- [75] J. Muller-Gerking, G. Pfurtscheller, and H. Flyvbjerg. Designing optimal spatial filters for single-trial EEG classification in a movement task. *Clinical Neurophysiology*, 110(5):787 – 798, 1999.
- [76] C.H. Nguyen and H. Mamitsuka. Latent feature kernels for link prediction on sparse graphs. *Neural Networks and Learning Systems, IEEE Transactions on*, 23(11):1793–1804, 2012.
- [77] Q. Novi, C. Guan, T.H. Dat, and P. Xue. Sub-band common spatial pattern (SBCSP) for brain-computer interface. In *Neural Engineering, 2007. CNE '07. 3rd International IEEE/EMBS Conference on*, pages 204 –207, may 2007.
- [78] P. Paatero and U. Tapper. Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*, 5(2):111–126, 1994.
- [79] E. Parkhomenko. *Sparse canonical correlation analysis*. PhD thesis, University of Toronto, 2008.
- [80] A. Pascual-Montano, J.M. Carazo, K. Kochi, D. Lehmann, and R.D. Pascual-Marqui. Nonsmooth nonnegative matrix factorization (nsnmf). *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 28(3):403–415, 2006.
- [81] X. Peng and Y. Wang. Geometric algorithms to large margin classifier based on affine hulls. *Neural Networks and Learning Systems, IEEE Transactions on*, 23(2):236–246, 2012.
- [82] A.H. Phan and A. Cichocki. Tensor decompositions for feature extraction and classification of high dimensional datasets. *IEICE, NOLTA*, 1(1):37–68, 2010.

- [83] A.H. Phan and A. Cichocki. Extended HALS algorithm for nonnegative Tucker decomposition and its applications for multiway analysis and classification. *Neurocomputing*, 74(11):1956 – 1969, 2011.
- [84] A.H. Phan and A. Cichocki. Fast and efficient algorithms for nonnegative tucker decomposition. In *Advances in Neural Networks - ISNN 2008*, volume 5264 of *Lecture Notes in Computer Science*, pages 772–782. Springer Berlin Heidelberg, 2008.
- [85] A.H. Phan and A. Cichocki. Local learning rules for nonnegative tucker decomposition. In *Neural Information Processing*, volume 5863 of *Lecture Notes in Computer Science*, pages 538–545. Springer Berlin Heidelberg, 2009.
- [86] G. Rätsch, T. Onoda, and K.R. Müller. Soft margins for adaboost. *Tech. Rep. NC-TR-1998-021, Royal Holloway College, University of London, UK*, 42(3):287–320, 1998.
- [87] J.D.M. Rennie. Maximum-margin logistic regression. <http://people.csail.mit.edu/jrennie/writing>, February 2005.
- [88] F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6):386–408, Nov. 1958.
- [89] L. Rueda and M. Herrera. Linear dimensionality reduction by maximizing the chernoff distance in the transformed space. *Pattern Recognition*, 41(10):3138–3152, 2008.
- [90] H. Sakano, T. Ohashi, A. Kimura, H. Sawada, and K. Ishiguro. Extended Fisher criterion based on auto-correlation matrix information. In *SSPR/SPR'12*, pages 409–416, 2012.
- [91] A. Sierra. High-order Fisher's discriminant analysis. *Pattern Recognition*, 2002.
- [92] M. Swanson. *Path Integrals and Quantum Processes*. Academic Press, CA, 1992.
- [93] D. Tao, X. Li, X. Wu, and S.J. Maybank. General averaged divergence analysis. In *Data Mining, 2007. ICDM 2007. Seventh IEEE International Conference on*, pages 302–311, 2007.
- [94] F.J. Theis, K. Stadlthanner, and T. Tanaka. First results on uniqueness of sparse non-negative matrix factorization. In *Proceedings of the 13th European Signal Processing Conference (EU-SIPCO '05)*, 2005.
- [95] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.
- [96] R. Tibshirani and P. Wang. Spatial smoothing and hot spot detection for CGH data using the fused lasso. *Biostatistics*, 9(1):18–29, 2008.
- [97] A.N. Tikhonov and V.Y. Arsenin. *Solutions of Ill-posed Problems*. Winston & Sons, 1977.

