

**STATISTICAL INFERENCE FOR HIGH DIMENSIONAL DATA WITH LOW
RANK STRUCTURE**

A Dissertation
Presented to
The Academic Faculty

By

Fan Zhou

In Partial Fulfillment
of the Requirements for the Degree
Doctor of Philosophy in the
School of Computational Science and Engineering

Georgia Institute of Technology

December 2018

Copyright © Fan Zhou 2018

STATISTICAL INFERENCE FOR HIGH DIMENSIONAL DATA WITH LOW RANK STRUCTURE

Approved by:

Dr. Vladimir Koltchinskii, Advisor
School of Mathematics
Georgia Institute of Technology

Dr. Edmond Chow
School of Computational Science
and Engineering
Georgia Institute of Technology

Dr. Hongyuan Zha
School of Computational Science
and Engineering
Georgia Institute of Technology

Dr. Sung Ha Kang
School of Mathematics
Georgia Institute of Technology

Dr. Mark A. Davenport
School of Electrical and Computer
Engineering
Georgia Institute of Technology

Dr. Mayya Zhilova
School of Mathematics
Georgia Institute of Technology

Date Approved: July 21, 2018

”You can’t connect the dots looking forward; you can only connect them looking backwards. So you have to trust that the dots will somehow connect in your future. You have to trust in something - your gut, destiny, life, karma, whatever. Because believing that the dots will connect down the road will give you the confidence to follow your heart even when it leads you off the well worn path; and that will make all the difference.”

Steve Jobs

To my parents

ACKNOWLEDGEMENTS

After the most precious five years of my life at Georgia Tech, so many aspects of me have been transformed. I would acknowledge the great teachers I met, the valuable friendship I made, and the unconditional support and love from my family I received that attribute such a transformation to the next level of my life.

First and foremost, I want to express my deepest gratitude to my advisor Prof. Vladimir Koltchinskii, a truly respectful scholar who patiently transformed me from a student to a researcher. Through those uncountable conversations, he not only introduced me lots of beauty in mathematics but also taught me how to be a scholar with integrity and a person with decency. Those qualities are more scarce I realize as I grow older, and they shaped me as who I am today and for sure I will benefit from them for the rest of my life.

I'd also like to thank Prof. Edmond Chow, the first teacher I met at Georgia Tech with whom I would develop a very good friendship later on. It is always nice to talk with Edmond! He introduced me lots of interesting research topics on the computational aspects and real world applications. He broadened my horizon as a researcher.

Most of my research is supported by NSF grants from Prof. Vladimir Koltchinskii's DMS-1509739 and Prof. Arkadi Nemirovski's CCF-1523768.

I want to thank my academic brother Dong Xia who gave me a lot of help and encouragement as a big brother along the path. We found so many shared interests and common attitudes towards research and life which eventually leads to a great friendship. I want to thank my friends Haiyu Zou, Xin Xing, Xin Wang, Haoyan Zhai, Jiangning Chen, Qianli Hu, Junxiong Jia, Jian Zu and Dawei He for those good old days we spent together in Atlanta. I also want to thank Guojing Cong from IBM Research who hosted me for a wonderful summer in New York.

Finally, I want to thank my parents especially my mother for their unconditional support and endless love. My gratitude to them is beyond words.

TABLE OF CONTENTS

Acknowledgments	v
List of Tables	ix
List of Figures	x
Chapter 1: Introduction to Nonparametric Estimation of Low Rank Matrix Valued Function	1
1 Introduction	1
2 Preliminaries	3
2.1 Notations	3
2.2 Matrix Completion and statistical learning setting	4
2.3 Matrix Valued Functions	5
Chapter 2: Optimal Estimation of Low Rank Matrix Valued Function	6
1 A local polynomial Lasso estimator	6
2 Global estimators and upper bounds on integrated risk	11
2.1 From localization to globalization	11
2.2 Bias reduction through higher order kernels	12
3 Lower bounds under matrix completion setting	17
4 Model selection	27

Chapter 3: Simulation Results of Nonparametric Estimation of Low Rank Matrix Valued Function	35
1 An ADMM Algorithm	35
2 Numerical results	36
2.1 Pointwise estimation simulation	38
2.2 Integrated risk estimation simulation	40
2.3 Simulation of model selection	41
Chapter 4: The ℓ_∞ Perturbation of HOSVD and Low Rank Tensor Denoising . .	44
1 Introduction	44
2 Preliminaries on Tensor and HOSVD	48
2.1 Notations	48
2.2 HOSVD and Eigengaps	49
3 Main Results	51
3.1 Second Order Spectral Analysis	51
3.2 Perturbation of Linear Forms of Singular Vectors	52
3.3 Low Rank Tensor Denoising ℓ_∞ Bound	55
4 Applications	57
4.1 High Dimensional Clustering	57
4.2 Subtensor Localization	59
4.3 Numerical Experiments	60
5 Proofs	61
5.1 Proof of Theorem 3.1	64
5.2 Proof of Theorem 3.2	76

5.3	Proof of Theorem 3.3	76
Chapter 5: Appendices		83
1	Proof of Lemma 1	85
2	Proof of Lemma 6	94
3	Proof of Lemma 3	96
4	Proof of Theorem 5.1	97
5	Proof of Lemma 5	104
References		113

LIST OF TABLES

3.1	Pointwise error rate comparison with different sample size n	39
3.2	Pointwise error rate comparison with different sample size n	40
3.3	Integrated error rate comparison with different sample size n	41
3.4	Integrated error rate comparison with different sample size n	41
3.5	Model Selection	42
3.6	Model Selection	42

LIST OF FIGURES

3.1	Pointwise risk convergence of ADMM Algorithm	37
3.2	Integrated risk convergence of ADMM Algorithm	37
3.3	$n = 1600$	38
3.4	$n = 6400$	38
3.5	$n = 25600$	38
3.6	$n = 102400$	38
3.7	$n = 409600$	39
3.8	$n = 1638400$	39
3.9	$n = 3276800$	39
3.10	True data	39
3.11	The point risk comparison at $t_0 = 0.5$	40
3.12	The integrated risk comparison with different sample size n	41
3.13	Model Selection on Grid \mathcal{H}	43
4.1	Comparison on ℓ_∞ -norm, ℓ_2 -norm and bias corrected ℓ_∞ norm in high dimension clustering and subtensor localization.	61

SUMMARY

As the big data era has come, a lot of machine learning problems involve data with very high dimension. However, the computational power is always limited. Such kind of practical issue motivates the works in this thesis. In the thesis, we study two major topics on statistical inference of high dimensional data with low rank structure occurred in many machine learning and statistics applications.

The first topic is about nonparametric estimation of low rank matrix valued function with applications in building dynamic recommender systems and recovering euclidean distance matrices in molecular biology. We propose an innovative nuclear norm penalized local polynomial estimator and establish an upper bound on its point-wise risk measured by Frobenius norm. Then we extend this estimator globally and prove an upper bound on its integrated risk measured by L_2 -norm. We also propose another new estimator based on bias-reducing kernels to study the case when the matrix valued function is not necessarily low rank and establish an upper bound on its risk measured by L_∞ -norm. We show that the obtained rates are all optimal up to some logarithmic factor in minimax sense. Finally, we propose an adaptive estimation procedure for practitioners based on Lepski's method and the penalized data splitting technique which is computationally efficient and can be easily implemented and parallelized. Most results in this work is in the paper [1].

The other topic is about spectral perturbation analysis of higher order singular value decomposition (HOSVD) of tensor under Gaussian noise. Given a tensor contaminated with Gaussian noise, we establish sharp upper bounds on the perturbation of linear forms of singular vectors of HOSVD. In particular, sharp upper bounds are proved for the component-wise perturbation of singular vectors. These results can be applied on sub-tensor localization and low rank tensor denoising. This work is a collaboration with Dong Xia and can be found in the paper [2].

CHAPTER 1
INTRODUCTION TO NONPARAMETRIC ESTIMATION OF LOW RANK
MATRIX VALUED FUNCTION

1 Introduction

Let $A : [0, 1] \rightarrow \mathbb{H}_m$ (the space of Hermitian matrices) be a matrix valued function. The goal is to study the problem of statistical estimation of A based on the regression model

$$\mathbb{E}(Y_j | \tau_j, X_j) = \langle A(\tau_j), X_j \rangle, \quad j = 1, \dots, n, \quad (1.1)$$

where τ_j are i.i.d. time design variables uniformly distributed in $[0, 1]$, X_j are i.i.d. matrix completion sampling matrices, Y_j are independent bounded random responses. Sometimes, it will be convenient to write model (1.1) in the form

$$Y_j = \langle A(\tau_j), X_j \rangle + \xi_j, \quad j = 1, \dots, n, \quad (1.2)$$

where the noise variables $\xi_j = Y_j - \mathbb{E}(Y_j | \tau_j, X_j)$ are independent and have zero means. In particular, we are interested in the case when A is low rank and satisfies certain smoothness condition. When $A(t) = A_0$ for some $A_0 \in \mathbb{H}_m$ and for any $t \in [0, 1]$, such problem coincides with the well known matrix completion/recovery problem which has drawn a lot of attention in the statistics community during the past few years, see [3, 4, 5, 6, 7, 8, 9, 10, 11, 12]. The low rank assumption in matrix completion/estimation problems has profound practical background. For instance, when [13] introduced their famous work on matrix factorization techniques for recommender systems, they considered temporal dynamics, see [14]. Another very common example is Euclidean distance matrix (EDM) which contains the distance information of a large set of points like molecules which are in

low dimensional spaces such as \mathbb{R}^2 or \mathbb{R}^3 . To be more specific, given m points p_1, \dots, p_m in \mathbb{R}^d , the EDM $D \in \mathbb{R}^{m \times m}$ formed by them has entries $D_{ij} = \|p_i - p_j\|_2^2$. Clearly, this matrix has rank at most $d + 1$ regardless of its size m . If $m \gg d$, then the recovery problem falls into the low rank realm. Similar topics in cases when points are fixed (see [15]) or in rigid motion (see [16]) have been studied. While points are moving in smooth trajectories, the EDMs are naturally high dimensional low rank matrix valued functions.

An appealing way to address the low rank issue in matrix completion/estimation is through nuclear norm minimization, see [17]. In section 1 of chapter 2, we inherit this idea and propose a local polynomial estimator (see [18]) with nuclear norm penalization:

$$\hat{S}^h = \arg \min_{S \in \mathbb{D}} \frac{1}{nh} \sum_{j=1}^n K\left(\frac{\tau_j - t_0}{h}\right) \left(Y_j - \left\langle \sum_{i=0}^{\ell} S_i p_i\left(\frac{\tau_j - t_0}{h}\right), X_j \right\rangle\right)^2 + \varepsilon \|S\|_1. \quad (1.3)$$

where $\mathbb{D} \subset \mathbb{H}_{(\ell+1)m}$ is a closed subset of block diagonal matrices with $S_j \in \mathbb{H}_m$ on its diagonal, and $\{p_i\}$ is a sequence of orthogonal polynomials with nonnegative weight function K . The solution to the convex optimization problem (1.3) induces a pointwise estimator of $A(t_0)$: $\hat{S}^h(t_0) := \sum_{i=0}^{\ell} \hat{S}_i^h p_i(0)$. We prove that under mild conditions, $\hat{S}^h(t_0)$ achieves a rate of $O\left(\left(\frac{mr \log n}{n}\right)^{2\beta/(2\beta+1)}\right)$ on the pointwise risk measured by $\frac{1}{m^2} \|\hat{S}^h(t_0) - A(t_0)\|_2^2$ over Hölder class $\Sigma(\beta, L)$ with low rank parameter r , where $\|\cdot\|_2$ denotes the Frobenius norm of a matrix. In section 2.1 of chapter 2, we propose a new global estimator \hat{A} based on the local results and prove that \hat{A} achieves a rate of $O\left(\left(\frac{mr \log n}{n}\right)^{2\beta/(2\beta+1)}\right)$ on the integrated risk measured by L_2 -norm, i.e. $\frac{1}{m^2} \int_0^1 \|\hat{A}(t) - A(t)\|_2^2 dt$. Then we study another naive kernel estimator \tilde{A} which can be used to estimate matrix valued functions which are not necessarily low rank. This estimator is associated with another popular approach to deal with low rank recovery which is called singular value thresholding, see [4, 9, 12]. We prove that \tilde{A} achieves a rate of $O\left(\left(\frac{m \log n}{n}\right)^{2\beta/(2\beta+1)}\right)$ measured by $\sup_{t \in [h, 1-h]} \frac{1}{m^2} \|\tilde{A}(t) - A(t)\|^2$, where $\|\cdot\|$ denotes the matrix operator norm. Note that those rates coincide with that of classical matrix recovery/estimation setting when the smoothness parameter $\beta \rightarrow \infty$.

An immediate question is whether the above rates are optimal. In section 3 of chapter 2, we prove that all the rates we established are optimal up to some logarithmic factor in the minimax sense, which essentially verified the effectiveness of our methodology.

As one may have noticed, there is an adaptation issue involved in (1.3). Namely, one needs to choose a proper bandwidth h and a proper order of degree ℓ of polynomials. Both parameters are closely related to the smoothness of A which is unknown to us in advance. In section 4 of chapter, we propose a model selection procedure based on Lepskii's method ([19]) and the work of [20] and [21]. We prove that this procedure adaptively selects an estimator that achieves a rate on the integrated risk measured by L_2 -norm which is the smallest among all candidates plus a negligible term. What is more important, such a procedure is computationally efficient, feasible in high dimensional setting, and can be easily parallelized.

The major contribution of our work is that we generalized the recent developments of matrix completion/estimation theory to low rank matrix valued function setting by proposing a new optimal estimation procedure. To our best knowledge, no one has ever thoroughly studied such problems from a theoretical point of view.

2 Preliminaries

In this section, we introduce some important definitions, basic facts, and notations for the convenience of presentation.

2.1 Notations

For any Hermitian matrices $A, B \in \mathbb{H}_m$, denote $\langle A, B \rangle = \text{tr}(AB)$ which is known as the Hilbert-Schmidt inner product. Denote $\langle A, B \rangle_{L_2(\Pi)} = \mathbb{E}\langle A, X \rangle \langle B, X \rangle$, where Π denotes the distribution of X . The corresponding norm $\|A\|_{L_2(\Pi)}^2$ is given by $\|A\|_{L_2(\Pi)}^2 = \mathbb{E}\langle A, X \rangle^2$.

We use $\|\cdot\|_2$ to denote the Hilbert-Schmidt norm (Frobenius norm or Schatten 2-norm) generated by the inner product $\langle \cdot, \cdot \rangle$; $\|\cdot\|$ to denote the operator norm (spectral norm) of

a matrix: the largest singular value; $\|\cdot\|_1$ to denote the trace norm (Schatten 1-norm or nuclear norm), i.e. the sum of singular values; $|A|$ to denote the nonnegative matrix with entries $|A_{ij}|$ corresponding to A .

We denote

$$\sigma^2 := \mathbb{E}\xi^2, \quad \sigma_X^2 := \left\| n^{-1} \sum_{j=1}^n \mathbb{E}X_j^2 \right\|, \quad U_X := \|\|X\|\|_{L_\infty}.$$

2.2 Matrix Completion and statistical learning setting

The matrix completion setting refers to that the random sampling matrices X_j are i.i.d. uniformly distributed on the following orthonormal basis \mathcal{X} of \mathbb{H}_m :

$$\mathcal{X} := \{E_{kj} : k, j = 1, \dots, m\},$$

where $E_{kk} := e_k \otimes e_k$, $k = 1, \dots, m$; $E_{jk} := \frac{1}{\sqrt{2}}(e_k \otimes e_j + e_j \otimes e_k)$, $1 \leq k < j \leq m$; $E_{kj} := \frac{i}{\sqrt{2}}(e_k \otimes e_j - e_j \otimes e_k)$, $1 \leq k < j \leq m$ with $\{e_j\}_{j=1}^m$ being the canonical basis of \mathbb{C}^m . The following identities are easy to check when the design matrices are under matrix completion setting:

$$\|A\|_{L_2(\Pi)}^2 = \frac{1}{m^2} \|A\|_2^2, \quad \sigma_X^2 \leq \frac{2}{m}, \quad U_X = 1. \quad (2.1)$$

The statistical learning setting refers to the bounded response case: there exists a constant a such that

$$\max_{j=1, \dots, n} |Y_j| \leq a, \quad a.s. \quad (2.2)$$

In this paper, we will consider model (1.1) under both matrix completion and statistical learning setting.

2.3 Matrix Valued Functions

Let $A : [0, 1] \rightarrow \mathbb{H}_m$ be a matrix valued function. One should notice that we consider the image space to be Hermitian matrix space for the convenience of presentation. Our methods and results can be readily extended to general rectangular matrix space. Now we define the rank of a matrix valued function. Let $\text{rank}_A(t) := \text{rank}(A(t))$, $\forall t \in [0, 1]$.

Definition 1. Let β and L be two positive real numbers. The **Hölder class** $\Sigma(\beta, L)$ on $[0, 1]$ is defined as the set of $\ell = \lfloor \beta \rfloor$ times differentiable functions $f : [0, 1] \rightarrow \mathbb{R}$ with derivative $f^{(\ell)}$ satisfying

$$|f^{(\ell)}(x) - f^{(\ell)}(x')| \leq L|x - x'|^{\beta-\ell}, \quad \forall x, x' \in [0, 1]. \quad (2.3)$$

In particular, we are interested in the following assumptions on matrix valued functions.

A1 Given a measurement matrix X and for some constant a ,

$$\sup_{t \in [0, 1]} |\langle A(t), X \rangle| \leq a.$$

A2 Given a measurement matrix X and for some constant a , the derivative matrices $A^{(k)}$ of A satisfy

$$\sup_{t \in (0, 1)} |\langle A^{(k)}(t), X \rangle| \leq a, \quad k = 1, \dots, \ell.$$

A3 The rank of $A, A', \dots, A^{(\ell)}$ are uniformly bounded by a constant r ,

$$\sup_{t \in [0, 1]} \text{rank}_{A^{(k)}}(t) \leq r, \quad k = 0, 1, \dots, \ell.$$

A4 Assume that for $\forall i, j$, a_{ij} is in the Hölder class $\Sigma(\beta, L)$.

CHAPTER 2

OPTIMAL ESTIMATION OF LOW RANK MATRIX VALUED FUNCTION

1 A local polynomial Lasso estimator

In this section, we study the estimation of matrix valued functions that are low rank. The construction of our estimator is inspired by localization of nonparametric least squares and nuclear norm penalization. The intuition of the localization technique originates from classical local polynomial estimators, see [18]. The intuition behind nuclear norm penalization is that whereas rank function counts the number of non-vanishing singular values, the nuclear norm sums their amplitude. The theoretical foundations behind nuclear norm heuristic for the rank minimization was proved by [17]. Instead of using the trivial basis $\{1, t, t^2, \dots, t^\ell\}$ to generate an estimator, we use orthogonal polynomials which fits our problem better. Let $\{p_i\}_{i=0}^\infty$ be a sequence of orthogonal polynomials with nonnegative weight function K compactly supported on $[-1, 1]$, then

$$\int_{-1}^1 K(u) p_i(u) p_j(u) du = \delta_{ij}.$$

There exist an invertible linear transformation $T \in \mathbb{R}^{(\ell+1) \times (\ell+1)}$ such that

$$(1, t, t^2/2!, \dots, t^\ell/\ell!)^T = T(p_0, p_1, \dots, p_\ell)^T.$$

Apparently, T is lower triangular. We denote $R(T) = \max_{1 \leq j \leq \ell+1} \sum_{i=1}^{\ell+1} |T_{ij}|$. Note that in some literature, $R(T)$ is denoted as $\|T\|_1$ as the matrix "column norm". Since we already used $\|\cdot\|_1$ to denote the nuclear norm, $R(T)$ is used to avoid any ambiguity. Denote

$$\mathbb{D} := \left\{ \text{Diag} \begin{bmatrix} S_0 & S_1 & \dots & S_{\ell-1} & S_\ell \end{bmatrix} \right\} \subset \mathbb{H}_{m(\ell+1)}$$

the set of block diagonal matrices with $S_k \in \mathbb{H}_m$ satisfying $|S_{ij}| \leq R(T)a$. With observations (τ_j, X_j, Y_j) , $j = 1, \dots, n$ from model (1.1), define \widehat{S}^h as

$$\widehat{S}^h = \arg \min_{S \in \mathbb{D}} \frac{1}{nh} \sum_{j=1}^n K\left(\frac{\tau_j - t_0}{h}\right) \left(Y_j - \left\langle \sum_{i=0}^{\ell} S_i p_i\left(\frac{\tau_j - t_0}{h}\right), X_j \right\rangle\right)^2 + \varepsilon \|S\|_1. \quad (1.1)$$

\widehat{S}^h naturally induces a local polynomial estimator of order ℓ around t_0 :

$$\widehat{S}^h(\tau) := \sum_{i=0}^{\ell} \widehat{S}_i^h p_i\left(\frac{\tau - t_0}{h}\right) \mathbb{I}\left\{\left|\frac{\tau - t_0}{h}\right| \leq 1\right\}. \quad (1.2)$$

The point estimate at t_0 is given by

$$\widehat{S}^h(t_0) := \sum_{i=0}^{\ell} \widehat{S}_i^h p_i(0). \quad (1.3)$$

In the following theorem, we establish an upper bound on the point-wise risk of $\widehat{S}^h(t_0)$ when $A(t)$ is in the Hölder class $\Sigma(\beta, L)$ with $\ell = \lfloor \beta \rfloor$.

Theorem 1.1. *Under model (1.1), let (τ_j, X_j, Y_j) , $j = 1, \dots, n$ be i.i.d. copies of the random triplet (τ, X, Y) with X uniformly distributed in \mathcal{X} , τ uniformly distributed in $[0, 1]$, X and τ are independent, and $|Y| \leq a$, a.s. for some constant $a > 0$. Let A be a matrix valued function satisfying A1, A2, A3, and A4. Denote $\Phi = \max_{i=0, \dots, \ell} \|\sqrt{K} p_i\|_{\infty}$, and $\ell = \lfloor \beta \rfloor$.*

Take

$$\widehat{h}_n = C_1 \left(\frac{(\ell^3 (\ell!)^2 \Phi^2 R(T)^2 a^2 m r \log n)}{L^2 n} \right)^{\frac{1}{2\beta+1}}, \quad \varepsilon = D \ell a \Phi \sqrt{\frac{\log 2m}{nm \widehat{h}_n}},$$

for some numerical constants C_1 and D . Then for any $\widehat{h}_n \leq t_0 \leq 1 - \widehat{h}_n$, the estimator defined in (1.3) satisfies with probability at least $1 - \frac{1}{n^{mr}}$,

$$\frac{1}{m^2} \left\| \widehat{S}^h(t_0) - A(t_0) \right\|_2^2 \leq C_1(a, \Phi, \ell, L) \left(\frac{mr \log n}{n} \right)^{\frac{2\beta}{2\beta+1}}, \quad (1.4)$$

where $C_1(a, \Phi, \ell, L)$ is a constant depending on a, Φ, ℓ and L .

One should notice that when $\beta \rightarrow \infty$, bound (1.4) coincides with similar result in classical matrix completion. In section 3, we prove that bound (1.4) is minimax optimal up to a logarithmic factor.

Proof. Firstly, we introduce a sharp oracle inequality of "locally integrated risk" of estimator (1.2) in the following lemma.

Lemma 1. Assume that the condition of Theorem 1.1 holds. Then there exist a numerical constants $D > 0$ such that for all

$$\varepsilon \geq D(\ell + 1)R(T)\Phi a \left(\sqrt{\frac{\log 2m}{nmh}} \vee \frac{(\log 2m)\Phi}{nh} \right),$$

and for arbitrary $\eta > 0$, the estimator (1.3) satisfies with probability at least $1 - e^{-\eta}$

$$\begin{aligned} & \frac{1}{h} \mathbb{E} K \left(\frac{\tau - t_0}{h} \right) \left\langle A(\tau) - \widehat{S}^h(\tau), X \right\rangle^2 \\ & \leq \inf_{S \in \mathbb{D}} \left\{ \frac{1}{h} \mathbb{E} K \left(\frac{\tau - t_0}{h} \right) \left\langle A(\tau) - S(\tau), X \right\rangle^2 \right. \\ & \quad \left. + \frac{D^2(\ell + 1)^2 \Phi^2 R(T)^2 a^2 (\text{rank}(S)m \log 2m + \eta)}{nh} \right\}. \end{aligned} \quad (1.5)$$

where $S(\tau) := \sum_{i=0}^{\ell} S_i p_i \left(\frac{\tau - t_0}{h} \right)$.

The proof of Lemma 1 can be derived from Theorem 19.1 in [22], see Appendix. To be more specific, one just needs to rewrite (1.1) as

$$\widehat{S}^h = \arg \min_{S \in \mathbb{D}} \frac{1}{n} \sum_{j=1}^n \left(\tilde{Y}_j - \left\langle S, \tilde{X}_j \right\rangle \right)^2 + \varepsilon \|S\|_1. \quad (1.6)$$

where

$$\tilde{X}_j = \text{Diag} \left[\sqrt{\frac{1}{h} K \left(\frac{\tau_j - t_0}{h} \right)} p_0 \left(\frac{\tau_j - t_0}{h} \right) X_j, \dots, \sqrt{\frac{1}{h} K \left(\frac{\tau_j - t_0}{h} \right)} p_\ell \left(\frac{\tau_j - t_0}{h} \right) X_j \right],$$

and $\tilde{Y}_j = \sqrt{\frac{1}{h} K \left(\frac{\tau_j - t_0}{h} \right)} Y_j$. Then Lemma 1 is just application of Theorem 19.1 in [22].

Consider

$$\begin{aligned}
& \frac{1}{h} \mathbb{E} K\left(\frac{\tau - t_0}{h}\right) \left\langle A(\tau) - \sum_{i=0}^{\ell} \widehat{S}_i^h p_i\left(\frac{\tau - t_0}{h}\right), X \right\rangle^2 \\
&= \frac{1}{h} \mathbb{E} K\left(\frac{\tau - t_0}{h}\right) \left\langle A(\tau) - \sum_{i=0}^{\ell} S_i p_i\left(\frac{\tau - t_0}{h}\right) + \sum_{i=0}^{\ell} (S_i - \widehat{S}_i^h) p_i\left(\frac{\tau - t_0}{h}\right), X \right\rangle^2 \\
&= \frac{1}{h} \mathbb{E} K\left(\frac{\tau - t_0}{h}\right) \left\langle \sum_{i=0}^{\ell} (S_i - \widehat{S}_i^h) p_i\left(\frac{\tau - t_0}{h}\right), X \right\rangle^2 \\
&+ \frac{1}{h} \mathbb{E} K\left(\frac{\tau - t_0}{h}\right) \left\langle A(\tau) - \sum_{i=0}^{\ell} S_i p_i\left(\frac{\tau - t_0}{h}\right), X \right\rangle^2 \\
&+ \frac{2}{h} \mathbb{E} K\left(\frac{\tau - t_0}{h}\right) \left\langle A(\tau) - \sum_{i=0}^{\ell} S_i p_i\left(\frac{\tau - t_0}{h}\right), X \right\rangle \left\langle \sum_{i=0}^{\ell} (S_i - \widehat{S}_i^h) p_i\left(\frac{\tau - t_0}{h}\right), X \right\rangle
\end{aligned} \tag{1.7}$$

Therefore, from (1.5) and (1.7), we have for any $S \in \mathbb{D}$

$$\begin{aligned}
& \frac{1}{h} \mathbb{E} K\left(\frac{\tau - t_0}{h}\right) \left\langle \sum_{i=0}^{\ell} (S_i - \widehat{S}_i^h) p_i\left(\frac{\tau - t_0}{h}\right), X \right\rangle^2 \\
&\leq \frac{2}{h} \mathbb{E} K\left(\frac{\tau - t_0}{h}\right) \left| \left\langle A(\tau) - \sum_{i=0}^{\ell} S_i p_i\left(\frac{\tau - t_0}{h}\right), X \right\rangle \left\langle \sum_{i=0}^{\ell} (S_i - \widehat{S}_i^h) p_i\left(\frac{\tau - t_0}{h}\right), X \right\rangle \right| \\
&+ \frac{D^2(\ell + 1)^2 \Phi^2 R(T)^2 a^2 (\text{rank}(S) m \log 2m + \eta)}{nh}. \\
&\leq \left(\frac{c^4}{c^2 - 1} \right) \frac{1}{h} \mathbb{E} K\left(\frac{\tau - t_0}{h}\right) \left\langle A(\tau) - \sum_{i=0}^{\ell} S_i p_i\left(\frac{\tau - t_0}{h}\right), X \right\rangle^2 \\
&+ \left(\frac{c^2}{c^2 - 1} \right) \left\{ \frac{D^2(\ell + 1)^2 \Phi^2 R(T)^2 a^2 (\text{rank}(S) m \log 2m + \eta)}{nh} \right\},
\end{aligned} \tag{1.8}$$

where we used the fact that for any positive constants a and b , $2ab \leq \frac{1}{c^2} a^2 + c^2 b^2$ for some $c > 1$. Take S such that

$$\sum_{i=0}^{\ell} S_i p_i\left(\frac{\tau - t_0}{h}\right) = A(t_0) + A'(t_0) h \left(\frac{\tau - t_0}{h}\right) + \dots + \frac{A^{(\ell)}(t_0) h^{\ell}}{\ell!} \left(\frac{\tau - t_0}{h}\right)^{\ell}. \tag{1.9}$$

Note that this is possible since the right hand side is a matrix valued polynomial of $\frac{\tau - t_0}{h}$ up to order ℓ , and $\text{span}\{p_0, p_1, \dots, p_{\ell}\} = \text{span}\{1, x, \dots, x^{\ell}\}$. Under the condition that all

entries of $A^{(k)}(t)$ are bounded by a , then entries of S_k are bounded by $R(T)a$. Thus, the corresponding $S \in \mathbb{D}$. Obviously, $\text{rank}(S_i) \leq (\ell+1-i)r$. Since $A \in \Sigma(\beta, L)$, we consider ℓ -th order Taylor expansion of A at t_0 to get

$$A(\tau) = A(t_0) + A'(t_0)(\tau - t_0) + \dots + \frac{\tilde{A}(\tau - t_0)^\ell}{\ell!}, \quad (1.10)$$

where \tilde{A} is the matrix with $\tilde{A}_{ij} = a_{ij}^{(\ell)}(t_0 + \alpha_{ij}(\tau - t_0))$ for some $\alpha_{ij} \in [0, 1]$. Then we apply the Taylor expansion (1.10) and identity (1.9) to get

$$\begin{aligned} & \frac{1}{h} \mathbb{E} K\left(\frac{\tau - t_0}{h}\right) \left\langle A(\tau) - \sum_{i=0}^{\ell} S_i p_i\left(\frac{\tau - t_0}{h}\right), X \right\rangle^2 \\ & \leq \frac{1}{h} \mathbb{E} K\left(\frac{\tau - t_0}{h}\right) \frac{1}{m^2} \left\| \frac{LU(\tau - t)^\beta}{\ell!} \right\|_2^2 \leq \frac{L^2 h^{2\beta}}{(\ell!)^2}. \end{aligned} \quad (1.11)$$

where U denotes the matrix with all entries being 1. The first inequality is due to $a_{ij} \in \Sigma(\beta, L)$, and the second is due to $|\tau - t_0| \leq h$. Under the condition that X is uniformly distributed in \mathcal{X} , and the orthogonality of $\{p_i\}_{i=0}^{\ell}$, it is easy to check that

$$\frac{1}{h} \mathbb{E} K\left(\frac{\tau - t_0}{h}\right) \left\langle \sum_{i=0}^{\ell} (S_i - \hat{S}_i^h) p_i\left(\frac{\tau - t_0}{h}\right), X \right\rangle^2 = \frac{1}{m^2} \sum_{i=0}^{\ell} \|\hat{S}_i^h - S_i\|_2^2 \quad (1.12)$$

Note that

$$\|\hat{S}^h(t_0) - S(t_0)\|_2^2 = \left\| \sum_{i=0}^{\ell} (\hat{S}_i^h - S_i) p_i(0) \right\|_2^2 \leq (\ell+1) \Phi^2 \sum_{i=0}^{\ell} \|\hat{S}_i^h - S_i\|_2^2, \quad (1.13)$$

where the second inequality is due to Cauchy-Schwarz inequality and p_i are uniformly bounded on $[-1, 1]$. Combining (1.8), (1.11), (1.12), and (1.13), we get with probability at least $1 - e^{-\eta}$

$$\begin{aligned} \frac{1}{m^2} \|\hat{S}^h(t_0) - A(t_0)\|_2^2 & \leq \left(\frac{c^4}{c^2 - 1} \right) \frac{2L^2 h^{2\beta}}{(\ell!)^2} \\ & \quad + \left(\frac{c^2}{c^2 - 1} \right) \left\{ \frac{D^2(\ell+1)^2 \Phi^2 R(T)^2 a^2 (\text{rank}(S)m \log 2m + \eta)}{nh} \right\}, \end{aligned}$$

By optimizing the right hand side with respect to h and take $\eta = mr \log n$, we take

$$\hat{h}_n = C \left(\frac{\ell^3 (\ell!)^2 \Phi^2 R(T)^2 a^2 mr \log n}{L^2 n} \right)^{\frac{1}{2\beta+1}}.$$

where C is a numerical constant. This completes the proof of the theorem. \square

2 Global estimators and upper bounds on integrated risk

In this section, we propose two global estimators and study their integrated risk measured by L_2 -norm and L_∞ -norm.

2.1 From localization to globalization

Firstly, we construct a global estimator based on (1.2). Take

$$\hat{h}_n = C_1 \left(\frac{\ell^3 (\ell!)^2 \Phi^2 R(T)^2 a^2 mr \log n}{L^2 n} \right)^{\frac{1}{2\beta+1}}, \quad M = \lceil 1/\hat{h}_n \rceil.$$

Without loss of generality, assume that M is even. Denote $\hat{S}_k^h(t)$ the local polynomial estimator around t_{2k-1} as in (1.2) by using orthogonal polynomials with $K = \mathbb{I}_{[-1,1]}$, where $t_{2k-1} = \frac{2k-1}{M}$, $k = 1, 2, \dots, M/2$ and \mathbb{I} is the indicator function. Denote

$$\hat{A}(t) = \sum_{k=1}^{M/2} \hat{S}_k^h(t) \mathbb{I}_{(t_{2k-1}-\hat{h}_n, t_{2k-1}+\hat{h}_n]}, \quad t \in (0, 1). \quad (2.1)$$

Note that the weight function K is not necessary to be $\mathbb{I}_{[-1,1]}$. It can be replaced by any K that satisfies $K \geq K_0 > 0$ on $[-1, 1]$. The following result characterizes the global performance of estimator (2.1) under matrix completion setting measured by L_2 -norm.

Theorem 2.1. *Assume that the conditions of Theorem 1.1 hold, and let \hat{A} be an estimator defined as in (2.1). Then with probability at least $1 - \frac{1}{n^{mr-1}}$,*

$$\frac{1}{m^2} \int_0^1 \|\hat{A}(t) - A(t)\|_2^2 dt \leq C_2(a, \Phi, \ell, L) \left(\frac{mr \log n}{n} \right)^{\frac{2\beta}{2\beta+1}}, \quad (2.2)$$

where $C_2(a, \Phi, \ell, L)$ is a constant depending on a, Φ, ℓ, L .

Compared with the integrated risk measured by L_2 -norm of real valued functions in Hölder class, the result in (2.2) has an excess $\log n$ term, which is introduced by the matrix Bernstein inequality, see [23]. In section 3, we show that bound (2.2) is minimax optimal up to a logarithmic factor.

Proof. It is easy to see that

$$\int_0^1 \|\widehat{A}(t) - A(t)\|_2^2 dt \leq \sum_{k=1}^{M/2} \int_{t_{2k-1}-\widehat{h}_n}^{t_{2k-1}+\widehat{h}_n} \|\widehat{S}_k^h(t) - A(t)\|_2^2 dt. \quad (2.3)$$

For each k ,

$$\frac{1}{m^2} \int_{t_{2k-1}-\widehat{h}_n}^{t_{2k-1}+\widehat{h}_n} \|\widehat{S}_k^h(t) - A(t)\|_2^2 dt = \mathbb{E}_{\tau, X} \mathbb{I}_{(t_{2k-1}-\widehat{h}_n, t_{2k-1}+\widehat{h}_n]} \left\langle A(\tau) - \widehat{S}^h(\tau), X \right\rangle^2$$

By (1.5), (1.11) and arguments used to prove Theorem 1.1, we have with probability at least $1 - \frac{1}{n^{mr}}$,

$$\frac{1}{m^2 \widehat{h}_n} \int_{t_{2k-1}-\widehat{h}_n}^{t_{2k-1}+\widehat{h}_n} \|\widehat{S}_k^h(t) - A(t)\|_2^2 dt \leq C_1(a, \Phi, \ell, L) \left(\frac{mr \log n}{n} \right)^{\frac{2\beta}{2\beta+1}}.$$

We take the union bound over k , from (2.3) we get with probability at least $1 - \frac{1}{n^{mr-1}}$,

$$\frac{1}{m^2} \int_0^1 \|\widehat{A}(t) - A(t)\|_2^2 dt \leq C_2(a, \Phi, \ell, L) \left(\frac{mr \log n}{n} \right)^{\frac{2\beta}{2\beta+1}}.$$

where $C_2(a, \Phi, \ell, L)$ is a constant depending on a, Φ, ℓ, L . □

2.2 Bias reduction through higher order kernels

If $A(t)$ is not necessarily low rank, we propose an estimator which is easy to implement and prove an upper bound on its risk measured by L_∞ -norm. Such estimators are related to another popular approach parallel to local polynomial estimators for bias reduction, namely,

using high order kernels to reduce bias. They can also be applied to another important technique of low rank estimation or approximation via singular value thresholding, see [4] and [12]. The estimator proposed by [9] is shown to be equivalent to soft singular value thresholding of such type of estimators.

The kernels we are interested in satisfy the following conditions:

K1 $K(\cdot)$ is symmetric, i.e. $K(u) = K(-u)$.

K2 $K(\cdot)$ is compactly supported on $[-1, 1]$.

K3 $R_K = \int_{-\infty}^{\infty} K^2(u) du < \infty$.

K4 $K(\cdot)$ is of order ℓ for some $\ell \in \mathbb{N}^*$.

K5 $K(\cdot)$ is Lipschitz continuous with $0 < L_K < \infty$.

Consider

$$\tilde{A}(t) = \frac{m^2}{nh} \sum_{j=1}^n K\left(\frac{\tau_j - t}{h}\right) Y_j X_j. \quad (2.4)$$

Note that when $K \geq 0$, (2.4) is the solution to the following optimization problem

$$\tilde{A}(t) = \arg \min_{S \in \mathbb{D}} \frac{1}{nh} \sum_{j=1}^n K\left(\frac{\tau_j - t}{h}\right) (Y_j - \langle S, X_j \rangle)^2. \quad (2.5)$$

In the following theorem we prove an upper bound on its global performance measured by L_∞ -norm over $\Sigma(\beta, L)$ which is much harder to obtain for matrix lasso problems.

Theorem 2.2. *Under model (1.1), let (τ_j, X_j, Y_j) , $j = 1, \dots, n$ be i.i.d. copies of the random triplet (τ, X, Y) with X uniformly distributed in \mathcal{X} , τ uniformly distributed in $[0, 1]$, X and τ are independent, and $|Y| \leq a$ a.s. for some constant $a > 0$; Kernel K satisfies K1-K5; $A(t)$ be any matrix valued function satisfying A1 and A4. Denote $\ell = \lfloor \beta \rfloor$. Take*

$$\tilde{h}_n := c_*(K) \left(\frac{a^2(\ell!)^2 m \log n}{2\beta L^2 n} \right)^{\frac{1}{2\beta+1}}, \quad (2.6)$$

Then with probability at least $1 - n^{-2}$, the estimator defined in (2.4) satisfies

$$\sup_{t \in [\tilde{h}_n, 1 - \tilde{h}_n]} \frac{1}{m^2} \|\tilde{A}(t) - A(t)\|^2 \leq C^*(K) \left(\frac{a^2(\ell!)^2 m \log n}{2\beta L^2 n} \right)^{\frac{2\beta}{2\beta+1}}, \quad (2.7)$$

where $C^*(K)$ and $c_*(K)$ are constants depending on K .

When the smoothness parameter β tends to infinity, bound (2.7) coincides with similar bounds in classical matrix completion, which is $O(\frac{m^3 \log n}{n})$. When m degenerates to 1, the bound coincides with that of real valued case, which is $O((\frac{\log n}{n})^{2\beta/(2\beta+1)})$. In section 3, we show that this bound is minimax optimal up to a logarithmic factor.

