

Coherent Pursuit and Boosting Learning

A DISSERTATION
SUBMITTED TO THE FACULTY OF THE GRADUATE SCHOOL
OF THE UNIVERSITY OF MINNESOTA
BY

Qi Yan

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE DEGREE OF
DOCTOR OF PHILOSOPHY

Professor Xiaotong Shen, Adviser

April 2015

ACKNOWLEDGEMENTS

I would like to thank to my Ph.D advisor, Professor Xiaotong Shen, for supporting me during these past few years. Xiaotong is the best advisor I have met. He is very supportive and is my primary resource for getting my science questions answered and was crucial in helping me crank out this thesis. Every time, when I felt frustrated in research, he encouraged me, provided insightful discussions about the research and inspired me to think in different angles. I am also thankful to him for training my scientific writing skills patiently. All the great advice and discussions which I had been benefited from in the Ph.D program will benefit for sure in the rest of my life.

I also have to thank the members of my Ph.D committee, Professors Galin Jones, Professor Charles Geyer from School of Statistics, University of Minnesota at Twin cities and Professor Wei Pan from Biostatistics, University of Minnesota at Twin cities for serving in my defense committee, spending time reviewing my thesis and providing instrumental suggestions and comments to my research.

I will forever be thankful to all faculty members and staff in School of Statistics, University of Minnesota at Twin cities. They are very friendly and helpful. I enjoyed the conversions and friendship with them.

Finally, I would like to thank my parents, my husband and my daughter for their unconditional love and support to make my life meaningful.

DEDICATION

I dedicate this thesis to my family, my husband, Gang and my daughter, Ella for their constant support. I love you all dearly.

ABSTRACT

In multi-response regression, pursuit of two different types of structures is essential to battle the curse of dimensionality. In this thesis, we seek a sparsest decomposition representation of a parameter matrix in terms of a sum of sparse and low rank matrices, among many overcomplete decompositions. On this basis, we propose a constrained method subject to two nonconvex constraints, respectively for sparseness and low-rank properties. Computationally, obtaining an exact global optimizer is rather challenging. To overcome the difficulty, we use an alternating directions method solving a low-rank subproblem and a sparseness subproblem alternatively, where we derive an exact solution to the low-rank subproblem, as well as an exact solution in a special case and an approximated solution generally through a surrogate of the L_0 -constraint and difference convex programming, for the sparse subproblem. Theoretically, we establish convergence rates of a global minimizer in the Hellinger-distance, providing an insight into why pursuit of two different types of decomposed structures is expected to deliver higher estimation accuracy than its counterparts based on either sparseness alone or low-rank approximation alone. Numerical examples are given to illustrate these aspects, in addition to an application to facial image recognition and multiple time series analysis.

In regression analysis, variables can often be combined into groups based on prior knowledge, such as genomic data, which can be naturally divided into biologically meaningful groups. Luan and Li (2008) and Yin et al (2012) utilize the group structure and propose a block coordinate descent procedure for group additive regression models and nonparametric additive models. Their simulation results demonstrate the good performance of the proposed algorithms in terms of support recovery and predic-

tion accuracy. However, none of them investigate the asymptotic properties of their methods. In this thesis, we generalize a smoothing spline based group L_2 Boosting algorithm and study the theoretical property for estimation of high-dimensional additive models with group variables.

Contents

List of Tables	vii
List of Figures	viii
1 Simultaneous pursuit of sparseness and rank structures for matrix decomposition	1
1.1 Introduction	1
1.2 Proposed method	5
1.2.1 Structure decomposition	5
1.2.2 Estimation	6
1.2.3 Method for nonconvex minimization	7
1.2.4 Computational properties	12
1.3 Theory	13
1.4 Numerical examples	16
1.4.1 Simulation I: Operating characteristics	17
1.4.2 Simulation II: Comparison	21
1.4.3 AR Face Database 20pt Markup	24
1.4.4 Greek Letters Image Reconstruction	27
1.4.5 US Macroeconomic Time Series	28
1.5 Appendix	32

2	Boosting for High-Dimensional Additive Models with Group Variables	45
2.1	Introduction	45
2.2	Models	47
2.2.1	Generalized Additive Models	47
2.2.2	GAM with Group Variables	48
2.3	Methods	48
2.3.1	G-GDB	49
2.3.2	GroupSpAM	50
2.3.3	The proposed algorithm	51
2.4	Consistency of Boosting	53
	References	66

List of Tables

1.1	Results of Simulation I. Algorithm 2 is used for computation. . . .	18
1.2	Results for Simulation I with fixed $k = 5$. Algorithm 2 is used for computation.	20
1.3	Results for Simulation II when $.1p$ nonzero are randomly chosen . Algorithm 1 is used for computation.	22
1.4	Results for Simulation II when $.3p$ nonzero are randomly chosen. Algorithm 1 is used for computation.	23
1.5	Economic indicators collected for U.S. macroeconomic time series. . .	29
1.6	Prediction errors of U.S. macroeconomic data for $K = 11$. Here “Low rank alone”, “Sparsity alone” and ”Ours” indicate our method for low rank pursuit only, for sparsity pursuit only and for simultaneous pursuit of low rank and sparsity. Algorithm 2 is used for computation.	30

List of Figures

1.1	The converted AR face image with markup points.	24
1.2	Extracted sparsity (first), low-rank (second) structures as well as the reconstructed image by the proposed method for AR face images; where the tuning parameters are set to $s_1 = 2500$, $s_2 = 5$	25
1.3	Extracted sparsity (first), low-rank (second) structures as well as the reconstructed image by the proposed method for AR face images; where the tuning parameters are set to $s_1 = 2100$, $s_2 = 10$	25
1.4	Original image (left) versus its noisy version (right).	27
1.5	Reconstructed images based on sparsity alone (first), low-rank alone (second) and our method (third). Algorithm 2 is used for computation.	28
1.6	Q-Q plots for each-fold in U.S. macroeconomic time series data example, where points on a straight line indicates non-departure from normality.	31

Chapter 1

Simultaneous pursuit of sparseness and rank structures for matrix decomposition

1.1 Introduction

In multivariate analysis, data as well as parameters are usually expressed in terms of a matrix form, as opposed to a vector representation in univariate analysis. This occurs frequently in multi-class classification (Amit et al., 2007), matrix completion (Cai et al., 2010; Jain et al., 2010), collaborative filtering (Srebro et al., 2005), computer vision (Wright, 2009), among others. In situations as such, it is essential to identify and employ certain lower-dimensional structures to battle the curse of dimensionality due to an increase in dimensionality from multivariate attributes. In this article, we explore rank and sparseness structures through matrix decomposition simultaneously in estimating large matrices through a novel notation of seeking a sparsest decomposition from a class of overcomplete decompositions.

Statistically, different structures have dramatically different interpretations. A low rank property of a matrix describes global information across different tasks, whereas sparseness concerns local information of specific task. For instance, for face

images, the global information corresponds to the overall shape of a face, but the local information characterizes specific facial expression such as laugh and cry. In linear time-invariant (LTI) system, a low rank property corresponds to a low-order LTI system and a sparseness property captures an LTI system with a sparse impulse response (Porat, 1997). In a high-dimensional situation, betting on one type of structure may not be adequate to battle the curse of dimensionality. In this article, we seek a sparsest decomposition for the purpose of dimension reduction, from a class of overcomplete decompositions into simpler sparse and low-rank components. Specifically, a matrix Θ is decomposed as $\Theta_1 + \Theta_2$, for a sparse Θ_1 and low-rank Θ_2 components, where Θ_1 and Θ_2 are chosen from many such decompositions, with a smallest effective degrees of freedom, leading to high accuracy of parameter estimation. Our objective is to reconstruct the parameter matrix by identifying a sparsest decomposition consisting of simpler components. Such a decomposition can be used to provide a simpler and more efficient description of a complex system in terms of its simpler components. This results in more efficient structure representations leading to higher accuracy of parameter estimation in high-dimensional data analysis.

In this dissertation, we consider a multi-response linear regression problem in which a random sample $(\mathbf{a}_i, \mathbf{z}_i)_{i=1}^n$ is observed with a k -dimensional response vector \mathbf{z}_i following

$$\mathbf{z}_i = \mathbf{a}_i^T \Theta + \epsilon_i, \quad E\epsilon_i = 0, \quad Cov(\epsilon_i) = \sigma^2 \mathbf{I}; \quad i = 1, \dots, n, \quad (1.1)$$

where \mathbf{a}_i is a p -dimensional design vector, is independent of random error ϵ_i , and \mathbf{I} is the identity matrix. Model (1.1) reduces to the univariate case when $k = 1$, and becomes a multivariate autoregressive model when $\mathbf{a}_i = \mathbf{z}_{i-1}$. Through matrix decomposition, we decompose a $p \times k$ regression parameter matrix Θ into a sum of a sparse matrix Θ_1 and a low rank matrix Θ_2 for structure exploration, that is,

$\Theta = \Theta_1 + \Theta_2$. Model (1.1) is expressible in a matrix form

$$\mathbf{Z} = \mathbf{A}\Theta + \mathbf{e}; \tag{1.2}$$

where $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)^T \in \mathbb{R}^{n \times k}$, $\mathbf{A} = (a_1, \dots, a_n)^T$ is a $n \times p$ matrix, and $\mathbf{e} = (\epsilon_1, \dots, \epsilon_n)^T \in \mathbb{R}^{n \times k}$ are the data, design and error matrices. In (1.1), we estimate Θ based on n paired observation vectors $(\mathbf{a}_i, \mathbf{z}_i)_{i=1}^n$, with prior knowledge that Θ_1 is sparse in the number of its nonzero entries, and $\text{rank } r(\Theta_2)$ is low relative to $\min(n, k, p)$. Our goal is to recover the parameter Θ by identifying Θ_1 and Θ_2 .

In the literature, the simultaneous exploration of rank and sparseness structures through matrix decomposition has received some attention, yet has not been well-studied. For robust principal component analysis (RPCA) where $\mathbf{A} = \mathbf{I}_{n \times p}$ is the $n \times p$ identity matrix with its diagonals and off-diagonals being one and zero, Yuan & Yang (2013) and Chandrasekaran et al. (2011) employed a linear combination of the L_1 sparsity regularization and the nuclear-norm regularization, and Zhou & Tao (2011) used a randomized projections based low rank approximations and thresholding for sparsity pursuit. Moreover, Wright et al. (2013) recovers the sparse and low-rank components by minimizing a linear combination of the L_1 -norm for sparsity and the nuclear-norm for low rank pursuit, while Waters et al. (2011) develops a greedy algorithm to pursue the sparse and low rank structures. For multiple task learning, Chen et al. (2010) studies sparse and low rank structures separately through convex regularization. In essence, most the existing literature focuses exclusively on a unique matrix decomposition of Θ with $\mathbf{A} = \mathbf{I}_{n \times p}$ or \mathbf{A} to be a set of random linear measurements, and without noise or with small noise that is essentially ignorable. For instance, Chandrasekaran et al. (2011) provided sufficient conditions for exact recovery of a convex relaxation method without noise; Wright et al. (2013) proved that recovering a target matrix is possible from a small set of randomly selected linear

measurements when the number of measurements is sufficiently large. Among these, Agarwal et al. (2012) considered a general \mathbf{A} and derived a theorem that bounds the Frobenius-norm error obtained through regularized convex relaxation under a "spikiness" condition that the max-norm of the low rank component $\|\Theta_2\|_{\max}$ is less than $\frac{\alpha}{\sqrt{pk}}$ for some fixed $\alpha > 0$.

In this dissertation, we consider a general design matrix \mathbf{A} and parameter matrices (Θ_1, Θ_2) , for regression analysis, where \mathbf{A} represents features of observations which is deterministic, and can be any matrix with n rows and p columns. Of particular interest is reconstruction of Θ in a high-dimensional situation in which (p, k) may exceed the sample size n . Computationally, we use an alternating direction method separating low-rank pursuit from sparsity pursuit alternatively, where an exact solution to the low-rank problem and that to the sparsity pursuit problem when $\mathbf{A} = \mathbf{I}_{n \times p}$ or an approximated solution for a general \mathbf{A} is obtained. In either case, the final solution is shown to be stationary without and with maximum block improvement (Chen et al., 2012) for $\mathbf{A} = \mathbf{I}_{n \times p}$ and a general \mathbf{A} . Theoretically, we establish error bound for the proposed method in the Hellinger-distance for reconstruction of Θ , based on which rates of convergence are obtained. Numerically, the proposed method compares favorably against two strong competitors in simulations.

This chapter is organized as follows. Section 2 develops a computational method through the alternating directions method and a closed-form solution for a rank problem. Section 3 investigates statistical properties of the proposed method, followed by simulation studies and a real data example in Section 4. Finally, technical proofs are contained in Section 5.

1.2 Proposed method

In this section, we explore a structure decomposition of a parameter matrix in the form $\Theta = \Theta_1 + \Theta_2$ under model (1.1), then develops computational methods in two situations and discuss their properties.

1.2.1 Structure decomposition

Due to non-uniqueness of such a decomposition under model (1.1), we seek one decomposition, among many overcomplete decompositions, that minimizes the effective degrees of freedom of Θ Efron (2004), defined as

$$\text{Eff}(\Theta) = \min_{\{\Theta = \Theta_1 + \Theta_2 : \|\Theta_1\|_0 \leq \max(0, p+k-2r(\Theta_2)-2)\}} \|\Theta_1\|_0 + (p+k-r(\Theta_2))r(\Theta_2),$$

where $\|\cdot\|_0$ is the L_0 -norm of a matrix, or the number of nonzero entries of the matrix, and $r(\cdot)$ denotes the rank of a matrix. In other words, we identify a decomposition minimizing the effective degrees of freedom $\text{Eff}(\Theta)$, among all candidate decompositions. Lemma 1.1 below says that the minimal of $\text{Eff}(\Theta)$ is unique in $(\|\Theta_1\|_0, r(\Theta_2))$ under the constraint that $\|\Theta_1\|_0 \leq \max(0, p+k-2r(\Theta_2)-2) \leq 2\max(p, k)$.

Lemma 1.1 *The minimizer of $\text{Eff}(\Theta)$ is unique with respect to $(\|\Theta_1\|_0, r(\Theta_2))$ if $\|\Theta_1\|_0 \leq \max(0, p+k-2r(\Theta_2)-2)$. Moreover,*

$$\text{Eff}(\Theta) \leq \min((p+k-r(\Theta))r(\Theta), \|\Theta\|_0).$$

Model (1.1) is identifiable with respect to Θ but may not be so in (Θ_1, Θ_2) even when \mathbf{A} is of full rank, due to non-uniqueness of a decomposition $\Theta = \Theta_1 + \Theta_2$.

1.2.2 Estimation

To pursue structures of low-rank and sparsity through matrix decomposition simultaneously, we propose a constrained likelihood method subject to two nonconvex constraints:

$$\min_{\Theta_1, \Theta_2} \|\mathbf{A}\Theta_1 + \mathbf{A}\Theta_2 - \mathbf{Z}\|_F^2, \quad \text{subject to} \quad \|\Theta_1\|_0 \leq s_1, \quad r(\Theta_2) \leq s_2, \quad (1.3)$$

where $\|\cdot\|_F$ is the Frobenius-norm defined as the L_2 -norm of all entries of a matrix, and s_1 and s_2 are integer-valued tuning parameters with $0 \leq s_1 \leq \max(p, k)$ and $1 \leq s_2 \leq \min(n, k, p)$ based on the consideration that the rank function and the sparsity measure are integer-valued.

When $\mathbf{A} = \mathbf{I}_{n \times p}$, (1.3) is simplified as

$$\min_{\Theta_1, \Theta_2} \|\mathbf{Z} - \Theta_1 - \Theta_2\|_F^2 \quad \text{subject to} \quad \|\Theta_1\|_0 \leq s_1, \quad r(\Theta_2) \leq s_2, \quad (1.4)$$

where a special structure may be taken into account to solve this nonconvex minimization.

When $\mathbf{A} \neq \mathbf{I}_{n \times p}$ is any matrix of full rank, the two constraints in (1.3) are either defined by the L_0 -function or the rank function, imposing computational challenges. To develop an efficient algorithm to solve (1.3), we approximate the $\|\Theta_1\|_0 = \sum_{i,j} I(|\theta_{ij}| \neq 0)$ by its computational surrogate—the truncated L_1 -function

$$\sum_{\theta_{ij} \in \Theta_1} \frac{1}{\tau} \min(|\theta_{ij}|, \tau)$$

(Shen et al., 2012) as $\tau \rightarrow 0^+$. This leads to a computational surrogate of (1.3):

$$\min_{\Theta_1, \Theta_2} f(\Theta_1, \Theta_2), \quad \text{subject to} \quad \frac{1}{\tau} \sum_{i,j} \min(|\theta_{ij}|, \tau) \leq s_1, \quad r(\Theta_2) \leq s_2, \quad (1.5)$$

where $f(\Theta_1, \Theta_2) = \|\mathbf{A}(\Theta_1 + \Theta_2) - \mathbf{Z}\|_F^2$ and τ is a nonnegative tuning parameter.

1.2.3 Method for nonconvex minimization

This section will develop computational strategies for (1.4) and (1.5) separately, based on blockwise coordinate decent as well as maximum block improvement (MBI, (Chen et al., 2012)). First, we separate the task of sparsity pursuit for Θ_1 from that of rank minimization for Θ_2 , where Θ_1 and Θ_2 correspond to two blocks for decent. Second, we apply MBI to assure that blockwise coordinate decent yields a stationary solution for nonconvex minimization, which would be otherwise impossible. In addition, for (1.5), we develop a gradient project method to permit fast computation of a constrained problem through the means of unconstrained optimization.

The strategy of blockwise coordinate decent proceeds as follows. For (1.4) and (1.5), we solve it in Θ_2 given Θ_1 and solve them in Θ_1 given Θ_2 , alternatively. In each step of alternating blocks, we proceed with the block giving the maximum block improvement.

Nonconvex minimization (1.4): a special case

For (1.4), when Θ_2 is held fixed, (1.4) has a global minimizer can be obtained through componentwise thresholding defined by the L_0 -function as follows:

$$\hat{\Theta}_1(\mathbf{Z}, \Theta_2) = \left(I \left\{ |z_{ij} - \theta_{ij}^{(2)}| > \lambda \right\} \cdot (z_{ij} - \theta_{ij}^{(2)}) \right)_{p \times k}, \quad (1.6)$$

where $\theta_{ij}^{(2)}$ is the ij th entry of Θ_2 and λ is any number between the s_1 th and $(s_1 + 1)$ th largest entries of $|\mathbf{Z} - \Theta_2|$.

When Θ_1 is held fixed, a global minimizer of (1.4) is

$$\hat{\Theta}_2(\mathbf{Z}, \Theta_1) = \mathbf{U} \mathbf{D}_{s_2} \mathbf{V}^T, \quad (1.7)$$

where \mathbf{U} and \mathbf{V} are given by singular value decomposition (SVD) of $\mathbf{Z} - \mathbf{\Theta}_1 = \mathbf{U}\mathbf{D}\mathbf{V}^T$ and \mathbf{D}_{s_2} is a diagonal matrix retaining the largest s_2 singular values of $\mathbf{Z} - \mathbf{\Theta}_1$ and truncating other singular values at zero.

Our algorithm for computing (1.4) is summarized.

Step 1.(Initialization) Supply a good initial estimate $(\hat{\mathbf{\Theta}}_1^{(0)}, \hat{\mathbf{\Theta}}_2^{(0)})$ in (1.4). Specify precision $\delta > 0$.

Step 2.(Iteration) At iteration m , update $\hat{\mathbf{\Theta}}_2^{(m)}$ in (1.7) with $\mathbf{\Theta}_1 = \hat{\mathbf{\Theta}}_1^{(m-1)}$. Then update $\hat{\mathbf{\Theta}}_1^{(m)}$ in (1.6) with $\mathbf{\Theta}_2 = \hat{\mathbf{\Theta}}_2^{(m)}$.

Step 3.(Stopping rule) Terminate if $|f(\hat{\mathbf{\Theta}}_1^{(m)}, \hat{\mathbf{\Theta}}_2^{(m)}) - f(\hat{\mathbf{\Theta}}_1^{(m-1)}, \hat{\mathbf{\Theta}}_2^{(m-1)})| \leq \delta$, where $f(\mathbf{\Theta}_1, \mathbf{\Theta}_2) = \|\mathbf{\Theta}_1 + \mathbf{\Theta}_2 - \mathbf{Z}\|_F^2$. Let m^* be the index at termination. The estimate is then $(\hat{\mathbf{\Theta}}_1^{(m^*)}, \hat{\mathbf{\Theta}}_2^{(m^*)})$.