- [98] L.R. Tucker. Implications of factor analysis of three-way matrices for measurement of change. In C. W. Harris, editor, *Problems in measuring change.*, pages 122–137. University of Wisconsin Press, Madison WI, 1963.
- [99] L.R. Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3):279–311, 1966.
- [100] M.H.V. Benthem and M.R. Keenan. Fast algorithm for the solution of large-scale non-negativity-constrained least squares problems. *Journal of chemometrics*, 18(10):441–450, 2004.
- [101] V. Vapnik and A. Chervonenkis. A note on one class of perceptrons. *Automation and Remote Control*, 25, 1964.
- [102] V. Vapnik and A. Lerner. Pattern recognition using generalized portrait method. *Automation and Remote Control*, 24, 1963.
- [103] V. Vapnik. *Statistical Learning Theory*. Wiley, New York, 1998.
- [104] S. Watanabe and N. Pakvasa. Subspace method of pattern recognition. In *Proceedings of 1st International Joint Conference of Pattern Recognition*, pages 25–32, 1973.
- [105] D.M. Witten, R. Tibshirani, and T. Hastie. A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*, 10(3):515–534, 2009.
- [106] W. Xu, X. Liu, and Y. Gong. Document clustering based on non-negative matrix factorization. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*, pages 267–273. ACM, 2003.
- [107] T. Yokota, A. Cichocki, and Y. Yamashita. Linked PARAFAC/CP tensor decomposition and its fast implementation for multi-block tensor analysis. In *Neural Information Processing*, volume 7665, pages 84–91. Springer, 2012.
- [108] T. Yokota, T. Wakahara, and Y. Yamashita. Heteroscedastic Gaussian based correction term for Fisher discriminant analysis and its kernel extension. In *Proceedings of IJCNN2013*, pages 662–669. IEEE, 2013.
- [109] T. Yokota and Y. Yamashita. Quadratically constrained maximum a posteriori estimation for binary classifier. In *Machine Learning and Data Mining in Pattern Recognition*, volume 6871, pages 1–15. Springer, 2011.
- [110] T. Yokota and Y. Yamashita. Support vector machines with weighted regularization. In *Neural Information Processing*, volume 7063, pages 471–480. Springer, 2011.

- [111] T. Yokota and Y. Yamashita. A quadratically constrained MAP classifier using the mixture of Gaussians models as a weight function. *IEEE Transactions on Neural Networks and Learning Systems*, 24(7):1127–1140, 2013.
- [112] T. Yokota, T. Wakahara, H. Sakano, and Y. Yamashita. Correction term of Fisher discriminant analysis based on normal distribution. *IEICE Tech. Rep, IBISML2012-185*, 112(495):31–36, 2012.
- [113] T. Yokota and Y. Yamashita. A new quadratically constrained criterion based on maximum a posteriori estimation for binary classifier. *IEICE Tech. Rep, IBISML2010-85*, 110(265):187–194, 2010.
- [114] R. Zdunek. Regularized active set least squares algorithm for nonnegative matrix factorization in application to raman spectra separation. In *Advances in Computational Intelligence*, pages 492–499. Springer, 2011.
- [115] R. Zdunek. Approximation of feature vectors in nonnegative matrix factorization with gaussian radial basis functions. In *Proceedings of the 19th international conference on Neural Information Processing - Volume Part I, ICONIP'12*, pages 616–623, Berlin, Heidelberg, 2012. Springer-Verlag.
- [116] R. Zdunek and A. Cichocki. Gibbs regularized nonnegative matrix factorization for blind separation of locally smooth signals. In *15th IEEE International Workshop on Nonlinear Dynamics of Electronic Systems (NDES 2007), Tokushima, Japan*, pages 317–320, 2007.
- [117] R. Zdunek and A. Cichocki. Blind image separation using nonnegative matrix factorization with gibbs smoothing. In *Neural information processing*, pages 519–528. Springer, 2008.
- [118] S. Zhang and T. Sim. Discriminant subspace analysis: a Fukunaga-Koontz approach. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 29(10):1732–1745, 2007.
- [119] M. Zhu and A.M. Martinez. Subclass discriminant analysis. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 28(8):1274–1286, aug. 2006.
- [120] H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005.
- [121] H. Zou, T. Hastie, and R. Tibshirani. Sparse principal component analysis. *Journal of computational and graphical statistics*, 15(2):265–286, 2006.
- [122] M. Kojima, S. Mizuno, T. Tsutiya, and H. Yabe. *Interior point methods*. Asakura Shoten, 2001. (in Japanese)
- [123] S. Akaho. *Multivariate kernel analysis -new approach for nonlinear data analysis-*. Iwanami Shoten, 2008. (in Japanese)