Proof. In this proof, we use $C(K)$ to denote any constant depending on K which may vary from place to place. This simplifies the representation while does no harm to the soundness of our proof. Consider

$$\sup_{t \in [\tilde{h}_n, 1 - \tilde{h}_n]} \|\tilde{A}(t) - A(t)\| \leq \sup_{t \in [\tilde{h}_n, 1 - \tilde{h}_n]} \|\tilde{A}(t) - \mathbb{E}\tilde{A}(t)\| + \sup_{t \in [\tilde{h}_n, 1 - \tilde{h}_n]} \|\mathbb{E}\tilde{A}(t) - A(t)\|. \quad (2.8)$$

The first term on the right hand side is recognized as the variance and the second is the bias. Firstly, we deal with the bias term. Denote $B(t_0) := \mathbb{E}\tilde{A}(t_0) - A(t_0)$, $t_0 \in [\tilde{h}_n, 1 - \tilde{h}_n]$. Recall from (1.2), $\mathbb{E}(\xi_j | \tau_j, X_j) = 0$ for any $t_0 \in [\tilde{h}_n, 1 - \tilde{h}_n]$ we have

$$\mathbb{E}\tilde{A}(t_0) = \mathbb{E} \frac{m^2}{nh} \sum_{j=1}^n K\left(\frac{\tau_j - t_0}{h}\right) (\langle A(\tau_j), X_j \rangle + \xi_j) X_j = \frac{m^2}{h} \mathbb{E} K\left(\frac{\tau - t_0}{h}\right) \langle A(\tau), X \rangle X.$$

By applying the Taylor expansion of $A(\tau)$ as in (1.10) and the fact that K is a kernel of order ℓ , we get

$$\mathbb{E}\tilde{A}(t_0) = \mathbb{E} \frac{m^2}{h} K\left(\frac{\tau - t_0}{h}\right) \langle A(t_0), X \rangle X + \mathbb{E} \frac{m^2}{h} K\left(\frac{\tau - t_0}{h}\right) \frac{(\tau - t_0)^\ell}{\ell!} \langle \tilde{A}, X \rangle X,$$

where \tilde{A} is the same as in (1.10). It is easy to check that the first term on the right hand side

is $A(t_0)$. Therefore we rewrite $B(t_0)$ as

$$B(t_0) = \mathbb{E} \frac{m^2}{h} K\left(\frac{\tau - t_0}{h}\right) \frac{(\tau - t_0)^\ell}{\ell!} \langle \tilde{A}, X \rangle X = \mathbb{E} \frac{m^2}{h} K\left(\frac{\tau - t_0}{h}\right) \frac{(\tau - t_0)^\ell}{\ell!} \langle \tilde{A} - A^{(\ell)}(t_0), X \rangle X,$$

where the second equality is due to the fact that each element of $A(t)$ is in $\Sigma(\beta, L)$ and K is a kernel of order ℓ . Then we can bound each element of matrix $B(t_0)$ as

$$\begin{aligned} |B_{ij}(t_0)| &\leq \int_0^1 \frac{1}{h} K\left(\frac{\tau - t_0}{h}\right) \frac{|\tau - t_0|^\ell}{\ell!} |a_{ij}^{(\ell)}(t_0 + \alpha(\tau - t_0)) - a_{ij}^{(\ell)}(t_0)| d\tau \\ &\leq L \int_0^1 |K(u)| \frac{|uh|^\beta}{\ell!} du \\ &\leq C(K) \frac{Lh^\beta}{\ell!}. \end{aligned}$$

Thus

$$\sup_{t \in [\tilde{h}_n, 1 - \tilde{h}_n]} \|B(t)\| \leq C(K) \frac{Lmh^\beta}{\ell!}. \quad (2.9)$$

On the other hand, for the variance term $\sup_{t \in [\tilde{h}_n, 1 - \tilde{h}_n]} \|\tilde{A}(t) - \mathbb{E}\tilde{A}(t)\|^2$, we construct a δ -net on the interval $[0, 1]$ with $\delta = 1/M$, and

$$M = n^2, \quad t_j = \frac{2j-1}{2M}, \quad j = 1, \dots, M.$$

Denote $S_n(t) := \tilde{A}(t) - \mathbb{E}\tilde{A}(t)$, then we have

$$\sup_{t \in [\tilde{h}_n, 1 - \tilde{h}_n]} \|S_n(t)\| \leq \sup_{t \in [0, 1]} \|S_n(t)\| \leq \max_i \|S_n(t_i)\| + \sup_{|t-t'| \leq \delta} \|S_n(t) - S_n(t')\|. \quad (2.10)$$

Next, we bound both terms on the right hand side respectively. For each t_i ,

$$S_n(t_i) = \frac{m^2}{nh} \sum_{j=1}^n \left(K\left(\frac{\tau_j - t_i}{h}\right) Y_j X_j - \mathbb{E} K\left(\frac{\tau_j - t_i}{h}\right) Y_j X_j \right).$$

The right hand side is a sum of zero mean random matrices, we apply the matrix Bernstein inequality, see [23]. Under the assumption of Theorem 2.2, one can easily check that with

probability at least $1 - e^{-\eta}$,

$$\|S_n(t_i)\| \leq C(K)m^2 \left(\sqrt{\frac{a^2(\eta + \log 2m)}{mnh}} \vee \frac{a(\eta + \log 2m)}{nh} \right).$$

Indeed, by setting $\bar{X} = \frac{m^2}{h}K\left(\frac{\tau-t}{h}\right)YX - \mathbb{E}\frac{m^2}{h}K\left(\frac{\tau-t}{h}\right)YX$, it is easy to check that $U_{\bar{X}} \lesssim \|K\|_{\infty}am^2/h$ and $\sigma_{\bar{X}}^2 \lesssim R_K a^2 m^3/h$. By taking the union bound over all i and setting $\eta = 4 \log n$, we get with probability at least $1 - n^{-2}$,

$$\max_i \|S_n(t_i)\|^2 \leq C(K) \frac{a^2 m^3 \log n}{nh},$$

As for the second term on the right hand side of (2.10), by the assumption that K is a Lipschitz function with Lipschitz constant L_K , we have

$$\begin{aligned} \sup_{|t-t'|\leq\delta} \|S_n(t) - S_n(t')\| &\leq \sup_{|t-t'|\leq\delta} \|(\tilde{A}(t) - \tilde{A}(t'))\| + \sup_{|t-t'|\leq\delta} \|\mathbb{E}(\tilde{A}(t) - \tilde{A}(t'))\| \\ &\leq \frac{L_K am^3}{n^2 h^2} + \frac{L_K am}{n^2 h^2}. \end{aligned}$$

Thus with probability at least $1 - n^{-2}$,

$$\sup_{t \in [\tilde{h}_n, 1-\tilde{h}_n]} \|S_n(t)\|^2 \leq C(K) \frac{a^2 m^3 \log n}{nh}$$

Together with the upper bound we get on the bias in (2.9), we have with probability at least $1 - n^{-2}$,

$$\sup_{t \in [\tilde{h}_n, 1-\tilde{h}_n]} \frac{1}{m^2} \|\tilde{A}(t) - A(t)\|^2 \leq C(K) \left(\frac{a^2 m \log n}{nh} + \frac{L^2 h^{2\beta}}{\ell!^2} \right).$$

Choose

$$\tilde{h}_n = C(K) \left(\frac{a^2 (\ell!)^2 m \log n}{2\beta L^2 n} \right)^{\frac{1}{2\beta+1}},$$

we get

$$\sup_{t \in [\tilde{h}_n, 1-\tilde{h}_n]} \frac{1}{m^2} \|\tilde{A}(t) - A(t)\|^2 \leq C(K) \left(\frac{a^2 (\ell!)^2 m \log n}{2\beta L^2 n} \right)^{\frac{2\beta}{2\beta+1}}.$$

□

3 Lower bounds under matrix completion setting

In this section, we prove the minimax lower bound of estimators (1.3), (2.1) and (2.4). In the realm of classical low rank matrix estimation, [11] studied the optimality issue measured by the Frobenius norm on the classes defined in terms of a "spikeness index" of the true matrix; [10] derived optimal rates in noisy matrix completion on different classes of matrices for the empirical prediction error; [9] established that the rates of the estimator they propose under matrix completion setting are optimal up to a logarithmic factor measured by the Frobenius norm. Based on the ideas of [9], standard methods to prove minimax lower bounds in real valued case in [24], and some fundamental results in coding theory, we establish the corresponding minimax lower bounds of (1.4), (2.2) and (2.7) which essentially shows that the upper bounds we get are all optimal up to some logarithmic factor.

For the convenience of representation, we denote by $\inf_{\hat{A}}$ the infimum over all estimators \hat{A} of A . We denote by $\mathcal{A}(r, a)$ the set of matrix valued functions satisfying A1, A2, A3, and A4. We denote by $\mathcal{P}(r, a)$ the class of distributions of random triplet (τ, X, Y) that satisfies model (1.1) with any $A \in \mathcal{A}(r, a)$.

Theorem 3.1. *Under model (1.1), let (τ_j, X_j, Y_j) , $j = 1, \dots, n$ be i.i.d. copies of the random triplet (τ, X, Y) with X uniformly distributed in \mathcal{X} , τ uniformly distributed in $[0, 1]$, X and τ are independent, and $|Y| \leq a$, a.s. for some constant $a > 0$; let A be any matrix valued function in $\mathcal{A}(r, a)$. Then there is an absolute constant $\eta \in (0, 1)$ such that for all $t_0 \in [0, 1]$*

$$\inf_{\hat{A}} \sup_{P_{\tau, X, Y}^A \in \mathcal{P}(r, a)} \mathbb{P}_{P_{\tau, X, Y}^A} \left\{ \frac{1}{m^2} \|\hat{A}(t_0) - A(t_0)\|_2^2 > C(\beta, L, a) \left(\frac{mr}{n} \right)^{\frac{2\beta}{2\beta+1}} \right\} \geq \eta. \quad (3.1)$$

where $C(\beta, L, a)$ is a constant depending on β , L and a .

Note that compared with the upper bound (1.4), the lower bound (3.1) matches it that

up to a logarithmic factor. As a consequence, it shows that the estimator (1.3) achieves a near optimal minimax rate of pointwise estimation. Although, the result of Theorem 3.1 is under bounded response condition, it can be readily extended to the case when the noise in (1.2) is Gaussian.

Proof. Without loss of generality, we assume that both m and r are even numbers. We introduce several notations which are key to construct the hypothesis set. For some constant $\gamma > 0$, denote

$$\mathcal{C} = \{\tilde{A} = (a_{ij}) \in \mathbb{C}^{\frac{m}{2} \times \frac{r}{2}} : a_{ij} \in \{0, \gamma\}, \forall 1 \leq i \leq m/2, 1 \leq j \leq r/2\},$$

and consider the set of block matrices

$$\mathcal{B}(\mathcal{C}) = \left\{ \begin{bmatrix} \tilde{A} & \tilde{A} & \dots & \tilde{A} & O \end{bmatrix} \in \mathbb{C}^{\frac{m}{2} \times \frac{m}{2}} : \tilde{A} \in \mathcal{C} \right\}, \quad (3.2)$$

where O denotes the $m/2 \times (m/2 - r\lfloor m/r \rfloor/2)$ zero matrix. Then we consider a subset of Hermitian matrices $\mathcal{S}_m \subset \mathbb{H}_m$,

$$\mathcal{S}_m = \left\{ \begin{bmatrix} \tilde{O} & \hat{A} \\ \hat{A}^* & \tilde{O} \end{bmatrix} \in \mathbb{C}^{m \times m} : \hat{A} \in \mathcal{B}(\mathcal{C}) \right\}. \quad (3.3)$$

An immediate observation is that for any matrix $A \in \mathcal{S}_m$, $\text{rank}(A) \leq r$.

Due to the Varshamov-Gilbert bound (see Lemma 2.9 in [24]), there exists a subset $\mathcal{A}^0 \subset \mathcal{S}_m$ with cardinality $\text{Card}(\mathcal{A}^0) \geq 2^{mr/32} + 1$ containing the zero $m \times m$ matrix $\mathbf{0}$ such that for any two distinct elements A_1 and A_2 of \mathcal{A}^0 ,

$$\|A_1 - A_2\|_2^2 \geq \frac{mr}{16} \lfloor \frac{m}{r} \rfloor \gamma^2 \geq \gamma^2 \frac{m^2}{32}. \quad (3.4)$$

Let $f_n(t)$ denote the function $f_n(t) := Lh_n^\beta f\left(\frac{t-t_0}{h_n}\right)$, $t \in [0, 1]$, where $h_n = c_0 \left(\frac{mr}{n}\right)^{\frac{1}{2\beta+1}}$,

with some constant $c_0 > 0$, and $f \in \Sigma(\beta, 1/2) \cap C^\infty$ and $\text{Supp}(f) = [-1/2, 1/2]$. Note that there exist functions f satisfying this condition. For instance, one can take

$$f(t) = \alpha e^{-\frac{1}{1-4u^2}} \mathbb{I}(|u| < 1/2), \quad (3.5)$$

for some sufficient small $\alpha > 0$. It is easy to check that $f_n(t) \in \Sigma(\beta, L)$ on $[0, 1]$.

We consider the following hypotheses of A at t_0 :

$$\mathcal{A}_0^\beta := \{\hat{A}(t) = Af_n(t), t \in [0, 1] : A \in \mathcal{A}^0\}.$$

The following claims are easy to check: firstly, any element in \mathcal{A}_0^β together with its derivative have rank uniformly bounded by r , and the difference of any two elements of \mathcal{A}_0^β satisfies the same property for fixed t_0 ; secondly, the entries of any element of \mathcal{A}_0^β together with its derivative are uniformly bounded by some constant for sufficiently small chosen γ ; finally, each element of $A(t) \in \mathcal{A}_0^\beta$ belongs to $\Sigma(\beta, L)$. Therefore, $\mathcal{A}_0^\beta \subset \mathcal{A}(r, a)$ with some chosen γ .

According to (3.4), for any two distinct elements $\hat{A}_1(t)$ and $\hat{A}_2(t)$ of \mathcal{A}_0^β , the difference between $\hat{A}_1(t)$ and $\hat{A}_2(t)$ at point t_0 is given by

$$\|\hat{A}_1(t_0) - \hat{A}_2(t_0)\|_2^2 \geq \frac{\gamma^2 L^2 c_0^{2\beta} f^2(0)}{32} m^2 \left(\frac{mr}{n}\right)^{\frac{2\beta}{2\beta+1}}. \quad (3.6)$$

On the other hand, we consider the joint distributions $P_{\tau, X, Y}^A$ such that $\tau \sim U[0, 1]$, $X \sim \Pi_0$ where Π_0 denotes the uniform distribution on \mathcal{X} , τ and X are independent, and

$$\mathbb{P}_A(Y|\tau, X) = \begin{cases} \frac{1}{2} + \frac{\langle A(\tau), X \rangle}{4a}, & Y = a, \\ \frac{1}{2} - \frac{\langle A(\tau), X \rangle}{4a}, & Y = -a. \end{cases}$$

One can easily check that as long as $A(\tau) \in \mathcal{A}_0^\beta$, such $P_{\tau, X, Y}^A$ belongs to the distribution

class $\mathcal{P}(r, a)$. We denote the corresponding n -product probability measure by \mathbb{P}_A . Then for any $A(\tau) \in \mathcal{A}_0^\beta$, the Kullback-Leibler Divergence between \mathbb{P}_0 and \mathbb{P}_A is

$$K(\mathbb{P}_0, \mathbb{P}_A) = n\mathbb{E}\left(p_0(\tau, X) \log \frac{p_0(\tau, X)}{p_A(\tau, X)} + (1 - p_0(\tau, X)) \log \frac{1 - p_0(\tau, X)}{1 - p_A(\tau, X)}\right),$$

where $p_A(\tau, X) = 1/2 + \langle A(\tau), X \rangle / 4a$. Note that $\mathbb{P}_A(Y = a | \tau, X) \in [1/4, 3/4]$ is guaranteed provided that $|\langle A(t), X \rangle| \leq a$. Thus by the inequality $-\log(1+u) \leq -u + u^2/2$, $\forall u > -1$, and the fact that $\mathbb{P}_A(Y = a | \tau, X) \in [1/4, 3/4]$, we have

$$K(\mathbb{P}_0, \mathbb{P}_A) \leq n\mathbb{E}2(p_0(\tau, X) - p_A(\tau, X))^2 \leq \frac{n}{8a^2}\mathbb{E}\langle A(\tau), X \rangle^2.$$

Recall that $A(\tau) = Af_n(\tau) \in \mathcal{A}_0^\beta$, by $\tau \sim U[0, 1]$ and $X \sim \Pi_0$, we have

$$K(\mathbb{P}_0, \mathbb{P}_A) \leq \frac{n}{8a^2} \frac{1}{m^2} L^2 \|f\|_2^2 h_n^{2\beta+1} m^2 \gamma^2 \leq \frac{L^2 \|f\|_2^2 c_0^{2\beta+1} \gamma^2}{8a^2} mr. \quad (3.7)$$

Therefore, provided the fact that $\text{Card}(\mathcal{A}^0) \geq 2^{mr/32} + 1$, together with (3.7), we have

$$\frac{1}{\text{Card}(\mathcal{A}_0^\beta) - 1} \sum_{A \in \mathcal{A}_0^\beta} K(\mathbb{P}_0, \mathbb{P}_A) \leq \alpha \log(\text{Card}(\mathcal{A}_0^\beta) - 1) \quad (3.8)$$

is satisfied for any $\alpha > 0$ if γ is chosen as a sufficiently small constant. In view of (3.6) and (3.8), the lower bound (3.1) follows from Theorem 2.5 in [24]. \square

Theorem 3.2. *Under model (1.1), let (τ_j, X_j, Y_j) , $j = 1, \dots, n$ be i.i.d. copies of the random triplet (τ, X, Y) with X uniformly distributed in \mathcal{X} , τ uniformly distributed in $[0, 1]$, X and τ are independent, and $|Y| \leq a$, a.s. for some constant $a > 0$; let A be any matrix valued function in $\mathcal{A}(r, a)$. Then there is an absolute constant $\eta \in (0, 1)$ such that*

$$\inf_{\hat{A}} \sup_{P_{\tau, X, Y}^A \in \mathcal{P}(r, a)} \mathbb{P}_{P_{\tau, X, Y}^A} \left\{ \frac{1}{m^2} \int_0^1 \|\hat{A}(t) - A(t)\|_2^2 dt > \tilde{C}(\beta, L, a) \left(\frac{mr}{n} \right)^{\frac{2\beta}{2\beta+1}} \right\} \geq \eta, \quad (3.9)$$

where $\tilde{C}(\beta, L, a)$ is a constant depending on L, β and a .

The lower bound in (3.9) matches the upper bound we get in (2.2) up to a logarithmic factor. Therefore, our estimator (2.1) achieves a near optimal minimax rate on the integrated risk measured by L_2 -norm up to a logarithmic factor. The result of Theorem 3.2 can be readily extended to the case when the noise in (1.2) is Gaussian.

Proof. Without loss of generality, we assume that both m and r are even numbers. Take a real number $c_1 > 0$, define

$$M = \left\lceil c_1 \left(\frac{n}{mr} \right)^{\frac{1}{2\beta+1}} \right\rceil, \quad h_n = \frac{1}{2M}, \quad t_j = \frac{2j-1}{2M},$$

and

$$\phi_j(t) = L h_n^\beta f\left(\frac{t - t_j}{h_n}\right), \quad j = 1, \dots, M, \quad t \in [0, 1],$$

where f is defined the same as in (3.5). Meanwhile, we consider the set of all binary sequences of length M : $\Omega = \{\omega = (\omega_1, \dots, \omega_M), \omega_i \in \{0, 1\}\} = \{0, 1\}^M$. By Varshamov-Gilbert bound, there exists a subset $\Omega_0 = \{\omega^0, \dots, \omega^N\}$ of Ω such that $\omega^0 = (0, \dots, 0) \in \Omega_0$, and $d(\omega^j, \omega^k) \geq \frac{M}{8}$, $\forall 0 \leq j < k \leq N$, and $N \geq 2^{\frac{M}{8}}$, where $d(\cdot, \cdot)$ denotes the Hamming distance of two binary sequences. Then we define a collection of functions based on Ω_0 : $\mathcal{E} = \left\{ f_\omega(t) = \sum_{j=1}^M \omega_j \phi_j(t) : \omega \in \Omega_0 \right\}$. From the result of Varshamov-Gilbert bound, we know that $S := \text{Card}(\mathcal{E}) = \text{Card}(\Omega_0) \geq 2^{\frac{M}{8}} + 1$. It is also easy to check that for all $f_\omega, f_{\omega'} \in \mathcal{E}$,

$$\begin{aligned} \int_0^1 (f_\omega(t) - f_{\omega'}(t))^2 dt &= \sum_{j=1}^M (\omega_j - \omega'_j)^2 \int_{\Delta_j} \phi_j^2(t) dt \\ &= L^2 h_n^{2\beta+1} \|f\|_2^2 \sum_{j=1}^M (\omega_j - \omega'_j)^2 \\ &\geq L^2 h_n^{2\beta} \|f\|_2^2 / 16, \end{aligned} \tag{3.10}$$

where $\Delta_j = [(j-1)/M, j/M]$.

In what follows, we combine two fundamental results in coding theory: one is Varshamov-Gilbert bound ([25, 26]) in its general form of a q -ary code, the other is the volume estimate of Hamming balls. Let $A_q(n, d)$ denote the largest size of a q -ary code of block length n with minimal Hamming distance d .

Proposition 3.3. The maximal size of a q -ary code of block length n with minimal Hamming distance $d = pn$, satisfies

$$A_q(n, d + 1) \geq q^{n(1-h_q(p))}, \quad (3.11)$$

where $p \in [0, 1 - 1/q]$, $h_q(p) = p \log_q(q - 1) - p \log_q p - (1 - p) \log_q(1 - p)$ is the q -ary entropy function.

We now have all the elements needed in hand to construct our hypotheses set. Denote $\Omega_1 = \{\omega^1, \dots, \omega^N\}$, which is a subset of Ω_0 without ω^0 . We then consider a subset \mathcal{E}_1 of \mathcal{E} which is given by $\mathcal{E}_1 := \left\{ f_\omega(t) = \sum_{j=1}^M \omega_j \phi_j(t) : \omega \in \Omega_1 \right\}$. Clearly, $S_1 := \text{Card}(\mathcal{E}_1) \geq 2^{M/8}$. Then we define a new collection of matrix valued functions as

$$\mathcal{C} = \left\{ \tilde{A} = (a_{ij}) \in \mathbb{C}^{\frac{m}{2} \times \frac{r}{2}} : a_{ij} \in \{\delta f_\omega : \omega \in \Omega_1\}, \delta \in \mathbb{C}, \forall 1 \leq i \leq m/2, 1 \leq j \leq r/2 \right\}.$$

Obviously, the collection \mathcal{C} is a S_1 -ary code of block length $mr/4$. Thus, we can apply the result of Proposition 3.3. It is easy to check that for $p = 1/4$, and $q \geq 4$

$$1 - h_q(p) = 1 - p \log_q \frac{q-1}{p} + (1-p) \log_q(1-p) \geq \frac{1}{4}. \quad (3.12)$$

In our case, $q = S_1 \geq 2^{M/8}$ and $n = mr/4$. If we take $p = 1/4$, we know that

$$A_{S_1}(mr/4, mr/16) \geq A_{S_1}(mr/4, mr/16 + 1) \geq S_1^{mr/16}. \quad (3.13)$$

In other words, (3.13) guarantees that there exists a subset $\mathcal{H}^0 \subset \mathcal{C}$ with $\text{Card}(\mathcal{H}^0) \geq$

$2^{Mmr/128}$ such that for any $A_1, A_2 \in \mathcal{H}^0$, the Hamming distance between A_1 and A_2 is at least $mr/16$. Now we define the building blocks of our hypotheses set

$$\mathcal{H} := \mathcal{H}^0 \cup \left\{ O_{\frac{m}{2} \times \frac{r}{2}} \right\},$$

where $O_{\frac{m}{2} \times \frac{r}{2}}$ is the $\frac{m}{2} \times \frac{r}{2}$ zero matrix. Obviously, \mathcal{H} has size $\text{Card}(\mathcal{H}) \geq 2^{Mmr/64} + 1$, and for any $A_1(t), A_2(t) \in \mathcal{H}$, the minimum Hamming distance is still greater than $mr/16$.

We consider the set of matrix valued functions

$$\mathcal{B}(\mathcal{H}) = \left\{ \begin{bmatrix} \tilde{A} & \tilde{A} & \dots & \tilde{A} & O \end{bmatrix} : \tilde{A} \in \mathcal{H} \right\},$$

where O denotes the $m/2 \times (m/2 - r\lfloor m/r \rfloor/2)$ zero matrix. Finally, our hypotheses set of matrix valued functions \mathcal{H}_m is defined as

$$\mathcal{H}_m = \left\{ \begin{bmatrix} \tilde{O} & \hat{A} \\ \hat{A}^* & \tilde{O} \end{bmatrix} \in \mathbb{C}^{m \times m} : \hat{A} \in \mathcal{B}(\mathcal{H}) \right\}.$$

By the definition of \mathcal{H}_m and similar to the arguments in proof of Theorem 3.1, it is easy to check that $\mathcal{H}_m \subset \mathcal{A}(r, a)$, and also

$$\text{Card}(\mathcal{H}_m) \geq 2^{Mmr/64} + 1. \quad (3.14)$$

Now we consider any two different hypotheses $A_j(t), A_k(t) \in \mathcal{H}_m$.

$$\int_0^1 \|A_j(t) - A_k(t)\|_2^2 dt \geq \gamma^2 \frac{mr}{16} 2^{\lfloor \frac{m}{r} \rfloor} \int_0^1 (f_\omega(t) - f_{\omega'}(t))^2 dt, \quad (3.15)$$

where $\omega \neq \omega'$. Based on (3.10), we have

$$\frac{1}{m^2} \int_0^1 \|A_j(t) - A_k(t)\|_2^2 dt \geq \frac{\gamma^2 L^2 h_n^{2\beta} \|f\|_2^2}{256} \geq c_* \left(\frac{mr}{n} \right)^{\frac{2\beta}{2\beta+1}}. \quad (3.16)$$

where c_* is a constant depending on $\|f\|_2$, L , c_1 and γ .

On the other hand, we repeat the same analysis on the Kullback-Leibler divergence $K(\mathbb{P}_0, \mathbb{P}_A)$ as in the proof of Theorem 3.1. One can get

$$K(\mathbb{P}_0, \mathbb{P}_A) \leq \frac{n}{8a^2} \mathbb{E} \langle A(\tau), X \rangle^2 \leq \frac{n}{8a^2} \gamma^2 \sum_{j=1}^M \int_0^1 \phi_j^2(\tau) d\tau \leq \frac{\gamma^2 c_1^{2\beta+1} L^2 M m r \|f\|_2^2}{8a^2}, \quad (3.17)$$

where $A(\tau) \in \mathcal{H}_m$. Combine (3.14) and (3.17) we know that

$$\frac{1}{\text{Card}(\mathcal{H}_m) - 1} \sum_{A(t) \in \mathcal{H}_m} K(\mathbb{P}_0, \mathbb{P}_A) \leq \alpha \log(\text{Card}(\mathcal{H}_m) - 1) \quad (3.18)$$

is satisfied for any $\alpha > 0$ if γ is chosen as a sufficiently small constant. In view of (3.16) and (3.18), the lower bound follows from Theorem 2.5 in [24]. \square

Now we consider the minimax lower bound on integrated risk measured by L_∞ -norm for general matrix valued functions without any rank information. Denote

$$\mathcal{A}(a) = \{A(t) \in \mathbb{H}_m, \forall t \in [0, 1] : |a_{ij}(t)| \leq a, a_{ij} \in \Sigma(\beta, L)\}.$$

We denote by $\mathcal{P}(a)$ the class of distributions of random triplet (τ, X, Y) that satisfies model (1.1) with any $A \in \mathcal{A}(a)$.

Theorem 3.4. *Under model (1.1), let (τ_j, X_j, Y_j) , $j = 1, \dots, n$ be i.i.d. copies of the random triplet (τ, X, Y) with X uniformly distributed in \mathcal{X} , τ uniformly distributed in $[0, 1]$, X and τ are independent, and $|Y| \leq a$, a.s. for some constant $a > 0$; let $A(t)$ be any matrix valued function in $\mathcal{A}(a)$. Then there exist an absolute constant $\eta \in (0, 1)$ such that*

$$\inf_{\hat{A}} \sup_{P_{\tau, X, Y}^A \in \mathcal{P}(a)} \mathbb{P}_{P_{\tau, X, Y}^A} \left\{ \sup_{t \in (0, 1)} \frac{1}{m^2} \|\hat{A}(t) - A(t)\|^2 > \bar{C}(\beta, L, a) \left(\frac{m \vee \log n}{n} \right)^{\frac{2\beta}{2\beta+1}} \right\} \geq \eta. \quad (3.19)$$

where $\bar{C}(\beta, L, a)$ is a constant depending on β , L and a .

Recall that in the real valued case, the minimax lower bound measured by sup norm of Hölder class is $O((\frac{\log n}{n})^{2\beta/(2\beta+1)})$. In our result (3.19), if the dimension m degenerates to 1, we get the same result as in real valued case and it is optimal. While the dimension m is large enough such that $m \gg \log n$, the lower bound (3.19) shows that the estimator (2.4) achieves a near minimax optimal rate up to a logarithmic factor.

Proof. Without loss of generality, assume that m is an even number. For some constant $\gamma > 0$, denote $\mathcal{V} = \left\{v \in \mathbb{C}^{\frac{m}{2}} : a_i \in \{0, \gamma\}, \forall 1 \leq i \leq m/2\right\}$. Due to the Varshamov-Gilbert bound (see Lemma 2.9 in [24]), there exists a subset $\mathcal{V}^0 \subset \mathcal{V}$ with cardinality $\text{Card}(\mathcal{V}^0) \geq 2^{m/16} + 1$ containing the zero vector $\mathbf{0} \in \mathbb{C}^{\frac{m}{2}}$, and such that for any two distinct elements v_1 and v_2 of \mathcal{V}^0 ,

$$\|v_1 - v_2\|_2^2 \geq \frac{m}{16}\gamma^2. \quad (3.20)$$

Consider the set of matrices

$$\mathcal{B}(\mathcal{V}) = \left\{ \begin{bmatrix} v & v & \dots & v \end{bmatrix} \in \mathbb{C}^{\frac{m}{2} \times \frac{m}{2}} : v \in \mathcal{V}^0 \right\}.$$

Clearly, $\mathcal{B}(\mathcal{V})$ is a collection of rank one matrices. Then we construct another matrix set \mathcal{V}_m ,

$$\mathcal{V}_m = \left\{ \begin{bmatrix} \tilde{O} & V \\ V^* & \tilde{O} \end{bmatrix} \in \mathbb{C}^{m \times m} : V \in \mathcal{B}(\mathcal{V}) \right\}$$

where \tilde{O} is the $m/2 \times m/2$ zero matrix. Apparently, $\mathcal{V}_m \subset \mathbb{H}_m$.

On the other hand, we define the grid on $[0, 1]$

$$M = \left\lceil c_2 \left(\frac{n}{m + \log n} \right)^{\frac{1}{2\beta+1}} \right\rceil, \quad h_n = \frac{1}{2M}, \quad t_j = \frac{2j-1}{2M},$$

and

$$\phi_j(t) = L h_n^\beta f\left(\frac{t - t_j}{h_n}\right), \quad j = 1, \dots, M, \quad t \in [0, 1]$$

where f is defined the same as in (3.5), and c_2 is some constant. Denote $\Phi := \left\{ \phi_j : j = 1, \dots, M \right\}$. We consider the following set of hypotheses: $\mathcal{A}_B^\beta := \{\hat{A}(t) = V\phi_j(t) : V \in \mathcal{V}_m, \phi_j \in \Phi\}$. One can immediately get that the size of \mathcal{A}_B^β satisfies

$$\text{Card}(\mathcal{A}_B^\beta) \geq (2^{m/16} + 1)M. \quad (3.21)$$

By construction, the following claims are obvious: any element $\hat{A}(t)$ of \mathcal{A}_B^β has rank at most 2; the entries of $\hat{A}(t) \in \mathcal{A}_B^\beta$ are uniformly bounded for some sufficiently small γ , and $\hat{A}_{ij}(t) \in \Sigma(\beta, L)$. Thus $\mathcal{A}_B^\beta \subset \mathcal{A}(a)$.

Now we consider the distance between two distinct elements $A(t)$ and $A'(t)$ of \mathcal{A}_B^β . An immediate observation is that

$$\sup_{t \in [0,1]} \|A(t) - A'(t)\|^2 \geq \frac{1}{4} \sup_{t \in [0,1]} \|A(t) - A'(t)\|_2^2,$$

due to the fact that $\forall t \in (0, 1)$, $\text{rank}(A(t) - A'(t)) \leq 4$. Then we turn to get lower bound on $\sup_{t \in (0,1)} \|A(t) - A'(t)\|_2^2$. Recall that by construction of \mathcal{A}_B^β , we have for any $A \neq A'$, $A(t) = A_1\phi_j(t)$, $A'(t) = A_2\phi_k(t)$, where $A_1, A_2 \in \mathcal{V}_m$. There are three cases need to be considered: 1). $A_1 \neq A_2$ and $j = k$; 2). $A_1 = A_2 \neq 0$ and $j \neq k$; 3). $A_1 \neq A_2$ and $j \neq k$.

For case 1,

$$\sup_{t \in [0,1]} \|A(t) - A'(t)\|_2^2 = \|A_1 - A_2\|_2^2 \|\phi_j\|_\infty^2 \geq \frac{m^2}{16} \gamma^2 L^2 h_n^{2\beta} \|f\|_\infty^2 \geq c_1^* m^2 \left(\frac{m + \log n}{n} \right)^{\frac{2\beta}{2\beta+1}},$$

where c_1^* is a constant depending on $\|f\|_\infty^2$, β , L and γ .

For case 2,

$$\sup_{t \in [0,1]} \|A(t) - A'(t)\|_2^2 = \|A_1\|_2^2 \|\phi_j - \phi_k\|_\infty^2 \geq \frac{m^2}{16} \gamma^2 L^2 h_n^{2\beta} \|f\|_\infty^2 \geq c_2^* m^2 \left(\frac{m + \log n}{n} \right)^{\frac{2\beta}{2\beta+1}},$$

where c_2^* is a constant depending on $\|f\|_\infty^2$, β , L and γ .

For case 3,

$$\sup_{t \in [0,1]} \|A(t) - A'(t)\|_2^2 \geq (\|A_1\|_2^2 \|\phi_j\|_\infty^2 \vee \|A_2\|_2^2 \|\phi_k\|_\infty^2) \geq \frac{m^2}{16} \gamma^2 L^2 h_n^{2\beta} \|f\|_\infty^2 \geq c_3^* m^2 \left(\frac{m + \log n}{n} \right)^{\frac{2\beta}{2\beta+1}},$$

where c_3^* is a constant depending on $\|f\|_\infty^2$, β , L and γ .

Therefore, by the analysis above we conclude that for any two distinct elements $A(t)$ and $A'(t)$ of \mathcal{A}_B^β ,

$$\sup_{t \in [0,1]} \|A(t) - A'(t)\|^2 \geq \frac{1}{4} \sup_{t \in [0,1]} \|A(t) - A'(t)\|_2^2 \geq c_* m^2 \left(\frac{m + \log n}{n} \right)^{\frac{2\beta}{2\beta+1}}, \quad (3.22)$$

where c_* is a constant depending on $\|f\|_\infty^2$, L , γ and β .

Meanwhile, we repeat the same analysis on the Kullback-Leibler divergence $K(\mathbb{P}_0, \mathbb{P}_A)$ as in the proof of Theorem 3.1. One can get that for any $A \in \mathcal{A}_B^\beta$, the Kullback-Leibler divergence $K(\mathbb{P}_0, \mathbb{P}_A)$ between \mathbb{P}_0 and \mathbb{P}_A satisfies

$$K(\mathbb{P}_0, \mathbb{P}_A) \leq \frac{n}{8a^2} \mathbb{E} |\langle A(\tau), X \rangle|^2 \leq \frac{n}{8a^2} \gamma^2 \int_0^1 \phi_j^2(\tau) d\tau \leq \frac{\gamma^2 c_2^{2\beta+1} L^2 (m + \log n) \|f\|_2^2}{8a^2}. \quad (3.23)$$

Combine (3.21) and (3.23) we know that

$$\frac{1}{\text{Card}(\mathcal{A}_B^\beta) - 1} \sum_{A \in \mathcal{A}_B^\beta} K(\mathbb{P}_0, \mathbb{P}_A) \leq \alpha \log(\text{Card}(\mathcal{A}_B^\beta) - 1) \quad (3.24)$$

is satisfied for any $\alpha > 0$ if γ is chosen as a sufficiently small constant. In view of (3.22) and (3.24), the lower bound follows from Theorem 2.5 in [24]. \square

4 Model selection

Despite the fact that estimators (1.3) and (2.1) achieve near optimal minimax rates in theory with properly chosen bandwidth h and order of degree ℓ , such parameters depend on

quantities like β and L which are unknown to us in advance. In this section, we propose an adaptive estimation procedure to choose h and ℓ adaptively. Two popular methods to address such problems are proposed in the past few decades. One is Lepskii's method, and the other is aggregation method. In the 1990s, many data-driven procedures for selecting the "best" estimator emerged. Among them, a series of papers stood out and shaped a method what is now called "Lepskii's method". This method has been described in its general form and in great detail in [19]. Later, [27] proposed a bandwidth selection procedure based on pointwise adaptation of a kernel estimator that achieves optimal minimax rate of point estimation over Hölder class, and [28] proposed a new bandwidth selector that achieves optimal rates of convergence over Besov classes with spatially inhomogeneous smoothness. The basic idea of Lepskii's method is to choose a bandwidth from a geometric grid to get an estimator not "very different" from those indexed by smaller bandwidths on the grid. Although Lepskii's method is shown to give optimal rates in pointwise estimation over Hölder class in [27], it has a major defect when applied to our problem: the procedure already requires a huge amount of computational cost when real valued functions are replaced by matrix valued functions. Indeed, with Lepskii's method, in order to get a good bandwidth, one needs to compare all smaller bandwidth with the target one, which leads to dramatically growing computational cost. Still, we have an extra parameter ℓ that needs to fit with h . As a result, we turn to aggregation method to choose a bandwidth from the geometric grid introduced by Lepskii's method, which is more computationally efficient for our problem. The idea of aggregation method can be briefly summarized as follows: one splits the data set into two parts; the first is used to build all candidate estimators and the second is used to aggregate the estimates to build a new one (aggregation) or select one (model selection) which is as good as the best candidate among all constructed. The model selection procedure we use was initially introduced by [20] in classical nonparametric estimation with bounded response. [21] generalized this method to the case where the noise can be unbounded but with a finite p -th moment for some $p > 2$. One can find a more

detailed review on such penalization methods in [29].

Firstly, we introduce the geometric grid created by [27] where to conduct our model selection procedure. Assume that the bandwidth being considered falls into the range $[h_{\min}, h_{\max}]$. Recall that the "ideal" bandwidth \hat{h}_n which is given as

$$\hat{h}_n = C_1 \left(\frac{\ell^3 (\ell! \Phi R(T) a)^2 m r \log n}{L^2 n} \right)^{\frac{1}{2\beta+1}}, \quad (4.1)$$

h_{\max}, h_{\min} can be chosen as

$$h_{\max} = C_1 \left(\frac{\ell^{*3} (\ell^*! \Phi R(T) a)^2 m r \log n}{L_*^2 n} \right)^{\frac{1}{2\beta_*+1}}, \quad h_{\min} = C_1 \left(\frac{\ell_*^3 (\ell_*! \Phi R(T) a)^2 m r \log n}{L^{*2} n} \right)^{\frac{1}{2\beta_*+1}},$$

where $[\beta_*, \beta^*]$ and $[L_*, L^*]$ are the possible ranges of β, L respectively. Obviously, β is the most important parameter among all. Note that when those ranges are not so clear, a natural upper bound of h_{\max} is 1, and a typical choice of h_{\min} can be set to $n^{-1/2}$. Denote

$$d(h) = \sqrt{1 \vee 2 \log \left(\frac{h_{\max}}{h} \right)}, \quad d_n = \sqrt{2 \log \left(\frac{h_{\max}}{h_{\min}} \right)}, \quad \alpha(h) = \frac{1}{\sqrt{d(h)}}.$$

Apparently, $d_n = O(\sqrt{\log n})$. Define grid \mathcal{H} inductively by

$$\mathcal{H} = \left\{ h_k \in [h_{\min}, h_{\max}] : h_0 = h_{\max}, h_{k+1} = \frac{h_k}{1 + \alpha(h_k)}, k = 0, 1, 2, \dots \right\}. \quad (4.2)$$

Note that the grid \mathcal{H} is a decreasing sequence and the sequence becomes denser as k grows.

We consider possible choices of ℓ_k for each h_k . A trivial candidate set is $\ell_k \in \mathcal{L} := \{ \lfloor \beta_* \rfloor, \lfloor \beta_* \rfloor + 1, \dots, \lfloor \beta^* \rfloor \} \subset \mathbb{N}^*$. If the size of this set is large, one can shrink it through the correspondence (4.1) between h and β , $\ell_k \leq \left\lfloor \frac{\log n^{-1} + \log m r \log 2m - 1}{\log h_k} \right\rfloor$. If $n \geq m^d$ for some $d > 1$, $\left\lfloor \frac{(1 - \frac{1}{d}) \log n^{-1} - 1}{\log h_k} \right\rfloor \leq \ell_k \leq \left\lfloor \frac{\log n^{-1} - 1}{\log h_k} \right\rfloor$, which indicates the more the data, the narrower the range. We denote the candidate set for ℓ as \mathcal{L} . Then the set

$$\tilde{\mathcal{H}} = \mathcal{H} \times \mathcal{L} := \{ (h, \ell) : h \in \mathcal{H}, \ell \in \mathcal{L} \}$$

indexed a countable set of candidate estimators. Once (h_k, ℓ_k) is fixed, one can take $\varepsilon_k = D(\ell_k + 1)R(T)\Phi a\sqrt{\frac{\log 2m}{nmh_k}}$.

Now we introduce our model selection procedure based on $\tilde{\mathcal{H}}$. We split the data (τ_j, X_j, Y_j) , $j = 1, \dots, 2n$, into two parts with equal size. The first part of the observations $\{(\tau_j, X_j, Y_j) : j \in \tilde{h}_n\}$ contains n data points, which are randomly drawn without replacement from the original data set. We construct a sequence of estimators \hat{A}^k , $k = 1, 2, \dots$ based on the training data set \tilde{h}_n through (2.1) for each pair in $\tilde{\mathcal{H}}$. Our main goal is to select an estimator \hat{A} among $\{\hat{A}^k\}$, which is as good as the one that has the smallest mean square error. We introduce an quantity π_k associated with each estimator \hat{A}^k which serves as a penalty term. We use the remaining part of the data set $\{(\tau_j, X_j, Y_j) : j \in \ell_n\}$ to perform the selection procedure:

$$k^* = \arg \min_k \frac{1}{n} \sum_{j \in \ell_n} (Y_j - \langle \hat{A}^k(\tau_j), X_j \rangle)^2 + \frac{\pi_k}{n}. \quad (4.3)$$

Denote $\hat{A}^* = \hat{A}^{k^*}$ as the adaptive estimator. In practice, we suggest one to rank all estimators \hat{A}^k according to the following rule: 1. pairs with bigger h always have smaller index; 2. if two pairs have the same h , the one with smaller ℓ has smaller index. Our selection procedure can be summarized in Algorithm 1:

Algorithm 1: Model Selection Procedure

1. Construct the geometric grid \mathcal{H} defined as in (4.2) and compute the candidate set $\tilde{\mathcal{H}}$;
 2. Equally split the data set (τ_j, X_j, Y_j) , $j = 1, \dots, N$ into two parts (\tilde{h}_n and ℓ_n) by randomly drawing without replacement;
 3. For each pair in $\tilde{\mathcal{H}}$, construct an estimator \hat{A}^k defined as in (2.1) using data set \tilde{h}_n ;
 4. Using the second data set ℓ_n to perform the selection rule defined as in (4.3).
-

The selection procedure described in Algorithm 1 have several advantages: firstly, it chooses a global bandwidth instead of a local one; secondly, since our selection procedure is only based on computations of entries of \hat{A}^k , no matrix computation is involved in the last step, which saves a lot of computational cost and can be easily applied to high dimensional problems; finally, step 3 and 4 can be easily parallelized. The following theorem shows that the integrated risk of \hat{A}^* measured by L_2 -norm can be bounded by the smallest one among all candidates plus an extra term of order n^{-1} which is negligible.