Nonconvex minimization (1.5): A general case

The problem of solving for $\mathbf{\Theta}_2$ in (1.5) given $\mathbf{\Theta}_1$ reduces to that of constrained rank minimization

$$\min_{\mathbf{\Theta}_2} \|\mathbf{A}\mathbf{\Theta}_2 - (\mathbf{Z} - \mathbf{A}\mathbf{\Theta}_1)\|_F^2 \quad \text{subject to} \quad r(\mathbf{\Theta}_2) \leq s_2, \quad (1.8)$$

provided that $\mathbf{\Theta}_1$ satisfies the sparsity constraint in (1.5). Now write $\mathbf{\Theta}_2 \equiv \mathbf{C}\mathbf{F}$, where \mathbf{C} and \mathbf{F} are $p \times r$ and $r \times k$ matrices with $r \leq s_2$, consisting of a basis of the column space and that of the row space of $\mathbf{\Theta}_2$, respectively. Note that $\{\mathbf{\Theta}_2 : r(\mathbf{\Theta}_2) \leq s_2\} = \{\mathbf{\Theta}_2 : \mathbf{\Theta}_2 = \mathbf{C}\mathbf{F}, r \leq s_2\}$. Then solving (1.8) is equivalent to that

$$\min_{\mathbf{C}, \mathbf{F}} \|\mathbf{A}(\mathbf{C}\mathbf{F}) - (\mathbf{Z} - \mathbf{A}\mathbf{\Theta}_1)\|_F^2, \quad (1.9)$$

An application of an argument of (Xing et al., 2012) yields a global minimizer of

(1.9), which has an analytic form

$$\hat{\Theta}_2(\Theta_1) = \hat{C}\hat{F}, \quad \hat{C} = \mathbf{V}\mathbf{D}^{-1}\mathbf{U}_w, \quad \hat{F} = \mathbf{D}_w\mathbf{V}_w^T, \quad (1.10)$$

where \mathbf{D} is a $r(\mathbf{A}) \times r(\mathbf{A})$ diagonal singular vector matrix based on SVD of $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T$, \mathbf{D}_w is also a diagonal matrix of s_2 leading singular values of $\mathbf{W} \equiv \mathbf{U}^T(\mathbf{Z} - \mathbf{A}\Theta_1)$ and $\mathbf{U}_w, \mathbf{V}_w$ are matrices consisting of the corresponding right and left singular vectors.

Note that computation involves only the first s_2 largest singular values. Therefore, we employ the randomized truncated SVD method (Halko et al., 2011), for efficient computation of a large problem. This amounts to a complexity of order $O(pk \log r)$, as compared to $O(\min(pk^2, p^2k))$ of a conventional SVD method (Golub & Van, 1996).

Solving for Θ_1 in (1.5) given Θ_2 , on the other hand, becomes the problem of sparsity pursuit. In particular, we solve, assuming that $r(\Theta_2) \leq s_2$,

$$\min_{\Theta_1} \|\mathbf{A}\Theta_1 - (\mathbf{Z} - \mathbf{A}\Theta_2)\|_F^2, \quad \text{subject to} \quad \frac{1}{\tau} \sum_{\theta_{ij} \in \Theta_1} \min(|\theta_{ij}|, \tau) \leq s_1, \quad (1.11)$$

which is solved iteratively by a difference of convex (DC) programming, constructing a convex set containing the original constrained set. The constraint in (1.5) is defined by $J(\Theta_1) = S_1(\Theta_1) - S_2(\Theta_1)$ with $S_1(\Theta_1) = \frac{1}{\tau} \sum |\theta_{ij}|$ and $S_2(\Theta_1) = \frac{1}{\tau} \sum \max(|\theta_{ij}| - \tau, 0)$ are convex in Θ_1 . Then a sequence of upper approximations of $J(\Theta_1)$ is constructed: At iteration step m by $J^{(m)}(\Theta_1) = \sum_{\theta_{ij} \in \Theta_1} \left(\frac{|\theta_{ij}|}{\tau} I(|\hat{\theta}_{ij}^{(m-1)}| \leq \tau) + I(|\hat{\theta}_{ij}^{(m-1)}| > \tau) \right)$. This yields a sequence of convex minimization subproblems with convex constraints: At iteration step m , we solve

$$\min_{\Theta_1} \|\mathbf{A}\Theta_1 - (\mathbf{Z} - \mathbf{A}\Theta_2)\|_F^2, \quad \text{subject to} \quad J^{(m)}(\Theta_1) \leq s_1. \quad (1.12)$$

For (1.12), we develop a gradient projection method. First, we generalize an l_1 -ball

result of (Liu & Ye, 2009) to (1.12).

Lemma 1.2 (*Projection*) For any set $K \subseteq \{1, 2, \dots, n\}$,

$$\mathbf{x}^* = \mathcal{T}_{K,z}(\mathbf{v}) = \underset{\mathbf{x} \in \mathbb{R}^n: \sum_{i \in K} |x_i| \leq z}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{x} - \mathbf{v}\|_2^2,$$

where $\mathcal{T}_{K,z} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is a projection operator defined by

$$\mathcal{T}_{K,z}(\mathbf{v})_i = \operatorname{sign}(v_i) \max(|v_i| - \lambda^*, 0)$$

where $\lambda^* = 0$ if $\sum_{i \in K} |v_i| \leq z$ or $i \notin K$ and $\lambda^* = \frac{\sum_{i \in K \setminus K_0} |v_i| - z}{|K| - |K_0|}$ otherwise, and $K_0 = \{j : \sum_{i \in K} \max(|v_i| - |v_j|, 0) - z > 0\}$.

Before solving (1.12), we simply extend the fast iterative shrinkage-thresholding (FISTA) algorithm (Beck & Teboulle, 2009) to solving (1.13).

Lemma 1.3 For any set K defined in Lemma 1.2, a global minimizer of

$$\min_{\mathbf{x} \in \mathbb{R}^n: \sum_{i \in K} |x_i| \leq z} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 \tag{1.13}$$

can be obtained by FISTA iteratively: At iteration step t :

$$\begin{aligned} \mathbf{x}^{(t)} &= \mathcal{T}_{K,z} \left(\mathbf{y}^{(t)} - \frac{1}{2L} \mathbf{A}^T (\mathbf{A}\mathbf{y}^{(t)} - \mathbf{b}) \right), \\ \rho_{t+1} &= \frac{1 + \sqrt{1 + 4\rho_t^2}}{2}, \\ \mathbf{y}^{(t+1)} &= \mathbf{x}^{(t)} + \left(\frac{\rho_t - 1}{\rho_{t+1}} \right) (\mathbf{x}^{(t)} - \mathbf{x}^{(k-1)}), \end{aligned}$$

where L is the largest singular value of \mathbf{A} .

Next we solve (1.12) using Lemma 1.3, which yields an analytic updating formula

in a matrix form.

Then a global minimizer of (1.12) is computed using an iterative scheme with respect to t as follows:

$$\begin{aligned} \mathbf{v}^{(1)} &= \hat{\Theta}_1^{(m,0)} = \hat{\Theta}_1^{(m-1)}, \quad \rho_1 = 1, \\ \hat{\Theta}_1^{(m,t)} &= \mathcal{T}_{K^{(m)}, z^{(m)}} \left(\mathbf{v}^{(t)} - \frac{1}{2\lambda_{\max}(\mathbf{A}^T \mathbf{A})} \mathbf{A}^T [\mathbf{A} \mathbf{v}^{(t)} - (\mathbf{Z} - \mathbf{A} \Theta_2)] \right), \\ \rho_{t+1} &= \frac{1 + \sqrt{1 + 4\rho_t^2}}{2}, \quad \mathbf{v}^{(t+1)} = \hat{\Theta}_1^{(m,t)} + \left(\frac{\rho_t - 1}{\rho_{t+1}} \right) (\hat{\Theta}_1^{(m,t)} - \hat{\Theta}_1^{(m,t-1)}), \end{aligned} \quad (1.14)$$

where $K^{(m)} = \{(i, j) : |\hat{\theta}_{ij}^{(m-1)}| \leq \tau\}$, $z^{(m)} = \tau(s_1 - \sum_{\theta_{ij} \in \Theta_1} I(|\hat{\theta}_{ij}^{(m-1)}| > \tau))$ and $\lambda_{\max}(\cdot)$ denotes the largest eigenvalue of a matrix.

The algorithm is summarized as follows.

Algorithm 2:

Step 1.(Initialization) Supply a good initial estimate $(\hat{\Theta}_1^{(0)}, \hat{\Theta}_2^{(0)})$ in (1.5). Specify precision $\delta > 0$.

Step 2.(Iteration) At iteration m , compute candidate $\hat{\Theta}_2$ in (1.10) with $\Theta_1 = \hat{\Theta}_1^{(m-1)}$ and candidate $\hat{\theta}_{ij} \in \hat{\Theta}_1$ in (1.14) with $\mathbf{A} \Theta_2 = \mathbf{A} \hat{\Theta}_2^{(m-1)}$.

Step 3.(Maximum block improvement) At each iteration m , determine which of the two candidates $(\hat{\Theta}_1, \hat{\Theta}_2^{(m-1)})$ and $(\hat{\Theta}_1^{(m-1)}, \hat{\Theta}_2)$ for updating according to the amounts of improvement. That is, update $(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)}) = (\hat{\Theta}_1, \hat{\Theta}_2^{(m-1)})$ if $f(\hat{\Theta}_1, \hat{\Theta}_2^{(m-1)}) \leq f(\hat{\Theta}_1^{(m-1)}, \hat{\Theta}_2)$; update $(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)}) = (\hat{\Theta}_1^{(m-1)}, \hat{\Theta}_2)$ otherwise.

Step 4.(Stopping rule) Terminate if $|f(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)}) - f(\hat{\Theta}_1^{(m-1)}, \hat{\Theta}_2^{(m-1)})| \leq \delta$. Denote by m^* the index at termination. The final estimate is

$$\hat{\Theta}_1 = \hat{\Theta}_1^{(m^*)}, \quad \hat{\Theta}_2 = \hat{\mathbf{C}} \hat{\mathbf{F}},$$

where $\hat{\mathbf{C}}$ and $\hat{\mathbf{F}}$ are defined in (1.10) with $\Theta_1 = \hat{\Theta}_1$.

1.2.4 Computational properties

This section discusses computational properties of **Algorithms 1** and **2**. For non-convex minimization, our methods may not guarantee a global minimizer for (1.3). However, the following lemma says that our solution of **Algorithms 1** and **2** yields a stationary point of the cost function. Note that the scheme of maximum block improvement is essential for the result of Lemma 1.5.

Lemma 1.4 *The minimal cost function $f(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)})$ in **Algorithm 1** is strictly decreasing in m before termination. Moreover, the solution is a stationary point of $f(\Theta_1, \Theta_2)$ in that $\theta_{ij}^{(*)} = \operatorname{argmin}_{\theta_{ij} \in \Theta_k; k=1,2} f((\Theta_1^*, \Theta_2^*) \setminus \theta_{ij})$, where $(\Theta_1, \Theta_2) \setminus \theta_{ij}$ is the set of parameters of (Θ_1, Θ_2) without one component θ_{ij} in Θ_1 or Θ_2 , and (Θ_1, Θ_2) satisfy the constraints in (1.5).*

Lemma 1.5 *If \mathbf{A} is of full rank, then $\hat{\Theta}_1$ computed from **Algorithm 2** satisfies the constraints in (1.12). Moreover, the minimal cost function $f(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)})$ is strictly decreasing in m before termination. Finally, if the solution $(\hat{\Theta}_1, \hat{\Theta}_2)$ satisfies (1.5) and it is a stationary point of $f(\Theta_1, \Theta_2)$ in that*

$$\theta_{ij}^{(*)} = \operatorname{argmin}_{\theta_{ij} \in \Theta_k; k=1,2} f((\Theta_1^*, \Theta_2^*) \setminus \theta_{ij}),$$

where $(\Theta_1, \Theta_2) \setminus \theta_{ij}$ is the set of parameters of (Θ_1, Θ_2) without one component θ_{ij} in Θ_1 or Θ_2 , and (Θ_1, Θ_2) satisfy the constraints in (1.5).

With regard to the computational complexity of **Algorithms 1** and **2**, the method of truncated SVD yields an approximated SVD with a complexity of $O(pk \log r + (p+k)r^2)$ operations (Halko et al., 2011). Sorting requires a complexity of $O(pk \log(pk))$. For FISTA, the convergence rate is $O(1/t^2)$ (Beck & Teboulle, 2009), where t is the number of iterations. Overall, the computational complexity of **Algorithm 1** is

$O(pk \log(pk) + (p+k)r^2)I_2$, while that of **Algorithm 2** is $O((pk \log r + (p+k)r^2 + I_1/\varepsilon^2)I_2)$, where ε denotes the precision specified in **Algorithm 2**, and I_1 and I_2 is the number of DC iteration and blockwise iteration, respectively. Based on our experience, I_1 and I_2 are about between 3 and 20.

1.3 Theory

This section drives a finite-sample probability error bound for reconstruction of the true Θ^0 by $\hat{\Theta}^{L_0}$, which is a global minimizer of (1.3) in that $\hat{\Theta}^{L_0} = \hat{\Theta}_1^{L_0} + \hat{\Theta}_2^{L_0}$. Note that existence of a global minimizer is assured by the fact that the cost function (1.3) is bounded blow by zero. Moreover, we will provide an insight into simultaneous pursuit of the low rank and sparsity structures through matrix decomposition by contrasting the proposed method with (s_1, s_2) against low rank approximation alone with $(s_1 = 0, s_2)$ and sparsity pursuit alone with $(s_1, s_2 = 0)$.

Let $\|\Theta\|_\infty = \max_i \sum_j |\theta_{ij}|$ and $\|\Theta\|_{\max} = \max_{ij} |\theta_{ij}|$ are the L_∞ -norm and max norm respectively. Before proceeding, we define a parameter space Λ as $\{\Theta = \Theta_1 + \Theta_2 : \|\Theta_1\|_0 \leq s_1, \|\Theta_1\|_{\max} \leq l_1, \Theta_2 = \mathbf{C}\mathbf{F}, \max(\|\mathbf{C}\|_\infty, \|\mathbf{F}^T\|_\infty) \leq l_2\}$, where $l_1, l_2 > 0$ are constant, \mathbf{C} is a $p \times s_2$ matrix, \mathbf{F} is a $s_2 \times k$ matrix, \mathbf{F}^T is the transport of \mathbf{F} and $s_2 > 0$ is an upper bound of $r(\Theta_2)$. Let $g(\Theta, \mathbf{Z})$ be the probability density of \mathbf{Z} with respect to dominating measure ν on Λ . Define the Hellinger distance between two densities as

$$h(\Theta, \Theta') = \frac{1}{2} \left(\int (g^{1/2}(\Theta, \mathbf{Z}) - g^{1/2}(\Theta', \mathbf{Z}))^2 d\nu \right)^{1/2}, \quad (1.15)$$

which will be used to measure estimation accuracy.

The following technical assumptions are made.

Assumption A: (Norm-relation) For any $\Theta, \Theta' \in \Lambda$ and any $\delta > 0$,

$$\int \sup_{\|\Theta - \Theta'\|_{\max} \leq \delta} (g^{1/2}(\Theta, y) - g^{1/2}(\Theta', y))^2 d\nu(y) \leq M^2 \delta^2,$$

where M might depend on p, k, s_1, s_2 and l_1, l_2 .

Assumption A specifies a norm relation between the metric $\|\cdot\|_{\max}$ over parameters and the Hellinger distance over the corresponding densities. This can be verified given a specific form of g .

Theorem 1.1 gives a probability error bound for $\hat{\Theta}^{L_0}$ under probability P under the true Θ^0 . Let (s_1^0, s_2^0) be the degree of sparsity and rank, as defined in $\text{Eff}(\Theta^0)$ in Lemma 1.1.

Theorem 1.1 Under **Assumptions A**, for any $\epsilon \geq \epsilon_{n,p,k}$

$$P\left(h(\hat{\Theta}^{L_0}, \Theta^0) \geq \epsilon\right) \leq 5 \exp(-c_1 n \epsilon^2),$$

$$\epsilon_{n,p,k} = \frac{C_{p,k}}{\sqrt{n}} \sqrt{\log\left(\frac{\sqrt{n}}{C_{p,k}}\right)} \text{ with}$$

$$C_{p,k} = c_2 \sqrt{\log(2^9 M c_4 (l_2^3 + l_1))} \sqrt{(p+k)s_2^0 + s_1^0} + c_2 \sqrt{s_1^0 \log\left(e \frac{(p+k-r(\Theta^0))r(\Theta^0)}{s_1^0}\right)}. \quad (1.16)$$

If $\log(r(\Theta^0)) \leq d s_2^0$ for some $d > 0$, then it can be simplified:

$$C_{p,k} = c_3 \sqrt{\log(M)} \sqrt{(p+k-s_2^0)s_2^0},$$

where $c_1 - c_3$ are positive constants and M is defined in **Assumption A**. Moreover, as $n, p, k \rightarrow \infty$, $h^2(\hat{\Theta}^{L_0}, \Theta^0) = O_p(\epsilon_{n,p,k}^2)$, and $Eh^2(\hat{\Theta}^{L_0}, \Theta^0) = O(\epsilon_{n,p,k}^2)$, where $O_p(\cdot)$ and E denote the stochastic order and the expectation under P .

Corollary 1.1 gives an order of $\epsilon_{n,p,k}$ in three extreme situations with M held fixed.

Corollary 1.1 *Suppose M in Assumptions A is a constant independent of (p, k, s_1, s_2) .*

(i) *When Θ^0 is extremely sparse, that is, $\|\Theta^0\|_0 \leq p + k - 2$, $C_{p,k}$ in (1.16) is no worse than*

$$O\left(\sqrt{\|\Theta^0\|_0 \log((p + k - r(\Theta^0))r(\Theta^0)/\|\Theta^0\|_0)}\right).$$

(ii) *When Θ^0 is a low-rank matrix, $C_{p,k}$ in (1.16) is no worse than*

$$O\left(\sqrt{(p + k - r(\Theta^0))r(\Theta^0)}\right).$$

(iii) *When Θ^0 is dense, say $\|\Theta^0\|_0 \geq cpk$ for a constant $0 < c \leq 1$, and of full rank, $C_{p,k}$ in (1.16) is*

$$O\left(\max\left(\sqrt{(p + k - s_2^0)s_2^0}, \sqrt{s_1^0 \log\left(\frac{pk}{s_1^0}\right)}\right)\right).$$

Then $C_{p,k}^L = O\left(\sqrt{(p + k - r(\Theta^0))r(\Theta^0)}\right)$.

Corollary 1.2 and Theorem 1.2 give a similar result under the Hellinger distance and the Kullback-Leibler distance, respectively, assuming that ϵ_i follows a normal distribution.

Corollary 1.2 *If ϵ_i in (1.1) follows $N(0, \sigma^2 I_{k \times k})$, $\|\mathbf{A}\|_\infty$ is bounded, then the results in Corollary 1.1 continue to hold.*

Theorem 1.2 *Under the same assumptions in Corollary 1.2, we have, for any $\epsilon \geq$*

$\epsilon_{n,p,k}$,

$$P\left(K(\Theta^0, \hat{\Theta}^{L_0}) \geq 4\epsilon^2\right) \leq 5 \exp(-c_1 n \epsilon^2).$$

where $K(\cdot, \cdot)$ is Kullback-Leibler distance under normality and $\epsilon_{n,p,k}$ and c_2 remain to be the same as in Theorem 1. As $n, p, k \rightarrow \infty$, $K(\Theta^0, \hat{\Theta}^{L_0}) = O_p(\epsilon_{n,p,k}^2)$ and $EK(\Theta^0, \hat{\Theta}^{L_0}) = O(\epsilon_{n,p,k}^2)$.

Theorem 1.3 gives an error bound for $\|\hat{\Theta}^{L_0} - \Theta^0\|_F^2$ under the normal assumption when $\mathbf{A} = \mathbf{I}_{n \times p}$.

Theorem 1.3 *Assume that $\mathbf{A} = \mathbf{I}_{n \times p}$ with $n = \max(p, k)$. Under the same assumptions in Corollary 2 with $\sigma = O(\frac{1}{\sqrt{\max(p, k)}})$, as $n, p, k \rightarrow \infty$, $\|\hat{\Theta}^{L_0} - \Theta^0\|_F^2 = O_p(C'_{p,k} \log(\frac{1}{C'_{p,k}}))$, where*

$$C'_{p,k} = \frac{\log(\max(p, k)) \cdot [(p+k)s_2^0 + s_1^0] + s_1^0 \log\left(e^{\frac{(p+k-r(\Theta^0))r(\Theta^0)}{s_1^0}}\right)}{\max(p^2, k^2)}$$

1.4 Numerical examples

This section examines operating characteristics of the proposed method through simulations, and demonstrates its effectiveness on applications in image reconstruction and in time series analysis. In the literature, it is known that the state-of-art methods are the low-rank approximation method subject to rank restriction as well as its regularized version, which outperforms the low-rank approximation method with the trace-norm (Xing et al., 2012; She, 2013; Zhou & Tao, 2011). In Section 1.4.1, we contrast our proposed method with pursuing low rank and sparsity structures through matrix decomposition simultaneously, with the former low rank approximation method subject to rank restriction (low-rank alone), as well as the method based

on sparsity pursuit alone (sparsity alone). Here **Algorithm 2** are used. Most importantly, in Section 1.4.2, we compare the proposed method using **Algorithm 1** with two strong competitors the method of Go Decomposition (GoDec, (Zhou & Tao, 2011)) and the method augmented Lagrange multipliers (ALM, (Lin et al., 2009)) when $\mathbf{A} = \mathbf{I}_{n \times p}$ in (1.2). In simulations, codes for ALM and GoDec are used at the authors' website, and the initial values for **Algorithms 1** and **2** are set to be the zero-matrix

1.4.1 Simulation I: Operating characteristics

The simulated example is generated as follows. First, a $n \times p$ design matrix \mathbf{A} is sampled with each entry being iid $N(0, 1)$. Second, the true Θ_1 is a $p \times k$ matrix with all diagonals one and two more non-zeros (2 and 2) being randomly chosen with equal probability, and the true Θ_2 is generated by multiplying a $p \times r$ matrix with a $r \times k$ matrix with each entry following $N(1, 1)$. Moreover, each entry of \mathbf{E} is iid $N(0, 0.25)$. Throughout the simulations, Θ_1 and Θ_2 are held fixed with different values of (n, p, k) .

