

# Florida State University Libraries

---

Electronic Theses, Treatises and Dissertations

The Graduate School

---

## Statistical Analysis for Complex Data by Generalized Indirect Dependency Learning and Slack Empirical Likelihood

Peng Zhao

FLORIDA STATE UNIVERSITY  
COLLEGE OF ARTS AND SCIENCES

STATISTICAL ANALYSIS FOR COMPLEX DATA BY GENERALIZED INDIRECT  
DEPENDENCY LEARNING AND SLACK EMPIRICAL LIKELIHOOD

By  
PENG ZHAO

A Dissertation submitted to the  
Department of Statistics  
in partial fulfillment of the  
requirements for the degree of  
Doctor of Philosophy

2020

Copyright © 2020 Peng Zhao. All Rights Reserved.

Peng Zhao defended this dissertation on July 7, 2020.  
The members of the supervisory committee were:

Yiyuan She  
Professor Directing Dissertation

Zhenghao Zhang  
University Representative

Fred W. Huffer  
Committee Member

Jonathan R. Bradley  
Committee Member

The Graduate School has verified and approved the above-named committee members, and certifies that the dissertation has been approved in accordance with university requirements.

To my parents.

# ACKNOWLEDGMENTS

First of all, I would like to express my sincere gratitude to my advisors Dr. Yiyuan She. His immense knowledge and deep views on statistics and optimization strongly support my Ph.D. study. Beyond that, he helped me tremendously on my English writings and presentations. In addition, his personalities, including great diligence, critical thinking and conscientiousness not only inspire me in my research, but also raise me up on my future career pursuit. Without his outstanding help, there is no possibility for me to finish my Ph.D. study.

I would also like to thank my committee members: Dr. Zhenghao Zhang, Dr. Fred Huffer and Dr. Bradley. Their instructive comments and strong encouragement are great motivations for me to write this dissertations during such a hard outbreaking period. Special thanks to Dr. Yun Yang, who led me to the word of statistical research and taught me a lot of knowledge in nonparametric Bayesian and implicit regularizations.

Moreover, I would also like to thank the people who helped me during my years of graduate study. Special thanks my peers Wenjing, Xiaoyang and Yuqing, the time we spent together is an unforgettable experience during my life. Thanks my friends Jiahui, Jingze and Dongrui, I still remember every time we hung out together for dinners. I would like to thank my roommates Yanpeng, Zihan and Zhiji for taking care of me during my daily life. Thanks Guanyu and Hou-Cheng for help in my research and postdoc applications. I would also thank to Dr. Niu and Pam for their support in building such a harmonious department.

Finally, I want express my deepest appreciation to my parents, who always provide me their love through my entire life.

# TABLE OF CONTENTS

List of Tables . . . . .	vii
List of Figures . . . . .	viii
Abstract . . . . .	ix
<b>1 Inference for Generalized Indirect Learning with Dependency (GIDL)</b>	<b>1</b>
1.1 Introduction . . . . .	1
1.2 Generalized Indirect Dependent Learning . . . . .	4
1.3 Asymptotic Analysis for the Dependency Matrix . . . . .	5
1.3.1 Main Results . . . . .	5
1.3.2 Equivalence to Pseudo-likelihood for Ising Models . . . . .	7
1.4 Generalized Indirect Dependency Learning with Predictors . . . . .	8
1.4.1 Asymptotic Analysis for Coefficients . . . . .	11
1.5 Algorithm Design . . . . .	13
1.6 Experiments . . . . .	16
1.6.1 Simulation Studies . . . . .	16
1.6.2 Real Data . . . . .	18
1.7 Summary . . . . .	19
1.8 Outlines of Proofs . . . . .	19
1.8.1 Proof of Theorem 1.3.1 . . . . .	19
1.8.2 Proof of Theorem 1.3.2 . . . . .	22
1.8.3 Proof of Theorem 1.4.1 . . . . .	24
<b>2 Slack Empirical Likelihood (SEL)</b>	<b>25</b>
2.1 Introduction . . . . .	25
2.2 Slack Empirical Likelihood . . . . .	26
2.2.1 Statistic-driven EL . . . . .	26
2.2.2 Introducing Slack Variables . . . . .	28
2.3 Constructing SEL from General Optimizations . . . . .	32
2.3.1 SEL for $\Theta$ Estimators . . . . .	33
2.3.2 SEL for General Penalties or Constraints . . . . .	36
2.3.3 SEL for Sparse Reduced Rank Regression . . . . .	37
2.4 Asymptotic Analysis for SEL . . . . .	38
2.4.1 SEL with Affine Inequality Constraints . . . . .	38
2.5 Analysis of Large $p$ SEL . . . . .	40
2.5.1 A Limiting Result for Large $p$ SEL for Regression . . . . .	41
2.5.2 Nonasymptotic Analysis of Large $p$ SEL for a Fixed $\lambda$ . . . . .	43
2.6 Experiments . . . . .	47
2.6.1 Simulation Studies . . . . .	47
2.6.2 Real Data . . . . .	49
2.7 Summary . . . . .	50

