STATISTICAL METHODS FOR HIGH DIMENSIONAL COUNT AND COMPOSITIONAL DATA
WITH APPLICATIONS TO MICROBIOME STUDIES

Yuanpei Cao

A DISSERTATION

in

Applied Mathematics and Computational Science

Presented to the Faculties of the University of Pennsylvania

in

Partial Fulfillment of the Requirements for the

Degree of Doctor of Philosophy

2016

Supervisor of Dissertation

_____

Hongzhe Li

Professor of Biostatistics

Graduate Group Chairperson

_____

Charles L. Epstein, Thomas A. Scott Professor of Mathematics

Dissertation Committee

Hongzhe Li, Professor of Biostatistics

T. Tony Cai, Dorothy Silberberg Professor and Professor of Statistics

Zongming Ma, Associate Professor of Statistics

STATISTICAL METHODS FOR HIGH DIMENSIONAL COUNT AND COMPOSITIONAL DATA

WITH APPLICATIONS TO MICROBIOME STUDIES

© COPYRIGHT

2016

Yuanpei Cao

# ACKNOWLEDGEMENT

I would like to first and foremost express my deepest gratitude to my advisor Professor Hongzhe Li for the continuous support of my Ph.D study and research. With his extensive knowledge, sharp thinking and enthusiasm about problems in biostatistics, he provided me with an excellent research atmosphere and patiently guided me through a lot of difficulties during my Ph.D study.

I would also like to thank my thesis committee members, Professor Tony Cai and Professor Zongming Ma. They gave me tremendous help and great suggestions on my scientific research and presentation skills.

I would like to express my gratitude to my collaborators, Professor Wei Lin and Professor Anru Zhang, for their stimulating discussions and inspiring comments.

My sincere thanks also go to Professor Charles Epstein, the graduate chair of Applied Mathematics and Computational Science, for offering me this wonderful opportunity to pursue graduate studies at Penn. I also want to thank my fellow graduate students and program coordinators for their continuous support.

Last but not the least, I am deeply grateful to my wife Xin Feng as well as my parents for their love and support all the time. Whenever I met with difficulties, they always gave me unconditional love and support.

# ABSTRACT

## STATISTICAL METHODS FOR HIGH DIMENSIONAL COUNT AND COMPOSITIONAL DATA WITH APPLICATIONS TO MICROBIOME STUDIES

Yuanpei Cao

Hongzhe Li

Next generation sequencing (NGS) technologies make the studies of microbiomes in very large-scale possible without cultivation *in vitro*. One approach to sequencing-based microbiome studies is to sequence specific genes (often the 16S rRNA gene) to produce a profile of diversity of bacterial taxa. Alternatively, the NGS-based sequencing strategy, also called shotgun metagenomics, provides further insights at the molecular level, such as species/strain quantification, gene function analysis and association studies. Such studies generate large-scale high-dimensional count and compositional data, which are the focus of this dissertation.

In microbiome studies, the taxa composition is often estimated based on the sparse counts of sequencing reads in order to account for the large variability in the total number of reads. The first part of this thesis deals with the problem of estimating the bacterial composition based on sparse count data, where a penalized likelihood of a multinomial model is proposed to estimate the composition by regularizing the nuclear norm of the compositional matrix. Under the assumption that the observed composition is approximately low rank, a nearly optimal theoretical upper bound of the estimation error under the Kullback-Leibler divergence and the Frobenius norm is obtained. Simulation studies demonstrate that the penalized likelihood-based estimator outperforms the commonly used naive estimator in term of the estimation error of the composition matrix and various bacterial diversity measures. An analysis of a microbiome dataset is used to illustrate the methods.

Understanding the dependence structure among microbial taxa within a community, including co-occurrence and co-exclusion relationships between microbial taxa, is another important problem in microbiome research. However, the compositional nature of the data complicates the investigation of the dependency structure since there are no known multivariate distributions that are flexible enough to model such a dependency. The second part of the thesis develops a composition-adjusted thresholding (COAT) method to estimate the sparse covariance matrix of the latent log-

basis components. The method is based on a decomposition of the variation matrix into a rank-2 component and a sparse component. The resulting procedure can be viewed as thresholding the sample centered log-ratio covariance matrix and hence is scalable to large covariance matrice estimations based on compositional data. The issue of the identifiability problem of the covariance parameters is rigorously characterized. In addition, rate of convergence under the spectral norm is derived and the procedure is shown to have theoretical guarantee on support recovery under certain assumptions. In the application to gut microbiome data, the COAT method leads to more stable and biologically more interpretable results when comparing the dependence structures of lean and obese microbiomes.

The third part of the thesis considers the two-sample testing problem for high-dimensional compositional data and formulates a testable hypothesis of compositional equivalence for the means of two latent log-basis vectors. A test for such a compositional equivalence through the centered log-ratio transformation of the compositions is proposed and is shown to have an asymptotic extreme value of type 1 distribution under the null. The power of the test against sparse alternatives is derived. Simulations demonstrate that the proposed tests can be significantly more powerful than existing tests that are applied to the raw and log-transformed compositional data. The usefulness of the proposed tests is illustrated by applications to test for differences in gut microbiome composition between lean and obese individuals and changes of gut microbiome between different time points during treatment in Crohn's disease patients.

## TABLE OF CONTENTS

# LIST OF TABLES

# CHAPTER 1

## COMPOSITION ESTIMATION FROM SPARSE COUNT DATA VIA A REGULARIZED LIKELIHOOD

In microbiome studies, taxa composition is often estimated based on the sequencing read counts in order to account for the large variability in the total number of observed reads across different samples. Due to sequencing depth, some rare microbial taxa might not be captured in the metagenomic sequencing, which results in many zero read counts. Naive composition estimation using count normalization therefore lead to many zero proportions, which underestimates the underlying compositions, especially for the rare taxa. Such an estimate of the composition can further lead to biased estimate of taxa diversity, and can also cause difficulty in downstream log-ratio based analysis for compositional data. In this paper, the observed counts are assumed to be sampled from a multinomial distribution, with the unknown composition being the probability parameter in a high dimensional positive simplex space. Under the assumption that the composition matrix is approximately low rank, a nuclear norm regularization-based likelihood estimation is developed to estimate underlying compositions of the samples. The theoretical upper bounds and the min-max lower bounds of the estimation errors measured by the Kullback-Leibler divergence and the Frobenius norm are established. Simulation studies demonstrate that the regularized maximum likelihood estimator outperforms the commonly used naive estimators. The methods are applied to an analysis of a human gut microbiome dataset.

## 1.1. Introduction

The human microbiome is the totality of all microbes at different body sites, whose contribution to human health and disease has increasingly been recognized. Recent studies have demonstrated that the microbiome composition varies across individuals due to different health and the environment status (The Human Microbiome Project Consortium, 2012a), and may be associated with complex diseases such as obesity, atherosclerosis, and Crohn's disease (Koeth et al., 2013; Lewis et al., 2015; Turnbaugh et al., 2009). With the development of next-generation sequencing technologies, the human microbiome organisms can be quantified by using direct DNA sequencing of either marker genes or the whole metagenomes. After aligning the sequence reads to the refer-

ence microbial genomes, the observed count data (e.g., 16S rRNA marker gene reads or shotgun metagenomic reads) depend on the amount of genetic material extracted from the community or the sequencing depth, and they provide a relative measure of the abundances of community components. In a microbiome study, these read counts are typically non-negative and over-dispersed, and contain a large number of zeros.

In order to account for the large variability in the total number of reads obtained, the taxa composition is often estimated based on the observed counts of sequencing reads. Due to sequencing depth, some rare microbial taxa might not be captured in the metagenomic sequencing, which results in zero read counts assigned to these taxa. Naive estimates of the taxa composition using count normalization therefore lead to many zeros due to under sampling, especially for rare taxa. Such a naive estimate of the composition can be biased and can lead to biased estimates of taxa diversity. It can also cause difficulty in downstream data analysis for compositional data. Since the pioneering work of Aitchison, (2003), several techniques have been proposed to deal with zeros (see Martın-Fernandez, Palarea-Albaladejo, and Olea, 2011 for an overview) in count data. One approach is to replace zero counts through a Bayesian-multiplicative model, followed by normalizing the count into the composition. Such a Bayesian method involves a Dirichlet prior distribution as the conjugate distribution of multinomial distribution and a multiplicative modification of the non-zero counts. The zero replacement results were determined by the parameterizations of the prior distribution. However, such a prior information cannot be easily obtained, and the subjective selection of the parameter may yield misleading results. Other approaches normalized the count first and treated zero compositions as the missing values. The missing part was then recovered by either non-parametric imputation or EM algorithms. However, the non-parametric imputation lacks theoretical guarantees for selecting a reasonable replacement value. The EM algorithm is not feasible when the number of taxa is very large, or every taxa contains at least one zero across the samples. In addition, the multivariate additive log-ratio (alr) normality assumption used in these methods is often violated in microbiome studies.

This paper addresses the problem of estimating the microbial compositions in positive simplex space from a high-dimensional sparse count dataset. We assume that the observed counts follow a Poisson-multinomial model where the read counts of the taxa in each individual follow a multinomial distribution with the underlying probability parameter given by a positive composition, and the

number of total count is a Poisson random variable. If the compositions across different individuals are treated as a matrix by combining them together, an approximately low rank structure on this matrix is indicated by recent observations on co-occurrence pattern (Faust et al., 2012) and various symbiotic relationships in microbial communities (Chaffron et al., 2010; Horner-Devine et al., 2007; Woyke et al., 2006). Motivated by much success in solving the matrix completion problem using nuclear norm minimization (Cao and Xie, 2016; Klopp et al., 2015; Lafond et al., 2014; Lu and Negahban, 2014; Negahban and Wainwright, 2012), this paper solves the problem of the composition estimation using a regularized maximum likelihood approach. However, it should be emphasized that the multinomial likelihood function in this framework has not been studied and the sampling scheme used in this article is also different from other matrix completion problems. The observed zero counts are the result of under sampling, rather than the random missingness assumed in the previous literature. We provide the asymptotic upper and min-max lower bounds of the resulting regularized estimator and show through simulations that the estimator recovers low-rank compositions accurately.

The rest of the paper is organized as follows. Section $1.2$ presents details of the proposed regularized likelihood approach when the underlying composition is approximately low-rank. The implementation is presented in Section 1.3. The theoretical properties of the estimators are analyzed in Section $1.4$, where the upper bounds for the estimation error measured by average Kullback-Leibler divergence and Frobenius norm are established. Simulation results are shown in Section $1.5$ to investigate the numerical performance of the proposed methods. A real data application to a human gut microbiome study is given in Section $1.6$.

## 1.2. Poisson-Multinomial Model for Microbiome Count Data and Penalized Estimation

In this section, we consider a Poisson-multinomial statistical model for composition estimation from the sparse count data observed in microbiome studies. The proposed procedure for composition estimation relies on a regularized maximum likelihood. We start by introducing some notation that will be used throughout the rest of the paper. For any integers $N > 0$, let $[N] = \{1, 2, \cdots, N\}$ be the set of integers ranging from 1 to $N$. We also denote $\mathbf{1}_n = (1, \ldots, 1)^\top \in \mathbb{R}^n$, $\mathbf{e}_i$ as the canonical basis with $i$-th entry one and others zero. For any vector $u \in \mathbb{R}^p$, we refer to it as a composition

vector if $u \geq 0$ and $\sum_{i=1}^{p} u_i = 1$. For any two composition vectors $u, v \in \mathbb{R}^p$, we can define the Kullback-Leibler (KL) divergence as

$$\mathrm{D}_{KL}(u, v) = \sum_{j=1}^{p} u_i \log \frac{u_i}{v_i}. \tag{1.1}$$

For any matrix $\mathbf{A} = (a_{ij})$, define its $L_1$, $L_\infty$, spectral, Frobenius, element-wise maximum, and nuclear norm respectively as $\|\mathbf{A}\|_1$, $\|\mathbf{A}\|_\infty$, $\|\mathbf{A}\|_2$, $\|\mathbf{A}\|_F$, $\|\mathbf{A}\|_{\max}$, and $\|\mathbf{A}\|_*$. Specifically, $\|\mathbf{A}\|_1 = \max_j \sum_i |a_{ij}|$, $\|\mathbf{A}\|_\infty = \max_i \sum_j |a_{ij}|$, $\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}$, $\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} a_{ij}^2}$, $\|\mathbf{A}\|_{\max} = \max_{i,j} |a_{ij}|$, and $\|\mathbf{A}\|_* = \sum_i \sigma_i(\mathbf{A})$, where $\lambda_{\max}(\cdot)$ denotes the largest eigenvalue and $\{\sigma_i(\cdot)\}$ denotes the set of singular values. For two matrices $\mathbf{A}$ and $\mathbf{B}$, let $\langle \mathbf{A}, \mathbf{B} \rangle = \mathrm{tr}(\mathbf{A}^T \mathbf{B}) = \sum_{i,j} a_{ij} b_{ij}$ be the trace inner product. Finally, for notational simplicity, we use $C_1, C_2, \ldots$ as generic symbols for constants whose values may vary from line to line.

Our starting point is a $n \times p$ matrix of counts $\mathbf{W}$ with element $W_{ij}$ representing the observed read count of taxon $j$ in individual $i$, where $i \in [n]$ and $j \in [p]$. For $i$-th individual, the simplest model for their count data $\mathbf{W}_i = (W_{i1}, W_{i2}, \cdots, W_{ip})$ is the multinomial model with its probability function given as

$$f_M(W_{i1}, W_{i2}, \cdots, W_{ip}; \mathbf{X}_i) = \binom{N_i}{\mathbf{W}_i} \prod_{j=1}^{p} X_{ij}^{* W_{ij}},$$

where $N_i = \sum_{j=1}^{p} W_{ij}$ and $\mathbf{X}_i^* = (X_{i1}^*, X_{i2}^*, \cdots, X_{ip}^*)$ are underlying bacterial composition with $\sum_{j=1}^{p} X_{ij}^* = 1, X_{ij}^* > 0$. The total taxa count $N_i$ is determined by the sequencing depth and can be treated as a Poisson random variable given by $N_i \sim \mathrm{Pois}(\nu_i)$, where $\nu_i$ is a positive parameter, but it is of less interest.

Our goal is to estimate $\mathbf{X}^* = (\mathbf{X}_1^{*T}, \mathbf{X}_2^{*T}, \cdots, \mathbf{X}_p^{*T})^T$ based on $\mathbf{W}$. The most natural estimate is obtained by the maximum likelihood estimation. Denote by $\mathcal{L}_N$ the (normalized) negative log-likelihood of the observations, ignoring the terms that do not depend on the compositions $\mathbf{X}^*$,

$$\mathcal{L}_N(\mathbf{X}) = -\frac{1}{N} \sum_{1 \leq i \leq n, 1 \leq j \leq p} W_{ij} \log X_{ij}^*, \tag{1.2}$$

where $N = \sum_{i=1}^{n} N_i = \sum_{ij} W_{ij}$ is the total number of the observed counts and $\mathbf{X}^*$ belongs to

4

positive simplex space $\mathcal{S} = \{\mathbf{X} \in \mathbb{R}^{n \times p} \mid \mathbf{X}\mathbf{1}_p = \mathbf{1}_n, \mathbf{X} > \mathbf{0}\}$. Without further constraints, minimizing (1.2) leads to the standard maximum likelihood estimate $\widehat{\mathbf{X}}$,

$$\widehat{X}_{ij} = \frac{W_{ij}}{\sum_{k=1}^{p} W_{ik}}, \quad i \in [n], \quad j \in [p].$$

However, as a consequence of under sampling when $N$ is not sufficiently large, the estimator $\widehat{\mathbf{X}}$ will contain a large number of zeros. These zeros underestimate the composition and cause difficulty in downstream log-ratio based compositional data analysis (Aitchison, 2003). For an arbitrary matrix $\mathbf{X}^*$ in positive simplex space, clearly there is no good way to recover a positive $\mathbf{X}^*$. However, in the metagenomic study, $\mathbf{X}^*$ could be approximately low-rank in the sense that the singular values decay gradually towards zero, which provides the possibility to recover $\mathbf{X}^*$ with high accuracy. In this paper, we propose a penalized estimator $\widehat{\mathbf{X}}$ based on a regularized maximum likelihood formulation:

$$\widehat{\mathbf{X}} = \arg \min_{\mathbf{X} \in \mathcal{S}(\alpha_x, \beta_x)} \mathcal{L}_N(\mathbf{X}) + \lambda \|\mathbf{X}\|_*, \tag{1.3}$$

where $\mathcal{S}(\alpha_x, \beta_x)$ is a bounded simplex space given by

$$\mathcal{S}(\alpha_x, \beta_x) = \left\{ \mathbf{X} \in \mathbb{R}^{n \times p} \mid \mathbf{X}\mathbf{1}_p = \mathbf{1}_n, \alpha_x/p \le X_{ij} \le \beta_x/p, \forall (i,j) \in [n] \times [p] \right\}.$$

Here $\lambda$ and $\alpha_x$ and $\beta_x$ are tuning parameters. The constrained element-wise lower bound guarantees the positive sign of the estimator. The element-wise upper bound constraint is only needed in the theory, while in practice, such a constraint is not required.

## 1.3. Optimization Algorithm and Tuning Parameter Selection

In this section we consider the implementation of the proposed estimator specified as (1.3). Specifically, we propose to solve the following constrained convex optimization:

$$\widehat{\mathbf{X}} = \arg \min_{\mathbf{X} \in \mathcal{S}(\alpha_X)} \mathcal{L}_N(\mathbf{X}) + \lambda \|\mathbf{X}\|_*, \tag{1.4}$$

$$\mathcal{S}(\alpha_X) = \left\{ \mathbf{X} \in \mathbb{R}^{n \times p} \mid \mathbf{X}\mathbf{1}_p = \mathbf{1}_n, X_{ij} \ge \alpha_X/p, \forall (i,j) \in [n] \times [p] \right\}.$$

Here $\mathcal{S}(\alpha_X)$ is a positive simplex space and $(\lambda, \alpha_X)$ is a pair of tuning parameters. Particularly, (1.4) is a nuclear norm minimization problem, which can be solved by either semidefinite programing via interior-point SDP solver, or first-order method via Templates for First-Order Conic Solvers (TFOC-S), see Becker, Candès, and Grant, 2011. However, the SDP solver computes the nuclear norm via a less efficient eigenvalue decomposition which does not scale well with high-dimensions $n$ and $p$. Besides, Nesterov's scheme used in TFOCS is not monotone in the objective function owing to the introduction of the momentum term, which often results in oscillations or overshoots along the trajectory of the iteration. In this article, we propose a more efficient algorithm based on the generalized accelerated proximal gradient method (Su, Boyd, and Candes, 2014). To adapt to the bounded simplex constraint $\mathcal{S}(\alpha_X)$, we develop a non-iterative projection scheme in the proposed algorithm.

### 1.3.1. Generalized Accelerated Proximal Gradient Method

Since $\mathcal{L}_N(\cdot)$ is convex and differentiable over the domain $\mathcal{S}(\alpha_X)$ and the nuclear norm is convex, the accelerated Nesterov's scheme can be formulated as follows. Given the count matrix $\mathbf{W}$, we first normalize it into the composition $\mathbf{X}$ by $X_{ij} = W_{ij} / \sum_{k=1}^{p} W_{ik}$ and initialize $\mathbf{Y}_0 = \mathbf{X}_0 = \mathbf{X}_{-1} = \mathbf{X}$, then update $\mathbf{X}_k$ and $\mathbf{Y}_k$ in the $k$th iteration as

$$\mathbf{X}_k = \underset{\mathbf{X} \in \mathcal{S}(\alpha_X)}{\arg\min} \frac{L_k}{2} \|\mathbf{X} - \mathbf{Y}_{k-1} + L_k^{-1} \nabla \mathcal{L}_N (\mathbf{Y}_{k-1})\|_F^2 + \lambda \|\mathbf{X}\|_*, \tag{1.5}$$

$$\mathbf{Y}_k = \mathbf{X}_k + \frac{k-1}{k+r-1} (\mathbf{X}_k - \mathbf{X}_{k-1}). \tag{1.6}$$

Here we provide the detailed explanation for (1.5) and (1.6).

- $L_k$ is the step size in the $k$-th iteration, which is chosen by line search strategy. Denote by $\mathcal{F}_L(\mathbf{X}, \mathbf{Y})$ the approximation error when approximating $\mathcal{L}_N(\mathbf{X})$ with its second order Taylor expansion around $\mathbf{Y}$ and using $L$ as the second order coefficient,

$$\mathcal{F}_L(\mathbf{X}, \mathbf{Y}) = \mathcal{L}_N(\mathbf{X}) - \mathcal{L}_N(\mathbf{Y}) - \langle \mathbf{X} - \mathbf{Y}, \nabla \mathcal{L}_N(\mathbf{Y}) \rangle - \frac{L}{2} \|\mathbf{X} - \mathbf{Y}\|_F^2.$$

In the $k$th iteration, given an initial parameter $L_k = L_{k-1}$, we repeated increasing it by $L_k = \gamma L_k$ for some scale parameter $\gamma > 1$ until the function $\mathcal{L}_N(\mathbf{X}_k)$ is dominated by its second order Taylor expansion around $\mathbf{Y}_{k-1}$, i.e., $\mathcal{F}_{L_k}(\mathbf{X}_k, \mathbf{Y}_{k-1}) \leq 0$.

6

- $\frac{k-1}{k+r-1}$ is the momentum term and $r$ is a friction parameter. In the standard accelerated gradient method, the friction parameter is set by $r = 3$, and this scheme exhibits the convergence rate $O(1/k^2)$ as long as the gradient function $\nabla \mathcal{L}_N$ is Lipschitz continuous with a constant Lipschitz coefficient (Nesterov, 1983, 2013). The Nesterov's scheme can be further generalized by setting a high friction rate, for example $r \geq 9/2$, and it succeeds in eliminating the overshooting and oscillation along the trajectory toward the minimizer and obtaining a $O(1/k^3)$ convergence rate (Su, Boyd, and Candes, 2014).

- The minimization of the objective function (1.5) can be solved by a form of Singular Value Thresholding (SVT) (Cai, Candès, and Shen, 2010):

$$\mathbf{X}_k = \Pi_{\mathcal{S}(\alpha_X)} \left( \mathcal{D}_{\lambda L_k^{-1}} \left( \mathbf{Y}_{k-1} - L_k^{-1} \nabla \mathcal{L}_N \left( \mathbf{Y}_{k-1} \right) \right) \right).$$

Here $\Pi_{\mathcal{S}(\alpha_X)} \left( \mathbf{X} \right)$ is Euclidean projection of $\mathbf{X}$ onto the positive simplex space $\mathcal{S}(\alpha_X)$ that we will discuss in Section 1.3.2. If $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ is the singular value decomposition (SVD), the soft-thresholding operator $\mathcal{D}_\tau$ can be defined as

$$\mathcal{D}_\tau \left( \mathbf{X} \right) = \mathbf{U} \mathcal{D}_\tau \left( \mathbf{\Sigma} \right) \mathbf{V}^T, \quad \mathcal{D}_\tau \left( \mathbf{\Sigma} \right) = \operatorname{diag} \left( \max \left\{ \sigma_i - \tau, 0 \right\} \right).$$

Combining these steps together, the generalized accelerated proximal gradient method is summarized in Algorithm 1, where $k_{\max}$ is the maximum number of iteration. The complexity of the algorithm are dominated by $\mathcal{O}(n^2 p + p^3)$, which is the cost of singular value decomposition. The convergence of Algorithm 1 cannot be easily established; however, the following proposition provides some insight.

**Proposition 1.** *Let $\mathbf{X}_k$ be the sequencing generated in the iteration of Algorithm 1. Denote by $f(\mathbf{X}) = \mathcal{L}_N(\mathbf{X}) + \lambda \|\mathbf{X}\|_*$. Suppose the Euclidean projection onto the simplex space $\Pi_{\mathcal{S}(\alpha_X)}$ does not influence the convergence rate, and the step size is always set by $L_k = \max\limits_{ij} W_{ij} p/(\alpha_X N)$. Then, for any friction parameter $r \geq 9/2$, we have,*

$$f(\mathbf{X}_k) - f(\mathbf{X}^\star) \leq C \sqrt{\frac{\max\limits_{ij} W_{ij}^3}{\min\limits_{\{ij | W_{ij} > 0\}} W_{ij}} \frac{p^3}{N^2} \frac{\|\mathbf{X}_0 - \mathbf{X}^\star\|_F^2}{k^3}},$$

*where $\mathbf{X}^\star$ is any minimizer of $f$ and $C$ only depends on $r$ and $\alpha_X$.*

Since the gradient function $\triangledown\mathcal{L}_N$ is Lipschitz continuous with the constant $L = \max\limits_{ij} W_{ij} p/(\alpha_X N)$ and the negative likelihood function $\mathcal{L}_N$ is $\mu-$strongly convex with $\mu = \min\limits_{\{ij|W_{ij}>0\}} W_{ij}/N$ on the constrained simplex space, it is not hard to prove Proposition 1 by applying Theorem 9 in Su, Boyd, and Candes, 2014. The parameters $L$ and $\mu$ vary with different observations, as a result, the rate of convergence shows an interesting dependency on the dimension $p$ and the observation count $\mathbf{W}$

---

**Algorithm 1** Generalized accelerated proximal gradient method

---

1: **Input**: Count $\mathbf{W}$ and its normalized composition $\mathbf{X}$
2: Initialize: $\mathbf{Y}_0 = \mathbf{X}_0 = \mathbf{X}_{-1} = \mathbf{X}$, $r \geq 9/2$, $\gamma > 1$, $L = 10^{-4}$, and $k_{\max} \in \mathbb{N}^+$
3: **for** $k = 1$ to $k_{\max}$ **do**
4:     $\mathbf{X}_k = \Pi_{\mathcal{S}(\alpha_X)}\left(\mathcal{D}_{\lambda/L}\left(\mathbf{Y}_{k-1} - (1/L)\triangledown\mathcal{L}_N\left(\mathbf{Y}_{k-1}\right)\right)\right)$
5:     **if** $\mathcal{F}_L(\mathbf{X}_k, \mathbf{Y}_{k-1}) \geq 0$, **then**
6:         $L = \gamma L$, go to Step 3
7:     **end if**
8:     Update $\mathbf{Y}_k = \mathbf{X}_k + \frac{k-1}{k+r-1}\left(\mathbf{X}_k - \mathbf{X}_{k-1}\right)$
9:     **if** $|\mathcal{F}_L(\mathbf{X}_k, \mathbf{Y}_{k-1})| < 10^{-5}$ **then**
10:         **return** $\mathbf{X}_k$
11:     **end if**
12: **end for**

---

### 1.3.2. Euclidean Projection onto the Simplex Space

The remaining part is to deal with the Euclidean projection onto the simplex space $\mathcal{S}(\alpha_X)$ in Algorithm 1. We introduce a non-iterative and efficient algorithm based on the standard KKT condition. Consider a one-dimensional simplex projection problem given by

$$\Pi_{\mathcal{S}(\alpha_X)}(\mathbf{y}) = \min_{\mathbf{x}\in\mathbb{R}^p} \quad \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|_2^2 \quad s.t. \quad \sum_{i=1}^{p} x_i = 1, \quad x_i \geq \alpha_X/p. \tag{1.7}$$

The following Proposition provides an implicit formulation for the minimizer $\mathbf{x}^\star$ to this optimization problem (1.7).

**Proposition 2.** *Suppose that $y_1 \geq y_2 \geq \cdots \geq y_p$, then the minimizer $\mathbf{x}^\star = (x_1, x_2, \cdots, x_p)^T$ is given by*

$$x_i = \max\{y_i + \mu, \alpha_x/p\}, \text{ for any } i \in [p],$$

*where $\mu = \rho^{-1}(1 - \alpha_X - \sum_{i=1}^{\rho} u_i) + \alpha_X/p$, and $\rho$ is the number of components in $\mathbf{x}^\star$ that are strictly*

*larger than $\alpha_X/p$. We establish the the following formulation for $\rho$,*

$$\rho = \max\left\{j \in [p] \ \middle| \ y_j + j^{-1}(1 - \alpha_X - \sum_{i=1}^{j} y_i) > 0\right\}.$$

In the multi-dimensional case that $\mathbf{Y} \in \mathbb{R}^{n \times p}$, we generalize the above simplex projection and summarize this non-iterative optimization procedure in Algorithm 2. The scheme is easy to implement and its complexity is $\mathcal{O}(np\log(p))$.

---

**Algorithm 2** Euclidean projection of a matrix onto the simplex space $\mathcal{S}(\alpha_X)$.

---

1: Input: $\mathbf{Y} \in \mathbb{R}^{n \times p}$ and $\mathcal{S}(\alpha_X)$
2: Sort each row of $\mathbf{Y}$ into $\mathbf{U}$: $U_{i1} \geq U_{i2} \cdots \geq U_{ip}, i \in [p]$.
3: Find vector $\rho = (\rho_1, \cdots, \rho_n)^T$ such that

$$\rho_i = \max\left\{j \in [p] \ \middle| \ U_{ij} + j^{-1}\left(1 - \alpha_X - \sum_{i=1}^{j} U_{ij}\right) > 0\right\}, i \in [n].$$

4: Define vector $\mu = (\mu_1, \cdots, \mu_n)^T$ by $\mu_i = \rho_i^{-1}\left(1 - \alpha_X - \sum_{j=1}^{\rho_i} U_{ij}\right) + \alpha_X/p, i \in [n]$.
5: Return $\mathbf{X}$ such that $X_{ij} = \max\left\{Y_{ij} + \mu_i, \alpha_X/p\right\}, (i,j) \in [n] \times [p]$.

---

### 1.3.3. Data Driven Selection of the Tuning Parameters

The proposed nuclear norm minimization involves the tuning parameters $\lambda$ and $\alpha_X$. We propose the following data-driven method for selecting these tunning parameters with a guaranteed performance. Given a selected parameter $\alpha_X$, we choose $\lambda = \lambda(\alpha_X, \widehat{\beta}_R)$ by plugging $\alpha_X$ and the estimated row probability parameter

$$\widehat{\beta}_R = n \cdot \max_{1 \leq i \leq n} \frac{\sum_{j=1}^{p} W_{ij}}{\sum_{k=1}^{n} \sum_{l=1}^{p} W_{kl}}$$

$$\lambda(\alpha_X, \widehat{\beta}_R) = \sqrt{\frac{32\left(\widehat{\beta}_R^2/n + (1 \vee \widehat{\beta}_R p/n)/\alpha_X\right) p \log(n+p)}{N}} \vee \frac{8(1/\alpha_X + \widehat{\beta}_R/(np)^{1/2})n \log(n+p)}{N}.$$

This choice of $\lambda$ is motivated by the theoretical results of Theorem 1 in the next Section.

It remains to find the estimated parameter $\alpha_X$, which can be selected using $K$-fold cross-validation as follows. Let $\mathbf{W}$ be the observed sample and let $T$ be a grid of positive real values. For each

$t \in T$, set

$$(\lambda, \alpha_X) = (\lambda(\alpha_X(t), \widehat{\beta}_R), \alpha_X(t)) = (\lambda((t \cdot \widehat{\alpha}_X), \widehat{\beta}_R), t \cdot \widehat{\alpha}_X),$$

where

$$\widehat{\alpha}_X = p \cdot \min_{1 \le i \le n, 1 \le j \le p} \frac{W_{ij}}{\sum_{k=1}^{n} \sum_{l=1}^{p} W_{kl}} \text{ and } \widehat{\beta}_R = n \cdot \max_{1 \le i \le n} \frac{\sum_{j=1}^{p} W_{ij}}{\sum_{k=1}^{n} \sum_{l=1}^{p} W_{kl}}.$$

We randomly split the rows of $\mathbf{W}$ into two groups of sizes $n_1 \sim \frac{(K-1)n}{K}$ and $n_2 \sim \frac{n}{K}$ for $I$ times. We used the second group with sample size $n_2$ as the testing set. In order to estimate the composition from the rows in testing set, we further randomly picked $1/K$ proportion of observed columns in each row from the second group and combined it with the first group as the training set. Denote by $\mathbf{W}^i$ be the selected testing set in the $i$th split and let $\mathbf{X}^i$ be its compositions through $X^i_{kl} = W^i_{kl}/\sum_{l=1}^{p} W^i_{kl}$. Denote by $\widehat{\mathbf{X}}^{(-i)}(\alpha_X(t))$ the estimator based on the training set. We consider the Kullback-Leibler divergence to evaluate the prediction error.

$$\widehat{R}(t) = \sum_{i=1}^{I} \mathrm{D}(\mathbf{X}^i, \widehat{\mathbf{X}}^{(-i)}(\alpha_X(t))).$$

We select $t^* = \arg\min_T \widehat{R}(t)$ and choose the tuning parameters $(\lambda(\alpha_X(t^*), \widehat{\beta}_R), \alpha_X(t^*))$. If $t^*$ is chosen on the boundary of $T$, we expand the range of $T$ and repeat the above procedure. With the chosen tuning parameters, we finally obtain estimate by solving (1.4) based on the full dataset.

## 1.4. Theoretical Properties

We prove that the proposed estimator $\widehat{\mathbf{X}}$ achieves the near optimal rate of convergence over a class of low-rank compositions. The regularization assumptions we need for theoretical analysis are formally stated as below.

**Condition 1.** Let $R_i = \nu_i / \sum_{j=1}^{p} \nu_j$ for $i \in [n]$, then there exist constants $(\alpha_R, \beta_R)$ such that, for any $i \in [n]$,

$$\alpha_R/n \le R_i \le \beta_R/n.$$

**Condition 2.** There exist constants $(\alpha_X, \beta_X)$ such that, for any $(i, j) \in [n] \times [p]$,

$$\alpha_X/p \le X_{ij} \le \beta_X/p.$$

Here $\mathbf{R} = (R_1, \cdots, R_n)^T$ represents the probability of observing an element from each row, and $\mathbf{X}$ represents the column probability. Conditions $1$ and $2$ are analogous to the incoherence conditions that are commonly assumed in the matrix completion literature. The element-wise upper bounds avoid the overly "spiky" situation that some rows or columns are sampled with very high probability. The element-wise lower bound on $\mathbf{R}$ helps to establish bounds in Frobenius norm, and the entry-wise bound on $\mathbf{X}^*$ ensure the gradient function of $\mathcal{L}_N(\mathbf{X})$ in (1.2) is Lipschitz continuous, which helps to effectively bound Frobenius norm in terms of Kullback-Leibler (KL) divergence and guarantee the feasibility of accelerated gradient descent algorithm in practice.

*1.4.1. Rate of Convergence*

To assess how close the estimator $\widehat{\mathbf{X}}$ from (1.3) to the real compositional matrix $\mathbf{X}^*$, we use average Kullback-Leibler divergence $\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}})$ and squared Frobenius norm $\|\mathbf{X}^* - \widehat{\mathbf{X}}\|_F^2$. Here $\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}})$ is defined as the sum of Kullback-Leibler (KL) divergence between rows of $\mathbf{X}^*$ and $\widehat{\mathbf{X}}$,

$$\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) = \sum_{i=1}^{n} \mathrm{D}_{KL}(\mathbf{X}_i^*, \widehat{\mathbf{X}}_i) = \sum_{i=1}^{n} \sum_{j=1}^{p} X_{ij}^* \log \frac{X_{ij}^*}{\widehat{X}_{ij}}.$$

The following theorem gives an upper bound on the loss of the proposed estimator $\mathbf{X}^*$ for the exactly low-rank composition matrix $\mathbf{X}$.

**Theorem 1.** *(Exactly low-rank matrices) Under Conditions $1$ and $2$, suppose that $N \geq c_0(n \vee p) \log(n + p)$ for some universal constant $c_0 > 0$, and the tuning parameter is selected as*

$$\lambda = 2 \left( \sqrt{\frac{C_1(n,p)p\log(n+p)}{N}} \vee \frac{C_2(n,p)p\log(n+p)}{N} \right), \tag{1.8}$$

*where $C_1(n,p) = 8\left(\beta_R^2/n + (1 \vee \beta_R p/n)/\alpha_X\right)$ and $C_2(n,p) = 4(1/\alpha_X + \beta_R/(np)^{1/2})$. If the composition $\mathbf{X}^*$ has rank at most $r$, then, with probability at least $1 - 3(n+p)^{-1}$, the estimate $\widehat{\mathbf{X}}$ in (1.3) satisfies*

$$\frac{1}{n}\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) \leq C_1 \left( \frac{(p+n)r\log(n+p)}{N} \right), \tag{1.9}$$

$$\frac{p}{n}\|\widehat{\mathbf{X}} - \mathbf{X}^*\|_F^2 \leq C_2 \left( \frac{(p+n)r\log(n+p)}{N} \right), \tag{1.10}$$

*for some constants $C_1$ and $C_2$ which only depend on $c_0, \alpha_X, \beta_X, \alpha_R$ and $\beta_R$.*

Theorem 1 states the rate of convergence for both KL divergence and Frobenius loss in terms of probability. With some additional mild assumptions, the same rate of convergence holds in expectation.

**Corollary 1.** Under the same conditions mentioned in Theorem 1, if $N$ further satisfies $N \leq c_1(n + p)^2 r \log(n + p)$, then, there exists some constants $C_1$ and $C_2$ only depending on $c_0, c_1, \alpha_X, \beta_X, \alpha_r$ and $\beta_R$, such that

$$\frac{1}{n}\mathbb{E}\,\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) \leq C_1 \frac{(p + n)r \log(n + p)}{N},$$
$$\frac{p}{n}\mathbb{E}\left\|\widehat{\mathbf{X}} - \mathbf{X}^*\right\|_F^2 \leq C_2 \frac{(p + n)r \log(n + p)}{N}.$$

We also have the corresponding lower bound that shows that the bound in Theorem 1 essentially cannot be improved.