The proposed method is trained with a training set, and the optimal tuning parameters, minimizing the prediction mean squares error over an independent tuning set, are obtained through a bisection search over integer values. Then a method's performance is examined over a test set. The training, tuning and testing data sizes are n , $4n$ and $2n$.

For parameter estimation, we employ the mean squares error to evaluate performance

$$\frac{1}{4n} \|\mathbf{A}(\hat{\Theta} - \Theta^0)\|_F^2. \quad (1.17)$$

For rank recovery, we calculate the absolute difference between an estimated rank \hat{r}

$p = 30, k = 20$									
n	Ours				Low-rank alone		Sparsity alone		
	$ \hat{r} - r_0 $	TP	FP	MSE	$ \hat{r} - r_0 $	MSE	TP	FP	MSE
50	0.00	1.00	0.00	1.54	15.69	7.79	0.12	0.14	4650.35
	(0.00)	(0.00)	(0.00)	(0.30)	(2.41)	(1.03)	(0.21)	(0.02)	(511.55)
100	0.00	1.00	0.00	0.51	16.94	2.16	0.13	0.05	4399.38
	(0.00)	(0.00)	(0.00)	(0.08)	(0.24)	(0.18)	(0.22)	(0.01)	(429.41)
$p = 20, k = 30$									
n	Ours				Low-rank alone		Sparsity alone		
	$ \hat{r} - r_0 $	TP	FP	MSE	$ \hat{r} - r_0 $	MSE	TP	FP	MSE
50	0.00	1.00	0.00	1.06	16.66	5.06	0.43	0.06	4276.25
	(0.00)	(0.00)	(0.00)	(0.17)	(0.76)	(0.62)	(0.28)	(0.01)	(508.06)
100	0.00	1.00	0.00	0.46	16.99	1.88	0.53	0.06	4087.58
	(0.00)	(0.00)	(0.00)	(0.05)	(0.10)	(0.16)	(0.20)	(0.01)	(406.97)
$p = 40, k = 30$									
n	Ours				Low-rank alone		Sparsity alone		
	$ \hat{r} - r_0 $	TP	FP	MSE	$ \hat{r} - r_0 $	MSE	TP	FP	MSE
50	0.00	1.00	0.00	4.08	1.88	19.39	0.09	0.20	12018.68
	(0.00)	(0.00)	(0.00)	(1.21)	(0.59)	(1.57)	(0.20)	(0.04)	(1422.84)
$p = 50, k = 20$									
n	Ours				Low-rank alone		Sparsity alone		
	$ \hat{r} - r_0 $	TP	FP	MSE	$ \hat{r} - r_0 $	MSE	TP	FP	MSE
100	0.00	1.00	0.00	0.95	16.86	5.05	0.04	0.03	11262.97
	(0.00)	(0.00)	(0.00)	(0.15)	(0.35)	(0.40)	(0.14)	(0.01)	(1003.69)
$p = 200, k = 100$									
n	Ours				Low-rank alone		Sparsity alone		
	$ \hat{r} - r_0 $	TP	FP	MSE	$ \hat{r} - r_0 $	MSE			
300	3.76	0.92	0.00	8.26	29.56	54.24	-	-	-
	(1.24)	(0.23)	(0.00)	(0.86)	(7.84)	(0.81)	(-)	(-)	(-)

Table 1.1: Results of Simulation I. **Algorithm 2** is used for computation.

and the true rank r_0 , that is $|\hat{r} - r_0|$. For sparsity pursuit, we define the true positive (TP) as a ratio of the true positive numbers of nonzero estimates over the number of nonzeros in the true model, and the false positive (FP) as a ratio of the false positive numbers of nonzero estimates over the number of zeros in the true model. Here “Low rank alone”, “Sparsity alone” and “Ours” indicate the low rank method subject to rank restriction, the sparsity pursuit method, and the proposed method

As indicated in Table 1.1, the proposed method performs favorably against its counterpart—the low rank approximation method subject to rank restriction and sparsity pursuit alone, across all situations with different values of n , p and k . Moreover, the proposed method enables to identify two structures through matrix decomposition simultaneously. In particular, it recovers the true rank of the matrix with nearly zero $|\hat{r} - r_0|$ -values as compared to relatively large $|\hat{r} - r_0|$ -values, ranging from 6.7 to 29.6, for its low-rank counterpart. At the same time, the proposed method has high true positives ranging from .92 to 1.00 and low false positives between 0.00 and 0.01, as compared to true positives ranging 0.04 to .44 and false positives between 0.03 and 0.20 of its counterpart based on sparsity pursuit. This suggests that pursuit of two types of structures is indeed advantageous than that of either one structure individually. This is mainly because these two structures are complementary to each other. As a result, higher parameter estimation accuracy, as measured by the MSE values, can be realized. In fact, the amount of improvement is large, which ranges from 147% to 1185400%. To see how each method performs as (n, p) increases, we fix $k = 5$.

As suggested by Table 1.2, the proposed method yields more stable performance than its two counterparts whose performance deteriorates rapidly, as the level of difficulty of a problem escalates when p and k increase.

n	p	Ours				Low-rank alone		Sparsity alone		
		$ \hat{r} - r_0 $	TP	FP	MSE	$ \hat{r} - r_0 $	MSE	TP	FP	MSE
50	20	0.00	1.00	0.002	0.58	2.00	0.84	0.433	0.08	570
		(0.00)	(0.00)	(0.006)	(0.14)	(0.00)	(0.18)	(0.30)	(0.02)	(73)
50	30	0.00	0.57	0.01	1.29	1.97	1.98	0.18	0.08	3772.33
		(0.00)	(0.17)	(0.01)	(0.32)	(0.17)	(0.42)	(0.27)	(0.01)	(542.38)
50	40	0.00	1.00	0.001	3.57	1.67	5.43	0.07	0.05	1998
		(0.00)	(0.00)	(0.003)	(1.58)	(0.60)	(1.73)	(0.18)	(0.01)	(257)
50	50	0.82	0.36	0.01	487.43	0.82	12255	0.05	0.03	3797
		(0.84)	(0.38)	(0.01)	(1081.68)	(0.81)	(79570)	(0.15)	(0.01)	(539)
100	20	0.00	1.00	0.01	0.23	2.00	0.32	0.53	0.08	541
		(0.00)	(0.00)	(0.03)	(0.05)	(0.00)	(0.05)	(0.21)	(0.02)	(58)
100	30	0.00	0.71	0.01	0.36	2.00	0.54	0.19	0.03	1461
		(0.00)	(0.25)	(0.01)	(0.05)	(0.00)	(0.08)	(0.21)	(0.01)	(147)
100	40	0.00	0.98	0.01	0.53	2.00	0.83	0.10	0.03	1929
		(0.00)	(0.10)	(0.02)	(0.08)	(0.00)	(0.11)	(0.20)	(0.01)	(179)

Table 1.2: Results for Simulation I with fixed $k = 5$. **Algorithm 2** is used for computation.

1.4.2 Simulation II: Comparison

To compare with ALM (Lin et al., 2009) and GoDec (Zhou & Tao, 2011) for RPCA, consider the case of $\mathbf{A} = \mathbf{I}_{n \times p}$ in (1.2) and $p = k$ as in these papers. GoDec minimizes

$$\min_{\Theta_1, \Theta_2} \|\mathbf{Z} - \Theta_1 - \Theta_2\|_F^2 \quad \text{subject to } \text{card}(\Theta_1) \leq s_1, \text{rank}(\Theta_2) \leq s_2, \quad (1.18)$$

where $\text{card}(\cdot)$ denotes the cardinality, and $s_j \geq 0$ are tuning parameters as in our case. Similarly, ALM that focuses on the non-noisy situation minimizes

$$\min_{\Theta_1, \Theta_2} \|\Theta_2\|_* + \lambda \sum_{\theta_{ij} \in \Theta_1} |\theta_{ij}|, \quad \text{subject to } \mathbf{Z} = \Theta_1 + \Theta_2, \quad (1.19)$$

where $\|\cdot\|_*$ is the nuclear-norm of a matrix.

Our simulation example remains the same as before except that the positions of nonzero elements in Θ_2 are randomly sampled with equal probability, in particular, $.1p$ and $.3p$ nonzeros are randomly chosen without replacement. For tuning, grid search is employed for GoDec in (1.18), with $1 \leq s_1 \leq (p + k)$ and $1 \leq s_2 \leq \min(p, k, 50)$; λ is fixed at $\frac{1}{\sqrt{p}}$ for (1.19).

From Table 1.3 and Table 1.4, it is evidenced that the proposed method outperforms ALM uniformly in terms of the MSE while being comparable to GoDec, in all the situations with different values of (p, k, σ) . Moreover, it always recovers the true rank of the matrix perfectly with $|\hat{r} - r_0| = 0$. Although ALM has comparable high TP values, its FP values are high as well in that they are at least 0.6488. As a result, ALM never captures the true rank.

p	k	σ	Method	$ \hat{r} - r_0 $	TP	FP	MSE
50	30	0.1	Ours	0.0000 (0.0000)	0.9940 (0.0343)	0.0000 (0.0002)	0.2366 (0.0251)
			ALM	13.0300 (0.6735)	1.000 (0.0000)	0.6488 (0.0082)	1.5057 (0.0576)
			GoDec	0.0000 (0.0000)	0.9940 (0.0342)	0.0000 (0.0001)	0.2363 (0.0245)
		1	Ours	0.0000 (0.0000)	0.0320 (0.0839)	0.0000 (0.0001)	2.5308 (0.2418)
			ALM	13.3900 (0.6651)	0.9280 (0.1223)	0.6540 (0.0080)	15.0569 (0.5758)
			GoDec	0.0000 (0.0000)	0.0300 (0.0823)	0.0001 (0.0003)	2.5537 (0.2523)
200	100	0.1	Ours	0.0000 (0.0000)	0.9770 (0.0337)	0.0000 (0.0000)	0.2345 (0.0169)
			ALM	54.3100 (0.7745)	1.0000 (0.0000)	0.7034 (0.0022)	4.9984 (0.0510)
			GoDec	0.0000 (0.0000)	0.9755 (0.0344)	0.0000 (0.0000)	0.2330 (0.0160)
		1	Ours	0.0000 (0.0000)	0.0075 (0.0206)	0.0000 (0.0000)	2.4469 (0.1387)
			ALM	54.2400 (0.7264)	0.9456 (0.0456)	0.7059 (0.0023)	49.9838 (0.5095)
			GoDec	0.0000 (0.0000)	0.0085 (0.0236)	0.0000 (0.0000)	2.4476 (0.1395)

Table 1.3: Results for Simulation II when $.1p$ nonzero are randomly chosen . **Algorithm 1** is used for computation.

p	k	σ	Method	$ \hat{r} - r_0 $	TP	FP	MSE
50	30	0.1	Ours	0.0000	0.9933	0.0002	0.2507
				(0.0000)	(0.0201)	(0.0003)	(0.0277)
			ALM	13.0000	1.0000	0.6472	1.5057
			(0.6195)	(0.0000)	(0.0079)	(0.0576)	
		GoDec	0.0000	0.9953	0.0001	0.2489	
			(0.0000)	(0.0171)	(0.0003)	(0.0271)	
1	Ours	0.0000	0.0373	0.0000	2.8870		
		(0.0000)	(0.0624)	(0.0001)	(0.2410)		
	ALM	13.37	0.9407	0.6531	15.0569		
	(0.6301)	(0.0621)	(0.0080)	(0.5758)			
GoDec	0.0000	0.0327	0.0001	2.8983			
	(0.0000)	(0.0653)	(0.0002)	(0.2504)			
200	100	0.1	Ours	0.0000	0.9867	0.0001	0.2495
				(0.0000)	(0.0164)	(0.0001)	(0.0198)
			ALM	54.3500	1.0000	0.7030	4.9984
			(0.6571)	(0.0000)	(0.0023)	(0.0510)	
		GoDec	0.0000	0.9882	0.0000	0.2479	
			(0.0000)	(0.0152)	(0.0001)	(0.0191)	
1	Ours	0.0000	0.0080	0.0000	2.8254		
		(0.0000)	(0.0122)	(0.0000)	(0.1402)		
	ALM	54.2200	0.9467	0.7054	49.9838		
	(0.6289)	(0.0297)	(0.0022)	(0.5095)			
GoDec	0.0000	0.0075	0.0000	2.8237			
	(0.0000)	(0.0135)	(0.0000)	(0.1409)			

Table 1.4: Results for Simulation II when $.3p$ nonzero are randomly chosen. **Algorithm 1** is used for computation.



Figure 1.1: The converted AR face image with markup points.

1.4.3 AR Face Database 20pt Markup

For face image reconstruction, we use a subset of *AR Face Data* for this experiment. The original image is available at http://www-prima.inrialpes.fr/FGnet/data/05-ARFace/markup_large.png, which is a colored one with size of $186 \times 200 \times 3$. To enable detailed testing, the image has been labeled with 20 facial features on the face. We convert the image into black and white and reduce it to size 171×180 . The target image is displayed in Figure 1.1.

Twenty one markup points around eyes, nose, mouth and cheeks, which are used to test face recognition or verification performance when the exact location of the face and features are known. To identify the locations, we extract sparse (Θ_1) and low-rank (Θ_2) structures for the face images as described by the matrix decomposition into Θ_1 and Θ_2 . For this purpose, \mathbf{A} in (1.3) is set to be the identity matrix of size 171×171 . Figure 1.2 and Figure 1.3 display two decomposed structures for the AR face images by the proposed method with different sparse and rank constraint parameters in (1.3).

As indicated in Figure 1.2 and Figure 1.3, the sparseness structure describes characteristics/detailed marks of the face, whereas the low-rank structure displays the rough outlook of the human face. This confirms our discussion regarding local and global features in the Introduction. Visually, both the first panels in Figure 1.2 and

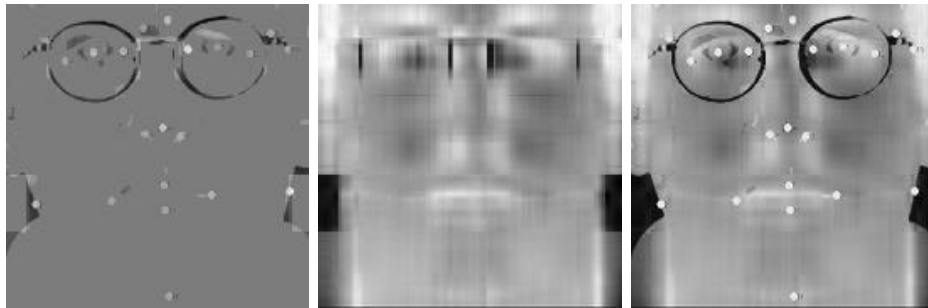


Figure 1.2: Extracted sparsity (first), low-rank (second) structures as well as the reconstructed image by the proposed method for AR face images; where the tuning parameters are set to $s_1 = 2500$, $s_2 = 5$.



Figure 1.3: Extracted sparsity (first), low-rank (second) structures as well as the reconstructed image by the proposed method for AR face images; where the tuning parameters are set to $s_1 = 2100$, $s_2 = 10$.

Figure 1.3 preserve at least 60% markup points, especially the points around nose two sides of face and lip. In other words, the sparsity structure captures most of markup points. Similarly, the second panels retain the overall look of the face. Most interestingly, this decomposition tends to remove the glasses from the human face.

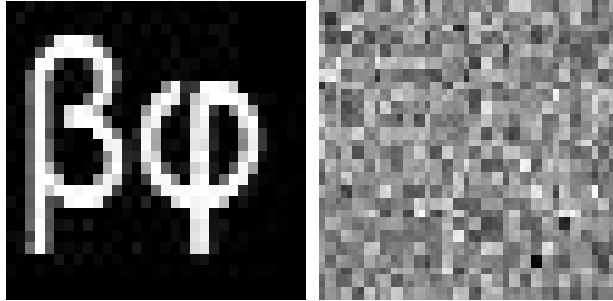


Figure 1.4: Original image (left) versus its noisy version (right).

1.4.4 Greek Letters Image Reconstruction

Now consider a 26×31 black-white image of two Greek letters β and ϕ , where its noisy version is obtained by adding noise $N(0, 1)$ after dividing the original matrix values by 100. The ratio of the maximum value of the image to the noise standard deviation is about 2.5. The images are displayed in Figure 1.4.

Our goal is reconstruction of the original image from its noise version, with a focus on restoration of detailed structures of the letters. Towards this end, we apply the proposed method and contrast with its counterpart based on sparse pursuit alone and low-rank approximations. Specifically, let \mathbf{A} to be the identity matrix of size 31×31 and Θ be a 31×26 parameter matrix in (1.3). For each method, grid search is performed for tuning, with $s_1 = (10, 20, 30, 50)$, $1 \leq s_2 \leq \min(p, k) = 26$ and $\tau = (0.05, 0.1, 0.2)$. For each method, the 10-*fold* cross-validation is employed. The reconstructed images are displayed in Figure 1.5.

Visually, the first two reconstructed images by the low-rank method and the sparsity method give the rough shape of two letters, but the letters β and ϕ not distinguishable with blurred segments in places, especially the right middle of β and the top of ϕ . By comparison, the third reconstructed image by our method enables to reconstruct the complete shape of these two letters, and yield the best quality of reconstruction.

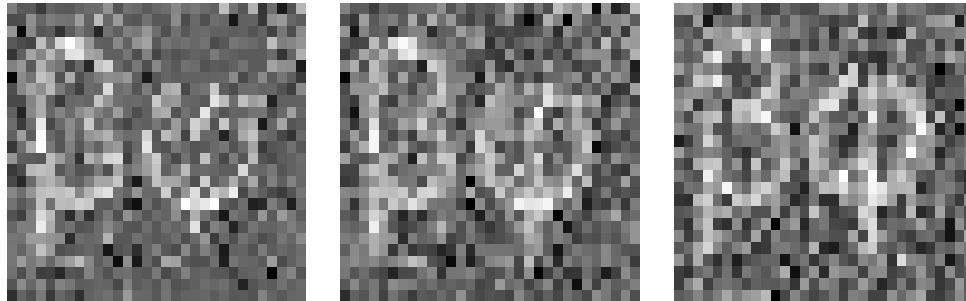


Figure 1.5: Reconstructed images based on sparsity alone (first), low-rank alone (second) and our method (third). **Algorithm 2** is used for computation.

1.4.5 US Macroeconomic Time Series

This subsection examines multiple time series data described in (Stock & Watson, 2012). The data measures 143 US macroeconomic variables quarterly over a time span from February 1, 1959 to November 1, 2008. These variables are categorized into 13 groups and are summarized in Table 1.5.

For data analysis, we consider time series starting from August 1, 1959 to November 1, 2008 due to incomplete initial observations. Our goal is one-step ahead forecasting, and contrast the proposed method with low-rank alone and sparsity alone in terms of forecasting accuracy. Using a multivariate autoregressive model, that is, $\mathbf{y}_t = \mathbf{y}_{t-1}^T \Theta + \epsilon_t$, we place it in the framework of (1.1), where \mathbf{y}_t is a vector that records the values of various macroeconomic variables at time point t , and ϵ_t follows normal distribution. In the presence of multiplicity and non-stationarity for economics data like this, we consider some transformations. For instance, log growth rates for quantity variables are differenced, nominal interest rates are differenced, as well as the logarithms of changes in rates of inflation for price series are differenced. See (Stock & Watson, 2012) for processing the data set. For this data set, $p = k = 143$ in (1.1) and the design matrix \mathbf{A} is specified by the time series, which can be written as $\mathbf{A} = (\mathbf{y}_{t_0}, \mathbf{y}_{t_0+1}, \dots, \mathbf{y}_{t_0+d-1})^T$.

A one-step ahead K -fold cross validation (CV) criterion is used for tuning the

Group	Description	Examples of series	# series
1	GDP component	GDP, consumption, investment	16
2	IP	IP, capacity utilization	14
3	Employment	Sectoral&total employment and hours	20
4	Unemployment rate	Unemployment rate, total and by duration	7
5	Housing	Housing starts, total and by region	6
6	Inventories	NAPM inventories, new orders	6
7	Prices	Price indexes, aggregate&disaggregate, commodity prices	37
8	Wages	Average hourly earning, unit labor cost	6
9	Interest rates	Treasuries, corporate, term spreads, public- private spreads	13
10	Money	M1, M2, business loans, consumer credit	7
11	Exchange rates	Average&selected trading partners	5
12	Stock prices	Various stock price indexes	5
13	Consumer expectations	Michigan consumer expectations	1

Table 1.5: Economic indicators collected for U.S. macroeconomic time series.

time series (Arlot & Celisse, 2010). In particular, for design matrix \mathbf{A} , at each fold i , we use observations i to $n - K + i - 1$ for training and the observation $n - K + i$ for tuning, where K is a pre-assigned integer and $K - 1$ indicates the number of folds. Note that the values of p and k are close to the sample size n for this time series. We therefore choose $K \leq 20$ to maintain adequate training samples.