Theorem 4.1. *Under model (1.1), let (τ_j, X_j, Y_j) , $j = 1, \dots, 2n$ be i.i.d. copies of the random triplet (τ, X, Y) with X uniformly distributed in \mathcal{X} , τ uniformly distributed in $[0, 1]$, X and τ are independent, and $|Y| \leq a$, a.s. for some constant $a > 0$; let A be a matrix valued function satisfying A1, A2, A3, and A4; let $\{\hat{A}^k\}$ be a sequence of estimators constructed from $\tilde{\mathcal{H}}$; let \hat{A}^* be the adaptive estimator selected through Algorithm 1. Then with probability at least $1 - \frac{1}{n^{mr}}$*

$$\frac{1}{m^2} \int_0^1 \|\hat{A}^*(t) - A(t)\|_2^2 dt \leq 3 \min_k \left\{ \frac{1}{m^2} \int_0^1 \|\hat{A}^k(t) - A(t)\|_2^2 dt + \frac{\pi_k}{n} \right\} + \frac{C(a)}{n}, \quad (4.4)$$

where $C(a)$ is a constant depending on a .

Recall that $\text{Card}(\mathcal{H}) = O(\log n)$, we can take $\pi_k = kmr$. Then $\pi_k \leq c_1 mr \log n$ uniformly for all k with some numerical constant c_1 . According to Lepskii's method that at least one candidate in \mathcal{H} gives the optimal bandwidth associated with the unknown smoothness parameter β , together with the result of Theorem 2.1, the following corollary follows from Theorem 4.1.

Proof. For any \hat{A}^k , denote the difference in empirical loss between \hat{A}^k and A by

$$r_n(\hat{A}^k, A) := \frac{1}{n} \sum_{j=1}^n (Y_j - \langle \hat{A}^k(\tau_j), X_j \rangle)^2 - \frac{1}{n} \sum_{j=1}^n (Y_j - \langle A(\tau_j), X_j \rangle)^2 = -\frac{1}{n} \sum_{j=1}^n U_j,$$

where $U_j = (Y_j - \langle A(\tau_j), X_j \rangle)^2 - (Y_j - \langle \hat{A}^k(\tau_j), X_j \rangle)^2$. It is easy to check that

$$U_j = 2(Y_j - \langle A(\tau_j), X_j \rangle)\langle \hat{A}^k(\tau_j) - A(\tau_j), X_j \rangle - \langle \hat{A}^k(\tau_j) - A(\tau_j), X_j \rangle^2. \quad (4.5)$$

We denote $r(\hat{A}^k, A) := \mathbb{E}\langle \hat{A}^k(\tau) - A(\tau), X \rangle^2$. The following concentration inequality developed by [30] to prove Bernstein's inequality is key to our proof.

Lemma 2. Let $U_j, j = 1, \dots, n$ be independent bounded random variables satisfying $|U_j - \mathbb{E}U_j| \leq M$ with $h = M/3$. Set $\bar{U} = n^{-1} \sum_{j=1}^n U_j$. Then for all $t > 0$

$$\mathbb{P}\left\{\bar{U} - \mathbb{E}\bar{U} \geq \frac{t}{n\varepsilon} + \frac{n\varepsilon \text{var}(\bar{U})}{2(1-c)}\right\} \leq e^{-t},$$

with $0 < \varepsilon h \leq c < 1$.

Firstly, we bound the variance of U_j . Under the assumption that $|Y|$ and $|\langle A(\tau), X \rangle|$ are bounded by a constant a , one can easily check that $h = 8a^2/3$. Given $\mathbb{E}(Y_j|\tau_j, X_j) = \langle A(\tau_j), X_j \rangle$, we know that the covariance between the two terms on the right hand side of (4.5) is zero. Conditionally on (τ, X) , the second order moment of the first term satisfies

$$4\mathbb{E}\sigma_{Y|\tau, X}^2 \langle \hat{A}^k(\tau_j) - A(\tau_j), X_j \rangle^2 \leq 4a^2 r(\hat{A}^k, A).$$

To see why, one can consider the random variable \tilde{Y} with the distribution $\mathbb{P}\{\tilde{Y} = a\} = \mathbb{P}\{\tilde{Y} = -a\} = 1/2$. The variance of Y is always bounded by the variance of \tilde{Y} which is a^2 under the assumption that $|Y_j|$ and $|\langle \hat{A}^k(\tau_j), X_j \rangle|$ are bounded by a constant $a > 0$. Similarly, we can get that the variance of the second term conditioned on (τ, X) is also bounded by $4a^2 \mathbb{E}\langle \hat{A}^k(\tau_j) - A(\tau_j), X_j \rangle^2$. As a result, $n\text{var}(\bar{U}) \leq 8a^2 r(\hat{A}^k, A)$. By the result of Lemma 2, we have for any \hat{A}^k with probability at least $1 - e^{-t}$

$$r(\hat{A}^k, A) - r_n(\hat{A}^k, A) < \frac{t}{n\varepsilon} + \frac{4a^2 \varepsilon r(\hat{A}^k, A)}{1-c}.$$

Set $t = \varepsilon\pi_k + \log 1/\delta$, we get with probability at least $1 - \delta/e^{\varepsilon\pi_k}$

$$(1 - \alpha)r(\hat{A}^k, A) < r_n(\hat{A}^k, A) + \frac{\pi_k}{n} + \frac{4a^2}{(1 - c)\alpha} \left(\frac{\log 1/\delta}{n} \right).$$

where $\alpha = 4a^2\varepsilon/(1 - c) < 1$. Denote

$$\tilde{k}^* = \arg \min_k \left\{ r(\hat{A}^k, A) + \frac{\pi_k}{n} \right\}.$$

By the definition of \hat{A}^* , we have with probability at least $1 - \delta/e^{\varepsilon\hat{\pi}^*}$

$$(1 - \alpha)r(\hat{A}^*, A) < r_n(\hat{A}^{\tilde{k}^*}, A) + \frac{\pi_{\tilde{k}^*}}{n} + \frac{4a^2}{(1 - c)\alpha} \left(\frac{\log 1/\delta}{n} \right). \quad (4.6)$$

where $\hat{\pi}^*$ is the penalty terms associated with \hat{A}^* .

Now we apply the result of Lemma 2 one more time and set $t = \log 1/\delta$, we get with probability at least $1 - \delta$

$$r_n(\hat{A}^{\tilde{k}^*}, A) \leq (1 + \alpha)r(\hat{A}^{\tilde{k}^*}, A) + \frac{4a^2}{(1 - c)\alpha} \frac{\log 1/\delta}{n}. \quad (4.7)$$

Apply the union bound of (4.6) and (4.7), we get with probability at least $1 - \delta(1 + e^{-\varepsilon\hat{\pi}^*})$

$$r(\hat{A}^*, A) \leq \frac{(1 + \alpha)}{(1 - \alpha)} \left(r(\hat{A}^{\tilde{k}^*}, A) + \frac{\pi_{\tilde{k}^*}}{n} \right) + \frac{4a^2}{(1 - c)\alpha(1 - \alpha)} \frac{\log 1/\delta}{n}.$$

By taking $\varepsilon = 3/32a^2$ and $c = \varepsilon h$,

$$r(\hat{A}^*, A) \leq 3 \left(r(\hat{A}^{\tilde{k}^*}, A) + \frac{\pi_{\tilde{k}^*}}{n} \right) + \frac{64a^2}{3} \frac{\log 1/\delta}{n}.$$

By taking $\delta = 1/n^{mr}$ and adjusting the constant, we have with probability at least $1 - 1/n^{mr}$

$$\frac{1}{m^2} \int_0^1 \|\hat{A}^*(t) - A(t)\|_2^2 dt \leq 3 \min_k \left\{ \frac{1}{m^2} \int_0^1 \|\hat{A}^k(t) - A(t)\|_2^2 dt + \frac{\pi_k}{n} \right\} + C(a) \frac{mr \log n}{n}$$

where $C(a)$ is a constant depending on a . □

Corollary 4.1. Assume that the conditions of Theorem 4.1 hold with $\pi_k = kmr$, and $n > mr \log n$. Then with probability at least $1 - \frac{1}{n^{mr-1}}$

$$\frac{1}{m^2} \int_0^1 \|\widehat{A}^*(t) - A(t)\|_2^2 dt \leq C(a, \Phi, \ell, L) \left(\frac{mr \log n}{n} \right)^{\frac{2\beta}{2\beta+1}} \quad (4.8)$$

where $C(a, \Phi, \ell, L)$ is a constant depending on a, Φ, ℓ , and L .

CHAPTER 3

SIMULATION RESULTS OF NONPARAMETRIC ESTIMATION OF LOW RANK MATRIX VALUED FUNCTION

1 An ADMM Algorithm

The alternating direction method of multipliers (ADMM) is a powerful algorithm to solve convex optimization problems, see a comprehensive introduction in [31]. The application of ADMM to matrix recovery problems can be found in [32, 33] and the references therein.

Recall that the major optimization problem to solve is

$$\hat{S}^h = \arg \frac{1}{nh} \sum_{j=1}^n K\left(\frac{\tau_j - t_0}{h}\right) \left(Y_j - \left\langle \sum_{i=0}^{\ell} S_i p_i\left(\frac{\tau_j - t_0}{h}\right), X_j \right\rangle\right)^2 + \varepsilon \|S\|_1. \quad (1.1)$$

where $\mathbb{D} \subset \mathbb{H}_{(\ell+1)m}$ is a closed subset of block diagonal matrices with $S_j \in \mathbb{H}_m$ on its diagonal, and $\{p_i\}$ is a sequence of orthogonal polynomials with nonnegative weight function K . Such a problem belongs to the standard form of optimization problems considered in ADMM applications, see [31]. To be more specific, the optimization problem involves the sum of two functions of the solution, i.e. a loss function and a penalization term. Instead of solving the original optimization problem, our ADMM algorithm presented below introduces a new variable such that the original problem to solve is transformed into the following one:

$$\hat{S}^h = \arg \min_{S=\bar{S}} \min_{S, \bar{S} \in \mathbb{D}} \frac{1}{nh} \sum_{j=1}^n K\left(\frac{\tau_j - t_0}{h}\right) \left(Y_j - \left\langle \sum_{i=0}^{\ell} S_i p_i\left(\frac{\tau_j - t_0}{h}\right), X_j \right\rangle\right)^2 + \varepsilon \|\bar{S}\|_1. \quad (1.2)$$

Then the corresponding augmented Lagrangian multipliers of function of (1.2) is defined

as

$$\begin{aligned} \hat{S}^h = \arg \min_{S=\bar{S}} \min_{S, \bar{S} \in \mathbb{D}} \frac{1}{nh} \sum_{j=1}^n K\left(\frac{\tau_j - t_0}{h}\right) \left(Y_j - \left\langle \sum_{i=0}^{\ell} S_i p_i\left(\frac{\tau_j - t_0}{h}\right), X_j \right\rangle\right)^2 \\ + \varepsilon \|\bar{S}\|_1 + \frac{\rho}{2} \|S - \bar{S}\|_2^2 + \langle Z, S - \bar{S} \rangle. \end{aligned} \quad (1.3)$$

where $\rho > 0$ is a constant and $Z \in \mathbb{D}$. Our ADMM algorithm is presented in Algorithm 2. It updates S and \bar{S} alternatively and the multiplier Z is updated by the difference between the iterates of S and \bar{S} . Note that in order to update $\bar{S}^{(k)}$, it is equivalent to solve the following optimization problem:

$$\bar{S}^{(k+1)} = \arg \min_{\bar{S} \in \mathbb{D}} \varepsilon \|\bar{S}\|_1 + \frac{\rho}{2} \|S^{(k+1)} - \bar{S} + \frac{Z^{(k)}}{\rho}\|_2^2. \quad (1.4)$$

It was proved by [9] that the solution to this problem has a simple form which can be obtained by soft thresholding of singular values of $S^{(k+1)} + \frac{Z^{(k)}}{\rho}$ as

$$\bar{S}^{(k+1)} = \sum_j (\sigma_j(\tilde{S}) - \frac{\varepsilon}{\rho})_+ u_j(\tilde{S}) v_j(\tilde{S})^T, \quad (1.5)$$

where $\tilde{S} = S^{(k+1)} + \frac{Z^{(k)}}{\rho}$, $x_+ = \max\{x, 0\}$, and $\sigma_j(\tilde{S})$, u_j , v_j are the singular values, left and right singular vectors of \tilde{S} respectively.

2 Numerical results

In this section, we present the numerical simulation results of our estimators (1.3) and (2.1), and simulation results of our model selection procedure in Algorithm 1. The underlying matrix valued function we create is in Hölder class $\Sigma(\beta, L)$ with $\beta = 3/2$, $L = 24$ and rank constraint $\text{rank}_A(t) \leq 3$. The orthogonal polynomial we choose is Chebyshev polynomials of the second kind.

Algorithm 2: ADMM Algorithm

Set up the values of $max_Iteration$ and tolerance $\varepsilon_{tol} > 0$; Initialize $S^{(0)}, \bar{S}^{(0)} \in \mathbb{D}$ and $Z^{(0)} = \mathbf{0}$; **while** $k < max_Iteration$ **do**

$$S^{(k+1)} = \arg \min_{S \in \mathbb{D}} \frac{1}{nh} \sum_{j=1}^n K\left(\frac{\tau_j - t_0}{h}\right) \left(Y_j - \left\langle \sum_{i=0}^{\ell} S_i p_i\left(\frac{\tau_j - t_0}{h}\right), X_j \right\rangle\right)^2 + \frac{\rho}{2} \|S -$$

$$\bar{S}^{(k)}\|_2^2 + \langle Z^{(k)}, S - \bar{S}^{(k)} \rangle;$$

$$\bar{S}^{(k+1)} = \arg \min_{\bar{S} \in \mathbb{D}} \varepsilon \|\bar{S}\|_1 + \frac{\rho}{2} \|S^{(k+1)} - \bar{S}\|_2^2 + \langle Z^{(k)}, S^{(k+1)} - \bar{S} \rangle;$$

$$Z^{(k+1)} = Z^{(k)} + \rho(S^{(k+1)} - \bar{S}^{(k+1)});$$

if $\|\bar{S}^{(k+1)} - \bar{S}^{(k)}\|_2^2 \leq \varepsilon_{tol}$ **or** $\|Z^{(k+1)} - Z^{(k)}\|_2^2 \leq \rho^2 \varepsilon_{tol}$ **then**

 Reaching the tolerance;

end

 Return $\bar{S}^{(k+1)}$. $k = k + 1$;

end

Return $\bar{S}^{(k+1)}$.

In Fig. 3.1 and Fig. 3.2 we plot the pointwise error at $t_0 = 0.5$ and integrated risk against the iteration number of Algorithm 2. As we can see, our ADMM algorithm converges really fast when ρ is small. We should also mention that according to our experiments, smaller ρ value gives faster convergence speed typically. But it doesn't mean that smaller ρ is always better. There is an optimal value of ρ that gives the best accuracy. One needs to tune this parameter in order to get the best accuracy and fairly good convergence speed.

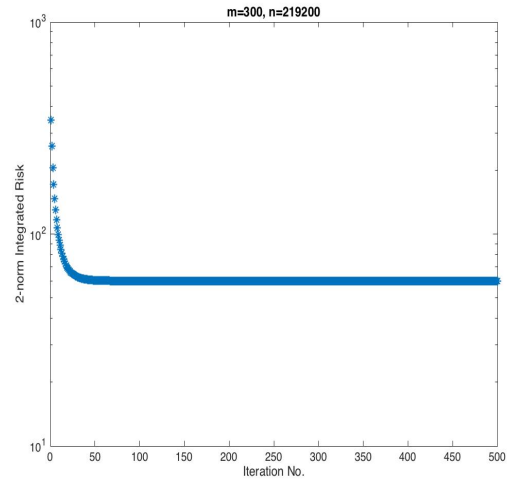
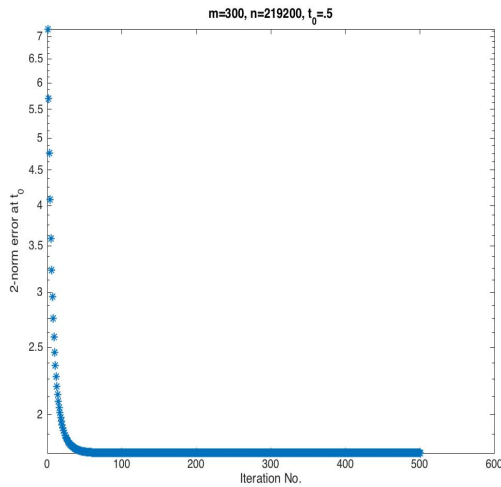


Figure 3.1: Pointwise risk convergence of ADMM Algorithm Figure 3.2: Integrated risk convergence of ADMM Algorithm

2.1 Pointwise estimation simulation

By plug in the optimal bandwidth in Theorem 1.1, we run our Algorithm 2 to solve the point estimator at $t_0 = 0.5$ with $m = 150$. Fig. 3.3 - Fig. 3.10 show different levels of recovery for the true underlying data matrix. As we can see, the recovery quality increases evidently as sample size n grows.

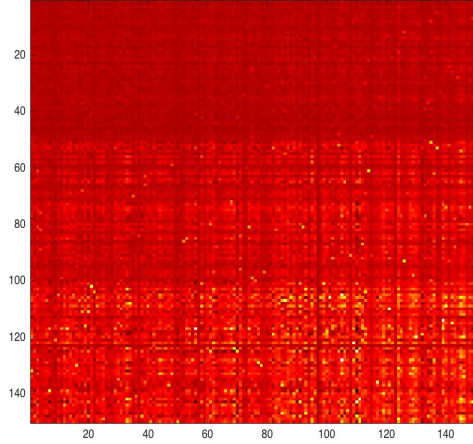


Figure 3.3: $n = 1600$

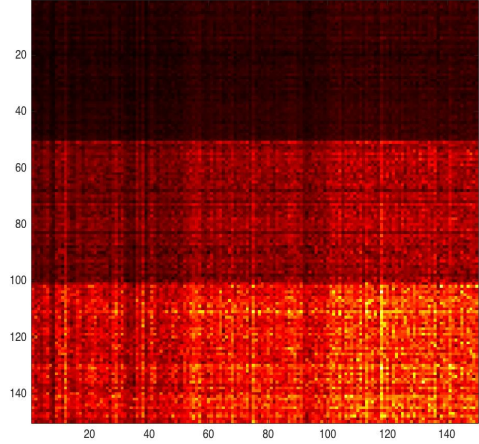


Figure 3.4: $n = 6400$

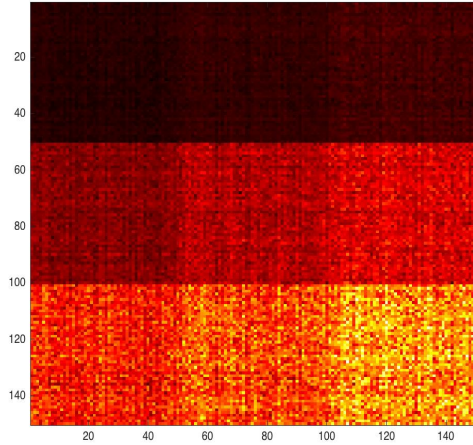


Figure 3.5: $n = 25600$

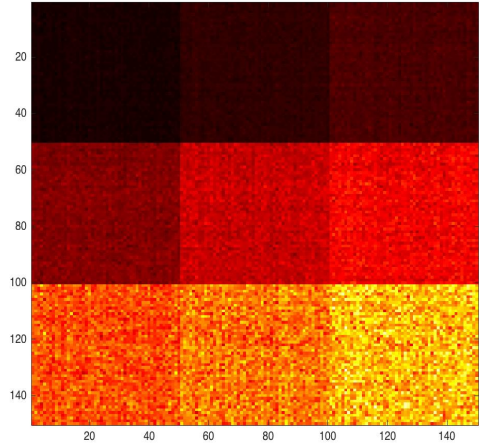


Figure 3.6: $n = 102400$

In table 3.1 and table 3.2 we display the comparison of pointwise risk measured by $\frac{1}{m^2} \|\hat{A}(t_0) - A(t_0)\|_2^2$ between our theoretical bounds proved in (1.4), (3.1) and our simula-

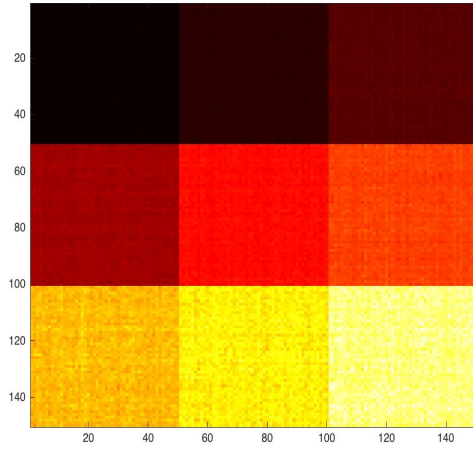


Figure 3.7: $n = 409600$

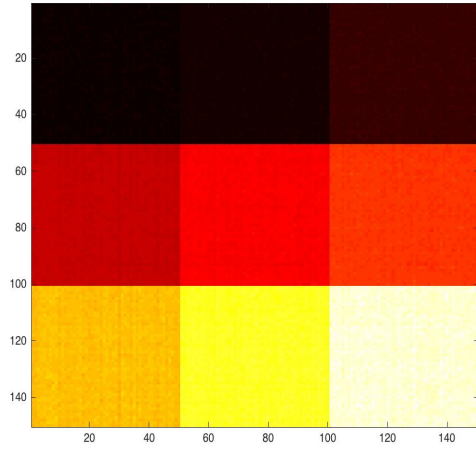


Figure 3.8: $n = 1638400$

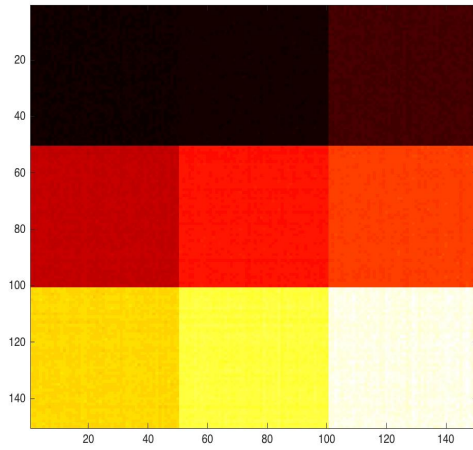


Figure 3.9: $n = 3276800$

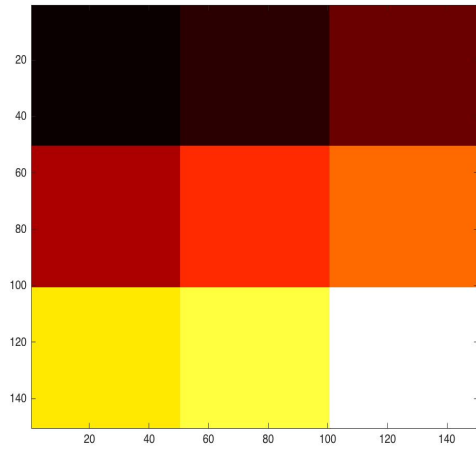


Figure 3.10: True data

tion results. The data is plotted in Fig. 3.11. As we can see, the simulation results match well with the minimax lower bound (3.1).

Sample size	800	1600	3200	6400	12800	25600
Theoretical upper bound	19.1780	11.4033	6.7805	4.0317	2.3973	1.4254
Minimax lower bound	5.1962	3.0897	1.8371	1.0924	0.6495	0.3862
Experimental error rate	7.2122	4.4569	1.9499	0.8600	0.5302	0.4329

Table 3.1: Pointwise error rate comparison with different sample size n

Sample size	51200	102400	204800	409600	819200	1638400
Theoretical upper bound	0.8476	0.5040	0.2997	0.1782	0.1059	0.0630
Minimax lower bound	0.2296	0.1365	0.0812	0.0483	0.0287	0.0171
Experimental error rate	0.2518	0.1156	0.0584	0.0466	0.0354	0.0194

Table 3.2: Pointwise error rate comparison with different sample size n

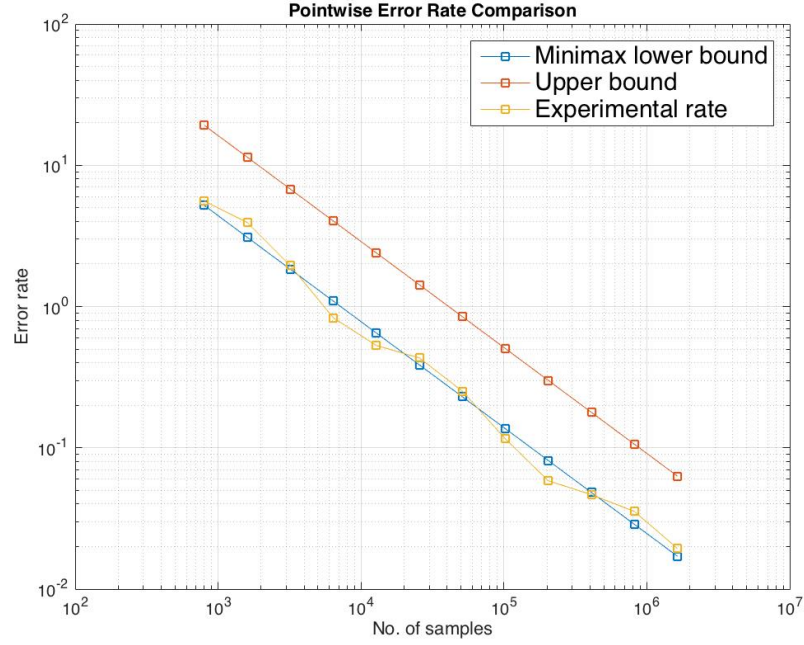


Figure 3.11: The point risk comparison at $t_0 = 0.5$

2.2 Integrated risk estimation simulation

In table 3.3 and table 3.4 we display the comparison of integrated risk measured by the L_2 -norm between the theoretical bounds proved in (2.2), (3.9) and our simulation results. Since $\beta = 3/2$ and $\ell = 1$, we use piecewise linear polynomials to approximate the underlying function. The data is plotted in Fig. 3.11. As we can see, the simulation results match well with the minimax lower bound (3.9).

Sample size	800	1600	3200	6400	12800	25600
Theoretical upper bound	1683.1	1000.8	595.1	353.8	210.4	125.1
Minimax lower bound	456.0163	271.1489	161.2261	95.8656	57.0020	33.8936
Experimental error rate	457.7443	293.3489	170.4948	106.8291	57.8282	37.2912

Table 3.3: Integrated error rate comparison with different sample size n

Sample size	51200	102400	204800	409600	819200	1638400
Theoretical upper bound	74.4	44.2	26.3	15.6	9.3	5.5
Minimax lower bound	20.1533	11.9832	7.1253	4.2367	2.5192	1.4979
Experimental error rate	19.1367	11.4798	8.0132	4.0110	2.9849	1.5030

Table 3.4: Integrated error rate comparison with different sample size n

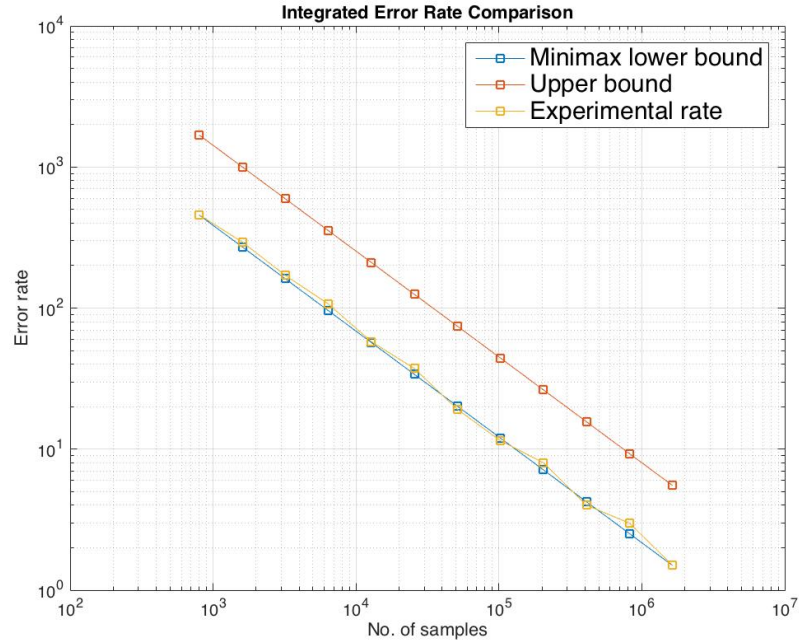


Figure 3.12: The integrated risk comparison with different sample size n

2.3 Simulation of model selection

Recall that in section 4, we emphasized that choosing a good bandwidth h is crucial to get better estimation. We developed Algorithm 1 to choose the optimal bandwidth. We implement Algorithm 1 in this section, and perform simulation with $m = 90$ and $n = 3200000$. We choose $h_{\max} = 1.0$ and $h_{\min} = 1/\sqrt{n}$ to construct the geometric grid \mathcal{H} as in (4.2). We

display the simulation results in table 3.5 and 3.6. To be more specific, we computed each global estimator as in (2.1) with each bandwidth on the \mathcal{H} . The corresponding integrated risks measured by L_2 -norm are displayed in second row and our model selection criterion computed as in (4.3) are displayed in the third row. The smaller value of the third row, the better. The data are plotted in Fig. 3.13. As we can see, our selector selects $\hat{h} = 0.0853$ with the smallest criterion value of 0.3490. The corresponding integrated risk is also the smallest among all candidates on the grid.

Bandwidth on grid \mathcal{H}	1.0000	0.5000	0.2602	0.1461
Integrated risk	68.1239	45.0275	1.0207	0.0657
Model selection criterion	5.8238	4.7442	1.0100	0.3862

Table 3.5: Model Selection

Bandwidth on grid \mathcal{H}	0.0853	0.0510	0.0311	0.0192	0.0121
Integrated risk	0.0333	0.04371	0.0538	0.0663	0.0807
Model selection criterion	0.3490	0.4821	0.6741	0.9771	1.3199

Table 3.6: Model Selection

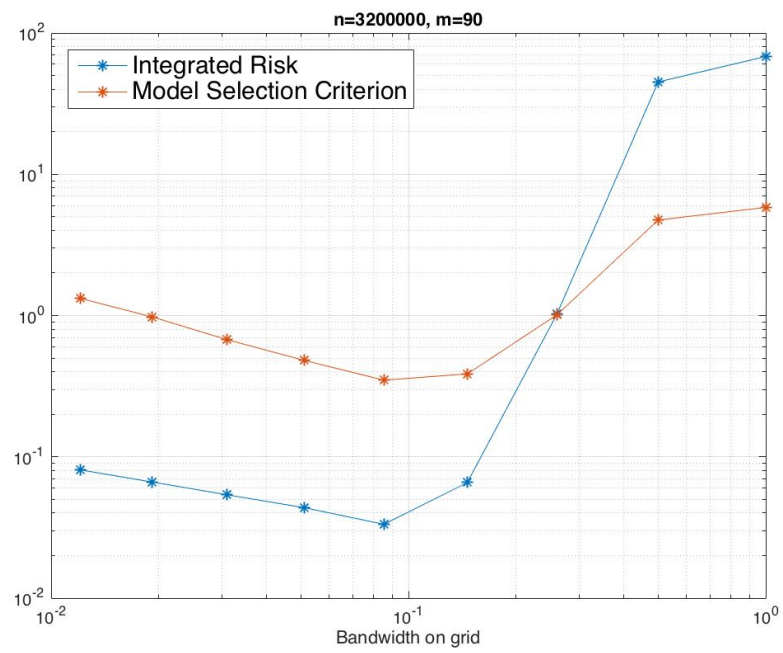


Figure 3.13: Model Selection on Grid \mathcal{H}

CHAPTER 4

THE ℓ_∞ PERTURBATION OF HOSVD AND LOW RANK TENSOR DENOISING

1 Introduction

A tensor is a mutliarray of more than 2 dimensions, which can be viewed as a higher order generalization of matrices. Data of tensor types has been widely available in many fields, such as image and video processing (see [34], [35], [36], [37], [38]); latent variable modeling (see [39], [40], [41]); genomic signal processing ([42], [43] and [44]) and references therein. It is demanding to handle these datasets in order to take the most advantages of the tensor structures. The task is challenging due to the highly non-convexity of tensor related optimization problems. For instance, computing the tensor operator norm is generally NP-hard while it can be implemented fast for matrices, see [45].

The higher order singular value decomposition (HOSVD) is one machinery to deal with tensors which generalizes the matrix SVD to higher order tensors, see [46],[47], [48], [49] and [50]. The conceptual simplicity and computational efficiency make HOSVD popular and successful on several applications including face recognition (see [38]), genomic signal processing (see [43]) and more examples in a survey paper [51]. Basically, the HOSVD unfolds a higher order tensor into matrices and treat it with standard matrix techniques to obtain the principal singular subspaces in each dimension, see more details in Section 2. Although HOSVD is appealing, there are several fundamental theoretical mysteries yet to be uncovered.

One important problem is to study the perturbation of HOSVD when stochastic noise is observed. The difficulty comes from both methodological and theoretical aspects. The computation of HOSVD is essentially reduced to matrix SVD which can be achieved efficiently. This naive estimator is actually statistically suboptimal and further power iterations

can lead to a minimax optimal estimator, see [52], [53], [54], [55] and references therein. Another intriguing phenomenon is on the signal-to-noise ratio (SNR) exhibiting distinct computational and statistical phase transitions, which do not exist for matrices. In particular, there is a gap on SNR between statistical optimality and computational optimality for HOSVD, see [53]. For introductory simplicity *, we consider the third-order tensors where an unknown tensor $\mathbf{A} \in \mathbb{R}^{d \times d \times d}$ with multilinear ranks (r, r, r) is planted in a noisy observation \mathbf{Y} with

$$\mathbf{Y} = \mathbf{A} + \mathbf{Z} \in \mathbb{R}^{d \times d \times d}$$

with $Z(i, j, k) \sim \mathcal{N}(0, \sigma^2)$ being i.i.d. for $i, j, k \in [d]$ and $[d] := \{1, \dots, d\}$. The signal strength $\underline{\Lambda}(\mathbf{A})$ is defined as the smallest nonzero singular values of matricizations of \mathbf{A} , see definitions in Section 3.3. Let $\mathbf{U}, \mathbf{V}, \mathbf{W} \in \mathbb{R}^{d \times r}$ denote the singular vectors of \mathbf{A} in the corresponding dimensions. It was proved (see [53] and [55]) that if the signal strength $\underline{\Lambda}(\mathbf{A}) \geq D_1 \sigma d^{3/4}$ for a large enough constant $D_1 > 0$, the following bound holds

$$\begin{aligned} r^{-1/2} \max \left\{ \|\widehat{\mathbf{U}}\widehat{\mathbf{U}}^\top - \mathbf{U}\mathbf{U}^\top\|_{\ell_2}, \|\widehat{\mathbf{V}}\widehat{\mathbf{V}}^\top - \mathbf{V}\mathbf{V}^\top\|_{\ell_2}, \|\widehat{\mathbf{W}}\widehat{\mathbf{W}}^\top - \mathbf{W}\mathbf{W}^\top\|_{\ell_2} \right\} \\ = O_p \left(\frac{\sigma d^{1/2}}{\underline{\Lambda}(\mathbf{A})} + \frac{\sigma d^{3/2}}{\underline{\Lambda}^2(\mathbf{A})} \right), \end{aligned}$$

where $\widehat{\mathbf{U}}, \widehat{\mathbf{V}}, \widehat{\mathbf{W}}$ represent the naive SVD obtained from noisy tensor \mathbf{Y} and $\|\cdot\|_{\ell_2}$ denotes the Euclidean norm. Power iterations (also called higher order orthogonal iterations, see [56]) can improve the estimate (denoted by $\widetilde{\mathbf{U}}, \widetilde{\mathbf{V}}, \widetilde{\mathbf{W}}$) to

$$\begin{aligned} r^{-1/2} \max \left\{ \|\widetilde{\mathbf{U}}\widetilde{\mathbf{U}}^\top - \mathbf{U}\mathbf{U}^\top\|_{\ell_2}, \|\widetilde{\mathbf{V}}\widetilde{\mathbf{V}}^\top - \mathbf{V}\mathbf{V}^\top\|_{\ell_2}, \|\widetilde{\mathbf{W}}\widetilde{\mathbf{W}}^\top - \mathbf{W}\mathbf{W}^\top\|_{\ell_2} \right\} \\ = O_p \left(\frac{\sigma d^{1/2}}{\underline{\Lambda}(\mathbf{A})} \right), \quad (1.1) \end{aligned}$$

which is minimax optimal (see [53]). Moreover, it is demonstrated in [53] via an assump-

*Results of this paper cover the general case where \mathbf{A} is $d_1 \times d_2 \times d_3$ with multilinear ranks (r_1, r_2, r_3) , and can be easily generalized to higher order tensors.

tion on hypergraphical planted clique detection that if $\underline{\Lambda}(\mathbf{A}) = o(\sigma d^{3/4})$, then all polynomial time algorithms produce trivial estimates of $\mathbf{U}, \mathbf{V}, \mathbf{W}$.

This work is focused on the estimation of linear forms of tensor singular vectors. More specifically, consider singular vectors $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_r) \in \mathbb{R}^{d \times r}$ and our goal is to estimate $\langle \mathbf{u}_j, \mathbf{x} \rangle$ for fixed $\mathbf{x} \in \mathbb{R}^d$ and $j = 1, \dots, r$. By choosing \mathbf{x} over the canonical basis vectors in \mathbb{R}^d , we end up with an estimation of \mathbf{u}_j whose componentwise perturbation bound can be attained. Unlike the ℓ_2 -norm perturbation bound, the ℓ_∞ bound can characterize the entrywise sign consistency and entrywise significance (i.e. entrywise magnitude) of singular vectors. The componentwise signs of singular vectors have been utilized in numerous applications, such as community detection (see [57], [58], [59] and [60]). The entrywise significance is useful in submatrix localizations, see [61], [62] and references therein. We show in Section 4 that ℓ_∞ bounds require a weaker condition than ℓ_2 bounds to guarantee exact clustering in high dimensions. Furthermore, it enables us to construct a low rank estimator of \mathbf{A} with a sharp bound on $\|\hat{\mathbf{A}} - \mathbf{A}\|_{\ell_\infty}$. To the best of our knowledge, ours is the first result concerning the low rank tensor denoising with sharp ℓ_∞ bound.

To better explain our results, Suppose that \mathbf{A} is an orthogonally decomposable third order tensor with (in particular, the CP decomposition of orthogonally decomposable tensors)

$$\mathbf{A} = \sum_{k=1}^r \lambda_k (\mathbf{u}_k \otimes \mathbf{v}_k \otimes \mathbf{w}_k), \quad \lambda_1 \geq \dots \geq \lambda_r > 0 \quad (1.2)$$

where $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_r)$, $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_r)$ and $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_r)$ are $d \times r$ orthonormal matrices. The k -th eigengap is written as $\bar{g}_k(\mathcal{M}_1(\mathbf{A})) = \bar{g}_k(\mathcal{M}_2(\mathbf{A})) = \bar{g}_k(\mathcal{M}_3(\mathbf{A})) = \min(\lambda_{k-1} - \lambda_k, \lambda_k - \lambda_{k+1})$ where $\mathcal{M}_j(\mathbf{A})$ represents the matricization (see Section 2) and we preset $\lambda_0 = +\infty$ and $\lambda_{r+1} = 0$. We show that if $\bar{g}_k(\mathcal{M}_1(\mathbf{A})\mathcal{M}_1^\top(\mathbf{A})) \geq D_1(\sigma \lambda_1 d^{1/2} + \sigma^2 d^{3/2})$, the following bound holds for any $\mathbf{x} \in \mathbb{R}^d$,

$$\left| \langle \hat{\mathbf{u}}_k, \mathbf{x} \rangle - (1 + b_k)^{1/2} \langle \mathbf{u}_k, \mathbf{x} \rangle \right| = O_p \left(\|\mathbf{x}\|_{\ell_2} \frac{\lambda_1 \sigma + d \sigma^2}{\bar{g}_k(\mathcal{M}_1(\mathbf{A})\mathcal{M}_1^\top(\mathbf{A}))} \right) = O_p \left(\frac{\|\mathbf{x}\|_{\ell_2}}{d^{1/2}} \right).$$

where $b_k \in [-1/2, 0]$ is an absolute constant which does not depend on \mathbf{x} .

If $r = 1$ (rank one spiked tensor PCA model, see [52]) such that $\underline{\Lambda}(\mathbf{A}) = \bar{g}_1(\mathcal{M}_1(\mathbf{A})) = \lambda_1$, we get

$$\left| \langle \hat{\mathbf{u}}_1, \mathbf{x} \rangle - (1 + b_1)^{1/2} \langle \mathbf{u}_1, \mathbf{x} \rangle \right| = O_p \left(\frac{\sigma}{\underline{\Lambda}(\mathbf{A})} + \frac{\sigma^2 d}{\underline{\Lambda}^2(\mathbf{A})} \right) \|\mathbf{x}\|_{\ell_2}.$$

By taking \mathbf{x} over the canonical basis vectors in \mathbb{R}^d , the above fact implies that

$$\|\hat{\mathbf{u}}_1 - (1 + b_1)^{1/2} \mathbf{u}_1\|_{\ell_\infty} = O_p \left(\left(\frac{\log d}{d} \right)^{1/2} \right)$$

under the eigengap condition $\lambda_1 \geq D_1 \sigma d^{3/4}$ which is a standard requirement in tensor PCA [†]. Moreover, a low rank estimator (denoted by $\hat{\mathbf{A}}$) is constructed under the same conditions such that

$$\|\hat{\mathbf{A}} - \mathbf{A}\|_{\ell_\infty} = O_p \left(\left(\frac{\sigma^2 d}{\lambda_1} + \sigma \right) (\|\mathbf{u}_1\|_{\ell_\infty} \|\mathbf{v}_1\|_{\ell_\infty} + \|\mathbf{u}_1\|_{\ell_\infty} \|\mathbf{w}_1\|_{\ell_\infty} + \|\mathbf{v}_1\|_{\ell_\infty} \|\mathbf{w}_1\|_{\ell_\infty}) \right)$$

implying that the ℓ_∞ bound is determined by the coherence $\max \{\|\mathbf{u}_1\|_{\ell_\infty}, \|\mathbf{v}_1\|_{\ell_\infty}, \|\mathbf{w}_1\|_{\ell_\infty}\}$.

Our main contribution is on the theoretical front. The HOSVD is essentially the standard SVD computed on an unbalanced matrix where the column size is much larger than the row size. The perturbation tools such as Wedin's $\sin \Theta$ theorem ([64]) characterize the ℓ_2 bounds through the larger dimension, even when the left singular space lies in a low dimensional space. At the high level, the HOSVD is connected to the one-sided spectral analysis, see [65], [66] and references therein, which provide sharp perturbation bounds in ℓ_2 -norm. There are recent bounds (see [67] and [68]) in ℓ_∞ -norm developed under additional constraint (incoherent singular spaces) and structural noise (sparse noise). To obtain a sharp ℓ_∞ -norm bound, we borrow the instruments invented by [69] and extensively applied in [63]. Our framework is built upon a second order method of estimating the singular

[†]We shall point out that a similar result on matrix SVD has appeared in [63] which is suboptimal for tensors. Indeed, the result in [63] is established under the eigengap condition $\lambda_1 \geq D_1 \sigma d$.

subspaces, which improves the eigengap requirement than the first order method. Similar techniques have been proposed for solving tensor completion ([70]) and tensor PCA ([55]). The success of this seemingly natural treatment hinges upon delicate dealing with the correlations among higher order terms. We benefit from these ℓ_∞ -norm spectral bound by proposing a low rank estimator for tensor denoising such that entrywise perturbation is guaranteed through the tensor incoherence.