**Theorem 2.** *Consider the matrix classes*

$$\mathbb{B}_0(r, \alpha, \beta) = \left\{\mathbf{X} \in \mathbb{R}^{n \times p} \big| rank(\mathbf{X}) \leq r, \mathbf{X}\mathbf{1}_p = \mathbf{1}_n, \alpha/p \leq X_{ij} \leq \beta/p, \text{ for any } (i, j) \in [n] \times [p]\right\}.$$

*If $2 \leq r \leq p/2$, there exists some constants $C_1$ and $C_2$ which only depend on $\alpha_X, \beta_X, \alpha_R, \beta_R$, such that*

$$\inf_{\widehat{\mathbf{X}}} \sup_{\substack{\mathbf{X}^* \in \mathbb{B}_0(r, \alpha_X, \beta_X) \\ \alpha_R \leq \mathbf{R}_i \leq \beta_R}} \frac{1}{n}\mathbb{E}\,\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) \geq C_1 \frac{(p + n)r}{N},$$
$$\inf_{\widehat{\mathbf{X}}} \sup_{\substack{\mathbf{X}^* \in \mathbb{B}_0(r, \alpha_X, \beta_X) \\ \alpha_R \leq \mathbf{R}_i \leq \beta_R}} \frac{p}{n}\mathbb{E}\left\|\widehat{\mathbf{X}} - \mathbf{X}^*\right\|_F^2 \geq C_2 \frac{(p + n)r}{N}.$$

In practice, the composition is typically approximately low-rank instead of exactly low-rank. In such case, we formalize the class of approximately low-rank matrices via the $l_q$-"ball" of matrices by

$$\mathbb{B}_q(\rho_q) = \left\{\mathbf{X} \in \mathbb{R}^{n \times p} \,\Big|\, \sum_{i=1}^{n \wedge p} |\sigma_i(\mathbf{X}^*)|^q \leq \rho_q, \right\}, \tag{1.11}$$

where $0 \leq q \leq 1$. In general, we obtain the following upper bound result.

**Theorem 3.** *(Approximately low-rank matrix): Under Conditions 1 and 2, suppose that $N \geq c_0(n \vee p) \log(n + p)$ for some constant $c_0 > 0$, and the tuning parameter is selected as (1.8). If*

*the composition $\mathbf{X}^*$ further belongs to a class of approximately low-rank matrices, Then, with the probability proceeding $1 - 3(n + p)^{-1}$, the estimator $\widehat{\mathbf{X}}$ in (1.3) satisfies*

$$\frac{1}{n} \mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) \leq C_1 \rho_q (p/n)^{\frac{q}{2}} \left( \frac{(n+p)\log(n+p)}{N} \right)^{1-\frac{q}{2}}, \tag{1.12}$$

$$\frac{p}{n} \|\widehat{\mathbf{X}} - \mathbf{X}^*\|_F^2 \leq C_2 \rho_q (p/n)^{\frac{q}{2}} \left( \frac{(n+p)\log(n+p)}{N} \right)^{1-\frac{q}{2}}, \tag{1.13}$$

*where constants $C_1$ and $C_2$ only depend on $c_0, \alpha_X, \beta_X, \alpha_R$ and $\beta_R$.*

### 1.4.2. Estimation of Diversity Index

Various microbial diversity meaures are aften used to quantify the composition of the microbial communities. Given $\mathbf{X} \in \mathbb{R}^p$ that represents $p$-bacteria composition across $n$ different individuals, three widely used measurements of microbial community diversity include

- Shannon's index $\quad \mathbf{H}_{\mathsf{sh}}(\mathbf{X}_i) = -\sum_{j=1}^p X_{ij} \log X_{ij}, 1 \leq i \leq n$;
- Simpson's index $\quad \mathbf{H}_{\mathsf{sp}}(\mathbf{X}_i) = \sum_{j=1}^p X_{ij}^2, 1 \leq i \leq n$;
- Bray-Curtis index $\quad \mathbf{H}_{\mathsf{bc}}(\mathbf{X}_i, \mathbf{X}_j) = \sum_{k=1}^p |X_{ik} - X_{jk}|/2, 1 \leq i, j \leq n$.

Here $\{\mathbf{H}_{\mathsf{sh}}(\mathbf{X}_i)\}_{i=1}^n$ and $\{\mathbf{H}_{\mathsf{sh}}(\mathbf{X}_i)\}_{i=1}^n$ are two vectors in which each component measures the richness and evenness of microbial community in an individual; $\{\mathbf{H}_{\mathsf{bc}}(\mathbf{X}_i, \mathbf{X}_j)\}_{i,j=1}^{n,p}$ is a matrix with each entry ranging from $[0, 1]$ that quantifies the dissimilarity between two individuals. Higher value of Bray-Curtis index indicates that two microbial communities are less likely to share similar taxa. The penalized likelihood estimator $\widehat{\mathbf{X}}$ from (1.3) can be used to estimate Shannon's, Simpson's and Bray-Curtis indices. The following Corollary provides the upper bound of these estimates.

**Corollary 2.** Under Conditions $1$ and $2$, suppose that $N \geq c_0(n \vee p) \log(n + p)$ for some constant $c_0$, and the tuning parameter is selected by (1.8). If the composition $\mathbf{X}^*$ has rank at most $r$, then the estimate $\widehat{\mathbf{X}}$ in (1.3) satisfies

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{H}_{\mathsf{sh}}(\widehat{\mathbf{X}}_i) - \mathbf{H}_{\mathsf{sh}}(\mathbf{X}_i^*))^2 = O_p \left( \frac{(n+p)(\log p)^2 r \log(n+p)}{N} \right),$$

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{H}_{\mathsf{sp}}(\widehat{\mathbf{X}}_i) - \mathbf{H}_{\mathsf{sp}}(\mathbf{X}_i^*))^2 = O_p \left( \frac{(n+p) r \log(n+p)}{p^2 N} \right),$$

$$\frac{1}{n^2} \sum_{1 \leq i < j \leq n} (\mathbf{H}_{\mathsf{bc}}(\widehat{\mathbf{X}}_i, \widehat{\mathbf{X}}_j) - \mathbf{H}_{\mathsf{bc}}(\mathbf{X}_i^*, \mathbf{X}_j^*))^2 = O_p \left( \frac{(n+p) r \log(n+p)}{N} \right).$$

If the composition $\mathbf{X}^*$ belongs the class of approximately low-rank matrices (1.11), then the estimate $\widehat{\mathbf{X}}$ in (1.3) satisfies

$$\frac{1}{n}\sum_{i=1}^{n}(\mathbf{H}_{\mathsf{sh}}(\widehat{\mathbf{X}}_i) - \mathbf{H}_{\mathsf{sh}}(\mathbf{X}_i^*))^2 = O_p\left(\rho_q(\log p)^2(p/n)^{\frac{q}{2}}\left(\frac{(n+p)\log(n+p)}{N}\right)^{1-\frac{q}{2}}\right),$$

$$\frac{1}{n}\sum_{i=1}^{n}(\mathbf{H}_{\mathsf{sp}}(\widehat{\mathbf{X}}_i) - \mathbf{H}_{\mathsf{sp}}(\mathbf{X}_i^*))^2 = O_p\left(\rho_q p^{\frac{q}{2}-2}/n^{\frac{q}{2}}\left(\frac{(n+p)\log(n+p)}{N}\right)^{1-\frac{q}{2}}\right),$$

$$\frac{1}{n^2}\sum_{1\le i<j\le n}(\mathbf{H}_{\mathsf{bc}}(\widehat{\mathbf{X}}_i, \widehat{\mathbf{X}}_j) - \mathbf{H}_{\mathsf{bc}}(\mathbf{X}_i^*, \mathbf{X}_j^*))^2 = O_p\left(\rho_q(p/n)^{\frac{q}{2}}\left(\frac{(n+p)\log(n+p)}{N}\right)^{1-\frac{q}{2}}\right).$$

## 1.5. Simulation studies

Simulations studies were performed to evaluate the proposed composition estimator $\widehat{\mathbf{X}}$ and to compare the results with the naive estimator $\widehat{\mathbf{X}}_s$ that replaces zero count with the maximum rounding error 0.5 (Aitchison, 2003) and transforms the counts into composition.

### 1.5.1. Simulation settings

Data $(\mathbf{X}^*, \mathbf{R})$ were generated as follows. The row probability vector $\{R_i\}_{i=1}^n$ was generated as the normalization of i.i.d entries $\{P_i\}_{i=1}^n$ uniformly drawn from Unif$[1, 10]$: $R_i = P_i/\sum_{k=1}^n P_k$. In order to generate the composition $\mathbf{X}^*$, we first generated a rank-$r$ matrix $\mathbf{Z}$ by $\mathbf{Z} = \mathbf{U}\mathbf{V}^T$, where $\mathbf{U} \in \mathbb{R}^{n\times r}$ and $\mathbf{V} \in \mathbb{R}^{p\times r}$. The components in $\mathbf{U}$ are the absolute values of i.i.d $N(0,1)$ normal random variables. $\mathbf{V} = \mathbf{V}_1 + \mathbf{V}_2$ is a spike matrix, where the diagonal elements of $\mathbf{V}_1$ are ones and off-diagonal entries are equal to 1 with the probability 0.3 and equal to 0 with the probability 0.7, and the entries of $\mathbf{V}_2$ are independent $N(0, 10^{-3})$ normal random variables. This procedure is repeated until we obtain a strict positive matrix $\mathbf{Z}$. The following two models are considered for $r$.

- Model 1 (Exactly low rank): $r = 20$.
- Model 2 (Approximately low rank): $r = n \wedge p$.

Then $\mathbf{X}^*$ was obtained through the normalization $X_{ij}^* = Z_{ij}/\sum_{k=1}^p Z_{ik}$, and count matrix $\mathbf{W}$ was generated as Mult$(\mathbf{R}\mathbf{X}^*, \gamma np)$, where $\gamma \in \{1, 2, 3, 4, 5\}$ was considered. We set the sample size and dimension as $n = p = 50, 100$, and 150, and repeated 50 simulations for each setting.

*1.5.2. Composition Estimate*

We applied the penalized maximum likelihood approach to simulated data in both low rank and approximately low rank cases. The tuning parameters $(\lambda, \alpha_X)$ in each estimator were chosen by five-fold cross-validation. For comparison, we calculated the naive estimators $\widehat{\mathbf{X}}_s$ that replaced zero counts by 0.5 and converted the counts into composition. Losses under squared Frobenius norm $\|\widehat{\mathbf{X}} - \mathbf{X}^*\|_F^2$ and Kullback-Leibler divergence $\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}})$ were used to measure the estimation performance.

The simulation results for Model 1 and 2 are summarized in Figures $1.1$ and $1.2$ respectively. We observed that the proposed estimator $\widehat{\mathbf{X}}$ resulted in uniformly smaller errors thatn those based on the naive estimator $\widehat{\mathbf{X}}_s$ in all settings, demonstrating the superiority of the penalized likelihood estimation. In addition, as expected, the difference in the loss of $\widehat{\mathbf{X}}$ and that of $\widehat{\mathbf{X}}_s$ got smaller as the total counts increased since the number of zeros decreased as more read counts were observed.



Figure 1.1: Frobenius norm error and Kullback-Leibler divergence between the estimated and the true compositions for different numbers of taxa $p$ in Model 1, where $\widehat{\mathbf{X}}$ is the proposed estimator and $\widehat{\mathbf{X}}_s$ is the estimator with simple zero replacement.

Figure 1.2: Frobenius norm error and Kullback-Leibler divergence between the estimated and the true compositions for different numbers of taxa $p$ in Model 2, where $\widehat{\mathbf{X}}$ is the proposed estimator and $\widehat{\mathbf{X}}_s$ is the estimator with simple zero replacement.

### 1.5.3. Diversity Index Estimate

To evaluate the ability to estimate the individual-level diversity and dispersion, we also calculated vector $L_2$ norm losses of the Shannon index and Simpson index, as well as the Frobenius norm error of Bray-Curtis index. The simulation results for both models are summarized in Figures $1.3$ and $1.4$. We see that the proposed estimator $\widehat{\mathbf{X}}$ uniformly outperformed the naive estimators $\widehat{\mathbf{X}}_s$ by a large margin.



Figure 1.3: Losses on different diversity indices between the estimated and the true compositions for different numbers of observed taxa $p$ in Model 1. Left panel: Shannon index; Middle panel:Simpson index; Right panel: Bray-Curtis index

Figure 1.4: Losses on different diversity indices between the estimated and the true compositions for different numbers of observed taxa $p$ in Model 2. Left panel: Shannon index; Middle panel:Simpson index; Right panel: Bray-Curtis index

## 1.6. Gut Microbiome Data Analysis

The gut microbiome plays an important role in regulating metabolic functions and immune homeostasis and exerts a profound influence on human health and disease. We applied the proposed method to a human gut microbiome dataset of a cross-sectional study of 98 healthy volunteers at the University of Pennsylvania (Wu et al., 2011). DNA from stool samples of these individuals were analyzed by 454/Roche pyrosequencing of 16S rRNA gene segments and yielded an average of 9265 reads per sample, with a standard deviation of 386, which led to identification of 3068 operational taxonomic units and 87 bacterial genera that were presented in at least one sample. Figure 1.5 show the proportions of zeros observed versus the size library sizes, indicating that many observed zeros are due to under sampling. It is therefore reasonably to assume that the true compositions of these rare genera are not zero.



Figure 1.5: Proportions of zeros observed versus the size library sizes, indicating that many observed zeros are due to under sampling.

Figure 1.6: Decay of singular values $d_{ii}$ from the SVD decomposition of $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$.

Figure 1.6 shows the decay of singular values $d_{ii}$ from the SVD decomposition of $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$, indicating that the approximate low-rank nature of the compositional matrix. We applied the proposed penalized likelihood method to estimate the positive compositions and used five-fold cross-validation to select the tuning parameters. As a comparison, we also replaced the count zeros by $0.5$ to obtain the naive estimator $\widehat{\mathbf{X}}_s$.

To illustrate the result, we define $\mathcal{M} = \{(i,j) \in [n] \times [p] \text{ such that } W_{ij} = 0\}$ as the set of zero counts $\mathbf{W}$. The top panel of Figure 2.1 shows the boxplots of the estimated compositions $\widehat{\mathbf{X}}$ except common genera *Bacteroides*, *Blautia* and *Roseburia* that have been observed in all individuals. Overall, we observed that the observed non-zero compositions had an effects in estimating the compositions with zeros counts and the estimated compositions in those with zero observations ($\mathcal{M}$) were almost always smaller than those with non-zero observations ($\mathcal{M}^c$). However, results from the simple zero replacement ($\widehat{\mathbf{X}}_s$) gave almost the same estimates for all samples/taxa in $\mathcal{M}$. The observed non-zero compositions almost had no effects in estimating the compositions with zero observed counts.

18

Figure 1.7: Boxplots of the estimated compositions for each genus for those with zero observations $(\mathcal{M})$ and those with non-zero observations $(\mathcal{M}^c)$. Top panel: proposed estimator $\widehat{\mathbf{X}}$; Bottom panel: estimator with zero-replacement $\widehat{\mathbf{X}}_s$.

## 1.7. Discussion

We have considered the problem of estimating the bacterial compositions based on sequencing data, particularly for those taxa with zero observed counts, one of the first step in any microbiome and metagenomic studies. We have developed a penalized likelihood estimation method for estimating the mcirobial abundances for these taxa with observed zero count. The estimate effectively

utilizes data across different individuals and across different taxa, which is in contrast to most of the available methods and has the flavor of shrinkage estimate. The estimation procedure makes two key assumptions. First, it assume that the true microbial compositions are always positive and the zero counts observed in metagenomic sequencing are due to under sampling. Our empirical data (Figure 1.5) seems to support this assumption. Second, it assumes that the true composition matrix has approximately low-rank structure. Under these assumptions, we have proposed a penalized likelihood estimation with a nuclear norm penalty function in order to obtain better estimate of the composition matrix. We have obtained the estimation upper bounds and also the min-max lower bounds and showed that our estimator is almost optimal. We have additionally obtained the upper bounds for the estimates of various commonly used diversity indices, including Shannon's index, Simpson's index and Brey-Curtis index. The resulting composition estimates can facilitate other downstream compositional data analysis, such as high dimensional regression analysis (Lin et al., 2014) and covariance estimation based on the composition data (Cao, Lin, and Li, 2016).

CHAPTER 2

LARGE COVARIANCE ESTIMATION FOR COMPOSITIONAL DATA VIA

COMPOSITION-ADJUSTED THRESHOLDING

In this chapter, we address the problem of covariance estimation for high-dimensional compositional data, and introduce a composition-adjusted thresholding (COAT) method under the assumption that the basis covariance matrix is sparse. Our method is based on a decomposition relating the compositional covariance to the basis covariance, which is approximately identifiable as the dimensionality tends to infinity. The resulting procedure can be viewed as thresholding the sample centered log-ratio covariance matrix and hence is scalable for large covariance matrices. We rigorously characterize the identifiability of the covariance parameters, derive rates of convergence under the spectral norm, and provide theoretical guarantees on support recovery. Simulation studies demonstrate that the COAT estimator outperforms some naive thresholding estimators that ignore the unique features of compositional data. We apply the proposed method to the analysis of a microbiome dataset in order to understand the dependence structure among bacterial taxa in the human gut.

## 2.1. Introduction

Compositional data, which represent the proportions or fractions of a whole, arise naturally in a wide range of applications; examples include geochemical compositions of rocks, household patterns of expenditures, species compositions of biological communities, and topic compositions of documents, among many others. This article is particularly motivated by the metagenomic analysis of microbiome data. The human microbiome is the totality of all microbes at various body sites, whose importance in human health and disease has increasingly been recognized. Recent studies have revealed that microbiome composition varies based on diet, health, and the environment (The Human Microbiome Project Consortium, 2012a), and may play a key role in complex diseases such as obesity, atherosclerosis, and Crohn's disease (Koeth et al., 2013; Lewis et al., 2015; Turnbaugh et al., 2009).

With the development of next-generation sequencing technologies, it is now possible to survey

the microbiome composition using direct DNA sequencing of either marker genes or the whole metagenomes. After aligning these sequence reads to the reference microbial genomes, one can quantify the relative abundances of microbial taxa. These sequencing-based microbiome studies, however, only provide a relative, rather than absolute, measure of the abundances of community components. The counts comprising these data (e.g., 16S rRNA gene reads or shotgun metagenomic reads) are set by the amount of genetic material extracted from the community or the sequencing depth, and analysis typically begins by normalizing the observed data by the total number of counts. The resulting fractions thus fall into a class of high-dimensional compositional data that we focus in this article. The high dimensionality refers to the fact that the number of taxa may be comparable to or much larger than the sample size.

An important question in metagenomic studies is to understand the co-occurrence and co-exclusion relationship between microbial taxa, which would provide valuable insights into the complex ecology of microbial communities (Faust et al., 2012). Standard correlation analysis from the raw proportions, however, can lead to spurious results due to the unit-sum constraint; the proportions tend to be correlated even if the absolute abundances are independent. Such undesired effects should be removed in an analysis in order to make valid inferences about the underlying biological processes. The compositional effects are further magnified by the low diversity of microbiome data, that is, a few taxa make up the overwhelming majority of the microbiome (Friedman and Alm, 2012).

Let $\mathbf{X} = (X_1, \ldots, X_p)^T$ be a composition of $p$ components (taxa) satisfying the simplex constraint

$$X_j > 0, \quad j = 1, \ldots, p, \quad \sum_{j=1}^{p} X_j = 1.$$

Owing to the difficulties arising from the simplex constraint, it has been a long-standing question how to appropriately model, estimate, and interpret the covariance structure of compositional data. The pioneering work of Aitchison, (1982, 2003) introduced several equivalent matrix specifications of compositional covariance structures via the log-ratios of components. Statistical methods based on these covariance models respect the unique features of compositional data and prove useful in a variety of applications such as geochemical analysis. A potential disadvantage of these models, however, is that they lack a direct interpretation in the usual sense of covariances and correlations; as a result, it is unclear how to impose certain structures such as sparsity in high dimensions, which

is crucial for our applications to microbiome data analysis.

Covariance matrix estimation is of fundamental importance in high-dimensional data analysis and has attracted much recent interest. It is well known that the sample covariance matrix performs poorly in high dimensions and regularization is thus indispensable. Bickel and Levina, (2008) and El Karoui, (2008) introduced regularized estimators by hard thresholding for large covariance matrices that satisfy certain notions of sparsity. Rothman, Levina, and Zhu, (2009) considered a more general class of thresholding functions, and Cai and Liu, (2011) proposed adaptive thresholding that adapts to the variability of individual entries. Exploiting a factor model structure, Fan, Fan, and Lv, (2008) proposed a factor-based method for high-dimensional covariance matrix estimation. Fan, Liao, and Mincheva, (2013) extended the work by considering a conditional sparsity structure and developed a POET method by thresholding principal orthogonal complements.

In this article, we address the problem of covariance estimation for high-dimensional compositional data. Let $\mathbf{W} = (W_1, \ldots, W_p)^T$ with $W_j > 0$ for all $j$ be a vector of latent variables, called the *basis*, that generate the observed data via the normalization

$$X_j = \frac{W_j}{\sum_{i=1}^p W_i}, \quad j = 1, \ldots, p. \tag{2.1}$$

Estimating the covariance structure of $\mathbf{W}$ has traditionally been considered infeasible owing to the apparent lack of identifiability. By exploring a decomposition relating the compositional covariance to the basis covariance, we find, however, that the nonidentifiability vanishes asymptotically as the dimensionality grows under certain sparsity assumptions. More specifically, define the *basis covariance matrix* $\mathbf{\Omega}_0 = (\omega_{ij}^0)_{p \times p}$ by

$$\omega_{ij}^0 = \mathrm{Cov}(Y_i, Y_j), \tag{2.2}$$

where $Y_j = \log W_j$. Then $\mathbf{\Omega}_0$ is approximately identifiable as long as it belongs to a class of large sparse covariance matrices.

The somewhat surprising "blessing of dimensionality" allows us to develop a simple, two-step method by first extracting a rank-2 component from the decomposition and then estimating the sparse component $\mathbf{\Omega}_0$ by thresholding the residual matrix. The resulting procedure can equivalently be viewed as thresholding the sample centered log-ratio covariance matrix, and hence is

optimization-free and scalable for large covariance matrices. We call our method *composition-adjusted thresholding* (COAT), which removes the "coat" of compositional effects from the covariance structure. We derive rates of convergence under the spectral norm and provide theoretical guarantees on support recovery. Simulation studies demonstrate that the COAT estimator outperforms some naive thresholding estimators that ignore the unique features of compositional data. We illustrate our method by analyzing a microbiome dataset in order to understand the dependence structure among bacterial taxa in the human gut.

The covariance relationship, which was due to Aitchison, (2003 sec. 4.11), has recently been exploited to develop algorithms for inferring correlation networks from metagenomic data (Ban, An, and Jiang, 2015; Fang et al., 2015; Friedman and Alm, 2012). Our contributions here are to turn the idea into a principled approach to sparse covariance matrix estimation and provide statistical insights into the issue of identifiability and the impacts of dimensionality. Our method also bears some resemblance to the POET method proposed by Fan, Liao, and Mincheva, (2013) in that underlying both methods is a low-rank plus sparse matrix decomposition. The rank-2 component in our method, however, arises from the covariance structure of compositional data rather than a factor model assumption. As a result, it can be obtained by simple algebraic operations without computing the principal components.

The rest of the article is organized as follows. Section 2 reviews a covariance relationship and addresses the issue of identifiability. Section 3 introduces the COAT methodology. Section 4 investigates the theoretical properties of the COAT estimator in terms of convergence rates and support recovery. Simulation studies and an application to human gut microbiome data are presented in Sections 5 and 6, respectively. We conclude the article with some discussion in Section 7 and relegate all proofs to the Appendix.

## 2.2. Identifiability of the Covariance Model

We first introduce some notation. Denote by $\| \cdot \|_1$, $\| \cdot \|_2$, $\| \cdot \|_F$, and $\| \cdot \|_{\max}$ the matrix $L_1$-norm, spectral norm, Frobenius norm, and entrywise $L_\infty$-norm, defined for a matrix $\mathbf{A} = (a_{ij})$ by $\|\mathbf{A}\|_1 = \max_j \sum_i |a_{ij}|$, $\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^T\mathbf{A})}$, $\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} a_{ij}^2}$, and $\|\mathbf{A}\|_{\max} = \max_{i,j} |a_{ij}|$, where $\lambda_{\max}(\cdot)$ denotes the largest eigenvalue.

In the latent variable covariance model (2.1) and (2.2), the basis covariance matrix $\boldsymbol{\Omega}_0$ is the parameter of interest. One of the matrix specifications of compositional covariance structures introduced by Aitchison, (2003) is the *variation matrix* $\mathbf{T}_0 = (\tau_{ij}^0)_{p \times p}$ defined by

$$\tau_{ij}^0 = \mathrm{Var}(\log(X_i/X_j)). \tag{2.3}$$

In view of the relationship (2.1), we can decompose $\tau_{ij}^0$ as

$$
\begin{aligned}
\tau_{ij}^0 &= \mathrm{Var}(\log W_i - \log W_j) \\
&= \mathrm{Var}(Y_i) + \mathrm{Var}(Y_j) - 2\,\mathrm{Cov}(Y_i, Y_j) \\
&= \omega_{ii}^0 + \omega_{jj}^0 - 2\omega_{ij}^0,
\end{aligned} \tag{2.4}
$$

or in matrix form,

$$\mathbf{T}_0 = \boldsymbol{\omega}_0 \mathbf{1}^T + \mathbf{1}\boldsymbol{\omega}_0^T - 2\boldsymbol{\Omega}_0, \tag{2.5}$$

where $\boldsymbol{\omega}_0 = (\omega_{11}^0, \ldots, \omega_{pp}^0)^T$ and $\mathbf{1} = (1, \ldots, 1)^T$. Corresponding to the many-to-one relationship between bases and compositions, the basis covariance matrix $\boldsymbol{\Omega}_0$ is unidentifiable from the decomposition (2.5), since $\boldsymbol{\omega}_0 \mathbf{1}^T + \mathbf{1}\boldsymbol{\omega}_0^T$ and $\boldsymbol{\Omega}_0$ are in general not orthogonal to each other (with respect to the usual Euclidean inner product). In fact, using the *centered log-ratio covariance matrix* $\boldsymbol{\Gamma}_0 = (\gamma_{ij}^0)_{p \times p}$ defined by

$$\gamma_{ij}^0 = \mathrm{Cov}\{\log(X_i/g(\mathbf{X})), \log(X_j/g(\mathbf{X}))\},$$

where $g(\mathbf{x}) = (\prod_{j=1}^p x_j)^{1/p}$ is the geometric mean of a vector $\mathbf{x} = (x_1, \ldots, x_p)^T$, we can similarly write

$$
\begin{aligned}
\tau_{ij}^0 &= \mathrm{Var}\{\log(X_i/g(\mathbf{X})) - \log(X_j/g(\mathbf{X}))\} \\
&= \mathrm{Var}\{\log(X_i/g(\mathbf{X}))\} + \mathrm{Var}\{\log(X_j/g(\mathbf{X}))\} - 2\,\mathrm{Cov}\{\log(X_i/g(\mathbf{X})), \log(X_j/g(\mathbf{X}))\} \\
&= \gamma_{ii}^0 + \gamma_{jj}^0 - 2\gamma_{ij}^0,
\end{aligned}
$$

or in matrix form,

$$\mathbf{T}_0 = \boldsymbol{\gamma}_0 \mathbf{1}^T + \mathbf{1}\boldsymbol{\gamma}_0^T - 2\boldsymbol{\Gamma}_0, \tag{2.6}$$

where $\gamma_0 = (\gamma_{11}^0, \ldots, \gamma_{pp}^0)^T$ and $\mathbf{1} = (1, \ldots, 1)^T$. Unlike (2.5), the following proposition shows that (2.6) is an orthogonal decomposition and hence the components $\gamma_0 \mathbf{1}^T + \mathbf{1}\gamma_0^T$ and $\Gamma_0$ are identifiable. In addition, by comparing the decompositions (2.5) and (2.6), we can bound the difference between $\Omega_0$ and its identifiable counterpart $\Gamma_0$ as follows.

**Proposition 3.** *The components $\gamma_0 \mathbf{1}^T + \mathbf{1}\gamma_0^T$ and $\Gamma_0$ in the decomposition* (2.6) *are orthogonal to each other. Moreover, for the covariance parameters $\Omega_0$ and $\Gamma_0$ in the decompositions* (2.5) *and* (2.6),

$$\|\Omega_0 - \Gamma_0\|_{\max} \leq 3p^{-1}\|\Omega_0\|_1.$$

Proposition 3 entails that the covariance parameter $\Omega_0$ is *approximately* identifiable as long as $\|\Omega_0\|_1 = o(p)$. In particular, suppose that $\Omega_0$ belongs to a class of sparse covariance matrices considered by Bickel and Levina, (2008),

$$\mathcal{U}(q, s_0(p), M) \equiv \left\{ \Omega \colon \Omega \succ 0, \max_j \omega_{jj} \leq M, \max_i \sum_{j=1}^p |\omega_{ij}|^q \leq s_0(p) \right\}, \qquad (2.7)$$

where $0 \leq q < 1$ and $\Omega \succ 0$ denotes that $\Omega$ is positive definite. Then

$$\|\Omega_0\|_1 = \max_i \sum_{j=1}^p |\omega_{ij}^0|^{1-q} |\omega_{ij}^0|^q \leq \max_i \sum_{j=1}^p (\omega_{ii}^0 \omega_{jj}^0)^{(1-q)/2} |\omega_{ij}^0|^q \leq M^{1-q} s_0(p),$$

and hence the parameters $\Omega_0$ and $\Gamma_0$ are asymptotically indistinguishable when $s_0(p) = o(p)$. This allows us to use $\Gamma_0$ as a proxy for $\Omega_0$ and greatly facilitates the development of new methodology and associated theory. The intuition behind the approximate identifiability under the sparsity assumption is that the rank-2 component $\omega_0 \mathbf{1}^T + \mathbf{1}\omega_0^T$ represents a global effect that spreads across all rows and columns, while the sparse component $\Omega_0$ represents a local effect that is confined to individual entries.

Also of interest is the *exact* identifiability of $\Omega_0$ over $L_0$-balls, which has been studied by Fang et al., (2015) and Ban, An, and Jiang, (2015). The following result provides a sufficient and necessary condition for the exact identifiability of $\Omega_0$ by confining it to an $L_0$-ball.

**Proposition 4.** *Suppose that $\mathbf{\Omega}_0$ belongs to the $L_0$-ball*

$$\mathcal{B}_0(s_e(p)) \equiv \left\{ \mathbf{\Omega} \colon \sum_{(i,j)\colon\, i<j} I(\omega_{ij} \neq 0) \leq s_e(p) \right\},$$

*where $p \geq 5$. Then there exist no two values of $\mathbf{\Omega}_0$ that correspond to the same $\mathbf{T}_0$ in (2.5) if and only if $s_e(p) < (p-1)/2$.*

A counterexample is provided in the proof of Proposition 4 to show that the sparsity conditions in Fang et al., (2015) and Ban, An, and Jiang, (2015), which are both at the order of $O(p^2)$, do not suffice. The identifiability condition in Proposition 4 essentially requires the average degree of the correlation network to be less than 1, which is too restrictive to be useful in practice. This illustrates the importance and necessity of introducing the notion of approximate identifiability.

## 2.3. A Sparse Covariance Estimator for Compositional Data

Suppose that $(\mathbf{W}_k, \mathbf{X}_k)$, $k = 1, \ldots, n$, are independent copies of $(\mathbf{W}, \mathbf{X})$, where the compositions $\mathbf{X}_k = (X_{k1}, \ldots, X_{kp})^T$ are observed and the bases $\mathbf{W}_k = (W_{k1}, \ldots, W_{kp})^T$ are latent. In Section 3.1, we rely on the decompositions (2.5) and (2.6) and Proposition 3 to develop an estimator of $\mathbf{\Omega}_0$, and in Section 3.2 discuss the selection of the tuning parameter.

### 2.3.1. Composition-Adjusted Thresholding

In view of Proposition 3, we wish to estimate the covariance parameter $\mathbf{\Omega}_0$ via the proxy $\mathbf{\Gamma}_0$. To this end, we first construct an empirical estimate of $\mathbf{\Gamma}_0$ and then apply adaptive thresholding to the estimate.

There are two equivalent ways to form the estimate of $\mathbf{\Gamma}_0$. Motivated by the decomposition (2.6), one can start with the sample counterpart $\widehat{\mathbf{T}} = (\hat{\tau}_{ij})_{p \times p}$ of $\mathbf{T}_0$ defined by

$$\hat{\tau}_{ij} = \frac{1}{n} \sum_{k=1}^{n} (\tau_{kij} - \bar{\tau}_{ij})^2,$$

where $\tau_{kij} = \log(X_{ki}/X_{kj})$ and $\bar{\tau}_{ij} = n^{-1}\sum_{k=1}^{n} \tau_{kij}$. A rank-2 component $\widehat{\boldsymbol{\alpha}}\mathbf{1}^T + \mathbf{1}\widehat{\boldsymbol{\alpha}}^T$ with $\widehat{\boldsymbol{\alpha}} = (\hat{\alpha}_1, \ldots, \hat{\alpha}_p)^T$ can be extracted from the decomposition (2.6) by projecting $\widehat{\mathbf{T}}$ onto the subspace

$\mathcal{A} \equiv \{\boldsymbol{\alpha}\mathbf{1}^T + \mathbf{1}\boldsymbol{\alpha}^T : \boldsymbol{\alpha} \in \mathbb{R}^p\}$, which is given by

$$\hat{\alpha}_i = \hat{\tau}_{i\cdot} - \frac{1}{2}\hat{\tau}_{\cdot\cdot},$$

where $\hat{\tau}_{i\cdot} = p^{-1}\sum_{j=1}^{p}\hat{\tau}_{ij}$ and $\hat{\tau}_{\cdot\cdot} = p^{-2}\sum_{i,j=1}^{p}\hat{\tau}_{ij}$. The residual matrix $\widehat{\boldsymbol{\Gamma}} = -(\widehat{\mathbf{T}} - \widehat{\boldsymbol{\alpha}}\mathbf{1}^T - \mathbf{1}\widehat{\boldsymbol{\alpha}}^T)/2$, with entries

$$\hat{\gamma}_{ij} = -\frac{1}{2}(\hat{\tau}_{ij} - \hat{\alpha}_i - \hat{\alpha}_j) = -\frac{1}{2}(\hat{\tau}_{ij} - \hat{\tau}_{i\cdot} - \hat{\tau}_{j\cdot} + \hat{\tau}_{\cdot\cdot}),$$

is then an estimate of $\boldsymbol{\Gamma}_0$. Alternatively, $\widehat{\boldsymbol{\Gamma}}$ can be obtained directly as the sample counterpart of $\boldsymbol{\Gamma}_0$ through the expression

$$\hat{\gamma}_{ij} = \frac{1}{n}\sum_{k=1}^{n}(\gamma_{ki} - \bar{\gamma}_i)(\gamma_{kj} - \bar{\gamma}_j), \tag{2.8}$$

where $\gamma_{kj} = \log(X_{kj}/g(\mathbf{X}_k))$ and $\bar{\gamma}_j = n^{-1}\sum_{k=1}^{n}\gamma_{kj}$.

Now applying adaptive thresholding to $\widehat{\boldsymbol{\Gamma}}$, we define the *composition-adjusted thresholding* (COAT) estimator

$$\widehat{\boldsymbol{\Omega}} = (\hat{\omega}_{ij})_{p \times p} \quad \text{with } \hat{\omega}_{ij} = S_{\lambda_{ij}}(\hat{\gamma}_{ij}), \tag{2.9}$$

where $S_\lambda(\cdot)$ is a general thresholding function and $\lambda_{ij} > 0$ are entry-dependent thresholds.

In this article, we consider a class of general thresholding functions $S_\lambda(\cdot)$ that satisfy the following conditions:

(i) $S_\lambda(z) = 0$ for $|z| \le \lambda$;

(ii) $|S_\lambda(z) - z| \le \lambda$ for all $z \in \mathbb{R}$.

These two conditions were assumed by Rothman, Levina, and Zhu, (2009) and Cai and Liu, (2011) along with another condition that is not required in our analysis. Examples of thresholding functions belonging to this class include the hard thresholding rule $S_\lambda(z) = zI(|z| \ge \lambda)$, the soft thresholding rule $S_\lambda(z) = \text{sgn}(z)(|z| - \lambda)_+$, and the adaptive lasso rule $S_\lambda(z) = z(1 - |\lambda/z|^\eta)_+$ for $\eta \ge 1$.

The performance of the COAT estimator depends critically on the choice of thresholds. Using entry-adaptive thresholds may in general improve the performance over applying a universal threshold.

To derive a data-driven choice of $\lambda_{ij}$, define

$$\theta_{ij} = \text{Var}\{(Y_i - \mu_i)(Y_j - \mu_j)\},$$

where $\mu_j = EY_j$. We take $\lambda_{ij}$ to be of the form

$$\lambda_{ij} = \lambda\sqrt{\hat{\theta}_{ij}}, \qquad (2.10)$$

where $\hat{\theta}_{ij}$ are estimates of $\theta_{ij}$, and $\lambda > 0$ is a tuning parameter to be chosen, for example, by cross-validation. We rewrite (2.8) as $\hat{\gamma}_{ij} = n^{-1}\sum_{k=1}^{n}\gamma_{kij}$, where $\gamma_{kij} = (\gamma_{ki} - \bar{\gamma}_i)(\gamma_{kj} - \bar{\gamma}_j)$. Then $\theta_{ij}$ can be estimated by

$$\hat{\theta}_{ij} = \frac{1}{n}\sum_{k=1}^{n}(\gamma_{kij} - \hat{\gamma}_{ij})^2.$$

### 2.3.2. Tuning Parameter Selection

The thresholds defined by (2.10) depend on the tuning parameter $\lambda$, which can be chosen through $V$-fold cross-validation. Denote by $\widehat{\boldsymbol{\Omega}}^{(-v)}(\lambda)$ the COAT estimate based on the training data excluding the $v$th fold, and $\widehat{\boldsymbol{\Gamma}}_v$ the residual matrix (or the sample centered log-ratio covariance matrix) based on the test data including only the $v$th fold. We choose the optimal value of $\lambda$ that minimizes the cross-validation error

$$\text{CV}(\lambda) = \frac{1}{V}\sum_{v=1}^{V}\|\widehat{\boldsymbol{\Omega}}^{(-v)}(\lambda) - \widehat{\boldsymbol{\Gamma}}^{(v)}\|_F^2.$$

With the optimal $\lambda$, we then compute the COAT estimate based on the full dataset as our final estimate. When the positive definiteness of the covariance estimate in finite samples is required for interpretation, we follow the approach of Fan, Liao, and Mincheva, (2013) and choose $\lambda$ in the range where the minimum eigenvalue of the COAT estimate is positive.