For tuning, the CV is optimized over a set of grids for $s_1 = (10, 20, 50, 100, 200)$, $1 \leq s_2 \leq \min(p, k)$ and $\tau = (0.02, 0.05, 0.1, 0.2)$. The results for $K = 11$ are reported in Table 1.6. The results for other K values are omitted due to similarity.

As suggested by Table 1.6, the proposed method outperforms its counterparts pursuing sparseness and low-rank alone. The amount of improvement over the low rank method and the sparsity method is 15% and 933%, respectively. The Q-Q plots

in Figure 1.6 indicate that the model assumption is adequate although some departure from normality has been detected. Overall, the proposed method performs reasonably well.

	Ours	Low-rank alone	Sparsity alone
$K = 11$	301.22	348.02	3111.89

Table 1.6: Prediction errors of U.S. macroeconomic data for $K = 11$. Here “Low rank alone”, “Sparsity alone” and ”Ours” indicate our method for low rank pursuit only, for sparsity pursuit only and for simultaneous pursuit of low rank and sparsity. **Algorithm 2** is used for computation.

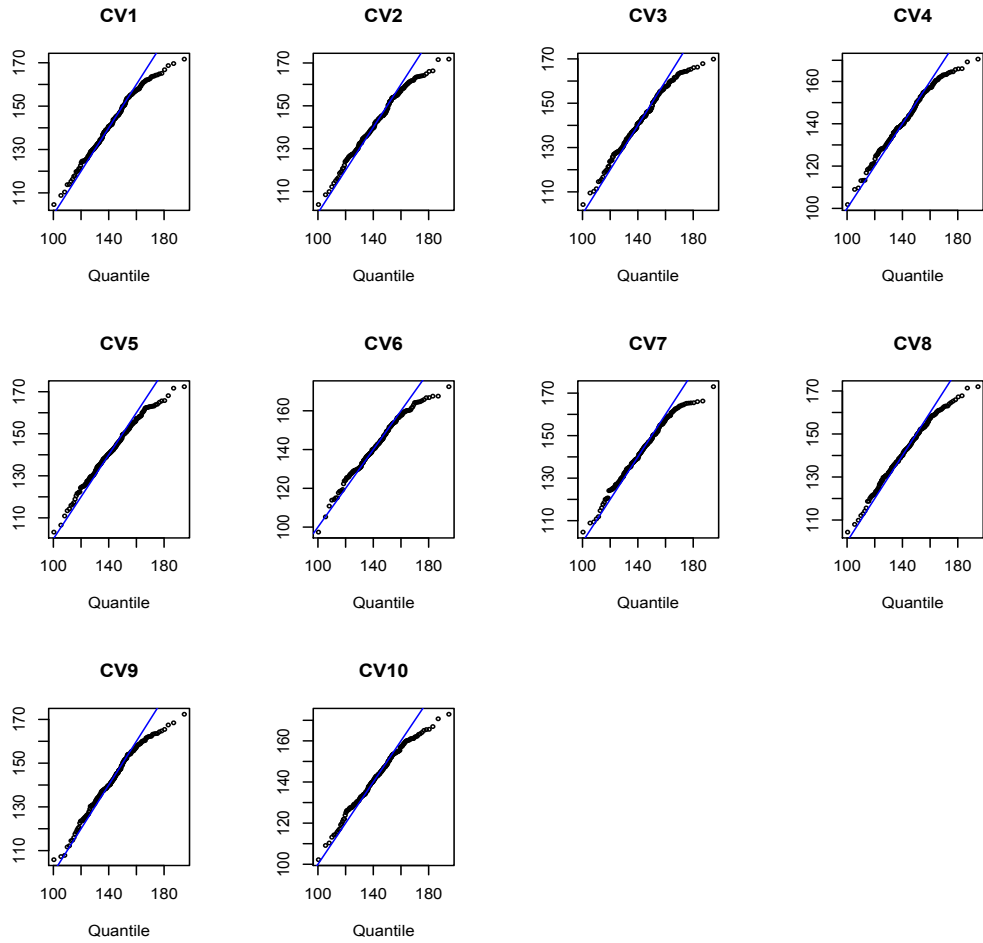


Figure 1.6: Q-Q plots for each-fold in U.S. macroeconomic time series data example, where points on a straight line indicates non-departure from normality.

1.5 Appendix

Proof of Lemma 1.1: Let $df(s, r) = s + (p + k - r)r$. By definition of the effective degrees of freedom, we obtain that

$$\text{Eff}(\Theta) \leq \min(df(0, r(\Theta^0)), df(\|\Theta^0\|_0, 0)).$$

To prove uniqueness in terms of (s, r) , suppose there exist $(\bar{s}, \bar{r}) \neq (\bar{s}', \bar{r}')$ such that $df(\bar{s}, \bar{r}) = df(\bar{s}', \bar{r}') = \min_{s,r} df(s, r)$. Without loss of generality, assume $\bar{r} = \bar{r}' - n_0 < \bar{r}'$, where $n_0 > 0$ is a positive integer. If $n_0 \leq \min(p, k) - \bar{r}$ and $\bar{r} < \min(p, k)$, then $\bar{s} + (p + k - \bar{r})\bar{r} = \bar{s}' + (p + k - \bar{r}')\bar{r}'$ implies that $\bar{s} = \bar{s}' + n_0(p + k - 2\bar{r} - n_0) \geq n_0(p + k - 2\bar{r} - n_0) > p + k - 2\bar{r} - 1$, which contradicts with the assumption that $s < p + k - 2r - 1$. Otherwise, if $\bar{r} = \min(p, k)$, \bar{s} must be zero. This completes the proof. ■

Proof of Lemma 1.2: Let $x_i = v_i$ for $i \notin K$. Then the problem reduces to the standard l_1 ball problem.

$$\underset{\sum_{i \in K} |x_i| \leq z}{\text{argmin}} \frac{1}{2} \sum_{i \in K} (x_i - v_i)^2.$$

The results follows by the proof of Theorem 1 of (Liu & Ye, 2009). ■

Proof of Lemma 1.3: It suffices to derive the basic step of ISTA in (Amit et al., 2007) for (1.13). Consider the following quadratic approximation of problem (1.13)

at a given point \mathbf{y} :

$$\min_{\mathbf{x} \in \mathbb{R}^n: \sum_{i \in K} |x_i| \leq z} Q_L(\mathbf{x}, \mathbf{y}) = \|\mathbf{A}\mathbf{y} - \mathbf{b}\|_2^2 + \langle \mathbf{x} - \mathbf{y}, \mathbf{A}^T(\mathbf{A}\mathbf{y} - \mathbf{b}) \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2, \quad (1.20)$$

where L is a Lipschitz constant of the function $\mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{b})$ with respect to \mathbf{x} . Solving (1.20) is equivalent to that of

$$\min_{\mathbf{x} \in \mathbb{R}^n: \sum_{i \in K} |x_i| \leq z} \left\| \mathbf{x} - \left(\mathbf{y} - \frac{1}{L} \mathbf{A}^T(\mathbf{A}\mathbf{y} - \mathbf{b}) \right) \right\|_2^2.$$

By Lemma 1.2, the solution is $\mathcal{T}_{K,z}(\mathbf{y} - \frac{1}{L} \mathbf{A}^T(\mathbf{A}\mathbf{y} - \mathbf{b}))$. The basic step of ISTA thus can be written as $\mathbf{x}^{(t)} = \mathcal{T}_{K,z}(\mathbf{x}^{(t-1)} - \frac{1}{L} \mathbf{A}^T(\mathbf{A}\mathbf{x}^{(t-1)} - \mathbf{b}))$. Then, Lemma 3 follows by taking L to be $\lambda_{\max}(\mathbf{A}^T \mathbf{A})$, where $\lambda_{\max}(\cdot)$ denotes the largest singular value. ■

Proof of Lemma 1.4: By (1.6) and (1.7), for any integer $m \geq 1$,

$$f(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)}) \geq f(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m+1)}) \geq f(\hat{\Theta}_1^{(m+1)}, \hat{\Theta}_2^{(m+1)}).$$

Meanwhile, it follows from (1.6) that

$$\begin{aligned} f(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)}) &= \|\mathbf{Z} - \hat{\Theta}_2^{(m)}\|_F^2 - \|\hat{\Theta}_1^{(m)}\|_F^2 \\ &\geq \|\mathbf{Z} - \hat{\Theta}_2^{(m+1)} - \hat{\Theta}_1^{(m)}\|_F^2 \\ &\geq \|\mathbf{Z} - \hat{\Theta}_2^{(m+1)}\|_F^2 - \|\hat{\Theta}_1^{(m)}\|_F^2. \end{aligned}$$

Therefore $\|\mathbf{Z} - \hat{\Theta}_2^{(m)}\|_F^2$ is lower bounded and decreasing in m . Moreover, by the monotone properties of $f(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)})$, $\|\hat{\Theta}_1^{(m)}\|_F^2$ converges as $m \rightarrow \infty$. Then there exists a subsequence $\{m_k\}$ such that $(\hat{\Theta}_1^{(m_k)}, \hat{\Theta}_2^{(m_k)}) \rightarrow (\hat{\Theta}_1^{(m^*)}, \hat{\Theta}_2^{(m^*)})$.

Let $R_{ij}(\Theta_1, \Theta_2) \in \operatorname{argmin}_{\theta_{ij} \in \Theta_1 \text{ or } \theta_{ij} \in \Theta_2} f((\Theta_1, \Theta_2) \setminus \theta_{ij})$. Let the cost function

for θ_{ij} to be $f_m(\theta_{ij}) = f((\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)}) \setminus \theta_{ij})$, where other components of $(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)})$ are held fixed. Then

$$\begin{aligned} f_{m_k}(R_{ij}(\hat{\Theta}_1^{(m^*)}, \hat{\Theta}_2^{(m^*)})) &\geq f_{m_k}(R_{ij}(\hat{\Theta}_1^{(m_k)}, \hat{\Theta}_2^{(m_k)})) \\ &\geq \min\left(f((\hat{\Theta}_1^{(m_k)}, \hat{\Theta}_2^{(m_k+1)})), f((\hat{\Theta}_1^{(m_k)}, \hat{\Theta}_2^{(m_k)}))\right) \\ &\geq f((\hat{\Theta}_1^{(m_k+1)}, \hat{\Theta}_2^{(m_k+1)})). \end{aligned}$$

As $m \rightarrow \infty$, by continuity of $f(\cdot)$, $f_{(m^*)}(R_{ij}(\hat{\Theta}_1^{(m^*)}, \hat{\Theta}_2^{(m^*)})) \geq f(\hat{\Theta}_1^{(m^*)}, \hat{\Theta}_2^{(m^*)})$, where the equality holds by the definition of R_{ij} . Hence, for each $\theta_{ij} \in \Theta_l$; $l = 1, 2$, $\hat{\theta}_{ij}^{(m^*)} = R_{ij}(\hat{\Theta}_1^{(m^*)}, \hat{\Theta}_2^{(m^*)})$ is the optimal componentwise solution. The results of Lemma 4 then follow. ■

Proof of Lemma 5: First we prove that $\hat{\Theta}_1^{(m)}$ satisfies

$$\sum_{\theta_{ij} \in \Theta_1} \left(\frac{|\hat{\theta}_{ij}^{(m)}|}{\tau} I(|\hat{\theta}_{ij}^{(m)}| \leq \tau) + I(|\hat{\theta}_{ij}^{(m)}| > \tau) \right) \leq s_1. \quad (1.21)$$

Toward this end, we rewrite the left side of (1.21) as

$$\begin{aligned} &\sum_{\theta_{ij} \in \Theta_1} \left(\frac{|\hat{\theta}_{ij}^{(m)}|}{\tau} I(|\hat{\theta}_{ij}^{(m-1)}| \leq \tau) + I(|\hat{\theta}_{ij}^{(m-1)}| > \tau) \right) \\ &+ \sum_{\theta_{ij} \in \Theta_1} \left(\frac{|\hat{\theta}_{ij}^{(m)}|}{\tau} I(|\hat{\theta}_{ij}^{(m)}| \leq \tau) + I(|\hat{\theta}_{ij}^{(m)}| > \tau) - \frac{|\hat{\theta}_{ij}^{(m)}|}{\tau} I(|\hat{\theta}_{ij}^{(m-1)}| \leq \tau) - I(|\hat{\theta}_{ij}^{(m-1)}| > \tau) \right) \\ &= \sum_{\theta_{ij} \in \Theta_1} \left(\frac{|\hat{\theta}_{ij}^{(m)}|}{\tau} I(|\hat{\theta}_{ij}^{(m-1)}| \leq \tau) + I(|\hat{\theta}_{ij}^{(m-1)}| > \tau) \right) + I_m, \end{aligned} \quad (1.22)$$

where $I_m = \sum_{\theta_{ij} \in \Theta_1} \frac{|\hat{\theta}_{ij}^{(m)}| - \tau}{\tau} \left(I(|\hat{\theta}_{ij}^{(m)}| \leq \tau) - I(|\hat{\theta}_{ij}^{(m-1)}| \leq \tau) \right)$. Note that it follows

from the DC construction that

$$\sum_{\theta_{ij} \in \Theta_1} \left(\frac{|\hat{\theta}_{ij}^{(m)}|}{\tau} I(|\hat{\theta}_{ij}^{(m-1)}| \leq \tau) + I(|\hat{\theta}_{ij}^{(m-1)}| > \tau) \right) \leq s_1.$$

Thus, to establish (1.21), we only need to prove $I_m \leq 0$. Rewrite I as

$$I_m = \begin{cases} 0 & \text{if } \min(|\hat{\theta}_{ij}^{(m)}|, |\hat{\theta}_{ij}^{(m-1)}|) > \tau \text{ or } \max(|\hat{\theta}_{ij}^{(m)}|, |\hat{\theta}_{ij}^{(m-1)}|) \leq \tau, \\ \sum_{\theta_{ij} \in \Theta_1} \left(\frac{|\hat{\theta}_{ij}^{(m)}|}{\tau} - 1 \right) & \text{if } |\hat{\theta}_{ij}^{(m)}| \leq \tau \text{ and } |\hat{\theta}_{ij}^{(m-1)}| > \tau, \\ -\sum_{\theta_{ij} \in \Theta_1} \left(\frac{|\hat{\theta}_{ij}^{(m)}|}{\tau} - 1 \right) & \text{if } |\hat{\theta}_{ij}^{(m)}| > \tau \text{ and } |\hat{\theta}_{ij}^{(m-1)}| \leq \tau, \end{cases}$$

implying that $I_m \leq 0$. Then, (1.21) follows.

For stationarity, note that it follows from (1.21) that

$$f(\hat{\Theta}_1^{(m-1)}, \hat{\Theta}_2^{(m-1)}) \geq f(\hat{\Theta}_1^{(m-1)}, \hat{\Theta}_2) \geq f(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)}),$$

where $\hat{\Theta}_2$ is defined in Step 2 of **Algorithm 2**.

Suppose that termination index m^* is infinite. Then we will prove that $\hat{\Theta}_1^{(m)} \rightarrow \hat{\Theta}_1^{(m^*)}$ as $m \rightarrow m^* = \infty$. When $m^* = \infty$, $\hat{\Theta}_1^{(m)}$ must be updated infinitely because $\hat{\Theta}_2^{(m)}$ is analytically solved. First consider, at step m , Θ_1 is updated whereas $\Theta_2 = \hat{\Theta}_2^{(m)}$. Denote by $\Lambda(\Theta_1, \Theta_2, \lambda^*)$ the dual problem of (1.12), where λ^* is the optimal Lagrange multiplier and $\Theta_2 = \hat{\Theta}_2^m$. Then

$$\begin{aligned} & f(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)}) - f(\hat{\Theta}_1^{(m+1)}, \hat{\Theta}_2^{(m+1)}) \\ &= \Lambda(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)}, \lambda^*) - \Lambda(\hat{\Theta}_1^{(m+1)}, \hat{\Theta}_2^{(m)}, \lambda^*) \\ & \quad - \lambda^* \sum_{\theta_{ij} \in \Theta_1} \left(\frac{|\hat{\theta}_{ij}^{(m)}|}{\tau} I(|\hat{\theta}_{ij}^{(m)}| \leq \tau) + I(|\hat{\theta}_{ij}^{(m)}| > \tau) - s_1 \right) \end{aligned}$$

The equality holds because $\hat{\Theta}_1^{(m+1)}$ is the global minimizer of a convex problem (1.12),

attaining at constraint boundaries, i.e

$$\sum_{\theta_{ij} \in \Theta_1} \left(\frac{|\hat{\theta}_{ij}^{(m+1)}|}{\tau} I(|\hat{\theta}_{ij}^{(m)}| \leq \tau) + I(|\hat{\theta}_{ij}^{(m)}| > \tau) - s_1 \right) = 0.$$

An application of the Taylor expansion to $\Lambda(\Theta_1, \hat{\Theta}_2^{(m)}, \lambda^*)$ at $\Theta_1 = \hat{\Theta}_1^{(m+1)}$ yields that

$$\begin{aligned} & f(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)}) - f(\hat{\Theta}_1^{(m+1)}, \hat{\Theta}_2^{(m+1)}) \\ &= \left\langle \frac{\partial \Lambda}{\partial \Theta_1}(\hat{\Theta}_1^{(m+1)}, \hat{\Theta}_2^{(m)}, \lambda^*), \hat{\Theta}_1^{(m)} - \hat{\Theta}_1^{(m+1)} \right\rangle \\ &+ \frac{1}{2} \langle \mathbf{A}(\hat{\Theta}_1^{(m)} - \hat{\Theta}_1^{(m+1)}), \mathbf{A}(\hat{\Theta}_1^{(m)} - \hat{\Theta}_1^{(m+1)}) \rangle \\ &- \lambda^* \sum_{\theta_{ij} \in \Theta_1} \left(\frac{|\hat{\theta}_{ij}^{(m)}|}{\tau} I(|\hat{\theta}_{ij}^{(m)}| \leq \tau) + I(|\hat{\theta}_{ij}^{(m)}| > \tau) - s_1 \right), \end{aligned}$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product. The first term in the right side of the equality is zero, because $\hat{\Theta}_1^{(m+1)}$ is the global minimizer and the third term is no less than zero by (1.21). Thus,

$$\begin{aligned} f(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)}) - f(\hat{\Theta}_1^{(m+1)}, \hat{\Theta}_2^{(m+1)}) &\geq \frac{1}{2} \langle \mathbf{A}(\hat{\Theta}_1^{(m)} - \hat{\Theta}_1^{(m+1)}), \mathbf{A}(\hat{\Theta}_1^{(m)} - \hat{\Theta}_1^{(m+1)}) \rangle \\ &\geq \frac{\lambda_{\min}(\mathbf{A}^T \mathbf{A})}{2} \|\hat{\Theta}_1^{(m)} - \hat{\Theta}_1^{(m+1)}\|_F^2, \quad (1.23) \end{aligned}$$

where $\lambda_{\min(\cdot)}$ is the smallest eigenvalue of a matrix. Therefore $f(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)})$ is lower bounded and decreasing in m , implying $f(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)})$ converges to some limit f^* as $m \rightarrow \infty$. By (1.23), convergence of $\hat{\Theta}_1^{(m)} \rightarrow \hat{\Theta}_1^{(m^*)}$ is established. Next consider the case in which Θ_2 is only updated finitely, say before step m_0 , using the same notation

with proof of Lemma 1.4, then for any $m > m_0$

$$f_m(R_{ij}(\hat{\Theta}_1^{(m^*)}, \hat{\Theta}_2^{(m^*)})) \geq f_m(R_{ij}(\hat{\Theta}_1^{(m)}, \hat{\Theta}_2^{(m)})) = f((\hat{\Theta}_1^{(m+1)}, \hat{\Theta}_2^{(m+1)})).$$

The second equality holds because the MBI is employed. As $m \rightarrow m^*$, by continuity of function f , $f_{(m^*)}(R_{ij}(\hat{\Theta}_1^{(m^*)}, \hat{\Theta}_2^{(m^*)})) \geq f(\hat{\Theta}_1^{(m^*)}, \hat{\Theta}_2^{(m^*)})$, where the equality holds by the definition of R_{ij} . Finally, we consider the case in which Θ_2 is updated infinitely. Then there is a subsequence $\{m_k\}$ such that $\hat{\Theta}_2^{(m_k)} \rightarrow \hat{\Theta}_2^{(m^*)}$. Similarly, $f_{m^*}(R_{ij}(\hat{\Theta}_1^{(m^*)}, \hat{\Theta}_2^{(m^*)})) = f(\hat{\Theta}_1^{(m^*)}, \hat{\Theta}_2^{(m^*)})$. Hence, for each $\theta_{ij} \in \Theta_l$, $l = 1, 2$, $\hat{\theta}_{ij}^{(m^*)} = R_{ij}(\hat{\Theta}_1^{(m^*)}, \hat{\Theta}_2^{(m^*)})$ is the optimal componentwise solution. The results of Lemma 1.5 then follow. ■

Let $\mathcal{B}_{S,r} = \{\Theta = \Theta_1^S + \Theta_2 : r(\Theta_2) = r\} \cap \Lambda$, a sub-parameter space with known sparsity structure S and rank r . Denote $H(\cdot, \Lambda)$ and $H^B(\cdot, \Lambda)$ to be the L_∞ entropy and bracketing Hellinger metric entropy for set Λ , respectively. The next two technical lemmas concern the size of the parameter space.