2.8	Outlines of Proofs . . . . .	50
2.8.1	Proof of Theorem 2.4.1 . . . . .	50
2.8.2	Proof of Theorem 2.5.1 . . . . .	53
2.8.3	Proof of Theorem 2.5.2 . . . . .	56
<b>Appendix</b>		
<b>A</b>	<b>Technical Lemmas</b>	<b>57</b>
	Bibliography . . . . .	61
	Biographical Sketch . . . . .	64

# LIST OF TABLES

1.1	Converge probability of the confidence interval constructed based on the Hotelling's $t$ -Squared statistics derived from GIDL . . . . .	17
1.2	Converge probability and length of the confidence interval comparisons between squared and Tukey's loss by GIDL . . . . .	18
2.1	Comparisons between SEL and exact Gaussian inference . . . . .	48
2.2	Comparisons between SEL and truncated Gaussian inference of the lasso . . . . .	49



# LIST OF FIGURES

2.1	Percentage of the type I error made by 100 hypothesis tests on for $H_0 : \beta = \beta^\circ$ v.s. the respective significant level by large $n$ SEL. . . . .	40
2.2	Percentage of the type I error made by 100 hypothesis tests on for $H_0 : \beta = \beta^\circ$ v.s. the respective significant level by conditional inference of SEL. . . . .	44

# ABSTRACT

Dependence is one of the most important concepts in probability and statistics. To detect and measure dependency between response variables and predictors, various models have been constructed, from simple models like least square regressions to complex ones like deep neural networks. However, less literature focuses on detecting the dependency structures among responses and incorporating this structure information to facilitate other analyses. Dependency among responses can appear when the responses are multivariate or the data are observed as groups. Markov Random Field (MRF) and Generalized Estimating Equations (GEE) are proposed to detecting the dependency structures and improve the efficiency of analysis of mean effect under these circumstances. However, these methods may not perform well when data is complex, such as non-Gaussian, heavy-tailed or skewed. In Chapter 1, we consider a generalized indirect learning with dependency (GIDL) framework to detect and apply the dependent structure between responses in both multivariate and grouped dependencies, when data are non-Gaussian, heavy-tailed or skewed. We focus on statistical analysis that covers asymptotic distribution and signal selection of dependency structures and asymptotic distribution of the coefficients of predictors with the assistant of structure information.

In Chapter 2, we develop a slack empirical likelihood (SEL) inference framework that is able to handle non-Gaussian, heavy-tailed or skewed types of data. Empirical likelihood is powerful because an inference problem is transformed into an optimization one. Besides, fewer distributional assumptions are required than traditional likelihood-based inference methods. Modern statistical models often implicitly put nonsmooth constraints through regularizations on the possible solutions for parameter estimations, like the restricted cone induced by  $\ell_1$  penalized regressions. Traditional empirical likelihood can not handle such nonsmooth constraints induced by regularizations when the solutions could appear at the boundary of constrained regions because the estimating equation is not sample additive. By carefully examining the directional derivatives at the nonsmooth region and introducing slack variables such that the modified estimating equations are sample additive, we extend the empirical likelihood framework such that the nonsmooth constraints can be handled by a joint optimization on parameters and the slack variable. We show some examples that our framework works for some traditional constrained empirical likelihood problems such that the

correct asymptotic distribution can be derived. Besides, some explorations on modern statistical problems including high dimensional regression with regularizations are provided.

# CHAPTER 1

## INFERENCE FOR GENERALIZED INDIRECT LEARNING WITH DEPENDENCY (GIDL)

### 1.1 Introduction

Big data arising in modern signal processing applications call for the need of analyzing multi-response systems, such as gene regulatory networks, multi-label classification, natural language processing, and stock markets, in the areas of biology, computer vision, neuroscience, economics and others. We assume a supervised setup with  $n$  samples of  $m$  responses and  $p$  predictors:  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_m] \in \mathbb{R}^{n \times m}$ ,  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_p] \in \mathbb{R}^{n \times p}$ . A common linear model minimizes  $\|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_F^2$  over  $\mathbf{B} \in \mathbb{R}^{p \times m}$ , where  $\|\cdot\|_F$  denotes the Frobenius norm.

Instead of naively treating all responses as independent random variables, this article studies how we capture the *between-response dependency* that cannot be ignored in real applications, and incorporate such information into a given loss.

Most studies focus on the type of regularizer in the presence of a number of responses. But we are concerned about how to make a more reasonable loss to include the multivariate dependency. Most matrix estimation problems are formulated using an *additive* loss, like squared Frobenius norm or KL divergence, as in matrix completion, recommender systems and nonnegative matrix factorization. From a statistical perspective, this amounts to assuming all responses are independent, a simplified treatment, but a true risk of distorting the data and a loss of information. How to capture various dependencies for *non-Gaussian* systems is quite a pressing challenge. The conventional probabilistic Markov random field models, in addition to being computationally intractable, have shown severe limitations in including both positive and negative correlations on say discrete data [29, 30, 7]. Given the losses to measure the accuracy on each response, which may or may not be the same, if simply adding them together is crude and inaccurate, this article studies how to make a suitable aggregation of marginal losses to incorporate higher-order statistics into the criterion.