We organize this chapter as follows. Tensor notations and preliminaries on HOSVD are explained in Section 2. Our main theoretical contributions are presented in Section 3 which includes the ℓ_∞ -norm bound on singular vector perturbation and the accuracy of a low rank tensor denoising estimator. In Section 4, we apply our theoretical results on applications including high dimensional clustering and sub-tensor localizations to manifest the advantages of utilizing ℓ_∞ bounds. The proofs are provided in Section 5.

2 Preliminaries on Tensor and HOSVD

2.1 Notations

We first review some notations which will be used through the paper. We use boldfaced upper-case letters to denote tensors or matrices, and use the same letter in normal font with indices to denote its entries. We use boldfaced lower-case letters to represent vectors, and the same letter in normal font with indices to represent its entries. For notationally simplicity, our main context is focused on third-order tensors, while our results can be easily generalized to higher order tensors.

Given a third-order tensor $\mathbf{A} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$, define a linear mapping $\mathcal{M}_1 : \mathbb{R}^{d_1 \times d_2 \times d_3} \mapsto \mathbb{R}^{d_1 \times (d_2 d_3)}$ such that

$$\mathcal{M}_1(\mathbf{A})(i_1, (i_2 - 1)d_3 + i_3) = A(i_1, i_2, i_3), \quad i_1 \in [d_1], i_2 \in [d_3], i_3 \in [d_3]$$

which is conventionally called the unfolding (or matricization) of tensor \mathbf{A} . The columns

of matrix $\mathcal{M}_1(\mathbf{A})$ are called the mode-1 fibers of \mathbf{A} . The corresponding matricizations $\mathcal{M}_2(\mathbf{A})$ and $\mathcal{M}_3(\mathbf{A})$ can be defined through a similar fashion. The multilinear ranks of \mathbf{A} are then defined by:

$$r_1(\mathbf{A}) := \text{rank}(\mathcal{M}_1(\mathbf{A})), \quad r_2(\mathbf{A}) := \text{rank}(\mathcal{M}_2(\mathbf{A})), \quad r_3(\mathbf{A}) := \text{rank}(\mathcal{M}_3(\mathbf{A}))$$

Note that $r_1(\mathbf{A}), r_2(\mathbf{A}), r_3(\mathbf{A})$ are unnecessarily equal with each other in general. We write $\mathbf{r}(\mathbf{A}) := (r_1(\mathbf{A}), r_2(\mathbf{A}), r_3(\mathbf{A}))$.

The marginal product $\times_1 : \mathbb{R}^{r_1 \times r_2 \times r_3} \times \mathbb{R}^{d_1 \times r_1} \mapsto \mathbb{R}^{d_1 \times r_2 \times r_3}$ is given by

$$\mathbf{C} \times_1 \mathbf{U} = \left(\sum_{j_1=1}^{r_1} C(j_1, j_2, j_3) U(i_1, j_1) \right)_{i_1 \in [d_1], j_2 \in [r_2], j_3 \in [r_3]},$$

and \times_2 and \times_3 are defined similarly. Therefore, we write the multilinear product of tensors

$\mathbf{C} \in \mathbb{R}^{r_1 \times r_2 \times r_3}$, $\mathbf{U} \in \mathbb{R}^{d_1 \times r_1}$, $\mathbf{V} \in \mathbb{R}^{d_2 \times r_2}$ and $\mathbf{W} \in \mathbb{R}^{d_3 \times r_3}$ as

$$\mathbf{C} \cdot (\mathbf{U}, \mathbf{V}, \mathbf{W}) = \mathbf{C} \times_1 \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W} \in \mathbb{R}^{d_1 \times d_2 \times d_3}.$$

We use $\|\cdot\|$ to denote the operator norm of matrices and $\|\cdot\|_{\ell_2}$ and $\|\cdot\|_{\ell_\infty}$ to denote ℓ_2 and ℓ_∞ norms of vectors, or vectorized matrices and tensors.

2.2 HOSVD and Eigengaps

For a tensor $\mathbf{A} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$ with multilinear ranks $\mathbf{r}(\mathbf{A}) = (r_1(\mathbf{A}), r_2(\mathbf{A}), r_3(\mathbf{A}))$, let $\mathbf{U} \in \mathbb{R}^{d_1 \times r_1(\mathbf{A})}$, $\mathbf{V} \in \mathbb{R}^{d_2 \times r_2(\mathbf{A})}$ and $\mathbf{W} \in \mathbb{R}^{d_3 \times r_3(\mathbf{A})}$ be the left singular vectors of $\mathcal{M}_1(\mathbf{A})$, $\mathcal{M}_2(\mathbf{A})$ and $\mathcal{M}_3(\mathbf{A})$ respectively, which can be computed efficiently via matricization followed by thin singular value decomposition. The higher order singular value decomposition (HOSVD) refers to the decomposition

$$\mathbf{A} = \mathbf{C} \times_1 \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W} \tag{2.1}$$

where the $r_1(\mathbf{A}) \times r_2(\mathbf{A}) \times r_3(\mathbf{A})$ core tensor \mathbf{C} is obtained by $\mathbf{C} := \mathbf{A} \times_1 \mathbf{U}^\top \times_2 \mathbf{V}^\top \times_3 \mathbf{W}^\top$.

Suppose that a noisy version of \mathbf{A} is observed:

$$\mathbf{Y} = \mathbf{A} + \mathbf{Z}$$

where $\mathbf{Z} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$ is a noise tensor with i.i.d. entries satisfying $Z(i, j, k) \sim \mathcal{N}(0, \sigma^2)$. By observing \mathbf{Y} , the goal is to estimate \mathbf{U} , \mathbf{V} and \mathbf{W} . An immediate solution is to compute HOSVD of \mathbf{Y} . To this end, let $\hat{\mathbf{U}} \in \mathbb{R}^{d_1 \times r_1}$, $\hat{\mathbf{V}} \in \mathbb{R}^{d_2 \times r_2}$, $\hat{\mathbf{W}} \in \mathbb{R}^{d_3 \times r_3}$ be the corresponding top singular vectors of $\mathcal{M}_1(\mathbf{Y})$, $\mathcal{M}_2(\mathbf{Y})$ and $\mathcal{M}_3(\mathbf{Y})$. The key factor characterizing the perturbation of $\hat{\mathbf{U}}$, $\hat{\mathbf{V}}$ and $\hat{\mathbf{W}}$ is the so-called eigengap.

Observe that the computing of $\hat{\mathbf{U}}$ is essentially via matrix SVD on $\mathcal{M}_1(\mathbf{A})$. It suffices to consider eigengaps for matrices. Given a rank r matrix $\mathbf{M} \in \mathbb{R}^{m_1 \times m_2}$ with SVD:

$$\mathbf{M} = \sum_{k=1}^r \lambda_k (\mathbf{g}_k \otimes \mathbf{h}_k)$$

where singular values $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$ and $\{\mathbf{g}_1, \dots, \mathbf{g}_r\}$ are the corresponding left singular vectors and $\{\mathbf{h}_1, \dots, \mathbf{h}_r\}$ are its corresponding right singular vectors. Introduce further $\lambda_0 = +\infty$ and $\lambda_{r+1} = 0$. The k -th eigengap of matrix \mathbf{M} is then defined by

$$\bar{g}_k(\mathbf{M}) := \min(\lambda_k - \lambda_{k+1}, \lambda_{k-1} - \lambda_k), \quad \forall 1 \leq k \leq r.$$

Recall that $\mathbf{U}, \hat{\mathbf{U}} \in \mathbb{R}^{d_1 \times r_1}$ are the top- r_1 left singular vectors of $\mathcal{M}_1(\mathbf{A})$ and $\mathcal{M}_1(\mathbf{Y})$ respectively. By Wedin's $\sin \Theta$ theorem ([64]),

$$\|\hat{\mathbf{U}}\hat{\mathbf{U}}^\top - \mathbf{U}\mathbf{U}^\top\| = O\left(\frac{\|\mathcal{M}_1(\mathbf{Z})\|}{\bar{g}_{r_1}(\mathcal{M}_1(\mathbf{A})\mathcal{M}_1^\top(\mathbf{A}))}\right), \quad (2.2)$$

which is generally suboptimal especially when $\mathcal{M}_1(\mathbf{Z}) \in \mathbb{R}^{d_1 \times (d_2 d_3)}$ is unbalanced such that $d_2 d_3 \gg d_1$. Sharper bounds in ℓ_2 -norm concerning one sided perturbation have been

derived in [65] and [66]. In this paper, we focus on the perturbation bound in ℓ_∞ -norm. To this end, write $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_r)$ and $\widehat{\mathbf{U}} = (\widehat{\mathbf{u}}_1, \dots, \widehat{\mathbf{u}}_r)$. We are interested in the perturbation of linear forms $\langle \widehat{\mathbf{u}}_k, \mathbf{x} \rangle$ for $\mathbf{x} \in \mathbb{R}^{d_1}$. Similar results can be obtained for singular vectors $\widehat{\mathbf{V}}$ and $\widehat{\mathbf{W}}$.

3 Main Results

3.1 Second Order Spectral Analysis

The ℓ_∞ -norm spectral perturbation for balanced matrices has been developed in [63]. Recall that \mathbf{u}_k denotes the k -th left singular vector of $\mathcal{M}_1(\mathbf{A})$ and $\widehat{\mathbf{u}}_k$ denotes the k -th left singular vector of $\mathcal{M}_1(\mathbf{Y})$ where $\mathcal{M}_1(\mathbf{A})$ is of size $d_1 \times (d_2 d_3)$. The operator norm $\|\mathcal{M}_1(\mathbf{Z})\|$ is determined by the larger dimension $(d_1 \vee d_2 d_3)$, see Section 5. It turns out that the machinery in [63] is suboptimal meaning that the eigengap requirement becomes $\bar{g}_k(\mathcal{M}_1(\mathbf{A})\mathcal{M}_1^\top(\mathbf{A})) \geq D_1 \sigma(d_1 \vee d_2 d_3)^{1/2}$, which shall be unnecessarily strong in view of the recent results in [66], [53] and [55].

In this paper, we conduct a second order spectral analysis for $\widehat{\mathbf{U}}$. Basically, the top left singular vectors of $\mathcal{M}_1(\mathbf{Y})$ are also the top eigenvectors of $\mathcal{M}_1(\mathbf{Y})\mathcal{M}_1^\top(\mathbf{Y})$. The second order method seeks the spectral perturbation on $\mathcal{M}_1(\mathbf{Y})\mathcal{M}_1^\top(\mathbf{Y})$ instead of on $\mathcal{M}_1(\mathbf{Y})$. Clearly,

$$\mathcal{M}_1(\mathbf{Y})\mathcal{M}_1^\top(\mathbf{Y}) = \mathcal{M}_1(\mathbf{A})\mathcal{M}_1^\top(\mathbf{A}) + \mathbf{\Gamma} \in \mathbb{R}^{d_1 \times d_1}$$

where $\mathbf{\Gamma} = \mathcal{M}_1(\mathbf{A})\mathcal{M}_1^\top(\mathbf{Z}) + \mathcal{M}_1(\mathbf{Z})\mathcal{M}_1^\top(\mathbf{A}) + \mathcal{M}_1(\mathbf{Z})\mathcal{M}_1^\top(\mathbf{Z})$. Note that \mathbf{U} are the leading eigenvectors of $\mathcal{M}_1(\mathbf{A})\mathcal{M}_1^\top(\mathbf{A})$ and $\widehat{\mathbf{U}}$ are the top- r_1 eigenvectors of $\mathcal{M}_1(\mathbf{Y})\mathcal{M}_1^\top(\mathbf{Y})$. Moreover, the following fact is obvious:

$$\bar{g}_{r_1}(\mathcal{M}_1(\mathbf{A})\mathcal{M}_1^\top(\mathbf{A})) \geq \bar{g}_{r_1}^2(\mathcal{M}_1(\mathbf{A})).$$

The advantage of our method comes from the observation that even though $\mathbb{E}\|\mathcal{M}_1(\mathbf{Z})\mathcal{M}_1^\top(\mathbf{Z})\|$

is of the order $\sigma^2(d_1 \vee d_2 d_3)$, the symmetric matrix $\mathcal{M}_1(\mathbf{Z})\mathcal{M}_1^\top(\mathbf{Z})$ is concentrated at $d_2 d_3 \sigma^2 \mathbf{I}_{d_1}$ such that (see more details in Section 5)

$$\|\mathcal{M}_1(\mathbf{Z})\mathcal{M}_1^\top(\mathbf{Z}) - \sigma^2 d_2 d_3 \mathbf{I}_{d_1}\| = O_p\left(\sigma^2 (d_1 d_2 d_3)^{1/2}\right).$$

Note that subtracting by an identity matrix does not change the eigen-structure. The second order method introduces the additional term $\mathcal{M}_1(\mathbf{A})\mathcal{M}_1^\top(\mathbf{Z})$ whose operator norm is bounded by $D_1 \sigma \sqrt{d_1} \|\mathcal{M}_1(\mathbf{A})\|$ with high probability, which creates a constraint on the condition number of $\mathcal{M}_1(\mathbf{A})$. Moreover, in order to characterize a sharp perturbation bound of linear forms $\langle \hat{\mathbf{u}}_k, \mathbf{x} \rangle$, we need to pay more attention to dealing with correlations among the higher order terms than the first order method in [63].

3.2 Perturbation of Linear Forms of Singular Vectors

In this section, we present our main theorem characterizing the perturbation of linear forms $\langle \hat{\mathbf{u}}_k, \mathbf{x} \rangle$ for any $\mathbf{x} \in \mathbb{R}^{d_1}$, where $\hat{\mathbf{u}}_k$ is the k -th left singular vector of $\mathcal{M}_1(\mathbf{Y})$. Our results have similar implications as the previous work [63], meaning that the bias $\mathbb{E}\hat{\mathbf{u}}_k - \mathbf{u}_k$ is well aligned with \mathbf{u}_k . Therefore, by correcting the bias term, we are able to obtain a sharper estimation of linear forms $\langle \mathbf{u}_k, \mathbf{x} \rangle$. To this end, denote the condition number of the matrix $\mathcal{M}_1(\mathbf{A})$ by

$$\kappa(\mathcal{M}_1(\mathbf{A})) = \frac{\lambda_{\max}(\mathcal{M}_1(\mathbf{A}))}{\lambda_{\min}(\mathcal{M}_1(\mathbf{A}))}$$

where $\lambda_{\max}(\cdot)$ and $\lambda_{\min}(\cdot)$ return the largest and smallest nonzero singular values.

Theorem 3.1. *Let[‡] $\mathbf{M} := \mathcal{M}_1(\mathbf{A})$ and $\delta(d_1, d_2, d_3) = \sigma d_1^{1/2} \|\mathbf{M}\| + \sigma^2 (d_1 d_2 d_3)^{1/2}$ and suppose $d_2 d_3 e^{-d_1/2} \leq 1$. There exist absolute constants $D_1, D_2 > 0$ such that the following fact holds. Let \mathbf{u}_k be \mathbf{M} 's k -th left singular vector with multiplicity 1. If $\bar{g}_k(\mathbf{M}\mathbf{M}^\top) \geq D_1 \delta(d_1, d_2, d_3)$, there exist a constant $b_k \in [-1/2, 0]$ with $|b_k| \leq \frac{\sqrt{2}\delta(d_1, d_2, d_3)}{\bar{g}_k(\mathbf{M}\mathbf{M}^\top)}$ such that for*

[‡]Observe that if we set $d_3 = 1$ and consider the case with $d_1 \ll d_2$, then Theorem 3.1 elaborates the one-sided spectral perturbation in ℓ_∞ -norm for unbalanced (or fat) matrices.

any \mathbf{x} , the following bound holds with probability at least $1 - e^{-t}$,

$$\begin{aligned} & \left| \langle \hat{\mathbf{u}}_k, \mathbf{x} \rangle - (1 + b_k)^{1/2} \langle \mathbf{u}_k, \mathbf{x} \rangle \right| \\ & \leq D_2 \left(t^{1/2} \frac{\sigma \|\mathbf{M}\| + \sigma^2 (d_2 d_3)^{1/2}}{\bar{g}_k(\mathbf{M}\mathbf{M}^\top)} + \frac{\sigma^2 d_1}{\bar{g}_k(\mathbf{M}\mathbf{M}^\top)} \left(\frac{\delta(d_1, d_2, d_3)}{\bar{g}_k(\mathbf{M}\mathbf{M}^\top)} \right) \right) \|\mathbf{x}\|_{\ell_2} \end{aligned} \quad (3.1)$$

for all $\log 8 \leq t \leq d_1$. In particular, if $\mathbf{x} = \pm \mathbf{u}_k$, then with the same probability,

$$\begin{aligned} & \left| |\langle \hat{\mathbf{u}}_k, \mathbf{u}_k \rangle| - 1 \right| \leq \left| \sqrt{1 + b_k} - 1 \right| \\ & + D_2 \left(t^{1/2} \frac{\sigma \|\mathbf{M}\| + \sigma^2 (d_2 d_3)^{1/2}}{\bar{g}_k(\mathbf{M}\mathbf{M}^\top)} + \frac{\sigma^2 d_1}{\bar{g}_k(\mathbf{M}\mathbf{M}^\top)} \left(\frac{\delta(d_1, d_2, d_3)}{\bar{g}_k(\mathbf{M}\mathbf{M}^\top)} \right) \right). \end{aligned}$$

It is easy to check that the condition $\bar{g}_k(\mathcal{M}_1(\mathbf{A})\mathcal{M}_1^\top(\mathbf{A})) \geq D_1 \delta(d_1, d_2, d_3)$ holds whenever

$$\bar{g}_k(\mathcal{M}_1(\mathbf{A})) \geq D_1 \left(\sigma (d_1 d_2 d_3)^{1/4} + \sigma d_1^{1/2} \kappa(\mathcal{M}_1(\mathbf{A})) \right).$$

If $\kappa(\mathcal{M}_1(\mathbf{A})) \leq \left(\frac{d_2 d_3}{d_1} \right)^{1/4}$, the above bound becomes $\bar{g}_k(\mathcal{M}_1(\mathbf{A})) \geq D_1 \sigma (d_1 d_2 d_3)^{1/4}$ which is a standard requirement in tensor SVD or PCA, see [53], [54] and [52]. By taking \mathbf{x} over the standard basis vectors in \mathbb{R}^{d_1} and choosing $t \geq D_3 \log d_1$, we end up with a ℓ_∞ -norm perturbation bound for empirical singular vector $\hat{\mathbf{u}}_k$.

Corollary 3.1. Under the conditions in Theorem 3.1, there exists a universal constant $D_1 > 0$ such that the following bound holds with probability at least $1 - \frac{1}{d_1}$,

$$\left\| \hat{\mathbf{u}}_k - (1 + b_k)^{1/2} \mathbf{u}_k \right\|_{\ell_\infty} \leq D_1 \left(\left(\frac{\log d_1}{d_1} \right)^{1/2} + \left(\frac{d_1}{d_2 d_3} \right)^{1/2} \right).$$

If $d_1 \asymp d_2 \asymp d_3 \asymp d$, we obtain

$$\mathbb{P} \left(\left\| \hat{\mathbf{u}}_k - (1 + b_k)^{1/2} \mathbf{u}_k \right\|_{\ell_\infty} \geq D_1 \left(\frac{\log d}{d} \right)^{1/2} \right) \leq \frac{1}{d}$$

which has an analogous form to the perturbation bound in [63] implying a famous delo-

calization phenomenon in random matrix theory, see [71] and [72] and references therein. The bias b_k is usually unknown and we borrow the idea in [63] to estimate b_k based on two independent samples, which happens in the application of tensor decomposition for gene expression, usually multiple independent copies are available, see [73].

Suppose that two independent noisy version of $\mathbf{A} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$ are observed with $\mathbf{Y}^{(1)} = \mathbf{A} + \mathbf{Z}^{(1)}$ and $\mathbf{Y}^{(2)} = \mathbf{A} + \mathbf{Z}^{(2)}$ where $\mathbf{Z}^{(1)}$ and $\mathbf{Z}^{(2)}$ have i.i.d. centered Gaussian entries with variance σ^2 . Let $\hat{\mathbf{u}}_k^{(1)}$ and $\hat{\mathbf{u}}_k^{(2)}$ denote the k -th left singular vector of $\mathcal{M}_1(\mathbf{Y}^{(1)})$ and $\mathcal{M}_1(\mathbf{Y}^{(2)})$ respectively. The signs of $\hat{\mathbf{u}}_k^{(1)}$ and $\hat{\mathbf{u}}_k^{(2)}$ are chosen such that $\langle \hat{\mathbf{u}}_k^{(1)}, \hat{\mathbf{u}}_k^{(2)} \rangle \geq 0$. Define the estimator of b_k by

$$\hat{b}_k := \langle \hat{\mathbf{u}}_k^{(1)}, \hat{\mathbf{u}}_k^{(2)} \rangle - 1.$$

Define the scaled version of empirical singular vector $\tilde{\mathbf{u}}_k := \frac{\hat{\mathbf{u}}_k}{(1+b_k)^{1/2}}$, which is not necessarily a unit vector.

Theorem 3.2. *Under the assumptions in Theorem 3.1, there exists an absolute constant $D_1 > 0$ such that for any $\mathbf{x} \in \mathbb{R}^{d_1}$, the follow bound holds with probability at least $1 - e^{-t}$ for all $t \geq 0$,*

$$|\hat{b}_k - b_k| \leq D_1 \left(t^{1/2} \frac{\sigma \|\mathbf{M}\| + \sigma^2 (d_2 d_3)^{1/2}}{\bar{g}_k(\mathbf{M}\mathbf{M}^\top)} + \frac{\sigma^2 d_1}{\bar{g}_k(\mathbf{M}\mathbf{M}^\top)} \left(\frac{\delta(d_1, d_2, d_3)}{\bar{g}_k(\mathbf{M}\mathbf{M}^\top)} \right) \right)$$

and

$$|\langle \tilde{\mathbf{u}}_k - \mathbf{u}_k, \mathbf{x} \rangle| \leq D_1 \left(t^{1/2} \frac{\sigma \|\mathbf{M}\| + \sigma^2 (d_2 d_3)^{1/2}}{\bar{g}_k(\mathbf{M}\mathbf{M}^\top)} + \frac{\sigma^2 d_1}{\bar{g}_k(\mathbf{M}\mathbf{M}^\top)} \left(\frac{\delta(d_1, d_2, d_3)}{\bar{g}_k(\mathbf{M}\mathbf{M}^\top)} \right) \right) \|\mathbf{x}\|_{\ell_2}$$

where $\mathbf{M} = \mathcal{M}_1(\mathbf{A})$.

Remark 1. If $d/2 \leq \min_k d_k \leq \max_k d_k \leq 2d$, we get

$$\mathbb{P} \left(\|\tilde{\mathbf{u}}_k - \mathbf{u}_k\|_{\ell_\infty} \geq D_1 \left(\frac{\log d}{d} \right)^{1/2} \right) \leq \frac{1}{d}.$$

Moreover, if $\text{rank}(\mathbf{A}) = (1, 1, 1)$, we can write $\|\tilde{\mathbf{u}}_1 - \mathbf{u}_1\|_{\ell_\infty} = O_p\left(\frac{\sigma \log^{1/2} d}{\underline{\Lambda}(\mathbf{A})} + \frac{\sigma^2 d \log^{1/2} d}{\underline{\Lambda}^2(\mathbf{A})}\right)$ where $\underline{\Lambda}(\mathbf{A}) = \lambda_{\min}(\mathcal{M}_1(\mathbf{A}))$. Note that $\|\tilde{\mathbf{u}}_1 - \mathbf{u}_1\|_{\ell_2} = O_p\left(\frac{\sigma d^{1/2}}{\underline{\Lambda}(\mathbf{A})} + \frac{\sigma^2 d^{3/2}}{\underline{\Lambda}^2(\mathbf{A})}\right)$, see [53]. Therefore, our one-sided SVD perturbation bound in ℓ_∞ -norm for a matrix $\mathbf{M} \in \mathbb{R}^{d_1 \times d_2}$ is optimal if it is ultra-fat such that $d_1^2 \leq d_2$.

3.3 Low Rank Tensor Denoising ℓ_∞ Bound

In this section, we consider low rank estimate of \mathbf{A} through projection of \mathbf{Y} . Let $\tilde{\mathbf{U}} = (\tilde{\mathbf{u}}_1, \dots, \tilde{\mathbf{u}}_{r_1}) \in \mathbb{R}^{d_1 \times r_1}$ be scaled singular vectors each of which is computed as in Theorem 3.2. Similarly, let $\tilde{\mathbf{V}} \in \mathbb{R}^{d_2 \times r_2}$ and $\tilde{\mathbf{W}} \in \mathbb{R}^{d_3 \times r_3}$ be the corresponding scaled singular vectors computed from $\mathcal{M}_2(\mathbf{Y})$ and $\mathcal{M}_3(\mathbf{Y})$. Define the low rank estimate

$$\tilde{\mathbf{A}} := \mathbf{Y} \times_1 \mathbf{P}_{\tilde{\mathbf{U}}} \times_2 \mathbf{P}_{\tilde{\mathbf{V}}} \times_3 \mathbf{P}_{\tilde{\mathbf{W}}}$$

where $\mathbf{P}_{\tilde{\mathbf{U}}}$ represents the scaled projector $\mathbf{P}_{\tilde{\mathbf{U}}} := \tilde{\mathbf{U}}\tilde{\mathbf{U}}^\top$. Clearly, $\text{rank}(\tilde{\mathbf{A}}) = (r_1, r_2, r_3)$ which serves as a low rank estimate of \mathbf{A} . We characterize the entrywise accuracy of $\tilde{\mathbf{A}}$, namely, the upper bound of $\|\tilde{\mathbf{A}} - \mathbf{A}\|_{\ell_\infty}$ in terms of the coherence of \mathbf{U} , \mathbf{V} and \mathbf{W} . Our $\|\tilde{\mathbf{A}} - \mathbf{A}\|_{\ell_\infty}$ bound relies on the simultaneous ℓ_∞ -norm perturbation bounds on $\tilde{\mathbf{u}}_k, \tilde{\mathbf{v}}_k, \tilde{\mathbf{w}}_k$. We shall need the following conditions on the eigengaps: for a large enough constant $D_1 > 0$,

$$\bar{g}_k(\mathcal{M}_1(\mathbf{A})\mathcal{M}_1^\top(\mathbf{A})) \geq D_1 \left(\sigma d_1^{1/2} \bar{\Lambda}(\mathbf{A}) + \sigma^2 (d_1 d_2 d_3)^{1/2} \right), \quad 1 \leq k \leq r_1, \quad (3.2)$$

$$\bar{g}_k(\mathcal{M}_2(\mathbf{A})\mathcal{M}_2^\top(\mathbf{A})) \geq D_1 \left(\sigma d_2^{1/2} \bar{\Lambda}(\mathbf{A}) + \sigma^2 (d_1 d_2 d_3)^{1/2} \right), \quad 1 \leq k \leq r_2, \quad (3.3)$$

$$\bar{g}_k(\mathcal{M}_3(\mathbf{A})\mathcal{M}_3^\top(\mathbf{A})) \geq D_1 \left(\sigma d_3^{1/2} \bar{\Lambda}(\mathbf{A}) + \sigma^2 (d_1 d_2 d_3)^{1/2} \right), \quad 1 \leq k \leq r_3, \quad (3.4)$$

where

$$\bar{\Lambda}(\mathbf{A}) := \max \left\{ \lambda_{\max}(\mathcal{M}_1(\mathbf{A})), \lambda_{\max}(\mathcal{M}_2(\mathbf{A})), \lambda_{\max}(\mathcal{M}_3(\mathbf{A})) \right\}.$$

Similarly, define

$$\underline{\Lambda}(\mathbf{A}) := \min \left\{ \lambda_{\min}(\mathcal{M}_1(\mathbf{A})), \lambda_{\min}(\mathcal{M}_2(\mathbf{A})), \lambda_{\min}(\mathcal{M}_3(\mathbf{A})) \right\}$$

and the overall eigengap

$$\bar{g}_{\min}(\mathbf{A}) := \min \left\{ \bar{g}_{k_1}^{1/2}(\mathcal{M}_1(\mathbf{A})\mathcal{M}_1^\top(\mathbf{A})), \bar{g}_{k_2}^{1/2}(\mathcal{M}_2(\mathbf{A})\mathcal{M}_2^\top(\mathbf{A})), \bar{g}_{k_3}^{1/2}(\mathcal{M}_3(\mathbf{A})\mathcal{M}_3^\top(\mathbf{A})), \right. \\ \left. , 1 \leq k_1 \leq r_1, 1 \leq k_2 \leq r_2, 1 \leq k_3 \leq r_3 \right\}.$$

By definition, it is clear that $\underline{\Lambda}(\mathbf{A}) \geq \bar{g}_{\min}(\mathbf{A})$.

Theorem 3.3. *Suppose conditions (3.2) (3.3) (3.4) hold and assume that for all $i \in [d_1], j \in [d_2], k \in [d_3]$,*

$$\|\mathbf{U}^\top \mathbf{e}_i\|_{\ell_2} \leq \mu_{\mathbf{U}} \sqrt{\frac{r_1}{d_1}}, \quad \|\mathbf{V}^\top \mathbf{e}_j\|_{\ell_2} \leq \mu_{\mathbf{V}} \sqrt{\frac{r_2}{d_2}}, \quad \|\mathbf{W}^\top \mathbf{e}_k\|_{\ell_2} \leq \mu_{\mathbf{W}} \sqrt{\frac{r_3}{d_3}}$$

for some constants $\mu_{\mathbf{U}}, \mu_{\mathbf{V}}, \mu_{\mathbf{W}} \geq 0$. Suppose that $\frac{d}{2} \leq \min_{1 \leq k \leq 3} d_k \leq \max_{1 \leq k \leq 3} d_k \leq 2d$ and $\frac{r}{2} \leq \min_{1 \leq k \leq 3} r_k \leq \max_{1 \leq k \leq 3} r_k \leq 2r$. There exists an absolute constant $D_2 > 0$ such that, with probability at least $1 - \frac{1}{d}$,

$$\|\tilde{\mathbf{A}} - \mathbf{A}\|_{\ell_\infty} \leq D_2 \sigma r^3 \left(\frac{\tilde{\kappa}(\mathbf{A})\sigma}{\bar{g}_{\min}(\mathbf{A})} + \frac{\tilde{\kappa}^2(\mathbf{A})}{d} \right) (\mu_{\mathbf{U}}\mu_{\mathbf{V}} + \mu_{\mathbf{U}}\mu_{\mathbf{W}} + \mu_{\mathbf{V}}\mu_{\mathbf{W}}) \log^{3/2} d$$

where $\tilde{\kappa}(\mathbf{A}) = \bar{\Lambda}(\mathbf{A})/\bar{g}_{\min}(\mathbf{A})$.

Remark 2. To highlight the contribution of Theorem 3.3, let $r = O(1)$ and $\tilde{\kappa}(\mathbf{A}) = O(1)$.

Note that if the coherence constants $\mu_{\mathbf{U}}, \mu_{\mathbf{V}}, \mu_{\mathbf{W}} = d^{(\frac{3}{4}-\varepsilon)/2}$ for $\varepsilon \in (0, 3/4)$, i.e., $\mathbf{U}, \mathbf{V}, \mathbf{W}$

can be almost spiked, under the minimal eigengap $\bar{g}_{\min}(\mathbf{A}) \gtrsim \sigma d^{3/4}$, then

$$\|\tilde{\mathbf{A}} - \mathbf{A}\|_{\ell_\infty} = O_p\left(\frac{\sigma}{d^\varepsilon} \log^{3/2} d\right)$$

It worths to point out that the minimax optimal bound of estimating \mathbf{A} in ℓ_2 -norm is $O(\sigma d^{1/2})$, see [53]. Theorem 3.3 is more interesting when \mathbf{A} is incoherent such that $\mu_{\mathbf{U}}, \mu_{\mathbf{V}}, \mu_{\mathbf{W}} = O(1)$. We conclude that

$$\|\tilde{\mathbf{A}} - \mathbf{A}\|_{\ell_\infty} = O_p\left(\left(\frac{\sigma^2}{\bar{g}_{\min}(\mathbf{A})} + \frac{\sigma}{d}\right) \log^{3/2} d\right) = O_p\left(\frac{\sigma}{d^{3/4}} \log^{3/2} d\right).$$

4 Applications

In this section, we review two applications of ℓ_∞ -norm. It is interesting to observe that it is unnecessary to estimate the bias b_k in these applications.

4.1 High Dimensional Clustering

Many statistical and machine learning tasks are associated with clustering high dimensional data, see [74], [75], [76], [77], [78] and references therein. We consider a two-class Gaussian mixture model such that each data point $\mathbf{y}_i \in \mathbb{R}^p$ can be represented by

$$\mathbf{y}_i = -\ell_i \boldsymbol{\beta} + (1 - \ell_i) \boldsymbol{\beta} + \boldsymbol{\varepsilon}_i \in \mathbb{R}^p$$

where the associated label $\ell_i \in \{0, 1\}$ for $i = 1, 2, \dots, n$ is unknown and the noise vector $\boldsymbol{\varepsilon}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$. The vector $\boldsymbol{\beta} \in \mathbb{R}^p$ is unknown with $p \gg n$.

Given the data matrix

$$\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)^\top \in \mathbb{R}^{n \times p},$$

the goal is to conduct bi-clustering. Let $n_k := \text{Card}(\{1 \leq i \leq n : \ell_i = k\})$ for $k = 0, 1$ such that $n_0 + n_1 = n$. Observe that $\mathbb{E}\mathbf{Y}$ has rank 1 and its leading left singular vector

$\mathbf{u} \in \mathbb{R}^n$ with

$$u(i) = \frac{1 - \ell_i}{n^{1/2}} - \frac{\ell_i}{n^{1/2}}, \quad 1 \leq i \leq n.$$

The signs of \mathbf{u} immediately produce the cluster membership. Moreover, the leading singular value of $\mathbb{E}\mathbf{Y}$ is $n^{1/2}\|\boldsymbol{\beta}\|_{\ell_2}$. Let $\hat{\mathbf{u}}$ denotes the leading left singular vector of \mathbf{Y} . By Corollary 3.1, if $\|\boldsymbol{\beta}\|_{\ell_2} \geq D_1(1 \vee (p/n)^{1/4})$ such that $|(1 + b_k)^{-1/2} - 1| \leq 1/2$, then

$$\mathbb{P}\left(\|\hat{\mathbf{u}} - (1 + b_k)^{1/2}\mathbf{u}\|_{\ell_\infty} \leq D_2\left(\frac{1}{\|\boldsymbol{\beta}\|_{\ell_2}} + \frac{(p/n)^{1/2}}{\|\boldsymbol{\beta}\|_{\ell_2}^2}\right)\left(\frac{1}{\|\boldsymbol{\beta}\|_{\ell_2}^2} + \sqrt{\frac{\log n}{n}}\right)\right) \geq 1 - \frac{1}{n}.$$

On this event, if $\|\boldsymbol{\beta}\|_{\ell_2} \geq D_1\left(n^{1/6} \vee p^{1/8} \vee (p \log(n)/n)^{1/4}\right)$

$$\begin{aligned} \|\hat{\mathbf{u}} - \mathbf{u}\|_{\ell_\infty} &\leq \|\hat{\mathbf{u}} - (1 + b_k)^{1/2}\mathbf{u}\|_{\ell_\infty} + |(1 + b_k)^{-1/2} - 1| \|\mathbf{u}\|_{\ell_\infty} \\ &\leq \|\hat{\mathbf{u}} - (1 + b_k)^{1/2}\mathbf{u}\|_{\ell_\infty} + \frac{1}{2n^{1/2}} \leq \frac{3}{4n^{1/2}} \end{aligned}$$

implying that if $\ell_i = \ell_j$, then $\text{sign}(\hat{u}(i)) = \text{sign}(\hat{u}(j))$ for all $1 \leq i, j \leq n$. The above analysis also implies that it is unnecessary to estimate b_k in this application, since scaling will not change the entrywise signs. Therefore, in order to guarantee exact clustering, the ℓ_∞ bound requires

$$\|\boldsymbol{\beta}\|_{\ell_2} \geq D_1\left(n^{1/6} \vee p^{1/8} \vee (p \log(n)/n)^{1/4}\right),$$

while the ℓ_2 bound in [66] requires

$$\|\boldsymbol{\beta}\|_{\ell_2} \geq D_1\left(n^{1/2} \vee p^{1/4} \vee (p/n)^{1/4}\right)$$

for exact clustering.

Remark 3. The above framework can be directly generalized to Gaussian mixture model with k -clusters. Suppose that the j -th cluster has mean vector $\boldsymbol{\beta}_j$ and size n_j , then without

loss of generality, the data matrix $\mathbf{Y} = \mathbf{M} + \mathbf{Z}$

$$\mathbf{M} = \left(\underbrace{\beta_1, \dots, \beta_1}_{n_1}, \dots, \underbrace{\beta_j, \dots, \beta_j}_{n_j}, \dots, \underbrace{\beta_k, \dots, \beta_k}_{n_k} \right)^\top \in \mathbb{R}^{N \times p}$$

with $N = \sum_{j=1}^k n_j$ and $\mathbf{Z} \in \mathbb{R}^{N \times p}$ having i.i.d. standard Gaussian entries. Observe that $\text{rank}(\mathbf{M}) \leq k$, it suffices to consider the top- k left singular vectors of \mathbf{M} . However, it requires nontrivial effort to investigate the eigengaps of \mathbf{M} without further assumptions on $\{\beta_j\}_{j=1}^k$. In the case that $n_j = n$ and β_1, \dots, β_k are mutually orthogonal such that $\|\beta_1\|_{\ell_2} \geq \dots \geq \|\beta_k\|_{\ell_2}$, then \mathbf{M} 's top- k singular values are $\lambda_j = \sqrt{n_j} \|\beta_j\|_{\ell_2}$, $1 \leq j \leq k$. Clearly, the non-zero entries of \mathbf{M} 's top- k left singular vectors provide the clustering membership. By Theorem 3.1, if $\Delta_j \geq C_1 \sqrt{k} \|\beta_1\|_{\ell_2} + C_2 (kp/n)^{1/2}$ where $\Delta_j = \min\{(\|\beta_j\|_{\ell_2}^2 - \|\beta_{j+1}\|_{\ell_2}^2), (\|\beta_{j-1}\|_{\ell_2}^2 - \|\beta_j\|_{\ell_2}^2)\}$, then

$$\|\hat{\mathbf{u}}_j - \sqrt{1 + b_j} \mathbf{u}_j\|_{\ell_\infty} = O_p \left(\left(\frac{\|\beta_1\|_{\ell_2}}{\Delta_j} + \frac{(p/n)^{1/2}}{\Delta_j} \right) \left(\frac{k^{3/2}}{\Delta_j} + \sqrt{\frac{k \log n}{n}} \right) \right)$$

for all $1 \leq j \leq k$.

4.2 Subtensor Localization

In gene expression association analysis (see [73], [79], [80] and [81]) and planted clique detection (see [82], [83] and [84]), the goal is equivalent to localizing a sub-tensor whose entries are statistically more significant than the others. One simple model characterizing this type of tensor data is as

$$\mathbf{Y} = \lambda \mathbf{1}_{C_1} \otimes \mathbf{1}_{C_2} \otimes \mathbf{1}_{C_3} + \mathbf{Z} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$$

with $C_k = \cup_{j=1}^{b_k} C_k^{(j)} \subset [d_k]$ where $\{C_k^{(1)}, \dots, C_k^{(b_k)}\}$ are disjoint subsets of $[d_k]$ for $k = 1, 2, 3$, i.e., there are b_k dense blocks in the k -th direction. Then, in total, there are $b_1 b_2 b_3$ dense blocks in $\mathbb{E}\mathbf{Y}$. The vector $\mathbf{1}_{C_k} \in \mathbb{R}^{p_k}$ is a zero-or-one vector whose entry equals

1 only when the index belongs to C_k . The noise tensor \mathbf{Z} has i.i.d. entries such that $Z(i, j, k) \sim \mathcal{N}(0, 1)$. Given the noisy observation of \mathbf{Y} , the goal is to localize the unknown subsets $\{C_1^{(j)}\}_{j=1}^{b_1}$, $\{C_2^{(j)}\}_{j=1}^{b_2}$ and $\{C_3^{(j)}\}_{j=1}^{b_3}$. The appealing scenario is $\lambda = O(1)$, since otherwise the signal is so strong that the problem can be easily solved by just looking at each entry. The tensor $\mathbb{E}\mathbf{Y}$ has rank 1 with leading singular value $\lambda|C_1|^{1/2}|C_2|^{1/2}|C_3|^{1/2}$ and corresponding singular vectors

$$\mathbf{u} = \frac{1}{|C_1|^{1/2}} \mathbf{1}_{C_1}, \quad \mathbf{v} = \frac{1}{|C_2|^{1/2}} \mathbf{1}_{C_2} \quad \text{and} \quad \mathbf{w} = \frac{1}{|C_3|^{1/2}} \mathbf{1}_{C_3},$$

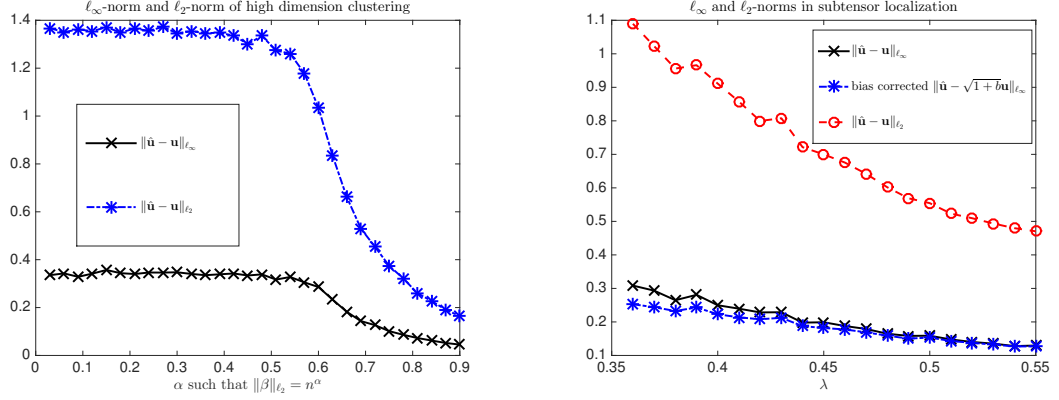
where $|C|$ denotes the cardinality of C . By Theorem 3.1, if $\lambda \geq D_1 \frac{(d_1 d_2 d_3)^{1/4}}{|C_1|^{1/2}|C_2|^{1/2}|C_3|^{1/2}}$ for a large enough constant $D_1 > 0$, then with probability at least $1 - \frac{1}{d_{\max}}$ where $d_{\max} := (d_1 \vee d_2 \vee d_3)$ and we assume $d_{\max} \leq D_1(d_1 d_2 d_3)^{1/2}$,

$$\begin{aligned} & \|\hat{\mathbf{u}} - (1 + b_1)^{1/2} \mathbf{u}\|_{\ell_\infty} \\ & \leq \frac{D_1}{\lambda|C_1|^{1/2}|C_2|^{1/2}|C_3|^{1/2}} + \frac{D_1(d_2 d_3)^{1/2}}{\lambda^2|C_1||C_2||C_3|} + \frac{D_1 d_1}{\lambda^2|C_1||C_2||C_3|} \left(\frac{(d_1 d_2 d_3)^{1/2}}{\lambda^2|C_1||C_2||C_3|} \right). \\ & \leq D_1 \left(\frac{1}{d_1^{1/2}} + \left(\frac{d_1}{d_2 d_3} \right)^{1/2} \right). \end{aligned} \tag{4.1}$$

If we let \hat{C}_1 denote the locations of entries of $\hat{\mathbf{u}}$ whose magnitudes are among the $|C_1|$ largest, it is straightforward to see that $\hat{C}_1 = C_1$ on the above event if $D_2|C_1|d_1 \leq d_2 d_3$ for a large enough constant $D_2 > 0$. Note that it is also unnecessary to estimate b_1 if we are only interested in the top- $|C_1|$ largest entries of $|\hat{\mathbf{u}}|$.