Lemma 1.6 *Suppose that Assumptions A is met.*

$$H^B(t, \mathcal{B}_{S,r}) \leq |S| \log(2Ml_1/t) + (p+k)r \log(2Ml_2^3/t),$$

where l_1, l_2 are constant and $M > 1$ is defined in Assumption A.

Lemma 1.7 *Suppose that Assumptions A is satisfied. If $s_1 = s_1^0$, $s_2 = s_2^0$, then*

$$\begin{aligned} H^B(t, \Lambda) \leq & 2(p+k)s_2^0 \log(2Ml_2^3\epsilon/t) + s_1^0 \log((1+2Ml_1)/t) \\ & + 2s_1^0 \log \left(e^{\frac{(p+k-r(\Theta^0))r(\Theta^0)}{s_1^0}} \right). \end{aligned}$$

Proof of Lemmas 1.6 and 1.7: For Lemma 1.6, note that $\Lambda = \cup_{|S| \leq s_1^0} \cup_{r \leq s_2^0} \mathcal{B}_{S,r}$. It suffices to calculate the entropy for each $\mathcal{B}_{S,r}$. Let

$$\Lambda_2 = \{(\Theta_1, \Theta_2) : \Theta_1, \Theta_2 \text{ satisfy conditions defined in } \Lambda\}.$$

For $\Theta = \Theta_1 + \Theta_2$ and $(\Theta_1, \Theta_2) \in \Lambda_2$, define $\mathcal{B}_\delta(\Theta_1, \Theta_2) = \{(\Theta'_1, \Theta'_2) \in \Lambda_2 : \|\Theta_1 - \Theta'_1\|_{\max} + \|\Theta_2 - \Theta'_2\|_{\max} \leq \delta\}$ to be the neighborhood of (Θ_1, Θ_2) . For any $\Theta' = \Theta'_1 + \Theta'_2$ with $(\Theta'_1, \Theta'_2) \in \mathcal{B}_\delta(\Theta_1, \Theta_2)$, by **Assumption A**,

$$\int \sup_{\mathcal{B}_\delta(\Theta_1, \Theta_2)} (g^{1/2}(\Theta, y) - g^{1/2}(\Theta', y))^2 d\nu(y) \leq M^2 \delta^2.$$

Combined the above with Lemma 2.1 of (Ossiander, 1987), we have

$$H^B(t, \mathcal{B}_{S,r}) \leq H(M^{-1}t, \mathcal{B}_{S,r}). \quad (1.24)$$

Since $\|\Theta_1\|_{\max}$ is bounded by l_1 , by constructing a $2t$ -net on $\mathcal{B}_{S,r}$ through the outer product of the t -nets on Θ_1^S and Θ_2 defined in the parameter space Λ , we can show that

$$H(M^{-1}t, \mathcal{B}_{S,r}) \leq |S| \log(2Ml_1/t) + H_r(M^{-1}t) \quad (1.25)$$

where $|S|$ is the number of nonzeros in Θ_1 and $H_r(M^{-1}t)$ is the entropy for Θ_2 with rank r . Let \mathbf{C} be a basis of column of Θ_2 , then there exists an $k \times r$ matrix \mathbf{F} such that $\Theta_2 = \mathbf{C}\mathbf{F}$. Hence

$$\|\Theta_2 - \Theta'_2\|_{\max} = \|\mathbf{C}\mathbf{F} - \mathbf{C}'\mathbf{F}'\|_{\max} \leq \|\mathbf{C}\|_{\infty} \|\mathbf{F} - \mathbf{F}'\|_{\max} + \|\mathbf{F}'^T\|_{\infty} \|\mathbf{C} - \mathbf{C}'\|_{\max}.$$

where $\|\Theta_{p \times k}\|_{\infty} = \max_{1 \leq i \leq p} \sum_{j=1}^k |\theta_{ij}|$ is the L_{∞} matrix-norm and $\|\Theta\|_{\max} = \max_{\theta_{ij} \in \Theta} |\theta_{ij}|$

is the max norm. Note that $\|\mathbf{C}\|_\infty$ and $\|\mathbf{F}^T\|_\infty$ are bounded by l_2 . This yields

$$H_r(M^{-1}t) \leq (p+k)r \log \frac{2l_2^3 M}{t}.$$

This, together with ((1.24)) and ((1.25)), implies Lemma 6.

For Lemma 1.7, note that

$$\begin{aligned} \exp(H^B(t, \Lambda)) &\leq \exp(H(M^{-1}t, \Lambda)) \\ &= \sum_{r=0}^{s_2^0} \sum_{|S|=0}^{s_1^0} \sum_{i=0}^{|S|} \binom{s_1^0}{i} \binom{(p+k-r(\Theta^0))r(\Theta^0)-s_1^0}{|S|-i} \exp(H(M^{-1}t, \mathcal{B}_{S,r})) \\ &\leq \binom{(p+k-r(\Theta^0))r(\Theta^0)}{s_1^0} \left(\sum_{|S|=0}^{s_1^0} \binom{s_1^0}{|S|} (2Ml_1/t)^{|S|} \right) \left(\sum_{r=0}^{s_2^0} (2Ml_2^3/t)^{(p+k)r} \right) \\ &\equiv \binom{(p+k-r(\Theta^0))r(\Theta^0)}{s_1^0} \times I \times II. \end{aligned}$$

Note that $\sum_{k=0}^n \binom{n}{k} a^k b^{n-k} = (a+b)^n$. Then $I = (1 + \frac{2Ml_1}{t})^{s_1^0}$ and $II \leq (s_2^0 + 1) \left(\frac{2Ml_2^3 \epsilon}{t} \right)^{(p+k)s_2^0}$. Thus,

$$\begin{aligned} H^B(t, \Lambda) &\leq \log \binom{(p+k-r(\Theta^0))r(\Theta^0)}{s_1^0} + \log(s_2^0 + 1) \\ &\quad + s_1^0 \log\left(1 + \frac{2Ml_1}{t}\right) + (p+k)s_2^0 \log\left(\frac{2Ml_2^3}{t}\right) \\ &\leq 2(p+k)s_2^0 \log(2Ml_2^3/t) + s_1^0 \log\left(\frac{1+2Ml_1}{t}\right) \\ &\quad + 2s_1^0 \log \left(e \frac{(p+k-r(\Theta^0))r(\Theta^0)}{s_1^0} \right), \end{aligned}$$

where e is the natural number and $0 < t < 1$. The last inequality follows Theorem 2.6 of (Stanica & Montgomery, 2001) that $\binom{b}{a} \leq \frac{b^{b+1/2}}{\sqrt{2\pi a^{a+1/2}(b-a)^{b-a+1/2}}} \leq \exp((a+1/2) \log(b/a) + a) \leq \exp(2a \log(b/a) + a)$ for any integer $0 < a < b$. This completes the proof. ■

Proof of Theorem 1.1: We apply a large deviation inequality in Theorem 2 of (Wong & Shen, 1995). To this end, we verify (1.2) there. By Lemma 1.7,

$$\begin{aligned}
& \int_{\epsilon^2/2^8}^{2^{1/2}\epsilon} (H^B(t/c_4, \Lambda))^{1/2} dt \\
& \leq \int_{\epsilon^2/2^8}^{2^{1/2}\epsilon} \sqrt{2(p+k)s_2^0 \log(2Ml_2^3 c_4/t) + s_1^0 \log((1+2Ml_1)c_4/t)} dt \\
& \quad + \int_{\epsilon^2/2^8}^{2^{1/2}\epsilon} \sqrt{2s_1^0 \log\left(e \frac{(p+k-r(\Theta^0))r(\Theta^0)}{s_1^0}\right)} dt \\
& \equiv I_1 + I_2,
\end{aligned}$$

for some constant $c_4 > 0$, say $c_4 = 10$. Then, for ϵ small,

$$\begin{aligned}
I_1 & \leq \sqrt{2}\epsilon \sqrt{2(p+k)s_2^0 \log(2^9 M l_2^3 c_4/\epsilon^2) + s_1^0 \log((1+2Ml_1)2^8 c_4/\epsilon^2)} \\
& \leq 2\epsilon \sqrt{(p+k)s_2^0 + s_1^0} \sqrt{\log(2^9 M c_4 (l_2^3 + l_1))} + 2 \log \frac{1}{\epsilon} \\
& \leq 2\sqrt{2}\epsilon \sqrt{\log(2^9 M c_4 (l_2^3 + l_1))} \sqrt{(p+k)s_2^0 + s_1^0} \cdot \sqrt{\log \frac{1}{\epsilon}}.
\end{aligned}$$

Similarly,

$$I_2 \leq 2\epsilon \sqrt{s_1^0 \log\left(e \frac{(p+k-r(\Theta^0))r(\Theta^0)}{s_1^0}\right)}.$$

Let $\epsilon_{n,p,k} = \frac{C_{p,k}}{\sqrt{n}} \log\left(\frac{C_{p,k}}{\sqrt{n}}\right)$ where

$$\begin{aligned}
C_{p,k} & = 2\sqrt{2}c_5^{-1} \sqrt{\log(2^9 M c_4 (l_2^3 + l_1))} \sqrt{(p+k)s_2^0 + s_1^0} \\
& \quad + 2c_5^{-1} \sqrt{s_1^0 \log\left(e \frac{(p+k-r(\Theta^0))r(\Theta^0)}{s_1^0}\right)}.
\end{aligned}$$

Then, for any $\epsilon \geq \epsilon_{n,p,k}$ and $c_5 = \frac{512}{(2/3)^{5/12}}$

$$\int_{\epsilon^2/2^8}^{2^{1/2}\epsilon} (H^B(t/c_4, \Lambda))^{1/2} dt \leq c_5^{-1} \sqrt{n} \epsilon^2.$$

By Theorem 2 of (Wong & Shen, 1995), $P\left(h(\hat{\Theta}^{L_0}, \Theta^0) \geq \epsilon\right) \leq 5 \exp(-c_1 n \epsilon^2)$, which yields $Eh^2(\hat{\Theta}^{L_0}, \Theta^0) = O(\epsilon_{n,p,k}^2)$ by using the fact that $h(\hat{\Theta}^{L_0}, \Theta^0) \leq 1$.

Consider a special situation when $\log(r(\Theta^0)) \leq ds_2^0$ for some constant $d > 0$ that is independent of p, k . Note that $s_1^0 < p + k - s_2^0$ and $p + k - r(\Theta^0) \leq p + k - s_2^0$. Then

$$\begin{aligned} s_1^0 \log\left(e^{\frac{(p+k-r(\Theta^0))r(\Theta^0)}{s_1^0}}\right) &\leq (p+k-s_2^0) \log\left(e^{\frac{(p+k-r(\Theta^0))r(\Theta^0)}{p+k-s_2^0}}\right) \\ &\leq (p+k-s_2^0) \log(er(\Theta^0)) \\ &\leq 2d(p+k-s_2^0)s_2^0. \end{aligned}$$

Thus, $I_1 + I_2$ is upper bounded by

$$2\sqrt{2}\epsilon \left(\sqrt{\log(2^9 M c_4 (l_2^3 + l_1))} + \sqrt{d} \right) \sqrt{(p+k)s_2^0 + s_1^0} \cdot \sqrt{\log \frac{1}{\epsilon}}.$$

Let $c_3 = 2\sqrt{2}c_5^{-1} \left(\sqrt{\log(2^9 c_4 (l_2^3 + l_1))} + \sqrt{d} \right) \sqrt{(p+k)s_2^0 + s_1^0}$. The result then follows. This completes the proof. ■

Proof of Corollary 1.1: If Θ^0 is sparse and $\|\Theta^0\|_0 \leq p+k-2$, then by the definition of effective degrees of freedom $s_0 = s_1^0 + (p+k-s_2^0)s_2^0 \leq \|\Theta^0\|_0$. This implies that

$$\begin{aligned} C_{p,k} &= O\left(\sqrt{\|\Theta^0\|_0}\right) + O\left(\sqrt{\|\Theta^0\|_0 \log\left(\frac{(p+k-r(\Theta^0))r(\Theta^0)}{\|\Theta^0\|_0}\right)}\right) \\ &= O\left(\sqrt{\|\Theta^0\|_0 \log\left(\frac{(p+k-r(\Theta^0))r(\Theta^0)}{\|\Theta^0\|_0}\right)}\right). \end{aligned}$$

The second inequality is because of nondecreasingness of \sqrt{x} and $\sqrt{x \log(a/x)}$ in x for $x \leq a/e$.

If Θ^0 is low-rank, we have

$$C_{p,k} = O \left(\sqrt{(p+k-r(\Theta^0))r(\Theta^0)} + \sqrt{s_1^0 \log((p+k-r(\Theta^0))r(\Theta^0)/s_1^0)} \right).$$

Note that $s_1^0 \log((p+k-r(\Theta^0))r(\Theta^0)/s_1^0) \leq \log((p+k-r(\Theta^0))r(\Theta^0)/e)$. The result follows.

If Θ^0 is dense and of full rank, then $(p+k-r(\Theta^0))r(\Theta^0)$ is of order $O(pk)$. Hence $C_{p,k}$ can be written as $O \left(\sqrt{(p+k-s_2^0)s_2^0} + \sqrt{s_1^0 \log(\frac{pk}{s_1^0})} \right)$. This completes the proof. ■

Proof of Corollary 1.2: It suffices to show the **Assumption A** is met. Let $f(\boldsymbol{\mu}_i, \mathbf{y}) = \frac{1}{(\sqrt{2\pi}\sigma)^k} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \boldsymbol{\mu}_i)^T(\mathbf{y} - \boldsymbol{\mu}_i)\right)$ for $i = 1, 2$. $\boldsymbol{\mu}_1 = \mathbf{a}^T \Theta$ and $\boldsymbol{\mu}_2 = \mathbf{a}^T \Theta'$. Then

$$\begin{aligned} & \int \sup_{\|\Theta - \Theta'\|_{\max} \leq \delta} (f^{1/2}(\boldsymbol{\mu}_1, \mathbf{y}) - f^{1/2}(\boldsymbol{\mu}_2, \mathbf{y}))^2 d\mathbf{y} \\ & \leq 2 - 2 \frac{1}{(\sqrt{2\pi}\sigma)^k} \int \inf_{\|\Theta - \Theta'\|_{\max} \leq \delta} \exp\left(-\frac{\|\mathbf{y} - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2}\|_2^2 + \frac{\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|_2^2}{2}}{2\sigma^2}\right) d\mathbf{y} \\ & \leq 2 - 2 \inf_{\|\Theta - \Theta'\|_{\max} \leq \delta} \exp\left(-\frac{\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|_2^2}{4\sigma^2}\right) \\ & \leq \frac{(\|\mathbf{a}\|_1)^2 \|\Theta - \Theta'\|_{\max}^2}{4\sigma^2} \leq \frac{(\|\mathbf{a}\|_1)^2 \delta^2}{4\sigma^2}. \end{aligned}$$

The second inequality follows from the invariance property of the normal distribution.

Corollary 1.2 follows when $\|\mathbf{a}\|_1$ is bounded. This completes the proof. ■

Proof of Theorem 1.2: After some calculations, we obtain that

$$\begin{aligned}
h^2(\Theta, \Theta^0) &= 2 \left(1 - \prod_{i=1}^n \frac{1}{(\sqrt{2\pi}\sigma)^k} \int \exp \left[-\frac{1}{4\sigma^2} (\|\mathbf{y}_i - \mathbf{a}_i^T \Theta\|^2 + \|\mathbf{y}_i - \mathbf{a}_i^T \Theta^0\|^2) \right] d\mathbf{y} \right) \\
&= 2 \left(1 - \prod_{i=1}^n \exp \left[-\frac{1}{8\sigma^2} \|\mathbf{a}_i^T (\Theta - \Theta^0)\|^2 \right] \right) \\
&= 2 \left(1 - \exp \left(-\frac{1}{8\sigma^2} \|\mathbf{A}(\Theta - \Theta^0)\|_F^2 \right) \right), \\
K(\Theta^0, \Theta) &= \frac{1}{2\sigma^2} \|\mathbf{A}(\Theta - \Theta^0)\|_F^2.
\end{aligned}$$

When $\epsilon < 1$,

$$\begin{aligned}
P(K(\Theta^0, \hat{\Theta}^{L_0}) \geq 4\epsilon^2) &= P \left(\frac{1}{8\sigma^2} \|\mathbf{A}(\hat{\Theta}^{L_0} - \Theta^0)\|_F^2 \geq \epsilon^2 \right) \\
&\leq P \left(\frac{1}{8\sigma^2} \|\mathbf{A}(\hat{\Theta}^{L_0} - \Theta^0)\|_F^2 \geq -\log \left(1 - \frac{\epsilon^2}{2} \right) \right) \\
&= P \left(2 \left(1 - \exp \left(-\frac{1}{8\sigma^2} \|\mathbf{A}(\hat{\Theta}^{L_0} - \Theta^0)\|_F^2 \right) \right) \geq \epsilon^2 \right) \\
&= P \left(h^2(\hat{\Theta}^{L_0}, \Theta^0) \geq \epsilon^2 \right).
\end{aligned}$$

For any $\epsilon \geq \epsilon_{n,p,k}$, it follows from Theorem 1 and Corollary 2 that

$$\begin{aligned}
EK(\Theta^0, \hat{\Theta}^{L_0}) &\leq EK(\Theta^0, \hat{\Theta}^{L_0}) I\{K(\Theta^0, \hat{\Theta}^{L_0}) \leq 4\epsilon^2\} + EK(\Theta^0, \hat{\Theta}^{L_0}) I\{K(\Theta^0, \hat{\Theta}^{L_0}) > 4\epsilon^2\} \\
&\leq 4\epsilon^2 + \left(EK^2(\Theta^0, \hat{\Theta}^{L_0}) \right)^{1/2} \left(P(K^2(\Theta^0, \hat{\Theta}^{L_0}) > 4\epsilon^2) \right)^{1/2}.
\end{aligned}$$

By the triangle inequality, $\|\mathbf{A}\Theta^0 - \mathbf{A}\hat{\Theta}^{L_0}\|_F - \|\epsilon\|_F \leq \|\mathbf{A}\Theta^0 + \epsilon - \mathbf{A}\hat{\Theta}^{L_0}\|_F$. Note that $\hat{\Theta}^{L_0}$ is a global minimizer of ((1.3)). Then $\|\mathbf{A}\Theta^0 + \epsilon - \mathbf{A}\hat{\Theta}^{L_0}\|_F \leq \|\epsilon\|_F$. Hence

$$K(\Theta^0, \hat{\Theta}^{L_0}) = \frac{1}{2\sigma^2} \|\mathbf{A}(\Theta^0 - \hat{\Theta}^{L_0})\|_F^2 \leq \frac{2}{\sigma^2} \|\epsilon\|_F^2.$$

Thus,

$$\begin{aligned} EK(\Theta^0, \hat{\Theta}^{L_0}) &\leq 4\epsilon^2 + \left(E \frac{4}{\sigma^4} \|\epsilon\|_F^4\right)^{1/2} P(K^2(\Theta^0, \hat{\Theta}^{L_0}) > 4\epsilon^2) \\ &\leq 4\epsilon^2 + 10 \exp(-c_1 n \epsilon^2 + \log \sqrt{3nk}). \end{aligned}$$

The results in Theorem 1.2 follow by letting $\epsilon = \epsilon_{n,p,k}$ and using the fact that $\log k \leq C_{p,k}^2$.