Denote the systematic (mean) component by  $\mathbf{M} = \mathbf{X}\mathbf{B}$  with  $\mathbf{m}_k$  the  $k$ th ( $1 \leq k \leq m$ ) column. As aforementioned, many multivariate methods simply utilize an additive loss to measure the discrepancy between  $\mathbf{Y}$  and  $\mathbf{M}$ :  $l_0(\mathbf{M}; \mathbf{Y}) = \sum_{k=1}^m l_k(\mathbf{m}_k, \mathbf{y}_k)$ , where  $l_k$  is the loss for the  $k$ th response (an instance is the use of a uniform loss  $\sum_{i,k} l(m_{i,k}, y_{i,k})$  like squared distance or KL divergence). Though convenient, the additive form is problematic either from the perspective of joint entropy or statistical modeling, since it implicitly assumes independence between all  $\mathbf{y}_k$ , while for large  $m$ , real life multi-response systems are almost surely dependent. For example, in multi-label classification problem, like in [25], where the authors considered a multi-label classification problem of music into emotions. Since a piece of music can belong to more than one class of emotions, the existence of correlation between different emotions is obvious. Even when the labels of emotions have already been clustered based on some psychological models [24], the labels left still enjoy some correlation: for example, a happy song is unlikely to be sad, while a relaxing song can make people feel quiet, too. Estimating with incorporation of the correlation structure can usually improve the efficiency. Another important points in multi-label classification is the mislabeled issue: since emotions are abstract and the same songs may provide inconsistent emotions to different people. The mislabeled data may cause a harmful effect to modeling if the loss functions are not properly chosen.

In longitudinal and survival analysis, it is very common to use the patients' information repeatedly collected over time (e.g., CD4 counts, blood pressures) to predict responses like severity of diseases, probability of recovery, time for death and so on. With the repeated observations over time, the correlation among values of the same response variable must be considered in the modeling to estimate efficiently. In addition, many responses have heavy-tails in distribution, for instance, the costs of health care are usually heavy-tailed [11], since people with rare disease tend to have a high cost in treatments.

There are mainly two kinds of methods to handle correlation structure in response variables when modeling the first order effect. The first one is the well-known Generalized Estimating Equations (GEE) [9]: Usually for non-Gaussian data, it is hard to model the joint distribution of the repeated observations, and GEE only assumes estimating equations for marginal distributions, avoiding the usage of joint distribution. GEE adds a working correlation matrix into the diagonal covariance matrix of the responses, and even if the working correlation matrix is misspecified, the

estimation of coefficients can still be consistent with improvement of the efficiency. However, GEE can be limited sometimes, the first problem of GEE is its non-robust to outliers, since all marginal losses for GEE are derived from exponential families. For example, by directly applying the logistic loss in the classification problem, the model could suffer a lot from the mislabeled issues. Another problem for GEE is the large sample size requirement for high dimensional estimation problems, for example, for penalized GEE proposed in [27], the number of variables  $p$  can only be in the polynomial order of the sample size  $n$ .

The other way that statisticians deal with the issue is by defining a joint Markov random field (MRF) and limiting the loss to the exponential family. This has met quite some obstacles in recent years. First, it has been realized MRF has severe limitations in modeling data dependencies. In particular, on discrete counts (prevalent in text, genomic sequencing, site-visit problems) the associated Poisson MRF can pick only negative conditional dependencies [29, 30]. Second, the partition function in MRF is often computationally intractable, and so people have to resort to various approximations of the joint likelihood [17, 6, 28]. Third, the exponential family places stringent restrictions on the marginal loss and becomes inappropriate when the data are fat-tailed, skewed, or truncated. For example, hinge loss, savage loss [12], Tukey/Hampel’s robust losses [5], and many other useful losses simply do not correspond to any probability distribution (let alone one in the exponential family).

To tackle the challenges, an intriguing and ambitious topic would be to aggregate customer-provided marginal losses in a clever way to incorporate and learn dependencies. This is different from defining an entire distribution for a large number of variables which may easily break down in applications due to limited modeling capability and data imperfections. One natural question is whether there is a unique framework to generalize the weighting strategy to losses which are not necessarily derived from probability distributions. It turns out our proposed iterative weighted algorithm can achieve this goal as long as the gradient of the loss function is Lipschitz continuous.

We study generalized indirect dependent learning framework that allows for customizing marginal losses. First, since our methods are based on marginal models, there is no explicit restrictions on how to insert the correlation structure in our framework, unlike the MRF in Poisson case. In addition, by allowing to apply robust type of losses into the marginal models, our framework can handle mislabeled or extreme data. Finally, with explicit loss functions, sparse induced penalty

can be added together, which makes our framework easily adapted to the high dimensional estimation problems.

The work is structured as follows: in Section 1.3, we provide analysis on asymptotic properties for the dependency matrix; Section 1.4 describes analysis on asymptotic properties for the coefficients of predictors; in section 1.5, we provide a computation framework for our approach.