4.3 Numerical Experiments

We present simulation results of experiments on the above applications. In high dimensional clustering, we randomly sample a vector $\boldsymbol{\beta} \in \mathbb{R}^p$ with $p = 1000$. Fixed $\boldsymbol{\beta}$, $n/2 = 50$ random vectors are sampled from distribution $\mathcal{N}(\boldsymbol{\beta}, \mathbf{I}_p)$ and $n/2 = 50$ random vectors are sampled from distribution $\mathcal{N}(-\boldsymbol{\beta}, \mathbf{I}_p)$. Then, we compare between the leading left singular



(a) High dimension clustering: a significant gap between ℓ_∞ -norm and ℓ_2 -norm. (b) Subtensor localization: the bias correction indeed improves the ℓ_∞ norm.

Figure 4.1: Comparison on ℓ_∞ -norm, ℓ_2 -norm and bias corrected ℓ_∞ norm in high dimension clustering and subtensor localization.

vector of \mathbf{Y} and leading left singular vector of $\mathbb{E}\mathbf{Y}$, i.e., $\|\hat{\mathbf{u}} - \mathbf{u}\|_\infty$ and $\|\hat{\mathbf{u}} - \mathbf{u}\|_{\ell_2}$, without bias correction. For each $\|\beta\|_{\ell_2} \in [n^{0.03}, n^{0.9}]$, the loss is reported by averaging 30 independent simulations. The results are displayed in Figure 4.1a where we can observe a significant gap between ℓ_2 -norm and ℓ_∞ -norm. It explains why ℓ_∞ -norm is more powerful for exact clustering than ℓ_2 -norm in this application.

In subtensor localization, we fix $d_1 = d_2 = d_3 = 100$ and $C_1 = C_2 = C_3 = [20]$, i.e., the subtensor is the bottom-left-front corner of $\mathbb{E}\mathbf{Y}$. The λ is varied from 0.36 to 0.55. For each λ , we report the average ℓ_∞ -norm, ℓ_2 -norm and bias corrected ℓ_∞ -norm, all from 30 independent simulations. It is interesting to observe that actually the bias correction indeed can improve the ℓ_∞ -norm when λ is small. The results are displayed in Figure 4.1b.

5 Proofs

For notational brevity, we write $A \lesssim B$ if there exists an absolute constant D_1 such that $A \leq D_1 B$. A similar notation would be \gtrsim and $A \asymp B$ means that $A \lesssim B$ and $A \gtrsim B$ simultaneously. If the constant D_1 depends on some parameter γ , we shall write \lesssim_γ , \gtrsim_γ and \asymp_γ .

Recall that the HOSVD is translated directly from SVD on $\mathcal{M}_1(\mathbf{A})$ and the matrix

perturbation model $\mathcal{M}_1(\mathbf{Y}) = \mathcal{M}_1(\mathbf{A}) + \mathcal{M}_1(\mathbf{Z})$. Without loss of generality, it suffices to focus on matrices with unbalanced sizes. In the remaining context, we write $\mathbf{A}, \mathbf{Z}, \mathbf{Y} \in \mathbb{R}^{m_1 \times m_2}$ instead of $\mathcal{M}_1(\mathbf{A}), \mathcal{M}_1(\mathbf{Z}), \mathcal{M}_1(\mathbf{Y}) \in \mathbb{R}^{m_1 \times m_2}$, where $m_1 = d_1$ and $m_2 = d_2 d_3$ such that $m_1 \lesssim m_2$. The second order spectral analysis begins with

$$\mathbf{Y}\mathbf{Y}^\top = \mathbf{A}\mathbf{A}^\top + \mathbf{\Gamma}, \quad \text{where} \quad \mathbf{\Gamma} = \mathbf{A}\mathbf{Z}^\top + \mathbf{Z}\mathbf{A}^\top + \mathbf{Z}\mathbf{Z}^\top.$$

Suppose that \mathbf{A} has the thin singular value decomposition

$$\mathbf{A} = \sum_{k=1}^{r_1} \lambda_k (\mathbf{u}_k \otimes \mathbf{h}_k) \in \mathbb{R}^{m_1 \times m_2}$$

where $\{\mathbf{h}_1, \dots, \mathbf{h}_{r_1}\} \subset \text{span}\{\mathbf{v}_j \otimes \mathbf{w}_k^\top : j \in [r_2], k \in [r_3]\}$ are the right singular vectors of \mathbf{A} . Moreover, $\mathbf{A}\mathbf{A}^\top$ admits the eigen-decomposition:

$$\mathbf{A}\mathbf{A}^\top = \sum_{k=1}^{r_1} \lambda_k^2 (\mathbf{u}_k \otimes \mathbf{u}_k).$$

In an identical fashion, denote the eigen-decomposition of $\mathbf{Y}\mathbf{Y}^\top$ by

$$\mathbf{Y}\mathbf{Y}^\top = \sum_{k=1}^{m_1} \hat{\lambda}_k^2 (\hat{\mathbf{u}}_k \otimes \hat{\mathbf{u}}_k).$$

Even though Theorem 3.1 and Theorem 3.2 are stated when the singular value λ_k has multiplicity 1, we present more general results in this section. Note that when there are repeated singular values, the singular vectors are not uniquely defined. In this case, let $\mu_1 > \mu_2 > \dots > \mu_s > 0$ be distinct singular values of \mathbf{A} with $s \leq r_1$. Denote $\Delta_k := \{j : \lambda_j = \mu_k\}$ for $1 \leq k \leq s$ and $\nu_k := \text{Card}(\Delta_k)$ the multiplicity of μ_k . Let $\mu_{s+1} = 0$ which is a trivial eigenvalue of $\mathbf{A}\mathbf{A}^\top$ with multiplicity $m_1 - r_1$. Then, the spectral decomposition

of $\mathbf{A}\mathbf{A}^\top$ can be represented as

$$\mathbf{A}\mathbf{A}^\top = \sum_{k=1}^{s+1} \mu_k^2 \mathbf{P}_k^{uu}$$

where the spectral projector $\mathbf{P}_k^{uu} := \sum_{j \in \Delta_k} \mathbf{u}_j \otimes \mathbf{u}_j$ which is uniquely defined. Correspondingly, define the empirical spectral projector based on eigen-decomposition of $\mathbf{Y}\mathbf{Y}^\top$,

$$\hat{\mathbf{P}}_k^{uu} := \sum_{j \in \Delta_k} \hat{\mathbf{u}}_j \otimes \hat{\mathbf{u}}_j.$$

We develop a sharp concentration bound for bilinear forms $\langle \hat{\mathbf{P}}_k^{uu} \mathbf{x}, \mathbf{y} \rangle$ for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{m_1}$. Observe that $\mathbf{Y}\mathbf{Y}^\top$ has an identical eigen-space as $\mathbf{Y}\mathbf{Y}^\top - m_2 \sigma^2 \mathbf{I}_{m_1}$. Let $\hat{\Gamma} := \Gamma - m_2 \sigma^2 \mathbf{I}_{m_1}$ and the spectral analysis shall be realized on $\mathbf{A}\mathbf{A}^\top + \hat{\Gamma}$.

Several preliminary facts are introduced as follows. It is clear that the k -th eigengap is $\bar{g}_k(\mathbf{A}\mathbf{A}^\top) := \min(\mu_{k-1}^2 - \mu_k^2, \mu_k^2 - \mu_{k+1}^2)$ for $1 \leq k \leq s$, where we set $\mu_0 = +\infty$. The proof of Lemma 3 is provided in the Appendix.

Lemma 3. For any deterministic matrix $\mathbf{B} \in \mathbb{R}^{m_3 \times m_2}$, the following bounds hold

$$\begin{aligned} \mathbb{E} \|\mathbf{B}\mathbf{Z}^\top\| &\lesssim \sigma \|\mathbf{B}\| \left(m_1^{1/2} + m_3^{1/2} + (m_1 m_3)^{1/4} \right) \\ \|\mathbb{E} \mathbf{Z}\mathbf{Z}^\top - m_2 \sigma^2 \mathbf{I}_{m_1}\| &\lesssim \sigma^2 (m_1 m_2)^{1/2}. \end{aligned} \quad (5.1)$$

For any $t > 0$, the following inequalities hold with probability at least $1 - e^{-t}$,

$$\begin{aligned} \|\mathbf{B}\mathbf{Z}^\top\| &\lesssim \sigma \|\mathbf{B}\| \left(m_1^{1/2} + m_3^{1/2} + (m_1 m_3)^{1/4} + t^{1/2} + (m_1 t)^{1/4} \right) \\ \|\mathbf{Z}\mathbf{Z}^\top - m_2 \sigma^2 \mathbf{I}_{m_1}\| &\lesssim \sigma^2 m_2^{1/2} (m_1^{1/2} + t^{1/2}). \end{aligned} \quad (5.2)$$

5.1 Proof of Theorem 3.1

To this end, define

$$\mathbf{C}_k^{uu} := \sum_{s \neq k} \frac{1}{\mu_s^2 - \mu_k^2} \mathbf{P}_s^{uu}$$

and

$$\mathbf{P}_k^{hh} := \sum_{j \in \Delta_k} \mathbf{h}_j \otimes \mathbf{h}_j.$$

Theorem 3.1 is decomposed of two separate components. Theorem 5.1 provides the concentration bound for $|\langle \mathbf{P}_k \mathbf{x}, \mathbf{y} \rangle - \mathbb{E} \langle \mathbf{P}_k \mathbf{x}, \mathbf{y} \rangle|$ by Gaussian isoperimetric inequality and the proof is postponed to the Appendix. In Theorem 5.2, we characterize the bias $\mathbb{E} \hat{\mathbf{P}}_k^{uu} - \mathbf{P}_k^{uu}$.

Theorem 5.1. *Let $\delta(m_1, m_2) := \mu_1 \sigma m_1^{1/2} + \sigma^2 (m_1 m_2)^{1/2}$ and suppose that $\bar{g}_k(\mathbf{A} \mathbf{A}^\top) \geq D_1 \delta(m_1, m_2)$ for a large enough constant $D_1 > 0$. Then, for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{m_1}$, there exists an absolute constant $D_2 > 0$ such that for all $\log 8 \leq t \lesssim m_1$, the following bound holds with probability at least $1 - e^{-t}$,*

$$|\langle \hat{\mathbf{P}}_k^{uu} \mathbf{x}, \mathbf{y} \rangle - \mathbb{E} \langle \hat{\mathbf{P}}_k^{uu} \mathbf{x}, \mathbf{y} \rangle| \leq D_2 t^{1/2} \left(\frac{\sigma \mu_1 + \sigma^2 m_2^{1/2}}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \right) \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}.$$

The following spectral representation formula is needed whose proof can be found in [69].

Lemma 4. The following bound holds

$$\|\hat{\mathbf{P}}_k^{uu} - \mathbf{P}_k^{uu}\| \leq \frac{4 \|\hat{\Gamma}\|}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)}.$$

Moreover, $\hat{\mathbf{P}}_k^{uu}$ can be represented as

$$\hat{\mathbf{P}}_k^{uu} - \mathbf{P}_k^{uu} = \mathbf{L}_k(\hat{\Gamma}) + \mathbf{S}_k(\hat{\Gamma})$$

where $\mathbf{L}_k(\widehat{\Gamma}) = \mathbf{P}_k^{uu} \widehat{\Gamma} \mathbf{C}_k^{uu} + \mathbf{C}_k^{uu} \widehat{\Gamma} \mathbf{P}_k^{uu}$ and

$$\|\mathbf{S}_k(\widehat{\Gamma})\| \leq 14 \left(\frac{\|\widehat{\Gamma}\|}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right)^2.$$

Theorem 5.2. *Let $\delta(m_1, m_2) := \mu_1 \sigma m_1^{1/2} + \sigma^2 (m_1 m_2)^{1/2}$ and suppose that $\bar{g}_k(\mathbf{A}\mathbf{A}^\top) \geq D_1 \delta(m_1, m_2)$ for a large enough constant $D_1 > 0$ and $m_2 e^{-m_1/2} \leq 1$. Then there exists an absolute constant $D_2 > 0$ such that*

$$\|\mathbb{E} \widehat{\mathbf{P}}_k^{uu} - \mathbf{P}_k^{uu} - \mathbf{P}_k^{uu} (\mathbb{E} \widehat{\mathbf{P}}_k^{uu} - \mathbf{P}_k^{uu}) \mathbf{P}_k^{uu}\| \leq D_2 \nu_k \frac{\sigma^2 m_1 + \sigma^2 m_2^{1/2} + \sigma \mu_1}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \left(\frac{\delta(m_1, m_2)}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right).$$

Proof of Theorem 3.1. Combining Theorem 5.1 and Theorem 5.2, we conclude that for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{m_1}$ with probability at least $1 - e^{-t}$ for all $\log 8 \leq t \leq m_1$,

$$\begin{aligned} & |\langle \widehat{\mathbf{P}}_k^{uu} \mathbf{x}, \mathbf{y} \rangle - \langle \mathbf{P}_k^{uu} \mathbf{x}, \mathbf{y} \rangle - \langle \mathbf{P}_k^{uu} (\mathbb{E} \widehat{\mathbf{P}}_k^{uu} - \mathbf{P}_k^{uu}) \mathbf{P}_k^{uu} \mathbf{x}, \mathbf{y} \rangle| \\ & \lesssim \left(t^{1/2} \frac{\sigma \mu_1 + \sigma^2 m_2^{1/2}}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} + \frac{\sigma^2 m_1 \delta(m_1, m_2)}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \right) \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} \end{aligned}$$

where we used the fact $\frac{\delta(m_1, m_2)}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \leq 1$ and $\nu_k = 1$. Since $\nu_k = 1$ such that $\mathbf{P}_k^{uu} = \mathbf{u}_k \otimes \mathbf{u}_k$ and $\widehat{\mathbf{P}}_k^{uu} = \widehat{\mathbf{u}}_k \otimes \widehat{\mathbf{u}}_k$, we can write

$$\mathbf{P}_k^{uu} (\mathbb{E} \widehat{\mathbf{P}}_k^{uu} - \mathbf{P}_k^{uu}) \mathbf{P}_k^{uu} = b_k \mathbf{P}_k^{uu}$$

where

$$b_k = \mathbb{E} \langle \widehat{\mathbf{u}}_k, \mathbf{u}_k \rangle^2 - 1 \in [-1, 0].$$

Moreover, a simple fact is $b_k \leq \mathbb{E} \|\widehat{\mathbf{P}}_k^{uu} - \mathbf{P}_k^{uu}\| \lesssim \frac{\delta(m_1, m_2)}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)}$ by Wedin's $\sin \Theta$ theorem ([64]). If $\bar{g}_k(\mathbf{A}\mathbf{A}^\top) \geq D \delta(m_1, m_2)$ for a large enough constant $D > 0$, we can ensure $b_k \in [-1/2, 0]$. Then, with probability at least $1 - e^{-t}$,

$$|\langle (\widehat{\mathbf{P}}_k^{uu} - (1 + b_k) \mathbf{P}_k^{uu}) \mathbf{x}, \mathbf{y} \rangle| \lesssim \left(t^{1/2} \frac{\sigma \mu_1 + \sigma^2 m_2^{1/2}}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} + \frac{\sigma^2 m_1 \delta(m_1, m_2)}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \right) \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}.$$

By choosing $\mathbf{x} = \mathbf{y} = \mathbf{u}_k$, we obtain for all $\log 8 \leq t \leq m_1$,

$$\mathbb{P}\left(\left|\langle \hat{\mathbf{u}}_k, \mathbf{u}_k \rangle^2 - (1 + b_k)\right| \gtrsim t^{1/2} \frac{\sigma\mu_1 + \sigma^2 m_2^{1/2}}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} + \frac{\sigma^2 m_1 \delta(m_1, m_2)}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)}\right) \leq e^{-t}.$$

Denote this event by \mathcal{E}_1 . Observe that if the constant $D > 0$ is large enough and $m_1 \ll m_2$, we conclude that on event \mathcal{E}_1 , $\langle \hat{\mathbf{u}}_k, \mathbf{u}_k \rangle^2 \geq \frac{1}{4}$. Then, on event \mathcal{E}_1 ,

$$\begin{aligned} & \left| \langle \hat{\mathbf{u}}_k, \mathbf{x} \rangle - \sqrt{1 + b_k} \langle \mathbf{u}_k, \mathbf{x} \rangle \right| \\ & \leq \left| \frac{1 + b_k}{\langle \hat{\mathbf{u}}_k, \mathbf{u}_k \rangle} - \sqrt{1 + b_k} \right| |\langle \mathbf{u}_k, \mathbf{x} \rangle| \\ & \quad + \frac{1}{|\langle \hat{\mathbf{u}}_k, \mathbf{u}_k \rangle|} \left| \langle \hat{\mathbf{u}}_k, \mathbf{u}_k \rangle \langle \hat{\mathbf{u}}_k, \mathbf{x} \rangle - (1 + b_k) \langle \mathbf{u}_k, \mathbf{x} \rangle \right| \\ & = \frac{\sqrt{1 + b_k} |1 + b_k - \langle \hat{\mathbf{u}}_k, \mathbf{u}_k \rangle^2| |\langle \mathbf{u}_k, \mathbf{x} \rangle|}{|\langle \hat{\mathbf{u}}_k, \mathbf{u}_k \rangle| (\sqrt{1 + b_k} + \langle \hat{\mathbf{u}}_k, \mathbf{u}_k \rangle)} + \frac{1}{|\langle \hat{\mathbf{u}}_k, \mathbf{u}_k \rangle|} \left| \langle (\hat{\mathbf{P}}_k^{uu} - (1 + b_k) \mathbf{P}_k^{uu}) \mathbf{u}_k, \mathbf{x} \rangle \right| \\ & \lesssim t^{1/2} \frac{\sigma\mu_1 + \sigma^2 m_2^{1/2}}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \|\mathbf{x}\|_{\ell_2} + \frac{\sigma^2 m_1}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \left(\frac{\delta(m_1, m_2)}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right) \|\mathbf{x}\|_{\ell_2}, \end{aligned}$$

which concludes the proof after replacing \mathbf{A} with $\mathcal{M}_1(\mathbf{A})$ and μ_1 with $\|\mathcal{M}_1(\mathbf{A})\|$. \square

Proof of Theorem 5.2. Recall the representation formula of $\hat{\mathbf{P}}_k^{uu}$ in Lemma 4 that

$$\mathbb{E} \hat{\mathbf{P}}_k^{uu} = \mathbf{P}_k^{uu} + \mathbb{E} \mathbf{S}_k(\hat{\mathbf{\Gamma}})$$

where $\hat{\mathbf{\Gamma}} := \mathbf{A}\mathbf{Z}^\top + \mathbf{Z}\mathbf{A}^\top + \mathbf{Z}\mathbf{Z}^\top - m_2\sigma^2\mathbf{I}_{m_1}$. To this end, define

$$\tilde{\mathbf{\Gamma}} := \hat{\mathbf{\Gamma}} - (\mathbf{Z}\mathbf{P}_k^{hh}\mathbf{Z}^\top - \nu_k\sigma^2\mathbf{I}_{m_1})$$

such that we can write $\mathbb{E} \hat{\mathbf{P}}_k^{uu} = \mathbf{P}_k^{uu} + \mathbb{E} \mathbf{S}_k(\tilde{\mathbf{\Gamma}}) + (\mathbb{E} \mathbf{S}_k(\hat{\mathbf{\Gamma}}) - \mathbb{E} \mathbf{S}_k(\tilde{\mathbf{\Gamma}}))$. We derive an upper bound on $\|\mathbb{E} \mathbf{S}_k(\tilde{\mathbf{\Gamma}}) - \mathbb{E} \mathbf{S}_k(\hat{\mathbf{\Gamma}})\|$ and the proof can be found in the Appendix. Lemma 5 implies that our analysis can be proceeded by replacing $\hat{\mathbf{\Gamma}}$ with $\tilde{\mathbf{\Gamma}}$.

Lemma 5. There exists a universal constant $D_1 > 0$ such that if $m_2 e^{-m_1/2} \leq 1$, then

$$\|\mathbb{E}\mathbf{S}_k(\tilde{\Gamma}) - \mathbb{E}\mathbf{S}_k(\hat{\Gamma})\| \leq D_1 \frac{\sigma\mu_1 + \sigma^2 m_1}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \left(\frac{\delta(m_1, m_2)}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right).$$

Let $\delta_t = \mathbb{E}\|\hat{\Gamma}\| + D_1\sigma\mu_1 t^{1/2} + D_2\sigma^2 m_2^{1/2} t^{1/2}$ for $0 < t \leq m_1$ to be determined later and large enough constants $D_1, D_2 > 0$ such that $\mathbb{P}(\|\hat{\Gamma}\| \geq \delta_t) \leq e^{-t}$. We write

$$\begin{aligned} & \mathbb{E}\hat{\mathbf{P}}_k^{uu} - \mathbf{P}_k^{uu} - \mathbf{P}_k^{uu}\mathbb{E}\mathbf{S}_k(\tilde{\Gamma})\mathbf{P}_k^{uu} \\ &= \mathbb{E}\mathbf{S}_k(\hat{\Gamma}) - \mathbb{E}\mathbf{S}_k(\tilde{\Gamma}) \\ &+ \mathbb{E}\left(\mathbf{P}_k^{uu}\mathbf{S}_k(\tilde{\Gamma})(\mathbf{P}_k^{uu})^\perp + (\mathbf{P}_k^{uu})^\perp\mathbf{S}_k(\tilde{\Gamma})\mathbf{P}_k^{uu} + (\mathbf{P}_k^{uu})^\perp\mathbf{S}_k(\tilde{\Gamma})(\mathbf{P}_k^{uu})^\perp\right)\mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) \\ &+ \mathbb{E}\left(\mathbf{P}_k^{uu}\mathbf{S}_k(\tilde{\Gamma})(\mathbf{P}_k^{uu})^\perp + (\mathbf{P}_k^{uu})^\perp\mathbf{S}_k(\tilde{\Gamma})\mathbf{P}_k^{uu} + (\mathbf{P}_k^{uu})^\perp\mathbf{S}_k(\tilde{\Gamma})(\mathbf{P}_k^{uu})^\perp\right)\mathbf{1}(\|\tilde{\Gamma}\| > \delta_t). \end{aligned} \tag{5.3}$$

We prove an upper bound for $\mathbb{E}\langle \mathbf{x}, (\mathbf{P}_k^{uu})^\perp\mathbf{S}_k(\tilde{\Gamma})\mathbf{P}_k^{uu}\mathbf{y} \rangle \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t)$ for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{m_1}$. Similar to the approach in [63], under the assumption $\|\tilde{\Gamma}\| \leq \delta_t$, $\mathbf{S}_k(\tilde{\Gamma})$ is represented in the following analytic form,

$$\mathbf{S}_k(\tilde{\Gamma}) = -\frac{1}{2\pi i} \oint_{\gamma_k} \sum_{r \geq 2} (-1)^r \left(\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma} \right)^r \mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) d\eta$$

where γ_k is a circle on the complex plane with center μ_k^2 and radius $\frac{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)}{2}$, and $\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta)$ is the resolvent of the operator $\mathbf{A}\mathbf{A}^\top$ with $\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) = (\mathbf{A}\mathbf{A}^\top - \eta\mathbf{I}_{m_1})^{-1}$ which can be explicitly written as

$$\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) := (\mathbf{A}\mathbf{A}^\top - \eta\mathbf{I}_{m_1})^{-1} = \sum_s \frac{1}{\mu_s^2 - \eta} \mathbf{P}_s^{uu}.$$

We also denote

$$\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) := \mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) - \frac{1}{\mu_k^2 - \eta} \mathbf{P}_k^{uu} = \sum_{s \neq k} \frac{1}{\mu_s^2 - \eta} \mathbf{P}_s^{uu}.$$

It is easy to check that

$$\begin{aligned}
& (\mathbf{P}_k^{uu})^\perp (\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^r \mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \mathbf{P}_k^{uu} \\
&= (\mathbf{P}_k^{uu})^\perp (\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^r \frac{1}{\mu_k^2 - \eta} \mathbf{P}_k^{uu} \\
&= \left(\frac{1}{(\mu_k^2 - \eta)^2} \sum_{s=2}^r (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{s-1} (\mathbf{P}_k^{uu} \tilde{\Gamma}) (\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{r-s} \mathbf{P}_k^{uu} \right) \\
&\quad + \frac{1}{\mu_k^2 - \eta} (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^r \mathbf{P}_k^{uu},
\end{aligned}$$

where we used the formula $(a + b)^r = b^r + \sum_{s=1}^r b^{s-1} a (a + b)^{r-s}$. As a result,

$$\begin{aligned}
& (\mathbf{P}_k^{uu})^\perp \mathbf{S}_k(\tilde{\Gamma}) \mathbf{P}_k^{uu} \\
&= - \sum_{r \geq 2} (-1)^r \frac{1}{2\pi i} \oint_{\gamma_k} \left(\frac{1}{(\mu_k^2 - \eta)^2} \sum_{s=2}^r (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{s-1} (\mathbf{P}_k^{uu} \tilde{\Gamma}) (\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{r-s} \mathbf{P}_k^{uu} \right. \\
&\quad \left. + \frac{1}{\mu_k^2 - \eta} (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^r \mathbf{P}_k^{uu} \right) d\eta. \tag{5.4}
\end{aligned}$$

For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{m_1}$, we shall derive an upper bound for

$$\mathbb{E} \left\langle \mathbf{x}, (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{s-1} (\mathbf{P}_k^{uu} \tilde{\Gamma}) (\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{r-s} \mathbf{P}_k^{uu} \mathbf{y} \right\rangle \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t), \quad s = 2, \dots, r.$$

Recall that $\text{rank}(\mathbf{P}_k^{uu}) = \nu_k$ and $\mathbf{P}_k^{uu} = \sum_{j \in \Delta_k} \mathbf{u}_j \otimes \mathbf{u}_j$. Then,

$$\begin{aligned}
& \left\langle \mathbf{x}, (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{s-1} (\mathbf{P}_k^{uu} \tilde{\Gamma}) (\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{r-s} \mathbf{P}_k^{uu} \mathbf{y} \right\rangle \\
&= \sum_{j \in \Delta_k} \left\langle \mathbf{x}, (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{s-1} (\mathbf{u}_j \otimes \mathbf{u}_j \tilde{\Gamma}) (\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{r-s} \mathbf{P}_k^{uu} \mathbf{y} \right\rangle \\
&= \sum_{j \in \Delta_k} \langle \tilde{\Gamma} (\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{r-s} \mathbf{P}_k^{uu} \mathbf{y}, \mathbf{u}_j \rangle \langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{s-2} \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma} \mathbf{u}_j, \mathbf{x} \rangle.
\end{aligned}$$

Observe that

$$\left| \langle \tilde{\Gamma} (\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{r-s} \mathbf{P}_k^{uu} \mathbf{y}, \mathbf{u}_j \rangle \right| \leq \|\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta)\|^{r-s} \|\tilde{\Gamma}\|^{r-s+1} \|\mathbf{y}\|_{\ell_2}$$

$$\leq \left(\frac{2}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right)^{(r-s)} \|\tilde{\Gamma}\|^{r-s+1} \|\mathbf{y}\|_{\ell_2}.$$

Therefore,

$$\begin{aligned} & \mathbb{E} \left\langle \mathbf{x}, (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{s-1} (\mathbf{P}_k^{uu} \tilde{\Gamma}) (\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{r-s} \mathbf{P}_k^{uu} \mathbf{y} \right\rangle \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) \\ &= \sum_{j \in \Delta_k} \mathbb{E} \left\langle \tilde{\Gamma} (\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{r-s} \mathbf{P}_k^{uu} \mathbf{y}, \mathbf{u}_j \right\rangle \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{s-1} \mathbf{u}_j, \mathbf{x} \right\rangle \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) \\ &\leq \sum_{j \in \Delta_k} \mathbb{E}^{1/2} \left| \left\langle \tilde{\Gamma} (\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{r-s} \mathbf{P}_k^{uu} \mathbf{y}, \mathbf{u}_j \right\rangle \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) \right|^2 \\ &\quad \times \mathbb{E}^{1/2} \left| \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{s-1} \mathbf{u}_j, \mathbf{x} \right\rangle \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) \right|^2 \\ &\leq \left(\frac{2\delta_t}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right)^{r-s} \delta_t \|\mathbf{y}\|_{\ell_2} \sum_{j \in \Delta_k} \mathbb{E}^{1/2} \left| \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{s-2} \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma} \mathbf{u}_j, \mathbf{x} \right\rangle \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) \right|^2. \end{aligned} \tag{5.5}$$

It then remains to bound, for each $j \in \Delta_k$,

$$\mathbb{E}^{1/2} \left| \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{s-2} \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma} \mathbf{u}_j, \mathbf{x} \right\rangle \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) \right|^2.$$

Recall that we can write

$$\tilde{\Gamma} = \mathbf{A}\mathbf{Z}^\top + \mathbf{Z}\mathbf{A}^\top + \mathbf{Z} \sum_{k' \neq k} \mathbf{P}_{k'}^{hh} \mathbf{Z}^\top - \sigma^2(m_2 - \nu_k) \mathbf{I}_{m_1}$$

and correspondingly

$$\tilde{\Gamma} \mathbf{u}_j = \mathbf{A}\mathbf{Z}^\top \mathbf{u}_j + \mathbf{Z}\mathbf{A}^\top \mathbf{u}_j + \mathbf{Z} \sum_{k' \neq k} \mathbf{P}_{k'}^{hh} \mathbf{Z}^\top \mathbf{u}_j - \sigma^2(m_2 - \nu_k) \mathbf{u}_j.$$

We write

$$\begin{aligned} & \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{s-2} \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma} \mathbf{u}_j, \mathbf{x} \right\rangle \\ &= \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma})^{s-2} \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \mathbf{Z}\mathbf{A}^\top \mathbf{u}_j, \mathbf{x} \right\rangle \end{aligned} \tag{5.6}$$

$$+ \left\langle \left(\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma} \right)^{s-2} \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \mathbf{A} \mathbf{Z}^\top \mathbf{u}_j, \mathbf{x} \right\rangle \quad (5.7)$$

$$+ \left\langle \left(\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma} \right)^{s-2} \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \left(\mathbf{Z} \sum_{k' \neq k} \mathbf{P}_{k'}^{hh} \mathbf{Z}^\top \mathbf{u}_j - \sigma^2(m_2 - \nu_k) \mathbf{u}_j \right), \mathbf{x} \right\rangle. \quad (5.8)$$

The upper bounds of (5.6), (5.7), and (5.8) shall be obtained separately via different representations.

Bound of $\mathbb{E}^{1/2} \left| \left\langle \left(\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma} \right)^{s-2} \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \mathbf{Z} \mathbf{A}^\top \mathbf{u}_j, \mathbf{x} \right\rangle \right|^2 \mathbf{1} \left(\|\tilde{\Gamma}\| \leq \delta_t \right)$. Observe that $\mathbf{A}^\top \mathbf{u}_j = \mu_k \mathbf{h}_j \in \mathbb{R}^{m_2}$ for $j \in \Delta_k$ such that

$$\mathbf{Z} \mathbf{A}^\top \mathbf{u}_j = \mu_k \mathbf{Z} \mathbf{h}_j = \mu_k \sum_{i=1}^{m_1} \langle \mathbf{z}_i, \mathbf{h}_j \rangle \mathbf{e}_i$$

where $\{\mathbf{e}_1, \dots, \mathbf{e}_{m_1}\}$ denote the canonical basis vectors in \mathbb{R}^{m_1} and $\{\mathbf{z}_1^\top, \dots, \mathbf{z}_{m_1}^\top\}$ denote the rows of \mathbf{Z} . Therefore,

$$\begin{aligned} & \left\langle \left(\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma} \right)^{s-2} \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \mathbf{Z} \mathbf{A}^\top \mathbf{u}_j, \mathbf{x} \right\rangle \\ &= \mu_k \sum_{i=1}^{m_1} \langle \mathbf{z}_i, \mathbf{h}_j \rangle \left\langle \left(\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma} \right)^{s-2} \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \mathbf{e}_i, \mathbf{x} \right\rangle. \end{aligned}$$

It is clear that $\langle \mathbf{z}_i, \mathbf{h}_j \rangle, i = 1, \dots, m_1$ are i.i.d. and $\langle \mathbf{z}_i, \mathbf{h}_j \rangle \sim \mathcal{N}(0, \sigma^2)$. Recall that $\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) = \sum_{k' \neq k} \frac{\mathbf{P}_{k'}^{uu}}{\mu_{k'}^2 - \eta}$, implying that $\left(\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\Gamma} \right)^{s-2} \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)$ can be viewed as a linear combination of operators

$$(\mathbf{P}_{t_1}^{uu} \tilde{\Gamma} \mathbf{P}_{t_2}^{uu}) (\mathbf{P}_{t_2}^{uu} \tilde{\Gamma} \mathbf{P}_{t_3}^{uu}) \dots (\mathbf{P}_{t_{s-2}}^{uu} \tilde{\Gamma} \mathbf{P}_{t_{s-1}}^{uu})$$

where $t_1, \dots, t_{s-1} \neq k$. For each $\mathbf{P}_{t_1}^{uu} \tilde{\Gamma} \mathbf{P}_{t_2}^{uu}$, we have

$$\mathbf{P}_{t_1}^{uu} \tilde{\Gamma} \mathbf{P}_{t_2}^{uu} = \mathbf{P}_{t_1}^{uu} \mathbf{A} \mathbf{Z}^\top \mathbf{P}_{t_2}^{uu} + \mathbf{P}_{t_1}^{uu} \mathbf{Z} \mathbf{A}^\top \mathbf{P}_{t_2}^{uu} + \mathbf{P}_{t_1}^{uu} \left(\mathbf{Z} \sum_{k' \neq k} \mathbf{P}_{k'}^{hh} \mathbf{Z}^\top \right) \mathbf{P}_{t_2}^{uu} - \sigma^2(m_2 - \nu_k) \mathbf{P}_{t_1}^{uu} \mathbf{P}_{t_2}^{uu}.$$

Clearly, $\mathbf{P}_{t_1}^{uu} \mathbf{A} \mathbf{Z}^\top$ is a function of random vectors $\mathbf{P}_{t_1}^{uu} \mathbf{A} \mathbf{z}_i, i = 1, \dots, m_1$; $\mathbf{Z} \mathbf{A}^\top \mathbf{P}_{t_2}^{uu}$ is a function of random vectors $\mathbf{P}_{t_2}^{uu} \mathbf{A} \mathbf{z}_i, i = 1, \dots, m_1$; $\mathbf{Z} \sum_{k' \neq k} \mathbf{P}_{k'}^{hh} \mathbf{Z}^\top = \mathbf{Z} \sum_{k' \neq k} (\mathbf{P}_{k'}^{hh})^2 \mathbf{Z}^\top$ is a function of random vectors $\mathbf{P}_{k'}^{hh} \mathbf{z}_i, i = 1, \dots, m_1$. The following facts are obvious

$$\mathbb{E} \langle \mathbf{z}_i, \mathbf{h}_j \rangle \mathbf{P}_{t_1}^{uu} \mathbf{A} \mathbf{z}_i = \mathbf{P}_{t_1}^{uu} \mathbf{A} (\mathbb{E} \mathbf{z}_i \otimes \mathbf{z}_i) \mathbf{h}_j = \sigma^2 \mathbf{P}_{t_1}^{uu} \mathbf{A} \mathbf{h}_j = \sigma^2 \mu_k \mathbf{P}_{t_1}^{uu} \mathbf{u}_j = \mathbf{0}, \quad \forall t_1 \neq k$$

and

$$\mathbb{E} \langle \mathbf{z}_i, \mathbf{h}_j \rangle \mathbf{P}_{k'}^{hh} \mathbf{z}_i = \mathbf{P}_{k'}^{hh} (\mathbb{E} \mathbf{z}_i \otimes \mathbf{z}_i) \mathbf{h}_j = \sigma^2 \mathbf{P}_{k'}^{hh} \mathbf{h}_j = \mathbf{0}, \quad \forall k' \neq k.$$

Since $\{\langle \mathbf{z}_i, \mathbf{h}_j \rangle, i = 1, \dots, m_1\}$ are Gaussian random variables and $\{\mathbf{P}_{t_1}^{uu} \mathbf{A} \mathbf{z}_i, \mathbf{P}_{k'}^{hh} \mathbf{z}_i, i = 1, \dots, m_1\}$ are (complex) Gaussian random vectors, uncorrelations indicate that $\{\langle \mathbf{z}_i, \mathbf{h}_j \rangle : i = 1, \dots, m_1\}$ are independent with $\{\mathbf{P}_{t_1}^{uu} \mathbf{A} \mathbf{z}_i, \mathbf{P}_{k'}^{hh} \mathbf{z}_i : t_1 \neq k, k' \neq k, i = 1, \dots, m_1\}$. We conclude that $\{\langle \mathbf{z}_i, \mathbf{h}_j \rangle : i = 1, \dots, m_1\}$ are independent with

$$\{\langle (\tilde{\mathbf{R}}_{\mathbf{A} \mathbf{A}^\top}(\eta) \tilde{\mathbf{\Gamma}})^{s-2} \tilde{\mathbf{R}}_{\mathbf{A} \mathbf{A}^\top}(\eta) \mathbf{e}_i, \mathbf{x} \rangle, i = 1, \dots, m_1\}.$$

To this end, define the complex random variables

$$\omega_i(\mathbf{x}) = \langle (\tilde{\mathbf{R}}_{\mathbf{A} \mathbf{A}^\top}(\eta) \tilde{\mathbf{\Gamma}})^{s-2} \tilde{\mathbf{R}}_{\mathbf{A} \mathbf{A}^\top}(\eta) \mathbf{e}_i, \mathbf{x} \rangle = \omega_i^{(1)}(\mathbf{x}) + \omega_i^{(2)}(\mathbf{x}) \text{Im} \in \mathbb{C}, \quad i = 1, \dots, m_1$$

where Im denotes the imaginary number. Then,

$$\begin{aligned} & \langle (\tilde{\mathbf{R}}_{\mathbf{A} \mathbf{A}^\top}(\eta) \tilde{\mathbf{\Gamma}})^{s-2} \tilde{\mathbf{R}}_{\mathbf{A} \mathbf{A}^\top}(\eta) \mathbf{Z} \mathbf{A}^\top \mathbf{u}_j, \mathbf{x} \rangle \\ &= \mu_k \sum_{i=1}^{m_1} \langle \mathbf{z}_i, \mathbf{h}_j \rangle \omega_i^{(1)}(\mathbf{x}) + \left(\mu_k \sum_{i=1}^{m_1} \langle \mathbf{z}_i, \mathbf{h}_j \rangle \omega_i^{(2)}(\mathbf{x}) \right) \text{Im} \\ &=: \kappa_1(\mathbf{x}) + \kappa_2(\mathbf{x}) \text{Im} \in \mathbb{C}. \end{aligned}$$

Conditioned on $\{\mathbf{P}_{t_1}^{uu} \mathbf{A} \mathbf{z}_i, \mathbf{P}_{k'}^{hh} \mathbf{z}_i : t_1 \neq k, k' \neq k, i = 1, \dots, m_1\}$, we get

$$\mathbb{E} \kappa_1^2(\mathbf{x}) = \mu_k^2 \sigma^2 \sum_{i=1}^{m_1} \left(\omega_i^{(1)}(\mathbf{x}) \right)^2$$

and

$$\mathbb{E} \kappa_1(\mathbf{x}) \kappa_2(\mathbf{x}) = \mu_k^2 \sigma^2 \sum_{i=1}^{m_1} \omega_i^{(1)}(\mathbf{x}) \omega_i^{(2)}(\mathbf{x})$$

implying that the centered Gaussian random vector $(\kappa_1(\mathbf{x}), \kappa_2(\mathbf{x}))$ has covariance matrix:

$$\left(\mu_k^2 \sigma^2 \sum_{i=1}^{m_1} \omega_i^{(k_1)}(\mathbf{x}) \omega_i^{(k_2)}(\mathbf{x}) \right)_{k_1, k_2=1,2}.$$

Finally,

$$\begin{aligned} & \mathbb{E}^{1/2} \left| \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\mathbf{\Gamma}})^{s-2} \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \mathbf{Z} \mathbf{A}^\top \mathbf{u}_j, \mathbf{x} \right\rangle \right|^2 \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| \leq \delta_t) \\ &= \mathbb{E}^{1/2} (\kappa_1^2(\mathbf{x}) + \kappa_2^2(\mathbf{x})) \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| \leq \delta_t) \\ &= \sigma \mu_k \mathbb{E}^{1/2} \left(\sum_{i=1}^{m_1} (\omega_i^{(1)}(\mathbf{x}))^2 + (\omega_i^{(2)}(\mathbf{x}))^2 \right) \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| \leq \delta_t) \\ &= \sigma \mu_k \mathbb{E}^{1/2} \sum_{i=1}^{m_1} |\omega_i(\mathbf{x})|^2 \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| \leq \delta_t). \end{aligned}$$