This completes the proof. ■

Proof of Theorem 1.3: Without loss of generality, assume $p \geq k$ and $n = p$. When $\sigma = O(1/\sqrt{p})$, by Theorem 2, we have

$$\begin{aligned} \|\hat{\Theta}^{L_0} - \Theta^0\|_F^2 &= 2\sigma^2 K(\Theta^0, \hat{\Theta}^{L_0}) = O_P\left(\frac{\epsilon_{n,p,k}^2}{p}\right) \\ &= O_P\left(\frac{C_{p,k}^2}{p^2} \log\left(\frac{\sqrt{p}}{C_{p,k}}\right)\right) \\ &= O_P\left(\frac{C_{p,k}^2}{p^2} \log\left(\frac{p^2}{C_{p,k}^2}\right)\right), \end{aligned} \quad (1.26)$$

where

$$C_{p,k} = O\left(\sqrt{\log(p)} \sqrt{(p+k)s_2^0 + s_1^0} + \sqrt{s_1^0 \log\left(e \frac{(p+k-r(\Theta^0))r(\Theta^0)}{s_1^0}\right)}\right). \quad (1.27)$$

(1.27) comes from the proof of Corollary 2 with M in **Assumption A** being $O(\sqrt{p})$. Thus,

$$\|\hat{\Theta}^{L_0} - \Theta^0\|_F^2 = O_P\left(C'_{p,k} \log\left(\frac{1}{C'_{p,k}}\right)\right)$$

with

$$C'_{p,k} = \frac{\log(p) \cdot [(p+k)s_2^0 + s_1^0] + s_1^0 \log\left(e \frac{(p+k-r(\Theta^0))r(\Theta^0)}{s_1^0}\right)}{p^2}.$$

This completes the proof. ■

Chapter 2

Boosting for High-Dimensional Additive Models with Group Variables

2.1 Introduction

Boosting was introduced in the machine learning literature by Schapire (1990), which came up with the first provable polynomial-time boosting algorithm. Freund (1995) developed a more efficient boosting algorithm for improving the accuracy of algorithms for learning binary concepts. His work provided an optimal upper bounds on the resources required for learning in Valiant's polynomial PAC (Probably Approximately Correct) learning framework. He also pointed out that the major drawback of this method is that its complexity is of order $O((\log \epsilon)^2/\epsilon)$ with ϵ being required accuracy. Freund and Schapire (1997) introduced the AdaBoost algorithm, which solved many of the practical difficulties of the earlier boosting algorithms. One of the main ideas of the algorithm is to maintain a distribution or set of weights over the training set. Initially, all weights are set equally, but on each iteration, the weights of incorrectly classified examples are increased so that the weaker learner is forced to focus on the hard examples in the training set (Freund and Schapire, 1999). After that, boosting gets more and more attention and has demonstrated great empirical success on a wide variety of especially high-dimensional prediction problems, including analysis of

microarray gene expression data (Dettling and Buhlmann (2003); Hothorn et al (2006); Li and Luan (2005)). Much has been written about the success of boosting as classifier. Friedman et al (2000) provided an elegant statistical justification of the boosting procedure and showed that boosting can be viewed as an approximation of additive modeling and maximum likelihood. This important insight opened a new perspective for using boosting in contexts other than classification. From the perspective of numerical optimization on function space, Friedman (2001) proposed a gradient descent boosting (GDB) procedure and demonstrated that such a procedure can be regarded as a stage-wise fitting of the additive models. Depending on the choice of the weak learner or base procedure, many different types of additive functions can be constructed. In the literature, many papers have explored the use of tree-based algorithms (Breiman et al, 1984) as weak learner, which renders final model as a linear combination of a large number of trees. It provides a natural nonparametric framework for modeling higher-order interactions.

Buhlmann and Yu (2003) proposed and studies the properties of boosting with L_2 loss for regression and classification. In particular, they proposed to use the component-wise linear least squares or component-wise univariate splines as based learners. Luan and Li (2008) proposed group additive regression models and a group L_2 Boosting (gL_2 Boost) procedure for identifying groups of genomic features that are related to clinical phenotypes. Yin et al (2012) considered a similar problem but in a nonparametric setting and presented a new method, called group sparse additive models (GroupSpAM), which generalized the l_1/l_2 norm as the sparsity-inducing penalty. Buhlmann (2006) proved that boosting with the squared error loss, L_2 Boosting, is consistent for very-high dimensional linear models, where the number of predictor variables is allowed to grow essentially as fast as $O(\exp(n))$ where n is the sample size under the assumption that the true underlying regression function is sparse in terms of the l_1 -norm of the regression coefficients. In this thesis, we extend the theoretical results of Buhlmann (2006) to the setting of high dimensional models with group variables and prove that the gL_2 Boosting in such a setting yields consistent estimates in high-dimensional context, when the number of the groups is allowed to grow essentially as fast as $O(\exp(n))$.

The rest of chapter is organized as follows. We first introduce the generalized additive models (GAM) with group variables. We then present the group gradient descent boosting (G-GDBoosting) procedure for fitting the GAM with smoothing spline as weak learners. Finally, we prove the consistency of the proposed method.

2.2 Models

Although regression models is important and useful, it often fails when the relationship between the predictors and the response is not linear or the effects of the predictors are not linear. In this section, we will review the generalized additive model (GAM) first then extend it to group variables.

2.2.1 Generalized Additive Models

In statistics, a generalized additive model (GAM) is assumed to be linear in terms of unknown smooth functions of the predictor variables. The challenge is to identify and characterize these nonlinear regression effects. Suppose we have n *i.i.d.* samples. The a traditional GAM has the form

$$g[\mu(X)] = \alpha + f_1(X_1) + f_2(X_2) + \cdots + f_p(X_p). \quad (2.1)$$

As usual, X_1, X_2, \dots, X_p are predictors; g is the link function, $\mu(X)$ is the expectation of the response; f_j 's are unspecified smooth functions, known as "nonparametric" functions. There are some examples:

- For Gaussian data, $g(\mu) = \mu$, the identity link. Then

$$E(Y|X) = \alpha + f_1(X_1) + f_2(X_2) + \cdots + f_p(X_p).$$

- For binary data, $g(\mu) = \text{logit}(\mu)$, the logit link. Then

$$\log\left(\frac{\mu(X)}{1 - \mu(X)}\right) = \alpha + f_1(X_1) + f_2(X_2) + \cdots + f_p(X_p).$$

- For Poisson count data, $g(\mu) = \log(\mu)$, the log link. Then

$$\log \mu(X) = \alpha + f_1(X_1) + f_2(X_2) + \cdots + f_p(X_p).$$

Different techniques are utilized to solve this kind of problem, including tree-based methods, spline smoother, kernel smoother and so on.

2.2.2 GAM with Group Variables

Consider X_k , $k = 1, \dots, p$, is not a single variable but a batch of variables. To distinguish with the traditional version, we use bold \mathbf{X}_k and $\mathbf{X}_k = (X_{k,1}, X_{k,2}, \dots, X_{k,p_k})$ be the collection of variables in the k th group. Consider a nonparametric GAM problem,

$$Y = F(\mathbf{X}) + \epsilon = f_1(\mathbf{X}_1) + f_2(\mathbf{X}_2) + \cdots + f_p(\mathbf{X}_p) + \epsilon, \quad (2.2)$$

where ϵ is the noise term and $f_k(\mathbf{X}_k)$ is the group effect as determined by the genomic data \mathbf{X}_k of the k th group. This model assumes additive effects of different groups on the response variable. If X_k is the vector of gene expression data of the p_k genes in the k th pathway, $f_k(\mathbf{X}_k)$ can be interpreted as the pathway activity (Luan and Li, 2008). When $p_k = 1$ for $k = 1, \dots, p$, it's reduced to the traditional GAM.

2.3 Methods

Wei and Li (2007) proposed a gradient descent boosting (GDB) procedure for fitting non-parametric pathways-based regression (NPR) models using regression trees as weak learners. Luan and Li (2008) pointed out that although trees are very flexible in modeling interac-

tions among variables, it's still difficult to interpret because the resulting model is a linear combination of many small trees. Throughout this chapter, we focus on the framework of GAM models using either linear regression or smoother as weak learners. Compared with tree-based methods, such a procedure leads to explicit expressions of the estimators. In this section, we will review G-GDB procedure (Luan and Li, 2008), GroupSpAM Yin et al (2012). Then, we propose a method using smoothing splines and reproducing kernel Hilbert space (RKHS) norms.

2.3.1 G-GDB

Luan and Li (2008) assume that the response Y is related to the predictors through an additive regression model:

$$Y = \sum_{k=1}^p f_k(\mathbf{X}_k) + \epsilon = \sum_{k=1}^p \sum_{l=1}^{p_k} \beta_{k,l} X_{k,l}.$$

For $k = 1, \dots, p$, let $\mathbf{H}_k = \mathbf{X}_k(\mathbf{X}_k' \mathbf{X}_k)^{-1} \mathbf{X}_k'$. Then, the algorithm is summarized below:

Algorithm 2.1:

1. Initialization. Let $\hat{U}^{(1)} = (Y_1, \dots, Y_n)^T$, $\hat{F}_0 = 0$ and $\hat{\mathbf{A}}_0 = 0$.
2. At each step,
 - (a) Compute

$$\hat{s}_m = \operatorname{argmin}_{1 \leq k \leq p} (\hat{U}^{(m)} - \mathbf{H}_k \hat{U}_m)^T (\hat{U}^{(m)} - \mathbf{H}_k \hat{U}_m),$$

$$\hat{\mathbf{A}}_m = \mathbf{I} - (\mathbf{I} - \rho \mathbf{H}_{\hat{s}_0}) \cdots (\mathbf{I} - \rho \mathbf{H}_{\hat{s}_{m-1}}),$$

where \mathbf{I} is the identity matrix of order n and ρ is the learning rate.

- (b) Let $\hat{F}_m = \hat{F}_{m-1} + \rho \mathbf{H}_{\hat{s}_m} \hat{U}_m$, $\hat{U}^{(m+1)} = Y - \hat{F}_m$ and $m = m + 1$. Go back to the last step until $m = M$.

When the inverse of matrix $\mathbf{X}'_k \mathbf{X}_k$ is singular or near singular, the above algorithm cannot be applied. To solve this problem, they proposed to apply a penalized least square regression as weak learners. Accordingly, they redefined \mathbf{H}_k as

$$\mathbf{H}_k^{(\lambda)} = \mathbf{X}_k (\mathbf{X}'_k \mathbf{X}_k + \lambda \mathbf{I})^{-1} \mathbf{X}'_k, \text{ for } k = 1, \dots, p,$$

where λ is a tuning parameter for L_2 -penalized estimation. The G-GDBoosting algorithm remains the same with \mathbf{H}_k being replaced by $\mathbf{H}_k^{(\lambda)}$.

2.3.2 GroupSpAM

Yin et al (2012) considered this problem in a nonparametric setting and utilized the group structure. They assume

$$E(Y|\mathbf{X}) = \sum_{k=1}^p \sum_{j=1}^{p_k} h_{k,j}(X_{k,j}).$$

Then, they generalized a l_1/l_2 norm to Hilbert spaces as the sparsity-inducing penalty and proposed an efficient block coordinate descent algorithm. Their optimization problem of GroupSpAM in the population setting is formulated as

$$\min_f \frac{1}{2} E \left[\left(Y - \sum_{k=1}^p \sum_{j=1}^{p_k} h_{k,j}(X_{k,j}) \right)^2 \right] + \lambda \sum_{k=1}^p d_k J_k(h),$$

where $J_k(h) = \sqrt{\frac{1}{n} \sum_{j=1}^{p_k} \|h_{k,j}\|^2} = \sqrt{\frac{1}{n} \sum_{j=1}^{p_k} E h_{k,j}^2}$ is the penalty.

Algorithm 2.2

1. Initialization. $\hat{h}_{k,j} = 0$ for any (k, j) ; pre-compute smoother matrices $\mathbf{S}_{k,j}$ for any (k, j) .
2. Cycle through $k = 1, \dots, p$:
 - (a) Compute the partial residual $\hat{\mathbf{q}}_k^{(m)} = y - \sum_{l \neq k} \sum_{j=1}^{p_l} \hat{h}_{l,j}^{(m)}$ and estimate the

penalty $J_k(h)$ by

$$\hat{w}_k = \sqrt{\frac{1}{n} \sum_{j=1}^{p_k} (\mathbf{S}_{k,j} \hat{\mathbf{q}}_k^{(m)})' \mathbf{S}_{k,j} \hat{\mathbf{q}}_k^{(m)}}.$$

If $\hat{w}_k \leq \lambda \sqrt{d_k}$, then $\hat{h}_{k,j}^{(m+1)} = 0$ for all $j = 1, \dots, p_k$. Otherwise,

$$\left(\hat{h}_{k,j}^{(m+1)} \right)_{1 \leq j \leq p_k} = \left(\hat{\mathbf{J}} + \frac{n \lambda d_k}{\sqrt{\sum_{j=1}^{p_k} \|\hat{h}_{k,j}^{(m)}\|^2}} \mathbf{I} \right)^{-1} \hat{\mathbf{Q}}_k^{(m)} \hat{\mathbf{q}}_k^{(m)},$$

where

$$\hat{\mathbf{Q}}_k^{(m)} = \begin{bmatrix} \mathbf{S}_{k,1} \hat{\mathbf{q}}_k^{(m)} \\ \mathbf{S}_{k,2} \hat{\mathbf{q}}_k^{(m)} \\ \vdots \\ \mathbf{S}_{k,p_k} \hat{\mathbf{q}}_k^{(m)} \end{bmatrix}, \quad \hat{\mathbf{J}}_k^{(m)} = \begin{bmatrix} \mathbf{I} & \mathbf{S}_{k,1} \hat{\mathbf{q}}_k^{(m)} & \cdots & \mathbf{S}_{k,1} \hat{\mathbf{q}}_k^{(m)} \\ \mathbf{S}_{k,2} \hat{\mathbf{q}}_k^{(m)} & \mathbf{I} & \cdots & \mathbf{S}_{k,2} \hat{\mathbf{q}}_k^{(m)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{S}_{k,p_k} \hat{\mathbf{q}}_k^{(m)} & \mathbf{S}_{k,p_k} \hat{\mathbf{q}}_k^{(m)} & \cdots & \mathbf{I} \end{bmatrix}.$$

3. Center each $\left(\hat{h}_{k,j}^{(m+1)} \right)_{1 \leq j \leq p_k}$ by subtracting its mean. Go back to the last step until $m = M$ or convergence.

2.3.3 The proposed algorithm

The above algorithm is not applicable when within each group, f_k is not additive. To solve this issue, we propose the following algorithm. Assume the functional space $\mathcal{F} = \oplus \mathcal{H}_k$ with \mathcal{H}_k being a subspace corresponding to f_k in (2.2). Let $E(Y|\mathbf{X}) = F(\mathbf{X}) = \sum_{k=1}^p f_k(\mathbf{X}_k)$. Denote the norm in the RKHS \mathcal{H}_k by $J(\cdot)$. A traditional smoothing spline type method finds $F \in \mathcal{F}$ to minimize

$$\frac{1}{n} \sum_{i=1}^n \{y_i - F(\mathbf{x}_i)\}^2 + \lambda_n \sum_{k=1}^p J^2(P^k F), \quad (2.3)$$

where $P^k F$ is the orthogonal projection of F onto \mathcal{H}_k and $\lambda_n \geq 0$. If λ_n is large, then $\|P^k F\|$ tends to be zero. It's well known that the solution F has the form $F(x) = \sum_{i=1}^n c_i R_\alpha(\mathbf{x}_i, \mathbf{x})$, where $c = (c_1, \dots, c_n)^T \in R^n$ and $R_\alpha = \sum_{k=1}^p \alpha_k R_k$, with R_k being the reproducing kernel of \mathcal{H}_k . However, when p is large, it motivates us to reduce the number of parameters at each step. Our proposed boosting procedure is to iteratively fit the residual using the covariates in each of the p groups and at each step, select an group that provides the best fit to the residuals as measured by the residuals sum of squares.

Before introducing the algorithm, we define the following notation:

$$\begin{aligned} \hat{U}^{(m)} &= (\hat{u}_1^{(m)}, \dots, \hat{u}_n^{(m)})^T, \text{ residual at } m\text{th step} \\ \mathbf{X}'_k &= (\mathbf{X}_{1,k}, \dots, \mathbf{X}_{n,k}), \text{ a matrix of } p_k \text{ by } n, k = 1, \dots, p, \\ R_k &= (R(\mathbf{X}_{i,k}, \mathbf{X}_{j,k}))_{n \times n}, \text{ a matrix of } n \text{ by } n, i, j = 1, \dots, n, k = 1, \dots, p, \\ \mathbf{Y} &= (Y_1, \dots, Y_n). \end{aligned}$$

The proposed Boosting algorithm is as follows:

1. Initialization. Let $\hat{U}^{(1)} = (Y_1, \dots, Y_n)^T$, $\hat{F}_0 = 0$.
2. For $m = 1$ to M :
 - (a) Compute

$$\begin{aligned} \hat{\alpha}_k^{(m)} &= (R_k + \lambda_n I)^{-1} \hat{U}^{(m)}, \text{ for } k = 1, \dots, p \\ \hat{s}_m &= \operatorname{argmin}_{1 \leq k \leq p} (\hat{U}^{(m)} - R_k \hat{\alpha}_k^{(m)})^T (\hat{U}^{(m)} - R_k \hat{\alpha}_k^{(m)}). \end{aligned}$$

- (b) Let $\hat{F}_m = \hat{F}_{m-1} + \rho R_{\hat{s}_m} \hat{\alpha}_{\hat{s}_m}^{(m)}$, $\hat{U}^{(m+1)} = Y - \hat{F}_m$ and $m = m + 1$. Go back to the last step until $m = M$.

2.4 Consistency of Boosting

In this section, we present the consistency of the group- L_2 -boosting in nonparametric models where the number of predictors is allowed to grow very fast as the sample size n increases.

Assume the data (Y_i, \mathbf{X}_i) , $i = 1, \dots, n$ are i.i.d. sample generated from the model (2.2). The goal is to learn the unknown function f through a generalized dictionary $\mathcal{D}_p = \{\mathcal{H}_k, k = 1, \dots, p\}$, composed of Hilbert spaces \mathcal{H}_k , $k = 1, \dots, p$ endowed with a common inner product $\langle g, g' \rangle = E(g(\mathbf{X})g'(\mathbf{X}))$ with $g, g' \in \mathcal{H}_k$. Define the linear space spanned by \mathcal{D}_p as $\mathcal{F}(\mathcal{D}_p) = \{F : F = \sum_{k=1}^p g_k \text{ where } g_k \in \mathcal{H}_k, \mathcal{H}_k \in \mathcal{D}_p\}$ and L_1 -ball as $\mathcal{F}(\mathcal{D}_p, W) = \{F : F = \sum_{k=1}^p g_k \text{ where } g_k \in \mathcal{H}_k, \mathcal{H}_k \in \mathcal{D}_p, \sum_{k=1}^p J(g_k) \leq W\}$.

Then, we introduce the sample version norm $\|\cdot\|_n$ and inner product $\langle \cdot, \cdot \rangle_n$ in R^n as

$$\|g\|_n^2 = \frac{1}{n} \sum_{i=1}^n g^2(\mathbf{X}_i), \quad \langle g, g' \rangle_n = \frac{1}{n} \sum_{i=1}^n g(\mathbf{X}_i)g'(\mathbf{X}_i).$$

Similarly, the population version norm and inner product are

$$\|g\|^2 = E g^2(X), \quad \langle g, g' \rangle = E(g(X)g'(X)).$$

And define

$$P_{\mathcal{H}}^{(n, \lambda)} g = \operatorname{argmin}_{g' \in \mathcal{H}} \|g - g'\|_n^2 + \lambda_n^2 J^2(g'),$$

$$P_{\mathcal{H}} g = \operatorname{argmin}_{g' \in \mathcal{H}} \|g - g'\|^2.$$

We make the following assumptions:

A1: Let $F \in \mathcal{F}(\mathcal{D}_{p_n}, W) = \{F : F_n = \sum_{k=1}^{p_n} g_k, g_k \in \mathcal{H}_k, \sum J(g_k) \leq W\}$, where the convex hull of \mathcal{H}_k is \mathcal{H}_k itself and for all $g \in \mathcal{H}_k$ satisfying $J(g) \leq W$, g is bounded by T , that is, $0 \leq g \leq T$ with T independent of n . In addition, $p_n = O(\exp(C_1 n^{1-\xi}))$, for some $0 < \xi < 1$.

A2: $H_{\infty}(\delta, \mathcal{H}_k(1)) \leq A\delta^{-r}$ with $0 < r < 2$ for all $k = 1, \dots, p$, where $\mathcal{H}_k(1) = \{g : g \in \mathcal{H}_k, \text{ and } J(g) \leq 1\}$.

A3: For any $g \in \mathcal{H}_j$, $P_{\mathcal{H}_k} g \in \mathcal{H}_k$ and $\sup_{k,j} J(P_{\mathcal{H}_k} g)/J(g) \leq K$. (For simplicity, we use the

$J(g)$ to represent the penalty in functional class that g belongs to)

A4: $E\epsilon^s \leq \infty$ for some $s > 4/\xi$

In order to analyze \hat{U}_m , we first define a population version group- L_2 -boosting algorithm, where the sequence $\{F_m\}_{m=0}^\infty$ satisfies the following condition:

$$F_0 = 0, F_m = F_{m-1} + h_m, \quad (2.4)$$

where $h_m = P_{\tilde{S}_m}(F - F_{m-1})$ satisfies $\|h_m\| \geq t_m \sup_j \|P_{S_j}(F - F_{m-1})\|$. The following lemma establishes the convergence rate of the population version algorithm, which generalize Theorem 5.1 in Temlyakov (2000).