We will use the following notations and symbols in the rest of the paper. For a matrix  $\mathbf{A} \in \mathbb{R}^{n \times m}$ , let  $\|\mathbf{A}\|_2$  and  $\|\mathbf{A}\|_F$  denote its Frobenius norm and spectral norm respectively. We denote  $\text{vec}$  as the vectorization operation, for example,  $\text{vec}(\mathbf{A})$  will convert a  $n \times m$  matrix to a  $nm$  dimensional column vector keeping all the consecutive orders of the column vectors of  $\mathbf{A}$ . Define the operator  $\text{vec}^{-1} : \mathbb{R}^{mn} \rightarrow \mathbb{R}^{m \times n}$  such that

$$\begin{aligned} \text{vec}^{-1}(\text{vec}(\mathbf{A})) &= \mathbf{A} \text{ for all } \mathbf{A} \in \mathbb{R}^{m \times n} \\ \text{vec}(\text{vec}^{-1}(\boldsymbol{\alpha})) &= \boldsymbol{\alpha} \text{ for all } \boldsymbol{\alpha} \in \mathbb{R}^{mn}. \end{aligned} \tag{1.1}$$

Let  $\otimes$  be the standard tensor product. Let  $\mathbb{S}_{++}^m$  denote the set of all positive definite matrices of dimension  $m \times m$ . The inner product of two matrices  $\mathbf{X}$  and  $\mathbf{Y}$  (of the same size) is defined as  $\langle \mathbf{X}, \mathbf{Y} \rangle = \text{Tr}(\mathbf{X}^T \mathbf{Y})$ . We use  $1_{\mathcal{A}}(t)$  to denote the indicator function, i.e.,  $1_{\mathcal{A}}(t) = 1$  if  $t \in \mathcal{A}$  and 0 otherwise. We use  $\mathbf{X}[\mathcal{I}, \mathcal{J}]$  to denote a submatrix of  $\mathbf{X}$  with rows and columns indexed by  $\mathcal{I}, \mathcal{J}$  respectively, and abbreviate  $\mathbf{X}[\mathcal{I}, 1:m]$  and  $\mathbf{X}[1:p, \mathcal{J}]$  to  $\mathbf{X}[\mathcal{I}, ]$  and  $\mathbf{X}[, \mathcal{J}]$ , respectively. When  $\mathcal{I} = \{i\}$  and  $\mathcal{J} = \{j\}$ , we also use  $X_{ij}$  to stand for  $X[\mathcal{I}, \mathcal{J}]$  for abbreviation. For a matrix  $\mathbf{B} = [\beta_1, \dots, \beta_p]^T$ , let  $\mathcal{J}(\mathbf{B}) = \{j : \beta_j \neq 0\}$  and  $J(\mathbf{B}) = |\mathcal{J}(\mathbf{B})| = \|\mathbf{B}\|_{2,0}$ .

## 1.2 Generalized Indirect Dependent Learning

Given an observation matrix  $\mathbf{Y} \in \mathbb{R}^{n \times m}$ , the target is to learn the column dependency structure  $\mathbf{W} \in \mathbb{S}_{++}^m$  when the pre-specified loss for the  $k$ th variable is  $l_k$  for  $k = 1, 2, \dots, m$ , which are not necessary the same. Let  $\mathbf{C} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{M} \in \mathbb{R}^{n \times m}$  where  $\mathbf{M} = \mathbf{1}\boldsymbol{\alpha}^T$  is the mean and already known,  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_m]^T$ . [23] considers the following generalized multivariate learning with dependency framework

$$\begin{aligned} \min_{(\mathbf{W}, \mathbf{C}) \in \mathcal{Z}} & \phi^{-1} \bar{l}(\mathbf{M} + \mathbf{C}(\mathbf{I} - \phi \mathbf{W})^{1/2}; \mathbf{Y}) + \frac{1}{2} \text{Tr}\{\mathbf{C}\mathbf{W}\mathbf{C}^T\} \\ & - \frac{n}{2} \log \det \mathbf{W} + P_{\mathbf{W}}(\mathbf{W}; \lambda_{\mathbf{W}}), \end{aligned} \tag{1.2}$$

where  $\mathcal{Z} = \mathcal{S}_{++}^m(\phi) \times \mathbb{R}^{n \times m}$  with  $\mathcal{S}_{++}^m(\phi) = \{\mathbf{W} \in \mathcal{S}_{++}^m : \mathbf{W} \preceq \mathbf{I}/\phi\}$  and  $\bar{l} = \sum_{k=1}^m l_k$ . Let  $\Theta = [\theta_1, \dots, \theta_m] = \mathbf{C}(\mathbf{I} - \phi\mathbf{W})^{1/2}$ . Assume that the  $l_k(\cdot)$  is scaled by  $l_k''(\alpha_k)$ , such that we have

$$\nabla \bar{l}(\Theta) = [\nabla l_1(\theta_1), \dots, \nabla l_m(\theta_m)] = \nabla \bar{l}(\mathbf{M}) + \mathbf{C}\mathbf{Z}^{1/2}\mathbf{D} + \Delta(\Theta), \quad (1.3)$$

with  $\mathbf{D} = \mathbf{I}$  and  $\Delta(\Theta)$  is the reminder term with respect to the expansion of  $\nabla \bar{l}(\Theta)$ .

In addition, the expansion of the loss  $\bar{l}$  can also be written as:

$$\bar{l}(\Theta) = \bar{l}(\mathbf{M}) + \langle \nabla \bar{l}(\mathbf{M}), \mathbf{C}\mathbf{Z}^{\frac{1}{2}} \rangle + \frac{1}{2}\mathbf{C}\mathbf{Z}\mathbf{C}^T + \delta(\Theta), \quad (1.4)$$

where  $\delta(\Theta)$  is the reminder term with respect to the expansion of  $\bar{l}(\Theta)$ .