Moreover,

$$\begin{aligned} \sum_{i=1}^{m_1} |\omega_i(\mathbf{x})|^2 &= \sum_{i=1}^{m_1} \left| \left\langle \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\mathbf{\Gamma}})^{s-2} \mathbf{x}, \mathbf{e}_j \right\rangle \right|^2 \leq \|\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\mathbf{\Gamma}})^{s-2} \mathbf{x}\|_{\ell_2}^2 \\ &\leq \|\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\|^{2(s-1)} \|\tilde{\mathbf{\Gamma}}\|^{2(s-2)} \|\mathbf{x}\|_{\ell_2}^2 \leq \left(\frac{2}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right)^{2(s-1)} \|\tilde{\mathbf{\Gamma}}\|^{2(s-2)} \|\mathbf{x}\|_{\ell_2}^2. \end{aligned}$$

As a result,

$$\begin{aligned} & \mathbb{E}^{1/2} \left| \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \tilde{\mathbf{\Gamma}})^{s-2} \tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta) \mathbf{Z} \mathbf{A}^\top \mathbf{u}_j, \mathbf{x} \right\rangle \right|^2 \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| \leq \delta_t) \\ &\leq \sigma \mu_k \mathbb{E}^{1/2} \left(\frac{2}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right)^{2(s-1)} \|\tilde{\mathbf{\Gamma}}\|^{2(s-2)} \|\mathbf{x}\|_{\ell_2}^2 \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| \leq \delta_t) \end{aligned}$$

$$\leq \frac{\sigma\mu_k}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \left(\frac{2\delta_t}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right)^{s-2} \|\mathbf{x}\|_{\ell_2}.$$

Bound of $\mathbb{E}^{1/2}$ $\left| \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\Gamma})^{s-2}\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\mathbf{A}\mathbf{Z}^\top\mathbf{u}_j, \mathbf{x} \right\rangle \right|^2 \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t)$. With a little abuse on the notations, we denote by $\mathbf{z}_1, \dots, \mathbf{z}_{m_2} \in \mathbb{R}^{m_1}$ the corresponding columns of \mathbf{Z} in this paragraph. Then,

$$\left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\Gamma})^{s-2}\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\mathbf{A}\mathbf{Z}^\top\mathbf{u}_j, \mathbf{x} \right\rangle = \sum_{i=1}^{m_2} \langle \mathbf{z}_i, \mathbf{u}_j \rangle \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\Gamma})^{s-2}\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\mathbf{A}\mathbf{e}_i, \mathbf{x} \right\rangle.$$

Similarly, $(\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\Gamma})^{s-2}\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)$ can be represented as linear combination of operators

$$(\mathbf{P}_{t_1}^{uu}\tilde{\Gamma}\mathbf{P}_{t_2}^{uu})(\mathbf{P}_{t_2}^{uu}\tilde{\Gamma}\mathbf{P}_{t_3}^{uu}) \dots (\mathbf{P}_{t_{s-2}}^{uu}\tilde{\Gamma}\mathbf{P}_{t_{s-1}}^{uu}), \quad t_1, \dots, t_{s-1} \neq k.$$

To this end, we write

$$\mathbf{P}_{t_1}^{uu}\tilde{\Gamma}\mathbf{P}_{t_2}^{uu} = \mathbf{P}_{t_1}^{uu}\mathbf{A}\mathbf{Z}^\top\mathbf{P}_{t_2}^{uu} + \mathbf{P}_{t_1}^{uu}\mathbf{Z}\mathbf{A}^\top\mathbf{P}_{t_2}^{uu} + \mathbf{P}_{t_1}^{uu}\left(\mathbf{Z}\sum_{k' \neq k} \mathbf{P}_{k'}^{hh}\mathbf{Z}^\top\right)\mathbf{P}_{t_2}^{uu} - \sigma^2(m_2 - \nu_k)\mathbf{P}_{t_1}^{uu}\mathbf{P}_{t_2}^{uu}.$$

Observe that $\mathbf{P}_{t_1}^{uu}\mathbf{A}\mathbf{Z}^\top\mathbf{P}_{t_2}^{uu}$, $\mathbf{P}_{t_1}^{uu}\mathbf{Z}\mathbf{A}^\top\mathbf{P}_{t_2}^{uu}$ and $\mathbf{P}_{t_1}^{uu}(\mathbf{Z}\sum_{k' \neq k} \mathbf{P}_{k'}^{hh}\mathbf{Z}^\top)\mathbf{P}_{t_2}^{uu}$ are functions of random vectors $\{\mathbf{P}_{t_1}^{uu}\mathbf{z}_i, \mathbf{P}_{t_2}^{uu}\mathbf{z}_i : t_1, t_2 \neq k, i = 1, \dots, m_2\}$. Moreover,

$$\mathbb{E}\langle \mathbf{z}_i, \mathbf{u}_j \rangle \mathbf{P}_{t_1}^{uu}\mathbf{z}_i = \mathbf{P}_{t_1}^{uu}(\mathbb{E}\mathbf{z}_i \otimes \mathbf{z}_i)\mathbf{u}_j = \sigma^2\mathbf{P}_{t_1}^{uu}\mathbf{u}_j = \mathbf{0}, \quad \forall t_1 \neq k$$

which implies that $\{\langle \mathbf{z}_i, \mathbf{u}_j \rangle : i = 1, \dots, m_2\}$ and $\left\{ \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\Gamma})^{s-2}\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\mathbf{A}\mathbf{e}_i, \mathbf{x} \right\rangle : i = 1, \dots, m_2 \right\}$ are independent. Following an identical analysis as above, we get

$$\mathbb{E}^{1/2} \left| \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\Gamma})^{s-2}\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\mathbf{A}\mathbf{Z}^\top\mathbf{u}_j, \mathbf{x} \right\rangle \right|^2 \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) \leq \frac{\sigma\mu_1}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \left(\frac{2\delta_t}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right)^{s-2} \|\mathbf{x}\|_{\ell_2}.$$

Bound of $\mathbb{E}^{1/2}$ $\left| \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\Gamma})^{s-2}\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)(\mathbf{Z}\sum_{k' \neq k} \mathbf{P}_{k'}^{hh}\mathbf{Z}^\top)\mathbf{u}_j, \mathbf{x} \right\rangle \right|^2 \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t)$. Note that we used the fact $\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\mathbf{u}_j = \mathbf{0}$ in (5.8). Again, let $\{\mathbf{z}_1, \dots, \mathbf{z}_{m_2}\} \subset \mathbb{R}^{m_1}$ denote the

corresponding columns of \mathbf{Z} . We write

$$\begin{aligned} & \langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\mathbf{\Gamma}})^{s-2}\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)(\mathbf{Z}\sum_{k'\neq k}\mathbf{P}_{k'}^{hh}\mathbf{Z}^\top)\mathbf{u}_j, \mathbf{x} \rangle \\ &= \sum_{i=1}^{m_2} \langle \mathbf{z}_i, \mathbf{u}_j \rangle \langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\mathbf{\Gamma}})^{s-2}\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\mathbf{Z}(\sum_{k'\neq k}\mathbf{P}_{k'}^{hh})\mathbf{e}_i, \mathbf{x} \rangle. \end{aligned}$$

In a similar fashion, we show that $(\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\mathbf{\Gamma}})^{s-2}\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\mathbf{Z}$ is a function of random vectors $\{\mathbf{P}_t^{uu}\mathbf{z}_i : t \neq k, i = 1, \dots, m_2\}$ which are independent with $\{\langle \mathbf{z}_i, \mathbf{u}_j \rangle : i = 1, \dots, m_2\}$. Then,

$$\begin{aligned} & \mathbb{E}^{1/2} \left| \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\mathbf{\Gamma}})^{s-2}\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)(\mathbf{Z}\sum_{k'\neq k}\mathbf{P}_{k'}^{hh}\mathbf{Z}^\top)\mathbf{u}_j, \mathbf{x} \right\rangle \right|^2 \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| \leq \delta) \\ & \leq \mathbb{E}^{1/2} \sigma^2 \|\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\|^{2(s-1)} \|\tilde{\mathbf{\Gamma}}\|^{2(s-2)} \|\mathbf{Z}\sum_{k'\neq k}\mathbf{P}_{k'}^{hh}\|^2 \|\mathbf{x}\|_{\ell_2}^2 \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| \leq \delta_t) \\ & \lesssim \frac{\sigma^2 m_2^{1/2}}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \left(\frac{\delta_t}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right)^{s-2} \|\mathbf{x}\|_{\ell_2}. \end{aligned}$$

where we used the fact $\mathbb{E}^{1/2} \|(\sum_{k'\neq k}\mathbf{P}_{k'}^{hh})\mathbf{Z}^\top\|^2 \lesssim \sigma m_2^{1/2}$ from Lemma 3.

Finalize the proof of Theorem. Combining the above bounds into (5.7), (5.6) and (5.8), we conclude that

$$\begin{aligned} & \mathbb{E}^{1/2} \left| \left\langle (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\mathbf{\Gamma}})^{s-2}\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\mathbf{\Gamma}}\mathbf{u}_j, \mathbf{x} \right\rangle \right|^2 \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| \leq \delta_t) \\ & \lesssim \frac{\sigma^2 m_2^{1/2} + \sigma \mu_1}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \left(\frac{2\delta_t}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right)^{s-2} \|\mathbf{x}\|_{\ell_2}. \end{aligned}$$

Continue from (5.5) and we end up with

$$\begin{aligned} & \mathbb{E} \langle \mathbf{x}, (\tilde{\mathbf{R}}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\mathbf{\Gamma}})^{s-1}(\mathbf{P}_k^{uu}\tilde{\mathbf{\Gamma}})(\mathbf{R}_{\mathbf{A}\mathbf{A}^\top}(\eta)\tilde{\mathbf{\Gamma}})^{r-s}\mathbf{P}_k^{uu}\mathbf{y} \rangle \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| \leq \delta_t) \\ & \lesssim \nu_k \delta_t \frac{\sigma^2 m_2^{1/2} + \sigma \mu_1}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \left(\frac{2\delta_t}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right)^{r-2} \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}. \end{aligned}$$

Plug the bounds into (5.4),

$$\begin{aligned}
& |\mathbb{E}\langle (\mathbf{P}_k^{uu})^\perp \mathbf{S}_k(\tilde{\Gamma}) \mathbf{P}_k^{uu} \mathbf{y}, \mathbf{x} \rangle \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) | \\
& \lesssim \sum_{r \geq 2} \frac{\pi \bar{g}_k(\mathbf{A} \mathbf{A}^\top)}{2\pi} \left(\frac{2}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \right)^2 (r-1) \nu_k \delta_t \frac{\sigma^2 m_2^{1/2} + \sigma \mu_1}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \left(\frac{2\delta_t}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \right)^{r-2} \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} \\
& \leq D_1 \nu_k \frac{\sigma^2 m_2^{1/2} + \sigma \mu_1}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} \sum_{r \geq 2} (r-1) \left(\frac{2\delta_t}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \right)^{r-1}
\end{aligned}$$

where we used the fact $\oint_{\gamma_k} (\tilde{\mathbf{R}}_{\mathbf{A} \mathbf{A}^\top}(\eta) \tilde{\Gamma})^r \mathbf{P}_k^{uu} d\eta = \mathbf{0}$. By the inequality $\sum_{r \geq 1} r q^r = \frac{q}{(1-q)^2}$, $\forall q < 1$ and the fact $D_1 \delta_t \leq \bar{g}_k(\mathbf{A} \mathbf{A}^\top)$ for some large constant $D_1 > 0$ and $t \leq m_1$,

we conclude with

$$\begin{aligned}
& |\mathbb{E}\langle (\mathbf{P}_k^{uu})^\perp \mathbf{S}_k(\tilde{\Gamma}) \mathbf{P}_k^{uu} \mathbf{y}, \mathbf{x} \rangle \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) | \\
& \lesssim \nu_k \frac{\sigma^2 m_2^{1/2} + \sigma \mu_1}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \left(\frac{2\delta_t}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \right) \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^{m_1}
\end{aligned}$$

implying that

$$\left\| \mathbb{E}(\mathbf{P}_k^{uu})^\perp \mathbf{S}_k(\tilde{\Gamma}) \mathbf{P}_k^{uu} \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) \right\| \lesssim \nu_k \frac{\sigma^2 m_2^{1/2} + \sigma \mu_1}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \left(\frac{2\delta_t}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \right).$$

The same bound holds for

$$\left\| \mathbb{E} \mathbf{P}_k^{uu} \mathbf{S}_k(\tilde{\Gamma}) (\mathbf{P}_k^{uu})^\perp \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) \right\| \quad \text{and} \quad \left\| \mathbb{E} (\mathbf{P}_k^{uu})^\perp \mathbf{S}_k(\tilde{\Gamma}) (\mathbf{P}_k^{uu})^\perp \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) \right\|,$$

following the same arguments. As a result,

$$\begin{aligned}
& \left\| \mathbb{E} \left((\mathbf{P}_k^{uu})^\perp \mathbf{S}_k(\tilde{\Gamma}) \mathbf{P}_k^{uu} + \mathbf{P}_k^{uu} \mathbf{S}_k(\tilde{\Gamma}) (\mathbf{P}_k^{uu})^\perp + (\mathbf{P}_k^{uu})^\perp \mathbf{S}_k(\tilde{\Gamma}) (\mathbf{P}_k^{uu})^\perp \right) \mathbf{1}(\|\tilde{\Gamma}\| \leq \delta_t) \right\| \\
& \lesssim \nu_k \frac{\sigma^2 m_2^{1/2} + \sigma \mu_1}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \left(\frac{2\delta_t}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \right). \tag{5.9}
\end{aligned}$$

By choosing $t = m_1$ such that $\mathbb{P}(\|\tilde{\Gamma}\| \geq \delta_{m_1}) \leq e^{-m_1/2}$, we get

$$\begin{aligned}
& \left\| \mathbb{E} \left((\mathbf{P}_k^{uu})^\perp \mathbf{S}_k(\tilde{\Gamma}) \mathbf{P}_k^{uu} + \mathbf{P}_k^{uu} \mathbf{S}_k(\tilde{\Gamma}) (\mathbf{P}_k^{uu})^\perp + (\mathbf{P}_k^{uu})^\perp \mathbf{S}_k(\tilde{\Gamma}) (\mathbf{P}_k^{uu})^\perp \right) \mathbf{1}(\|\tilde{\Gamma}\| > \delta_{m_1}) \right\| \\
& \leq \mathbb{E} \left\| \left((\mathbf{P}_k^{uu})^\perp \mathbf{S}_k(\tilde{\Gamma}) \mathbf{P}_k^{uu} + \mathbf{P}_k^{uu} \mathbf{S}_k(\tilde{\Gamma}) (\mathbf{P}_k^{uu})^\perp + (\mathbf{P}_k^{uu})^\perp \mathbf{S}_k(\tilde{\Gamma}) (\mathbf{P}_k^{uu})^\perp \right) \mathbf{1}(\|\tilde{\Gamma}\| > \delta_{m_1}) \right\| \\
& \leq \mathbb{E} \|\mathbf{S}_k(\tilde{\Gamma})\| \mathbf{1}(\|\tilde{\Gamma}\| > \delta_{m_1}) \leq \mathbb{E}^{1/2} \|\mathbf{S}_k(\tilde{\Gamma})\|^2 \mathbb{P}^{1/2}(\|\tilde{\Gamma}\| > \delta_{m_1}) \\
& \lesssim \left(\frac{\delta_{m_1}}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right)^2 \mathbb{P}^{1/2}(\|\tilde{\Gamma}\| > \delta_{m_1}) \lesssim \left(\frac{\delta_{m_1}}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right)^2 e^{-m_1/2},
\end{aligned}$$

which is clearly dominated by (5.9). Substitute the above bounds into (5.3) and we get

$$\begin{aligned}
\left\| \mathbb{E} \hat{\mathbf{P}}_k^{uu} - \mathbf{P}_k^{uu} - \mathbf{P}_k^{uu} \mathbf{S}_k(\tilde{\Gamma}) \mathbf{P}_k^{uu} \right\| & \leq \|\mathbb{E} \mathbf{S}_k(\tilde{\Gamma}) - \mathbf{S}_k(\hat{\Gamma})\| + D_1 \nu_k \frac{\sigma^2 m_2^{1/2} + \sigma \mu_1}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \left(\frac{2\delta(m_1, m_2)}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right) \\
& \leq D_2 \nu_k \frac{\sigma^2 m_2^{1/2} + \sigma^2 m_1 + \sigma \mu_1}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \left(\frac{2\delta(m_1, m_2)}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right).
\end{aligned}$$

□

5.2 Proof of Theorem 3.2

The proof of Theorem 3.2 is identical to the proof of Corollary 1.5 in [63] and will be skipped here.

5.3 Proof of Theorem 3.3

It suffices to prove the upper bound of $|\tilde{A}(i, j, k) - A(i, j, k)|$ for $i \in [d_1], j \in [d_2], k \in [d_3]$.

To this end, denote by \mathbf{e}_i the i -th canonical basis vectors. Observe that

$$\begin{aligned}
\langle \tilde{\mathbf{A}} - \mathbf{A}, \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle &= \langle \mathbf{A} \times_1 \mathbf{P}_{\tilde{\mathbf{U}}} \times_2 \mathbf{P}_{\tilde{\mathbf{V}}} \times_3 \mathbf{P}_{\tilde{\mathbf{W}}} - \mathbf{A}, \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle \\
&\quad + \langle \mathbf{Z} \times_1 \mathbf{P}_{\tilde{\mathbf{U}}} \times_2 \mathbf{P}_{\tilde{\mathbf{V}}} \times_3 \mathbf{P}_{\tilde{\mathbf{W}}}, \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle.
\end{aligned}$$

Some preliminary facts shall be concluded from Theorem 3.1. By Theorem 3.2, there exists an event \mathcal{E}_2 with $\mathbb{P}(\mathcal{E}_2) \geq 1 - \frac{1}{d^2}$ on which

$$\|\mathbf{e}_i^\top (\tilde{\mathbf{U}} - \mathbf{U})\|_{\ell_2} \leq r^{1/2} \|\mathbf{e}_i^\top (\tilde{\mathbf{U}} - \mathbf{U})\|_{\ell_\infty} \lesssim \frac{\sigma \bar{\Lambda}(\mathbf{A}) r^{1/2} + \sigma^2 d r^{1/2}}{\bar{g}_{\min}^2(\mathbf{A})} \log^{1/2} d$$

and

$$\|\tilde{\mathbf{U}}^\top \mathbf{U} - \mathbf{I}_{r_1}\| \leq \|\tilde{\mathbf{U}}^\top \mathbf{U} - \mathbf{I}_{r_1}\|_{\text{F}} \lesssim r \|\tilde{\mathbf{U}}^\top \mathbf{U} - \mathbf{I}_{r_1}\|_{\ell_\infty} \lesssim \frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 d r}{\bar{g}_{\min}^2(\mathbf{A})} \log^{1/2} d.$$

The following decomposition is straightforward,

$$\begin{aligned} & \mathbf{A} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}}, \mathbf{P}_{\tilde{\mathbf{V}}}, \mathbf{P}_{\tilde{\mathbf{W}}}) - \mathbf{A} \\ &= \mathbf{A} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\mathbf{W}}) + \mathbf{A} \cdot (\mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\tilde{\mathbf{V}}} - \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\mathbf{W}}) \\ &+ \mathbf{A} \cdot (\mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\tilde{\mathbf{W}}} - \mathbf{P}_{\mathbf{W}}) + \mathbf{A} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\tilde{\mathbf{V}}} - \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\mathbf{W}}) \\ &+ \mathbf{A} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\tilde{\mathbf{W}}} - \mathbf{P}_{\mathbf{W}}) + \mathbf{A} \cdot (\mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\tilde{\mathbf{V}}} - \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\tilde{\mathbf{W}}} - \mathbf{P}_{\mathbf{W}}) \\ &+ \mathbf{A} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\tilde{\mathbf{V}}} - \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\tilde{\mathbf{W}}} - \mathbf{P}_{\mathbf{W}}) \end{aligned}$$

Recall that $\mathbf{A} = \mathbf{C} \cdot (\mathbf{U}, \mathbf{V}, \mathbf{W})$ and we get

$$\begin{aligned} & \left\langle \mathbf{A} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\mathbf{W}}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \right\rangle \\ &= \mathbf{e}_i^\top \left(\tilde{\mathbf{U}} (\tilde{\mathbf{U}}^\top \mathbf{U}) - \mathbf{U} \right) \mathcal{M}_1(\mathbf{C}) (\mathbf{V} \otimes \mathbf{W})^\top (\mathbf{e}_j \otimes \mathbf{e}_k). \end{aligned}$$

Observe that

$$\mathbf{e}_i^\top \left(\tilde{\mathbf{U}} (\tilde{\mathbf{U}}^\top \mathbf{U}) - \mathbf{U} \right) = \mathbf{e}_i^\top (\tilde{\mathbf{U}} - \mathbf{U}) (\tilde{\mathbf{U}}^\top \mathbf{U}) + \mathbf{e}_i^\top \mathbf{U} (\tilde{\mathbf{U}}^\top \mathbf{U} - \mathbf{I}_{r_1})$$

implying that on event \mathcal{E}_2 ,

$$\begin{aligned}
& \left\| \mathbf{e}_i^\top \left(\tilde{\mathbf{U}}(\tilde{\mathbf{U}}^\top \mathbf{U}) - \mathbf{U} \right) \right\|_{\ell_2} \\
& \leq \|(\tilde{\mathbf{U}} - \mathbf{U})^\top \mathbf{e}_i\|_{\ell_2} \|\tilde{\mathbf{U}}^\top \mathbf{U}\| + \|\tilde{\mathbf{U}}^\top \mathbf{U} - \mathbf{I}_{r_1}\| \|\mathbf{U}^\top \mathbf{e}_i\|_{\ell_2} \\
& \lesssim \frac{\sigma \bar{\Lambda}(\mathbf{A}) r^{1/2} + \sigma^2 d r^{1/2}}{\bar{g}_{\min}^2(\mathbf{A})} \log^{1/2} d + \|\mathbf{U}^\top \mathbf{e}_i\|_{\ell_2} \frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 d r}{\bar{g}_{\min}^2(\mathbf{A})} \log^{1/2} d \\
& \lesssim \frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 d r}{\bar{g}_{\min}^2(\mathbf{A})} \log^{1/2} d,
\end{aligned}$$

where we used the facts $\|\tilde{\mathbf{U}}^\top \mathbf{U}\| \leq \|\tilde{\mathbf{U}}\| \|\mathbf{U}\| \leq (1 + b_k)^{-1/2} = O(1)$ and

$$\|\mathbf{U}^\top \mathbf{e}_i\|_{\ell_2} = \langle \mathbf{U} \mathbf{U}^\top, \mathbf{e}_i \otimes \mathbf{e}_i \rangle^{1/2} \leq 1.$$

Therefore, on event \mathcal{E}_2 ,

$$\begin{aligned}
& \left| \langle \mathbf{A} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\mathbf{W}}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle \right| \\
& \lesssim \bar{\Lambda}(\mathbf{A}) \left(\frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 d r}{\bar{g}_{\min}^2(\mathbf{A})} \log^{1/2} d \right) \|\mathbf{V}^\top \mathbf{e}_j\|_{\ell_2} \|\mathbf{W}^\top \mathbf{e}_k\|_{\ell_2}.
\end{aligned}$$

Similar bounds hold for

$$\left| \langle \mathbf{A} \cdot (\mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\tilde{\mathbf{V}}} - \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\mathbf{W}}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle \right| \quad \text{and} \quad \left| \langle \mathbf{A} \cdot (\mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\tilde{\mathbf{W}}} - \mathbf{P}_{\mathbf{W}}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle \right|.$$

Following the same method, we can show that on event \mathcal{E}_2 ,

$$\begin{aligned}
& \left| \langle \mathbf{A} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\tilde{\mathbf{V}}} - \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\mathbf{W}}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle \right| \\
& \lesssim \bar{\Lambda}(\mathbf{A}) \left(\frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 d r}{\bar{g}_{\min}^2(\mathbf{A})} \log^{1/2} d \right)^2 \|\mathbf{W}^\top \mathbf{e}_k\|_{\ell_2}
\end{aligned}$$

and

$$\left| \langle \mathbf{A} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\tilde{\mathbf{V}}} - \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\tilde{\mathbf{W}}} - \mathbf{P}_{\mathbf{W}}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle \right|$$

$$\lesssim \bar{\Lambda}(\mathbf{A}) \left(\frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 dr}{\bar{g}_{\min}^2(\mathbf{A})} \log^{1/2} d \right)^3.$$

We conclude that on event \mathcal{E}_2 ,

$$\begin{aligned} & \left| \langle \mathbf{A} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}}, \mathbf{P}_{\tilde{\mathbf{V}}}, \mathbf{P}_{\tilde{\mathbf{W}}}) - \mathbf{A}, \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle \right| \\ & \lesssim \bar{\Lambda}(\mathbf{A}) \left(\frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 dr}{\bar{g}_{\min}^2(\mathbf{A})} \log^{1/2} d \right) \left(\|\mathbf{V}^\top \mathbf{e}_j\|_{\ell_2} \|\mathbf{W}^\top \mathbf{e}_k\|_{\ell_2} \right. \\ & \quad \left. + \|\mathbf{U}^\top \mathbf{e}_i\|_{\ell_2} \|\mathbf{W}^\top \mathbf{e}_k\|_{\ell_2} + \|\mathbf{U}^\top \mathbf{e}_i\|_{\ell_2} \|\mathbf{V}^\top \mathbf{e}_j\|_{\ell_2} \right) \\ & \quad + \bar{\Lambda}(\mathbf{A}) \left(\frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 dr}{\bar{g}_{\min}^2(\mathbf{A})} \log^{1/2} d \right)^2 \left(\|\mathbf{V}^\top \mathbf{e}_j\|_{\ell_2} + \|\mathbf{U}^\top \mathbf{e}_i\|_{\ell_2} + \|\mathbf{W}^\top \mathbf{e}_k\|_{\ell_2} \right) \\ & \quad + \bar{\Lambda}(\mathbf{A}) \left(\frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 dr}{\bar{g}_{\min}^2(\mathbf{A})} \log^{1/2} d \right)^3. \end{aligned}$$

Recall that for all $i \in [d_1], j \in [d_2], k \in [d_3]$

$$\|\mathbf{U}^\top \mathbf{e}_i\|_{\ell_2} \leq \mu_{\mathbf{U}} \sqrt{\frac{r}{d}}, \quad \|\mathbf{V}^\top \mathbf{e}_j\|_{\ell_2} \leq \mu_{\mathbf{V}} \sqrt{\frac{r}{d}}, \quad \|\mathbf{W}^\top \mathbf{e}_k\|_{\ell_2} \leq \mu_{\mathbf{W}} \sqrt{\frac{r}{d}}$$

and conditions (3.2) (3.3) (3.4) imply

$$\frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 dr}{\bar{g}_{\min}^2(\mathbf{A})} \log^{1/2} d \lesssim r \left(\frac{\log d}{d} \right)^{1/2}.$$

We end up with a simpler bound on event \mathcal{E}_2 ,

$$\begin{aligned} & \left| \langle \mathbf{A} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}}, \mathbf{P}_{\tilde{\mathbf{V}}}, \mathbf{P}_{\tilde{\mathbf{W}}}) - \mathbf{A}, \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle \right| \\ & \lesssim \sigma r^3 \left(\frac{\sigma \tilde{\kappa}(\mathbf{A})}{\bar{g}_{\min}(\mathbf{A})} + \frac{\tilde{\kappa}^2(\mathbf{A})}{d} \right) (\mu_{\mathbf{U}} \mu_{\mathbf{V}} + \mu_{\mathbf{U}} \mu_{\mathbf{W}} + \mu_{\mathbf{V}} \mu_{\mathbf{W}}) \log^{3/2} d \end{aligned} \tag{5.10}$$

where $\tilde{\kappa}(\mathbf{A}) = \bar{\Lambda}(\mathbf{A}) / \bar{g}_{\min}(\mathbf{A})$.

Next, we prove the upper bound of $\left| \langle \mathbf{Z} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}}, \mathbf{P}_{\tilde{\mathbf{V}}}, \mathbf{P}_{\tilde{\mathbf{W}}}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle \right|$ and we proceed

with the same decomposition. Observe that

$$\begin{aligned}\langle \mathbf{Z} \cdot (\mathbf{P}_\mathbf{U}, \mathbf{P}_\mathbf{V}, \mathbf{P}_\mathbf{W}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle &= \langle \mathbf{Z}, (\mathbf{P}_\mathbf{U} \mathbf{e}_i) \otimes (\mathbf{P}_\mathbf{V} \mathbf{e}_j) \otimes (\mathbf{P}_\mathbf{W} \mathbf{e}_k) \rangle \\ &\sim \mathcal{N}\left(0, \sigma^2 \|\mathbf{P}_\mathbf{U} \mathbf{e}_i\|_{\ell_2}^2 \|\mathbf{P}_\mathbf{V} \mathbf{e}_j\|_{\ell_2}^2 \|\mathbf{P}_\mathbf{W} \mathbf{e}_k\|_{\ell_2}^2\right)\end{aligned}$$

The standard concentration inequality of Gaussian random variables yields that with probability at least $1 - \frac{1}{d^2}$,

$$\begin{aligned}|\langle \mathbf{Z} \cdot (\mathbf{P}_\mathbf{U}, \mathbf{P}_\mathbf{V}, \mathbf{P}_\mathbf{W}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle| &\lesssim \sigma \|\mathbf{U}^\top \mathbf{e}_i\|_{\ell_2} \|\mathbf{V}^\top \mathbf{e}_j\|_{\ell_2} \|\mathbf{W}^\top \mathbf{e}_k\|_{\ell_2} \log^{1/2} d \\ &\lesssim \sigma \left(\frac{r}{d}\right)^{3/2} \mu_\mathbf{U} \mu_\mathbf{V} \mu_\mathbf{W} \log^{1/2} d.\end{aligned}$$

Similarly, with probability at least $1 - \frac{1}{d^2}$,

$$\begin{aligned}&|\langle \mathbf{Z} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_\mathbf{U}, \mathbf{P}_\mathbf{V}, \mathbf{P}_\mathbf{W}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle| \\ &= |\mathbf{e}_i^\top (\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_\mathbf{U}) \mathcal{M}_1(\mathbf{Z}) (\mathbf{V} \otimes \mathbf{W}) ((\mathbf{V}^\top \mathbf{e}_j) \otimes (\mathbf{W}^\top \mathbf{e}_k))| \\ &\leq \|(\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_\mathbf{U}) \mathbf{e}_i\|_{\ell_2} \|\mathcal{M}_1(\mathbf{Z}) (\mathbf{V} \otimes \mathbf{W})\| \|\mathbf{V}^\top \mathbf{e}_j\|_{\ell_2} \|\mathbf{W}^\top \mathbf{e}_k\|_{\ell_2} \\ &\lesssim \sigma d^{1/2} \|(\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_\mathbf{U}) \mathbf{e}_i\|_{\ell_2} \|\mathbf{V}^\top \mathbf{e}_j\|_{\ell_2} \|\mathbf{W}^\top \mathbf{e}_k\|_{\ell_2}\end{aligned}$$

where we used Lemma 3 for the upper bound of $\|\mathcal{M}_1(\mathbf{Z}) (\mathbf{V} \otimes \mathbf{W})\|$. Moreover, since $\mu_\mathbf{U} \geq 1$,

$$\begin{aligned}\|(\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_\mathbf{U}) \mathbf{e}_i\|_{\ell_2} &\leq \|(\tilde{\mathbf{U}} - \mathbf{U}) \mathbf{e}_i\|_{\ell_2} + \|\tilde{\mathbf{U}} - \mathbf{U}\|_{\ell_2} \|\mathbf{U}^\top \mathbf{e}_i\|_{\ell_2} \\ &\lesssim \frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 d r}{\bar{g}_{\min}^2(\mathbf{A})} \mu_\mathbf{U} \log^{1/2} d.\end{aligned}$$

Denote the above event by \mathcal{E}_3 . On $\mathcal{E}_2 \cap \mathcal{E}_3$,

$$|\langle \mathbf{Z} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_\mathbf{U}, \mathbf{P}_\mathbf{V}, \mathbf{P}_\mathbf{W}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle| \lesssim \frac{\sigma r}{d^{1/2}} \left(\frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 d r}{\bar{g}_{\min}^2(\mathbf{A})} \right) \mu_\mathbf{U} \mu_\mathbf{V} \mu_\mathbf{W} \log^{1/2} d.$$

Similar bounds can be attained for

$$|\langle \mathbf{Z} \cdot (\mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\tilde{\mathbf{V}}} - \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\mathbf{W}}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle| \quad \text{and} \quad |\langle \mathbf{Z} \cdot (\mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\tilde{\mathbf{W}}} - \mathbf{P}_{\mathbf{W}}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle|.$$

In an identical fashion, on event $\mathcal{E}_2 \cap \mathcal{E}_3$,

$$\begin{aligned} & |\langle \mathbf{Z} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\tilde{\mathbf{V}}} - \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\mathbf{W}}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle| \\ & \lesssim \sigma r^{1/2} \left(\frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 dr}{\bar{g}_{\min}^2(\mathbf{A})} \right)^2 \mu_{\mathbf{U}} \mu_{\mathbf{V}} \mu_{\mathbf{W}} \log d. \end{aligned}$$

and

$$\begin{aligned} & |\langle \mathbf{Z} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}} - \mathbf{P}_{\mathbf{U}}, \mathbf{P}_{\tilde{\mathbf{V}}} - \mathbf{P}_{\mathbf{V}}, \mathbf{P}_{\tilde{\mathbf{W}}} - \mathbf{P}_{\mathbf{W}}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle| \\ & \lesssim \sigma d^{1/2} \left(\frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 dr}{\bar{g}_{\min}^2(\mathbf{A})} \right)^3 \mu_{\mathbf{U}} \mu_{\mathbf{V}} \mu_{\mathbf{W}} \log^{3/2} d. \end{aligned}$$

Observe by conditions (3.2) (3.3) (3.4) that

$$\frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 dr}{\bar{g}_{\min}^2(\mathbf{A})} \lesssim \frac{r}{d^{1/2}}.$$

We conclude on event $\mathcal{E}_2 \cap \mathcal{E}_3$ with

$$|\langle \mathbf{Z} \cdot (\mathbf{P}_{\tilde{\mathbf{U}}}, \mathbf{P}_{\tilde{\mathbf{V}}}, \mathbf{P}_{\tilde{\mathbf{W}}}), \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle| \lesssim \frac{\sigma r^2}{d^{1/2}} \left(\frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 dr}{\bar{g}_{\min}^2(\mathbf{A})} \right) \mu_{\mathbf{U}} \mu_{\mathbf{V}} \mu_{\mathbf{W}} \log^{3/2} d \quad (5.11)$$

By combining (5.10) and (5.11), we get on event $\mathcal{E}_2 \cap \mathcal{E}_3$,

$$\begin{aligned} & |\langle \tilde{\mathbf{A}} - \mathbf{A}, \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \rangle| \\ & \lesssim \sigma r^3 \left(\frac{\sigma \tilde{\kappa}(\mathbf{A})}{\bar{g}_{\min}(\mathbf{A})} + \frac{\tilde{\kappa}^2(\mathbf{A})}{d} \right) (\mu_{\mathbf{U}} \mu_{\mathbf{V}} + \mu_{\mathbf{U}} \mu_{\mathbf{W}} + \mu_{\mathbf{V}} \mu_{\mathbf{W}}) \log^{3/2} d \\ & + \frac{\sigma r^2}{d^{1/2}} \left(\frac{\sigma \bar{\Lambda}(\mathbf{A}) r + \sigma^2 dr}{\bar{g}_{\min}^2(\mathbf{A})} \right) \mu_{\mathbf{U}} \mu_{\mathbf{V}} \mu_{\mathbf{W}} \log^{3/2} d \end{aligned}$$

$$\lesssim \sigma r^3 \left(\frac{\sigma \tilde{\kappa}(\mathbf{A})}{\bar{g}_{\min}(\mathbf{A})} + \frac{\tilde{\kappa}^2(\mathbf{A})}{d} \right) (\mu_{\mathbf{U}}\mu_{\mathbf{V}} + \mu_{\mathbf{U}}\mu_{\mathbf{W}} + \mu_{\mathbf{V}}\mu_{\mathbf{W}}) \log^{3/2} d,$$

where the last inequality is due to fact $\bar{g}_{\min}(\mathbf{A}) \gtrsim \sigma d^{3/4}$ and $\max \{ \mu_{\mathbf{U}}, \mu_{\mathbf{V}}, \mu_{\mathbf{W}} \} \lesssim \sqrt{d}$.

CHAPTER 5
APPENDICES

Appendices

1 Proof of Lemma 1

The proof of Lemma 1 follows from a similar approach introduced by [22].

Proof. For any $S \in \mathbb{H}_m$ of rank r , $S = \sum_{j=1}^r \lambda_j (e_j \otimes e_j)$, where λ_j are non-zero eigenvalues of S (repeated with their multiplicities) and $e_j \in \mathbb{C}^m$ are the corresponding orthonormal eigenvectors. Denote $\text{sign}(S) := \sum_{j=1}^r \text{sign}(\lambda_j) (e_j \otimes e_j)$. Let $\mathcal{P}_L, \mathcal{P}_L^\perp$ be the following orthogonal projectors in the space $(\mathbb{H}_m, \langle \cdot, \cdot \rangle)$:

$$\mathcal{P}_L(A) := A - P_{L^\perp} A P_{L^\perp}, \quad \mathcal{P}_L^\perp(A) := P_{L^\perp} A P_{L^\perp}, \quad \forall A \in \mathbb{H}_m$$

where P_L denotes the orthogonal projector on the linear span of $\{e_1, \dots, e_r\}$, and P_{L^\perp} is its orthogonal complement. Clearly, this formulation provides a decomposition of a matrix A into a "low rank part" $\mathcal{P}_L(A)$ and a "high rank part" $\mathcal{P}_L^\perp(A)$ if $\text{rank}(S) = r$ is small. Given $b > 0$, define the following cone in the space \mathbb{H}_m :

$$\mathcal{K}(\mathbb{D}; L; b) := \{A \in \mathbb{D} : \|\mathcal{P}_L^\perp A\|_1 \leq b \|\mathcal{P}_L(A)\|_1\}$$

which consists of matrices with a "dominant" low rank part if S is low rank.

Firstly, we can rewrite (1.1) as

$$\hat{S}^h = \arg \min_{S \in \mathbb{D}} \frac{1}{n} \sum_{j=1}^n \left(\tilde{Y}_j - \langle S, \tilde{X}_j \rangle \right)^2 + \varepsilon \|S\|_1. \quad (1.1)$$

where $\tilde{X}_j = \text{Diag} \left[\sqrt{\frac{1}{h} K \left(\frac{\tau_j - t_0}{h} \right)} p_0 \left(\frac{\tau_j - t_0}{h} \right) X_j, \sqrt{\frac{1}{h} K \left(\frac{\tau_j - t_0}{h} \right)} p_1 \left(\frac{\tau_j - t_0}{h} \right) X_j, \dots, \sqrt{\frac{1}{h} K \left(\frac{\tau_j - t_0}{h} \right)} p_\ell \left(\frac{\tau_j - t_0}{h} \right) X_j \right]$ and $\tilde{Y}_j = \sqrt{\frac{1}{h} K \left(\frac{\tau_j - t_0}{h} \right)} Y_j$.

Denote the loss function as

$$\mathcal{L}(\tilde{Y}; \langle S(\tau), \tilde{X} \rangle) := \left(\tilde{Y}_j - \langle S, \tilde{X}_j \rangle \right)^2,$$

and the risk

$$P\mathcal{L}(\tilde{Y}; \langle S(\tau), \tilde{X} \rangle) := \mathbb{E}\mathcal{L}(\tilde{Y}; \langle S(\tau), \tilde{X} \rangle) = \sigma^2 + \mathbb{E}\frac{1}{h}K\left(\frac{\tau - t_0}{h}\right)(Y - \langle S(\tau), X \rangle)^2$$

Since \hat{S}^h is a solution of the convex optimization problem (1.1), there exists a $\hat{V} \in \partial\|\hat{S}^h\|_1$, such that for $\forall S \in \mathbb{D}$ (see [85] Chap. 2)

$$\frac{2}{n} \sum_{j=1}^n \left(\langle \hat{S}^h, \tilde{X}_j \rangle - \tilde{Y}_j \right) \langle \hat{S}^h - S, \tilde{X}_j \rangle + \varepsilon \langle \hat{V}, \hat{S}^h - S \rangle \leq 0.$$

This implies that, for all $S \in \mathbb{D}$,

$$\begin{aligned} & \mathbb{E}\mathcal{L}'(\tilde{Y}; \langle \hat{S}, \tilde{X} \rangle) \langle \hat{S}^h - S, \tilde{X} \rangle + \varepsilon \langle \hat{V}, \hat{S}^h - S \rangle \\ & \leq \mathbb{E}\mathcal{L}'(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle) \langle \hat{S}^h - S, \tilde{X} \rangle - \frac{2}{n} \sum_{j=1}^n (\langle \hat{S}^h, \tilde{X}_j \rangle - \tilde{Y}_j) \langle \hat{S}^h - S, \tilde{X}_j \rangle. \end{aligned} \quad (1.2)$$

where \mathcal{L}' denotes the partial derivative of $\mathcal{L}(y; u)$ with respect to u . One can easily check that for $\forall S \in \mathbb{D}$,

$$\mathbb{E}\mathcal{L}'(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle) \langle \hat{S}^h - S, \tilde{X} \rangle \geq \mathbb{E}(\mathcal{L}(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle) - \mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle)) + \|\hat{S}^h - S\|_{L_2(\tilde{\Pi})}^2. \quad (1.3)$$

where $\tilde{\Pi}$ denotes the distribution of \tilde{X} . If $\mathbb{E}\mathcal{L}(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle) \leq \mathbb{E}\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle)$ for $\forall S \in \mathbb{D}$, then the oracle inequality in Lemma 1 holds trivially. So we assume that $\mathbb{E}\mathcal{L}(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle) > \mathbb{E}\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle)$ for some $S \in \mathbb{D}$. Thus, inequalities (1.2) and (1.3) imply that

$$\begin{aligned} & \mathbb{E}\mathcal{L}(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle) + \|\hat{S}^h - S\|_{L_2(\tilde{\Pi})}^2 + \varepsilon \langle \hat{V}, \hat{S}^h - S \rangle \\ & \leq \mathbb{E}\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle) + \mathbb{E}\mathcal{L}'(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle) \langle \hat{S}^h - S, \tilde{X} \rangle - \frac{2}{n} \sum_{j=1}^n (\langle \hat{S}^h, \tilde{X}_j \rangle - \tilde{Y}_j) \langle \hat{S}^h - S, \tilde{X}_j \rangle. \end{aligned} \quad (1.4)$$

According to the well known representation of subdifferential of nuclear norm, see [86]

Sec. A.4, for any $V \in \partial\|S\|_1$, we have

$$V := \text{sign}(S) + \mathcal{P}_L^\perp(W), \quad W \in \mathbb{H}_m, \|W\| \leq 1.$$

By the duality between nuclear norm and operator norm

$$\langle \mathcal{P}_L^\perp(W), \hat{S}^h - S \rangle = \langle \mathcal{P}_L^\perp(W), \hat{S}^h \rangle = \langle W, \mathcal{P}_L^\perp(\hat{S}^h) \rangle = \|\mathcal{P}_L^\perp(\hat{S}^h)\|_1.$$

Therefore, by the monotonicity of subdifferentials of convex function $\|\cdot\|_1$, for any $V := \text{sign}(S) + \mathcal{P}_L^\perp(W) \in \partial\|S\|_1$, we have

$$\langle V, \hat{S}^h - S \rangle = \langle \text{sign}(S), \hat{S}^h - S \rangle + \|\mathcal{P}_L^\perp(\hat{S}^h)\|_1 \leq \langle \hat{V}, \hat{S}^h - S \rangle, \quad (1.5)$$

we can use (1.5) to change the bound in (1.4) to get

$$\begin{aligned} & \mathbb{E}\mathcal{L}(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle) + \|S - \hat{S}^h\|_{L_2(\tilde{\Pi})}^2 + \varepsilon \|\mathcal{P}_L^\perp(\hat{S}^h)\|_1 \\ & \leq \mathbb{E}\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle) + \varepsilon \langle \text{sign}(S), S - \hat{S}^h \rangle + \mathbb{E}\mathcal{L}'(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle) \langle \hat{S}^h - S, \tilde{X} \rangle \\ & \quad - \frac{2}{n} \sum_{j=1}^n (\langle \hat{S}^h, \tilde{X}_j \rangle - \tilde{Y}_j) \langle \hat{S}^h - S, \tilde{X}_j \rangle. \end{aligned} \quad (1.6)$$

For the simplicity of representation, we use the following notation to denote the empirical process:

$$\begin{aligned} & (P - P_n)(\mathcal{L}'(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle)) \langle \hat{S}^h - S, \tilde{X} \rangle := \\ & \mathbb{E}\mathcal{L}'(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle) \langle \hat{S}^h - S, \tilde{X} \rangle - \frac{2}{n} \sum_{j=1}^n (\langle \hat{S}^h, \tilde{X}_j \rangle - \tilde{Y}_j) \langle \hat{S}^h - S, \tilde{X}_j \rangle. \end{aligned} \quad (1.7)$$

The following part of the proof is to derive an upper bound on the empirical process (1.7).