Lemma 2.1 *Suppose $F \in \mathcal{F}(\mathcal{D}_p, T)$. Then, $\|F - F_m\| \leq T(1 + \sum_{k=1}^m t_k^2)^{-t_m/[2(2+t_m)]}$.*

Proof: Let $a_m = \|F - F_m\|^2$, $y_m = \|h_m\|$, $b_0 = T$, and $b_m = b_{m-1} + y_m$. Because $\|F - F_m\|^2 = \|F - F_{m-1}\|^2 - \|P_{S_m^*}(F - F_{m-1})\|^2$, we have $a_m = a_{m-1} - y_m^2$. In addition, note that

$$\|F - F_m\| = \left\| \sum_{j=1}^{\infty} g_j - \sum_{k=1}^m h_k \right\| \leq \sum_{j=1}^{\infty} \|g_j\| + \sum_{k=1}^m \|h_k\| = b_m$$

and

$$\begin{aligned} \|F - F_m\|^2 &= \sum_{j=1}^{\infty} \langle g_j, F - F_m \rangle - \sum_{k=1}^m \langle h_k, F - F_m \rangle \\ &\leq \sum_{j=1}^{\infty} \sup_{S \in \mathcal{D}} \|P_S(F - F_m)\| \|g_j\| + \sum_{k=1}^m \sup_{S \in \mathcal{D}} \|P_S(F - F_m)\| \|h_k\| \\ &= \sup_{S \in \mathcal{D}} \|P_S(F - F_m)\| b_m. \end{aligned}$$

Consequently, $\sup_{S \in \mathcal{D}} \|P_S(F - F_{m-1})\| \geq \|f - f_{m-1}\|^2 / b_{m-1}$, which implies $y_m \geq t_m a_{m-1} / b_m$.

Therefore, the sequence $\{a_m\}$, $\{b_m\}$, and $\{y_m\}$ can be characterized by the following equa-

tions:

$$\begin{aligned} b_m &= b_{m-1} + y_m, \\ a_m &= a_{m-1} - y_m^2, \\ y_m &\geq t_m a_{m-1} / b_m. \end{aligned}$$

The convergence rate of $\{a_m\}$ then follows as in the proof of Theorem 5.1 in Temlyakov (2000). ■

Now we define the sample version of the weak greedy algorithm, where the sequence $\{\hat{f}_m\}_{m=1}^\infty$ satisfy the following conditions:

$$\hat{F}_0 = 0, \hat{F}_m = \hat{F}_{m-1} + \hat{h}_m, \tag{2.5}$$

where $\hat{h}_m = P_{\hat{S}_m}^{(n,\lambda)}(F + \epsilon - \hat{F}_{m-1})$ satisfies $\|\hat{h}_m\|_n \geq t_m \sup_j \|P_{S_j}^{(n,\lambda)}(F + \epsilon - F_{m-1})\|_n$.

We will apply it to a semipopulation version \tilde{F}_m which extends its originally definition made by Buhlmann (2006). That is, the sequence $\{\tilde{F}_m\}_{m=0}^\infty$ satisfy the following conditions:

$$\tilde{F}_0 = 0, \tilde{F}_m = \tilde{F}_{m-1} + P_{\hat{s}_m}(F - \tilde{F}_{m-1})$$

where \hat{s}_m is selected from the sample version in the proposed algorithm above.

In following, we prove four lemmas which show the consistency of the semipopulation version and the convergence rate of the difference between semipopulation version and sample version. For arbitrary, fixed step-size $0 < \rho \leq 1$, we can then use exactly the same reasoning in Section 6.3 in Buhlmann (2006) to show its consistency.

Lemma 2.2 *Under the assumptions A1, A2 and A3, A4*

$$\begin{aligned}\delta_{n,1} &= \sup_{\mathcal{H}_j, \mathcal{H}_k \in \mathcal{D}_{p_n}} \sup_{g \in \mathcal{H}_j, g' \in \mathcal{H}_k} |\langle g, g' \rangle_n - \langle g, g' \rangle| / J(g)J(g') = O_p(n^{-\xi/2}) \\ \delta_{n,2} &= \sup_{\mathcal{H}_k \in \mathcal{D}_{p_n}} \sup_{g \in \mathcal{H}_k} |\langle \epsilon, g \rangle_n| / J(g) = O_p(n^{-\xi/2})\end{aligned}$$

Proof: Without loss of generality, assume $J(g) = J(g') = 1$. Consider a class $\mathcal{H}_{j,k} = \{h = gg' : g \in \mathcal{H}_j(1), g' \in \mathcal{H}_k(1)\}$. Note that for any $h_1 = g_1g'_1$ and $h_2 = g_2g'_2$, $\|h_1 - h_2\|_\infty \leq \|g_1(g'_1 - g'_2)\|_\infty + \|g'_2(g_1 - g_2)\|_\infty \leq T\|g'_1 - g'_2\|_\infty + T\|(g_1 - g_2)\|_\infty$. Therefore, one can show that $H_\infty(2T\delta, \mathcal{H}_{j,k}) \leq 2A\delta^{-r}$.

Let $M = \sqrt{nt}$, $t = Kn^{-\xi/2}$, $\psi(M, n) = Mn^{1/2} \frac{M/n^{1/2}}{2(1+M/3n^{1/2})}$ with sufficiently large K and $0 < \xi < 1$. Define t_0 by $H_\infty(t_0, \mathcal{H}_{j,k}) = 1/4 h\psi(M, n)$. Then

$$At_0^{-r} \geq 1/4 hMn^{1/2} \frac{M/n^{1/2}}{2(1+M/3n^{1/2})} \geq 1/4 hMn^{1/2} \frac{M/n^{1/2}}{2(1+1)}.$$

And

$$M > 2^8 h^{-3/2} I_\infty\left(\frac{hM}{64n^{1/2}}, C(h, A, r)M^{-2/r}\right) > 2^8 h^{-3/2} I_\infty\left(\frac{hM}{64n^{1/2}}, t_0\right)$$

with a given $0 < h < 1$.

It follows from Theorem 2.1 of Kenneth (1984)

$$\begin{aligned}& P^* \left[\sup_{g \in \mathcal{H}_j(1), g' \in \mathcal{H}_k(1)} |\langle g, g' \rangle_n - \langle g, g' \rangle| > T^2 t \right] \\ &= P^* \left[\sup_{gg' \in \mathcal{H}_{j,k}} \left| \frac{1}{n} \sum_{i=1}^n g(X_i)g'(X_i) - Eg(X)g'(X) \right| > T^2 t \right] \\ &\leq 5 \exp(-(1-h)\psi(M, n)) \\ &= 5 \exp\left(- (1-h) \frac{nt^2}{2(1+t/3)}\right),\end{aligned}$$

which yields

$$\begin{aligned} & P^* \left[\sup_{\mathcal{H}_j, \mathcal{H}_k \in \mathcal{D}_{p_n}} \sup_{g \in \mathcal{H}_j(1), g' \in \mathcal{H}_k(1)} |\langle g, g' \rangle_n - \langle g, g' \rangle| > T^2 t \right] \\ & \leq p_n^2 \times 5 \exp(-2(1-h) \frac{nt^2}{2(1+t/3)}) = o(1) \end{aligned}$$

with sufficiently large K . The desired result follows.

Consider a truncated version of ϵ defined as $\tilde{\epsilon} = \text{sign}(\epsilon) \min\{|\epsilon|, M_n\}$. For any $S \in \mathcal{D}_{p_n}$,

$$\begin{aligned} & P \left(\sup_{1 \leq k \leq p_n} \sup_{g \in \mathcal{H}_k(1)} |\langle g, \epsilon \rangle_n| > 3t \right) \\ & \leq P \left(\sup_{1 \leq k \leq p_n} \sup_{g \in \mathcal{H}_k(1)} |\langle g, \tilde{\epsilon} \rangle_n - \langle g, \epsilon \rangle_n| > t \right) + P \left(\sup_{1 \leq k \leq p_n} \sup_{g \in \mathcal{H}_k(1)} |\langle g, \tilde{\epsilon} \rangle_n - \langle g, \epsilon \rangle_n| > t \right) \\ & \quad + P \left(\sup_{1 \leq k \leq p_n} \sup_{g \in \mathcal{H}_k(1)} |\langle g, \tilde{\epsilon} \rangle| > t \right) \\ & = I + II + III. \end{aligned}$$

To bound I , let $Z_i = (X_i, \tilde{\epsilon}_i)$ and let $h_g(Z_i) = g(X_i)\tilde{\epsilon}_i/M_n$. Define $\mathcal{G}_k = \{h_g : g \in \mathcal{H}_k(1)\}$. Note that $\|h_g - h_{g'}\|_\infty \leq \|g - g'\|_\infty$. Similarly to the proof of (i), $H_\infty(\delta, \mathcal{G}_k) \leq A\delta^{-r}$, and $\sup_{h_g \in \mathcal{G}_k} \text{Var}(h_g(Z_i)/T) \leq E|\tilde{\epsilon}_i|^2/M_n^2 \triangleq \alpha_n = O(M_n^{-2})$. Let $M = \sqrt{nt}/M_n, t = Kn^{-\xi/2}, M_n = n^{\xi/4}, \psi(M, n, \alpha_n) = Mn^{1/2} \frac{M/n^{1/2}\alpha_n}{2(1+M/3n^{1/2}\alpha_n)}$ with sufficiently large K . Define t_0 by $H_\infty(t_0, \mathcal{G}_k) = 1/4 h\psi(M, n, \alpha_n)$. Then

$$At_0^{-r} \geq 1/4 hMn^{1/2} \frac{M/n^{1/2}}{2(\alpha_n + M/3n^{1/2})} \geq 1/4 hMn^{1/2} \frac{M/n^{1/2}}{4M/n^{1/2}}.$$

And

$$M \geq 2^8 h^{-3/2} I_\infty \left(\frac{hM}{64n^{1/2}}, C(h, A, r)(M\sqrt{n})^{-1/r} \right) \text{ with a given } 0 < h < 1.$$

By Theorem 2.1 of Kenneth (1984)

$$\begin{aligned}
P^*[\sup_{g \in \mathcal{H}_j(1)} |\langle g, \tilde{\epsilon} \rangle_n - \langle g, \tilde{\epsilon} \rangle| > Tt] &= P^*[\sup_{h_g \in \mathcal{G}_k} |1/n \sum_{i=1}^n h_g(Z_i) - Eh_g(Z_i)| > Tt/M_n] \\
&\leq 5 \exp(-(1-h)\psi(M, n, \alpha_n)) \\
&= 5 \exp(-(1-h) \frac{n(t/M_n)^2}{2(\alpha_n + t/3M_n)}),
\end{aligned}$$

which leads to

$$P^*[\sup_{1 \leq k \leq p_n} \sup_{g \in \mathcal{H}_j(1)} |\langle g, \tilde{\epsilon} \rangle_n - \langle g, \tilde{\epsilon} \rangle| > Tt] \leq 5p_n \exp(-(1-h) \frac{n(t/M_n)^2}{2(\alpha_n + t/3M_n)})$$

for $t = Kn^{-\xi/2}$ and $M_n = n^{\xi/4}$ with sufficiently large K .

To bound *II*, note that

$$\begin{aligned}
P(\sup_{1 \leq k \leq p_n} \sup_{g \in \mathcal{H}_j(1)} |\langle g, \tilde{\epsilon} \rangle_n - \langle g, \epsilon \rangle_n| > t) &= P(\sup_{1 \leq k \leq p_n} \sup_{g \in \mathcal{H}_j(1)} |\langle g, \tilde{\epsilon} - \epsilon \rangle_n| > t) \\
&\leq P(\sup_i |\epsilon_i| > M_n) \leq nP(|\epsilon| > M_n) \leq nM_n^{-s} E|\epsilon|^s = O(n^{1-s\xi/4}) = o(1).
\end{aligned}$$

To bound *III*, note that

$$\begin{aligned}
\sup_{g \in \mathcal{H}_k(1)} |\langle g, \tilde{\epsilon} \rangle| &= \sup_{g \in \mathcal{H}_k(1)} |Eg(X)(\tilde{\epsilon} - \epsilon)| \\
&\leq TE|\tilde{\epsilon} - \epsilon| \leq E|\epsilon|I[|\epsilon| > M_n] \\
&\leq 2(E(|\epsilon|^s))^{1/s} (P|\epsilon| > M_n)^{1-1/s} = O(M_n^{1-s}).
\end{aligned}$$

This implies that *III* = 0 for sufficiently large n . ■

Lemma 2.3 *Let $\lambda = O(\max(\delta_{n,1}, \delta_{n,2})^{1/2})$,*

$$\begin{aligned}\delta_{n,3} &= \sup_{j,k} \sup_{g \in \mathcal{H}_j(1)} \|P_{\mathcal{H}_k}^{(n,\lambda)} g - P_{\mathcal{H}_k} g\|_n = O_P(\lambda + \lambda^{-1} \delta_{n,1}), \\ \sup_{\mathcal{H} \in \mathcal{D}_{pn}} \|P_{\mathcal{H}}^{(n,\lambda)} \epsilon\| &= O_P(\lambda^{-1} \delta_{n,2}), \\ C_n &= \sup_{j,k} \sup_{g \in \mathcal{H}_j(1)} J(P_{\mathcal{H}_k}^{(n,\lambda)} g) + \sup_{\mathcal{H} \in \mathcal{D}_{pn}} J(P_{\mathcal{H}_k}^{(n,\lambda)} \epsilon) = O_P(1).\end{aligned}$$

Proof:

$$\begin{aligned}P_{\mathcal{H}_k}^{(n,\lambda)} g &= \operatorname{argmin}_{g' \in \mathcal{H}_k} \|g - g'\|_n^2 + \lambda^2 J^2(g') \\ &= \operatorname{argmin}_{g' \in \mathcal{H}_k} \|g - P_{\mathcal{H}_k} g + P_{\mathcal{H}_k} g - g'\|_n^2 + \lambda^2 J^2(g') \\ &= \operatorname{argmin}_{g' \in \mathcal{H}_k} \|P_{\mathcal{H}_k} g - g'\|_n^2 + 2\langle g - P_{\mathcal{H}_k} g, P_{\mathcal{H}_k} g - g' \rangle_n + \lambda^2 J^2(g').\end{aligned}$$

Then

$$\begin{aligned}\|P_{\mathcal{H}_k} g - P_{\mathcal{H}_k}^{(n,\lambda)} g\|_n^2 + 2\langle g - P_{\mathcal{H}_k} g, P_{\mathcal{H}_k} g - P_{\mathcal{H}_k}^{(n,\lambda)} g \rangle_n + \lambda^2 J^2(P_{\mathcal{H}_k}^{(n,\lambda)} g) \\ \leq \lambda^2 J^2(P_{\mathcal{H}_k} g) \leq K^2 \lambda^2.\end{aligned}\tag{2.6}$$

Using Lemma 2.2 and $\langle g - P_{\mathcal{H}_k} g, f \rangle = 0$ for any $f \in \mathcal{H}_k$, the second term in the left side of inequality above is

$$\begin{aligned}|\langle g - P_{\mathcal{H}_k} g, P_{\mathcal{H}_k} g - P_{\mathcal{H}_k}^{(n,\lambda)} g \rangle_n| &\leq \delta_{n,1} J(g - P_{\mathcal{H}_k} g) J(P_{\mathcal{H}_k} g - P_{\mathcal{H}_k}^{(n,\lambda)} g) \\ &\leq \delta_{n,1} (1 + K) (K + J(P_{\mathcal{H}_k}^{(n,\lambda)} g))\end{aligned}\tag{2.7}$$

Combining (2.6) and (2.7), we get

$$0 \leq \|P_{\mathcal{H}_k} g - P_{\mathcal{H}_k}^{(n,\lambda)} g\|_n^2 \leq K^2 \lambda^2 + \delta_{n,1} (1 + K) (K + J(P_{\mathcal{H}_k}^{(n,\lambda)} g)) - \lambda^2 J^2(P_{\mathcal{H}_k}^{(n,\lambda)} g).\tag{2.8}$$

Solving (2.8) yields

$$J(P_{\mathcal{H}_k}^{(n,\lambda)} g) \leq \frac{(1+K) + \sqrt{(1+K)^2 - 4\lambda^2(K^2\lambda^2\delta_{n,1}^{-2} + (1+K)K\delta_{n,1}^{-1})}}{2\lambda^2\delta_{n,1}^{-1}},$$

and

$$\|P_{\mathcal{H}_k} g - P_{\mathcal{H}_k}^{(n,\lambda)} g\|_n^2 \leq K^2\lambda^2 + K(1+K)\delta_{n,1} + \lambda^{-2}\delta_{n,1}^2(K+1)^2.$$

When $\lambda = O(\delta_{n,1}^{1/2})$,

$$\begin{aligned} J(P_{\mathcal{H}_k}^{(n,\lambda)} g) &= O_p(1), \\ \|P_{\mathcal{H}_k} g - P_{\mathcal{H}_k}^{(n,\lambda)} g\|_n &= O_p(\lambda + \lambda^{-1}\delta_{n,1}). \end{aligned}$$

The desired result follows.

Similarly,

$$\begin{aligned} P_{\mathcal{H}_k}^{(n,\lambda)} \epsilon &= \operatorname{argmin}_{g' \in \mathcal{H}_k} \|\epsilon - g'\|_n^2 + \lambda^2 J^2(g') \\ &= \operatorname{argmin}_{g' \in \mathcal{H}_k} \|\epsilon\|_n^2 + \|g'\|_n^2 - 2\langle \epsilon, g' \rangle_n + \lambda^2 J^2(g') \\ &= \operatorname{argmin}_{g' \in \mathcal{H}_k} \|g'\|_n^2 - 2\langle \epsilon, g' \rangle_n + \lambda^2 J^2(g'). \end{aligned}$$

Again using Lemma 2.2, note that $\langle \epsilon, f \rangle = 0$ for any $f \in \mathcal{H}_k$, we have

$$\|P_{\mathcal{H}_k}^{(n,\lambda)} \epsilon\|_n^2 \leq \delta_{n,2} J(P_{\mathcal{H}_k}^{(n,\lambda)} \epsilon) - \lambda^2 J^2(P_{\mathcal{H}_k}^{(n,\lambda)} \epsilon).$$

Solving it yields

$$\begin{aligned} J(P_{\mathcal{H}_k}^{(n,\lambda)} \epsilon) &= O_p(\delta_{n,2}\lambda^{-2}) = O_p(1), \\ \|P_{\mathcal{H}_k}^{(n,\lambda)} \epsilon\|_n &= O_p(\delta_{n,2}\lambda^{-1}). \end{aligned}$$

■

Lemma 2.4 *Under the assumptions A1, A2, A3 and A4 and λ defined above*

- (i) $\sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|P_{\mathcal{H}}^{(n,\lambda)} \hat{U}_m - P_{\mathcal{H}} \hat{U}_m\|_n \leq (W+1)(C_n+1)^m \delta_{n,3}$
- (ii) $\sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|P_{\mathcal{H}}^{(n,\lambda)} \hat{U}_m - P_{\mathcal{H}} \hat{U}_m\| \leq (W+1)(C_n+1)^m (2(1+K)\delta_{n,1}^{1/2} + \delta_{n,3})$

Proof: (i)

$$\begin{aligned}
& \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|P_{\mathcal{H}}^{(n,\lambda)} \hat{U}_m - P_{\mathcal{H}} \hat{U}_m\|_n \\
&= \|P_{\mathcal{H}}^{(n,\lambda)} (F - \hat{F}_m + \epsilon) - P_{\mathcal{H}} (F - \hat{F}_m + \epsilon)\|_n \\
&\leq \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|P_{\mathcal{H}}^{(n,\lambda)} F - P_{\mathcal{H}} F\|_n + \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|P_{\mathcal{H}}^{(n,\lambda)} \hat{F}_m - P_{\mathcal{H}} \hat{F}_m\|_n + \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|P_{\mathcal{H}}^{(n,\lambda)} \epsilon - P_{\mathcal{H}} \epsilon\|_n \\
&= I + II + III.
\end{aligned}$$

As for I , by Lemma 2.3,

$$\begin{aligned}
I &\leq \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \sum_{j=1}^{p_n} \|P_{\mathcal{H}}^{(n,\lambda)} g_j - P_{\mathcal{H}} g_j\|_n \\
&\leq \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \sum_{j=1}^{p_n} J(g_j) \|P_{\mathcal{H}}^{(n,\lambda)} \frac{g_j}{J(g_j)} - P_{\mathcal{H}} \frac{g_j}{J(g_j)}\|_n \\
&\leq \sum_{j=1}^{p_n} J(g_j) \delta_{n,3} \leq W \delta_{n,3}.
\end{aligned}$$

As for II ,

$$\begin{aligned}
II &\leq \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \sum_{j=1}^m \|P_{\mathcal{H}}^{(n,\lambda)} \hat{h}_j - P_{\mathcal{H}} \hat{h}_j\|_n \\
&\leq \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \sum_{j=1}^m J(\hat{h}_j) \|P_{\mathcal{H}}^{(n,\lambda)} \frac{\hat{h}_j}{J(\hat{h}_j)} - P_{\mathcal{H}} \frac{\hat{h}_j}{J(\hat{h}_j)}\|_n \\
&\leq \sum_{j=1}^m J(\hat{h}_j) \delta_{n,3}.
\end{aligned}$$