Define

$$\hat{\Sigma}^n = \frac{(\nabla \bar{l}(\mathbf{M}) + \Delta^*)^T (\nabla \bar{l}(\mathbf{M}) + \Delta^*)}{n}, \quad (1.5)$$

where  $\Delta^* = \Delta(\Theta^*)$ . Let  $\Sigma^*$  be the inverse of the true dependency structure:  $\Sigma^* = \mathbf{W}^{*-1}$ . Suppose that the true dependency structure  $\mathbf{W}^*$  is sparse with nonzero index set  $\mathcal{J}_W^*$  such that  $\mathbf{W}_{\mathcal{J}_W^{*c}} = \mathbf{0}$  and  $W_{ij} \neq 0$  for all  $(i, j) \in \mathcal{J}_W^*$ .

Some remarkable results are shown in [23]: first, as shown in Theorem 1 in this paper, when the loss is the squared loss, then the high order term  $\Delta^* = \mathbf{0}$  such that the optimization problem for optimization (1.2) degenerates to Gaussian graph leaning. In addition,  $\hat{\Sigma}^n$  plays a role as the sample covariance matrix, where the higher order term  $\Delta^*$  is incorporated. Finally, unlike the GEE which often fails in high large  $p$  learning, as stated in Theorem 4 in [23], even for the case  $p \gg n$ , the above optimization can still deliver a reasonable estimating results.

## 1.3 Asymptotic Analysis for the Dependency Matrix

### 1.3.1 Main Results

In this part, we show our main theorem. To be more specific, we show that under certain regularity conditions, the signal part of the true dependency matrix can be identified, in addition, the signals follow a multivariate normal distribution, which can be directly applied into inference via Hotelling's T-squared statistics. The following conditions are necessary in our main theorem:

1. Effective noise with respect to the mean:  $\mathbb{E}(-\nabla \bar{l}(\mathbf{M})) = \mathbf{0}$ ;
2. Effective noise with respect to  $\mathbf{W}^*$ :  $\mathbb{E}\{-\hat{\Sigma}^n + \mathbf{W}^{*-1}\} = \mathbf{0}$ ;



3. The second order moment of the effective noise with respect to the true signal part of  $\mathbf{W}^*$ :

$$\mathbb{E} \left\{ (\hat{\Sigma}_{\mathcal{J}_W^*}^n - \Sigma_{\mathcal{J}_W^*}^*) \otimes (\hat{\Sigma}_{\mathcal{J}_W^*}^n - \Sigma_{\mathcal{J}_W^*}^*) \right\}$$

is positive-definite and finite;

4. Let  $\sigma_{\max}(\mathbf{W}^*)$  and  $\sigma_{\min}(\mathbf{W}^*)$  be the maximal and minimal eigenvalues of  $\mathbf{W}^*$ , then there exist positive constants  $\gamma_{\min}$  and  $\gamma_{\max}$  such that  $\gamma_{\min} < \sigma_{\min}(\mathbf{W}^*) < \sigma_{\max}(\mathbf{W}^*) < \gamma_{\max}$ ;

5. There exists an open set  $\Omega \subset \mathcal{Z}$  that contains true  $(\mathbf{W}^*, \mathbf{C}^*)$  that for almost all  $\mathbf{y}$  all  $(\mathbf{W}, \mathbf{C}) \in \Omega$ , we have  $\|\Delta(\Theta)\|_F = o_p(1)$  and  $\|\delta(\Theta)/d\mathbf{W}\|_F = o_p(n^{1/2})$ .

6. Properties of the penalty:

$$\begin{aligned} \lim_{n \rightarrow \infty} \max(P_W''(\mathbf{W}_{\mathcal{J}_W^*}^*; \lambda_W))/n &= 0; \\ \lim_{n \rightarrow \infty} \lim_{t \rightarrow 0^+} P_W'(t; \lambda_W)/\lambda_W &> 0; \\ \max(P_W'(\mathbf{W}_{\mathcal{J}_W^*}^*; \lambda_W))/n &= O_p(n^{-1/2}); \quad \text{and} \end{aligned}$$

there exist constants  $C_1, C_2$  such that when  $t_1, t_2 > C_1 \lambda_n$ , we have  $|P_W''(t_1; \lambda_n) - P_W''(t_2; \lambda_n)| \leq C_2 |t_1 - t_2|$ ;

7. Choices of the tuning parameter:

$$\lim_{n \rightarrow \infty} \lambda_W/n = 0, \quad \lim_{n \rightarrow \infty} \lambda_W/\sqrt{n} \rightarrow \infty,$$

Then we have the following theorem:

**Theorem 1.3.1.** *Suppose the above conditions 1–7 provided are satisfied, let  $a_n = \max(P_W'(\mathbf{W}_{\mathcal{J}_W^*}^*; \lambda_W)/n)$ , with probability tending to 1, we have:*