## 2.4. Theoretical Properties

In this section, we investigate the asymptotic properties of the COAT estimator. As a distinguishing feature of our theoretical analysis, we assume neither the exact identifiability of the parameters nor that the degree of (approximate) identifiability is dominated by the statistical error. Instead, the degree of identifiability enters our analysis and shows up in the resulting rate of convergence. Such

theoretical analysis is rare in the literature, but is extremely relevant for latent variable models in the presence of nonidentifiability and is of theoretical interest in its own right. We introduce our assumptions in Section 4.1, and present our main results on rates of convergence and support recovery in Section 4.2.

### 2.4.1. Assumptions

Recall that $Y_j = \log W_j$, $\mu_j = EY_j$, and $\theta_{ij} = \mathrm{Var}\{(Y_i - \mu_i)(Y_j - \mu_j)\}$, and define $Y_{kj} = \log W_{kj}$. Without loss of generality, assume $\mu_j = 0$ for all $j$ throughout this section. We need to impose the following moment conditions on the log-basis $\mathbf{Y} = (Y_1, \ldots, Y_p)^T$.

**Condition 3.** There exists a constant $\alpha > 0$ such that $\max_j E \exp(\alpha Y_j^2) \leq 2$.

**Condition 4.** The basis covariance matrix $\mathbf{\Omega}_0$ belongs to the class $\mathcal{U}(q, s_0(p), M)$ defined by (2.7), where $0 \leq q < 1$, $s_0(p) = o(p)$, and $\log p = o(n^{1/5})$.

**Condition 5.** There exists a constant $\tau > 0$ such that $\min_{i,j} \theta_{ij} \geq \tau$.

**Condition 6.** There exists a sequence $s_1(p) = o(p)$ such that

$$\max_{i,j,\ell} \left| \sum_{m=1}^{p} EY_i Y_j Y_\ell Y_m \right| \leq s_1(p).$$

Conditions 1–3 are similar to those commonly assumed in the covariance estimation literature; see, for example, Cai and Liu, (2011). Condition 3 requires that the variables $Y_j$s be uniformly sub-Gaussian; the definition we use here is among several equivalent ways of defining sub-Gaussianity (Boucheron, Lugosi, and Massart, 2013 sec. 2.3), and is most convenient for our technical analysis. Condition 4 imposes some restrictions on the dimensionality and sparsity of the basis covariance matrix $\mathbf{\Omega}_0$. It is worth mentioning that the sparsity level condition $s_0 = o(p)$ is so weak that it suffices to guarantee only approximate identifiability but allows the degree of nonidentifiability to be large relative to the statistical error. Condition 5 is essential for methods based on adaptive thresholding. Condition 6 arises from identifiability considerations in estimating the variances $\theta_{ij}$. In particular, if $\mathbf{Y}$ is multivariate normal, then Condition 6 is implied by the assumptions $\mathbf{\Omega}_0 \in \mathcal{U}(q, s_0(p), M)$ and $s_0(p) = o(p)$ in Condition 4, since from Isserlis' theorem (Isserlis, 1918) we have

$$\max_{i,j,\ell} \left| \sum_{m=1}^{p} EY_i Y_j Y_\ell Y_m \right| \leq \max_{i,j,\ell} \sum_{m=1}^{p} \left( |\omega_{ij}^0||\omega_{\ell m}^0| + |\omega_{i\ell}^0||\omega_{jm}^0| + |\omega_{im}^0||\omega_{j\ell}^0| \right) \leq 3M^{2-q} s_0(p).$$

*2.4.2. Main Results*

We are now in a position to state our main results. The following theorem gives the rate of convergence under the spectral norm for the COAT estimator.

**Theorem 4** (Rate of convergence). *Under Conditions 3–6, if the tuning parameter $\lambda$ in (2.10) is chosen to be*

$$\lambda = C_1 \sqrt{\frac{\log p}{n}} + C_2 \frac{s_0(p)}{p} \tag{2.11}$$

*for sufficiently large $C_1, C_2 > 0$, then the COAT estimator $\widehat{\boldsymbol{\Omega}}$ in (2.9) satisfies*

$$\|\widehat{\boldsymbol{\Omega}} - \boldsymbol{\Omega}_0\|_2 = O_p \left\{ s_0(p) \left( \sqrt{\frac{\log p}{n}} + \frac{s_0(p)}{p} \right)^{1-q} \right\}$$

*uniformly on $\mathcal{U}(q, s_0(p), M)$.*

The rate of convergence provided by Theorem 4 exhibits an interesting decomposition: the term $s_0(p)\{(\log p)/n\}^{(1-q)/2}$ represents the estimation error due to estimating $\boldsymbol{\Gamma}_0$, while the term $s_0(p)(s_0(p)/p)^{1-q}$ accounts for the approximation error due to using $\boldsymbol{\Gamma}_0$ as a proxy for $\boldsymbol{\Omega}_0$. In particular, if the approximation error is dominated by the estimation error, then the COAT estimator attains the minimax optimal rate under the spectral norm over $\mathcal{U}(q, s_0(p), M)$ (Cai and Zhou, 2012). It is important to note that the dimensionality $p$ appears in both terms where it plays opposite roles. We observe a "curse of dimensionality" in the first term, where the growth of dimensionality contributes a logarithmic factor to the estimation error. In contrast, a "blessing of dimensionality" is reflected by the second term in that a diverging dimensionality shrinks the approximation error toward zero at a power rate.

The insights gained from Theorem 4 have important implications for compositional data analysis. In the analysis of many compositional datasets, the dimensionality often depends on the taxonomic level to be examined. For example, in metagenomic studies, the dimensionality may range from only a few taxa at the phylum level to thousands of taxa at the operational taxonomic unit (OTU) level. Suppose, for simplicity, that the magnitudes of correlation signals are of about the same order across different taxonomic levels. Then Theorem 4 indicates a tradeoff between an accurate estimation of the covariance structure with low dimensionality and a sensible interpretation in terms of the basis components with high dimensionality. This tradeoff thus suggests the need to

analyze compositional data at relatively finer taxonomic levels when a latent variable interpretation is desired.

The proof of Theorem 4 relies on a series of concentration inequalities that take the approximation error term into account, which can be found in the Appendix. As a consequence of these inequalities, we obtain the following result regarding the support recovery property of the COAT estimator. Here the support of $\Omega_0$ refers to the set of all indices $(i, j)$ with $\omega_{ij}^0 \neq 0$.

**Theorem 5** (Support recovery)**.** *Under Conditions 3–6, if the tuning parameter $\lambda$ in* (2.10) *is chosen as in* (2.11)*, then the COAT estimator $\widehat{\Omega}$ in* (2.9) *satisfies*

$$P\left(\hat{\omega}_{ij} = 0 \text{ for all } (i, j) \text{ with } \omega_{ij}^0 = 0\right) \to 1. \tag{2.12}$$

*Moreover, if in addition*

$$\min_{(i,j)\colon \omega_{ij}^0 \neq 0} |\omega_{ij}^0| / \sqrt{\theta_{ij}} \geq C\lambda \tag{2.13}$$

*for some constant $C > 3/2$, then*

$$P\left(\operatorname{sgn}(\hat{\omega}_{ij}) = \operatorname{sgn}(\omega_{ij}^0) \text{ for all } (i, j)\right) \to 1. \tag{2.14}$$

Theorem 5 parallels the support recovery results in Rothman, Levina, and Zhu, (2009) and Cai and Liu, (2011). However, owing to the extra term $s_0(p)/p$ in the expression of $\lambda$, the assumption (2.13) requires in addition that no correlation signals fall below the approximation error. In other words, exact support recovery will break down if any correlation signal is confounded by the compositional effect.

## 2.5. Simulation Studies

We conducted simulation studies to compare the numerical performance of the COAT estimator $\widehat{\Omega}$ with that of the oracle thresholding estimator $\widehat{\Omega}_o$, which knew the latent basis components and applied the thresholding procedure to the sample covariance matrix of the log-basis $\mathbf{Y}$. We also include in our comparison two naive thresholding estimators $\widehat{\Omega}_c$ and $\widehat{\Omega}_l$, which are based on the sample covariance matrices of the composition $\mathbf{X}$ and its logarithm $\log \mathbf{X}$, respectively. Note that $\widehat{\Omega}_o$ is the ideal estimator that the COAT estimator attempts to mimic, whereas both $\widehat{\Omega}_c$ and $\widehat{\Omega}_l$ ignore

the unique features of compositional data and thus are expected to perform poorly.

### 2.5.1. Simulation Settings

The data $(\mathbf{W}_k, \mathbf{X}_k)$, $k = 1, \ldots, n$, were generated as follows. We first generated $\mathbf{Y}_k$ in two different ways:

(i) $\mathbf{Y}_k$ are independent from the multivariate normal distribution $N_p(\mu, \mathbf{\Omega}_0)$;

(ii) $\mathbf{Y}_k = \mu + \mathbf{F}\mathbf{U}_k/\sqrt{10}$, where $\mathbf{F}\mathbf{F}^T = \mathbf{\Omega}_0$ and the components of $\mathbf{U}_k$ are independent gamma variables with shape parameter 10 and scale parameter 1, so that $\mathrm{Var}(\mathbf{Y}_k) = \mathbf{\Omega}_0$. Here the matrix $\mathbf{F}$ is obtained by computing the singular value decomposition $\mathbf{\Omega}_0 = \mathbf{Q}\mathbf{S}\mathbf{Q}^T$ and letting $\mathbf{F} = \mathbf{Q}\mathbf{S}^{1/2}$.

Then $\mathbf{W}_k = (W_{k1}, \ldots, W_{kp})^T$ and $\mathbf{X}_k = (X_{k1}, \ldots, X_{kp})^T$ were obtained through the transformations $W_{kj} = e^{Y_{kj}}$ and $X_{kj} = W_{kj}/\sum_{i=1}^{p} W_{ki}$, $j = 1, \ldots, p$. Hence, in Case (i), $\mathbf{W}_k$ and $\mathbf{X}_k$ follow multivariate log-normal and logistic normal distributions (Aitchison and Shen, 1980), respectively; the distributions of $\mathbf{W}_k$ and $\mathbf{X}_k$ in Case (ii) can similarly be viewed as a type of multivariate log-gamma and logistic-gamma distributions.

In both cases, we took the the components of $\mu$ randomly from the uniform distribution on $[0, 10]$, in order to reflect the fact that compositional data arising from metagenomic studies are often heterogeneous. The following two models for the covariance matrix $\mathbf{\Omega}_0$ were considered:

- Model 1 (Identity covariance): $\mathbf{\Omega}_0 = \mathbf{I}_p$.
- Model 2 (Sparse covariance): $\mathbf{\Omega}_0 = \mathrm{diag}(\mathbf{A}_1, \mathbf{A}_2)$, where $\mathbf{A}_1 = \mathbf{B} + \varepsilon\mathbf{I}_{p_1}$, $\mathbf{A}_2 = 4\mathbf{I}_{p_2}$, $p_1 = \lfloor 2\sqrt{p} \rfloor$, $p_2 = p - p_1$, and $\mathbf{B}$ is a symmetric matrix whose lower triangular entries are independent from the uniform distribution on $[-1, -0.5] \cup [0.5, 1]$ with probability 0.2 and equal to 0 with probability 0.8. We set $\varepsilon = \max(-\lambda_{\min}(\mathbf{B}), 0) + 0.01$ to ensure that $\mathbf{A}_1$ is positive definite, where $\lambda_{\min}(\cdot)$ denotes the smallest eigenvalue.

Model 1 is an extreme but illustrative case intended for comparing the distributions of spurious correlations under different transformations. The setting of Model 2 is typical in the covariance estimation literature and similar to that in Cai and Liu, (2011). We set the sample size $n = 100$ and the dimension $p = 50$, 100, and 200, and repeated 100 simulations for each setting.

Figure 2.1: Boxplots of sample correlations with simulated data under different transformations in Model 1.

### 2.5.2. Spurious Correlations

The boxplots of sample correlations with simulated data under different transformations in Model 1 are shown in Figure 2.1. Clearly, the sample centered log-ratio (clr) correlations are centered around zero and have a similar distribution to that of the sample correlations of $\mathbf{Y}$; the resemblance tends to increase as the dimension $p$ grows. This trend is consistent with Proposition 3 and provides numerical evidence for the validity of the centered log-ratio covariance matrix $\mathbf{\Gamma}_0$ as a proxy for $\mathbf{\Omega}_0$. In fact, from the proof of Proposition 3 we have, when $\mathbf{\Omega}_0 = \mathbf{I}_p$,

$$\|\mathbf{\Omega}_0 - \mathbf{\Gamma}_0\|_{\max} = \max_{i,j} |\omega_{i\cdot}^0 + \omega_{j\cdot}^0 - \omega_{\cdot\cdot}^0| = p^{-1}.$$

In contrast, the phenomenon of spurious correlations is observed on both $\log \mathbf{X}$ and $\mathbf{X}$. The sample correlations of $\log \mathbf{X}$ exhibit a severe upward bias, while the sample correlations of $\mathbf{X}$ contain many outliers that would be detected as signals by a thresholding procedure with threshold level close to 1. Moreover, the spurious correlations seem to become worse with gamma-related distributions where the components of the composition have more heterogeneous means.

### 2.5.3. Performance Comparisons

We applied the COAT method with hard and soft thresholding rules to simulated data in Model 2. For comparison, we also applied the thresholding procedure to the sample covariance matrices of $\mathbf{Y}$, $\log\mathbf{X}$, and $\mathbf{X}$, resulting in the estimators $\widehat{\boldsymbol{\Omega}}_o$, $\widehat{\boldsymbol{\Omega}}_l$, and $\widehat{\boldsymbol{\Omega}}_c$, respectively. The tuning parameter $\lambda$ in each thresholding estimator was chosen by tenfold cross-validation. Losses under the matrix $L_1$-norm, spectral norm, and Frobenius norm were used to measure the estimation performance, while the true positive rate and false positive rate were employed to assess the quality of support recovery.

The simulation results for Model 2 with normal- and gamma-related distributions are summarized in Tables 2.1 and 2.2, respectively. We see that the COAT estimator $\widehat{\boldsymbol{\Omega}}$ performs almost equally well as the ideal estimator $\widehat{\boldsymbol{\Omega}}_o$, and outperforms the naive thresholding estimators $\widehat{\boldsymbol{\Omega}}_l$ and $\widehat{\boldsymbol{\Omega}}_c$ by a large margin. In particular, the estimation losses of $\widehat{\boldsymbol{\Omega}}_l$ are disastrously large in the gamma setting, in agreement with the severe bias observed in Figure 2.1. The estimation losses of $\widehat{\boldsymbol{\Omega}}_c$ do not change much across different thresholding rules and distributions, since all entries of the estimate are very small relative to the true values. Both $\widehat{\boldsymbol{\Omega}}_l$ and $\widehat{\boldsymbol{\Omega}}_c$ show inferior performance in terms of true and false positive rates, indicating that they are not model selection consistent. Comparisons between hard and soft thresholding rules suggest that the former is more conservative in selecting false positives and results in a more parsimonious model, whereas the latter strikes a balance between true and false positives due to the shrinkage effect.

To further compare the support recovery performance without selecting a threshold level, we plot the receiver operating characteristic (ROC) curves for all methods in Figure 2.2. Note that hard and soft thresholding rules lead to the same ROC curve for each method. We observe that the ROC curves for $\widehat{\boldsymbol{\Omega}}$ and $\widehat{\boldsymbol{\Omega}}_o$ are almost indistinguishable and uniformly dominate those for $\widehat{\boldsymbol{\Omega}}_l$ and $\widehat{\boldsymbol{\Omega}}_c$, demonstrating the superiority of the COAT method. Of the two naive thresholding estimators, $\widehat{\boldsymbol{\Omega}}_l$ tends to outperform $\widehat{\boldsymbol{\Omega}}_c$ when the threshold level is high, since the former is less influenced by the high spurious correlations as reflected in Figure 2.1.

Table 2.1: Means (standard errors) of various performance measures for four methods with hard and soft thresholding rules with normal-related distributions over 100 replications

| | Hard | | | | Soft | | | |
|---|---|---|---|---|---|---|---|---|
| $p$ | $\widehat{\Omega}$ | $\widehat{\Omega}_o$ | $\widehat{\Omega}_l$ | $\widehat{\Omega}_c$ | $\widehat{\Omega}$ | $\widehat{\Omega}_o$ | $\widehat{\Omega}_l$ | $\widehat{\Omega}_c$ |
| Matrix $L_1$-norm loss | | | | | | | | |
| 50 | 4.09 (0.05) | 4.02 (0.05) | 11.72 (1.51) | 6.91 (0.00) | 4.34 (0.05) | 4.10 (0.05) | 18.73 (0.64) | 6.91 (0.00) |
| 100 | 5.46 (0.04) | 5.50 (0.05) | 7.85 (1.13) | 8.07 (0.00) | 5.50 (0.05) | 5.40 (0.05) | 27.10 (1.18) | 8.07 (0.00) |
| 200 | 8.07 (0.04) | 8.10 (0.04) | 8.36 (0.04) | 10.93 (0.00) | 7.72 (0.06) | 7.66 (0.06) | 22.61 (1.13) | 10.93 (0.00) |
| Spectral norm loss | | | | | | | | |
| 50 | 2.32 (0.02) | 2.22 (0.02) | 7.23 (0.99) | 4.91 (0.00) | 2.49 (0.02) | 2.40 (0.02) | 10.23 (0.42) | 4.92 (0.00) |
| 100 | 2.89 (0.02) | 2.90 (0.02) | 4.50 (0.74) | 5.46 (0.00) | 3.01 (0.02) | 2.98 (0.02) | 13.93 (0.70) | 5.46 (0.00) |
| 200 | 3.55 (0.02) | 3.55 (0.02) | 3.68 (0.02) | 6.43 (0.00) | 3.93 (0.02) | 3.89 (0.02) | 9.28 (0.60) | 6.43 (0.00) |
| Frobenius norm loss | | | | | | | | |
| 50 | 5.63 (0.03) | 5.50 (0.03) | 11.47 (1.01) | 26.00 (0.00) | 8.37 (0.03) | 7.99 (0.03) | 15.18 (0.39) | 26.01 (0.00) |
| 100 | 8.70 (0.04) | 8.66 (0.03) | 11.39 (0.81) | 38.39 (0.00) | 13.11 (0.04) | 12.87 (0.04) | 24.18 (0.70) | 38.39 (0.00) |
| 200 | 12.03 (0.03) | 12.05 (0.03) | 12.97 (0.05) | 55.78 (0.00) | 20.48 (0.03) | 20.32 (0.03) | 27.06 (0.68) | 55.78 (0.00) |
| True positive rate | | | | | | | | |
| 50 | 0.65 (0.01) | 0.67 (0.01) | 0.70 (0.01) | 0.76 (0.02) | 0.94 (0.00) | 0.95 (0.00) | 0.93 (0.00) | 0.94 (0.00) |
| 100 | 0.59 (0.00) | 0.59 (0.00) | 0.59 (0.01) | 0.46 (0.02) | 0.91 (0.00) | 0.91 (0.00) | 0.87 (0.00) | 0.92 (0.00) |
| 200 | 0.60 (0.00) | 0.60 (0.00) | 0.60 (0.00) | 0.36 (0.02) | 0.83 (0.00) | 0.84 (0.00) | 0.87 (0.00) | 0.89 (0.00) |
| False positive rate | | | | | | | | |
| 50 | 0.00 (0.00) | 0.00 (0.00) | 0.15 (0.03) | 0.44 (0.03) | 0.11 (0.00) | 0.09 (0.00) | 0.53 (0.01) | 0.61 (0.01) |
| 100 | 0.00 (0.00) | 0.00 (0.00) | 0.02 (0.01) | 0.11 (0.01) | 0.06 (0.00) | 0.06 (0.00) | 0.41 (0.01) | 0.59 (0.01) |
| 200 | 0.00 (0.00) | 0.00 (0.00) | 0.00 (0.00) | 0.07 (0.01) | 0.03 (0.00) | 0.03 (0.00) | 0.18 (0.01) | 0.54 (0.01) |

Table 2.2: Means (standard errors) of various performance measures for four methods with hard and soft thresholding rules with gamma-related distributions over 100 replications

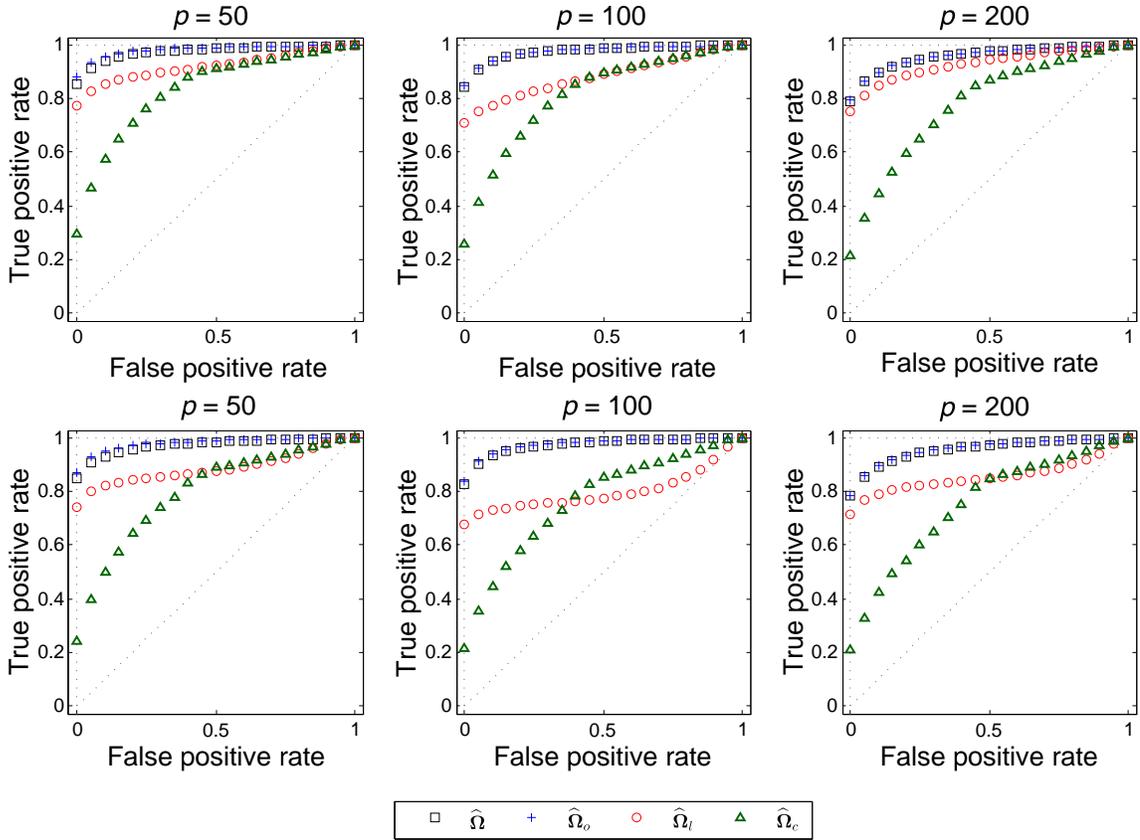| | Hard | | | | Soft | | | |
|---|---|---|---|---|---|---|---|---|
| $p$ | $\widehat{\Omega}$ | $\widehat{\Omega}_o$ | $\widehat{\Omega}_l$ | $\widehat{\Omega}_c$ | $\widehat{\Omega}$ | $\widehat{\Omega}_o$ | $\widehat{\Omega}_l$ | $\widehat{\Omega}_c$ |
| Matrix $L_1$-norm loss | | | | | | | | |
| 50 | 4.15 (0.07) | 4.09 (0.06) | 92.60 (1.85) | 6.91 (0.00) | 4.34 (0.06) | 4.11 (0.06) | 72.77 (1.45) | 6.91 (0.00) |
| 100 | 5.45 (0.04) | 5.44 (0.04) | 159.43 (4.91) | 8.07 (0.00) | 5.68 (0.05) | 5.58 (0.05) | 124.90 (3.18) | 8.07 (0.00) |
| 200 | 8.09 (0.05) | 7.99 (0.05) | 256.12 (11.01) | 10.93 (0.00) | 7.98 (0.07) | 7.95 (0.07) | 200.10 (5.37) | 10.93 (0.00) |
| Spectral norm loss | | | | | | | | |
| 50 | 2.50 (0.05) | 2.38 (0.05) | 68.27 (1.51) | 4.92 (0.00) | 2.53 (0.02) | 2.43 (0.02) | 51.83 (1.17) | 4.92 (0.00) |
| 100 | 3.25 (0.05) | 3.19 (0.05) | 111.79 (3.66) | 5.46 (0.00) | 3.07 (0.02) | 3.03 (0.02) | 83.24 (2.42) | 5.46 (0.00) |
| 200 | 3.86 (0.03) | 3.87 (0.02) | 170.37 (7.79) | 6.43 (0.00) | 3.94 (0.02) | 3.91 (0.02) | 122.81 (4.05) | 6.43 (0.00) |
| Frobenius norm loss | | | | | | | | |
| 50 | 6.17 (0.06) | 5.96 (0.06) | 70.52 (1.46) | 25.98 (0.00) | 8.82 (0.03) | 8.45 (0.04) | 54.44 (1.12) | 25.99 (0.00) |
| 100 | 9.40 (0.06) | 9.32 (0.06) | 117.87 (3.51) | 38.38 (0.00) | 13.92 (0.03) | 13.67 (0.04) | 90.22 (2.30) | 38.38 (0.00) |
| 200 | 13.55 (0.08) | 13.54 (0.09) | 185.38 (7.65) | 55.78 (0.00) | 21.64 (0.04) | 21.45 (0.04) | 140.56 (3.83) | 55.78 (0.00) |
| True positive rate | | | | | | | | |
| 50 | 0.65 (0.01) | 0.68 (0.01) | 0.99 (0.00) | 0.76 (0.02) | 0.94 (0.01) | 0.95 (0.00) | 0.99 (0.00) | 0.93 (0.00) |
| 100 | 0.60 (0.00) | 0.61 (0.00) | 0.97 (0.01) | 0.39 (0.02) | 0.91 (0.00) | 0.92 (0.00) | 0.93 (0.01) | 0.89 (0.01) |
| 200 | 0.60 (0.00) | 0.61 (0.00) | 0.94 (0.01) | 0.28 (0.02) | 0.84 (0.00) | 0.84 (0.00) | 0.93 (0.00) | 0.88 (0.01) |
| False positive rate | | | | | | | | |
| 50 | 0.00 (0.00) | 0.00 (0.00) | 0.98 (0.01) | 0.48 (0.03) | 0.12 (0.00) | 0.11 (0.00) | 0.95 (0.00) | 0.72 (0.01) |
| 100 | 0.00 (0.00) | 0.00 (0.00) | 0.94 (0.02) | 0.10 (0.01) | 0.07 (0.00) | 0.07 (0.00) | 0.92 (0.01) | 0.65 (0.01) |
| 200 | 0.00 (0.00) | 0.00 (0.00) | 0.86 (0.03) | 0.06 (0.01) | 0.04 (0.00) | 0.04 (0.00) | 0.86 (0.02) | 0.61 (0.01) |

Figure 2.2: ROC curves for four methods in Model 2 with normal-related distribution (top panel) and gamma-related distribution (bottom panel).

## 2.6. Gut Microbiome Data Analysis

The gut microbiome plays a critical role in energy extraction from the diet and interacts with the immune system to exert a profound influence on human health and disease. Despite an emerging interest in characterizing the ecology of human-associated microbial communities, the complex interactions among microbial taxa remain poorly understood (Coyte, Schluter, and Foster, 2015). We now illustrate the proposed method by applying it to a human gut microbiome dataset described by Wu et al., (2011), which was collected from a cross-sectional study of 98 healthy individuals at the University of Pennsylvania. DNA from stool samples of these subjects were analyzed by 454/Roche pyrosequencing of 16S rRNA gene segments, resulting in an average of 9265 reads per sample, with a standard deviation of 3864. Taxonomic assignment yielded 3068 operational taxonomic units, which were further combined into 87 genera that appeared in at least one sample. Demographic information, including body mass index (BMI), was also collected from the subjects. We are interested in identifying and comparing the correlation structures among bacterial genera between lean and obese subjects. We therefore divided the dataset into a lean group ($\mathrm{BMI} <$ $25$, $n = 63$) and an obese group ($\mathrm{BMI} \geq 25$, $n = 35$), and focused on the $p = 40$ bacterial genera that appeared in at least four samples in each group. The count data were transformed into compositions after zero counts were replaced by 0.5.

We applied the COAT method with the soft thresholding rule to each group, and used tenfold cross-validation to select the tuning parameter. The resulting estimate was represented by a correlation network among the bacterial genera with each edge representing a nonzero correlation. To assess the stability of support recovery, we further generated 100 bootstrap samples for each group and repeated the thresholding procedure on each sample. The stability of the correlation network was measured by the average proportion of edges reproduced by each bootstrap replicate. Finally, we retained only the edges in the correlation network that were reproduced in at least 80 bootstrap replicates. The numbers of positive and negative correlations and the stability of correlation networks are reported in Table 2.3; the results for the two naive thresholding estimators $\widehat{\Omega}_l$ and $\widehat{\Omega}_c$ are also included for comparison. We see that the COAT method achieves the highest stability among the three methods and has the most edges passing the stability test. The correlation network identified by $\widehat{\Omega}_l$ has substantially fewer negative correlations than the other two methods, which is likely due to the severe upward bias observed in Figure 2.1. The correlation network identified by $\widehat{\Omega}_c$ is
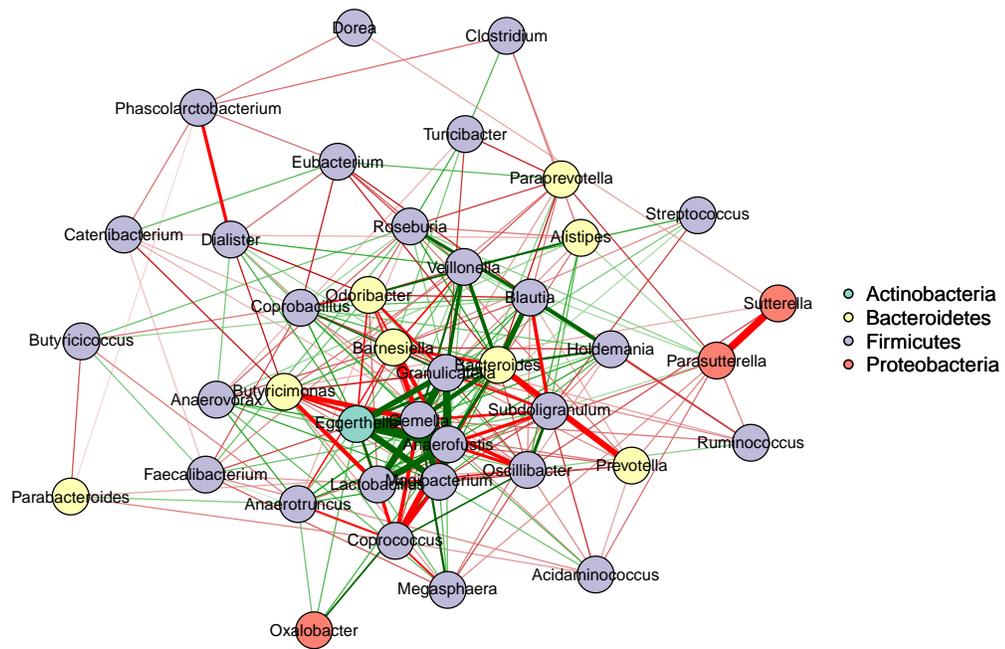
Table 2.3: Numbers of positive and negative correlations and stability of correlation networks for three methods applied to the gut microbiome data

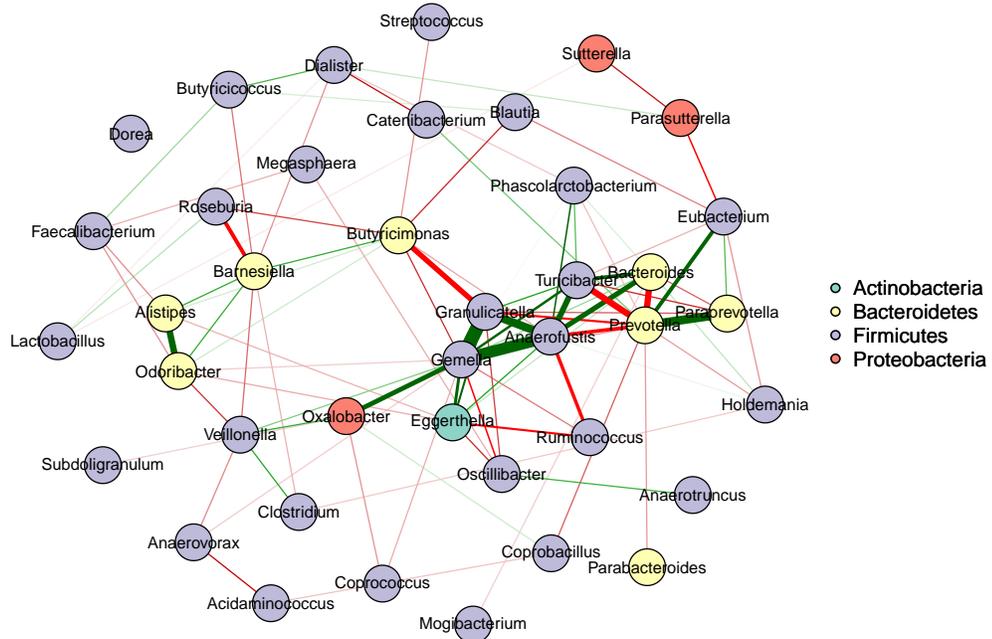| | Lean | | | Obese | | |
|---|---|---|---|---|---|---|
| | $\widehat{\Omega}$ | $\widehat{\Omega}_l$ | $\widehat{\Omega}_c$ | $\widehat{\Omega}$ | $\widehat{\Omega}_l$ | $\widehat{\Omega}_c$ |
| Positive correlations | 111 | 108 | 119 | 41 | 34 | 31 |
| Negative correlations | 134 | 55 | 95 | 55 | 11 | 43 |
| Network stability | 0.83 | 0.68 | 0.67 | 0.87 | 0.62 | 0.54 |

the least stable.

The correlation networks identified by the COAT method for the two groups are displayed in Figure 2.3. Clearly, the networks for the lean and obese groups show markedly different architecture, indicating that the obese microbiome is less modular with less complex interactions between the modules. This phenomenon has been demonstrated by previous studies and is possibly due to adaptation of the microbiome to low-diversity environments (Greenblum, Turnbaugh, and Borenstein, 2012). Table 2.3 and Figure 2.3 also suggest that the gut microbial network tends to contain more competitive (negative) interactions than cooperative (positive) ones, which seems consistent with the recent finding that the ecological stability of the gut microbiome can be attributed to the benefits from limiting positive feedbacks and dampening cooperative networks (Coyte, Schluter, and Foster, 2015).

A closer inspection of the correlation networks identifies *Bacteroides* and *Prevotella* as two key genera of the gut microbiome. The abundances of these two genera are well known to distinguish two gut microbial enterotypes, which are strongly associated with long-term dietary patterns (Arumugam et al., 2011; Wu et al., 2011). The negative correlations between *Bacteroides* and *Prevotella* ($-0.404$ in the lean group and $-0.296$ in the obese group) are well explained by the diet-dependent enterotypes and the within-body separation of the two genera (Jordán et al., 2015). Moreover, recent studies have suggested several keystone species belonging to the genus *Bacteroides*, through which the structure of gut microbial communities may be influenced by small perturbations (Fisher and Mehta, 2014). Also, the Firmicutes-enriched microbiome has been found to hold greater metabolic potential than the Bacteroidetes-enriched microbiome for more efficient energy harvest from the diet (Turnbaugh et al., 2006). Figure 2.3 seems to support these findings, in view of the central position of *Bacteroides* in the networks and its strong correlations with a few

(a) Lean



(b) Obese

Figure 2.3: Correlation networks identified by the COAT method for the lean and obese groups in the gut microbiome data. Positive and negative correlations are displayed in green and red, respectively. The thickness of edges indicates the magnitude of correlations.

genera belonging to the Firmicutes. Such patterns, however, are less clearly seen in the correlation networks identified by the other two methods.

## 2.7. Discussion

Understanding the dependence structure among microbial taxa within a community, including co-occurrence and co-exclusion relationships between microbial taxa, is an important problem in microbiome research. Such structures provide biological insights into the community dynamics and factors that change the community structures. To overcome the difficulties arising from the unit-sum constraint of the observed compositional data, we have developed a COAT method to estimate the sparse covariance matrix of the latent log-basis components. Our method is based on a decomposition of the variation matrix into a rank-2 component and a sparse component. The resulting procedure is equivalent to thresholding the sample centered log-ratio covariance matrix, and thus is optimization-free and scalable for high-dimensional data.

Our simulation results demonstrate that the COAT method performs almost as well as the oracle thresholding estimator that knew the latent basis components, and outperforms some naive thresholding estimators by a large margin. These improvements are more pronounced when the basis components have a skewed distribution, as is often observed in microbiome studies. In the application to gut microbiome data, the COAT method leads to more stable and biologically more interpretable results for comparing the dependence structures of lean and obese microbiomes.