Note that

$$\begin{aligned}
J(\hat{h}_j) &= J(P_{\hat{S}_j}^{(n,\lambda)}(F + \epsilon - \hat{F}_{j-1})) \\
&\leq \sup_{\mathcal{H} \in \mathcal{D}_n} J(P_{\mathcal{H}}^{(n,\lambda)}(F + \epsilon - \hat{F}_{j-1})) \\
&\leq \sup_{\mathcal{H} \in \mathcal{D}_n} J(P_{\mathcal{H}}^{(n,\lambda)}F) + \sup_{\mathcal{H} \in \mathcal{D}_n} J(P_{\mathcal{H}}^{(n,\lambda)}\epsilon) + \sup_{\mathcal{H} \in \mathcal{D}_n} J(P_{\mathcal{H}}^{(n,\lambda)}\hat{F}_{j-1}) \\
&\leq (W+1)C_n + \sum_{k=1}^{j-1} J(\hat{h}_{k-1}) \sup_{\mathcal{H} \in \mathcal{D}_n} J(P_{\mathcal{H}}^{(n,\lambda)} \frac{\hat{h}_{k-1}}{J(\hat{h}_{k-1})}) \\
&\leq (W+1)C_n + \sum_{k=1}^{j-1} C_n J(\hat{h}_{k-1}). \tag{2.9}
\end{aligned}$$

Let $a_j = J(\hat{h}_j)$, $a_1 = (W+1)C_n$. They have such relationship $a_j \leq a_1 + \sum_{i=1}^{j-1} C_n a_i$. Then $a_j \leq (C_n + 1)^{j-1} a_1$, i.e. $J(\hat{h}_j) \leq (C_n + 1)^{j-1} (W+1)C_n$. Hence,

$$II \leq \sum_{j=1}^m (C_n + 1)^{j-1} (W+1)C_n \delta_{n,3} \leq (W+1)\delta_{n,3}((C_n + 1)^m - 1).$$

As for *III*, by Lemma 2.3

$$III = \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|P_{\mathcal{H}}^{(n,\lambda)}\epsilon\|_n \leq \delta_{n,3}. \tag{2.10}$$

(ii)

$$\begin{aligned}
&\sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|P_{\mathcal{H}}^{(n,\lambda)}\hat{U}_m - P_{\mathcal{H}}\hat{U}_m\| \\
&\leq \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \left(\|P_{\mathcal{H}}^{(n,\lambda)}\hat{U}_m - P_{\mathcal{H}}\hat{U}_m\|_n + \delta_{n,1}^{1/2} J(P_{\mathcal{H}}^{(n,\lambda)}\hat{U}_m - P_{\mathcal{H}}\hat{U}_m) \right) \\
&\leq I + II + III + \delta_{n,1}^{1/2} \left(\sup_{\mathcal{H} \in \mathcal{D}_n} J(P_{\mathcal{H}}^{(n,\lambda)}\hat{U}_m) + \sup_{\mathcal{H} \in \mathcal{D}_n} J(P_{\mathcal{H}}\hat{U}_m) \right) \\
&\leq I + II + III + \delta_{n,1}^{1/2} (J(P_{\mathcal{H}}^{(n,\lambda)}(F + \epsilon - \hat{F}_m)) + J(P_{\mathcal{H}}(F + \epsilon - \hat{F}_m))). \tag{2.11}
\end{aligned}$$

From (2.9), the last item above is

$$\begin{aligned}
&\leq \delta_{n,1}^{1/2} [(W+1)((C_n+1)^m - 1) + (W+1)C_n + (KW + K(W+1))((C_n+1)^m - 1)] \\
&\leq [(W+1)(K+1)(C_n+1)^m + (W+1)C_n - (K+W+1)]\delta_{n,1}^{1/2} \\
&\leq 2(W+1)(K+1)(C_n+1)^m\delta_{n,1}^{1/2}. \tag{2.12}
\end{aligned}$$

Lemma 2.4 follows from (2.11) and (2.12). ■

Lemma 2.5 *Let $\delta_n = 2(1+K)\delta_{n,1}^{1/2} + \delta_{n,3}$,*

$$\begin{aligned}
&\max\{\|\tilde{F}_m - \hat{F}_m\|, \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|P_{\mathcal{H}}^{(n,\lambda)}(F + \epsilon - \hat{F}_m)\|_n - \|P_{\mathcal{H}}(F - \tilde{F}_m)\|\} \\
&\leq 2 \sum_{k=1}^m (W+1)(C_n+1)^k \cdot 2^{m-k}\delta_n + 2\delta_n(W+1)(C_n+1)^m.
\end{aligned}$$

Proof: Let $A_{n,m} = \|\tilde{F}_m - \hat{F}_m\|$. We have

$$\begin{aligned}
A_{n,m} &= \|\tilde{F}_{m-1} + P_{\hat{S}_m}(F - \tilde{F}_{m-1}) - (\hat{F}_{m-1} + \rho P_{\hat{S}_m}^{(n,\lambda)}(F + \epsilon - \hat{F}_{m-1}))\| \\
&\leq A_{n,m-1} + \|P_{\hat{S}_m}(F - \tilde{F}_{m-1}) - P_{\hat{S}_m}^{(n,\lambda)}(F + \epsilon - \hat{F}_{m-1})\| \\
&\leq 2A_{n,m-1} + \|(P_{\hat{S}_m}^{(n,\lambda)} - P_{\hat{S}_m})\hat{U}_{m-1}\| \\
&\leq 2A_{n,m-1} + (W+1)(C_n+1)^m(2(1+K)\delta_{n,1}^{1/2} + \delta_{n,3}).
\end{aligned}$$

Let $\delta_n = 2(1+K)\delta_{n,1}^{1/2} + \delta_{n,3}$, note that for $m = 0$, $A_{n,m} = 0$, it can be proved recursively that $A_{n,m} \leq \sum_{k=1}^m (W+1)(C_n+1)^k \cdot 2^{m-k}\delta_n$.

For the second part, note

$$\begin{aligned}
&\sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|P_{\mathcal{H}}^{(n,\lambda)}(F + \epsilon - \hat{F}_m)\|_n - \|P_{\mathcal{H}}(F - \tilde{F}_m)\| \\
&\leq \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|P_{\mathcal{H}}(\hat{F}_m - \tilde{F}_m)\| + \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|P_{\mathcal{H}}^{(n,\lambda)}(F + \epsilon - \hat{F}_m)\|_n - \|P_{\mathcal{H}}(F - \tilde{F}_m)\| \\
&\leq A_{n,m} + II.
\end{aligned}$$

The second term in the above equality is

$$\begin{aligned}
II &\leq \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \left| \|P_{\mathcal{H}}^{(n,\lambda)}(F + \epsilon - \hat{F}_m)\|_n - \|P_{\mathcal{H}}(F - \hat{F}_m)\| \right| \\
&\quad + \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \left| \|P_{\mathcal{H}}(F - \hat{F}_m)\| - \|P_{\mathcal{H}}(F - \tilde{F}_m)\| \right| \\
&\leq \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \left(\|(P_{\mathcal{H}}^{(n,\lambda)} - P_{\mathcal{H}})\hat{U}_m\| + \delta_{n,1}^{1/2} J(P_{\mathcal{H}}^{(n,\lambda)}\hat{U}_m) \right) + \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|\hat{F}_m - \tilde{F}_m\| \\
&\leq \sup_{\mathcal{H} \in \mathcal{D}_{p_n}} \|(P_{\mathcal{H}}^{(n,\lambda)} - P_{\mathcal{H}})\hat{U}_m\| + \delta_{n,1}^{1/2}(W+1)[(C_n+1)^m - 1 + C_n] + A_{n,m}.
\end{aligned}$$

The second inequality follows Lemma 2.2 and the fact $\sqrt{x+y} \leq \sqrt{x} + \sqrt{y}$ as $x \geq 0, y \geq 0$. Hence,

$$\begin{aligned}
I + II &\leq 2A_{n,m} + (W+1)(C_n+1)^m \delta_n + \delta_{n,1}^{1/2}(W+1)[(C_n+1)^m - 1 + C_n] \\
&\leq 2 \sum_{k=1}^m (W+1)(C_n+1)^k \cdot 2^{m-k} \delta_n + 2\delta_n(W+1)(C_n+1)^m \triangleq \xi_{n,m}. \quad (2.13)
\end{aligned}$$

We first show the convergence of \tilde{F} . Let $\delta_n = \max_j \delta_{n,j}$. Note that $\hat{F}_m = \sum_{k=1}^m \hat{h}_m$ with $\hat{h}_m = P_{\hat{S}_m}^{(n,\lambda)}(F + \epsilon - \hat{F}_{m-1})$. By Lemma 2.5,

$$\begin{aligned}
\|\tilde{h}_m\| &= \|P_{\hat{S}_m}(F - \tilde{F}_{m-1})\| \geq \|P_{\hat{S}_m}^{(n,\lambda)}(F + \epsilon - \hat{F}_{m-1})\|_n - \xi_{m,n} \\
&= \|\hat{h}_m\|_n - \xi_{m,n} \\
&\geq t \sup_{S \in \mathcal{D}} \|P_{\hat{S}_m}^{(n,\lambda)}(F + \epsilon - \hat{F}_{m-1})\|_n - \xi_{m,n} \\
&\geq t(\sup_{S \in \mathcal{D}} \|P_{\hat{S}_m}(F - \tilde{F}_{m-1})\|) - (1+t)\xi_{m,n}. \quad (2.14)
\end{aligned}$$

Note that

$$\sup_{S \in \mathcal{D}} \|P_S(F - \tilde{F}_m)\| = \sup_{S \in \mathcal{D}} \|P_S(F - \tilde{F}_{m-1})\| + \sup_{S \in \mathcal{D}} \|P_S P_{\hat{S}_m}(F - \tilde{F}_{m-1})\| \leq 2 \sup_{S \in \mathcal{D}} \|P_S(F - \tilde{F}_{m-1})\|$$

and

$$\begin{aligned} \frac{\xi_{n,m}}{\xi_{n,m-1}} &= \frac{\sum_{k=1}^m (C_n + 1)^k 2^{m-k} + (C_n + 1)^m}{\sum_{k=1}^{m-1} (C_n + 1)^k 2^{m-1-k} + (C_n + 1)^{m-1}} \\ &= 2 + \frac{2(C_n + 1)^m - 2(C_n + 1)^{m-1}}{\xi_{n,m-1}} > 2. \end{aligned}$$

It implies that $\sup_{S \in \mathcal{D}} \|P_S(F - \tilde{F}_{m-1})\| / [(1+t)\xi_{n,m}]$ is a decreasing sequence with respect to m . Let $m_n = o(\log n)$. Define $B_n = \{\omega : \sup_{S \in \mathcal{D}} \|P_S(F - \tilde{F}_{m-1})\| / [(1+t)\xi_{n,m_n}] > 2/t\}$. Therefore, for all $m \leq m_n$, $\sup_{S \in \mathcal{D}} \|P_S(F - \tilde{F}_{m-1})\| / [(1+t)\xi_{n,m_n}] > 2/t$, which together with (2.14), lead to $\|\tilde{h}_m\| \geq t/2 \sup_{S \in \mathcal{D}} \|P_S(F - \tilde{F}_{m-1})\|$ for all $m \leq m_n$ on B_n . It then follows that $\{\tilde{F}_m\}_{m_n}$ satisfies condition (2.4) with $t_m = t/2$, and by Lemma 2.1, $\|F - \tilde{F}_{m_n}\| \leq T(1 + m_n t^2/4)^{-t/(8+2t)} = o(1)$ on B_n . On the other hand, on B_n^C , $\|F - \tilde{F}_{m_n}\|^2 \leq \sup_{S \in \mathcal{D}} \|P_S(F - \tilde{F}_{m-1})\| (\sum_{k=1}^{p_n} \|g_k\| + \sum_{k=1}^{m_n} \|\tilde{h}\|) \leq 2(1+t)\xi_{n,m_n}(1+m_n)T \leq 8(W+1)(C_n+1)^m \cdot 2^m \delta_n \cdot (1+t)(1+m_n)T = o_P(1)$.

For the consistency of \hat{F}_{m_n} , since $\|F - \hat{F}_{m_n}\| \leq \|F - \tilde{F}_{m_n}\| + \|\tilde{F}_{m_n} - \hat{f}_{m_n}\|$, by Lemma 2.5, $\|F - \hat{F}_{m_n}\| = o_P(1)$. ■

References

- Agarwal, A., Negahban, S. and Wainwright, M. (2012). Noisy matrix decomposition via convex relaxation: optimal rates in high dimensions. *The Annals of Statistics*, **Vol. 40**, No. 2, 1171–1197.
- Akaike, H (1973). Information theory and an extension of the maximum likelihood principle. *International Symposium on Information Theory*, eds. B.N. Petrov and F. Csaki, Budapest: Akademia Kiado, 267-281.
- Amit, Y., Fink, M., Srebro, N., and Ullman, S. (2007). Uncovering shared structures in multiclass classification. *Proceedings of the 24th Annual International Conference on Machine Learning*, 17–24.
- Arlot, S. and Celisse, A. (2010). A survey of cross-validation procedures for model selection. *Statist. Surv.*, **4**, 40–79.
- Beck, Amir and Teboulle, Marc (2009). A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, **Vol. 2**, No. 1, 183-202.
- Breiman, L. Friedman, J., Olshen, R. and Stone, C. (1984). Classification and Regression Trees. Belmont, CA: Wadsworth.
- Bhlmann, P. and Yu, B. (2003). Boosting with the L2 loss: regression and classification. *Journal of the American Statistical Association*, **98**, 324-339.
- Bhlmann, P. (2006). Boosting for high-dimensional linear models. *The Annals of Statistics*, **Vol. 34**, No. 2, 559-583.

- Bunea, F., and She, Y. (2011). Optimal selection of reduced rank estimators of high-dimensional matrices. *Annals of Statistics*, **39(2)**, 1282-1309.
- Cai, J.F., Candès, E.J. and Shen, Z. (2010). A singular value thresholding algorithm for matrix completion. Arxiv preprint arXiv:0810.3286.
- Cai, T. T., Liu, W. and Luo, X. (2011). A constrained l_1 minimization approach to sparse precision matrix estimation. *J. Amer. Statist. Assoc.*, **106**, 594-607.
- Candes, E., Li. X., Ma, U., and Wright, J. (2009). Robust principal component analysis. *Journal of ACM*, **58(1)**, 1-37.
- Candes, E.J. and Recht, B. (2009). Exact matrix completion via convex optimization. *Found of Comput. Math.*, **9**, 717-772.
- Chandrasekaran, V., Sanghavi, S., Parrilo, P.A., and Willsky, A.S. (2011). Rank-sparsity incoherence for matrix decomposition, *SIAM J. Optim.*, **21**, 572-596.
- Chen, B. , He, S., Li, Z. and Zhang, S. (2012). Maximum block improvement and polynomial optimization., *SIAM Journal on Optimization*, **22**, 87-107.
- Chen, J., Liu, J., and Ye, J. (2010). Learning incoherent sparse and low-rank patterns from multiple tasks. *The Sixteenth ACM SIGKDD International Conference On Knowledge Discovery and Data Mining*. (SIGKDD 2010).
- Detting, M. and Bhlmann, P. (2006). Boosting for tumor classification with gene expression data. *Bioinformatics*, **Vol. 19**, No. 9, 1061-1069.
- Efron, B. (2004). The estimation of prediction error: Covariance penalties and cross-validation. *Journal of the American Statistical Association*, **99**, 619-642. (with discussion).
- Freund, Y. (1995). Boosting a weak learning algorithm by majority. *Information and Computation*, **121(2)**, 256-285.

- Freund, Y. and Schapire R. (1996). Experiments with a new boosting algorithm. *Machine Learning: Proceedings of the Thirteenth International Conference*, 148-156.
- Freund, Y. and Schapire R. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, **55(1)**, 148-156.
- Freund, Y. and Schapire R. (1999). A short introduction to Boosting. *Journal of Japanese Society for Artificial Intelligence*, **14(5)**, 771-780.
- Friedman, J. (2001). Greedy function approximation: a gradient boosting machine. *The Annals of Statistics*, **29**, 1189-1232.
- Friedman, J., Hastie, T., and Tibshirani, R. (2000). Additive logistic regression: A statistical view of boosting. *The Annals of Statistics*, **28**, 337-374.
- Ganesh, A., Min, K., Wright, J. and Ma, Y. (2012). Principal component pursuit with reduced linear measurements. *International Symposium on Information Theory*.
- Golub, G. and Van Loan, C.(1996). Matrix Computations. Third edition. *London: The Johns Hopkins University Press*.
- Halko, N., Martinsson P. G. and Tropp, J. A.. (2011). Finding structure with randomness: stochastic algorithms for constructing approximate matrix decompositions. *SIAM Rev.*, **53(2)**, 217-288.
- Hothorn, T., Buhlmann, P., Dudoit, S., Molinaro, A. and Van Der Lann, M. J. (2006). Survival ensembles. *Biostatistics*, **7**, 355-373.
- Hurvich, C., Simonoff, J. and Tsai C. (1998). Smoothing parameter selection in non-parametric regression using an improved Akaike information criterion. *Journal of Royal Statistical Society, Series B*, **60**, 271-293.

- Jain, P., Meka, R. and Dhillon, I. (2010). Guaranteed rank minimization via singular value projection. *Advances in Neural Information Processing Systems*, **23**, 937–945.
- Kenneth, S. (1984). Probability inequalities for empirical processes and a law of the iterated logarithm. *The Annals of Probability*, **Vol 12**, No 2, 1041-1067.
- Kolmogorov, A.N. and Tihomirov, V.M. (1959). ϵ -entropy and ϵ -capacity of sets in function spaces. *Uspekhi Mat. Nauk.* **14** 3-86. [In Russian. English translation, *Ameri. Math. Soc. Transl.* **2**, **17**, 277-364.(1961)].
- Li, H. and Luan, Y. (2005). Boosting proportional hazards models using smoothing splines, with applications to high-dimensional microarray data. *Bioinformatics*, **21**, 2403-2409.
- Lin, Z., Chen, M., Wu, L. and Ma, Y. (2009). The augmented Lagrange multiplier method for exact recovery of corrupted low-rank matrices. *UIUC Technical Report UILU-ENG-09-2215*.
- Liu, J., and Ye, J. (2009). Efficient Euclidean projections in linear time. *The Twenty-Sixth International Conference on Machine Learning*.
- Luan, Y. and Li, H. (2008). Group additive regression models for genomic data analysis. *Biostatistics*, **9(1)**, 100-113.
- Negahban, S. and Wainwright, M.J. (2011). Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *Annals of Statistics*, **39(2)**, 1069-1097.
- Ossiander, M. (1987). A central limit theory under metric entropy with L_2 bracketing. *Ann. Probab.* **15** 897-919.
- Porat, B. (1997). *A Course in Digital Signal Processing*, New York: John Wiley.
- Schapire, R. (1990). The strength of weak learnability. *Machine Learning*, **5(2)**, 197-227.
- She, Y. (2013). Reduced rank vector generalized linear models for feature extraction. *Statistics and Its Interface*, **Vol 6**, 197-209.

- Shen, X., Pan, W., and Zhu, Y. (2012). Likelihood-based selection and sharp parameter estimation. *Journal of the American Statistical Association*, **107**, 223-232.
- Srebro, N., Rennie, J.D.M., and Jaakkola, T.S. (2005). Maximum-margin matrix factorization. *Advances in neural information processing system*, **17**, 1329–1336.
- Stanica, P., and Montgomery, A.P. (2001). Good lower and upper bounds on binomial coefficients. *J. Ineq. in Pure. Appl. Math.*, **2**, art 30.
- Stock, J. H. and Watson, M. W. (2012). Generalized shrinkage methods for forecasting using many predictors. *Journal of Business and Economic Statistics*. To appear.
- Temlyakov, V. (2000). Weak Greedy Algorithm. *Advances in Computational Mathematics*, **12**, 213-227.
- Van der Vaart, A. and Wellner, J. (2000). Weak Convergence and Empirical Processes with Application to Statistics. *Springer*, New York.
- Waters, A.E., Sankaranarayanan, A.C. and Baraniuk, R.G.(2011). SpaRCS: Recovering low-rank and sparse matrices from compressive measurements. *Neural Information Processing Systems*, Granada, Spain.
- Wei, Z. and Li, H. (2007). Nonparametric pathway-based regression models for analysis of genomic data. *Biostatistics*, **8(2)**, 265-284.
- Wong, W.H., and Shen, X. (1995). Probability inequalities for likelihood ratios and convergence rates of sieve MLEs. *Ann. Statist.*, **23**, 339-362.
- Wright, J., Ganesh, A., Min, K., and Ma, Y. (2013). Compressive principal component pursuit. *Information and Inference*.
- Wright, J., Ganesh, A., Rao, S., and Ma, Y. (2009). Robust principal component analysis: Exact recovery of corrupted low-rank matrices via convex optimization. *Advances in Neural Information Processing Systems*.

- Xing, S., Zhu, Y., Shen, X., and Ye, J. (2012). Optimal exact rank minimization for noisy data. *Proceeding the 18th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, Beijing, China.
- Yin, J. , Chen, Xi and Xing, E. (2012). Group Sparse Additive Models. *ICML*, 2012.
- Yuan, X. and Yang, J. (2013). Sparse and low-rank matrix decomposition via alternating direction methods. *Pacific Journal of Optimization*, **9(1)**, 167-180.
- Zhou, T. and Tao, D. (2011). GoDec: Randomized low-rank & sparse matrix decomposition in noisy case. *Proceedings of the 28th International Conference on Machine Learning*, Bellevue, WA, USA.