Before we start with the derivation, let us present several vital ingredients that will be used

in the following literature. For a given $S \in \mathbb{D}$ and for $\delta_1, \delta_2, \delta_3, \delta_4 \geq 0$, denote

$$\mathcal{A}(\delta_1, \delta_2) := \{A \in \mathbb{D} : A - S \in \mathcal{K}(\mathbb{D}; L; b), \|A - S\|_{L_2(\tilde{\Pi})} \leq \delta_1, \|\mathcal{P}_L^\perp A\|_1 \leq \delta_2\},$$

$$\tilde{\mathcal{A}}(\delta_1, \delta_2, \delta_3) := \{A \in \mathbb{D} : \|A - S\|_{L_2(\tilde{\Pi})} \leq \delta_1, \|\mathcal{P}_L^\perp A\|_1 \leq \delta_2, \|\mathcal{P}_L(A - S)\|_1 \leq \delta_3\},$$

$$\check{\mathcal{A}}(\delta_1, \delta_4) := \{A \in \mathbb{D} : \|A - S\|_{L_2(\tilde{\Pi})} \leq \delta_1, \|A - S\|_1 \leq \delta_4\},$$

and

$$\alpha_n(\delta_1, \delta_2) := \sup\{|(P - P_n)(\mathcal{L}'(\tilde{Y}; \langle A, \tilde{X} \rangle))\langle A - S, \tilde{X} \rangle| : A \in \mathcal{A}(\delta_1, \delta_2)\},$$

$$\tilde{\alpha}_n(\delta_1, \delta_2, \delta_3) := \sup\{|(P - P_n)(\mathcal{L}'(\tilde{Y}; \langle A, \tilde{X} \rangle))\langle A - S, \tilde{X} \rangle| : A \in \tilde{\mathcal{A}}(\delta_1, \delta_2, \delta_3)\},$$

$$\check{\alpha}_n(\delta_1, \delta_4) := \sup\{|(P - P_n)(\mathcal{L}'(\tilde{Y}; \langle A, \tilde{X} \rangle))\langle A - S, \tilde{X} \rangle| : A \in \check{\mathcal{A}}(\delta_1, \delta_4)\}.$$

Given the definitions above, Lemma 6 below shows upper bounds on the three quantities

$\alpha_n(\delta_1, \delta_2)$, $\tilde{\alpha}_n(\delta_1, \delta_2, \delta_3)$, $\check{\alpha}_n(\delta_1, \delta_4)$. The proof of Lemma 6 can be found in section 2.

Denote

$$\Xi := n^{-1} \sum_{j=1}^n \varepsilon_j \tilde{X}_j \tag{1.8}$$

where ε_j are i.i.d. Rademacher random variables.

Lemma 6. Suppose $0 < \delta_k^- < \delta_k^+$, $k = 1, 2, 3, 4$. Let $\eta > 0$ and

$$\bar{\eta} := \eta + \sum_{k=1}^2 \log([\log_2(\frac{\delta_k^+}{\delta_k^-})] + 2) + \log 3,$$

$$\tilde{\eta} := \eta + \sum_{k=1}^3 \log([\log_2(\frac{\delta_k^+}{\delta_k^-})] + 2) + \log 3,$$

$$\check{\eta} := \eta + \sum_{k=1, k=4} \log([\log_2(\frac{\delta_k^+}{\delta_k^-})] + 2) + \log 3.$$

Then with probability at least $1 - e^{-\eta}$, for all $\delta_k \in [\delta_k^-, \delta_k^+]$, $k=1,2,3$

$$\alpha_n(\delta_1, \delta_2) \leq \frac{C_1(\ell+1)R(T)\Phi a}{\sqrt{h}} \left\{ \mathbb{E}\|\Xi\|(\sqrt{\text{rank}(S)}m\delta_1 + \delta_2) + \frac{2(\ell+1)R(T)\Phi a\bar{\eta}}{n\sqrt{h}} + \delta_1 \sqrt{\frac{\bar{\eta}}{n}} \right\} \quad (1.9)$$

$$\tilde{\alpha}_n(\delta_1, \delta_2, \delta_3) \leq \frac{C_2(\ell+1)R(T)\Phi a}{\sqrt{h}} \left\{ \mathbb{E}\|\Xi\|(\delta_2 + \delta_3) + \frac{2(\ell+1)R(T)\Phi a\tilde{\eta}}{n\sqrt{h}} + \delta_1 \sqrt{\frac{\tilde{\eta}}{n}} \right\} \quad (1.10)$$

$$\check{\alpha}_n(\delta_1, \delta_4) \leq \frac{C_3(\ell+1)R(T)\Phi a}{\sqrt{h}} \left\{ \mathbb{E}\|\Xi\|\delta_4 + \frac{2(\ell+1)R(T)\Phi a\check{\eta}}{n\sqrt{h}} + \delta_1 \sqrt{\frac{\check{\eta}}{n}} \right\} \quad (1.11)$$

where C_1 , C_2 , and C_3 are numerical constants.

Since both \hat{S}^h and S are in \mathbb{D} , by the definition of $\tilde{\alpha}$ and $\check{\alpha}$, we have

$$(P - P_n)(\mathcal{L}'(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle)) \langle \hat{S}^h - S, \tilde{X} \rangle \leq \tilde{\alpha}(\|\hat{S}^h - S\|_{L_2(\tilde{\Pi})}; \|\mathcal{P}_L^\perp \hat{S}^h\|_1; \|\mathcal{P}_L(\hat{S}^h - S)\|_1), \quad (1.12)$$

and

$$(P - P_n)(\mathcal{L}'(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle)) \langle \hat{S}^h - S, \tilde{X} \rangle \leq \check{\alpha}(\|\hat{S}^h - S\|_{L_2(\tilde{\Pi})}; \|\hat{S}^h - S\|_1), \quad (1.13)$$

If $\hat{S}^h - S \in \mathcal{K}(\mathbb{D}; L; b)$, by the definition of α , we have

$$(P - P_n)(\mathcal{L}'(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle)) \langle \hat{S}^h - S, \tilde{X} \rangle \leq \alpha(\|\hat{S}^h - S\|_{L_2(\tilde{\Pi})}; \|\mathcal{P}_L^\perp \hat{S}^h\|_1), \quad (1.14)$$

Assume for a while that

$$\|\hat{S}^h - S\|_{L_2(\tilde{\Pi})} \in [\delta_1^-, \delta_1^+], \quad \|\mathcal{P}_L^\perp \hat{S}^h\|_1 \in [\delta_2^-, \delta_2^+], \quad \|\mathcal{P}_L(\hat{S}^h - S)\|_1 \in [\delta_3^-, \delta_3^+]. \quad (1.15)$$

By the definition of subdifferential, for any $\widehat{V} \in \partial \|\widehat{S}^h\|_1$,

$$\langle \widehat{V}, S - \widehat{S}^h \rangle \leq \|S\|_1 - \|\widehat{S}^h\|_1.$$

Then we apply (1.13) in bound (1.4) and use the upper bound on $\check{\alpha}_n(\delta_1, \delta_4)$ of Lemma 6, and get with probability at least $1 - e^{-\eta}$,

$$\begin{aligned} & P(\mathcal{L}(\tilde{Y}; \langle \widehat{S}^h, \tilde{X} \rangle)) + \|\widehat{S}^h - S\|_{L_2(\tilde{\Pi})}^2 \\ & \leq P(\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle)) + \varepsilon(\|S\|_1 - \|\widehat{S}^h\|_1) + \check{\alpha}_n(\|\widehat{S}^h - S\|_{L_2(\tilde{\Pi})}, \|\widehat{S}^h - S\|_1) \\ & \leq P(\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle)) + \varepsilon(\|S\|_1 - \|\widehat{S}^h\|_1) \\ & \quad + \frac{C_3(\ell+1)R(T)\Phi a}{\sqrt{h}} \left\{ \mathbb{E}\|\Xi\| \|\widehat{S}^h - S\|_1 + \frac{2(\ell+1)R(T)\Phi a \check{\eta}}{n\sqrt{h}} + \|\widehat{S}^h - S\|_{L_2(\tilde{\Pi})} \sqrt{\frac{\check{\eta}}{n}} \right\}. \end{aligned} \quad (1.16)$$

Assuming that

$$\varepsilon > \frac{C(\ell+1)R(T)\Phi a}{\sqrt{h}} \mathbb{E}\|\Xi\|, \quad (1.17)$$

where $C = C_1 \vee 4C_2 \vee C_3$. From (1.16)

$$P(\mathcal{L}(\tilde{Y}; \langle \widehat{S}^h, \tilde{X} \rangle)) \leq P(\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle)) + 2\varepsilon\|S\|_1 + \frac{C_3(\ell+1)^2 R(T)^2 \Phi^2 a^2 \check{\eta}}{nh}. \quad (1.18)$$

We now apply the upper bound on $\check{\alpha}_n(\|\widehat{S}^h - S\|_{L_2(\tilde{\Pi})}, \|\mathcal{P}_L^\perp \widehat{S}^h\|_1, \|\mathcal{P}_L(\widehat{S}^h - S)\|_1)$ to (1.6) and get with probability at least $1 - e^{-\eta}$,

$$\begin{aligned} & P(\mathcal{L}(\tilde{Y}; \langle \widehat{S}^h, \tilde{X} \rangle)) + \|\widehat{S}^h - S\|_{L_2(\tilde{\Pi})}^2 + \varepsilon\|\mathcal{P}_L^\perp(\widehat{S}^h)\|_1 \\ & \leq P(\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle)) + \varepsilon\|\mathcal{P}_L(\widehat{S}^h - S)\|_1 + \check{\alpha}_n(\|\widehat{S}^h - S\|_{L_2(\tilde{\Pi})}, \|\mathcal{P}_L^\perp \widehat{S}^h\|_1, \|\mathcal{P}_L(\widehat{S}^h - S)\|_1) \\ & \leq P(\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle)) + \varepsilon\|\mathcal{P}_L(\widehat{S}^h - S)\|_1 \\ & \quad + \frac{C_2(\ell+1)R(T)\Phi a}{\sqrt{h}} \left\{ \mathbb{E}\|\Xi\| (\|\mathcal{P}_L^\perp \widehat{S}^h\|_1 + \|\mathcal{P}_L(\widehat{S}^h - S)\|_1) \right\} + \frac{C_2(\ell+1)^2 R(T)^2 \Phi^2 a^2 \check{\eta}}{nh}, \end{aligned} \quad (1.19)$$

where the first inequality is due to the fact that

$$|\langle \text{sign}(S), S - \widehat{S}^h \rangle| = |\langle \text{sign}(S), \mathcal{P}_L(S - \widehat{S}^h) \rangle| \leq \|\text{sign}(S)\| \|\mathcal{P}_L(S - \widehat{S}^h)\|_1 \leq \|\mathcal{P}_L(S - \widehat{S}^h)\|_1.$$

With assumption (1.17) holds, we get from (1.19)

$$\begin{aligned} & P\mathcal{L}(\tilde{Y}; \langle \widehat{S}^h, \tilde{X} \rangle) + \varepsilon \|\mathcal{P}_L^\perp(\widehat{S}^h)\|_1 \\ & \leq P\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle) + \frac{5\varepsilon}{4} \|\mathcal{P}_L(\widehat{S}^h - S)\|_1 + \frac{\varepsilon}{4} \|\mathcal{P}_L^\perp(\widehat{S}^h)\|_1 + \frac{C_2(\ell+1)^2 R(T)^2 \Phi^2 a^2 \tilde{\eta}}{nh}. \end{aligned} \quad (1.20)$$

If the following is satisfied:

$$\frac{C_2(\ell+1)^2 R(T)^2 \Phi^2 a^2 \tilde{\eta}}{nh} \geq \frac{5\varepsilon}{4} \|\mathcal{P}_L(\widehat{S}^h - S)\|_1 + \frac{\varepsilon}{4} \|\mathcal{P}_L^\perp(\widehat{S}^h)\|_1, \quad (1.21)$$

we can just conclude that

$$P(\mathcal{L}(\tilde{Y}; \langle \widehat{S}^h, \tilde{X} \rangle)) \leq P(\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle)) + \frac{C_2(\ell+1)^2 R(T)^2 \Phi^2 a^2 \tilde{\eta}}{nh}, \quad (1.22)$$

which is sufficient to meet the bound of Lemma 1. Otherwise, by the assumption that

$P(\mathcal{L}(\tilde{Y}; \langle \widehat{S}^h, \tilde{X} \rangle)) > P(\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle))$, one can easily check that

$$\|\mathcal{P}_L^\perp(\widehat{S}^h - S)\|_1 \leq 5 \|\mathcal{P}_L(\widehat{S}^h - S)\|_1,$$

which implies that $\widehat{S}^h - S \in \mathcal{K}(\mathbb{D}; L; 5)$. This fact allows us to use the bound on $\alpha_n(\delta_1, \delta_2)$

of Lemma 6. We get from (1.6)

$$\begin{aligned} & P(\mathcal{L}(\tilde{Y}; \langle \widehat{S}^h, \tilde{X} \rangle)) + \|\widehat{S}^h - S\|_{L_2(\tilde{\Pi})}^2 + \varepsilon \|\mathcal{P}_L^\perp(\widehat{S}^h)\|_1 \\ & \leq P(\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle)) + \varepsilon \langle \text{sign}(S), S - \widehat{S}^h \rangle \\ & + \frac{C_1(\ell+1)R(T)\Phi a}{\sqrt{h}} \mathbb{E} \|\Xi\| (\sqrt{\text{rank}(S)} m \|\widehat{S}^h - S\|_{L_2(\tilde{\Pi})} + \|\mathcal{P}_L^\perp(\widehat{S}^h)\|_1) + \frac{C_1(\ell+1)^2 R(T)^2 \Phi^2 a^2 \tilde{\eta}}{nh}. \end{aligned} \quad (1.23)$$

By applying the inequality

$$|\langle \text{sign}(S), \hat{S}^h - S \rangle| \leq m \sqrt{\text{rank}(S)} \|\hat{S}^h - S\|_{L_2(\tilde{\Pi})},$$

and the assumption (1.17), we have with probability at least $1 - e^{-\eta}$,

$$P(\mathcal{L}(\tilde{Y}; \langle \hat{S}^h, \tilde{X} \rangle)) \leq P(\mathcal{L}(\tilde{Y}; \langle S, \tilde{X} \rangle)) + \varepsilon^2 m^2 \text{rank}(S) + \frac{C_1(\ell + 1)^2 R(T)^2 \Phi^2 a^2 \bar{\eta}}{nh}. \quad (1.24)$$

To sum up, the bound of Lemma 1 follows from (1.18), (1.22) and (1.24) provided that condition (1.17) and condition (1.15) hold.

We still need to specify $\delta_k^-, \delta_k^+, k = 1, 2, 3, 4$ to establish the bound of the theorem. By the definition of \hat{S}^h , we have

$$P_n(\mathcal{L}(\tilde{Y}; \langle X, \hat{S}^h \rangle)) + \varepsilon \|\hat{S}^h\|_1 \leq P_n(\mathcal{L}(\tilde{Y}; \langle X, 0 \rangle)) \leq Q,$$

implying that $\|\hat{S}^h\|_1 \leq \frac{Q}{\varepsilon}$. Next, $\|\mathcal{P}_L^\perp \hat{S}^h\|_1 \leq \|\hat{S}^h\|_1 \leq \frac{Q}{\varepsilon}$ and $\|\mathcal{P}_L(\hat{S}^h - S)\|_1 \leq 2\|\hat{S}^h - S\|_1 \leq \frac{2Q}{\varepsilon} + 2\|S\|_1$. Finally, we have $\|\hat{S}^h - S\|_{L_2(\tilde{\Pi})} \leq 2a$. Thus, we can take $\delta_1^+ := 2a$, $\delta_2^+ := \frac{Q}{\varepsilon}$, $\delta_3^+ := \frac{2Q}{\varepsilon} + 2\|S\|_1$, $\delta_4^+ := \frac{Q}{\varepsilon} + \|S\|_1$. With these choices, $\delta_k^+, k = 1, 2, 3, 4$ are upper bounds on the corresponding norms in condition (1.15). We choose $\delta_1^- := \frac{a}{\sqrt{n}}$, $\delta_2^- := \frac{a^2}{n\varepsilon} \wedge \frac{\delta_2^+}{2}$, $\delta_3^- := \frac{a^2}{n\varepsilon} \wedge \frac{\delta_3^+}{2}$, $\delta_4^- := \frac{a^2}{n\varepsilon} \wedge \frac{\delta_4^+}{2}$. Let $\eta^* := \eta + 3 \log(B \log_2(\|S\|_1 \vee n \vee \varepsilon \vee a^{-1} \vee Q))$. It is easy to verify that $\bar{\eta} \vee \tilde{\eta} \vee \tilde{\eta} \leq \eta^*$. for a proper choice of numerical constant B in the definition of η^* . When condition (1.15) does not hold, which means at least one of the numbers $\delta_k^-, k = 1, 2, 3, 4$ we chose is not a lower bound on the corresponding norm, we can still use the bounds

$$\begin{aligned} & (P - P_n)(\mathcal{L}'(Y; \langle \hat{S}^h, \tilde{X} \rangle)) \langle \hat{S}^h - S, \tilde{X} \rangle \\ & \leq \tilde{\alpha} (\|\hat{S}^h - S\|_{L_2(\tilde{\Pi})} \vee \delta_1^-; \|\mathcal{P}_L^\perp \hat{S}^h\|_1 \vee \delta_2^-; \|\mathcal{P}_L(\hat{S}^h - S)\|_1 \vee \delta_3^-), \end{aligned} \quad (1.25)$$

and

$$(P - P_n)(\mathcal{L}'(Y; \langle \hat{S}^h, \tilde{X} \rangle)) \langle \hat{S}^h - S, \tilde{X} \rangle \leq \check{\alpha}(\|\hat{A}(t)^\varepsilon - S\|_{L_2(\tilde{\Pi})} \vee \delta_1^-; \|\hat{S}^h - S\|_1 \vee \delta_4^-), \quad (1.26)$$

instead of (1.12), (1.13). In the case when $\hat{S}^h - S \in \mathcal{K}(\mathbb{D}; L; 5)$, we can use the bound

$$(P - P_n)(\mathcal{L}'(Y; \langle \hat{S}^h, \tilde{X} \rangle)) \langle \hat{S}^h - S, \tilde{X} \rangle \leq \alpha(\|\hat{S}^h - S\|_{L_2(\tilde{\Pi})} \vee \delta_1^-; \|\mathcal{P}_L^\perp \hat{S}^h\|_1 \vee \delta_2^-), \quad (1.27)$$

instead of bound (1.14). Then one can repeat the arguments above with only minor modifications. By the adjusting the constants, the result of Lemma 1 holds.

The last thing we need to specify is the size of ε which controls the nuclear norm penalty. Recall that from condition (1.17), the essence is to control $\mathbb{E}\|\Xi\|$. Here we use a simple but powerful noncommutative matrix Bernstein inequalities. The original approach was introduced by [87]. Later, the result was improved by [23] based on the classical result of [88]. We give the following lemma which is a direct consequence of the result proved by [23], and we omit the proof here.

Lemma 7. Under the model (1.1), Ξ is defined as in (1.8) with τ_j are i.i.d. uniformly distributed in $[0,1]$, and ε_j are i.i.d. Rademacher random variables, and X_j are i.i.d uniformly distributed in \mathcal{X} . Then for any $\eta > 0$, with probability at least $1 - e^{-\eta}$

$$\|\Xi\| \leq 4 \left(\sqrt{\frac{(\eta + \log 2m)}{nm}} \vee \frac{(\eta + \log 2m)\Phi}{n\sqrt{h}} \right),$$

and by integrating its exponential tail bounds

$$\mathbb{E}\|\Xi\| \leq C \left(\sqrt{\frac{\log 2m}{nm}} \vee \frac{(\log 2m)\Phi}{n\sqrt{h}} \right)$$

where C is a numerical constant.

Together with (1.17), we know for some numerical constant $D > 0$,

$$\varepsilon \geq D \frac{\Phi a(\ell+1)R(T)}{\sqrt{h}} \left(\sqrt{\frac{\log 2m}{nm}} \vee \frac{(\log 2m)\Phi}{n\sqrt{h}} \right).$$

which completes the proof of Lemma 1. \square

2 Proof of Lemma 6

Proof. We only prove the first bound in detail, and proofs of the rest two bounds are similar with only minor modifications. By Talagrand's concentration inequality [89], and its Bousquet's form [90], with probability at least $1 - e^{-\eta}$,

$$\alpha_n(\delta_1, \delta_2) \leq 2\mathbb{E}\alpha_n(\delta_1, \delta_2) + \frac{24(\ell+1)^2 R(T)^2 \Phi^2 a^2 \eta}{nh} + \frac{12(\ell+1)R(T)\Phi a \delta_1 \sqrt{\eta}}{\sqrt{nh}}. \quad (2.1)$$

By standard Rademacher symmetrization inequalities, see [86], Sec. 2.1, we can get

$$\mathbb{E}\alpha_n(\delta_1, \delta_2) \leq 4\mathbb{E} \sup \left\{ \left| \frac{1}{n} \sum_{j=1}^n \varepsilon_j (\langle A, \tilde{X}_j \rangle - \tilde{Y}_j) \langle A - S, \tilde{X}_j \rangle \right| : A \in \mathcal{A}(\delta_1, \delta_2) \right\}, \quad (2.2)$$

where $\{\varepsilon_j\}$ are i.i.d. Rademacher random variables independent of $\{(\tau_j, X_j, \tilde{Y}_j)\}$. Then we consider the function $f(u) = (u - y + v)u$, where $|y| \leq \frac{2\Phi a}{\sqrt{h}}$ and $|v|, |u| \leq \frac{2(\ell+1)R(T)\Phi a}{\sqrt{h}}$. Clearly, this function has a Lipschitz constant $\frac{6(\ell+1)R(T)\Phi a}{\sqrt{h}}$. Thus by comparison inequality, see [86], Sec. 2.2, we can get

$$\begin{aligned} & \mathbb{E} \sup \left\{ \left| n^{-1} \sum_{j=1}^n \varepsilon_j (\langle A, \tilde{X}_j \rangle - \tilde{Y}_j) \langle A - S, \tilde{X}_j \rangle \right| : A \in \mathcal{A}(\delta_1, \delta_2) \right\} \\ & \leq \frac{6(\ell+1)R(T)\Phi a}{\sqrt{h}} \mathbb{E} \sup \left\{ n^{-1} \left| \sum_{j=1}^n \varepsilon_i \langle A - S, \tilde{X}_j \rangle \right| : A \in \mathcal{A}(\delta_1, \delta_2) \right\}. \end{aligned} \quad (2.3)$$

As a consequence of (2.2) and (2.3), we have

$$\mathbb{E}\alpha_n(\delta_1, \delta_2) \leq \frac{12(\ell+1)R(T)\Phi a}{\sqrt{h}} \mathbb{E} \sup \left\{ n^{-1} \left| \sum_{j=1}^n \varepsilon_i \langle A-S, \tilde{X}_j \rangle \right| : A \in \mathcal{A}(\delta_1, \delta_2) \right\}. \quad (2.4)$$

The next step is to get an upper bound on $\left| n^{-1} \sum_{j=1}^n \varepsilon_i \langle A-S, \tilde{X}_j \rangle \right|$. Recall that $\Xi := n^{-1} \sum_{j=1}^n \varepsilon_j \tilde{X}_j$, then we have $n^{-1} \sum_{j=1}^n \varepsilon_i \langle A-S, \tilde{X}_j \rangle = \langle A-S, \Xi \rangle$. One can check that

$$\begin{aligned} |\langle A-S, \Xi \rangle| &\leq |\langle \mathcal{P}_L(A-S), \mathcal{P}_L \Xi \rangle| + |\langle \mathcal{P}_L^\perp(A-S), \Xi \rangle| \\ &\leq \|\mathcal{P}_L \Xi\|_2 \|\mathcal{P}_L(A-S)\|_2 + \|\Xi\| \|\mathcal{P}_L^\perp A\|_1 \\ &\leq m \sqrt{2\text{rank}(S)} \|\Xi\| \|A-S\|_{L_2(\tilde{\Pi})} + \|\Xi\| \|\mathcal{P}_L^\perp A\|_1. \end{aligned}$$

The second line of this inequality is due to Hölder's inequality and the third line is due to the facts that $(A-S) \in \mathcal{K}(\mathbb{D}; L; 5)$, $\text{rank}(\mathcal{P}_L(\Xi)) \leq 2\text{rank}(S)$, $\|\mathcal{P}_L \Xi\|_2 \leq 2\sqrt{\text{rank}(\mathcal{P}_L(\Xi))} \|\Xi\|$, and $\|A-S\|_{L_2(\tilde{\Pi})}^2 = \frac{1}{m^2} \|A-S\|_2^2$. Therefore,

$$\begin{aligned} &\frac{12(\ell+1)R(T)\Phi a}{\sqrt{h}} \mathbb{E} \sup \left\{ \left| n^{-1} \sum_{j=1}^n \varepsilon_i \langle A-S, \tilde{X}_j \rangle \right| : A \in \mathcal{A}(\delta_1, \delta_2) \right\} \\ &\leq \frac{12(\ell+1)R(T)\Phi a}{\sqrt{h}} \mathbb{E} \|\Xi\| (2\sqrt{2\text{rank}(S)} m \delta_1 + \delta_2). \end{aligned} \quad (2.5)$$

It follows from (2.1), (2.4) and (2.5) that with probability at least $1 - e^{-\eta}$,

$$\begin{aligned} \alpha_n(\delta_1, \delta_2) &\leq \left(\frac{12(\ell+1)R(T)\Phi a}{\sqrt{h}} \mathbb{E} \|\Xi\| (\sqrt{\text{rank}(S)} m \delta_1 + \delta_2) \right) \\ &\quad + \frac{24(\ell+1)^2 R(T)^2 \Phi^2 a^2 \eta}{nh} + \frac{12(\ell+1)R(T)\Phi a \delta_1 \sqrt{\eta}}{\sqrt{nh}}. \end{aligned}$$

Now similar to the approach in [22], we make this bound uniform in $\delta_k \in [\delta_k^-, \delta_k^+]$. Let $\delta_k^{j_k} = \delta_k^+ 2^{-j_k}$, $j_k = 0, \dots, [\log_2(\delta_k^+ / \delta_k^-)] + 1$, $k = 1, 2$. By the union bound, with probability

at least $1 - e^{-\eta}/3$, for all $j_k = 0, \dots, [\log_2(\delta_k^+/\delta_k^-)] + 1$, $k = 1, 2$, we have

$$\begin{aligned} \alpha_n(\delta_1, \delta_2) &\leq \left(\frac{12(\ell+1)R(T)\Phi a}{\sqrt{h}} \mathbb{E}\|\Xi\|(\sqrt{\text{rank}(S)}m\delta_1^{j_1} + \delta_2^{j_2}) \right) \\ &\quad + \frac{24(\ell+1)^2 R(T)^2 \Phi^2 a^2 \eta}{nh} + \frac{12(\ell+1)R(T)\Phi a\delta_1^{j_1}\sqrt{\eta}}{\sqrt{nh}}. \end{aligned}$$

which implies that for all $\delta_k \in [\delta_k^-, \delta_k^+]$, $k = 1, 2$,

$$\begin{aligned} \alpha_n(\delta_1, \delta_2) &\leq \left(\frac{12(\ell+1)R(T)\Phi a}{\sqrt{h}} \mathbb{E}\|\Xi\|(\sqrt{\text{rank}(S)}m\delta_1 + \delta_2) \right) \\ &\quad + \frac{24(\ell+1)^2 R(T)^2 \Phi^2 a^2 \bar{\eta}}{nh} + \frac{12(\ell+1)R(T)\Phi a\delta_1\sqrt{\bar{\eta}}}{\sqrt{nh}}. \end{aligned}$$

The proofs of the second and the third bounds are similar to this one, we omit the repeated arguments. \square

3 Proof of Lemma 3

Let $\mathbf{z}_i \in \mathbb{R}^{m_1}$, $i = 1, \dots, m_2$ denote the columns of \mathbf{Z} . Then, we write

$$\mathbf{Z}\mathbf{Z}^\top - \sigma^2 m_2 \mathbf{I}_{m_1} = \sum_{i=1}^{m_2} (\mathbf{z}_i \otimes \mathbf{z}_i - \sigma^2 \mathbf{I}_{m_1}).$$

Similarly, let $\tilde{\mathbf{z}}_j \in \mathbb{R}^{m_1}$, $j = 1, \dots, m_1$ denote the rows of \mathbf{Z} and observe that $\|\mathbf{B}\mathbf{Z}^\top\| = \|\mathbf{B}\mathbf{Z}^\top \mathbf{Z} \mathbf{B}^\top\|^{1/2}$ and

$$\mathbf{B}\mathbf{Z}^\top \mathbf{Z} \mathbf{B}^\top = \sum_{j=1}^{m_1} \left((\mathbf{B}\tilde{\mathbf{z}}_j) \otimes (\mathbf{B}\tilde{\mathbf{z}}_j) - \sigma^2 \mathbf{B}\mathbf{B}^\top \right).$$

The inequalities (5.7) and (5.2) are on the concentration of sample covariance operator, where a sharp bound has been derived in [91] and will be skipped here.

4 Proof of Theorem 5.1

Since $\mathbb{E}\hat{\Gamma} = \mathbf{0}$, we immediately get $\mathbb{E}\mathbf{L}_k(\hat{\Gamma}) = \mathbf{0}$. Then,

$$\langle \mathbf{x}, \hat{\mathbf{P}}_k^{uu} \mathbf{y} \rangle - \mathbb{E} \langle \mathbf{x}, \hat{\mathbf{P}}_k^{uu} \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{L}_k(\hat{\Gamma}) \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{S}_k(\hat{\Gamma}) \mathbf{y} \rangle - \mathbb{E} \langle \mathbf{x}, \mathbf{S}_k(\hat{\Gamma}) \mathbf{y} \rangle.$$

Lemma 8. For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{m_1}$, there exists an absolute constant $D_1 > 0$ such that for all $0 \leq t \leq m_1$, with probability at least $1 - e^{-t}$,

$$|\langle \mathbf{x}, \mathbf{L}_k(\hat{\Gamma}) \mathbf{y} \rangle| \leq D_1 t^{1/2} \left(\frac{\sigma \mu_1 + \sigma^2 m_2^{1/2}}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \right) \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}.$$

Proof. Recall that

$$\hat{\Gamma} = \mathbf{A} \mathbf{Z}^\top + \mathbf{Z} \mathbf{A}^\top + \mathbf{Z} \mathbf{Z}^\top - m_2 \sigma^2 \mathbf{I}_{m_1}.$$

Then, we write $\langle \mathbf{x}, \mathbf{L}_k(\hat{\Gamma}) \mathbf{y} \rangle$ as

$$\begin{aligned} \langle \mathbf{x}, \mathbf{L}_k(\hat{\Gamma}) \mathbf{y} \rangle &= \langle \hat{\Gamma} \mathbf{P}_k^{uu} \mathbf{x}, \mathbf{C}_k^{uu} \mathbf{y} \rangle + \langle \hat{\Gamma} \mathbf{C}_k^{uu} \mathbf{x}, \mathbf{P}_k^{uu} \mathbf{y} \rangle \\ &= \langle (\mathbf{A} \mathbf{Z}^\top + \mathbf{Z} \mathbf{A}^\top + \mathbf{Z} \mathbf{Z}^\top - m_2 \sigma^2 \mathbf{I}_{m_1}) \mathbf{P}_k^{uu} \mathbf{x}, \mathbf{C}_k^{uu} \mathbf{y} \rangle \\ &\quad + \langle (\mathbf{A} \mathbf{Z}^\top + \mathbf{Z} \mathbf{A}^\top + \mathbf{Z} \mathbf{Z}^\top - m_2 \sigma^2 \mathbf{I}_{m_1}) \mathbf{C}_k^{uu} \mathbf{x}, \mathbf{P}_k^{uu} \mathbf{y} \rangle. \end{aligned}$$

It suffices to consider the following terms separately for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{m_1}$:

$$\langle \mathbf{Z} \mathbf{A}^\top \mathbf{x}, \mathbf{y} \rangle, \quad \langle \mathbf{A} \mathbf{Z}^\top \mathbf{x}, \mathbf{y} \rangle, \quad \langle (\mathbf{Z} \mathbf{Z}^\top - m_2 \sigma^2 \mathbf{I}_{m_1}) \mathbf{x}, \mathbf{y} \rangle.$$

It is straightforward to check that $\langle \mathbf{Z} \mathbf{A}^\top \mathbf{x}, \mathbf{y} \rangle$ is a normal random variable with zero mean and variance

$$\mathbb{E} \langle \mathbf{Z} \mathbf{A}^\top \mathbf{x}, \mathbf{y} \rangle^2 = \mathbb{E} \langle \mathbf{Z}, \mathbf{y} \otimes (\mathbf{A}^\top \mathbf{x}) \rangle^2 = \sigma^2 \|\mathbf{y} \otimes (\mathbf{A}^\top \mathbf{x})\|_{\ell_2}^2 = \sigma^2 \|\mathbf{y}\|_{\ell_2}^2 \|\mathbf{A}^\top \mathbf{x}\|_{\ell_2}^2,$$

where we used the fact that \mathbf{Z} is a $m_1 \times m_2$ matrix with i.i.d. $\mathcal{N}(0, \sigma^2)$ entries. Therefore,

$$\mathbb{E} \langle \mathbf{Z} \mathbf{A}^\top \mathbf{P}_k^{uu} \mathbf{x}, \mathbf{C}_k^{uu} \mathbf{y} \rangle^2 \leq \frac{\sigma^2 \mu_k^2}{\bar{g}_k^2(\mathbf{A} \mathbf{A}^\top)} \|\mathbf{x}\|_{\ell_2}^2 \|\mathbf{y}\|_{\ell_2}^2,$$

where we used the facts $\|\mathbf{C}_k\| \leq \frac{1}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)}$ and $\|\mathbf{A}^\top \mathbf{P}_k^{uu}\| \leq \mu_k$. By the standard concentration inequality of Gaussian random variables, we get for all $t \geq 0$,

$$\mathbb{P} \left(\left| \langle \mathbf{Z} \mathbf{A}^\top \mathbf{P}_k^{uu} \mathbf{x}, \mathbf{C}_k^{uu} \mathbf{y} \rangle \right| \geq 2t^{1/2} \frac{\sigma \mu_k}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} \right) \leq e^{-t}.$$

Similarly, for all $t \geq 0$,

$$\mathbb{P} \left(\left| \langle \mathbf{Z} \mathbf{A}^\top \mathbf{C}_k^{uu} \mathbf{x}, \mathbf{P}_k^{uu} \mathbf{y} \rangle \right| \geq 2t^{1/2} \frac{\sigma \mu_1}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} \right) \leq e^{-t}.$$

We next turn to the bound of $|\langle (\mathbf{Z} \mathbf{Z}^\top - m_2 \sigma^2 \mathbf{I}_{m_1}) \mathbf{P}_k^{uu} \mathbf{x}, \mathbf{C}_k^{uu} \mathbf{y} \rangle|$. Recall that $\mathbf{P}_k^{uu} \mathbf{C}_k^{uu} = \mathbf{0}$ implying that it suffices to consider $\langle \mathbf{Z} \mathbf{Z}^\top \mathbf{P}_k^{uu} \mathbf{x}, \mathbf{C}_k^{uu} \mathbf{y} \rangle$. Let $\mathbf{z}_1, \dots, \mathbf{z}_{m_2} \in \mathbb{R}^{m_1}$ denote the columns of \mathbf{Z} such that $\mathbf{z}_i \in \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_{m_1})$ for $1 \leq i \leq m_2$. Write

$$\langle \mathbf{Z} \mathbf{Z}^\top (\mathbf{P}_k^{uu} \mathbf{x}), \mathbf{C}_k^{uu} \mathbf{y} \rangle = \sum_{i=1}^{m_2} \langle \mathbf{z}_i, \mathbf{P}_k^{uu} \mathbf{x} \rangle \langle \mathbf{z}_i, \mathbf{C}_k^{uu} \mathbf{y} \rangle.$$

Observe that $\mathbb{E}(\mathbf{P}_k^{uu} \mathbf{z}_i) \otimes (\mathbf{C}_k^{uu} \mathbf{z}_i) = \mathbf{0}$ implying that $\langle \mathbf{z}_i, \mathbf{P}_k^{uu} \mathbf{x} \rangle$ is independent of $\langle \mathbf{z}_i, \mathbf{C}_k^{uu} \mathbf{y} \rangle$.

By concentration inequalities of Gaussian random variables, for all $t \geq 0$,

$$\mathbb{P} \left(\left| \langle \mathbf{Z} \mathbf{Z}^\top (\mathbf{P}_k^{uu} \mathbf{x}), \mathbf{C}_k^{uu} \mathbf{y} \rangle \right| \geq 2t^{1/2} \|\mathbf{y}\|_{\ell_2} \frac{\sigma \left(\sum_{i=1}^{m_2} \langle \mathbf{z}_i, \mathbf{P}_k^{uu} \mathbf{x} \rangle^2 \right)^{1/2}}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \left| \left\{ \langle \mathbf{z}_i, \mathbf{P}_k^{uu} \mathbf{x} \rangle : i = 1, \dots, m_2 \right\} \right| \right) \leq e^{-t}.$$

By [92, Prop 5.16], the following bound holds with probability at least $1 - e^{-t}$,

$$\left| \sum_{i=1}^{m_2} \langle \mathbf{z}_i, \mathbf{P}_k^{uu} \mathbf{x} \rangle^2 - \sigma^2 m_2 \|\mathbf{x}\|_{\ell_2}^2 \right| \lesssim \sigma \left(m_2^{1/2} t^{1/2} + t \right) \|\mathbf{x}\|_{\ell_2}.$$

If $t \lesssim m_1 \leq m_2$, we conclude that there exists an absolute constant $D_1 > 0$ such that

$$\mathbb{P}\left(\left|\langle \mathbf{Z}\mathbf{Z}^\top (\mathbf{P}_k^{uu} \mathbf{x}), \mathbf{C}_k^{uu} \mathbf{y} \rangle\right| \geq D_1 \frac{\sigma^2 m_2^{1/2} t^{1/2}}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}\right) \leq e^{-t}.$$

To sum up, for all $0 \leq t \lesssim m_1$, the following bound holds with probability at least $1 - e^{-t}$,

$$|\langle \mathbf{x}, \mathbf{L}_k(\hat{\Gamma}) \mathbf{y} \rangle| \lesssim t^{1/2} \left(\frac{\sigma \mu_1 + \sigma^2 m_2^{1/2}}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right) \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}$$

which concludes the proof. \square

It remains to derive the upper bound of $|\langle \mathbf{x}, \mathbf{S}_k(\hat{\Gamma}) \mathbf{y} \rangle - \mathbb{E} \langle \mathbf{x}, \mathbf{S}_k(\hat{\Gamma}) \mathbf{y} \rangle|$. The following lemma is due to [69].

Lemma 9. Let $\delta(m_1, m_2) := \sigma \mu_1 m_1^{1/2} + \sigma^2 (m_1 m_2)^{1/2}$ and suppose that $\delta(m_1, m_2) \leq \frac{1-\gamma}{2(1+\gamma)} \bar{g}_k(\mathbf{A}\mathbf{A}^\top)$ for some $\gamma \in (0, 1)$. There exists a constant $D_\gamma > 0$ such that, for all symmetric $\hat{\Gamma}_1, \hat{\Gamma}_2 \in \mathbb{R}^{m_1 \times m_1}$ satisfying the condition $\max \{ \|\hat{\Gamma}_1\|, \|\hat{\Gamma}_2\| \} \leq (1 + \gamma) \delta(m_1, m_2)$,

$$\|\mathbf{S}_k(\hat{\Gamma}_1) - \mathbf{S}_k(\hat{\Gamma}_2)\| \leq D_\gamma \frac{\delta(m_1, m_2)}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \|\hat{\Gamma}_1 - \hat{\Gamma}_2\|.$$

Define function $\varphi(\cdot) : \mathbb{R}_+ \mapsto [0, 1]$ such that $\varphi(t) = 1$ for $0 \leq t \leq 1$ and $\varphi(t) = 0$ for $t \geq (1 + \gamma)$ and φ is linear in between. Then, function φ is Lipschitz on \mathbb{R}_+ with constant $\frac{1}{\gamma}$. To illustrate the dependence of $\hat{\Gamma}$ on \mathbf{Z} , we write $\hat{\Gamma}(\mathbf{Z})$ instead of $\hat{\Gamma}$. To this end, fix $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{m_1}$ and constants $\delta_1, \delta_2 > 0$ and define the function

$$F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}) := \left\langle \mathbf{x}, \mathbf{S}_k(\hat{\Gamma}(\mathbf{Z})) \mathbf{y} \right\rangle \varphi\left(\frac{\|\hat{\Gamma}(\mathbf{Z})\|}{\delta_1}\right) \varphi\left(\frac{\|\mathbf{Z}\|}{\delta_2}\right).$$

where we view \mathbf{Z} as a point in $\mathbb{R}^{m_1 \times m_2}$ rather than a random matrix.