We have provided conditions for the approximate and exact identifiability of the covariance parameters, and have established rates of convergence and support recovery guarantees for the COAT estimator. The rate of convergence includes an extra term of $O_p(s_0(p)(s_0(p)/p)^{1-q})$ in addition to the usual minimax optimal rate of convergence for sparse covariance estimation. The extra term represents an approximation error due to using $\Gamma_0$ as a proxy for $\Omega_0$, which vanishes under mild assumptions as the dimensionality increases.

The proposed methodology may be extended in several ways. First, it would be possible to develop a joint optimization procedure based on the decomposition (2.5). For example, one may consider

42

the regularized estimator

$$\widehat{\boldsymbol{\Omega}}_{\mathrm{reg}} = \arg\min_{\boldsymbol{\Omega}}\{\|\widehat{\mathbf{T}} - \boldsymbol{\omega}\mathbf{1}^T - \mathbf{1}\boldsymbol{\omega}^T + 2\boldsymbol{\Omega}\|_F^2 + P_\lambda(\boldsymbol{\Omega})\},$$

where $\boldsymbol{\omega} = \mathrm{diag}(\boldsymbol{\Omega})$ and $P_\lambda(\cdot)$ is a sparsity-inducing penalty function. The COAT estimator can be viewed as a one-step approximation to $\widehat{\boldsymbol{\Omega}}_{\mathrm{reg}}$ with appropriately chosen penalty function and initial value $\widehat{\boldsymbol{\Omega}} = \mathbf{0}$. Solving the full optimization problem is computationally more expensive but is expected to improve on the performance of the COAT estimator. Another worthwhile extension would be to deal with zero counts directly. One may, in principle, combine the ideas presented here with models that account for sampling and structural zeros. The issues of identifiability and computational feasibility are the major concerns with such extensions.

# CHAPTER 3

## TWO-SAMPLE MEAN TESTS FOR HIGH-DIMENSIONAL COMPOSITIONAL DATA

Motivated by microbiome and metagenomic research, in this chapter, we consider a two-sample testing problem for high-dimensional compositional data and formulate a testable hypothesis of compositional equivalence for the means of two latent log-basis vectors. We propose a test for such compositional equivalence through the centered log-ratio transformation of the compositions. The asymptotic null distribution of the test statistic is derived and the power of the test against sparse alternatives is studied. A modified test for paired observations is also developed. Simulations show that the proposed tests can be significantly more powerful than existing tests that are applied to the raw and log-transformed compositional data. The usefulness of the proposed tests is illustrated by applications to test for differences in gut microbiome composition between lean and obese individuals and changes between different time points during treatment in Crohn's disease patients.

Compositional data, which belong to the unit simplex sample space, are ubiquitous in many scientific disciplines such as geology, economics, genomics, and machine learning. This paper is motivated by microbiome and metagenomic research, where the relative abundances of hundreds to thousands of bacterial taxa on a few tens to hundreds of individuals are available for analysis (The Human Microbiome Project Consortium, 2012b). Due to varying amounts of DNA generating materials across different samples, sequencing read counts are often normalized to relative abundances; the resulting data are therefore compositional (Li, 2015). One fundamental problem in microbiome data analysis is to test whether two populations have the same microbiome composition, which can be viewed as a two-sample mean testing problem for high-dimensional compositional data. Owing to the key feature that the components of a composition must sum to one, applying standard multivariate statistical methods intended for unconstrained data directly to compositional data may result in inappropriate or misleading inferences (Aitchison, 2003).

Various methods for compositional data analysis have been developed in the literature since the seminal work of Aitchison, 1982. Most existing methods for the two-sample mean testing problem, however, deal only with the low-dimensional setting where the dimensionality is fixed or much s-

maller than the sample size; see, e.g., the generalized likelihood ratio tests discussed in Aitchison, (2003 §7.5). In this paper, we consider the two-sample mean testing problem for high-dimensional compositional data, where compositions in the $(p-1)$-dimensional unit simplex $\mathcal{S}^{p-1}$ are thought of as arising from latent basis vectors in the $p$-dimensional positive orthant $\mathbb{R}_+^p$. In microbiome studies, the basis components may represent the true abundances of bacterial taxa in a microbial community such as the gut of a healthy individual (Li, 2015). To circumvent the nonidentifiability issue associated with the basis vectors, we formulate a testable hypothesis of compositional equivalence for the means of two log-basis vectors. We then propose a test for such compositional equivalence through the centered log-ratio transformation of the compositions. The proposed test thus honors the principle of scale invariance, which is crucial for compositional data analysis. We emphasize that we are adopting the extrinsic analysis approach, which leads to biologically meaningful interpretations and is in contrast to intrinsic analysis where no such basis exists and interest focuses on the composition itself (Aitchison, 1982).

Development of tests for the equality of two population means in high-dimensional settings has received much attention recently; see, e.g., Bai and Saranadasa, (1996), Srivastava, (2009), Srivastava, (2009), Chen and Qin, (2010) and Cai, Liu, and Xia, (2014). These high-dimensional tests, however, are not directly applicable to compositional data because the required regularity conditions are generally not met. For example, the covariance matrix of compositional variables is singular, thereby violating the usual assumptions on the eigenvalues of the covariance matrix such as Condition 1 in Cai, Liu, and Xia, (2014). Our assumptions are instead made on the latent log-basis vectors, which are free of the simplex constraint. We show that, under mild conditions, the centered log-ratio transformed variables satisfy certain desired properties, which in turn guarantee the validity of the proposed test. The asymptotic null distribution of the test statistic is then derived and the power of the test against sparse alternatives is investigated.

The proposed test is further extended to the setting with paired observations or repeatedly measured compositions. Extensive simulations and applications to two microbiome datasets are provided to illustrate the proposed methodology. All proofs are given in the Appendix.

## 3.1. A testable hypothesis of compositional equivalence

Denote by $X^{(k)} = (X_1^{(k)}, \ldots, X_{n_k}^{(k)})^\mathrm{T}$ the observed $n_k \times p$ data matrices for group $k$ ($k = 1, 2$), where $X_i^{(k)}$ represent compositions that lie in the $p - 1$-simplex $\mathcal{S}^{p-1} = \{(x_1, \ldots, x_p) : x_j > 0 \, (j = 1, \ldots, p), \sum_{j=1}^p x_j = 1\}$, and the two-sample sizes $n_1$ and $n_2$ are supposed to be comparable, that is, the ratio $n_1/n_2$ is always constant. Let $n = \max(n_1, n_2)$. We assume that the compositional variables arise from a vector of latent variables, which we call the basis. For microbiome data, the basis components may refer to the true abundances of bacterial taxa in a microbial community. Denote by $W^{(k)} = (W_1^{(k)}, \ldots, W_{n_k}^{(k)})^\mathrm{T}$ the $n_k \times p$ matrices of unobserved bases, which generate the observed compositional data via the normalization

$$X_{ij}^{(k)} = W_{ij}^{(k)} \Big/ \sum_{\ell=1}^p W_{i\ell}^{(k)} \quad (i = 1, \ldots, n_k; \; j = 1, \ldots, p; \; k = 1, 2),$$

where $X_{ij}^{(k)}$ and $W_{ij}^{(k)} > 0$ are the $j$th components of $X_i^{(k)}$ and $W_i^{(k)}$, respectively.

Denote by $Z_i^{(k)} = \log W_i^{(k)}$ the log-basis vectors, where the logarithm applies componentwise. Suppose that $Z_1^{(k)}, \ldots, Z_{n_k}^{(k)}$ are independent and identically distributed from a distribution with mean $\mu_k$ and covariance matrix $\Omega$ ($k = 1, 2$). One might attempt to test the hypotheses

$$H_0 : \mu_1 = \mu_2 \quad \text{versus} \quad H_1 : \mu_1 \neq \mu_2. \tag{3.1}$$

These hypotheses, however, are not testable through the observed compositional data $X^{(k)}$ ($k = 1, 2$). Clearly, a basis is determined by its composition only up to a multiplicative factor, and the set of bases giving rise to a composition $x \in \mathcal{S}^{p-1}$ forms the equivalence class $\mathcal{W}(x) = \{(tx_1, \ldots, tx_p) : t > 0\}$ (Aitchison, 2003 p. 32). As an immediate consequence, a log-basis vector is determined by the resulting composition only up to an additive constant, and the set of log-basis vectors corresponding to $x$ constitutes the equivalence class $\mathcal{Z}(x) = \{(\log x_1 + c, \ldots, \log x_p + c) : c \in \mathbb{R}\}$. We therefore introduce the following definition.

**Definition 1.** Two log-basis vectors $z_1$ and $z_2$ are said to be compositionally equivalent if their components differ by a constant $c \in \mathbb{R}$, i.e., $z_1 = z_2 + c1_p$, where $1_p$ is the $p$-vector of 1s.

Now, instead of testing the hypotheses in (3.1), we propose to test

$$H_0 : \mu_1 = \mu_2 + c1_p \text{ for some } c \in \mathbb{R} \quad \text{versus} \quad H_1 : \mu_1 \neq \mu_2 + c1_p \text{ for all } c \in \mathbb{R}, \qquad (3.2)$$

which are testable using only the observed compositional data. Clearly, $H_0$ in (3.1) implies $H_0$ in (3.2), so that rejecting the latter would lead to rejection of the former. Note, however, that $H_0$ in (3.2) neither implies nor is implied by $E(X_1^{(1)}) = E(X_1^{(2)})$ or $E(\log X_1^{(1)}) = E(\log X_1^{(2)})$. We do not consider the latter two hypotheses because they are not scale invariant, whereas we will derive in the next section an equivalent form of $H_0$ in (3.2), from which its scale invariance is obvious. Moreover, these two hypotheses do not allow us to obtain biological interpretations in terms of the true underlying abundances.

## 3.2. The centered log-ratio transformation and a test for compositional equivalence

### 3.2.1. The centered log-ratio transformation and an equivalent hypothesis

The unit-sum constraint entails that compositional variables must not vary independently, making many covariance-based multivariate analysis methods inapplicable. Aitchison, 1982 proposed to relax the constraint by performing statistical analysis through log-ratios. Among various forms of log-ratio transformations, the centered log-ratio transformation possesses some attractive features and has been widely used in practice. For the observed compositional data $X^{(k)}$ ($k = 1, 2$), the centered log-ratios are defined by

$$Y_{ij}^{(k)} = \log\{X_{ij}^{(k)}/g(X_i^{(k)})\} \quad (i = 1, \ldots, n_k; \ j = 1, \ldots, p; \ k = 1, 2), \qquad (3.3)$$

where $g(x) = (\prod_{i=1}^{p} x_i)^{1/p}$ denotes the geometric mean of a vector $x = (x_1, \ldots, x_p)^{\mathrm{T}}$. The relationship (3.3) can be expressed in matrix form as

$$Y_i^{(k)} = G \log X_i^{(k)} \quad (i = 1, \ldots, n_k; \ k = 1, 2), \qquad (3.4)$$

where $Y_i^{(k)}$ are the centered log-ratio vectors, $G = I_p - p^{-1}1_p 1_p^{\mathrm{T}}$, and $I_p$ is the $p \times p$ identity matrix.

Let $\nu_k = E(Y_1^{(k)})$ ($k = 1, 2$). In view of (3.4) and the scale invariance of the centered log-ratios, we

have

$$\nu_k = E(G \log X_1^{(k)}) = E(G \log W_1^{(k)}) = GE(\log W_1^{(k)}) = GE(Z_1^{(k)}) = G\mu_k \quad (k = 1, 2).$$

Note that the matrix $G$ has rank $p - 1$ and hence a null space of dimension 1, $\mathcal{N}(G) \equiv \{x \in \mathbb{R}^p : Gx = 0\} = \{c1_p : c \in \mathbb{R}\}$. As a result, $\nu_1 = \nu_2$ if and only if $\mu_1 = \mu_2 + c1_p$ for some $c \in \mathbb{R}$. Therefore, testing the hypotheses in (3.2) is equivalent to testing

$$H_0 : \nu_1 = \nu_2 \quad \text{versus} \quad H_1 : \nu_1 \neq \nu_2. \tag{3.5}$$

Despite this equivalence, the hypotheses in (3.2) are meaningful only when the bases exist, which is the case in microbiome studies. On the other hand, the hypotheses in (3.5) concern only the compositions through the centered log-ratios, from which its scale invariance and testability using the observed compositional data are evident.

### 3.2.2. A test for compositional equivalence

A natural test statistic for testing $H_0$ in (3.5), and hence $H_0$ in (3.2), would be based on the differences $\bar{Y}_j^{(1)} - \bar{Y}_j^{(2)}$, where $\bar{Y}_j^{(k)} = \sum_{i=1}^{n_k} Y_{ij}^{(k)}/n_k$ are the sample means of the centered log-ratios. Moreover, it is well known that tests based on maximum type statistics are generally more powerful than those based on sum-of-squares type statistics against sparse alternatives (Cai, Liu, and Xia, 2014). Since in microbiome studies we are mainly interested in the sparse setting where only a small number of taxa may have different mean abundances between the two groups, we consider the test statistic

$$M_n = \frac{n_1 n_2}{n_1 + n_2} \max_{1 \leq j \leq p} \frac{(\bar{Y}_j^{(1)} - \bar{Y}_j^{(2)})^2}{\hat{\gamma}_{jj}}, \tag{3.6}$$

where $\hat{\gamma}_{jj} = \sum_{k=1}^{2} \sum_{i=1}^{n_k} (Y_{ij}^{(k)} - \bar{Y}_j^{(k)})^2/(n_1 + n_2)$ are the pooled sample variances.

The asymptotic behavior of $M_n$ will be investigated in the next section. Specifically, under suitable conditions on the log-basis variables $Z_{1j}^{(k)}$, we will show that the centered log-ratio transformed variables $Y_{1j}^{(k)}$ are only weakly dependent and satisfy certain concentration properties. As a result, the null distribution of $M_n - 2 \log p + \log \log p$ is asymptotically a type I extreme value distribution. The test defined by $\Phi_\alpha = I(M_n \geq q_\alpha + 2 \log p - \log \log p)$, where $q_\alpha = -\log \pi - 2 \log \log(1 - \alpha)^{-1}$ is the $(1 - \alpha)$-quantile of the type I extreme value distribution, is then an asymptotic $\alpha$-level test for

testing $H_0$ in (3.2) or (3.5).

Although $M_n$ is similar to the test statistic $M_I$ defined in Cai, Liu, and Xia, 2014, their theoretical analyses are radically different, since our assumptions are not imposed on the observed variables. Besides, the test statistic based on a linear transformation by the precision matrix proposed by Cai, Liu, and Xia, 2014 is not considered here, because the covariance matrix of $Y_1^{(k)}$ is singular and its precision matrix is not well defined.

## 3.3. Theoretical Results for the CLR Transformation-based Global Test

### 3.3.1. Covariance and correlation of CLR-transformed compositions

The relationships between the covariance and correlation for the true log of the abundance $\log W_d$ and the log-ratio vectors $Y_d$, $(d = 1, 2)$ are first studied. Let the centered log-ratio covariance matrix be $\Gamma = (\gamma_{i,j}) := \text{cov}(y_{i,d}^\star, y_{j,d}^\star)$ for $1 \leq i \leq j \leq p$ and $d = 1, 2$. Then $\Gamma$ is related to $\Omega$ via the following relationship (Aitchison, 2003),

$$\Omega \to \Gamma : \ \Gamma = \mathrm{G}\Omega\mathrm{G}^{\mathrm{T}}. \tag{3.7}$$

In the high dimensional setting when $p \gg 0$, $G = I_p - p^{-1}J_p \approx I_p$, suggesting that the covariance and correlation structure of $\log W_d$ and $Y_d$ are similar. Such relationships serve as the basis for the theoretical validity of the testing procedure. Denote the correlation of $\log W_d$ and $Y_d$ as $R$ and $R^{\mathrm{clr}}$, i.e., $R = (r_{i,j}) = \text{corr}(\log w_{i,d}^\star, \log w_{j,d}^\star)$ and $R^{\mathrm{clr}} = (r_{i,j}^{\mathrm{clr}}) = \text{corr}(y_{i,d}^\star, y_{j,d}^\star)$. The following assumptions are made on the correlation matrix $R$.

**Condition 1.** $\max_{1 \leq i < j \leq p} |r_{i,j}| \leq r_1 < 1$ for some constant $0 < r_1 < 1$.

**Condition 2.** $\max_{1 \leq j \leq p} \sum_{i=1}^p r_{i,j}^2 \leq r_2 < \infty$ for some constant $r_2 > 0$.

**Condition 3.** $0 < 1/\tau \leq \omega_{i,i} \leq \tau < \infty$, for any $i = 1, \cdots, p$, where $\tau > 0$ is a constant.

Condition $1$ is mild since $\Omega$ is non-singular. Both Conditions $2$ and $3$ are standard assumptions encountered in high dimensional settings. Condition $2$ guarantees weak correlations among majority of the variables, which is reasonable in the context of microbiome study as only a small number of bacterial species in human microbiome may have strong cooperative and competitive relationships. Condition $3$ assumes a uniform variance. Under these conditions, the following properties of the correlation and covariance matrices hold.

**Proposition 1.** Under Condition 2, let $r_3 := \max_{1 \leq j \leq p} \sum_{i=1}^p |r_{i,j}|$, then

$$r_3 = \max_{1 \leq j \leq p} \sum_{i=1}^p |r_{i,j}| \leq p^{1/2} \cdot \max_{1 \leq j \leq p} (\sum_{i=1}^p r_{i,j}^2)^{1/2} = O(p^{1/2}).$$

**Proposition 2.** Under Conditions 1, 2 and 3, from Proposition 1, the difference between $\Omega$ and $\Gamma$ is bounded by

$$\|\Omega - \Gamma\|_{\max} \leq 3p^{-1} \cdot \max_{1 \leq i \leq p} \omega_{i,i} \cdot \max_{1 \leq j \leq p} \sum_{i=1}^p |r_{i,j}| = o(1),$$

which implies $\Omega$ and $\Gamma$ are approximately identical as $p \to \infty$. Therefore, for sufficiently large $p$,

$$\min_{1 \leq i \leq p} \gamma_{i,i} \geq \min_{1 \leq i \leq p} \omega_{i,i} - \|\Omega - \Gamma\|_{\max} \geq 1/(2\tau). \tag{3.8}$$

Propositions 1 and 2 bound the difference between $R$ and $R^{\mathrm{clr}}$, which is given in the following theorem.

**Theorem 6.** *Suppose Conditions 1, 2 and 3 hold respectively for the correlation matrix $R$ and the covariance matrix $\Omega$, as $p$ is sufficiently large,*

$$\|R^{\mathrm{clr}} - R\|_{\max} = o(1), \ \max_{1 \leq j \leq p} \left| \sum_{i=1}^p (r_{i,j}^{\mathrm{clr}})^2 - \sum_{i=1}^p r_{i,j}^2 \right| = O(1). \tag{3.9}$$

Equations (3.9), combined with Conditions 1 and 2 guarantee that a similar correlation structure holds for $R^{\mathrm{clr}}$.

**Corollary 1.** Suppose Condition 1, 2 and 3 hold respectively on $R$ and $\Omega$, then, for sufficiently large $p$, there exists some constant $r_4 > 0$ and $r_5 > 0$, such that

$$\max_{1 \leq i < j \leq p} r_{i,j}^{\mathrm{clr}} \leq r_4 < 1, \ \text{and} \ \max_{1 \leq j \leq p} \sum_{i=1}^p (r_{i,j}^{\mathrm{clr}})^2 \leq r_5 < \infty. \tag{3.10}$$

### 3.3.2. Tail distribution of the CLR-transformed compositions

This section investigates the concentration property the CLR-transformed variable $Y_d$ based on the assumptions for the random vector $\log W_d$. Specifically,

**Condition 4.** (Sub-Gaussian-type tails). Random vector $\log W_d^{\star}$ $(d = 1, 2)$ follows sub-Gaussian-type tails (Cai, Liu, and Xia, 2014), if $\log p = o(n_d^{1/4})$ and there exist some constants $\eta > 0$ and $K > 0$ such that,

$$E(\exp(\eta(\log W_{i,d}^{\star} - \mu_{i,d})^2/\omega_{i,i})) \leq K, \text{ for } 1 \leq i \leq p.$$

**Condition 5.** (Polynomial-type tails). Random vector $\log W_d^{\star}$ $(d = 1, 2)$ follow polynomial-type tails (Cai, Liu, and Luo, 2011), if, for some constant $\gamma_0 > 0$, $p = O(n^{\gamma_0})$ and for some constants $\epsilon > 0$ and $K > 0$ such that,

$$E\left|(\log W_{i,d}^{\star} - \mu_{i,d})/\omega_{i,i}^{1/2}\right|^{4\gamma_0 + 4 + \epsilon} \leq K, \text{ for } 1 \leq i \leq p.$$

Conditions $4$ and $5$ are assumed for $\log W_d^{\star}$. The tail distribution of the CLR-transformed observation $Y_d$ $(d = 1, 2)$ has the following probability inequality, as well as the rate of convergence of its sample variance.

**Theorem 7.** *Suppose that the correlation/variance structure (Conditions $2$ and $3$), and tail probability of its distribution (Condition $4$ (or $5$)) hold. Then, there exists $\tau_n = o(n_d^{1/2}/(\log p)^{3/2})$ such that,*

$$\mathrm{pr}\big(\max_{1 \leq k \leq n_d, 1 \leq i \leq p} |y_{k,i,d} - \nu_{i,d}| / \gamma_{i,i}^{1/2} \leq \tau_n\big) \to 1, \text{ as } n_d, p \to \infty. \tag{3.11}$$

*In addition, uniformly in $1 \leq i \leq p$, the rate of convergence of the pooled sample variance $\gamma_{i,i}$ is*

$$|\widehat{\gamma}_{i,i} - \gamma_{i,i}| = O_p\left\{(\log p/n)^{1/2}\right\} \gamma_{i,i}. \tag{3.12}$$

### 3.3.3. Asymptotic null distribution of $M_n$ and power analysis

The following theorem states that although no assumptions are made directly on $Y_d^\star$ $(d = 1, 2)$, Conditions $1$-$5$ guarantee that the asymptotic null distribution of $M_n$ follows a type I extreme value distribution and its size is effectively controlled.

**Theorem 8.** *Suppose that equations* $(3.10)$, $(3.11)$ *and* $(3.12)$ *hold, which are guaranteed by Conditions* $1 - 3$ *and* $4$ *(or* $5$*). Under the null hypothesis* $H_0 : \nu_1 = \nu_2$, *for any* $t \in \mathbb{R}$,

$$\mathrm{pr}(M_n - 2\log p + \log\log p \leq t) \to \exp\left\{-\pi^{-1/2}\exp(-t/2)\right\}, \text{ as } n_1, n_2, p \to \infty. \tag{3.13}$$

*Besides, if the* $\alpha-$*level test* $\Phi_\alpha$ *is defined by* $\Phi_\alpha = I(M_n \geq q_\alpha + 2\log p - \log\log p)$, *the probability of type I error is controlled by*

$$\mathrm{pr}(\text{\textit{Type I error}}) = \mathrm{pr}_{\mathrm{H}_0}(\Phi_\alpha = 1) \leq -\log(1-\alpha) + o(1), \text{ for any } 0 < \alpha < 1. \tag{3.14}$$

To study the power of the test defined in Theorem (8), consider the alternative hypothesis

$$H_1 : \nu_1 - \nu_2 \in S(k_p) \text{ with } k_p = p^r, 0 \leq r < 1, \tag{3.15}$$

where the non-zero support is randomly and uniformly drawn from $\{1, \cdots, p\}$ with the magnitude given by

$$\max_{1 \leq i \leq p} |\nu_{i,1} - \nu_{i,2}|/\gamma_{i,i}^{1/2} = \{2\beta\log p(1/n_1 + 1/n_2)\}^{1/2}, \text{ where } \beta \in (0,1). \tag{3.16}$$

This hypothesis $H_1$ can be rephrased using the parameters for the basis counts. For $\forall e \in S(k_p)$ with the support index defined as $\mathrm{S}$,

$$\nu_1 - \nu_2 = e \iff \begin{cases} \exists c^\star \in \mathbb{R}, \ s.t. \ \mu_{\mathrm{S},1} = \mu_{\mathrm{S},2} + c^\star \times 1_{p^r}, \\ \forall c \in \mathbb{R}, \ \mu_{\mathrm{S^c},1} \neq \mu_{\mathrm{S^c},2} + c \times 1_{p-p^r}, \end{cases}$$

where $\mu_{\mathrm{S},d}$ and $\mu_{\mathrm{S^c},d}$ are the sub-vectors of $\mu$ corresponding to the support $\mathrm{S}$ and its complement. Since the observed composition is in high dimensional space, it is reasonably to assume the means

of the basis count from two groups compositionally differ only in a small number of the coordinates. In addition, under Conditions $1$, $2$ and $3$, equation $(3.16)$ is equivalent to

$$\max_{i \in \mathbb{S}} |\mu_{i,1} - \mu_{i,2} - c^\star| / \omega_{i,i}^{1/2} = \{2\beta \log p(1/n_1 + 1/n_2)\}^{1/2} \cdot (1 + o(1)).$$

The following theorem provides the results on test power under the alternative specified in (3.15).

**Theorem 9.** *Under $H_1$ given by* $(3.15)$, *for some $\epsilon > 0$, we have*

$$\lim_{p \to \infty} \mathrm{pr}_{H_1}(\Phi_\alpha = 1) = 1, \text{ if } \beta \geq (1 - \sqrt{r})^2 + \epsilon, \tag{3.17}$$

$$\overline{\lim_{p \to \infty}} \, \mathrm{pr}_{H_1}(\Phi_\alpha = 1) \leq \alpha, \text{ if } \beta < (1 - \sqrt{r})^2. \tag{3.18}$$

## 3.4. Two-sample Test for Paired Observations

Test of compositional equality for paired observations $\{(X_{i,1}, X_{i,2})\}_{i=1}^n$, where $X_{i,1}$ and $X_{i,2}$ are two $p$-dimensional compositional observations on a subject $i$ before and after a treatment, requires slight modification. Suppose that $\{(W_{i,1}, W_{i,2})\}_{i=1}^n$ are the corresponding $p-$variate basis vectors. Let $\widetilde{W}_i = (\widetilde{w}_{i,1}, \cdots, \widetilde{w}_{i,p})^\mathrm{T}$ be the element-wise ratio of $i^\mathrm{th}$ samples $W_{i,1}$ and $W_{i,2}$: $\widetilde{w}_{i,j} = w_{i,j,1}/w_{i,j,2}$ for $1 \leq j \leq p$, and $\log \widetilde{W}_i := \log W_{i,1} - \log W_{i,2}$. For $i = 1 \cdots, n$, let $\widetilde{Y}_i = (\widetilde{y}_{i,1}, \cdots, \widetilde{y}_{i,p})^\mathrm{T}$ be the corresponding CLR-transformed random vector. Through $(3.4)$ and the principle of scale invariance that $G \log W_{i,d} = G \log X_{i,d}$, it can be written by the observations $X_d$ as

$$\widetilde{Y}_i := G \log \widetilde{W}_i = G \log W_{i,1} - G \log W_{i,2} = G \log X_{i,1} - G \log X_{i,2}, \text{ for } i = 1, \cdots, n.$$

The compositional equality null hypothesis can be written in term of the mean of the difference of the centered log-ratio variables,

$$H_0 : \widetilde{\nu} = 0 \quad \textit{vs} \quad H_1 : \widetilde{\nu} \in S(k_p) \text{ with } k_p = p^r, \ 0 \leq r < 1, \tag{3.19}$$

where $\widetilde{\nu} = E\widetilde{Y}_i$ $(i = 1, \cdots, n)$. The non-zero locations in $S(k_p)$ are randomly uniformly drawn from $\{1, \cdots, p\}$ and the magnitude of support in $\widetilde{\nu}$ is given by $(3.16)$, where $\nu_{i,1} - \nu_{i,2}$ and $\gamma_{i,i}$ is

respectively replaced by $\widetilde{\nu}_i$ and $\widetilde{\gamma}_{i,i} = \mathrm{var}(\widetilde{y}_{k,i})$ $(k = 1, \cdots, n)$.

The following test statistic is proposed for testing the null (3.19),

$$\widetilde{M}_n = \max_{1 \leq i \leq p} \frac{\bar{\widetilde{y}}_i^2}{\widehat{\gamma}_{i,i}/n},$$

where $\bar{\widetilde{y}}_i = n^{-1} \sum_{j=1}^n \widetilde{y}_{j,i}$ is the sample mean of $i^{\mathrm{th}}$ variable in CLR-transformed observation $\widetilde{Y} = (\widetilde{Y}_1, \cdots, \widetilde{Y}_n)^{\mathrm{T}}$, and $\widehat{\Gamma} = (\widehat{\gamma}_{i,j}) = (n^{-1} \sum_{l=1}^n (y_{l,i} - \bar{\widetilde{y}}_i)(y_{l,j} - \bar{\widetilde{y}}_j))$ is its sample covariance. Let $\widetilde{W}$ be drawn from the distribution of the variable $\widetilde{W}^\star = (\widetilde{w}_1^\star, \cdots, \widetilde{w}_p^\star)^{\mathrm{T}}$. Suppose Conditions 1, 2 and 3 hold for the correlation matrix $\widetilde{R} = (\widetilde{r}_{i,j}) = \mathrm{corr}(\log \widetilde{w}_i^\star, \log \widetilde{w}_j^\star)$ and covariance matrix $\widetilde{\Omega} = (\widetilde{\omega}_{i,j}) = \mathrm{cov}(\log \widetilde{w}_i^\star, \log \widetilde{w}_j^\star)$. In addition, assume the tail distribution of $\log \widetilde{W}^\star$ follows sub-Gaussian type (Condition 4) or the polynomial type (Condition 5), then the asymptotic distribution under the null hypothesis and the power of the test can be written out in the same ways as those for the two-sample case.

## 3.5. Simulation Studies

### 3.5.1. Simulation settings and performance evaluation

Simulation studies were conducted to evaluate the numerical performance of the proposed test $\Phi_\alpha$ based on the CLR-transformed data and to compare with the test $M_I$ in (Cai, Liu, and Xia, 2014) when applied to the compositional data $X_d$, the logarithm of the compositional data $\log X_d$ and the logarithm of the true basis count $\log W_d$ $(d = 1, 2)$. The results based on $\log W_d$ are considered as an oracle procedure to test the difference between $\mu_1$ and $\mu_2$.

To simulate the data, the basis counts and the compositional data were generated as the following. Two $n \times p$ data matrices $V_d = (V_{1,d}, \cdots, V_{n,d})^{\mathrm{T}} = (v_{i,j,d})$ $(d = 1, 2)$ were first generated from a multivariate normal distribution $N_p(\mu_d, \Omega)$ or a Gamma multivariate model, where $V_{i,d} = FU_d + \mu_d$. Here the matrix $F$ is generated by calculating the singular value decomposition $\Omega = QSQ^T$ and setting $F = QS^{1/2}$, and the components of $U_d$ are i.i.d standardized Gamma(10,1) random variables. The data $(W_d, X_d)$ were then generated through $w_{i,j,d} = \exp(v_{i,j,d})$ and $x_{i,j,d} = \exp(v_{i,j,d})/\sum_{k=1}^p \exp(v_{i,k,d})$ $(d = 1, 2)$. Thus, $X_d$ followed a logistic-normal distribution (Aitchison and Shen, 1980) or a type of logistic-gamma distribution.

The parameters $(\mu_d, \Omega)$ were set as follows. In both cases, we picked the components of $\mu_2$ randomly from the uniform distribution on $[0, 10]$. Under the null hypothesis, $\mu_1 = \mu_2$. Under the alternative hypothesis, $\mu_1 = \mu_2 + \delta$, the support set $S = \{l_1, \cdots, l_m : l_1 < l_2 < \cdots < l_m\}$ of $\delta$, with cardinality $m$, was randomly and uniformly selected from $\{1, \cdots, p\}$. For the elements in the support $S$, four different magnitudes were considered:

M 1: $\mu_{1_j,1} = \pm(\alpha_1 \log p/n)^{1/2}$, with equal probability and $m = \lfloor \beta_1 p \rfloor$.

M 2: $\mu_{1_j,1} = \pm(\alpha_2 \log p/n)^{1/2}$, with equal probability and $m = \lfloor p^{1/2} \rfloor$.

M 3: $\mu_{1_j,1}$ is uniformly drawn from $[-(\alpha_3 \log p/n)^{1/2}, (\alpha_3 \log p/n)^{1/2}]$, $m = \lfloor \beta_2 p \rfloor$.

M 4: $\mu_{1_j,1}$ is uniformly drawn from $[-(\alpha_4 \log p/n)^{1/2}, (\alpha_4 \log p/n)^{1/2}]$, $m = \lfloor p^{1/2} \rfloor$.

Denote by $D = (d_{i,j})$ the diagonal matrix with diagonal elements $d_{i,i} = \mathrm{unif}(1,3)$, and let $\lambda_{\min}(\cdot)$ be the smallest eigenvalue. Four different covariance structures of the log basis were considered as follows.

Model 1: (Bandable $\Omega$): $\Omega = (\omega_{i,j})$ where $\omega_{i,j} = 0.6^{|i-j|}$ for $1 \leq i, j \leq p$.

Model 2: (Sparse $\Omega$): $\Omega = \mathrm{diag}(A_1, A_2)$, where $A_1 = B + \varepsilon I_{p_1}$, $A_2 = 4I_{p_2}$, $p_1 = \lfloor \sqrt{p} \rfloor$, $p_2 = p - p_1$, and $B$ is a symmetric matrix where lower triangular entries are independent from the uniform distribution on $[-1, -0.5] \cup [0.5, 1]$ with probability 0.2 and equal to 0 with probability 0.8, and $\varepsilon = \max(-\lambda_{\min}(B), 0) + 0.05$.

Model 3: (Sparse $\Omega$): $\Sigma = (\sigma_{i,j})$ where $\sigma_{i,j} = 0.6^{|i-j|}$ for $1 \leq i, j \leq p$. $\Omega = D^{1/2} \Sigma^{-1} D^{1/2}$.

Model 4: (Non-sparse $\Omega$): $\Omega^\star = (\omega_{i,j}^\star)$ where $\omega_{i,i}^\star = 1$, $\omega_{i,j}^\star = 0.8$ for $2(k-1)+1 \leq i \neq j \leq 2k$, where $k = 1, \cdots, \lfloor p/2 \rfloor$, and $\omega_{i,j}^\star = 0$ otherwise, $\Omega = D^{1/2} \Omega^\star D^{1/2} + E + \delta I_p$ with $\delta = \left| \lambda_{\min}(D^{1/2} \Omega^\star D^{1/2} + E) \right| + 0.05$, where $E$ is a symmetric matrix with the lower triangular components generated independently from the uniform distribution on $[-0.2, 0.2]$ with probability 0.3 and equal to 0 with probability 0.7.

*3.5.2. Simulation results*

The sample size and the dimension were set as $(n, p) = (100, 50), (100, 100)$ and $(100, 200)$ and the simulations were repeated 1000 times under each setting. The empirical size and power of the

proposed test $\Phi_\alpha$ based on CLR$X_d$, $\log W_d$, $\log X_d$ and $X_d$ $(d = 1, 2)$ under the four different models and various null and alternative mean vectors are summarized in Table $3.1$ and $3.2$, altogether with the parameters $\alpha_i$ and $\beta_j$ representing the magnitude of the signals and the size of the support.

The results showed that the proposed test $\Phi_\alpha$ based on CLR-transformation had similar empirical size and power to the oracle test based on the true log-count data $\log W_d$. The empirical test size of the proposed test $\Phi_\alpha$ was close to, and controlled in most settings by the nominal level of $0.05$. However, the test of Cai, Liu, and Xia, 2014, when applied directly to the observed compositional data $X_d$, was conservative. In addition, the performance of the test of Cai, Liu, and Xia, 2014 when applied to $\log X_d$ was not stable. As seen from Table $3.1$, the empirical size and power of the test of Cai, Liu, and Xia, 2014 when applied to $\log X_d$ was close to our proposed test under the log-normal distribution setting. However, under the multivariate log-Gamma model as shown in Table $3.2$, the power of the proposed test based on CLR$X_d$ uniformly outperformed the test based on $\log X_d$. As expected, for both log-normal and log-Gamma distributions, the power of the proposed test depended on the magnitude of the signal, the size of the support and the dimension $p$. In majority of the cases, the empirical power increased when the magnitude of the signal, or the size of the support increased.