1. *There is a local minimizer  $\widehat{\mathbf{W}}$  such that  $\|\widehat{\mathbf{W}} - \mathbf{W}^*\|_F = O(n^{-1/2} + a_n)$ ;*
2. *The above defined local minimizer  $\widehat{\mathbf{W}}$  satisfies  $\widehat{\mathbf{W}}_{\mathcal{J}_W^{*c}} = \mathbf{0}$ ;*
3. *Asymptotic normality:  $\sqrt{n}(\text{vec } \widehat{\mathbf{W}}_{\mathcal{J}_W^*} - \text{vec } \mathbf{W}_{\mathcal{J}_W^*}^*)$  converges to a normal distribution with mean zero and covariance matrix:*

$$\{\mathbf{W}_{\mathcal{J}_W^*}^* \otimes \mathbf{W}_{\mathcal{J}_W^*}^*\} \lim_{n \rightarrow \infty} \text{Cov}(\text{vec}[-\sqrt{n}\{\hat{\Sigma}_{\mathcal{J}_W^*}^n - \Sigma_{\mathcal{J}_W^*}^*\}])\{\mathbf{W}_{\mathcal{J}_W^*}^* \otimes \mathbf{W}_{\mathcal{J}_W^*}^*\}, \quad (1.6)$$

where recall that we define  $\hat{\Sigma}^n = (\nabla \bar{l}(\mathbf{M}) + \Delta^*)^T (\nabla \bar{l}(\mathbf{M}) + \Delta^*)/n$  and  $\Sigma^* = \mathbf{W}^{*-1}$ .

### 1.3.2 Equivalence to Pseudo-likelihood for Ising Models

In this part, we show the equivalence of selection between pseudo-likelihood model and GIDL in the statistical sense when the sample size  $n \rightarrow \infty$ . Consider the Ising model:

$$p(\mathbf{y}, \mathbf{W}) = \exp \left[ \sum_{(s,t) \in E} W_{st} y_s y_t - \Psi(\mathbf{W}) \right], \quad (1.7)$$

where  $\mathbf{y} = (y_1, \dots, y_m)^T \in \{-1, 1\}^m$ , the vertices  $\mathbf{V} = \{1, 2, \dots, m\}$ , the edges  $E$  captures the conditional independence between variables of  $\mathbf{y}$ , and  $\Psi(\mathbf{W}) = \log \sum_{\mathbf{y}_s, \mathbf{y}_t \in \{-1, 1\}} \exp\{\sum_{(s,t) \in E} W_{st} y_s y_t\}$  is the normalizing constant. Due to the infeasible computation of the normalizing constant. The pseudo-likelihood approach is often considered [6], which is the direct aggregate of all conditional likelihoods. By conditional likelihood of  $y_s$  given all others and normalizing, given the index  $s$ , let  $\{-s\}$  be the index set of all other response variables except of  $s$ , we have:

$$p(y_s, \mathbf{W} | \mathbf{y}_{\{-s\}}) = \frac{\exp(\sum_{t \neq s} W_{st} y_s y_t)}{\exp(\sum_{t \neq s} W_{st} y_t) + \exp(-\sum_{t \neq s} W_{st} y_t)}. \quad (1.8)$$

Then the log-pseudo-likelihood is defined by

$$\begin{aligned} l(\mathbf{W} | \mathbf{y}) &= \sum_{s=1}^m \log p(y_s, \mathbf{W} | \mathbf{y}_{\{-s\}}) \\ &= \sum_{s=1}^m \left[ \sum_{t \neq s} y_s y_t W_{st} - \log \cosh(\sum_{t \neq s} W_{st} y_t) - \log 2 \right]. \end{aligned} \quad (1.9)$$

Then for data matrix  $\mathbf{Y} \in \{-1, 1\}^{n \times m}$ , we have the log-pseudo-likelihood:

$$l(\mathbf{W} | \mathbf{Y}) = \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^m \left[ \sum_{t \neq s} Y_{is} Y_{it} W_{st} - \log \cosh(\sum_{t \neq s} W_{st} Y_{it}) - \log 2 \right]. \quad (1.10)$$

To minimize  $-l(\mathbf{W} | \mathbf{Y})$ , the first order derivative can be calculated:

$$-\frac{\partial l}{\partial W_{sk}} = -\frac{1}{n} \sum_{i=1}^n \left[ Y_{is} Y_{ik} - Y_{ik} \tanh(\sum_{t \neq s} W_{st} Y_{it}) \right]. \quad (1.11)$$

Note that

$$\tanh(\sum_{t \neq s} W_{st} Y_{it}) = \frac{\exp(\sum_{t \neq s} W_{st} Y_{it}) - \exp(-\sum_{t \neq s} W_{st} Y_{it})}{\exp(\sum_{t \neq s} W_{st} Y_{it}) + \exp(-\sum_{t \neq s} W_{st} Y_{it})} = \mathbb{E}(Y_{is} | \{Y_{it}, t \neq s\}) \quad (1.12)$$

is the conditional expectation. Consistency of pseudo-likelihood when  $n \rightarrow \infty$  can be shown by the fact that  $Y_{ik}$  is uncorrelated with  $Y_{is} - \mathbb{E}(Y_{is} | \{Y_{it}, t \neq s\})$  and the identity  $\mathbb{E}(Y_{is} - \mathbb{E}(Y_{is} | \{Y_{it}, t \neq s\})) = 0$ .