Lemma 10. For any $\delta_1 \leq \frac{1-\gamma}{2(1+\gamma)} \bar{g}_k(\mathbf{A}\mathbf{A}^\top)$ for some $\gamma \in (0, 1)$ and $\delta_2 > 0$, there exists an

absolute constant $C_\gamma > 0$ such that

$$|F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}_1) - F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}_2)| \leq C_\gamma \frac{\delta_1}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \left(\mu_1 + \delta_2 + \frac{\delta_1}{\delta_2} \right) \|\mathbf{Z}_1 - \mathbf{Z}_2\| \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}$$

Proof of Lemma 10. Since $\varphi\left(\frac{\|\hat{\Gamma}(\mathbf{Z})\|}{\delta_1}\right)\varphi\left(\frac{\|\mathbf{Z}\|}{\delta_2}\right) \neq 0$ only if $\|\hat{\Gamma}(\mathbf{Z})\| \leq (1 + \gamma)\delta_1$ and $\|\mathbf{Z}\| \leq (1 + \gamma)\delta_2$, Lemma 4 implies that

$$|F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z})| = \left| \left\langle \mathbf{x}, \mathbf{S}_k(\hat{\Gamma}(\mathbf{Z}))\mathbf{y} \right\rangle \varphi\left(\frac{\|\hat{\Gamma}(\mathbf{Z})\|}{\delta_1}\right) \varphi\left(\frac{\|\mathbf{Z}\|}{\delta_2}\right) \right| \leq 14(1 + \gamma)^2 \frac{\delta_1^2}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)}.$$

Case 1. If $\max\{\|\hat{\Gamma}(\mathbf{Z}_1)\|, \|\hat{\Gamma}(\mathbf{Z}_2)\|\} \leq (1 + \gamma)\delta_1$ and $\max\{\|\mathbf{Z}_1\|, \|\mathbf{Z}_2\|\} \leq (1 + \gamma)\delta_2$.

By the Lipschitzity of function φ , Lemma 9 and definition of $\hat{\Gamma}(\mathbf{Z})$, it is easy to check

$$\begin{aligned} & |F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}_1) - F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}_2)| \\ & \leq \|\mathbf{S}_k(\hat{\Gamma}(\mathbf{Z}_1)) - \mathbf{S}_k(\hat{\Gamma}(\mathbf{Z}_2))\| \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} \\ & \quad + \frac{14(1 + \gamma)^2 \delta_1}{\gamma \bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \|\hat{\Gamma}(\mathbf{Z}_1) - \hat{\Gamma}(\mathbf{Z}_2)\| \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} + \frac{14(1 + \gamma)^2 \delta_1^2}{\delta_2 \gamma \bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \|\mathbf{Z}_1 - \mathbf{Z}_2\| \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} \\ & \leq D_\gamma \frac{\delta_1}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \|\hat{\Gamma}(\mathbf{Z}_1) - \hat{\Gamma}(\mathbf{Z}_2)\| \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} + \frac{14(1 + \gamma)^2 \delta_1^2}{\delta_2 \gamma \bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \|\mathbf{Z}_1 - \mathbf{Z}_2\| \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} \\ & \leq D_\gamma \frac{\delta_1}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \left(\mu_1 + \delta_2 + \frac{\delta_1}{\delta_2} \right) \|\mathbf{Z}_1 - \mathbf{Z}_2\| \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}. \end{aligned}$$

Case 2. If $\|\hat{\Gamma}(\mathbf{Z}_1)\| \leq (1 + \gamma)\delta_1$, $\|\hat{\Gamma}(\mathbf{Z}_2)\| \geq (1 + \gamma)\delta_1$ and $\max\{\|\mathbf{Z}_1\|, \|\mathbf{Z}_2\|\} \leq (1 + \gamma)\delta_2$. Since $\|\hat{\Gamma}(\mathbf{Z}_2)\| \geq (1 + \gamma)\delta_1$, we have $\varphi\left(\frac{\|\hat{\Gamma}(\mathbf{Z}_2)\|}{\delta_1}\right) = 0$ and $F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}_2) = 0$.

Then,

$$\begin{aligned} & |F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}_1) - F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}_2)| \\ & = \left| \left\langle \mathbf{x}, \mathbf{S}_k(\hat{\Gamma}(\mathbf{Z}_1))\mathbf{y} \right\rangle \varphi\left(\frac{\|\hat{\Gamma}(\mathbf{Z}_1)\|}{\delta_1}\right) \varphi\left(\frac{\|\mathbf{Z}_1\|}{\delta_2}\right) \right| \\ & = \left| \left\langle \mathbf{x}, \mathbf{S}_k(\hat{\Gamma}(\mathbf{Z}_1))\mathbf{y} \right\rangle \varphi\left(\frac{\|\hat{\Gamma}(\mathbf{Z}_1)\|}{\delta_1}\right) \varphi\left(\frac{\|\mathbf{Z}_1\|}{\delta_2}\right) - \left\langle \mathbf{x}, \mathbf{S}_k(\hat{\Gamma}(\mathbf{Z}_1))\mathbf{y} \right\rangle \varphi\left(\frac{\|\hat{\Gamma}(\mathbf{Z}_2)\|}{\delta_1}\right) \varphi\left(\frac{\|\mathbf{Z}_1\|}{\delta_2}\right) \right| \end{aligned}$$

$$\begin{aligned}
&\leq \|\mathbf{S}_k(\widehat{\Gamma}(\mathbf{Z}_1))\| \frac{1}{\delta_1 \gamma} \|\widehat{\Gamma}(\mathbf{Z}_1) - \widehat{\Gamma}(\mathbf{Z}_2)\| \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} \\
&\leq \frac{(1+\gamma)^2 \delta_1^2}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top) \delta_1 \gamma} (2\mu_1 + 2(1+\gamma)\delta_2) \|\mathbf{Z}_1 - \mathbf{Z}_2\| \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} \\
&\leq D_\gamma \frac{\delta_1}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} (\mu_1 + \delta_2) \|\mathbf{Z}_1 - \mathbf{Z}_2\| \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}.
\end{aligned}$$

Case 3. If $\|\widehat{\Gamma}(\mathbf{Z}_1)\| \leq (1+\gamma)\delta_1$, $\|\widehat{\Gamma}(\mathbf{Z}_2)\| \geq (1+\gamma)\delta_1$, $\|\mathbf{Z}_1\| \leq (1+\gamma)\delta_2$, $\|\mathbf{Z}_2\| \geq (1+\gamma)\delta_2$. It can be proved similarly as *Case 2*.

Case 4. If $\|\widehat{\Gamma}(\mathbf{Z}_1)\| \leq (1+\gamma)\delta_1$, $\|\widehat{\Gamma}(\mathbf{Z}_2)\| \geq (1+\gamma)\delta_1$, $\|\mathbf{Z}_1\| \geq (1+\gamma)\delta_2$, $\|\mathbf{Z}_2\| \geq (1+\gamma)\delta_2$. It is a trivial case since $F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}_1) = F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}_2) = 0$.

Case 5. If $\max\{\|\widehat{\Gamma}(\mathbf{Z}_1)\|, \|\widehat{\Gamma}(\mathbf{Z}_2)\|\} \leq (1+\gamma)\delta_1$, $\|\mathbf{Z}_1\| \leq (1+\gamma)\delta_2$, $\|\mathbf{Z}_2\| \geq (1+\gamma)\delta_2$. Again, we have $F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}_2) = 0$. Then,

$$\begin{aligned}
&|F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}_1) - F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}_2)| \\
&= \left| \left\langle \mathbf{x}, \mathbf{S}_k(\widehat{\Gamma}(\mathbf{Z}_1)) \mathbf{y} \right\rangle \varphi\left(\frac{\|\widehat{\Gamma}(\mathbf{Z}_1)\|}{\delta_1}\right) \varphi\left(\frac{\|\mathbf{Z}_1\|}{\delta_2}\right) - \left\langle \mathbf{x}, \mathbf{S}_k(\widehat{\Gamma}(\mathbf{Z}_2)) \mathbf{y} \right\rangle \varphi\left(\frac{\|\widehat{\Gamma}(\mathbf{Z}_2)\|}{\delta_1}\right) \varphi\left(\frac{\|\mathbf{Z}_2\|}{\delta_2}\right) \right| \\
&\leq \|\mathbf{S}_k(\widehat{\Gamma}(\mathbf{Z}_1))\| \frac{1}{\delta_2 \gamma} \|\mathbf{Z}_1 - \mathbf{Z}_2\| \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} \leq \frac{(1+\gamma)^2 \delta_1^2}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top) \delta_2 \gamma} \|\mathbf{Z}_1 - \mathbf{Z}_2\| \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} \\
&\leq D_\gamma \frac{\delta_1}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \frac{\delta_1}{\delta_2} \|\mathbf{Z}_1 - \mathbf{Z}_2\| \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}.
\end{aligned}$$

All the other cases shall be handled similarly and we conclude the proof. \square

Note that $\|\mathbf{Z}_1 - \mathbf{Z}_2\| \leq \|\mathbf{Z}_1 - \mathbf{Z}_2\|_{\ell_2}$, Lemma 10 indicates that $F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z})$ is Lipschitz with constant

$$D_\gamma \frac{\delta_1}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \left(\mu_1 + \delta_2 + \frac{\delta_1}{\delta_2} \right) \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}.$$

Lemma 11. Let $\delta(m_1, m_2) := \sigma \mu_1 m_1^{1/2} + \sigma^2 (m_1 m_2)^{1/2}$ and suppose that $\mathbb{E}\|\widehat{\Gamma}\| \leq \frac{1-\gamma}{2} \bar{g}_k(\mathbf{A}\mathbf{A}^\top)$ for some $\gamma \in (0, 1)$. There exists some constant D_γ such that for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{m_1}$ and all

$\log 8 \leq t \leq m_1$, the following inequality holds with probability at least $1 - e^{-t}$,

$$|\langle \mathbf{x}, \mathbf{S}_k(\widehat{\Gamma})\mathbf{y} \rangle - \mathbb{E}\langle \mathbf{x}, \mathbf{S}_k(\widehat{\Gamma})\mathbf{y} \rangle| \leq D_\gamma t^{1/2} \frac{\sigma\mu_1 + \sigma^2 m_2^{1/2}}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \left(\frac{\delta(m_1, m_2)}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right) \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}.$$

Proof of Lemma 11. Choose $\delta_1 = \delta_1(m_1, m_2)$ and $\delta_2 = \delta_2(m_1, m_2)$ as follows where $\log 8 \leq t \leq m_1$ is to be determined:

$$\begin{aligned} \delta_1(m_1, m_2) &:= \delta_1(m_1, m_2, t) := \mathbb{E}\|\widetilde{\Gamma}\| + D_1 t^{1/2} (\sigma\mu_1 + \sigma^2 m_2^{1/2}) \\ \delta_2(m_1, m_2) &:= \delta_2(m_1, m_2, t) := \mathbb{E}\|\mathbf{Z}\| + D_2 \sigma t^{1/2} \end{aligned}$$

and the constants $D_1, D_2 > 0$ are chosen such that $\mathbb{P}(\|\widehat{\Gamma}\| \geq \delta_1(m_1, m_2, t)) \leq e^{-t}$ and $\mathbb{P}(\|\mathbf{Z}\| \geq \delta_2(m_1, m_2, t)) \leq e^{-t}$. Let $M := \text{Med}(\langle \mathbf{x}, \mathbf{S}_k(\widehat{\Gamma})\mathbf{y} \rangle)$ denote its median.

Case 1. If $D_1 t^{1/2} (\mu_1 \sigma + \sigma^2 m_2^{1/2}) \leq \frac{\gamma}{4} \bar{g}_k(\mathbf{A}\mathbf{A}^\top)$. Then, $\delta_1 \leq (1 - \frac{\gamma}{2}) \frac{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)}{2} = \frac{1 - 2\gamma'}{1 + 2\gamma'} \frac{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)}{2}$ for some $\gamma' \in (0, 1/2)$. By Lemma 10, $F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\cdot)$ satisfies the Lipschitz condition. By definition of $F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z})$, we have $F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}) = \langle \mathbf{x}, \mathbf{S}_k(\widehat{\Gamma})\mathbf{y} \rangle$ on the event $\{\|\widehat{\Gamma}\| \leq \delta_1, \|\mathbf{Z}\| \leq \delta_2\}$. By Lemma 3 and $t \geq \log 8$,

$$\begin{aligned} &\mathbb{P}\left\{F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}) \geq M\right\} \\ &\geq \mathbb{P}\left\{F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}) \geq M, \quad \|\widehat{\Gamma}\| \leq \delta_1, \quad \|\mathbf{Z}\| \leq \delta_2\right\} \\ &\geq \mathbb{P}\left\{\langle \mathbf{x}, \mathbf{S}_k(\widehat{\Gamma})\mathbf{y} \rangle \geq M\right\} - \mathbb{P}\{\|\widehat{\Gamma}\| \leq \delta_1, \|\mathbf{Z}\| \leq \delta_2\} \\ &\geq \mathbb{P}\left\{\langle \mathbf{x}, \mathbf{S}_k(\widehat{\Gamma})\mathbf{y} \rangle \geq M\right\} - \mathbb{P}\left\{\|\widehat{\Gamma}\| \leq \delta_1\right\} - \mathbb{P}\left\{\|\mathbf{Z}\| \leq \delta_2\right\} \\ &\geq \frac{1}{2} - \frac{1}{8} - \frac{1}{8} = 1/4, \end{aligned}$$

and similarly,

$$\mathbb{P}\left\{F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}) \leq M\right\} \geq 1/4.$$

It follows from Gaussian isoperimetric inequality (see [63, Lemma 2.6]) and Lemma 10 that with some constant $D_\gamma > 0$, for all $t \geq \log 8$ with probability at least $1 - e^{-t}$,

$$|F_{\delta_1, \delta_2, \mathbf{x}, \mathbf{y}}(\mathbf{Z}) - M| \leq D_\gamma \frac{\sigma \delta_1 t^{1/2}}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \left(\mu_1 + \delta_2 + \frac{\delta_1}{\delta_2} \right) \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}.$$

Since $t \leq m_1 \leq m_2$, it is easy to check by Lemma 3 that $\delta_1 \asymp \sigma \mu_1 m_1^{1/2} + \sigma^2 (m_1 m_2)^{1/2}$ and $\delta_2 \asymp \sigma m_2^{1/2}$. Moreover, $\mathbb{P}\{\|\hat{\Gamma}\| \leq \delta_1, \|\mathbf{Z}\| \leq \delta_2\} \geq 1 - 2e^{-t}$. As a result, with probability at least $1 - e^{-3t}$,

$$|\langle \mathbf{x}, \mathbf{S}_k(\hat{\Gamma})\mathbf{y} \rangle - M| \leq D_\gamma \frac{\sigma \mu_1 t^{1/2} + \sigma^2 m_2^{1/2} t^{1/2}}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \left(\frac{\delta(m_1, m_2)}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right) \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}. \quad (4.1)$$

Case 2. If $D_1 t^{1/2} (\sigma \mu_1 + \sigma^2 m_2^{1/2}) > \frac{\gamma}{4} \bar{g}_k(\mathbf{A}\mathbf{A}^\top)$. It implies that

$$\mathbb{E}\|\hat{\Gamma}\| \leq D_1 \frac{(1 - \gamma)}{\gamma} t^{1/2} (\sigma \mu_1 + \sigma^2 m_2^{1/2}),$$

and $\delta_1 \leq D_\gamma t^{1/2} (\sigma \mu_1 + \sigma^2 m_2^{1/2})$. By Lemma 3 and Lemma 4, with probability at least $1 - e^{-t}$,

$$|\langle \mathbf{x}, \mathbf{S}_k(\hat{\Gamma})\mathbf{y} \rangle| \leq \|\mathbf{S}_k(\hat{\Gamma})\| \leq D_\gamma t \frac{(\sigma \mu_1 + \sigma^2 m_2^{1/2})^2}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2},$$

which immediately yields that

$$M \leq D_\gamma \frac{(\sigma \mu_1 + \sigma^2 m_2^{1/2})^2}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}.$$

The above inequalities imply that with probability at least $1 - e^{-t}$ for $\log 8 \leq t \leq m_1$,

$$\begin{aligned} |\langle \mathbf{x}, \mathbf{S}_k(\hat{\Gamma})\mathbf{y} \rangle - M| &\leq D_\gamma t \frac{(\sigma \mu_1 + \sigma^2 m_2^{1/2})^2}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2} \\ &\leq D_\gamma \frac{\sigma \mu_1 t^{1/2} + \sigma^2 m_2^{1/2} t^{1/2}}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \left(\frac{\delta(m_1, m_2)}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \right) \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}. \end{aligned} \quad (4.2)$$

Therefore, bounds (4.1) and (4.2) hold in both cases. The rest of the proof is quite standard by integrating the exponential tails and will be skipped here, see [63]. \square

Proof of Theorem 5.1. By Lemma 8 and Lemma 11, if $D_1\delta(m_1, m_2) \leq \bar{g}_k(\mathbf{A}\mathbf{A}^\top)$ for a large enough constant $D_1 > 0$ such that $\gamma \leq 1/2$, we conclude that for all $\log 8 \leq t \leq m_1$, with probability at least $1 - 2e^{-t}$,

$$|\langle \mathbf{x}, \hat{\mathbf{P}}_k \mathbf{y} \rangle| \leq Dt^{1/2} \frac{\sigma\mu_1 + \sigma^2 m_2^{1/2}}{\bar{g}_k(\mathbf{A}\mathbf{A}^\top)} \|\mathbf{x}\|_{\ell_2} \|\mathbf{y}\|_{\ell_2}$$

which concludes the proof after adjusting the constant D accordingly. \square

5 Proof of Lemma 5

Observe that for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{m_1}$ with $\|\mathbf{x}\|_{\ell_2} = \|\mathbf{y}\|_{\ell_2} = 1$ and $\delta_t = \mathbb{E}\|\hat{\mathbf{\Gamma}}\| + D_1\sigma\mu_1 t^{1/2} + D_2\sigma^2 m_2^{1/2} t^{1/2}$ with $t \leq m_1$ and some $\gamma \in (0, 1/2]$,

$$\begin{aligned} \left| \mathbb{E} \langle \mathbf{x}, (\mathbf{S}_k(\tilde{\mathbf{\Gamma}}) - \mathbf{S}_k(\hat{\mathbf{\Gamma}})) \mathbf{y} \rangle \right| &\leq \mathbb{E} \left\| \mathbf{S}_k(\tilde{\mathbf{\Gamma}}) - \mathbf{S}_k(\hat{\mathbf{\Gamma}}) \right\| \\ &= \mathbb{E} \left\| \mathbf{S}_k(\tilde{\mathbf{\Gamma}}) - \mathbf{S}_k(\hat{\mathbf{\Gamma}}) \right\| \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| \leq (1 + \gamma)\delta_t) \mathbf{1}(\|\hat{\mathbf{\Gamma}}\| \leq (1 + \gamma)\delta_t) \\ &\quad + \mathbb{E} \left\| \mathbf{S}_k(\tilde{\mathbf{\Gamma}}) - \mathbf{S}_k(\hat{\mathbf{\Gamma}}) \right\| \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| \leq (1 + \gamma)\delta_t) \mathbf{1}(\|\hat{\mathbf{\Gamma}}\| > (1 + \gamma)\delta_t) \\ &\quad + \mathbb{E} \left\| \mathbf{S}_k(\tilde{\mathbf{\Gamma}}) - \mathbf{S}_k(\hat{\mathbf{\Gamma}}) \right\| \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| > (1 + \gamma)\delta_t) \mathbf{1}(\|\hat{\mathbf{\Gamma}}\| \leq (1 + \gamma)\delta_t) \\ &\quad + \mathbb{E} \left\| \mathbf{S}_k(\tilde{\mathbf{\Gamma}}) - \mathbf{S}_k(\hat{\mathbf{\Gamma}}) \right\| \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| > (1 + \gamma)\delta_t) \mathbf{1}(\|\hat{\mathbf{\Gamma}}\| > (1 + \gamma)\delta_t) \end{aligned}$$

where the constants $D_1, D_2 > 0$ are chosen such that $\max \{ \mathbb{P}(\|\tilde{\mathbf{\Gamma}}\| \geq \delta_t), \mathbb{P}(\|\hat{\mathbf{\Gamma}}\| \geq \delta_t) \} \leq e^{-t}$. By Lemma 9,

$$\begin{aligned} &\mathbb{E} \left\| \mathbf{S}_k(\tilde{\mathbf{\Gamma}}) - \mathbf{S}_k(\hat{\mathbf{\Gamma}}) \right\| \mathbf{1}(\|\tilde{\mathbf{\Gamma}}\| \leq (1 + \gamma)\delta_t) \mathbf{1}(\|\hat{\mathbf{\Gamma}}\| \leq (1 + \gamma)\delta_t) \\ &\leq D_\gamma \frac{\delta_t}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \mathbb{E} \|\tilde{\mathbf{\Gamma}} - \hat{\mathbf{\Gamma}}\| \leq D_\gamma \frac{\delta_t}{\bar{g}_k^2(\mathbf{A}\mathbf{A}^\top)} \mathbb{E} \|\mathbf{Z}\mathbf{P}_k^{hh}\mathbf{Z}^\top - \nu_k\sigma^2\mathbf{I}_{m_1}\|. \end{aligned}$$

By writing $\mathbf{P}_k^{hh} := \sum_{j \in \Delta_k} \mathbf{h}_j \otimes \mathbf{h}_j$, we obtain

$$\begin{aligned} \mathbf{Z} \mathbf{P}_k^{hh} \mathbf{Z}^\top - \sigma^2 \nu_k \mathbf{I}_{m_1} &= \sum_{j \in \Delta_k} (\mathbf{Z} \mathbf{h}_j) \otimes (\mathbf{Z} \mathbf{h}_j) - \sigma^2 \nu_k \mathbf{I}_{m_1} \\ &= \nu_k \left(\frac{1}{\nu_k} \sum_{j \in \Delta_k} (\mathbf{Z} \mathbf{h}_j) \otimes (\mathbf{Z} \mathbf{h}_j) - \sigma^2 \mathbf{I}_{m_1} \right). \end{aligned}$$

where $\nu_k = \text{Card}(\Delta_k)$. The vectors $\mathbf{Z} \mathbf{h}_j \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_{m_1})$ and $\{\mathbf{Z} \mathbf{h}_j : \dots, j \in \Delta_k\}$ are independent. By [91],

$$\mathbb{E} \left\| \frac{1}{\nu_k} \sum_{j \in \Delta_k} (\mathbf{Z} \mathbf{h}_j) \otimes (\mathbf{Z} \mathbf{h}_j) - \sigma^2 \mathbf{I}_{m_1} \right\| \lesssim \sigma^2 \left(\sqrt{\frac{m_1}{\nu_k}} \vee \frac{m_1}{\nu_k} \right).$$

Since $\nu_k \leq m_1$, we conclude with

$$\begin{aligned} \mathbb{E} \left\| \mathbf{S}_k(\tilde{\Gamma}) - \mathbf{S}_k(\hat{\Gamma}) \right\| \mathbf{1} \left(\|\tilde{\Gamma}\| \leq (1 + \gamma) \delta_t \right) \mathbf{1} \left(\|\hat{\Gamma}\| \leq (1 + \gamma) \delta_t \right) \\ \lesssim \gamma \frac{\delta_t}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \left(\frac{m_1 \sigma^2}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \right). \end{aligned} \quad (5.1)$$

Choose $t = m_1$, by Lemma 4 and Lemma 3,

$$\begin{aligned} \mathbb{E} \left\| \mathbf{S}_k(\tilde{\Gamma}) - \mathbf{S}_k(\hat{\Gamma}) \right\| \mathbf{1} \left(\|\hat{\Gamma}\| \leq (1 + \gamma) \delta_{m_1} \right) \mathbf{1} \left(\|\tilde{\Gamma}\| > (1 + \gamma) \delta_{m_1} \right) \\ \leq D_\gamma \frac{\delta_{m_1}^2}{\bar{g}_k^2(\mathbf{A} \mathbf{A}^\top)} \mathbb{E} \frac{\|\tilde{\Gamma}\|^2}{\bar{g}_k^2(\mathbf{A} \mathbf{A}^\top)} \mathbf{1} \left(\|\tilde{\Gamma}\| > (1 + \gamma) \delta_{m_1} \right) \\ \lesssim \gamma \frac{\delta_{m_1}^2}{\bar{g}_k^4(\mathbf{A} \mathbf{A}^\top)} e^{-m_1/2} \mathbb{E}^{1/2} \|\tilde{\Gamma}\|^4 \lesssim \frac{\delta_{m_1}^4}{\bar{g}_k^4(\mathbf{A} \mathbf{A}^\top)} e^{-m_1/2} \\ \lesssim \frac{\delta(m_1, m_2)}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \left(\frac{\sigma \mu_1 + \sigma^2 m_1}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \right) \end{aligned}$$

which is clearly dominated by (5.1) for $t = m_1$ and $m_2 e^{-m_1/2} \leq 1$. The other terms are bounded in a similar fashion. To sum up, we obtain

$$\|\mathbb{E} \mathbf{S}_k(\tilde{\Gamma}) - \mathbb{E} \mathbf{S}_k(\hat{\Gamma})\| \lesssim \frac{\sigma \mu_1 + \sigma^2 m_1}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \left(\frac{\delta(m_1, m_2)}{\bar{g}_k(\mathbf{A} \mathbf{A}^\top)} \right).$$

REFERENCES

- [1] F. Zhou, “Nonparametric estimation of low rank matrix valued function,” *ArXiv preprint arXiv:1802.06292*, 2018.
- [2] D. Xia and F. Zhou, “The ℓ_∞ perturbation of hosvd and low rank tensor denoising,” *ArXiv preprint arXiv:1707.01207*, 2017.
- [3] E. J. Candès and B. Recht, “Exact matrix completion via convex optimization,” *Foundations of Computational mathematics*, vol. 9, no. 6, p. 717, 2009.
- [4] J.-F. Cai, E. J. Candès, and Z. Shen, “A singular value thresholding algorithm for matrix completion,” *SIAM Journal on Optimization*, vol. 20, no. 4, pp. 1956–1982, 2010.
- [5] E. J. Candès and Y. Plan, “Matrix completion with noise,” *Proceedings of the IEEE*, vol. 98, no. 6, pp. 925–936, 2010.
- [6] E. J. Candès and T. Tao, “The power of convex relaxation: Near-optimal matrix completion,” *IEEE Transactions on Information Theory*, vol. 56, no. 5, pp. 2053–2080, 2010.
- [7] V. Koltchinskii, “Von Neumann entropy penalization and low-rank matrix estimation,” *The Annals of Statistics*, pp. 2936–2973, 2011.
- [8] D. Gross, “Recovering low-rank matrices from few coefficients in any basis,” *IEEE Transactions on Information Theory*, vol. 57, no. 3, pp. 1548–1566, 2011.
- [9] V. Koltchinskii, K. Lounici, and A. B. Tsybakov, “Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion,” *The Annals of Statistics*, vol. 39, no. 5, pp. 2302–2329, 2011.
- [10] A. Rohde and A. B. Tsybakov, “Estimation of high-dimensional low-rank matrices,” *The Annals of Statistics*, vol. 39, no. 2, pp. 887–930, 2011.
- [11] S. Negahban and M. J. Wainwright, “Restricted strong convexity and weighted matrix completion: Optimal bounds with noise,” *Journal of Machine Learning Research*, vol. 13, no. May, pp. 1665–1697, 2012.
- [12] S. Chatterjee, “Matrix estimation by universal singular value thresholding,” *The Annals of Statistics*, vol. 43, no. 1, pp. 177–214, 2015.

- [13] Y. Koren, R. Bell, and C. Volinsky, “Matrix factorization techniques for recommender systems,” *Computer*, vol. 42, no. 8, 2009.
- [14] Y. Koren, “Collaborative filtering with temporal dynamics,” *Communications of the ACM*, vol. 53, no. 4, pp. 89–97, 2010.
- [15] L. Zhang, G. Wahba, and M. Yuan, “Distance shrinkage and euclidean embedding via regularized kernel estimation,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 78, no. 4, pp. 849–867, 2016.
- [16] A. Singer and M. Cucuringu, “Uniqueness of low-rank matrix completion by rigidity theory,” *SIAM Journal on Matrix Analysis and Applications*, vol. 31, no. 4, pp. 1621–1641, 2010.
- [17] B. Recht, M. Fazel, and P. A. Parrilo, “Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization,” *SIAM review*, vol. 52, no. 3, pp. 471–501, 2010.
- [18] J. Fan and I. Gijbels, *Local polynomial modelling and its applications: Monographs on statistics and applied probability 66*. CRC Press, 1996, vol. 66.
- [19] O. V. Lepskii, “On a problem of adaptive estimation in gaussian white noise,” *Theory of Probability & Its Applications*, vol. 35, no. 3, pp. 454–466, 1991.
- [20] A. R. Barron, “Complexity regularization with application to artificial neural networks,” *Nonparametric Functional Estimation and Related Topics*, vol. 335, pp. 561–576, 1991.
- [21] M. Wegkamp, “Model selection in nonparametric regression,” *The Annals of Statistics*, vol. 31, no. 1, pp. 252–273, 2003.
- [22] V. Koltchinskii, “Sharp oracle inequalities in low rank estimation,” in *Empirical Inference*, Springer, 2013, pp. 217–230.
- [23] J. A. Tropp, “User-friendly tail bounds for sums of random matrices,” *Foundations of computational mathematics*, vol. 12, no. 4, pp. 389–434, 2012.
- [24] A. B. Tsybakov, *Introduction to nonparametric estimation. revised and extended from the 2004 french original. translated by vladimir zaiats*, 2009.
- [25] E. N. Gilbert, “A comparison of signalling alphabets,” *Bell Labs Technical Journal*, vol. 31, no. 3, pp. 504–522, 1952.
- [26] R. R. Varshamov, “Estimate of the number of signals in error correcting codes,” in *Dokl. Akad. Nauk SSSR*, vol. 117, 1957, pp. 739–741.

- [27] O. V. Lepski and V. G. Spokoiny, “Optimal pointwise adaptive methods in nonparametric estimation,” *The Annals of Statistics*, pp. 2512–2546, 1997.
- [28] O. V. Lepski, E. Mammen, and V. G. Spokoiny, “Optimal spatial adaptation to inhomogeneous smoothness: An approach based on kernel estimates with variable bandwidth selectors,” *The Annals of Statistics*, pp. 929–947, 1997.
- [29] V. Koltchinskii, “Local rademacher complexities and oracle inequalities in risk minimization,” *The Annals of Statistics*, vol. 34, no. 6, pp. 2593–2656, 2006.
- [30] C. C. Craig, “On the Tchebychef inequality of Bernstein,” *The Annals of Mathematical Statistics*, vol. 4, no. 2, pp. 94–102, 1933.
- [31] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, *et al.*, “Distributed optimization and statistical learning via the alternating direction method of multipliers,” *Foundations and Trends® in Machine learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [32] Z. Lin, M. Chen, and Y. Ma, “The augmented lagrange multiplier method for exact recovery of corrupted low-rank matrices,” *ArXiv preprint arXiv:1009.5055*, 2010.
- [33] C. Chen, B. He, and X. Yuan, “Matrix completion via an alternating direction method,” *IMA Journal of Numerical Analysis*, vol. 32, no. 1, pp. 227–245, 2012.
- [34] J. Liu, P. Musialski, P. Wonka, and J. Ye, “Tensor completion for estimating missing values in visual data,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 208–220, 2013.
- [35] C.-F. Westin, S. E. Maier, H. Mamata, A. Nabavi, F. A. Jolesz, and R. Kikinis, “Processing and visualization for diffusion tensor mri,” *Medical image analysis*, vol. 6, no. 2, pp. 93–108, 2002.
- [36] T Hildebrand and P Rügsegger, “A new method for the model-independent assessment of thickness in three-dimensional images,” *Journal of microscopy*, vol. 185, no. 1, pp. 67–75, 1997.
- [37] N. Li and B. Li, “Tensor completion for on-board compression of hyperspectral images,” in *Image Processing (ICIP), 2010 17th IEEE International Conference on*, IEEE, 2010, pp. 517–520.
- [38] M Vasilescu and D. Terzopoulos, “Multilinear analysis of image ensembles: Tensor-faces,” *Computer Vision?ECCV 2002*, pp. 447–460, 2002.
- [39] A. Anandkumar, R. Ge, D. Hsu, S. M. Kakade, and M. Telgarsky, “Tensor decompositions for learning latent variable models,” *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 2773–2832, 2014.

- [40] A. Cichocki, D. Mandic, L. De Lathauwer, G. Zhou, Q. Zhao, C. Caiafa, and H. A. Phan, "Tensor decompositions for signal processing applications: From two-way to multiway component analysis," *IEEE Signal Processing Magazine*, vol. 32, no. 2, pp. 145–163, 2015.
- [41] A. T. Chaganty and P. Liang, "Spectral experts for estimating mixtures of linear regressions," in *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, 2013, pp. 1040–1048.
- [42] L. Omberg, G. H. Golub, and O. Alter, "A tensor higher-order singular value decomposition for integrative analysis of DNA microarray data from different studies," *Proceedings of the National Academy of Sciences*, vol. 104, no. 47, pp. 18 371–18 376, 2007.
- [43] C. Muralidhara, A. M. Gross, R. R. Gutell, and O. Alter, "Tensor decomposition reveals concurrent evolutionary convergences and divergences and correlations with structural motifs in ribosomal rna," *PloS one*, vol. 6, no. 4, e18768, 2011.
- [44] S. P. Ponnappalli, M. A. Saunders, C. F. Van Loan, and O. Alter, "A higher-order generalized singular value decomposition for comparison of global mrna expression from multiple organisms," *PloS one*, vol. 6, no. 12, e28072, 2011.
- [45] C. J. Hillar and L.-H. Lim, "Most tensor problems are NP-hard," *Journal of the ACM (JACM)*, vol. 60, no. 6, p. 45, 2013.
- [46] Q. Zheng and R. Tomioka, "Interpolating convex and non-convex tensor decompositions via the subspace norm," in *Advances in Neural Information Processing Systems*, 2015, pp. 3106–3113.
- [47] L. De Lathauwer, B. De Moor, and J. Vandewalle, "A multilinear singular value decomposition," *SIAM journal on Matrix Analysis and Applications*, vol. 21, no. 4, pp. 1253–1278, 2000.
- [48] G. Bergqvist and E. G. Larsson, "The higher-order singular value decomposition: Theory and an application [lecture notes]," *IEEE Signal Processing Magazine*, vol. 27, no. 3, pp. 151–154, 2010.
- [49] J. Chen and Y. Saad, "On the tensor svd and the optimal low rank orthogonal approximation of tensors," *SIAM Journal on Matrix Analysis and Applications*, vol. 30, no. 4, pp. 1709–1734, 2009.
- [50] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," *SIAM review*, vol. 51, no. 3, pp. 455–500, 2009.

- [51] E. Acar and B. Yener, “Unsupervised multiway data analysis: A literature survey,” *IEEE transactions on knowledge and data engineering*, vol. 21, no. 1, pp. 6–20, 2009.
- [52] E. Richard and A. Montanari, “A statistical model for tensor PCA,” in *Advances in Neural Information Processing Systems*, 2014, pp. 2897–2905.
- [53] A. Zhang and D. Xia, “Tensor svd: Statistical and computational limits,” *IEEE Transactions on Information Theory*, to appear, 2018+.
- [54] S. B. Hopkins, J. Shi, and D. Steurer, “Tensor principal component analysis via sum-of-square proofs,” in *COLT*, 2015, pp. 956–1006.
- [55] T. Liu, M. Yuan, and H. Zhao, “Characterizing spatiotemporal transcriptome of human brain via low rank tensor decomposition,” *ArXiv preprint arXiv:1702.07449*, 2017.
- [56] L. De Lathauwer, B. De Moor, and J. Vandewalle, “On the best rank-1 and rank-(r_1, r_2, \dots, r_n) approximation of higher-order tensors,” *SIAM Journal on Matrix Analysis and Applications*, vol. 21, no. 4, pp. 1324–1342, 2000.
- [57] L. Florescu and W. Perkins, “Spectral thresholds in the bipartite stochastic block model,” *ArXiv preprint arXiv:1506.06737*, 2015.
- [58] M. E. Newman, “Detecting community structure in networks,” *The European Physical Journal B-Condensed Matter and Complex Systems*, vol. 38, no. 2, pp. 321–330, 2004.
- [59] P. Mitra, “Entrywise bounds for eigenvectors of random graphs,” *Electronic journal of combinatorics*, vol. 16, no. 1, R131, 2009.
- [60] J. Jin, “Fast community detection by score,” *The Annals of Statistics*, vol. 43, no. 1, pp. 57–89, 2015.
- [61] T. T. Cai, T. Liang, and A. Rakhlin, “Computational and statistical boundaries for submatrix localization in a large noisy matrix,” *ArXiv preprint arXiv:1502.01988*, 2015.
- [62] Z. Ma and Y. Wu, “Computational barriers in minimax submatrix detection,” *The Annals of Statistics*, vol. 43, no. 3, pp. 1089–1116, 2015.
- [63] V. Koltchinskii and D. Xia, “Perturbation of linear forms of singular vectors under gaussian noise,” in *High Dimensional Probability VII*, Springer, 2016, pp. 397–423.

- [64] P. Wedin, “Perturbation bounds in connection with singular value decomposition,” *BIT Numerical Mathematics*, vol. 12, no. 1, pp. 99–111, 1972.
- [65] R. Wang, “Singular vector perturbation under gaussian noise,” *SIAM Journal on Matrix Analysis and Applications*, vol. 36, no. 1, pp. 158–177, 2015.
- [66] T. T. Cai and A. Zhang, “Rate-optimal perturbation bounds for singular subspaces with applications to high-dimensional statistics,” *ArXiv preprint arXiv:1605.00353*, 2016.
- [67] J. Fan, W. Wang, and Y. Zhong, “An ℓ_∞ eigenvector perturbation bound and its application to robust covariance estimation,” *ArXiv preprint arXiv:1603.03516*, 2016.
- [68] J. Cape, M. Tang, and C. E. Priebe, “The two-to-infinity norm and singular subspace geometry with applications to high-dimensional statistics,” *ArXiv preprint arXiv:1705.10735*, 2017.
- [69] V. Koltchinskii and K. Lounici, “Asymptotics and concentration bounds for bilinear forms of spectral projectors of sample covariance,” in *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, Institut Henri Poincaré, vol. 52, 2016, pp. 1976–2013.
- [70] D. Xia and M. Yuan, “On polynomial time methods for exact low rank tensor completion,” *ArXiv preprint arXiv:1702.06980*, 2017.
- [71] M. Rudelson, R. Vershynin, *et al.*, “Delocalization of eigenvectors of random matrices with independent entries,” *Duke Mathematical Journal*, vol. 164, no. 13, pp. 2507–2538, 2015.
- [72] V. Vu and K. Wang, “Random weighted projections, random quadratic forms and random eigenvectors,” *Random Structures & Algorithms*, vol. 47, no. 4, pp. 792–821, 2015.
- [73] V. Hore, A. Viñuela, A. Buil, J. Knight, M. I. McCarthy, K. Small, and J. Marchini, “Tensor decomposition for multiple-tissue gene expression experiments,” *Nature Genetics*, vol. 48, no. 9, pp. 1094–1100, 2016.
- [74] A. McCallum, K. Nigam, and L. H. Ungar, “Efficient clustering of high-dimensional data sets with application to reference matching,” in *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, 2000, pp. 169–178.
- [75] L. Parsons, E. Haque, and H. Liu, “Subspace clustering for high dimensional data: A review,” *Acm Sigkdd Explorations Newsletter*, vol. 6, no. 1, pp. 90–105, 2004.

- [76] J. Fan and Y. Fan, “High dimensional classification using features annealed independence rules,” *Annals of statistics*, vol. 36, no. 6, p. 2605, 2008.
- [77] T. Hastie, R. Tibshirani, and J. Friedman, “Unsupervised learning,” in *The elements of statistical learning*, Springer, 2009, pp. 485–585.
- [78] J. H. Friedman, “Regularized discriminant analysis,” *Journal of the American statistical association*, vol. 84, no. 405, pp. 165–175, 1989.
- [79] Q. Xiong, N. Ancona, E. R. Hauser, S. Mukherjee, and T. S. Furey, “Integrating genetic and gene expression evidence into genome-wide association analysis of gene sets,” *Genome research*, vol. 22, no. 2, pp. 386–397, 2012.
- [80] M. Kolar, S. Balakrishnan, A. Rinaldo, and A. Singh, “Minimax localization of structural information in large noisy matrices,” in *Advances in Neural Information Processing Systems*, 2011, pp. 909–917.
- [81] A. Ben-Dor, B. Chor, R. Karp, and Z. Yakhini, “Discovering local structure in gene expression data: The order-preserving submatrix problem,” *Journal of computational biology*, vol. 10, no. 3-4, pp. 373–384, 2003.
- [82] S. C. Brubaker and S. S. Vempala, “Random tensors and planted cliques,” in *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, Springer, 2009, pp. 406–419.
- [83] A. Anandkumar, R. Ge, D. Hsu, and S. Kakade, “A tensor spectral approach to learning mixed membership community models,” in *Conference on Learning Theory*, 2013, pp. 867–881.
- [84] L. Gauvin, A. Panisson, and C. Cattuto, “Detecting the community structure and activity patterns of temporal networks: A non-negative tensor factorization approach,” *PloS one*, vol. 9, no. 1, e86028, 2014.
- [85] J.-P. Aubin and I. Ekeland, *Applied nonlinear analysis*. Courier Corporation, 2006.
- [86] V. Koltchinskii, “Oracle inequalities in empirical risk minimization and sparse recovery problems,” 2011.
- [87] R. Ahlswede and A. Winter, “Strong converse for identification via quantum channels,” *IEEE Transactions on Information Theory*, vol. 48, no. 3, pp. 569–579, 2002.
- [88] E. H. Lieb, “Convex trace functions and the wigner-yanase-dyson conjecture,” *Advances in Mathematics*, vol. 11, no. 3, pp. 267–288, 1973.

- [89] M. Talagrand, “New concentration inequalities in product spaces,” *Inventiones mathematicae*, vol. 126, no. 3, pp. 505–563, 1996.
- [90] O. Bousquet, “A bennett concentration inequality and its application to suprema of empirical processes,” *Comptes Rendus Mathematique*, vol. 334, no. 6, pp. 495–500, 2002.
- [91] V. Koltchinskii and K. Lounici, “Concentration inequalities and moment bounds for sample covariance operators,” *Bernoulli*, vol. 23, no. 1, pp. 110–133, 2017.
- [92] R. Vershynin, “Introduction to the non-asymptotic analysis of random matrices,” *ArXiv preprint arXiv:1011.3027*, 2010.