Table 3.1: Empirical size and power of the tests based on 1000 replications with $\alpha = 0.05$ and $n = 100$ for basis generated from log-normal distributions. Model 1: $\alpha_1 = \alpha_2 = 3$, $\alpha_3 = \alpha_4 = 10$, $\beta_1 = \beta_2 = 0.05$. Model 2: $\alpha_1 = \alpha_2 = 10$, $\alpha_3 = \alpha_4 = 20$, $\beta_1 = 0.15$, $\beta_2 = 0.2$. Model 3: $\alpha_1 = \alpha_2 = 10$, $\alpha_3 = \alpha_4 = 20$, $\beta_1 = 0.05$, $\beta_2 = 0.1$. Model 4: $\alpha_1 = \alpha_2 = 10$, $\alpha_3 = \alpha_4 = 20$, $\beta_1 = 0.15$, $\beta_2 = 0.2$.

| | Model 1 | | | Model 2 | | | Model 3 | | | Model 4 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $p$=50 | $p$=100 | $p$=200 | $p$=50 | $p$=100 | $p$=200 | $p$=50 | $p$=100 | $p$=200 | $p$=50 | $p$=100 | $p$=200 |
| **Null** | | | | | | | | | | | | |
| $\mathrm{CLR}X$ | 0.049 | 0.051 | 0.042 | 0.041 | 0.047 | 0.051 | 0.050 | 0.051 | 0.048 | 0.053 | 0.047 | 0.047 |
| $\log W$ | 0.048 | 0.047 | 0.047 | 0.051 | 0.043 | 0.052 | 0.047 | 0.052 | 0.049 | 0.050 | 0.046 | 0.045 |
| $\log X$ | 0.042 | 0.055 | 0.049 | 0.045 | 0.043 | 0.050 | 0.037 | 0.060 | 0.054 | 0.060 | 0.040 | 0.039 |
| $X$ | 0.021 | 0.020 | 0.019 | 0.007 | 0.004 | 0.000 | 0.007 | 0.004 | 0.001 | 0.013 | 0.010 | 0.002 |
| **M1: $m = \beta_1 p$ with fixed magnitude** | | | | | | | | | | | | |
| $\mathrm{CLR}X$ | 0.365 | 0.698 | 0.903 | 0.791 | 0.974 | 0.994 | 0.203 | 0.527 | 0.822 | 0.964 | 0.994 | 1.000 |
| $\log W$ | 0.342 | 0.669 | 0.895 | 0.823 | 0.970 | 0.994 | 0.219 | 0.556 | 0.827 | 0.968 | 0.997 | 1.000 |
| $\log X$ | 0.285 | 0.647 | 0.876 | 0.657 | 0.929 | 0.992 | 0.143 | 0.415 | 0.766 | 0.932 | 0.993 | 1.000 |
| $X$ | 0.119 | 0.284 | 0.457 | 0.143 | 0.151 | 0.067 | 0.017 | 0.013 | 0.047 | 0.580 | 0.422 | 0.301 |
| **M2: $m = \sqrt{p}$ with fixed magnitude** | | | | | | | | | | | | |
| $\mathrm{CLR}X$ | 0.797 | 0.894 | 0.967 | 0.811 | 0.829 | 0.903 | 0.551 | 0.858 | 0.944 | 0.873 | 0.991 | 0.990 |
| $\log W$ | 0.753 | 0.902 | 0.957 | 0.840 | 0.840 | 0.905 | 0.589 | 0.885 | 0.951 | 0.871 | 0.988 | 0.989 |
| $\log X$ | 0.731 | 0.874 | 0.961 | 0.652 | 0.765 | 0.907 | 0.468 | 0.728 | 0.857 | 0.810 | 0.969 | 0.988 |
| $X$ | 0.408 | 0.453 | 0.584 | 0.139 | 0.034 | 0.067 | 0.105 | 0.062 | 0.098 | 0.324 | 0.284 | 0.205 |
| **M3: $m = \beta_2 p$ with varied magnitude** | | | | | | | | | | | | |
| $\mathrm{CLR}X$ | 0.265 | 0.900 | 0.966 | 0.538 | 0.792 | 0.993 | 0.753 | 0.568 | 0.896 | 0.888 | 0.995 | 1.000 |
| $\log W$ | 0.264 | 0.890 | 0.962 | 0.586 | 0.781 | 0.993 | 0.753 | 0.591 | 0.896 | 0.901 | 0.994 | 1.000 |
| $\log X$ | 0.219 | 0.866 | 0.965 | 0.425 | 0.737 | 0.985 | 0.655 | 0.503 | 0.837 | 0.818 | 0.985 | 1.000 |
| $X$ | 0.099 | 0.459 | 0.639 | 0.086 | 0.045 | 0.074 | 0.119 | 0.028 | 0.048 | 0.444 | 0.379 | 0.368 |
| **M4: $m = \sqrt{p}$ with varied magnitude** | | | | | | | | | | | | |
| $\mathrm{CLR}X$ | 0.992 | 0.957 | 1.000 | 0.520 | 0.946 | 0.791 | 0.624 | 0.795 | 0.921 | 0.467 | 0.668 | 0.999 |
| $\log W$ | 0.992 | 0.958 | 1.000 | 0.554 | 0.953 | 0.783 | 0.632 | 0.811 | 0.923 | 0.491 | 0.653 | 1.000 |
| $\log X$ | 0.985 | 0.94 | 1.000 | 0.390 | 0.883 | 0.747 | 0.537 | 0.680 | 0.862 | 0.450 | 0.586 | 0.994 |
| $X$ | 0.843 | 0.622 | 0.946 | 0.027 | 0.143 | 0.033 | 0.131 | 0.053 | 0.056 | 0.248 | 0.082 | 0.237 |

Table 3.2: Empirical size and power of tests based on 1000 replications with $\alpha = 0.05$ and $n = 100$ for basis generated from log-Gamma models. Model 1: $\alpha_1 = \alpha_2 = 3$, $\alpha_3 = \alpha_4 = 10$, $\beta_1 = \beta_2 = 0.05$. Model 2: $\alpha_1 = \alpha_2 = 10$, $\alpha_3 = \alpha_4 = 20$, $\beta_1 = 0.15$, $\beta_2 = 0.2$. Model 3: $\alpha_1 = \alpha_2 = 10$, $\alpha_3 = \alpha_4 = 20$, $\beta_1 = 0.05$, $\beta_2 = 0.1$. Model 4: $\alpha_1 = \alpha_2 = 10$, $\alpha_3 = \alpha_4 = 20$, $\beta_1 = 0.15$, $\beta_2 = 0.2$.

| | Model 1 | | | Model 2 | | | Model 3 | | | Model 4 | | |
| | $p$=50 | $p$=100 | $p$=200 | $p$=50 | $p$=100 | $p$=200 | $p$=50 | $p$=100 | $p$=200 | $p$=50 | $p$=100 | $p$=200 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Null** | | | | | | | | | | | | |
| $\text{CLR}X$ | 0.044 | 0.041 | 0.053 | 0.048 | 0.049 | 0.057 | 0.048 | 0.048 | 0.047 | 0.040 | 0.048 | 0.051 |
| $\log W$ | 0.040 | 0.046 | 0.052 | 0.049 | 0.049 | 0.058 | 0.050 | 0.048 | 0.044 | 0.041 | 0.050 | 0.049 |
| $\log X$ | 0.043 | 0.049 | 0.034 | 0.041 | 0.045 | 0.056 | 0.038 | 0.034 | 0.035 | 0.034 | 0.033 | 0.037 |
| $X$ | 0.025 | 0.025 | 0.010 | 0.010 | 0.003 | 0.000 | 0.010 | 0.000 | 0.003 | 0.011 | 0.004 | 0.002 |
| **M1: $m = \beta_1 p$ with fixed magnitude** | | | | | | | | | | | | |
| $\text{CLR}X$ | 0.335 | 0.677 | 0.927 | 0.778 | 0.966 | 0.985 | 0.217 | 0.496 | 0.840 | 0.964 | 0.999 | 1.000 |
| $\log W$ | 0.326 | 0.668 | 0.924 | 0.816 | 0.969 | 0.988 | 0.237 | 0.534 | 0.844 | 0.972 | 0.999 | 1.000 |
| $\log X$ | 0.215 | 0.470 | 0.453 | 0.489 | 0.807 | 0.943 | 0.094 | 0.249 | 0.461 | 0.913 | 0.782 | 0.987 |
| $X$ | 0.082 | 0.158 | 0.103 | 0.070 | 0.092 | 0.021 | 0.025 | 0.003 | 0.009 | 0.531 | 0.105 | 0.160 |
| **M2: $m = \sqrt{p}$ with fixed magnitude** | | | | | | | | | | | | |
| $\text{CLR}X$ | 0.767 | 0.863 | 0.969 | 0.817 | 0.849 | 0.885 | 0.448 | 0.857 | 0.846 | 0.886 | 0.995 | 0.990 |
| $\log W$ | 0.736 | 0.866 | 0.956 | 0.843 | 0.849 | 0.888 | 0.482 | 0.883 | 0.843 | 0.891 | 0.995 | 0.986 |
| $\log X$ | 0.713 | 0.761 | 0.548 | 0.501 | 0.600 | 0.801 | 0.188 | 0.514 | 0.384 | 0.615 | 0.691 | 0.369 |
| $X$ | 0.418 | 0.378 | 0.130 | 0.058 | 0.006 | 0.026 | 0.008 | 0.038 | 0.003 | 0.089 | 0.059 | 0.001 |
| **M3: $m = \beta_2 p$ with varied magnitude** | | | | | | | | | | | | |
| $\text{CLR}X$ | 0.863 | 0.998 | 0.993 | 0.564 | 0.808 | 0.986 | 0.111 | 0.428 | 0.998 | 0.934 | 1.000 | 1.000 |
| $\log W$ | 0.830 | 0.998 | 0.994 | 0.601 | 0.801 | 0.989 | 0.104 | 0.430 | 0.998 | 0.917 | 1.000 | 1.000 |
| $\log X$ | 0.560 | 0.982 | 0.732 | 0.037 | 0.606 | 0.936 | 0.065 | 0.118 | 0.798 | 0.727 | 0.972 | 0.895 |
| $X$ | 0.218 | 0.726 | 0.161 | 0.042 | 0.01 | 0.022 | 0.012 | 0.006 | 0.036 | 0.081 | 0.219 | 0.047 |
| **M4: $m = \sqrt{p}$ with varied magnitude** | | | | | | | | | | | | |
| $\text{CLR}X$ | 0.934 | 0.968 | 1.000 | 0.862 | 0.622 | 0.942 | 0.207 | 0.406 | 0.890 | 0.490 | 0.646 | 1.000 |
| $\log W$ | 0.923 | 0.976 | 1.000 | 0.877 | 0.641 | 0.943 | 0.220 | 0.411 | 0.888 | 0.513 | 0.633 | 0.999 |
| $\log X$ | 0.811 | 0.932 | 0.851 | 0.509 | 0.412 | 0.759 | 0.065 | 0.199 | 0.498 | 0.251 | 0.237 | 0.583 |
| $X$ | 0.470 | 0.572 | 0.233 | 0.057 | 0.006 | 0.049 | 0.004 | 0.012 | 0.004 | 0.031 | 0.010 | 0.004 |

## 3.6. Real Data Analysis

### 3.6.1. Application to a cross-sectional study of diet

Gut microbiome plays an important role human metabolism in order to maintain human health. Wu et al., 2011 reported a cross-sectional study to investigate the association between long-term dietary patterns and gut microbiome composition. The gut microbiome composition data were collected from 98 healthy individuals at the University of Pennsylvania, together with demographic data including body mass indexes (BMI). From each healthy subject, DNAs collected from stool samples were analysed by 454/Roche pyrosequencing of 16S rRNA gene segments from the V1-V2 region. An average of 9265 reads per sample were yielded, with a standard deviation of 3864, by denoising the pyrosequences. These reads were further grouped into 87 bacterial genera that were observed in at least one sample. Since the number of sequencing reads varied greatly across samples, the count data were converted into compositional data by dividing the total number of reads, where the maximum rounding error 0.5 was used to replace zero counts Aitchison, 2003.

One important question was to test whether obese and lean individuals had the same gut micro-biome composition, where obese (n=24) and lean group (n=25) were defined based on whether the BMI was in the upper or lower quartile. At a nominal level of $0.05$, the CLR based test indicated a significant difference in bacterial genus compositions $(p = 0.009)$. This was consistent with the previous finding in human and mice gut microbiome studies that obesity was associated with changes in the relative abundance of bacterial taxa (Bäckhed et al., 2004). Figure 3.1 presents the bar plots of the CLR transformation of genus composition, showing that the abundance of *Acidaminococcus* was clearly different between the obese and lean groups. *Acidaminococcus* is a genus in the phylum *Firmicutes* that was found to contribute to the change of energy balance and subsequent weight gain by holding great metabolic potential for efficient energy harvest from the diet (Turnbaugh et al., 2006). As a comparison, tests based the compositions $X_d$ and the logarithm of the compositions $\log X_d$ $(d = 1, 2)$ resulted $p = 0.542$ and $p = 0.100$, respectively. This may be due to the fact that the conditions for these tests to be valid did not hold for $X_d$ and $\log X_d$. The analysis was also performed for the compositions of 51 relatively commonly genera that were observed in at least 6 samples in the whole dataset. The $p$-value from our proposed test and the test statistics $M_I$ when applied to $X$ and $\log X$ was $0.007$, $0.064$ and $0.388$, respectively.
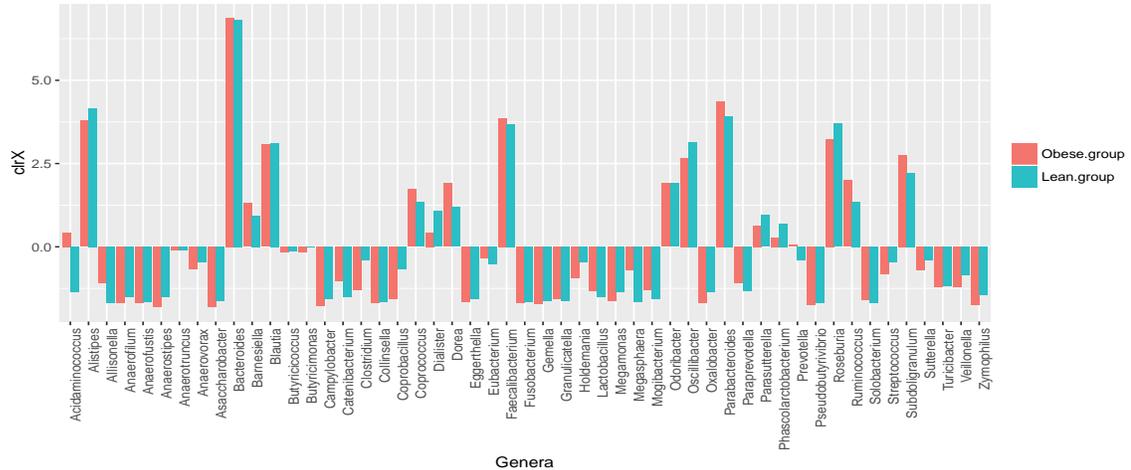
Figure 3.1: Bar plots of CLR$X$ of 51 bacterial genera composition in obese and lean samples.

### 3.6.2. Application to a microbiome study of Crohn's disease

Crohn's disease, characterized by abnormal composition of the intestinal dysbiosis, is one type of inflammatory bowel disease. The etiology of such a dysbiosis remained unknown. Lewis et al., 2015 recently reported a study to examine the gut microbiota composition among a cohort of 90 children with Crohn's disease at the University of Pennsylvania. Of these patients, 47 received an anti-tumor necrosis factor (anti-TNF) treatment. For each sample, fecal sample was collected at four time points: baseline, 1, 4, and 8 weeks after the treatment. Compositions of the bacterial genera were measured using shotgun metagenomic sequencing and the MetaPhlAn program (Segata et al., 2012). A total of 52 bacterial genera were identified that appeared in more than 5% of the samples. In addition, zeros were replaced with half of the non-zero minimum composition observed in the data (Aitchison, 2003).

To assess the effect of anti-TNF on fecal microbiome, testing whether there were significant changes in the overall microbiome compositions over the four time points during the anti-TNF treatment was performed using the proposed test for repeated measured data. The $p$-values of three pairs, including baseline *v.s.* week 1, week 1 *v.s.* week 4, week 4 *v.s.* week 8) were 0.0087, 0.493 and 0.449, respectively. The results indicated a significant change of gut microbiome composition within one week after the anti-TNF treatment, but were relatively stable during the rest of the treatment. It was interesting to note that these patients clinically responded to anti-TNF treatment within a week

after initiation of the treatment (Lewis et al., 2015), indicating that gut microbiome may play a role in reducing gut inflammation.

## 3.7. Discussion

This paper has proposed statistical tests for compositional equality of the log basis abundances based on the observed compositional data in high dimensional settings. The key assumptions of the test are certain dependency structure and tail distributions of the logarithm of the basis counts. Different from many existing two-sample mean tests, no sampling assumptions are made directly on the observed compositional data, rather these assumptions are made on the basis counts that are not directly observable. This overcomes the difficulty of modeling the compositional data in simplex.

Since the test statistic (3.6) takes the maximum of the normalized difference of the CLR transformed data, it is most powerful when the difference of the mean vectors is sparse. These assumptions are most likely met in analysis of microbiome data. The proposed test using the CLR transformation is powerful even when the true basis counts do not follow log-normal distributions. Our simulation results have showed that proposed test can be used to investigate the difference of the basis abundances and outperformed those naive tests based directly on the compositions or logarithm of the compositions.

# AAppendices

## A.1. Additional Lemmas and Technical Proofs for Chapter 1

### A.1.1. Proof of Theorems 1 and 3.

We prove a more general theorem first then move back to the proof of Theorem 1 and 3.

**Theorem 10.** *Under Conditions* $1$ *and* $2$*, suppose that* $N \geq c_0(n \vee p)\log(n+p)$ *for some constant* $c_0 > 0$*. Consider any solution* $\widehat{\mathbf{X}}$ *to the optimization problem* (1.3) *using regularization parameter selected by* (1.8)*. Then, with probability at least* $1 - 3(n+p)^{-1}$*, for each* $r \in \{1, 2, \cdots, n \wedge p\}$*, the average KL divergence satisfies*

$$\frac{1}{n} \mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) \leq \max \left\{ C_1 \sqrt{\frac{\log(n+p)}{N}}, \frac{C_2(n \vee p)r\log(n+p)}{N}, \right.$$
$$\left. \left( C_3 \sqrt{\frac{p\,(n \vee p)\log(n+p)}{nN}} \right) \sum_{i=r+1}^{n \wedge p} \sigma_i(\mathbf{X}^*) \right\}, \tag{A.1}$$

*where constants* $C_1, C_2$*, and* $C_3$ *only depend on* $c_0, \alpha_X, \beta_X, \alpha_R$ *and* $\beta_R$*.*

**Remark 1.** The rate of convergence provided by Theorem $10$ exhibits an interesting decomposition: besides the first term $O\left(\sqrt{\frac{\log(n+p)}{N}}\right)$, the rate $O\left(\frac{(n \vee p)r\log(n+p)}{N}\right)$ represents the estimation error corresponding to a rank-$r$ matrix, while the rest term $O\left(\left(\sqrt{\frac{p(n \vee p)\log(n+p)}{nN}}\right) \sum_{i=r+1}^{n \wedge p} \sigma_i(\mathbf{X}^*)\right)$ accounts for the approximation error due to using $r$ as a proxy for the rank of $\mathbf{X}^*$. When $\mathbf{X}^*$ is exactly a rank-$r$ matrix, this approximation error vanishes. When $\mathbf{X}^*$ is approximately low-rank, the value of $r$ can be optimally chosen to obtain the sharpest bound, which is presented in the following corollary.

*Proof.* ( Proof of Theorem 10 ): Note that, we can rewrite $\mathbf{W}$ by

$$\mathbf{W} = \sum_{k=1}^{N} \mathbf{E}_k, \quad \text{where} \quad \mathbf{E}_k \text{ are i.i.d copies of } \mathbf{E},$$

$$\mathbf{E} = \mathbf{e}_i \mathbf{e}_j \text{ with probability } \mathbf{\Pi}_{ij}, \quad 1 \leq i \leq n, 1 \leq j \leq p, \quad \mathbf{\Pi} = \mathbf{R}\mathbf{X}^* \in \mathbb{R}^{N \times p},$$

we rewrite the negative log-likelihood function (1.2) as

$$\mathcal{L}_N\left(\mathbf{X}\right) = -\frac{1}{N}\sum_{k=1}^{N}\log\langle\mathbf{X},\mathbf{E}_k\rangle = -\frac{1}{N}\sum_{i=1}^{n}\sum_{j=1}^{p}\mathbf{W}_{ij}\log\mathbf{X}_{ij}. \tag{A.2}$$

Then the estimator $\widehat{\mathbf{X}}$ associated with any optimal solution to the convex optimization (1.3) satisfies

$$\mathcal{L}_N(\widehat{\mathbf{X}}) + \lambda\|\widehat{\mathbf{X}}\|_* \le \mathcal{L}_N(\mathbf{X}^*) + \lambda\|\mathbf{X}^*\|_*$$

$$\Rightarrow \quad \mathcal{L}_N(\widehat{\mathbf{X}}) - \mathcal{L}_N(\mathbf{X}^*) = \frac{1}{N}\sum_{k=1}^{N}\langle\log\mathbf{X}^* - \log\widehat{\mathbf{X}},\mathbf{E}_k\rangle \le \lambda\left(\|\mathbf{X}^*\|_* - \|\widehat{\mathbf{X}}\|_*\right). \tag{A.3}$$

To derive a lower bound for $\frac{1}{N}\sum_{k=1}^{N}\langle\log\mathbf{X}^* - \log\widehat{\mathbf{X}},\mathbf{E}_k\rangle$, we first present the following lemma.

**Lemma 1.** *For all* $\mathbf{X}$ *in the constraint set* $\mathcal{C}(\alpha_X,\beta_X)$ *defined below,*

$$\mathcal{C}(\alpha_X,\beta_X) = \left\{\mathbf{X}\in\mathcal{S}(\alpha_x,\beta_x)\ \middle|\ \mathrm{D}(\mathbf{X}^*,\mathbf{X}) \ge n\log(\beta_X/\alpha_X)\sqrt{\frac{512\log(n+p)}{\log(4)\alpha_R^2 N}}\right\},$$

*with the probability proceeding* $1 - 2\left(n+p\right)^{-1}$, *we have*

$$\frac{1}{N}\sum_{k=1}^{N}\langle\log\mathbf{X}^* - \log\mathbf{X},\mathbf{E}_k\rangle \ge \frac{1}{2}\sum_{i,j}R_i X_{ij}^*\log\frac{X_{ij}^*}{X_{ij}} - E(n,p,r),$$

*where*

$$E(n,p,r) = \frac{1024\beta_X^2 npr}{\alpha_X^3\alpha_R}\left(\sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n}\vee\frac{\beta_X}{p}\right)}{N}} + \frac{C_0(n,p)}{N}\right)^2$$

$$+ \frac{16p}{\alpha_X}\left(\sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n}\vee\frac{\beta_X}{p}\right)}{N}} + \frac{C_0(n,p)}{N}\right)\sum_{i=r+1}^{n\wedge p}\sigma_i\left(\mathbf{X}^*\right).$$

*and* $C_0(n,p) = 4(3 + 2\log(n+p))$.

Given Lemma 1, we consider in two cases based on whether $\widehat{\mathbf{X}}\in\mathcal{C}(\alpha_X,\beta_X)$ or not.

- Case 1: If $\widehat{\mathbf{X}} \notin \mathcal{C}(\alpha_X, \beta_X)$, then

$$\frac{1}{n} \mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) < \log(\beta_X/\alpha_X) \sqrt{\frac{512 \log(n+p)}{\log(4)\alpha_R^2 N}}. \tag{A.4}$$

- Case 2: If $\widehat{\mathbf{X}} \in \mathcal{C}(\alpha_X, \beta_X)$, using the assumption $\min_i R_i \geq \alpha_R/n$ and applying Lemma 1, we obtain, with the probability proceeding $1 - 2(n+p)^{-1}$,

$$\frac{1}{N} \sum_{k=1}^{N} \langle \log \mathbf{X}^* - \log \widehat{\mathbf{X}}, \mathbf{E}_k \rangle \geq \frac{1}{2} \sum_{i,j} R_i X_{ij}^* \log \frac{X_{ij}^*}{X_{ij}} - E(n,p,r) \tag{A.5}$$

$$\geq \frac{\alpha_R}{2n} \mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) - E(n,p,r). \tag{A.6}$$

Another important element for establishing the error is to upper bound $\|\mathbf{X}^*\|_* - \|\widehat{\mathbf{X}}\|_*$ in (A.3). Denote $\boldsymbol{\Delta} = \widehat{\mathbf{X}} - \mathbf{X}^*$, $\nabla \mathcal{L}_N(\mathbf{X}) = -\frac{1}{N} \sum_{k=1}^{N} \langle \mathbf{X}, \mathbf{E}_k \rangle^{-1} \mathbf{E}_k$ as the gradient function of $\mathcal{L}_N(\mathbf{X})$. By Taylor's expansion of $\mathcal{L}_N$, there exists $\xi = (\xi_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}$ such that

$$\mathcal{L}_N(\hat{\mathbf{X}}) - \mathcal{L}_N(\mathbf{X}^*) - \langle \nabla \mathcal{L}_N(\mathbf{X}^*), \boldsymbol{\Delta} \rangle = \frac{1}{2N} \sum_{k=1}^{N} \frac{\langle \boldsymbol{\Delta}, E_k \rangle^2}{\langle \xi, E_k \rangle^2}, \quad \xi_{ij} \text{ is between } \hat{\mathbf{X}}_{ij}, \mathbf{X}_{ij}^*.$$

Adding $\langle \nabla \mathcal{L}_N(\mathbf{X}^*), \boldsymbol{\Delta} \rangle$ on both-hand sides of (A.3) and using Taylor's expansion, we obtain

$$\frac{1}{2N} \sum_{k=1}^{N} \frac{\langle \boldsymbol{\Delta}, E_k \rangle^2}{\langle \xi, E_k \rangle^2} = \mathcal{L}_N(\widehat{\mathbf{X}}) - \mathcal{L}_N(\mathbf{X}^*) - \langle \nabla \mathcal{L}_N(\mathbf{X}^*), \boldsymbol{\Delta} \rangle$$

$$\leq -\langle \nabla \mathcal{L}_N(\mathbf{X}^*), \boldsymbol{\Delta} \rangle + \lambda \left( \|\mathbf{X}^*\|_* - \|\widehat{\mathbf{X}}\|_* \right)$$

$$= -\langle \nabla \mathcal{L}_N(\mathbf{X}^*) + \mathbf{R} \mathbf{1}_n \mathbf{1}_p^T, \boldsymbol{\Delta} \rangle + \lambda \left( \|\mathbf{X}^*\|_* - \|\widehat{\mathbf{X}}\|_* \right)$$

$$\leq \|\nabla \mathcal{L}_N(\mathbf{X}^*) + \mathbf{R} \mathbf{1}_n \mathbf{1}_p^T\|_2 \|\boldsymbol{\Delta}\|_* + \lambda \left( \|\mathbf{X}^*\|_* - \|\widehat{\mathbf{X}}\|_* \right). \tag{A.7}$$

The third line comes from the identity $\langle \mathbf{R} \mathbf{1}_n \mathbf{1}_p^T, \boldsymbol{\Delta} \rangle = \langle \mathbf{R} \mathbf{1}_n, \boldsymbol{\Delta} \mathbf{1}_p \rangle = \langle \mathbf{R} \mathbf{1}_n, \mathbf{0}_n \rangle = 0$ and the forth inequality is the Hölder's inequality between the nuclear norm and operator norm. To further upper bound the nuclear norm $\|\boldsymbol{\Delta}\|_*$, we state three useful technical results.

**Lemma 2.** *With probability at least $1 - (n+p)^{-1}$, we have*

$$\|\nabla \mathcal{L}_N(\mathbf{X}^*) + \mathbf{R} \mathbf{1}_n \mathbf{1}_p^T\|_2 \leq \left\{ \sqrt{\frac{C_1(n,p)p \log(n+p)}{N}} \vee \frac{C_2(n,p)p \log(n+p)}{N} \right\},$$

*where* $C_1(n,p) = 8\left(\beta_R^2/n + (1 \vee \beta_R p/n)/\alpha_X\right)$ *and* $C_2(n,p) = 4(1/\alpha_X + \beta_R/(np)^{1/2})$.

Based on Lemma 2, with probability proceeding $1 - (n+p)^{-1}$, the selected tuning parameter $\lambda$ satisfies $\lambda \geq 2\|\nabla\mathcal{L}_N(\mathbf{X}^*) + \mathbf{R}\mathbf{1}_n\mathbf{1}_p^T\|_2$. As a result, we can use the following lemmas to upper bound the nuclear norm $\|\boldsymbol{\Delta}\|_*$.

**Lemma 3.** *If* (A.7) *holds and* $\lambda \geq 2\|\nabla\mathcal{L}_N(\mathbf{X}^*) + \mathbf{R}\mathbf{1}_n\mathbf{1}_p^T\|_2$, *we have the following upper bound for the nuclear norm of* $\boldsymbol{\Delta}$ *:*

$$\|\boldsymbol{\Delta}\|_* \leq 4\sqrt{2r}\|\boldsymbol{\Delta}\|_F + 4\sum_{i=r+1}^{n\wedge p}\sigma_i(\mathbf{X}^*). \tag{A.8}$$

In addition, Frobenius norm $\|\boldsymbol{\Delta}\|_F$ can be effectively bounded in terms of $\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}})$ as follows.

**Lemma 4.** *Under Condition* 2, *for any estimator* $\widehat{\mathbf{X}} \in \mathcal{S}(\alpha_X, \beta_X)$, *we have*

$$\frac{\alpha_X^2}{\beta_X p}\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) \leq \|\widehat{\mathbf{X}} - \mathbf{X}^*\|_F^2 \leq \frac{\beta_X^2}{\alpha_X p}\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}). \tag{A.9}$$

By applying Lemma 3 and 4, we obtain the upper bound of $\|\mathbf{X}^*\|_* - \|\widehat{\mathbf{X}}\|_*$ as

$$\|\mathbf{X}^*\|_* - \|\widehat{\mathbf{X}}\|_* \leq \|\mathbf{X}^* - \widehat{\mathbf{X}}\|_* \leq 4\sqrt{2r}\|\mathbf{X}^* - \widehat{\mathbf{X}}\|_F + 4\sum_{i=r+1}^{n\wedge p}\sigma_i(\mathbf{X}^*)$$

$$\leq 4\sqrt{\frac{2\beta_X^2 r}{\alpha_X p}\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}})} + 4\sum_{i=r+1}^{n\wedge p}\sigma_i(\mathbf{X}^*). \tag{A.10}$$

Therefore, combining (A.3), (A.5) and (A.10), if $\widehat{\mathbf{X}} \in \mathcal{C}(\alpha_X, \beta_X)$, we obtain,

$$\frac{\alpha_R}{2n}\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) \leq 4\lambda\left(\sqrt{\frac{2\beta_X^2 r}{\alpha_X p}\mathrm{D}(\mathbf{X}^*, \mathbf{X})} + \sum_{i=r+1}^{n\wedge p}\sigma_i(\mathbf{X}^*)\right) + E(n,p,r)$$

$$\leq \max\left\{8\lambda\sqrt{\frac{2\beta_X^2 r}{\alpha_X p}\mathrm{D}(\mathbf{X}^*, \mathbf{X})}, 8\lambda\sum_{i=r+1}^{n\wedge p}\sigma_i(\mathbf{X}^*) + 2E(n,p,r)\right\}$$

with probability at least $1 - 3(n+p)^{-1}$. The above equation yields

$$\frac{1}{n}\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) \leq \max\left\{\frac{512\beta_X^2\lambda^2 nr}{\alpha_R^2\alpha_X p}, \frac{4}{\alpha_r}\left(4\lambda\sum_{i=r+1}^{n\wedge p}\sigma_i(\mathbf{X}^*) + E(n,p,r)\right)\right\}. \tag{A.11}$$

Note that, in the formulation of $E(n,p,r)$, since $\log(n+p) \geq \log(2) > 0.6$, we have $\sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n} \vee \frac{\beta_X}{p}\right)}{N}} =$

$\sqrt{\frac{4(3+2\log(n+p))\left(\frac{\beta_R}{n}\vee\frac{\beta_X}{p}\right)}{N}} < \sqrt{\frac{28(\beta_R\vee\beta_X)\log(n+p)}{(n\wedge p)N}}$. In addition, under the assumption that $N \geq$

$c_0(n\vee p)\log(n+p)$, we have $\frac{C_0(n,p)}{N} < \frac{28}{\sqrt{c_0}}\sqrt{\frac{\log(n+p)}{(n\wedge p)N}}$. Consequently, we obtain $\sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n}\vee\frac{\beta_X}{p}\right)}{N}} +$

$\frac{C_0(n,p)}{N} \leq c_0'\sqrt{\frac{\log(n+p)}{(n\wedge p)N}}$ with $c_0' = \sqrt{28(\beta_R\vee\beta_X)} + \frac{28}{\sqrt{c_0}}$. Therefore, we can further upper bound

$E(n,p,r)$ by

$$E(n,p,r) \leq \frac{1024\beta_X^2 c_0'^{\,2}}{\alpha_X^3\alpha_R^2} \cdot \frac{(p\vee n)r\log(n+p)}{N} + \frac{16c_0'}{\alpha_X}\sqrt{\frac{p^2\log(n+p)}{(n\wedge p)N}}\sum_{i=r+1}^{n\wedge p}\sigma_i(\mathbf{X}^*). \quad \text{(A.12)}$$

We complete the proof by combining (A.4), (A.11) and (A.12),

$$\frac{1}{n}\mathrm{D}(\mathbf{X}^*,\widehat{\mathbf{X}}) \leq \max\left\{ c_1\sqrt{\frac{\log(n+p)}{N}}, \frac{c_2\lambda^2 nr}{p}, \right.$$
$$\left. \left(c_3\lambda + c_4\sqrt{\frac{p^2\log(n+p)}{(n\wedge p)N}}\right)\sum_{i=r+1}^{n\wedge p}\sigma_i(\mathbf{X}^*) + \frac{c_5(n\vee p)r\log(n+p)}{N} \right\}, \quad \text{(A.13)}$$

where constants $(c_1,c_2,c_3,c_4,c_5)$ are given by $c_1 = \log(\beta_X/\alpha_X)\sqrt{512/\log(4)\alpha_R^2}$, $c_2 = 512\beta_X^2/(\alpha_R^2\alpha_X)$,

$c_3 = 16/\alpha_R$, $c_4 = 16c_0'/\alpha_X$ and $c_5 = 1024\beta_X^2 c_0'^{\,2}/(\alpha_X^3\alpha_R^2)$. We also observe that, under the assumption $N \geq c_0(n\vee p)\log(n+p)$, the upper bound of selected tuning parameter $\lambda$ is given by

$$\lambda = 2\left(\sqrt{\frac{8\left(\beta_R^2/n + (1\vee\beta_R p/n)/\alpha_X\right)p\log(n+p)}{N}} \vee \frac{4(1/\alpha_X + \beta_R/(np)^{1/2})p\log(n+p)}{N}\right)$$
$$\leq 2\left(\sqrt{\frac{8\beta_R(\alpha_X\beta_R+1)}{\alpha_X}} \vee \frac{4(\alpha_X\beta_R+1)}{\sqrt{\alpha_X c_0}}\right)\sqrt{\frac{p(n\vee p)\log(n+p)}{nN}}. \quad \text{(A.14)}$$

Denote by $c_6 = 2\left(\sqrt{\frac{8\beta_R(\alpha_X\beta_R+1)}{\alpha_X}} \vee \frac{4(\alpha_X\beta_R+1)}{\sqrt{\alpha_X c_0}}\right)$. The proof for Theorem 10 is completed by plug-

ging (A.14) into (A.13):

$$\frac{1}{n}\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}})$$

$$\leq \max \left\{ c_1 \sqrt{\frac{\log(n+p)}{N}}, \frac{c_2 c_6^2 (n \vee p) r \log(n+p)}{N}, \right.$$

$$\left. (c_3 c_6 + c_4) \sqrt{\frac{p(n \vee p) \log(n+p)}{nN}} \sum_{i=r+1}^{n \wedge p} \sigma_i(\mathbf{X}^*) + \frac{c_5 (n \vee p) r \log(n+p)}{N} \right\}$$

$$\leq \max \left\{ c_1 \sqrt{\frac{\log(n+p)}{N}}, \frac{(c_2 c_6^2 + 2c_5)(n \vee p) r \log(n+p)}{N}, \right.$$

$$\left. 2(c_3 c_6 + c_4) \sqrt{\frac{p(n \vee p) \log(n+p)}{nN}} \sum_{i=r+1}^{n \wedge p} \sigma_i(\mathbf{X}^*) \right\}.$$

*Proof.* Proof of Theorem 1: By applying Theorem 10, when $\mathrm{rank}(\mathbf{X}^*) \leq r$, $\sum_{i=r+1}^{n \wedge p} \sigma_i(\mathbf{X}^*)$ vanishes in (A.1), and it yields (1.9). Besides, (1.10) can be obtained by applying Lemma $4$ to (1.9).

*Proof.* Proof of Theorem 3: We first focus on the proof of (1.12). If the composition $\mathbf{X}^* \in \mathbb{B}_q(\rho_q)$, we set $r = \max \left\{ i \mid \sigma_i(\mathbf{X}^*) > \tau \right\}$ for some fixed thresholding $\tau > 0$. By using this choice of $r$, we obtain

$$r\tau^q \leq \sum_{j=1}^{r} \sigma_j(\mathbf{X}^*)^q \leq \rho_q,$$

which implies $r \leq \tau^{-q} \rho_q$. In addition,

$$\sum_{i=r+1}^{n \wedge p} \sigma_i(\mathbf{X}^*) = \tau \sum_{i=r+1}^{n \wedge p} \frac{\sigma_i(\mathbf{X}^*)}{\tau} \leq \tau \left( \sum_{i=r+1}^{n \wedge p} \frac{\sigma_i(\mathbf{X}^*)}{\tau} \right)^q \leq \tau^{1-q} \rho_q.$$

We substitute the above relations into the upper bound (A.1) in Theorem $10$,

$$\frac{1}{n}\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) \leq \max \left\{ C_1 \sqrt{\frac{\log(n+p)}{N}}, \frac{C_2(n \vee p) \log(n+p) \tau^{-q} \rho_q}{N}, \right.$$

$$\left. C_3 \sqrt{\frac{p(n \vee p) \log(n+p)}{nN}} \tau^{1-q} \rho_q \right\}.$$

Since the rate of dominating terms are $O\left( \frac{(n \vee p) \log(n+p) \tau^{-q} \rho_q}{N} \right)$ and $O\left( \sqrt{\frac{p(n \vee p) \log(n+p)}{nN}} \tau^{1-q} \rho_q \right)$, we

can set $\tau = \frac{C_2}{C_3} \cdot \sqrt{\frac{(n \vee p)n \log(n+p)}{pN}}$ as to obtain the sharpest bound. As a result, we obtain (1.12),

$$\frac{1}{n} \mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) \leq C_2^{1-q} C_3^q \rho_q(p/n)^{\frac{q}{2}} \left( \frac{(n \vee p) \log(n+p)}{N} \right)^{1-\frac{q}{2}} .$$

Finally, (1.13) can be obtained by applying Lemma 4 to (1.12).