**Theorem 1.3.2.** *Suppose the data matrix  $\mathbf{Y} \in \{-1, 1\}^{n \times m}$  is  $n$  independent observations generated from density  $P(\mathbf{y}, \mathbf{W}^*)$  defined in equation (1.7). Let  $\mathbf{W}^{pse}$  be the estimator through pseudo-likelihood approach. If the conditions in Theorem 1.3.1 are satisfied. Then as  $n \rightarrow \infty$ , if for any index  $s, k$ , the following condition holds*

$$\frac{1}{n} \sum_{i=1}^n (Y_{ik} \mathbf{Y}_{i\{-s\}} (\mathbf{Y}_{\{-s\}}^T \mathbf{Y}_{\{-s\}})^{-1} \mathbf{Y}_{\{-s\}}^T \mathbf{Y}_s)^3 \rightarrow 0. \quad (1.13)$$

*Then there exists a stationary point for our GIDL method, denoted as  $\mathbf{W}^{gidl}$ , such that*

$$P(\text{nzoff}(\mathbf{W}^{gidl}) = \text{nzoff}(\mathbf{W}^{pse})) \rightarrow 1, \quad (1.14)$$

*where  $\text{nzoff}(\mathbf{W})$  is the index set for the non-zero components on the off-diagonal components of  $\mathbf{W}$ .*

**Remark 1.3.1.** *The regularity condition (1.13) puts a high order constraint on the graph structure of the Ising model. Similar constraints on graph structure are also required in other literature, for example in [10], the graph with singleton separator sets is required.*

## 1.4 Generalized Indirect Dependency Learning with Predictors

Although we formulate the problem by introducing the dependency between columns of response matrix, we show that through converting the problem to the vectorized case. The clustered data type of problem can also be handled. Let us first illustrate our motivation from multivariate Gaussian case: for multiple responses case, let  $\mathbf{y} = \text{vec}(\mathbf{Y})$ ,  $\mathbf{L} = \mathbf{I}_{m \times m} \otimes \mathbf{X}$  and  $\boldsymbol{\beta} = \text{vec}(\mathbf{B})$ , then the multivariate Gaussian distribution can be represented as  $\mathbf{y} \sim \mathcal{N}(\mathbf{L}\boldsymbol{\beta}, \mathbf{S})$ , where  $\mathbf{S} \in \mathbb{S}_{++}^{mn}$  captures the covariance structure of the column vectors of  $\mathbf{Y}$ . For example, for multiple response problems,  $\mathbf{S} = \boldsymbol{\Sigma} \otimes \mathbf{I}_{n \times n}$ , where  $\boldsymbol{\Sigma}$  is the covariance structure between response variables. Now suppose the loss is quadratic for  $\mathbf{y} - \mathbf{L}\boldsymbol{\beta}$ , then we have the estimating equation  $\mathbf{L}^T(\mathbf{L}\boldsymbol{\beta} - \mathbf{y}) = \mathbf{0}$ . In the Gaussian case we can directly add the weighting matrix into the estimation equations:

$$\mathbf{L}^T \mathbf{S}^{-1}(\mathbf{L}\boldsymbol{\beta} - \mathbf{y}) = \mathbf{0}. \quad (1.15)$$

Note that when  $\mathbf{S} = \boldsymbol{\Sigma} \otimes \mathbf{I}_{n \times n}$ ,

$$\begin{aligned} \mathbf{L}^T \mathbf{S}^{-1}(\mathbf{L}\boldsymbol{\beta} - \mathbf{y}) &= (\mathbf{I} \otimes \mathbf{X}^T)(\boldsymbol{\Sigma}^{-1} \otimes \mathbf{I})(\mathbf{L}\boldsymbol{\beta} - \mathbf{y}) \\ &= (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{X}^T)(\mathbf{L}\boldsymbol{\beta} - \mathbf{y}) \\ &= (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{X}^T) \text{vec}(\mathbf{X}\mathbf{B} - \mathbf{Y}) \\ &= \text{vec}(\mathbf{X}^T(\mathbf{X}\mathbf{B} - \mathbf{Y})\boldsymbol{\Sigma}^{-1}) = \mathbf{0}, \end{aligned} \quad (1.16)$$

which is reduced to the ordinary estimating equation for multivariate Gaussian case.

For clustered data case,  $\mathbf{Y} \in \mathbb{R}^{n \times m}$ , where  $n$  is the number of observation and  $m$  is the number of groups.  $\mathbf{X} \in \mathbb{R}^{n \times p}$  is the predictor matrix and  $\beta_0 \in \mathbb{R}^p$  is the coefficient. Now the coefficient is just a vector because there are only measurement errors within each cluster and sharing the same treatment effect  $\beta_0$ . Following the same idea, let  $\mathbf{W} \in \mathbb{R}^{n \times n}$  is the dependency matrix for each cluster. By vectorization, let  $\mathbf{L} = \mathbf{I}_{m \times m} \otimes \mathbf{X}$  and  $\beta = \mathbf{I}_{m \times 1} \otimes \beta_0$ , then the distribution can also be specified if likelihood is normal:  $\mathbf{y} \sim \mathcal{N}(\mathbf{L}\beta, \mathbf{S})$ , where  $\mathbf{S} = (\mathbf{I}_{m \times m} \otimes \mathbf{W})^{-1}$ . That is to say, the differences between clustered data case and multiple response case is that for clustered data, the dependency is introduced among rows rather than columns. In addition, if we formulate the coefficient  $\beta$  into a  $\mathbb{R}^{p \times m}$  matrix, then all columns should be the same. For convenience, we show all the following results in the vectorized multiple responses case, but similar analysis technique can be applied on clustered data.