*A.1.2. Proof of Theorem 2.*

We discuss the proof for Theorem 2 under two different scenarios.

- If $n \geq p$, we randomly generate $M$ copies of i.i.d. Rademachar random matrices: $\mathbf{B}_1, \cdots, \mathbf{B}_M \in \mathbb{R}^{n \times (r-1)}$. Since $(\mathbf{B}_{k,ij} - \mathbf{B}_{l,ij})^2$ has the following probability distribution

$$\mathbb{P}\left((\mathbf{B}_{k,ij} - \mathbf{B}_{l,ij})^2 = x\right) = \begin{cases} \frac{1}{2}, & x = 4, \\ \frac{1}{2}, & x = 0, \end{cases}$$

based on Bernstein's inequality,

$$\mathbb{P}\left(\|\mathbf{B}_k - \mathbf{B}_l\|_F^2 \leq n(r-1)\right) = \mathbb{P}\left( \sum_{i=1}^{n} \sum_{j=1}^{r-1} (\mathbf{B}_{k,ij} - \mathbf{B}_{l,ij})^2 - 2n(r-1) \leq -n(r-1) \right)$$

$$\leq \exp\left( -\frac{(n(r-1))^2/2}{4n(r-1) + \frac{2}{3}n(r-1)} \right) \leq \exp\left( -\frac{1}{10} n(r-1) \right).$$

Therefore, whenever $M \leq \exp(n(r-1)/20)$, there is a positive probability that

$$\min_{1 \leq k < l \leq M} \left\{ \|\mathbf{B}_k - \mathbf{B}_l\|_F^2 \right\} \geq n(r-1), \tag{A.15}$$

which means we can find such fixed $\mathbf{B}_1, \ldots, \mathbf{B}_M \in \{-1, 1\}^{n \times (r-1)}$ such that (A.15) holds. For the rest of proof, we assume $\mathbf{B}_1, \ldots, \mathbf{B}_M$ are such fixed matrices while $M = \lfloor \exp(n(r-1)/20) \rfloor$. Note that $r - 1 \leq p/2$, we consider the following set of random rank-$r$ matrices,

$$\mathbf{X}_k = \begin{bmatrix} \frac{1}{p} & \cdots & \frac{1}{p} \\ \vdots & & \vdots \\ \frac{1}{p} & \cdots & \frac{1}{p} \end{bmatrix}_{n \times p} + \begin{bmatrix} r-1 & r-1 & p-2r+2 \\ \nu \mathbf{B}_k & -\nu \mathbf{B}_k & 0 \end{bmatrix},$$

where $0 < \nu < \frac{\beta_X - 1}{p} \wedge \frac{1 - \alpha_X}{p}$ is a to-be-determined constant. Then for $k \neq l$,

$$\|\mathbf{X}_k - \mathbf{X}_l\|_F^2 = 2\nu^2 \|\mathbf{B}_k - \mathbf{B}_l\|_F^2 \geq 2\nu^2 n(r-1).$$

$$\mathrm{D}(\mathbf{X}_k, \mathbf{X}_l) \overset{\text{Lemma 4}}{\geq} cp\|\mathbf{X}_k - \mathbf{X}_l\|_F^2 \geq 2c\nu^2 np(r-1),$$

where $c = \frac{\alpha_X}{2\beta_X^2}$. We also fix $\mathbf{R} = \left(\frac{1}{n}, \ldots, \frac{1}{n}\right)^\top$, i.e., the uniform distribution on each row. Suppose $P_k \sim \mathrm{Mult}(N; \frac{1}{n}\mathbf{X}_k)$, i.e., the multinomial distribution corresponding to composition $\mathbf{X}_k$ and $\mathbf{R}$. Based on Lemma 4, we have

$$\begin{aligned}
\mathrm{D}_{KL}(P_k, P_l) =& N \sum_{i=1}^{n} \sum_{j=1}^{p} \frac{\mathbf{X}_{k,ij}}{n} \log\left(\frac{\mathbf{X}_{k,ij}/n}{\mathbf{X}_{l,ij}/n}\right) \leq N \sum_{i=1}^{n} \sum_{j=1}^{p} Cnp\left(\frac{\mathbf{X}_{k,ij}}{n} - \frac{\mathbf{X}_{l,ij}}{n}\right)^2 \\
\leq& \frac{CNp}{n} \|\mathbf{B}_k - \mathbf{B}_l\|_F^2 \leq C\nu^2 Np(r-1),
\end{aligned}$$

where $C = \frac{\beta_X}{2\alpha_X^2}$. By generalized Fano's lemma (Yu, 1997),

$$\inf_{\widehat{\mathbf{X}}} \sup_{\mathbf{X} \subseteq \{\mathbf{X}_1, \cdots, \mathbf{X}_M\}} \mathbb{E} \left\|\widehat{\mathbf{X}} - \mathbf{X}\right\|_F^2 \geq \nu^2 n(r-1)\left(1 - \frac{C\nu^2 Np(r-1) + \log 2}{\log(M)}\right).$$

We further set $\nu^2 = c_\nu n/(Np)$ for some small constant $c_\nu > 0$ such that $(c_\nu Cn(r-1) + \log(2))/(n(r-1)/20) < 1/2$, then the lower bound above becomes

$$\inf_{\widehat{\mathbf{X}}} \sup_{\mathbf{X} \subseteq \{\mathbf{X}_1, \cdots, \mathbf{X}_M\}} \mathbb{E} \left\|\widehat{\mathbf{X}} - \mathbf{X}\right\|_F^2 \geq \frac{c_\nu n^2(r-1)}{2Np},$$

which implies

$$\inf_{\widehat{\mathbf{X}}} \sup_{\mathbf{X} \subseteq \{\mathbf{X}_1, \cdots, \mathbf{X}_M\}} \frac{p}{n} \mathbb{E} \left\|\widehat{\mathbf{X}} - \mathbf{X}\right\|_F^2 \geq \frac{c_\nu n(r-1)}{2N} = \frac{c_\nu(r-1)(n \vee p)}{2N}.$$

We can similarly derive that, for some constant $c'$,

$$\inf_{\widehat{\mathbf{X}}} \sup_{\mathbf{X} \subseteq \{\mathbf{X}_1, \cdots, \mathbf{X}_M\}} \frac{1}{n} \mathrm{D}\left(\mathbf{X}, \widehat{\mathbf{X}}\right) \geq c' \frac{(r-1)(n \vee p)}{N}.$$

Note that if $r \geq 2$, $r - 1 \geq r/2$, the lower bound result has been finally shown.

- If $n < p$, the proof is essentially the same as the case of $n \geq p$. Here we construct $M$ copies of i.i.d. Rademachar random matrices: $\mathbf{B}_1, \ldots, \mathbf{B}_M \in \mathbb{R}^{(r-1) \times \lfloor p/2 \rfloor}$, and the following set of random rank-$r$ matrices,

$$
\mathbf{X}_k = \begin{bmatrix} \frac{1}{p} & \cdots & \frac{1}{p} \\ \vdots & & \vdots \\ \frac{1}{p} & \cdots & \frac{1}{p} \end{bmatrix}_{n \times p} + \begin{array}{c} r-1 \\ n-r+1 \end{array} \begin{array}{ccc} \lfloor p/2 \rfloor & \lfloor p/2 \rfloor & p-2\lfloor p/2 \rfloor \\ \begin{bmatrix} \nu \mathbf{B}_k & -\nu \mathbf{B}_k & 0 \\ 0 & 0 & 0 \end{bmatrix} \end{array}.
$$

We omit the rest of the proof as it is essentially the same as the part for $n \geq p$.

*A.1.3. Proof of Corollary 2*

Using first order Taylor's expansion on the function $f(x) = x \log(x) = x_0 \log(x_0) + (\log(\xi) + 1)(x - x_0)$ for some $\xi$ between $x$ and $x_0$, we have

$$
\frac{1}{(\log p)^2 n} \sum_{i=1}^{n} (\mathbf{H}_{\mathsf{sh}}(\widehat{\mathbf{X}}_i) - \mathbf{H}_{\mathsf{sh}}(\mathbf{X}_i^*))^2
$$

$$
= \frac{1}{(\log p)^2 n} \sum_{i=1}^{n} \left( \sum_{j=1}^{p} X_{ij}^* \log X_{ij}^* - \widehat{X}_{ij} \log \widehat{X}_{ij} \right)^2
$$

$$
= \frac{1}{(\log p)^2 n} \sum_{i=1}^{n} \left( \sum_{j=1}^{p} \log \xi_{ij} (X_{ij}^* - \widehat{X}_{ij}) + \sum_{j=1}^{p} (X_{ij}^* - \widehat{X}_{ij}) \right)^2
$$

$$
= \frac{1}{(\log p)^2 n} \sum_{i=1}^{n} \left( \sum_{j=1}^{p} \log \xi_{ij} (X_{ij}^* - \widehat{X}_{ij}) \right)^2
$$

$$
\leq \frac{1}{(\log p)^2 n} \sum_{i=1}^{n} \left( \left( \sum_{j=1}^{p} (\log \xi_{ij})^2 \right) \left( \sum_{j=1}^{p} (X_{ij}^* - \widehat{X}_{ij})^2 \right) \right)
$$

$$
\leq \frac{(\log(p/\alpha_X))^2 p}{(\log p)^2 n} \|\widehat{\mathbf{X}} - \mathbf{X}^*\|_F^2,
$$

where $\xi_{ij}$ is between $X_{ij}^*$ and $\widehat{X}_{ij}$. In addition, using the Taylor's expansion on $f(x) = x^2 = x_0^2 + 2\xi(x - x_0)$, we have

$$\frac{p^2}{n} \sum_{i=1}^{n} (\mathbf{H}_{\mathsf{sp}}(\widehat{\mathbf{X}}_i) - \mathbf{H}_{\mathsf{sp}}(\mathbf{X}_i^*))^2 = \frac{p^2}{n} \sum_{i=1}^{n} \left( \sum_{j=1}^{p} \widehat{X}_{ij}^2 - X_{ij}^{*\,2} \right)^2$$

$$= \frac{p^2}{n} \sum_{i=1}^{n} \left( \sum_{j=1}^{p} 2\xi_{ij}(\widehat{X}_{ij} - X_{ij}^*) \right)^2$$

$$\leq \frac{4p^2}{n} \sum_{i=1}^{n} \left( \left( \sum_{j=1}^{p} \xi_{ij}^2 \right) \left( \sum_{j=1}^{p} (\widehat{X}_{ij} - X_{ij}^*)^2 \right) \right)$$

$$\leq \frac{4p\beta_X^2}{n} \|\widehat{\mathbf{X}} - \mathbf{X}^*\|_F^2.$$

We further bound the Frobenius norm loss of Bray-Curtis index by

$$\frac{1}{n^2} \sum_{1 \leq i < j \leq n} (\mathbf{H}_{\mathsf{bc}}(\widehat{\mathbf{X}}_i, \widehat{\mathbf{X}}_j) - \mathbf{H}_{\mathsf{bc}}(\mathbf{X}_i^*, \mathbf{X}_j^*))^2$$

$$= \frac{1}{4n^2} \sum_{i=1}^{n} \sum_{j=1}^{n} \left( \sum_{k=1}^{p} |\widehat{X}_{ik} - \widehat{X}_{jk}| - |X_{ik}^* - X_{jk}^*| \right)^2$$

$$\leq \frac{1}{4n^2} \sum_{i=1}^{n} \sum_{j=1}^{n} \left( \sum_{k=1}^{p} |\widehat{X}_{ik} - X_{ik}^*| + |\widehat{X}_{jk} - X_{jk}^*| \right)^2$$

$$\leq \frac{p}{4n^2} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{p} \left( |\widehat{X}_{ik} - X_{ik}^*| + |\widehat{X}_{jk} - X_{jk}^*| \right)^2$$

$$\leq \frac{p}{2n^2} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{p} \left( |\widehat{X}_{ik} - X_{ik}^*|^2 + |\widehat{X}_{jk} - X_{jk}^*|^2 \right)$$

$$= \frac{p}{n} \|\widehat{\mathbf{X}} - \mathbf{X}^*\|_F^2,$$

where we used the triangle inequalities that $|\widehat{x} - \widehat{y}| \leq |\widehat{x} - x| + |x - y| + |\widehat{y} - y|$ and $|x - y| \leq |\widehat{x} - x| + |\widehat{x} - \widehat{y}| + |\widehat{y} - y|$ in the second inequality, and Cauchy-Schwarz inequality in the third and forth inequalities. We complete the proof by applying Theorem $1$ and Theorem $3$ to the above inequalities.

*A.1.4. Proof of Proposition $2$*

It suffices to show that, $u_j + j^{-1}(1 - \epsilon p - \sum_{i=1}^{j} u_i) > 0$ for any $j \in [\rho]$ and $u_j + j^{-1}(1 - \epsilon p - \sum_{i=1}^{j} u_i) \leq 0$ for $j > \rho$.

71

- Case 1: If $j < \rho$, we note that $\mu = \rho^{-1}(1 - \epsilon p - \sum_{i=1}^{\rho} u_i) + \epsilon$, so we have

$$u_j + j^{-1}(1 - \epsilon p - \sum_{i=1}^{j} u_i) = j^{-1}\left(ju_j - \epsilon p - \left(\sum_{i=1}^{\rho} u_i - \sum_{i=j+1}^{\rho} u_i\right)\right)$$

$$= j^{-1}\left(j(u_j + \mu - \epsilon) + \sum_{i=j+1}^{\rho}(u_i + \mu - \epsilon)\right).$$

Using KKT condition that $u_i + \mu = x_i > \epsilon$, we obtain $u_j + j^{-1}(1 - \epsilon p - \sum_{i=1}^{j} u_i) > 0$ for $j \in [\rho]$.

- Case 2: If $j = \rho$, it is apparent that $u_\rho + \rho^{-1}(1 - \epsilon p - \sum_{i=1}^{\rho} u_i) = x_\rho - \epsilon > 0$.

- Case 3: Otherwise, $j > \rho$, then $u_j + \mu - \epsilon < 0$. According to similar argument, we obtain

$$u_j + j^{-1}(1 - \epsilon p - \sum_{i=1}^{j} u_i) = j^{-1}\left(\rho(u_j + \mu - \epsilon) + \sum_{i=\rho+1}^{j}(u_j - u_i)\right) < 0.$$

*A.1.5. Proof of technical lemmas*

**Proof of Lemma** 1

For notational simplicity, we let $\nu = n\log(\beta_X/\alpha_X)\sqrt{\frac{512\log(n+p)}{\log(4)\alpha_R^2 N}}$ and $D_{\mathbf{R}}(\mathbf{X}^*, \mathbf{X}) = \sum_{i,j} R_i X_{ij}^* \log\frac{X_{ij}^*}{X_{ij}}$. The main lines of this proof are in the same spirit as Lemma 3 in (Negahban and Wainwright, 2012), but sampling scheme and the constraint set are quite different. We use a standard peeling argument to prove the probability of the following "bad" event is small

$$\mathcal{B} = \left\{\exists \mathbf{X} \in \mathcal{C} \text{ such that } \left|\frac{1}{N}\sum_{i=1}^{N}\langle\log\mathbf{X}^* - \log\mathbf{X}, \mathbf{E}_i\rangle - D_{\mathbf{R}}(\mathbf{X}^*, \mathbf{X})\right| \geq \frac{1}{2}D_{\mathbf{R}}(\mathbf{X}^*, \mathbf{X}) - E(n, p, r)\right\}.$$

We separate the constraint set $\mathcal{C}$ into pieces and focus on a sequences of small sets $\mathcal{C}_l$,

$$\mathcal{C}_l = \left\{\mathbf{X} \in \mathcal{S} \,\middle|\, 2^{l-1}\nu \leq D(\mathbf{X}^*, \mathbf{X}) \leq 2^l\nu\right\}, \quad l \in \mathbb{N}^+.$$

As $\mathcal{C} \in \bigcup_{l=1}^{\infty} \mathcal{C}_l$, $\mathbf{X} \in \mathcal{C}$ implies that $\mathbf{X} \in \mathcal{C}_l$ with some $l$, and $D_R(\mathbf{X}^*, \mathbf{X}) \geq \frac{\alpha_R}{n}D(\mathbf{X}^*, \mathbf{X})$ under Condition 1, it suffices to estimate the probability of the following events and then apply the union bound.

$$\mathcal{B}_l = \left\{\exists \mathbf{X} \in \mathcal{C}_l \text{ such that } \left|\frac{1}{N}\sum_{i=1}^{N}\langle\log\mathbf{X}^* - \log\mathbf{X}, \mathbf{E}_i\rangle - D_{\mathbf{R}}\right| \geq \frac{2^l\nu\alpha_R}{4n} - E(n, p, r)\right\}, \quad l \in \mathbb{N}^+.$$

72

since $\mathcal{C}_l \subseteq \mathcal{C}\left(2^l \nu\right) := \left\{ \mathbf{X} \in \mathcal{S} \mid D_{\mathbf{R}}\left(\mathbf{X}^*, \mathbf{X}\right) \leq 2^l \nu \right\}$ that is defined in (A.19), we can establish the upper bound of the probability of event $\mathcal{B}$ by using the union bound and the fact that $e^x \geq x$, and applying Lemma 7,

$$
\begin{aligned}
\mathbb{P}\left(\mathcal{B}\right) &\leq \sum_{l=1}^{k} \mathbb{P}\left(\mathcal{B}_l\right) \\
&\leq \sum_{l=1}^{k} \exp\left(-\frac{4^l \alpha_R^2 N v^2}{512(n \log(\beta_X/\alpha_X))^2}\right) \\
&\leq \sum_{l=1}^{\infty} \exp\left(-\frac{\log(4)\alpha_R^2 N v^2 l}{512(n \log(\beta_X/\alpha_X))^2}\right) \\
&= \frac{\exp\left(-\frac{\log(4)\alpha_R^2 N v^2}{512(n \log(\beta_X/\alpha_X))^2}\right)}{1 - \exp\left(-\frac{\log(4)\alpha_R^2 N v^2}{512(n \log(\beta_X/\alpha_X))^2}\right)}.
\end{aligned}
$$

The proof is completed by plugging $\nu = n \log(\beta_X/\alpha_X)\sqrt{\frac{512 \log(n+p)}{\log(4)\alpha_R^2 N}}$.

**Proof of Lemma** 2

Let $\mathbf{Y}_i = -\langle \mathbf{X}, \mathbf{E}_i \rangle^{-1}\mathbf{E}_i + \mathbf{R}\mathbf{1}_n\mathbf{1}_p^T$, then $\nabla \mathcal{L}_N\left(\mathbf{X}^*\right) + \mathbf{R}\mathbf{1}_n\mathbf{1}_p^T = \frac{1}{N}\sum_{i=1}^N \mathbf{Y}_i$ and $\mathbb{E}\mathbf{Y}_i = -\sum_{jk}\frac{R_j X_{jk}^*}{X_{jk}^*}\mathbf{e}_j\mathbf{e}_k^T + \mathbf{R}\mathbf{1}_n\mathbf{1}_p^T = \mathbf{0}$. We note that, under Conditions 1 and 2, using Weyl's inequality, we have

$$
\begin{aligned}
\|\mathbf{Y}_i\|_2 &= \left\|-\langle \mathbf{E}_i, \mathbf{X}^* \rangle^{-1}\mathbf{E}_i + \mathbf{R}\mathbf{1}_n\mathbf{1}_p^T\right\|_2 \\
&\leq \max_{j,k}\left\|X_{jk}^{*-1}\mathbf{e}_j\mathbf{e}_k^T\right\|_2 + \left\|\mathbf{R}\mathbf{1}_n\mathbf{1}_p^T\right\|_2 \\
&\leq \max_{j,k}X_{jk}^{*-1} + \left(p\sum_{i=1}^n R_i^2\right)^{1/2} \\
&\leq p/\alpha_X + (\beta_R^2 p/n)^{1/2}.
\end{aligned}
$$

We also observe that

$$
\mathbb{E}\mathbf{Y}_i^T\mathbf{Y}_i = \sum_{ij}\frac{R_i X_{ij}^*}{X_{ij}^{*2}}\mathbf{e}_j\mathbf{e}_j^T - \|\mathbf{R}\|_F^2\mathbf{1}_p\mathbf{1}_p^T \quad \text{and} \quad \mathbb{E}\mathbf{Y}_i\mathbf{Y}_i^T = \sum_{ij}\frac{R_i X_{ij}^*}{X_{ij}^{*2}}\mathbf{e}_i\mathbf{e}_i^T - p\mathbf{R}\mathbf{1}_n\left(\mathbf{R}\mathbf{1}_n\right)^T.
$$

Hence, under Conditions $1$ and $2$, we apply Weyl's inequality to $\mathbb{E}\mathbf{Y}_i^T\mathbf{Y}_i$ and $\mathbb{E}\mathbf{Y}_i\mathbf{Y}_i^T$ and obtain

$$\|\mathbb{E}\mathbf{Y}_i^T\mathbf{Y}_i\|_2 \leq \left\|\sum_{ij}\frac{R_i}{X_{ij}^*}\mathbf{e}_j\mathbf{e}_j^T\right\|_2 + \|\mathbf{R}\|_F^2\,\|\mathbf{1}_p\mathbf{1}_p^T\|_2 = \max_{1\leq j\leq p}\sum_{i=1}^{n}\frac{R_i}{X_{ij}^*} + p\sum_{i=1}^{n}R_i^2 \leq p/\alpha_X + \beta_R^2 p/n,$$

$$\|\mathbb{E}\mathbf{Y}_i\mathbf{Y}_i^T\|_2 \leq \left\|\sum_{ij}\frac{R_i}{X_{ij}^*}\mathbf{e}_i\mathbf{e}_i^T\right\|_2 + \left\|p\mathbf{R}\mathbf{1}_n\left(\mathbf{R}\mathbf{1}_n\right)^T\right\|_2 = \max_{1\leq i\leq n}\sum_{j=1}^{p}\frac{R_i}{X_{ij}^*} + p\sum_{i=1}^{n}R_i^2 \leq \beta_R p^2/(\alpha_X n) + \beta_R^2 p/n.$$

Denote by $M = p/\alpha_X + (\beta_R^2 p/n)^{1/2}$ and $\sigma^2 = Np\left(\beta_R^2/n + (1\vee\beta_R p/n)/\alpha_X\right)$, then applying Lemma $5$, with any $t > 0$, we have

$$\mathbb{P}\left(\|\nabla\mathcal{L}_N\left(\mathbf{X}^*\right) + \mathbf{R}\mathbf{1}_n\mathbf{1}_p^T\|_2 \geq t\right)$$
$$\leq (n+p)\left\{\exp\left(-\frac{Nt^2}{4p\left(\beta_R^2/n + (1\vee\beta_R p/n)/\alpha_X\right)}\right) \vee \exp\left(-\frac{Nt}{2(p/\alpha_X + (\beta_R^2 p/n)^{1/2})}\right)\right\}.$$

We complete the proof by setting $t = \left\{\sqrt{\frac{8\left(\beta_R^2/n + (1\vee\beta_R p/n)/\alpha_X\right)p\log(n+p)}{N}} \vee \frac{4(1/\alpha_X + \beta_R/(np)^{1/2})p\log(n+p)}{N}\right\}.$

**Proof of Lemma $3$**

We observe that (A.7) is essentially equivalent to (B.2) in Lemma 1 from Negahban and Wainwright, 2011. Therefore, following their results, under the assumption $\lambda \geq 2\|\nabla\mathcal{L}_N(\mathbf{X}^*) + \mathbf{R}\mathbf{1}_n\mathbf{1}_p^T\|_2$, for each constant $r \leq n \wedge p$, there exists an orthogonal decomposition $\boldsymbol{\Delta} = \boldsymbol{\Delta}' + \boldsymbol{\Delta}''$, where the rank of $\boldsymbol{\Delta}'$ is less than $2r$ and $\boldsymbol{\Delta}''$ satisfies

$$\|\boldsymbol{\Delta}''\|_\star \leq 3\|\boldsymbol{\Delta}'\|_\star + 4\sum_{i=r+1}^{n\wedge p}\sigma_i\left(\mathbf{X}^\star\right) \text{ and } \|\boldsymbol{\Delta}\|_F^2 = \|\boldsymbol{\Delta}'\|_F^2 + \|\boldsymbol{\Delta}''\|_F^2.$$

Using the triangle inequality and $\|\boldsymbol{\Delta}'\|_* \leq \sqrt{2r}\|\boldsymbol{\Delta}'\|_F \leq \sqrt{2r}\|\boldsymbol{\Delta}\|_F$, we obtain

$$\|\boldsymbol{\Delta}\|_* \leq \|\boldsymbol{\Delta}'\|_* + \|\boldsymbol{\Delta}''\|_* \leq 4\|\boldsymbol{\Delta}'\|_* + 4\sum_{i=r+1}^{n\wedge p}\sigma_i\left(\mathbf{X}^*\right) \leq 4\sqrt{2r}\|\boldsymbol{\Delta}\|_F + 4\sum_{i=r+1}^{n\wedge p}\sigma_i\left(\mathbf{X}^*\right),$$

which completes the proof.

**Proof of Lemma** 4

Using Taylor expansion on the function $f(x) = \log(x)$, we rewrite KL divergence $\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}})$ as

$$\mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) = \sum_{i,j} -X_{ij}^* \log \frac{\widehat{X}_{ij}}{X_{ij}^*} = \sum_{i,j} -(\widehat{X}_{ij} - X_{ij}^*) + \frac{X_{ij}^*}{\xi_{ij}^2}(\widehat{X}_{ij} - X_{ij}^*)^2 = \sum_{i,j} \frac{X_{ij}^*}{\xi_{ij}^2}(\widehat{X}_{ij} - X_{ij}^*)^2,$$

where we use $\sum_{ij} X_{ij}^* = \sum_{ij} \widehat{X}_{ij} = n$ in the third equality, and $\xi_{ij}$ is a quantity between $\widehat{X}_{ij}$ and $X_{ij}^*$. Since both $X_{ij}^*$ and $\widehat{X}_{ij}$ are uniformly bounded by $[\alpha_X/p, \beta_X/p]$ for any $(i,j)$, we complete the proof by

$$\mathrm{D}(\mathbf{X}^*, \mathbf{X}) = \sum_{i,j} \frac{X_{ij}^*}{\xi_{ij}^2}(\widehat{X}_{ij} - X_{ij}^*)^2 \leq \frac{\beta_X/p}{(\alpha_X/p)^2} \sum_{i,j}(\widehat{X}_{ij} - X_{ij}^*)^2 = \frac{\beta_X p}{\alpha_X^2}\|\mathbf{X}^* - \mathbf{X}\|_F^2,$$

$$\mathrm{D}(\mathbf{X}^*, \mathbf{X}) = \sum_{i,j} \frac{X_{ij}^*}{\xi_{ij}^2}(\widehat{X}_{ij} - X_{ij}^*)^2 \geq \frac{\alpha_X/p}{(\beta_X/p)^2} \sum_{i,j}(\widehat{X}_{ij} - X_{ij}^*)^2 = \frac{\alpha_X p}{\beta_X^2}\|\mathbf{X}^* - \mathbf{X}\|_F^2.$$

**Concentration Inequalities**

**Lemma 5.** *Let $\mathbf{Y}_i$ be independent $n \times p$ zero-mean random matrices such that $\|\mathbf{Y}_i\|_2 \leq M$ and define $\sigma^2 = \max\left\{\sum_{i=1}^N \|\mathbb{E}\mathbf{Y}_i^T\mathbf{Y}_i\|_2, \sum_{i=1}^N \|\mathbb{E}\mathbf{Y}_i\mathbf{Y}_i^T\|_2\right\}$. For all $t > 0$, we have*

$$\mathbb{P}\left[\|\frac{1}{N}\sum_{i=1}^N \mathbf{Y}_i\|_2 \geq t\right] \leq (n+p)\left\{\exp\left(-N^2t^2/(4\sigma^2)\right) \vee \exp\left(-Nt/(2M)\right)\right\}. \tag{A.16}$$

*In addition, the expected spectral norm satisfies*

$$\left(\mathbb{E}\|\frac{1}{N}\sum_{i=1}^N \mathbf{Y}_i\|_2^2\right)^{1/2} \leq \sqrt{\frac{C_0(n,p)\sigma^2}{N}} + \frac{C_0(n,p)M}{N}, \tag{A.17}$$

*where the dimension constant $C_0(n,p) = 4(3 + 2\log(n+p))$.*

*Proof.* The proof of concentration inequality (A.16) follows, for example, Theorem 1.6 of (Tropp, 2011); see also Theorem 3.2 of (Recht, 2011). The proof of inequality (A.17) follows, for example, Theorem 1 of (Tropp, 2015).

**Lemma 6.** *Let $n \times p$ random matrices $\{\mathbf{E}_i\}_{i=1}^N$ be i.i.d with distribution $\mathbf{\Pi}$ on $\{\mathbf{e}_i(n)\mathbf{e}_j^T(p), (i,j) \in$*

$[n] \times [p]\}$ *and* $\{\varepsilon_i\}_{i=1}^N$ *is an i.i.d Rademacher sequence. Under Conditions* $1$ *and* $2$*, we have the upper bound*

$$\mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \varepsilon_i \mathbf{E}_i \right\|_2 \leq \sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n} \vee \frac{\beta_X}{p}\right)}{N}} + \frac{C_0(n,p)}{N}. \tag{A.18}$$

*Proof.* We establish this bound by applying Lemma $5$. Let $\mathbf{Y}_i = \varepsilon_i \mathbf{E}_i$, we first calculate the terms $M$ and $\sigma^2$ involved in Lemma $5$. According to the definition of $\varepsilon_i$ and $\mathbf{E}_i$,

$$M = \max_{1 \leq i \leq N} \|\mathbf{Y}_i\|_2 = \max_{1 \leq i \leq N} \|\varepsilon_i \mathbf{E}_i\|_2 = \max_{1 \leq i \leq N} \|\varepsilon_i \mathbf{e}_{k(i)} \mathbf{e}_{j(i)}^T\|_2 = 1.$$

Also note that,

$$\mathbb{E}\left[\mathbf{Y}_i^T \mathbf{Y}_i\right] = \mathbb{E}\left[\varepsilon_i^2 \mathbf{E}_i^T \mathbf{E}_i\right] = \sum_{k,j} R_k X_{kj}^* \mathbf{e}_j \mathbf{e}_j^T \text{ and } \mathbb{E}\left[\mathbf{Y}_i^T \mathbf{Y}_i\right] = \mathbb{E}\left[\varepsilon_i^2 \mathbf{E}_i \mathbf{E}_i^T\right] = \sum_{k,j} R_k X_{kj}^* \mathbf{e}_k \mathbf{e}_k^T.$$

We observe that, under Conditions $1$ and $2$,

$$\|\sum_{k,j} R_k X_{kj}^* \mathbf{e}_j \mathbf{e}_j^T\|_2 = \max_j \sum_{k=1}^n R_k X_{kj}^* \leq \sum_{k=1}^n R_k \cdot \max_{k,j} X_{kj}^* \leq \beta_X / p,$$

$$\|\sum_{k,j} R_k X_{kj}^* \mathbf{e}_k \mathbf{e}_k^T\|_2 = \max_k \sum_{j=1}^p R_k X_{kj}^* \leq \sum_{j=1}^p X_{kj}^* \cdot \max_k R_k \leq \beta_R / n.$$

As a result, $\sigma^2 = N\left(\frac{\beta_X}{p} \vee \frac{\beta_R}{n}\right)$. By applying Jensen's inequality and inequality (A.17), we obtain

$$\mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \epsilon_i \mathbf{E}_i \right\|_2 \leq \left( \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \epsilon_i \mathbf{E}_i \right\|_2^2 \right)^{1/2} \leq \sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n} \vee \frac{\beta_X}{p}\right)}{N}} + \frac{C_0(n,p)}{N}.$$

**Lemma 7.** *We define a constraint set* $\mathcal{C}(T)$ *with some constant* $T > 0$*,*

$$\mathcal{C}(T) = \left\{ \mathbf{X} \in \mathcal{S}(\alpha_x, \beta_x) \,\middle|\, \mathrm{D}(\mathbf{X}^*, \widehat{\mathbf{X}}) \leq T \right\}. \tag{A.19}$$

*And denote by $Z_T$ the function on the constraint set $\mathcal{C}(T)$*

$$Z_T = \sup_{\mathbf{X} \in \mathcal{C}(T)} \left| \frac{1}{N} \sum_{i=1}^{N} \langle \log \mathbf{X}^* - \log \mathbf{X}, \mathbf{E}_i \rangle - \sum_{ij} R_i X_{ij}^* \log \frac{X_{ij}^*}{X_{ij}} \right|,$$

*where $\{\mathbf{E}_i\}_{i=1}^{N}$ is i.i.d with distribution $\mathbf{\Pi} = \mathbf{R}\mathbf{X}^*$ on $\{\mathbf{e}_i(n)\mathbf{e}_j^T(p), (i,j) \in [n] \times [p]\}$. We suppose $\{R_i\}_{i=1}^{n}$ and $\mathbf{X}$ satisfy Condition $1$ and $2$ respectively. If $\mathbf{X}$ further satisfies $\|\mathbf{X}^* - \mathbf{X}\|_* \leq 4\sqrt{2r}\|\mathbf{X}^* - \mathbf{X}\|_F + 4\sum_{i=r+1}^{n \wedge p} \sigma_i(\mathbf{X}^*)$, then*

$$\mathbb{P}\left( Z_T \geq \frac{\alpha_R T}{4n} + E(n,p,r) \right) \leq \exp\left( -\frac{\alpha_R^2 N T^2}{512(n\log(\beta_X/\alpha_X))^2} \right),$$

*where*

$$E(n,p,r) = \frac{1024\beta_X^2 npr}{\alpha_X^3 \alpha_R} \left( \sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n} \vee \frac{\beta_X}{p}\right)}{N}} + \frac{C_0(n,p)}{N} \right)^2$$

$$+ \frac{16p}{\alpha_X} \left( \sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n} \vee \frac{\beta_X}{p}\right)}{N}} + \frac{C_0(n,p)}{N} \right) \sum_{i=r+1}^{n \wedge p} \sigma_i(\mathbf{X}^*).$$

*Proof.* Under Condition $2$, $\sup_{\mathbf{X} \in \mathcal{S}(\alpha_X, \beta_X)} \|\log X^* - \log X\|_\infty \leq \log(\beta_X/\alpha_X)$, we obtain the following concentration inequality by a version of Hoeffding's inequality due to Theorem 14.2 of (Bühlmann and Van De Geer, 2011),

$$\mathbb{P}(Z_T - \mathbb{E}Z_T \geq \alpha_R T/(8n)) \leq \exp\left( -\frac{\alpha_R^2 N T^2}{512(n\log(\beta_X/\alpha_X))^2} \right). \tag{A.20}$$

It remains to upper bound the quantity $\mathbb{E}Z_T$. By using a standard symmetrization argument, we obtain

$$\mathbb{E}(Z_T) = \mathbb{E} \sup_{\mathbf{X} \in \mathcal{C}(T)} \left| \frac{1}{N} \sum_{i=1}^{N} \langle \log \mathbf{X}^* - \log \mathbf{X}, \mathbf{E}_i \rangle - \sum_{ij} R_i X_{ij}^* \log \frac{X_{ij}^*}{X_{ij}} \right|$$

$$\leq 2\mathbb{E}\left( \sup_{\mathbf{X} \in \mathcal{C}(T)} \left| \frac{1}{N} \sum_{i=1}^{N} \varepsilon_i \langle \log \mathbf{X}^* - \log \mathbf{X}, \mathbf{E}_i \rangle \right| \right)$$

$$= 2\mathbb{E}\left( \sup_{\mathbf{X} \in \mathcal{C}(T)} \left| \frac{1}{N} \sum_{i=1}^{N} \varepsilon_i \sum_{j,k} \mathbb{I}_{(\mathbf{E}_i = \mathbf{e}_j(n)\mathbf{e}_k(p)^T)} \log \frac{X_{jk}^*}{X_{jk}} \right| \right),$$

where $\{\varepsilon_i\}_{i=1}^N$ is an i.i.d Rademacher sequence. We notice that, for any number $i \in [N]$ and any $t$ that satisfies $X_{jk}^* + t \geq \alpha_X/p$ with $\forall(j,k) \in [n] \times [p]$, the function

$$\phi_i(t) = \frac{\alpha_X}{p} \sum_{j,k} \mathbb{I}_{(\mathbf{E}_i = \mathbf{e}_j(n)\mathbf{e}_k(p)^T)} \log \frac{X_{jk}^*}{X_{jk}^* + t}$$

is a contraction with $\phi_i(0) = 0$. Then the contraction principle from Theorem 4.12 in (Ledoux and Talagrand, 2013), together with Hölder's inequality between nuclear and operator norm, yields

$$\mathbb{E}(Z_T) \leq \frac{4p}{\alpha_X} \mathbb{E}\left( \sup_{\mathbf{X} \in \mathcal{C}(T)} \left| \frac{1}{N} \sum_{i=1}^N \langle \mathbf{X}^* - \mathbf{X}, \varepsilon_i \mathbf{E}_i \rangle \right| \right)$$

$$\leq \frac{4p}{\alpha_X} \sup_{\mathbf{X} \in \mathcal{C}(T)} \|\mathbf{X}^* - \mathbf{X}\|_* \mathbb{E}\left\| \frac{1}{N} \sum_{i=1}^N \varepsilon_i \mathbf{E}_i \right\|_2. \tag{A.21}$$

We bound $\mathbb{E}\left\| \frac{1}{N} \sum_{i=1}^N \varepsilon_i \mathbf{E}_i \right\|_2$ by applying Lemma 6,

$$\mathbb{E}\left\| \frac{1}{N} \sum_{i=1}^N \varepsilon_i \mathbf{E}_i \right\|_2 \leq \sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n} \vee \frac{\beta_X}{p}\right)}{N}} + \frac{C_0(n,p)}{N}. \tag{A.22}$$

Under the assumption that $\|\mathbf{X}^* - \mathbf{X}\|_* \leq 4\sqrt{2r}\|\mathbf{X}^* - \mathbf{X}\|_F + 4\sum_{i=r+1}^{n \wedge p} \sigma_i(\mathbf{X}^*)$, applying Lemma 4, we can bound $\|\mathbf{X}^* - \mathbf{X}\|_*$ by

$$\sup_{\mathbf{X} \in \mathcal{C}(T)} \|\mathbf{X}^* - \mathbf{X}\|_* \leq 4\sqrt{2r} \sup_{\mathbf{X} \in \mathcal{C}(T)} \|\mathbf{X}^* - \mathbf{X}\|_F + 4\sum_{i=r+1}^{n \wedge p} \sigma_i(\mathbf{X}^*)$$

$$\leq \sup_{\mathbf{X} \in \mathcal{C}(T)} 4\sqrt{\frac{2\beta_X^2 r}{\alpha_X p} \mathrm{D}(\mathbf{X}^*, \mathbf{X})} + 4\sum_{i=r+1}^{n \wedge p} \sigma_i(\mathbf{X}^*)$$

$$\leq 4\sqrt{\frac{2\beta_X^2 rT}{\alpha_X p}} + 4\sum_{i=r+1}^{n \wedge p} \sigma_i(\mathbf{X}^*). \tag{A.23}$$

Combining inequalities (A.21), (A.22) and (A.23) yields,

$$\mathbb{E}(Z_T) \leq \frac{16p}{\alpha_X} \left( \sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n} \vee \frac{\beta_X}{p}\right)}{N}} + \frac{C_0(n,p)}{N} \right) \left( \sqrt{\frac{2\beta_X^2 rT}{\alpha_X p}} + \sum_{i=r+1}^{n \wedge p} \sigma_i(\mathbf{X}^*) \right).$$

Finally, using

$$\frac{\alpha_R T}{8n} + \frac{16p}{\alpha_X} \left( \sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n} \vee \frac{\beta_X}{p}\right)}{N}} + \frac{C_0(n,p)}{N} \right) \left( \sqrt{\frac{2\beta_X^2 rT}{\alpha_X p}} + \sum_{i=r+1}^{n \wedge p} \sigma_i\left(\mathbf{X}^*\right) \right)$$

$$\leq \frac{\alpha_R T}{4n} + \frac{1024\beta_X^2 npr}{\alpha_X^3 \alpha_R} \left( \sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n} \vee \frac{\beta_X}{p}\right)}{N}} + \frac{C_0(n,p)}{N} \right)^2$$

$$+ \frac{16p}{\alpha_X} \left( \sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n} \vee \frac{\beta_X}{p}\right)}{N}} + \frac{C_0(n,p)}{N} \right) \sum_{i=r+1}^{n \wedge p} \sigma_i\left(\mathbf{X}^*\right)$$

and concentration inequality (A.20), we achieve at

$$\mathbb{P}\left( Z_T \geq \frac{\alpha_R T}{4n} + E(n,p,r) \right) \leq \exp\left( -\frac{\alpha_R^2 N T^2}{512(n\log(\beta_X/\alpha_X))^2} \right),$$

with

$$E(n,p,r) = \frac{1024\beta_X^2 npr}{\alpha_X^3 \alpha_R} \left( \sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n} \vee \frac{\beta_X}{p}\right)}{N}} + \frac{C_0(n,p)}{N} \right)^2$$

$$+ \frac{16p}{\alpha_X} \left( \sqrt{\frac{C_0(n,p)\left(\frac{\beta_R}{n} \vee \frac{\beta_X}{p}\right)}{N}} + \frac{C_0(n,p)}{N} \right) \sum_{i=r+1}^{n \wedge p} \sigma_i\left(\mathbf{X}^*\right).$$

## A.2. Additional Lemmas and Technical Proofs for Chapter 2

### A.2.1. Proof of Proposition 3

Using the fact that the centered log-ratio covariance matrix $\mathbf{\Gamma}_0$ is symmetric and has all zero row sums (Aitchison, 2003 Property 4.6), we have

$$\mathrm{tr}\{(\boldsymbol{\gamma}_0 \mathbf{1}^T + \mathbf{1}\boldsymbol{\gamma}_0^T)^T \mathbf{\Gamma}_0\} = \mathrm{tr}(\boldsymbol{\gamma}_0^T \mathbf{\Gamma}_0 \mathbf{1}) + \mathrm{tr}(\boldsymbol{\gamma}_0 \mathbf{1}^T \mathbf{\Gamma}_0) = 0,$$

that is, the components $\boldsymbol{\gamma}_0 \mathbf{1}^T + \mathbf{1}\boldsymbol{\gamma}_0^T$ and $\mathbf{\Gamma}_0$ are orthogonal to each other.

To show the desired inequality, by the identity (4.35) of Aitchison, (2003), we have

$$\omega_{ij}^0 - \gamma_{ij}^0 = \omega_{ij}^0 - (\omega_{ij}^0 - \omega_{i\cdot}^0 - \omega_{\cdot j}^0 + \omega_{\cdot\cdot}^0) = \omega_{i\cdot}^0 + \omega_{\cdot j}^0 - \omega_{\cdot\cdot}^0.$$

Therefore,

$$\|\mathbf{\Omega}_0 - \mathbf{\Gamma}_0\|_{\max} \le \max_{i,j}(|\omega_{i\cdot}^0| + |\omega_{\cdot j}^0| + |\omega_{\cdot\cdot}^0|) \le 3p^{-1}\|\mathbf{\Omega}_0\|_1.$$

### A.2.2. Proof of Proposition 4

We first claim that if $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_p)^T \ne \mathbf{0}$, then the matrix $\mathbf{A} \equiv \boldsymbol{\alpha}\mathbf{1}^T + \mathbf{1}\boldsymbol{\alpha}^T$ has at least $p-1$ nonzero upper-triangular entries. To prove this, without loss of generality, assume $\alpha_1 \ne 0$ and that the last $q$ entries of the first row of $\mathbf{A}$ are zero, where $0 \le q \le p-1$; that is, $\alpha_1 + \alpha_j \ne 0$ for $1 \le j \le p-q$, and $\alpha_1 + \alpha_{p-q+1} = \cdots = \alpha_1 + \alpha_p = 0$. The latter implies $\alpha_{p-q+1} = \cdots = \alpha_p = -\alpha_1 \ne 0$, which gives rise to $\binom{q}{2} = q(q-1)/2$ nonzero entries at positions $(i,j)$ with $p - q + 1 \le i < j \le p$. Putting these pieces together, we obtain that the number of nonzero upper-triangular entries in $\mathbf{A}$ is at least

$$f(q) \equiv p - q - 1 + \frac{q(q-1)}{2} \ge f(1) = f(2) = p - 2.$$

To show that the lower bound $p-2$ is not attainable, note that if there are only $p-2$ nonzero upper-triangular entries, then $q = 1$ or $2$, and we have $\alpha_2 + \alpha_p = \cdots = \alpha_{p-2} + \alpha_p = 0$, which implies $\alpha_2 = \cdots = \alpha_{p-2} = -\alpha_p = \alpha_1 \ne 0$. Since $p \ge 5$, this gives rise to at least one nonzero entry at positions $(i,j)$ with $2 \le i < j \le p - 2$, which is a contradiction.

Now suppose $s_e(p) < (p-1)/2$ and that $\mathbf{\Omega}_1$ and $\mathbf{\Omega}_2$ in $\mathcal{B}_0(s_e(p))$ lead to $\mathbf{T}_1 = \mathbf{T}_2$, that is,

$$(\boldsymbol{\omega}_1 - \boldsymbol{\omega}_2)\mathbf{1}^T + \mathbf{1}(\boldsymbol{\omega}_1 - \boldsymbol{\omega}_2)^T = 2(\mathbf{\Omega}_1 - \mathbf{\Omega}_2).$$

Note that the right-hand side has fewer than $p-1$ nonzero upper-triangular entries. Then it follows from the above claim that $\mathbf{\Omega}_1 = \mathbf{\Omega}_2$.

We prove the other direction by showing that, if $s_e(p) \ge (p-1)/2$, then there exist $\mathbf{\Omega}_1$ and $\mathbf{\Omega}_2$ in

$\mathcal{B}_0(s_e(p))$ with $\mathbf{\Omega}_1 \neq \mathbf{\Omega}_2$ that lead to $\mathbf{T}_1 = \mathbf{T}_2$. Indeed, let

$$
\mathbf{\Omega}_1 = \begin{pmatrix} 1+c & c\mathbf{1}_{p_1}^T & \mathbf{0}_{p_2}^T \\ c\mathbf{1}_{p_1} & \mathbf{I} & \mathbf{0} \\ \mathbf{0}_{p_2} & \mathbf{0} & \mathbf{I} \end{pmatrix}, \quad \mathbf{\Omega}_2 = \begin{pmatrix} 1-c & \mathbf{0}_{p_1}^T & -c\mathbf{1}_{p_2}^T \\ \mathbf{0}_{p_1} & \mathbf{I} & \mathbf{0} \\ -c\mathbf{1}_{p_2} & \mathbf{0} & \mathbf{I} \end{pmatrix},
$$

where $p_1 = \lfloor (p-1)/2 \rfloor$, $p_2 = p - 1 - p_1$, and $0 < |c| < 1$. Then it is easy to verify that

$$
\mathbf{T}_1 = \mathbf{T}_2 = \begin{pmatrix} 0 & (2-c)\mathbf{1}_{p_1}^T & (2+c)\mathbf{1}_{p_2}^T \\ (2-c)\mathbf{1}_{p_1} & 2(\mathbf{1}_{p_1}\mathbf{1}_{p_1}^T - \mathbf{I}) & 2\mathbf{1}_{p_1}\mathbf{1}_{p_2}^T \\ (2+c)\mathbf{1}_{p_2} & 2\mathbf{1}_{p_2}\mathbf{1}_{p_1}^T & 2(\mathbf{1}_{p_2}\mathbf{1}_{p_2}^T - \mathbf{I}) \end{pmatrix}.
$$

This completes the proof.

*A.2.3. Concentration Inequalities*

To prepare for the proofs of Theorems 1 and 2, we first establish some useful concentration in-equalities. For notational simplicity, the constants $C_1, C_2, \ldots$ below may vary from line to line.

**Lemma 8.** *Under Condition 3, there exist constants $C_1, C_2 > 0$ such that*

$$
P\left( \max_j \left| \frac{1}{n} \sum_{k=1}^n Y_{kj} \right| \geq t \right) \leq C_1 p e^{-C_2 n t^2} \tag{A.24}
$$

*and*

$$
P\left( \max_{i,j} \left| \frac{1}{n} \sum_{k=1}^n Y_{ki} Y_{kj} - EY_i Y_j \right| \geq t \right) \leq C_1 p^2 e^{-C_2 n t^2} \tag{A.25}
$$

*for sufficiently small $t > 0$. Moreover, if $\log p = o(n^{1/5})$, then there exists a constant $C_3 > 0$ such that*

$$
P\left( \max_{i,j,\ell,m} \left| \frac{1}{n} \sum_{k=1}^n Y_{ki} Y_{kj} Y_{k\ell} Y_{km} - EY_i Y_j Y_\ell Y_m \right| \geq \varepsilon \right) = O(p^{-C_3}) \tag{A.26}
$$

*for every constant $\varepsilon > 0$.*

*Proof.* Inequalities (A.24) and (A.25) follow, for example, from Exercise 2.27 of Boucheron, Lugosi, and Massart, (2013); see also Bickel and Levina, (2008).

To prove (A.26), let $Z_{kijlm} = Y_{ki} Y_{kj} Y_{k\ell} Y_{km}$ and $Z_{ijlm} = Y_i Y_j Y_\ell Y_m$. Note first that, by Condition 3

and the sub-Gaussian tail bound, for any $K > 0$ and $i, j, \ell, m$,

$$P(|Z_{ijlm}| > K) \le 4P(|Y_j| > K^{1/4}) \le 8e^{-\alpha\sqrt{K}/8}.$$

Hence,

$$
\begin{aligned}
E|Z_{ijlm}|I(|Z_{ijlm}| > K) &= \int_0^\infty P(|Z_{ijlm}|I(|Z_{ijlm}| > K) > z)\, dz \\
&= KP(|Z_{ijlm}| > K) + \int_K^\infty P(|Z_{ijlm}| > z)\, dz \\
&\le 8Ke^{-\alpha\sqrt{K}/8} + \int_K^\infty 8e^{-\alpha\sqrt{z}/8}\, dz \\
&= \frac{8}{\alpha^2}(\alpha^2 K + 16\alpha\sqrt{K} + 128)e^{-\alpha\sqrt{K}/8},
\end{aligned}
$$

which is less than $\varepsilon/4$ if we choose $K$ sufficiently large. Then we have

$$
\begin{aligned}
P&\left( \max_{i,j,\ell,m} \left| \frac{1}{n}\sum_{k=1}^n Z_{kijlm} - EZ_{ijlm} \right| \ge \varepsilon \right) \\
&\le P\left( \max_{i,j,\ell,m} \left| \frac{1}{n}\sum_{k=1}^n Z_{kijlm}I(|Z_{kijlm}| \le K) - EZ_{ijlm}I(|Z_{ijlm}| \le K) \right| \ge \frac{\varepsilon}{2} \right) \\
&\quad + P\left( \max_{i,j,\ell,m} \left| \frac{1}{n}\sum_{k=1}^n Z_{kijlm}I(|Z_{kijlm}| > K) \right| \ge \frac{\varepsilon}{4} \right) \\
&\equiv T_1 + T_2.
\end{aligned}
$$

By Hoeffding's inequality and the union bound,

$$T_1 \le 2p^4 \exp\left( -\frac{n\varepsilon^2}{8K^2} \right).$$

Also, by Condition 1 and the sub-Gaussian tail bound,

$$T_2 \le P\left( \max_{k,i,j,\ell,m} |Z_{kijlm}| > K \right) \le P\left( \max_{k,j} |Y_{kj}| > K^{1/4} \right) \le 2npe^{-\alpha\sqrt{K}/8}.$$

Combining both terms, choosing $K = C^2(\log p + \log n)^2$ with $C > 8/\alpha$, and noting $\log p = o(n^{1/5})$,

we arrive at

$$P \left( \max_{i,j,\ell,m} \left| \frac{1}{n} \sum_{k=1}^{n} Z_{kijlm} - EZ_{ijlm} \right| \geq \varepsilon \right)$$

$$\leq 2p^4 \exp \left( -\frac{n\varepsilon^2}{8C^4(\log p + \log n)^4} \right) + 2(np)^{1-C\alpha/8}$$

$$= O(p^{-C_3})$$

for some $C_3 > 0$. This proves (A.26) and completes the proof.

**Lemma 9.** *Under Conditions 3–6, there exist constants $C_1, C_2, C_3 > 0$ such that*

$$P \left( \max_{i,j} |\hat{\theta}_{ij} - \theta_{ij}| \geq \varepsilon \right) = O(p^{-C_3}) \tag{A.27}$$

*and*

$$P \left( \max_{i,j} |\hat{\gamma}_{ij} - \omega_{ij}^0| / \sqrt{\hat{\theta}_{ij}} \geq C_1 \sqrt{\frac{\log p}{n}} + C_2 \frac{s_0(p)}{p} \right) = O(p^{-C_3}) \tag{A.28}$$

*for every constant $\varepsilon > 0$.*

*Proof.* We first prove (A.27). Define

$$\tilde{\theta}_{ij} = \frac{1}{n} \sum_{k=1}^{n} (\gamma_{ki}\gamma_{kj} - \tilde{\gamma}_{ij})^2,$$

where $\tilde{\gamma}_{ij} = n^{-1} \sum_{k=1}^{n} \gamma_{ki}\gamma_{kj}$. We then write

$$\hat{\theta}_{ij} - \tilde{\theta}_{ij} = \frac{1}{n} \sum_{k=1}^{n} \{ (\gamma_{ki}\gamma_{kj} - \tilde{\gamma}_{ij}) - \gamma_{ki}\bar{\gamma}_j - \gamma_{kj}\bar{\gamma}_i + 2\bar{\gamma}_i\bar{\gamma}_j \}^2 - \frac{1}{n} \sum_{k=1}^{n} (\gamma_{ki}\gamma_{kj} - \tilde{\gamma}_{ij})^2$$

$$= \frac{2}{n} \sum_{k=1}^{n} (\gamma_{ki}\gamma_{kj} - \tilde{\gamma}_{ij})(-\gamma_{ki}\bar{\gamma}_j - \gamma_{kj}\bar{\gamma}_i + 2\bar{\gamma}_i\bar{\gamma}_j) + \frac{1}{n} \sum_{k=1}^{n} (-\gamma_{ki}\bar{\gamma}_j - \gamma_{kj}\bar{\gamma}_i + 2\bar{\gamma}_i\bar{\gamma}_j)^2. \tag{A.29}$$

Note that, by definition, $\gamma_{kj} = Y_{kj} - \bar{Y}_k$, where $\bar{Y}_k = p^{-1} \sum_{j=1}^{p} Y_{kj}$. Define $\gamma_j = Y_j - \bar{Y}$, where $\bar{Y} = p^{-1} \sum_{j=1}^{p} Y_j$. Since $Y_j$ are uniformly sub-Gaussian by Condition 3, $\gamma_j$ are also uniformly sub-Gaussian. Using a truncation argument similar to that for proving (A.26), we can show that

$$P \left( \max_{i,j} \left| \frac{1}{n} \sum_{k=1}^{n} \gamma_{ki}^2 \gamma_{kj} - E\gamma_i^2 \gamma_j \right| \geq C_1 \right) = O(p^{-C_3})$$

for some $C_1, C_3 > 0$. The sub-Gaussian tails imply also that $E\gamma_i^2|\gamma_j| \leq \frac{1}{2}(E\gamma_i^4 + E\gamma_j^2) = O(1)$. Combining these two pieces yields

$$P\left(\max_{i,j}\left|\frac{1}{n}\sum_{k=1}^{n}\gamma_{ki}^2\gamma_{kj}\right| \geq C_1\right) = O(p^{-C_3}).$$

It follows from Lemma 8 that

$$P\left(\max_j|\bar{\gamma}_j| \geq C_1\sqrt{\frac{\log p}{n}}\right) = O(p^{-C_3}).$$

The above two inequalities together imply

$$P\left(\max_{i,j}\left|\frac{1}{n}\sum_{k=1}^{n}\gamma_{ki}^2\gamma_{kj}\bar{\gamma}_j\right| \geq C_1\sqrt{\frac{\log p}{n}}\right) = O(p^{-C_3}). \tag{A.30}$$

We can similarly bound the other terms in (A.29) and obtain

$$P\left(\max_{i,j}|\hat{\theta}_{ij} - \tilde{\theta}_{ij}| \geq C_1\sqrt{\frac{\log p}{n}}\right) = O(p^{-C_3}). \tag{A.31}$$

Next, write

$$\tilde{\theta}_{ij} - \theta_{ij} = \frac{1}{n}\sum_{k=1}^{n}(\gamma_{ki}\gamma_{kj} - \tilde{\gamma}_{ij})^2 - \mathrm{Var}(Y_iY_j)$$

$$= \frac{1}{n}\sum_{k=1}^{n}\gamma_{ki}^2\gamma_{kj}^2 - EY_i^2Y_j^2 - \{\tilde{\gamma}_{ij}^2 - (\omega_{ij}^0)^2\}$$

$$\equiv T_1 + T_2.$$

To bound the term $T_1$, we further write

$$T_1 = \frac{1}{n}\sum_{k=1}^{n}\{(Y_{ki} - \bar{Y}_k)(Y_{kj} - \bar{Y}_k)\}^2 - EY_i^2Y_j^2$$

$$= \frac{1}{n}\sum_{k=1}^{n}\left(Y_{ki}Y_{kj} - Y_{ki}\bar{Y}_k - Y_{kj}\bar{Y}_k + \bar{Y}_k^2\right)^2 - EY_i^2Y_j^2$$

$$= \frac{1}{n}\sum_{k=1}^{n}Y_{ki}^2Y_{kj}^2 - EY_i^2Y_j^2 + \frac{2}{n}\sum_{k=1}^{n}Y_{ki}Y_{kj}(-Y_{ki}\bar{Y}_k - Y_{kj}\bar{Y}_k + \bar{Y}_k^2)$$

$$+ \frac{1}{n}(-Y_{ki}\bar{Y}_k - Y_{kj}\bar{Y}_k + \bar{Y}_k^2)^2.$$

Consider the event $A_1$ on which

$$\max_{i,j,\ell,m} \left| \frac{1}{n} \sum_{k=1}^{n} Y_{ki} Y_{kj} Y_{k\ell} Y_{km} - E Y_i Y_j Y_\ell Y_m \right| \leq \varepsilon_1.$$

Then, on $A_1$, we have

$$\left| \frac{1}{n} \sum_{k=1}^{n} Y_{ki}^2 Y_{kj}^2 - E Y_i^2 Y_j^2 \right| \leq \varepsilon_1.$$

To bound the next term in $T_1$, we write

$$\frac{1}{n} \sum_{k=1}^{n} Y_{ki}^2 Y_{kj} \bar{Y}_k = \frac{1}{n} \sum_{k=1}^{n} Y_{ki}^2 Y_{kj} \bar{Y}_k - E Y_i^2 Y_j \bar{Y} + E Y_i^2 Y_j \bar{Y}$$

$$= \frac{1}{p} \sum_{\ell=1}^{p} \left( \frac{1}{n} \sum_{k=1}^{n} Y_{ki}^2 Y_{kj} Y_{k\ell} - E Y_i^2 Y_j Y_\ell \right) + \frac{1}{p} \sum_{\ell=1}^{p} E Y_i^2 Y_j Y_\ell,$$

which, on $A_1$ and by Condition 6, is bounded by $\varepsilon_1 + s_1(p)/p$. We can similarly bound the other terms in $T_1$ and obtain, on $A_1$,

$$|T_1| \leq 16\varepsilon_1 + 15 s_1(p)/p. \tag{A.32}$$

To bound the term $T_2$, note that

$$\tilde{\gamma}_{ij} - \omega_{ij}^0 = \frac{1}{n} \sum_{k=1}^{n} (Y_{ki} - \bar{Y}_k)(Y_{kj} - \bar{Y}_k) - E Y_i Y_j$$

$$= \frac{1}{n} \sum_{k=1}^{n} Y_{ki} Y_{kj} - E Y_i Y_j + \frac{1}{n} \sum_{k=1}^{n} (-Y_{ki} \bar{Y}_k - Y_{kj} \bar{Y}_k + \bar{Y}_k^2). \tag{A.33}$$

Consider the event $A_2$ on which

$$\max_{i,j} \left| \frac{1}{n} \sum_{k=1}^{n} Y_{ki} Y_{kj} - E Y_i Y_j \right| \leq \varepsilon_2.$$

To bound the next term in (A.33), we write

$$\frac{1}{n} \sum_{k=1}^{n} Y_{ki} \bar{Y}_k = \frac{1}{n} \sum_{k=1}^{n} Y_{ki} \bar{Y}_k - E Y_i \bar{Y} + E Y_i \bar{Y}$$

$$= \frac{1}{p} \sum_{j=1}^{p} \left( \frac{1}{n} \sum_{k=1}^{n} Y_{ki} Y_{kj} - E Y_i Y_j \right) + \frac{1}{p} \sum_{j=1}^{p} \omega_{ij}^0,$$

which, on $A_2$ and by Condition 4, is bounded by $\varepsilon_2 + M^{1-q} s_0(p)/p$. We can similarly bound the

85

other terms in (A.33) and obtain, on $A_2$,

$$|\tilde{\gamma}_{ij} - \omega_{ij}^0| \le 4\varepsilon_2 + 3M^{1-q}s_0(p)/p. \tag{A.34}$$

Note also that, on $A_2$,

$$|\tilde{\gamma}_{ij} + \omega_{ij}^0| \le |\tilde{\gamma}_{ij} - \omega_{ij}^0| + 2|\omega_{ij}^0| \le 4\varepsilon_2 + 3M^{1-q}s_0(p)/p + 2M.$$

Hence, on $A_2$, we have

$$|T_2| = |\tilde{\gamma}_{ij} - \omega_{ij}^0||\tilde{\gamma}_{ij} + \omega_{ij}^0| \le (4\varepsilon_2 + 3M^{1-q}s_0(p)/p)(4\varepsilon_2 + 3M^{1-q}s_0(p)/p + 2M). \tag{A.35}$$

Finally, it follows from Lemma 8 that the event $A_1 \cap A_2$ occurs with probability at least $1 - O(p^{-C_3})$ for all constants $\varepsilon_1, \varepsilon_2 > 0$ and some constant $C_3 > 0$. Combining (A.31), (A.32), and (A.35) and noting $\log p = o(n)$, $s_0(p) = o(p)$, and $s_1(p) = o(p)$, we arrive at (A.27).

It remains to prove (A.28). We first write

$$\begin{aligned}
\hat{\gamma}_{ij} - \tilde{\gamma}_{ij} &= \frac{1}{n}\sum_{k=1}^{n}(\gamma_{ki} - \bar{\gamma}_i)(\gamma_{kj} - \bar{\gamma}_j) - \frac{1}{n}\sum_{k=1}^{n}\gamma_{ki}\gamma_{kj} \\
&= \frac{1}{n}\sum_{k=1}^{n}(-\gamma_{ki}\bar{\gamma}_i - \gamma_{kj}\bar{\gamma}_j + \bar{\gamma}_i\bar{\gamma}_j).
\end{aligned}$$

Using arguments similar to those for proving (A.30), we can show that

$$P\left(\max_{i,j}\left|\frac{1}{n}\sum_{k=1}^{n}\gamma_{ki}\bar{\gamma}_j\right| \ge C_1\sqrt{\frac{\log p}{n}}\right) = O(p^{-C_3}).$$

We can similarly bound the other two terms and obtain

$$P\left(\max_{i,j}|\hat{\gamma}_{ij} - \tilde{\gamma}_{ij}| \ge C_1\sqrt{\frac{\log p}{n}}\right) = O(p^{-C_3}).$$

Taking $\varepsilon_2 = C_1\sqrt{(\log p)/n}$ in (A.34), we have

$$P\left(\max_{i,j}|\tilde{\gamma}_{ij} - \omega_{ij}^0| \ge C_1\sqrt{\frac{\log p}{n}} + C_2\frac{s_0(p)}{p}\right) = O(p^{-C_3}).$$

86

The above two inequalities together imply

$$P\left(\max_{i,j} |\hat{\gamma}_{ij} - \omega_{ij}^0| \geq C_1 \sqrt{\frac{\log p}{n}} + C_2 \frac{s_0(p)}{p}\right) = O(p^{-C_3}). \tag{A.36}$$

From Condition 5 and (A.27) with $\varepsilon_2 = \tau/2$, it follows that $|\hat{\theta}_{ij}| \geq \tau/2$ with probability at least $1 - O(p^{-C_3})$. This, together with (A.36), implies (A.28) and completes the proof.

*A.2.4. Proof of Theorem 4*

By the triangle inequality, we have

$$\|\widehat{\Omega} - \Omega_0\|_1 \leq \sum_{j=1}^{p} |S_{\lambda_{ij}}(\omega_{ij}^0) - \omega_{ij}^0| + \sum_{j=1}^{p} |S_{\lambda_{ij}}(\hat{\gamma}_{ij}) - S_{\lambda_{ij}}(\omega_{ij}^0)|. \tag{A.37}$$

Using Conditions (i) and (ii) that define a general thresholding function, the first term above is bounded by

$$\sum_{j=1}^{p} |\omega_{ij}^0| I(|\omega_{ij}^0| \leq \lambda_{ij}) + \sum_{j=1}^{p} \lambda_{ij} I(|\omega_{ij}^0| > \lambda_{ij})$$

$$= \sum_{j=1}^{p} |\omega_{ij}^0|^q |\omega_{ij}^0|^{1-q} I(|\omega_{ij}^0| \leq \lambda_{ij}) + \sum_{j=1}^{p} \lambda_{ij}^q \lambda_{ij}^{1-q} I(|\omega_{ij}^0| > \lambda_{ij})$$

$$\leq \sum_{j=1}^{p} |\omega_{ij}^0|^q \lambda_{ij}^{1-q}.$$

On the other hand, the second term in (A.37) is bounded by

$$2\sum_{j=1}^{p} |\hat{\gamma}_{ij}| I(|\hat{\gamma}_{ij}| > \lambda_{ij}, |\omega_{ij}^0| \leq \lambda_{ij}) + 2\sum_{j=1}^{p} |\omega_{ij}^0| I(|\hat{\gamma}_{ij}| \leq \lambda_{ij}, |\omega_{ij}^0| > \lambda_{ij})$$

$$+ \sum_{j=1}^{p} |S_{\lambda_{ij}}(\hat{\gamma}_{ij}) - S_{\lambda_{ij}}(\omega_{ij}^0)| I(|\hat{\gamma}_{ij}| > \lambda_{ij}, |\omega_{ij}^0| > \lambda_{ij})$$

$$\equiv T_1 + T_2 + T_3.$$

To bound the term $T_1$, we write

$$\frac{T_1}{2} \leq \sum_{j=1}^{p} |\hat{\gamma}_{ij} - \omega_{ij}^0| I(|\hat{\gamma}_{ij}| > \lambda_{ij}, |\omega_{ij}^0| \leq \lambda_{ij}/2)$$

$$+ \sum_{j=1}^{p} |\hat{\gamma}_{ij} - \omega_{ij}^0| I(|\hat{\gamma}_{ij}| > \lambda_{ij}, \lambda_{ij}/2 < |\omega_{ij}^0| \leq \lambda_{ij}) + \sum_{j=1}^{p} |\omega_{ij}^0| I(|\hat{\gamma}_{ij}| > \lambda_{ij}, |\omega_{ij}^0| \leq \lambda_{ij})$$

$$\equiv T_4 + T_5 + T_6.$$

Consider the event $B_1$ on which $|\hat{\gamma}_{ij} - \omega_{ij}^0| \leq \lambda_{ij}/2$ for all $i, j$. On $B_1$, we have

$$T_4 \leq \sum_{j=1}^{p} |\hat{\gamma}_{ij} - \omega_{ij}^0| I(|\hat{\gamma}_{ij} - \omega_{ij}^0| > \lambda_{ij}/2) = 0,$$

$$T_5 \leq \sum_{j=1}^{p} \left(\frac{\lambda_{ij}}{2}\right)^q \left(\frac{\lambda_{ij}}{2}\right)^{1-q} I(|\hat{\gamma}_{ij}| > \lambda_{ij}, \lambda_{ij}/2 < |\omega_{ij}^0| \leq \lambda_{ij}) \leq \frac{1}{2^{1-q}} \sum_{j=1}^{p} |\omega_{ij}^0|^q \lambda_{ij}^{1-q},$$

and

$$T_6 \leq \sum_{j=1}^{p} |\omega_{ij}^0|^q \lambda_{ij}^{1-q}.$$

Combining these pieces yields

$$T_1 \leq 2 \left(1 + \frac{1}{2^{1-q}}\right) \sum_{j=1}^{p} |\omega_{ij}^0|^q \lambda_{ij}^{1-q} \leq 4 \sum_{j=1}^{p} |\omega_{ij}^0|^q \lambda_{ij}^{1-q}.$$

We can similarly bound the terms $T_2$ and $T_3$ on $B_1$:

$$T_2 \leq 2 \sum_{j=1}^{p} \left(|\hat{\gamma}_{ij} - \omega_{ij}^0| + |\hat{\gamma}_{ij}|\right) I(|\hat{\gamma}_{ij}| \leq \lambda_{ij}, |\omega_{ij}^0| > \lambda_{ij})$$

$$\leq 2 \sum_{j=1}^{p} \left(\frac{\lambda_{ij}}{2} + \lambda_{ij}\right) I(|\hat{\gamma}_{ij}| \leq \lambda_{ij}, |\omega_{ij}^0| > \lambda_{ij}) \leq 3 \sum_{j=1}^{p} |\omega_{ij}^0|^q \lambda_{ij}^{1-q},$$

$$T_3 \leq \sum_{j=1}^{p} \left(|\hat{\gamma}_{ij} - \omega_{ij}^0| + |S_{\lambda_{ij}}(\hat{\gamma}_{ij}) - \hat{\gamma}_{ij}| + |S_{\lambda_{ij}}(\omega_{ij}^0) - \omega_{ij}^0|\right) I(|\hat{\gamma}_{ij}| > \lambda_{ij}, |\omega_{ij}^0| > \lambda_{ij})$$

$$\leq \sum_{j=1}^{p} \left(\frac{\lambda_{ij}}{2} + \lambda_{ij} + \lambda_{ij}\right) I(|\hat{\gamma}_{ij}| > \lambda_{ij}, |\omega_{ij}^0| > \lambda_{ij}) \leq \frac{5}{2} \sum_{j=1}^{p} |\omega_{ij}^0|^q \lambda_{ij}^{1-q}.$$

Collecting all terms, we obtain, on $B_1$,

$$\|\widehat{\Omega} - \Omega_0\|_1 \leq \frac{21}{2} \sum_{j=1}^{p} |\omega_{ij}^0|^q \lambda_{ij}^{1-q}. \tag{A.38}$$

Next, we consider the event $B_2$ on which $|\hat{\theta}_{ij} - \theta_{ij}| \leq \tau$ for all $i, j$. From Condition 5 we have, on $B_2$,

$$\hat{\theta}_{ij} \leq |\hat{\theta}_{ij} - \theta_{ij}| + \theta_{ij} \leq \tau + \theta_{ij} \leq 2\theta_{ij}. \tag{A.39}$$

Note that, by Condition 3,

$$\theta_{ij} \leq EY_i^2 Y_j^2 \leq \frac{1}{2}(EY_i^4 + EY_j^4) \leq \frac{2}{\alpha^2}. \tag{A.40}$$

Taking $\lambda_{ij} = \lambda\sqrt{\hat{\theta}_{ij}}$ with $\lambda = C_1\sqrt{(\log p)/n} + C_2 s_0(p)/p$ in (A.38) and applying (A.39) and (A.40), we obtain, on $B_1 \cap B_2$,

$$\|\widehat{\boldsymbol{\Omega}} - \boldsymbol{\Omega}_0\|_1 \leq \frac{21}{2}\sum_{j=1}^{p} |\omega_{ij}^0|^q \lambda^{1-q}\left(\frac{2}{\alpha}\right)^{1-q} \leq \frac{21}{\alpha}s_0(p)\left(C_1\sqrt{\frac{\log p}{n}} + C_2\frac{s_0(p)}{p}\right)^{1-q}.$$

We conclude the proof by noting that the event $B_1 \cap B_2$ occurs with probability $1 - O(p^{-C_3})$ by Lemma 9 and that the spectral norm is bounded by the matrix $L_1$-norm.

*A.2.5. Proof of Theorem 5*

It follows from Condition (i) and (A.28) that

$$P\left(\hat{\omega}_{ij} \neq 0, \omega_{ij}^0 = 0 \text{ for some } i, j\right) \leq P\left(\max_{i,j} |\hat{\gamma}_{ij} - \omega_{ij}^0| \geq \lambda_{ij}\right)$$
$$= P\left(\max_{i,j} |\hat{\gamma}_{ij} - \omega_{ij}^0|/\sqrt{\hat{\theta}_{ij}} \geq C_1\sqrt{\frac{\log p}{n}} + C_2\frac{s_0(p)}{p}\right) = O(p^{-C_3}),$$

which proves (2.12).

To prove (2.14), note that, by Condition (ii),

$$P\left(\text{sgn}(\hat{\omega}_{ij}) \neq \text{sgn}(\omega_{ij}^0), \omega_{ij}^0 \neq 0 \text{ for some } i, j\right) \leq P\left(|\hat{\gamma}_{ij} - \omega_{ij}^0| \geq |\omega_{ij}^0| - \lambda_{ij} \text{ for some } i, j\right).$$

Also, by taking $\varepsilon = 3\tau/4$ in (A.27), we have, with probability $1 - O(p^{-C_3})$,

$$\left|\sqrt{\hat{\theta}_{ij}} - \sqrt{\theta_{ij}}\right| = \frac{|\hat{\theta}_{ij} - \theta_{ij}|}{\sqrt{\hat{\theta}_{ij}} + \sqrt{\theta_{ij}}} \leq \frac{3\tau/4}{\sqrt{\tau/4} + \sqrt{\tau}} = \frac{\sqrt{\tau}}{2},$$

and hence

$$|\omega_{ij}^0| - \lambda_{ij} \geq C\lambda\sqrt{\theta_{ij}} - \lambda\left(\sqrt{\hat{\theta}_{ij}} - \sqrt{\theta_{ij}} + \sqrt{\theta_{ij}}\right)$$

$$\geq (C-1)\lambda\sqrt{\tau} - \lambda\frac{\sqrt{\tau}}{2} = \left(C - \frac{3}{2}\right)\lambda\sqrt{\tau}$$

for all $i, j$. Now applying (A.36) yields

$$P\left(\mathrm{sgn}(\hat{\omega}_{ij}) \neq \mathrm{sgn}(\omega_{ij}^0), \omega_{ij}^0 \neq 0 \text{ for some } i, j\right) = O(p^{-C_3}),$$

which, together with (2.12), proves the result.

## A.3. Additional Lemmas and Technical Proofs for Chapter 3

### A.3.1. Preliminary Lemmas

Suppose $\log W_d$ is drawn from the distribution of $\log W_d^\star = (\log w_{1,d}^\star, \cdots, \log w_{p,d}^\star)$ with the covariance $\Omega = (\omega_{i,j}) = \mathrm{cov}(\log w_{i,d}^\star, \log w_{j,d}^\star)$. Let $\hat{\Omega}_d = (\hat{\omega}_{i,j,d})$ as the sample covariance of $\log W_d$ $(d = 1, 2)$,

$$\hat{\omega}_{i,j,d} = \frac{1}{n}\sum_{k=1}^{n_d}(\log w_{k,i,d} - \frac{1}{n_d}\sum_{l=1}^{n_d}\log w_{l,i,d})(\log w_{k,j,d} - \frac{1}{n_d}\sum_{l=1}^{n_d}\log w_{l,j,d}).$$

For notational simplicity, the constants $C_1, C_2, \cdots$ below may vary from line to line.

**Lemma 10.** *If the tail distribution of $\log W_d^\star$ follows Condition $4$, for any $t > 0$, we have*

$$\mathrm{pr}(\max_{i,j}|\hat{\omega}_{i,j,d} - \omega_{i,j}|/(\omega_{i,i}\omega_{j,j})^{1/2} \geq t) \leq C_1 p\exp(-C_2 n_d t/2) + p^2 C_3 \exp(-C_4 n_d t^2/4), \quad \text{(A.41)}$$

*where $C_1, C_2, C_3$ and $C_4$ are constants that do not depend on $p$ and $n_d$.*

*Besides, if tail distribution of $\log W_d^\star$ $(d = 1, 2)$ follows Condition $5$. Let*

$$\theta := \max_{i,j} E\left\{(\log w_{i,d}^\star - \mu_{i,d})(\log w_{j,d}^\star - \mu_{j,d}) - \omega_{i,j}\right\}^2/(\omega_{i,i}\omega_{j,j})$$

*which is a bounded constant depending only on $\gamma_0, \epsilon, K$ under the condition $5$, then for any $M > 0$,*

*we have*

$$\mathrm{pr}(\max_{i,j} |\widehat{\omega}_{i,j,d} - \omega_{i,j}| / (\omega_{i,i}\omega_{j,j})^{1/2} \geq \{(\theta+1)(5+M)\log p/n_d\}^{1/2}) = O(n_d^{-\epsilon/8} + p^{-M/2}). \quad \text{(A.42)}$$

**Lemma 11.** *Let $M_n^\star = \frac{n_1 n_2}{n_1 + n_2} \max_{1 \leq i \leq p} \frac{(\bar{y}_{i,1} - \bar{y}_{i,2})^2}{\gamma_{i,i}}$. If equations* (3.10) *and* (3.11) *hold, under the null hypothesis $H_0 : \nu_1 = \nu_2$, for any $t \in \mathbb{R}$ and $\epsilon_{n,p}^{(1)} = o(1)$, we have*

$$\mathrm{pr}(M_n^\star \leq (t + 2\log p - \log\log p)(1 + \epsilon_{n,p}^{(1)}(\log p)^{-1})) \to \exp\left\{-\pi^{-1/2}\exp(-t/2)\right\}, \quad \text{(A.43)}$$

*as $n_1, n_2, p \to \infty$. We define the corresponding $\alpha-$level test $\Phi_\alpha^\star$ by $\Phi_\alpha^\star = I(M_n^\star \geq (q_\alpha + 2\log p - \log\log p)(1 + \epsilon_{n,p}^{(1)}(\log p)^{-1}))$, where $q_\alpha$ is the $1 - \alpha$ quantile of the Type I extreme value distribution function, then*

$$\mathrm{pr}_{H_0}(\Phi_\alpha = 1) \leq -\log(1-\alpha) + o(1). \quad \text{(A.44)}$$

*In addition, under $H_1 : \nu_1 - \nu_2 \in S(k_p)$ with $k_p = p^r, 0 \leq r < 1$, for some $\epsilon > 0$, we have*

$$\lim_{p \to \infty} \mathrm{pr}_{H_1}(\Phi_\alpha^\star = 1) = 1. \text{ if } \beta \geq (1 - \sqrt{r})^2 + \epsilon, \quad \text{(A.45)}$$

$$\overline{\lim}_{p \to \infty} \mathrm{pr}_{H_1}(\Phi_\alpha^\star = 1) \leq \alpha, \text{ if } \beta < (1 - \sqrt{r})^2. \quad \text{(A.46)}$$

*A.3.2. Proof of Theorem 6*

Denote $\omega_{i,-} = p^{-1} \sum_{k=1}^p \omega_{i,k}$ ($i = 1, \cdots, p$) and $\omega_{-,-} = p^{-2} \sum_{k,l=1}^p \omega_{k,l}$, which will be used frequently in the following a few sections.

We first observed that

$$r_{i,j}^{\mathrm{clr}} = \frac{\omega_{i,j} + \epsilon_1}{\{(\omega_{i,i} + \epsilon_2)(\omega_{j,j} + \epsilon_3)\}^{1/2}}, \ 1 \leq i < j \leq p,$$

where $\epsilon_1 = -\omega_{i,-} - \omega_{j,-} + \omega_{-,-}$,   $\epsilon_2 = -2\omega_{i,-} + \omega_{-,-}$ and $\epsilon_3 = -2\omega_{j,-} + \omega_{-,-}$. Under Condition $1, 2$ and $3$, as a result of Proposition $2$, $|\epsilon_i| \leq 3\tau r_3/p, i = 1, 2, 3$. Therefore, by using this inequality

91

and $1/\tau \leq \omega_{i,i} \leq \tau$, when $p$ is sufficiently large, we have

$$r_{i,j}^{\mathrm{clr}} = \frac{\omega_{i,j} + \epsilon_1}{(\omega_{i,i}\omega_{j,j})^{1/2}} \cdot \frac{(\omega_{i,i}\omega_{j,j})^{1/2}}{\{(\omega_{i,i} + \epsilon_2)(\omega_{j,j} + \epsilon_3)\}^{1/2}}$$

$$\leq (r_{i,j} + |\epsilon_1|\tau)/\{(1 - |\epsilon_2|\tau)(1 - |\epsilon_3|\tau)\}^{1/2}$$

$$\leq r_{i,j} + O(r_3/p).$$

Similarly, we obtain $r_{i,j}^{\mathrm{clr}} \geq r_{i,j} - O(r_3/p)$. According to Proposition 1, $r_3 = O(p^{1/2})$, thus, uniformly in $1 \leq i < j \leq p$, we have that $\left|r_{i,j}^{\mathrm{clr}} - r_{i,j}\right| = o(1)$. Using $|r_{i,j}^{clr} - r_{i,j}| \leq Cr_3/p$ and Taylor's expansion that $(r_{i,j}^{clr})^2 = r_{i,j}^2 + 2r_{i,j}(r_{i,j}^{clr} - r_{i,j}) + (r_{i,j}^{clr} - r_{i,j})^2$, we have

$$\left|\sum_{i=1}^{p}((r_{i,j}^{clr})^2 - r_{i,j}^2)\right| \leq \sum_{i=1}^{p}\left|(r_{i,j}^{clr})^2 - r_{i,j}^2\right|$$

$$\leq \sum_{i=1}^{p}(2|r_{i,j}||r_{i,j}^{clr} - r_{i,j}| + |r_{i,j}^{clr} - r_{i,j}|^2)$$

$$\leq (2C + C^2)r_3^2/p.$$

As a consequence,

$$\sum_{i=1}^{p}(r_{i,j}^{\mathrm{clr}})^2 \leq \sum_{i=1}^{p}r_{i,j}^2 + O(r_3^2/p), \quad \sum_{i=1}^{p}(r_{i,j}^{\mathrm{clr}})^2 \geq \sum_{i=1}^{p}r_{i,j}^2 - O(r_3^2/p).$$

The proof of Theorem $6$ is completed, as $r_3 = O(p^{1/2})$.

*A.3.3. Proof of Theorem $7$*

Proof of equation (3.11): Note that,

$$|y_{k,i,d} - \nu_{i,d}| = \left|\log w_{k,i,d} - p^{-1}\sum_{j=1}^{p}\log w_{k,j,d} - \left(\mu_{i,d} - p^{-1}\sum_{j=1}^{p}\mu_{j,d}\right)\right|.$$

Under Conditions $1$, $2$ and $3$, by using the property (3.8), it follows that, uniformly in $1 \leq k \leq n_d$ and $1 \leq i \leq p$, $|y_{k,i,d} - \nu_{i,d}|/\gamma_{i,i}^{1/2} \leq 2\sqrt{2}\tau \max_{k,i}|\log w_{k,i,d} - \mu_{i,d}|/\omega_{i,i}^{1/2}$, $(d = 1, 2)$. Thus, if $\log W_d^{\star}$ $(d = 1, 2)$ follow Condition $4$ (sub-Gaussian-type tails), let $\tau_n = \left\{8\tau^2((M+1)\log p + \log n_d)/\eta\right\}^{1/2}$

with some constant $M > 0$, then, as $n_d, p \to \infty$,

$$\text{pr}\Big(\max_{1 \le k \le n_d, 1 \le i \le p} |y_{k,i,d} - \nu_{i,d}| / \gamma_{i,i}^{1/2} \ge \tau_n\Big) \le \text{pr}\Big(\max_{1 \le k \le n_d, 1 \le i \le p} |\log w_{k,i,d} - \mu_{i,d}| / \omega_{i,i}^{1/2} \ge \tau_n/(2\sqrt{2}\tau)\Big)$$

$$\le K n_d p \exp(-\eta \tau_n^2/(8\tau^2)) \to 0.$$

By using $p \gg n$ and $\log p = o(n_d^{1/4})$, we can verify that $\tau_n = o(n_d^{1/2}/(\log p)^{3/2})$. If $\log W_d^\star$ follow Condition $5$ (polynomial-type tails), let $\tau_n = 2\sqrt{2}\tau n_d^{1/4}$, then as $n_d, p \to \infty$,

$$\text{pr}\Big(\max_{1 \le k \le n_d, 1 \le i \le p} |y_{k,i,d} - \nu_{i,d}| / \gamma_{i,i}^{1/2} \ge \tau_n\Big) \le K n_d p (\tau_n/(2\sqrt{2}\tau))^{-4\gamma_0 - 4 - \epsilon} \to 0.$$

Under the condition that $p = O(n_d^{\gamma_0})$, it can be verified that $\tau_n = o(n_d^{1/2}/(\log p)^{3/2})$.

Proof of equation $(3.12)$: Without the loss of generality, assume $\mu_1 = \mu_2 = 0$. Denote by $\widehat{\gamma}_{i,i,d}$ as the sample covariance of $Y_{i,i}^{(d)}$. Using the notation defined in §A.3.2 and the equation (3.7), we obtain

$$|\widehat{\gamma}_{i,i,d} - \gamma_{i,i}| = |(\widehat{\omega}_{i,i,d} - 2\widehat{\omega}_{i,-,d} + \widehat{\omega}_{-,-}) - (\omega_{i,i} - 2\omega_{i,-} + \omega_{-,-})| \le 4 \max_{i,j} |\widehat{\omega}_{i,j,d} - \omega_{i,j}|.$$

Therefore, under Condition $2$ and $3$, if $\log W_d^\star$ $(d = 1, 2)$ follow Condition $4$, as a result of the inequality (A.41) and (3.8), uniformly in $1 \le i \le p$ for any $t > 0$,

$$\text{pr}(|\widehat{\gamma}_{i,i,d} - \gamma_{i,i}|/\gamma_{i,i} \ge t) \le \text{pr}\Big(4 \max_{i,j} \frac{|\widehat{\omega}_{i,j,d} - \omega_{i,j}|}{(\omega_{i,i}\omega_{j,j})^{1/2}} \cdot \frac{(\omega_{i,i}\omega_{j,j})^{1/2}}{\min_k \gamma_{k,k}} \ge t\Big)$$

$$\le \text{pr}\Big(\max_{i,j} |\widehat{\omega}_{i,j,d} - \omega_{i,j}|/(\omega_{i,i}\omega_{j,j})^{1/2} \ge t/(8\tau^2)\Big)$$

$$\le 2p \exp(-C_1 n_d t/(16\tau^2)) + p^2 C_2 \exp(-C_3 n_d t^2/(256\tau^4)).$$

Since $\widehat{\gamma}_{i,i} = \frac{n_1}{n_1+n_2}\widehat{\gamma}_{i,i,1} + \frac{n_2}{n_1+n_2}\widehat{\gamma}_{i,i,2}$, and $n_1$ and $n_2$ are comparable, it yields that, for some constants $C_4, C_5$ and $C_6$, we have

$$\text{pr}(|\widehat{\gamma}_{i,i} - \gamma_{i,i}|/\gamma_{i,i} \ge t) \le 4p \exp(-C_4 n t/\tau^2) + C_5 p^2 \exp(-C_6 n t^2/\tau^4).$$

It implies $|\widehat{\gamma}_{i,i} - \gamma_{i,i}| = O_p\Big\{(\log p/n)^{1/2}\Big\} \gamma_{i,i}$. Similarly, under Condition $2$ and $3$. If $\log W_d^\star$ $(d = 1, 2)$

follow Condition $5$, by use of (A.42), uniformly in $1 \leq i \leq p$ for any $M > 0$,

$$\mathrm{pr}(|\widehat{\gamma}_{i,i,d} - \gamma_{i,i}|/\gamma_{i,i} \geq \left\{64\tau^4(\theta + 1)(5 + M)\log p/n_d\right\}^{1/2})$$

$$\leq \mathrm{pr}(\max_{i,j}|\widehat{\omega}_{i,j,d} - \omega_{i,j}|/(\omega_{i,i}\omega_{j,j})^{1/2} \geq \{(\theta + 1)(5 + M)(\log p/n_d)\}^{1/2}) = O(n_d^{-\epsilon/8} + p^{-M/2}).$$

As a result, $|\widehat{\gamma}_{i,i} - \gamma_{i,i}| = O_p\left\{(\log p/n)^{1/2}\right\}\gamma_{i,i}$, which completes the proof.

### A.3.4. Proof of Theorem 8 & 9

**Proof 1.** We only prove (3.13), the proofs of (3.14), (3.17) and (3.18) are similar. Proof of (3.13): In the event $\{|\widehat{\gamma}_{i,i}/\gamma_{i,i} - 1| \leq C_1(\log p/n)^{1/2}\}$, we have $|M_n - M_n^\star| \leq C_2 M_n^\star(\log p/n)^{1/2}$. Since $(\log p/n)^{1/2} = o(1/\log p)$ under Conditions $4$ and $5$, the proof is completed by applying Theorem $7$ and Lemma $11$.

### A.3.5. Proof of Technical Lemmas

**Lemma 12.** (*Bonferroni Inequality*) Let $A = \bigcup_{l=1}^{p} A_l$. For any $k < [p/2]$, we have

$$\sum_{l=1}^{2k}(-1)^{l-1}E_l \leq \mathrm{pr}(A) \leq \sum_{l=1}^{2k-1}(-1)^{l-1}E_l,$$

where $E_t = \sum_{1 \leq i_1 < \cdots < i_l \leq p}\mathrm{pr}(A_{i_1} \cap \cdots \cap A_{i_l})$.

**Lemma 13.** Let $(Z_1, \cdots, Z_p)^{\mathrm{T}}$ be a zero-mean multivariate normal random vector with covariance matrix $\Omega = (\omega_{i,j})_{1 \leq i,j \leq p}$ and diagonal $\omega_{i,i} = 1$ for $1 \leq i \leq p$. Suppose that $\max_{1 \leq i < j \leq p}|\omega_{i,j}| \leq r_1 < 1$ and $\max_j \sum_{i=1}^{p}\omega_{i,j}^2 \leq r_2 < \infty$, where $r_1$ and $r_2$ are some constants. For any fixed integer $d^\star \geq 1$, let the vector $N_{d^\star} = (Z_{k_1}, \cdots, Z_{k_{d^\star}})$ with the index $\{k_1, \cdots, k_{d^\star}\}$. Then for any $\epsilon_{n,p} = o(1)$ and any $t \in \mathbb{R}$,

$$\sum_{1 \leq k_1 < \cdots < k_{d^\star} \leq p}\mathrm{pr}(|N_{d^\star}|_{min} \geq (t + 2\log p - \log\log p)^{1/2} \pm \epsilon_{n,p}(\log p)^{-1/2})$$

$$= (\pi^{-1/2}\exp(-t/2))^{d^\star}/d^\star! \cdot (1 + o(1)).$$

$$\tag{A.47}$$

*A.3.6. Proof of Lemma* A10

Proof of (A.41): Without loss of generality, we assume $\mu_1 = \mu_2 = 0$ and $\omega_{i,i} = 1, i = 1 \cdots, p$, then it suffices to prove that, for $d = 1, 2$,

$$\text{pr}(\max_i \left| \frac{1}{n_d} \sum_{k=1}^{n_d} \log w_{k,i,d} \right| \geq t) \leq C_1 p \exp(-C_2 n_d t^2), \tag{A.48}$$

$$\text{pr}(\max_{i,j} \left| \frac{1}{n_d} \sum_{k=1}^{n} \log w_{k,i,d} \log w_{k,j,d} - E \log w_{i,d}^\star \log w_{j,d}^\star \right| \geq t) \leq p^2 C_3 \exp(-C_4 n_d t^2). \tag{A.49}$$

Proof of inequality (A.48) and (A.49) are the results from Exercise 2.27 in (Boucheron, Lugosi, and Massart, 2013); Also see (Bickel and Levina, 2008).

Proof of (A.42): See Theorem 1(b) in (Cai, Liu, and Luo, 2011).

*A.3.7. Proof of Lemma* A11

Proof of (A.43): Let $t_p = (t + 2 \log p - \log \log p)(1 + \epsilon_{n,p}^{(1)} (\log p)^{-1})$, according to Lemma B12, for any $m \in \mathbb{Z}, 0 < m < q/2$,

$$\sum_{d^\star=1}^{2m} (-1)^{d^\star - 1} \sum_{1 \leq k_1 < \cdots < k_d^\star \leq p} \text{pr}(\bigcap_{j=1}^{d} E_{k_j}) \leq \text{pr}(\frac{n_1 n_2}{n_1 + n_2} \max_{1 \leq i \leq p} \frac{(\bar{y}_{i,1} - \bar{y}_{i,2})^2}{\gamma_{i,i}} \geq t_p)$$

$$\leq \sum_{d^\star=1}^{2m-1} (-1)^{d^\star - 1} \sum_{1 \leq k_1 < \cdots < k_{d^\star} \leq p} \text{pr}(\bigcap_{j=1}^{d^\star} E_{k_j}),$$

where $E_{k_j} = \left\{ \left| \frac{\bar{y}_{k_j,1} - \bar{y}_{k_j,2}}{\sqrt{\gamma_{k_j,k_j}(1/n_1 + 1/n_2)}} \right| \geq t_p^{1/2} \right\}$, $j = 1, \cdots, d^\star$. For a fixed $d^\star$, define the vectors $(T_1, \cdots, T_{n_1}, T_{n_1+1}, \cdots, T_{n_1+n_2})$ as

$$T_l = (\frac{1}{n_1} \cdot \frac{y_{l,k_1,1} - \nu_{k_1,1}}{\sqrt{\gamma_{k_1,k_1}(1/n_1 + 1/n_2)}}, \cdots, \frac{1}{n_1} \cdot \frac{y_{l,k_d,1} - \nu_{k_d,1}}{\sqrt{\gamma_{k_d,k_d}(1/n_1 + 1/n_2)}})^\text{T}, \, l = 1, \cdots, n_1,$$

$$T_{l+n_1} = (-\frac{1}{n_2} \cdot \frac{y_{l,k_1,2} - \nu_{k_1,2}}{\sqrt{\gamma_{k_1,k_1}(1/n_1 + 1/n_2)}}, \cdots, -\frac{1}{n_2} \cdot \frac{y_{l,k_d,2} - \nu_{k_d,2}}{\sqrt{\gamma_{k_d,k_d}(1/n_1 + 1/n_2)}})^\text{T}, \, l = 1, \cdots, n_2.$$

Note that we define $|\alpha|_{min} = \min_{1 \le i \le d} |\alpha_i|$ for any vector $\alpha \in \mathbb{R}^{d^*}$. Then we obtain

$$\text{pr}(\bigcap_{j=1}^{d^*} E_{k_j}) = \text{pr}(\left|\sum_{k=1}^{n_1+n_2} T_k\right|_{min} \ge t_p^{1/2}).$$

Denote the multivariate normal random vector $N_{d^*} = (N_{k_1}, \cdots, N_{k_{d^*}})$ with $EN_{d^*} = 0$ and $\text{cov}(N_{d^*}) = n_1\text{cov}(T_1) + n_2\text{cov}(T_{n_1+1}) = R$. Let $\lambda = \epsilon_{n,p}^{(2)}(\log p)^{-1/2}d^{*1/2}$. Note that under the event that $\left\{\max_{1 \le k \le n_d, 1 \le i \le p}\left|(y_{k,i,d} - \nu_{i,d})/\gamma_{i,i}^{1/2}\right| \le \tau_n\right\}$ $(d = 1, 2)$, we have $|T_{k,i}| \le \tau_n/n^{1/2}$ for any $1 \le k \le n_1 + n_2$, $1 \le i \le p$, then by using Theorem 1 in (Zaitsev, 1987), we have

$$\text{pr}(\left|\sum_{k=1}^{n_1+n_2} T_k\right|_{min} \ge t_p^{1/2}) \le \text{pr}(|N_{d^*}|_{min} \ge t_p^{1/2} - \lambda/d^{*1/2}) + C_1 d^{*5/2}\exp(-\frac{\lambda n^{1/2}}{C_2 d^{*5/2}\tau_n})$$

$$= \text{pr}(|N_{d^*}|_{min} \ge t_p^{1/2} - \epsilon_{n,p}^{(2)}(\log p)^{-1/2}) + C_1 d^{*5/2}\exp(-\frac{\epsilon_{n,p}^{(2)}n^{1/2}}{C_2 d^{*2}(\log p)^{1/2}\tau_n}).$$

Since $\lim_{n,p \to \infty}(\log p)^{3/2}\tau_n/n^{1/2} \to 0$, we set $\epsilon_{n,p}^{(2)} = MC_6 d^{*2}(\log p)^{3/2}\tau_n/n^{1/2}$ which is $o(1)$ for any constant $M > 0$, and

$$C_1 d^{*5/2}\exp(-\frac{\epsilon_{n,p}^{(2)}n^{1/2}}{C_2 d^{*2}(\log p)^{1/2}\tau_n}) = O(p^{-M}). \tag{A.50}$$

By Taylor's expansion, we also observed that $t_p^{1/2} = (t + 2\log p - \log\log p)^{1/2} + C_3\epsilon_{n,p}^{(1)}(\log p)^{-1/2}$. Therefore, setting $\epsilon_{n,p}^{(3)} = C_3\epsilon_{n,p}^{(1)} - \epsilon_{n,p}^{(2)}$ and using equation (3.11) and (A.50), we have

$$\text{pr}(\max_{1 \le i \le p}\frac{(\bar{y}_{i,1} - \bar{y}_{i,2})^2}{\gamma_{i,i}(1/n_1 + 1/n_2)} \ge t_p)$$

$$\le \sum_{d^*=1}^{2m-1}(-1)^{d^*-1}\sum_{1 \le k_1 < \cdots < k_{d^*} \le p}\text{pr}(|N_{d^*}|_{min} \ge (t + 2\log p - \log\log p)^{1/2} + \epsilon_{n,p}^{(3)}(\log p)^{-1/2}) + o(1).$$

$$\tag{A.51}$$

Similarly, we also obtain

$$\text{pr}(\max_{1 \le i \le p}\frac{(\bar{y}_{i,1} - \bar{y}_{i,2})^2}{\gamma_{i,i}(1/n_1 + 1/n_2)} \ge t_p)$$

$$\ge \sum_{d^*=1}^{2m}(-1)^{d^*-1}\sum_{1 \le k_1 < \cdots < k_{d^*} \le p}\text{pr}(|N_{d^*}|_{min} \ge (t + 2\log p - \log\log p)^{1/2} + \epsilon_{n,p}^{(3)}(\log p)^{-1/2}) - o(1).$$

$$\tag{A.52}$$

Combining (A.51) and (A.52) and applying Lemma B13 yields

$$\sum_{d^\star=1}^{2m-1} (-1)^{d^\star-1}(\pi^{-1/2}\exp(-t/2))^{d^\star}/d^\star! \leq \lim_{n_1,n_2,p\to\infty} \mathrm{pr}\Big(\max_{1\leq i\leq p} \frac{(\bar{y}_{i,1}-\bar{y}_{i,2})^2}{\gamma_{i,i}(1/n_1+1/n_2)} \geq t_p\Big)$$

$$\leq \sum_{d^\star=1}^{2m} (-1)^{d^\star-1}(\pi^{-1/2}\exp(-t/2))^{d^\star}/d^\star!.$$

As $p \to \infty$, the proof is completed by letting $m \to \infty$ on both hand sides.

Proof of (A.44): Follow the same procedure as the proof of (A.43). Let $d^\star = 1$ in (A.51), by use of $-\pi^{-1/2}\exp(-q_\alpha/2) = \log(1-\alpha)$, we finally obtain

$$\mathrm{pr}_{H_0}(\Phi_\alpha^\star = 1) = \mathrm{pr}(M_n^\star \geq (q_\alpha + 2\log p - \log\log p)(1 + \epsilon_{n,p}^{(1)}(\log p)^{-1}))$$

$$\leq \sum_{i=1}^{p} \mathrm{pr}(|N(0,1)| \geq (q_\alpha + 2\log p - \log\log p)^{1/2} - \epsilon_{n,p}^{(3)}(\log p)^{-1/2}) + o(1)$$

$$\leq \frac{2p\exp\left\{-((q_\alpha + 2\log p - \log\log p)^{1/2} - \epsilon_{n,p}^{(3)}(\log p)^{-1/2})^2/2\right\}}{(2\pi)^{1/2}\left\{(q_\alpha + 2\log p - \log\log p)^{1/2} - \epsilon_{n,p}^{(3)}(\log p)^{-1/2}\right\}} + o(1)$$

$$= -\log(1-\alpha) \cdot (1 + o(1)) + o(1).$$

Proof of (A.45) and (A.46): See the proof of Proposition 2 in supplement to (Cai, Liu, and Xia, 2014).

*A.3.8. Proof of Lemma* 13

See the proof of Lemma 6 in Cai, Liu, and Xia, 2014.

# BIBLIOGRAPHY

Aitchison, J (1982). The statistical analysis of compositional data. *Journal of the Royal Statistical Society* Series B, 44, 139–177.

Aitchison, J (2003). *The statistical analysis of compositional data*. Caldwell, NJ: Blackburn Press.

Aitchison, J and Shen, SM (1980). Logistic-normal distributions: some properties and uses. *Biometrika* 67, 261–272.

Arumugam, M, Raes, J, Pelletier, E, Le Paslier, D, Yamada, T, Mende, DR, Fernandes, GR, Tap, J, Bruls, T, Batto, J-M, Bertalan, M, Borruel, N, Casellas, F, Fernandez, L, Gautier, L, Hansen, T, Hattori, M, Hayashi, T, Kleerebezem, M, Kurokawa, K, Leclerc, M, Levenez, F, Manichanh, C, Nielsen, HB, Nielsen, T, Pons, N, Poulain, J, Qin, J, Sicheritz-Ponten, T, Tims, S, Torrents, D, Ugarte, E, Zoetendal, EG, Wang, J, Guarner, F, Pedersen, O, Vos, WM de, Brunak, S, Doré, J, MetaHIT Consortium, Weissenbach, J, Ehrlich, SD, and Bork, P (2011). Enterotypes of the human gut microbiome. *Nature* 473, 174–180.

Bäckhed, F, Ding, H, Wang, T, Hooper, LV, Koh, GY, Nagy, A, Semenkovich, CF, and Gordon, JI (2004). The gut microbiota as an environmental factor that rgulates fat storage. *Proceedings of the National Academy of Sciences of the United States of America* 101.44, 15718–15723.

Bai, Z and Saranadasa, H (1996). Effect of high dimension: by an example of a two sample problem. *Statst. Sinica* 6, 311–29.

Ban, Y, An, L, and Jiang, H (2015). Investigating microbial co-ocurrence patterns based on metagenomic compositional data. *Bioinformatics* 31, 3322–3329.

Becker, SR, Candès, EJ, and Grant, MC (2011). Templates for convex cone problems with applications to sparse signal recovery. *Mathematical programming computation* 3.3, 165–218.

Bickel, PJ and Levina, E (2008). Covariance regularization by thresholding. *The Annals of Statistics* 36, 2577–2604.

Boucheron, S, Lugosi, G, and Massart, P (2013). *Concentration inequalities: a nonasymptotic theory of independence*. Oxford: Oxford University Press.

Bühlmann, P and Van De Geer, S (2011). *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media.

Cai, J-F, Candès, EJ, and Shen, Z (2010). A singular value thresholding algorithm for matrix completion. *SIAM Journal on Optimization* 20.4, 1956–1982.

Cai, TT and Zhou, HH (2012). Optimal rates of convergence for sparse covariance matrix estimation. *The Annals of Statistics* 40, 2389–2420.

Cai, T and Liu, W (2011). Adaptive thresholding for sparse covariance matrix estimation. *Journal of the American Statistical Association* 106, 672–684.

Cai, T, Liu, W, and Luo, X (2011). A constrained $\ell_1$ minimization approach to sparse precision matrix estimation. *J. Am. Statist. Assoc.* 106, 594–607.

Cai, T, Liu, W, and Xia, Y (2014). Two-sample test of high dimensional means under dependence. *J. R. Statist. Soc.* B 76, 349–72.

Cao, Y and Xie, Y (2016). Poisson matrix recovery and completion. *IEEE Transactions on Signal Processing* 64.6, 1609–1620.

Cao, Y, Lin, W, and Li, H (2016). Large covariance estimation for compositional data via composition-adjusted thresholding. *arXiv preprint arXiv:1601.04397*.

Chaffron, S, Rehrauer, H, Pernthaler, J, and Mering, C von (2010). A global network of coexisting microbes from environmental and whole-genome sequence data. *Genome research* 20.7, 947–959.

Chen, SX and Qin, Y-L (2010). A two-sample test for high-dimensional data with applications to gene-set testing. *Ann. Statist.* 38, 808–35.

Coyte, KZ, Schluter, J, and Foster, KR (2015). The ecology of the microbiome: networks, competition, and stability. *Science* 350, 663–666.

El Karoui, N (2008). Operator norm consistent estimation of large-dimensional sparse covariance matrices. *The Annals of Statistics* 36, 2717–2756.

Fan, J, Fan, Y, and Lv, J (2008). High dimensional covariance matrix estimation using a factor model. *Journal of Econometrics* 147, 186–197.

Fan, J, Liao, Y, and Mincheva, M (2013). Large covariance estimation by thresholding principal orthogonal complements. *Journal of the Royal Statistical Society* Series B, 75, 603–680.

Fang, H, Huang, C, Zhao, H, and Deng, M (2015). CCLasso: correlation inference for compositional data through Lasso. *Bioinformatics* 31, 3172–3180.

Faust, K, Sathirapongsasuti, JF, Izard, J, Segata, N, Gevers, D, Raes, J, and Huttenhower, C (2012). Microbial co-occurrence relationships in the human microbiome. *PLoS Comput Biol* 8.7, e1002606.

Fisher, CK and Mehta, P (2014). Identifying keystone species in the human gut microbiome from metagenomic timeseries using sparse linear regression. *PLoS ONE* 9, e102451.

Friedman, J and Alm, EJ (2012). Inferring correlation networks from genomic survey data. *PLoS Computational Biology* 8, e1002687.

Greenblum, S, Turnbaugh, PJ, and Borenstein, E (2012). Metagenomic systems biology of the human gut microbiome reveals topological shifts associated with obesity and inflammatory bowel disease. *Proceedings of the National Academy of Sciences* 109, 594–599.

Horner-Devine, MC, Silver, JM, Leibold, MA, Bohannan, BJ, Colwell, RK, Fuhrman, JA, Green, JL, Kuske, CR, Martiny, JB, Muyzer, G, et al. (2007). A comparison of taxon co-occurrence patterns for macro-and microorganisms. *Ecology* 88.6, 1345–1353.

Isserlis, L (1918). On a formula for the product-moment coefficient of any order of a normal frequency distribution in any number of variables. *Biometrika* 12, 134–139.

Jordán, F, Lauria, M, Scotti, M, Nguyen, T-P, Praveen, P, Morine, M, and Priami, C (2015). Diversity of key players in the microbial ecosystems of the human body. *Scientific Reports* 5, 15920.

Klopp, O, Lafond, J, Moulines, r, and Salmon, J (2015). Adaptive multinomial matrix completion. *Electron. J. Statist.* 9.2, 2950–2975.

Koeth, RA, Wang, Z, Levison, BS, Buffa, JA, Org, E, Sheehy, BT, Britt, EB, Fu, X, Wu, Y, Li, L, Smith, JD, DiDonato, JA, Chen, J, Li, H, Wu, GD, Lewis, JD, Warrier, M, Brown, JM, Krauss, RM, Tang, WHW, Bushman, FD, Lusis, AJ, and Hazen, SL (2013). Intestinal microbiota metabolism of L-carnitine, a nutrient in red meat, promotes atherosclerosis. *Nature Medicine* 19, 576–585.

Lafond, J, Klopp, O, Moulines, E, and Salmon, J (2014). Probabilistic low-rank matrix completion on finite alphabets. In: *Advances in Neural Information Processing Systems 27*. Ed. by Z Ghahramani, M Welling, C Cortes, N Lawrence, and K Weinberger. Curran Associates, Inc., 1727–1735.

Ledoux, M and Talagrand, M (2013). *Probability in Banach Spaces: isoperimetry and processes*. Springer Science & Business Media.

Lewis, JD, Chen, EZ, Baldassano, RN, Otley, AR, Griffiths, AM, Lee, D, Bittinger, K, Bailey, A, Friedman, ES, Hoffmann, C, Albenberg, L, Sinha, R, Compher, C, Gilroy, E, Nessel, L, Grant, A, Chehoud, C, Li, H, Wu, GD, and Bushman, FD (2015). Inflammation, antibiotics, and diet as environmental stressors of the gut microbiome in pediatric crohn's disease. *Cell Host & Microbe* 18, 489–500.

Li, H (2015). Microbiome, metagenomics, and high-dimensional compositional data analysis. *Ann. Rev. Statist. Appl.* 2, 73–94.

Lin, W, Shi, P, Feng, R, and Li, H (2014). Variable selection in regression with compositional covariates. *Biometrika* 101, 785–797.

Lu, Y and Negahban, SN (2014). Individualized rank aggregation using nuclear norm regularization. *arXiv preprint arXiv:1410.0860*.

Martın-Fernandez, JA, Palarea-Albaladejo, J, and Olea, RA (2011). Dealing with zeros. *Compositional data analysis: Theory and applications*, 43–58.

Negahban, S and Wainwright, MJ (2011). Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *The Annals of Statistics*, 1069–1097.

Negahban, S and Wainwright, MJ (2012). Restricted strong convexity and weighted matrix completion: optimal bounds with noise. *The Journal of Machine Learning Research* 13.1, 1665–1697.

Nesterov, Y (1983). "A method of solving a convex programming problem with convergence rate O (1/k2)". In: *Soviet Mathematics Doklady*. Vol. 27. 2, 372–376.

Nesterov, Y (2013). *Introductory lectures on convex optimization: a basic course*. Vol. 87. Springer Science & Business Media.

Recht, B (2011). A simpler approach to matrix completion. *The Journal of Machine Learning Research* 12, 3413–3430.

Rothman, AJ, Levina, E, and Zhu, J (2009). Generalized thresholding of large covariance matrices. *Journal of the American Statistical Association* 104, 177–186.

Segata, N, Waldron, L, Ballarini, A, Narasimhan, V, Jousson, O, and Huttenhower, C (2012). Metagenomic microbial community profiling using unique clade-specific marker genes. *Nature methods* 9, 811–814.

Srivastava, MS (2009). A test for the mean vector with fewer observations than the dimension under non-normality. *J. Mult. Anal.* 100, 518–32.

Su, W, Boyd, S, and Candes, E (2014). "A differential equation for modeling Nesterovs accelerated gradient method: theory and insights". In: *Advances in Neural Information Processing Systems*, 2510–2518.

The Human Microbiome Project Consortium (2012a). A framework for human microbiome research. *Nature* 486, 215–221.

The Human Microbiome Project Consortium (2012b). Structure, function and diversity of the healthy human microbiome. *Nature* 486, 207–14.

Tropp, JA (2011). *User-friendly tail bounds for matrix martingales*. Tech. rep. DTIC Document.

Tropp, JA (2015). The expected norm of a sum of independent random matrices: An elementary approach. *arXiv preprint arXiv:1506.04711*.

Turnbaugh, PJ, Ley, RE, Mahowald, MA, Magrini, V, Mardis, ER, and Gordon, JI (2006). An obesity-associated gut microbiome with increased capacity for energy harvest. *Nature* 444, 1027–1031.

Turnbaugh, PJ, Hamady, M, Yatsunenko, T, Cantarel, BL, Duncan, A, Ley, RE, Sogin, ML, Jones, WJ, Roe, BA, Affourtit, JP, Egholm, M, Henrissat, B, Heath, AC, Knight, R, and Gordon, JI (2009). A core gut microbiome in obese and lean twins. *Nature* 457, 480–484.

Woyke, T, Teeling, H, Ivanova, NN, Huntemann, M, Richter, M, Gloeckner, FO, Boffelli, D, Anderson, IJ, Barry, KW, Shapiro, HJ, et al. (2006). Symbiosis insights through metagenomic analysis of a microbial consortium. *Nature* 443.7114, 950–955.

Wu, GD, Chen, J, Hoffmann, C, Bittinger, K, Chen, Y-Y, Keilbaugh, SA, Bewtra, M, Knights, D, Walters, WA, Knight, R, Sinha, R, Gilroy, E, Gupta, K, Baldassano, R, Nessel, L, Li, H, Bushman, FD, and Lewis, JD (2011). Linking long-term dietary patterns with gut microbial enterotypes. *Science* 334, 105–108.

Yu, B (1997). Assouad, fano, and le cam. In: *Festschrift for Lucien Le Cam*. Springer, 423–435.

Zaitsev, AY (1987). On the Gaussian approximation of convolutions under multidimensional analogues of SN Bernstein's inequality conditions. *Probability theory and related fields* 74, 535–566.