Now if we have a response data matrix  $\mathbf{Y}$  which is not multivariate Gaussian, and the loss  $l$  we want to evaluate is also not  $\ell_2$  type, even not necessarily derived from a probability distribution. Like in traditional weighted least squares, we prespecify an estimated inverse covariance matrix  $\mathbf{W} \in \mathbb{S}_{++}^m$  of  $\mathbf{Y}$ , which captures the conditional independence between the response variables ( $\mathbf{S}^{-1}$  in multivariate Gaussian case). In order to incorporate the covariance structure into estimation beyond the Gaussian case, we consider an auxiliary vector  $\mathbf{c}$  as random effects, with a design matrix  $(\mathbf{I} - \phi\mathbf{W} \otimes \mathbf{I})^{1/2}$ . Given the estimated  $\mathbf{W}$ , let  $\boldsymbol{\theta} = \mathbf{L}\beta + (\mathbf{I} - \phi\mathbf{W} \otimes \mathbf{I})^{1/2}\mathbf{c}$ ,  $\boldsymbol{\Theta} = \text{vec}^{-1}(\boldsymbol{\theta})$ , and  $\mathbf{C} = \text{vec}^{-1}(\mathbf{c})$ , we first consider the multi-dependency learning framework to without regularizations on the estimated coefficients of predictors:

$$\min_{\beta, \mathbf{c}, \mathbf{W}} \phi^{-1}l(\boldsymbol{\theta}; \mathbf{y}) + \frac{1}{2}\mathbf{c}^T(\mathbf{W} \otimes \mathbf{I})\mathbf{c} - \frac{n}{2} \log \det \mathbf{W} + P_W(\mathbf{W}; \lambda_W). \quad (1.17)$$

Denote  $F(\beta, \mathbf{W}, \mathbf{c})$  as the objective function, by taking  $\nabla_{\mathbf{c}}F = \mathbf{0}$ ,  $\nabla_{\beta}F = \mathbf{0}$ , we can obtain the estimation equation:

$$\phi^{-1}(\mathbf{I} - \phi\mathbf{W} \otimes \mathbf{I})^{1/2}\nabla l(\boldsymbol{\theta}) + (\mathbf{W} \otimes \mathbf{I})\mathbf{c} = \mathbf{0}, \quad (1.18)$$

$$\phi^{-1}\mathbf{L}^T\nabla l(\boldsymbol{\theta}) = \mathbf{0}, \quad (1.19)$$

For the Gaussian case, our GIDL framework can help to estimate  $\beta$  with the structure information  $\mathbf{W}$ . In fact, when  $l$  in equation (1.17) is the quadratic loss, plugging the loss into the gradient

equation (1.18), we have

$$\phi^{-1}(\mathbf{I} - \phi \mathbf{W} \otimes \mathbf{I})^{1/2}(\mathbf{L}\boldsymbol{\beta} - \mathbf{y} + (\mathbf{I} - \phi \mathbf{W} \otimes \mathbf{I})^{1/2}\mathbf{c}) + (\mathbf{W} \otimes \mathbf{I})\mathbf{c} = \mathbf{0},$$

where  $\mathbf{c} = (\mathbf{I} - \phi \mathbf{W} \otimes \mathbf{I})^{1/2}(\mathbf{y} - \mathbf{L}\boldsymbol{\beta})$  can be obtained. Plugging the form of  $\mathbf{c}$  back into the optimization problem,

$$\begin{aligned} & \min_{\boldsymbol{\beta}, \mathbf{c}, \mathbf{W}} \phi^{-1}l(\boldsymbol{\theta}; \mathbf{y}) + \frac{1}{2}\mathbf{c}^T(\mathbf{W} \otimes \mathbf{I})\mathbf{c} + P_{\mathbf{W}}(\mathbf{W}; \lambda_{\mathbf{W}}) \\ &= \min_{\boldsymbol{\beta}, \mathbf{W}} \phi^{-1}\|(\phi \mathbf{W} \otimes \mathbf{I})(\mathbf{y} - \mathbf{L}\boldsymbol{\beta})\|_2^2 + \frac{1}{2}(\mathbf{y} - \mathbf{L}\boldsymbol{\beta})^\top (\mathbf{I} - \phi \mathbf{W} \otimes \mathbf{I})(\mathbf{W} \otimes \mathbf{I})(\mathbf{y} - \mathbf{L}\boldsymbol{\beta}) + P_{\mathbf{W}}(\mathbf{W}; \lambda_{\mathbf{W}}) \\ &= \min_{\boldsymbol{\beta}, \mathbf{W}} \frac{1}{2}\{(\mathbf{y} - \mathbf{L}\boldsymbol{\beta})^\top (\mathbf{W} \otimes \mathbf{I}^{\frac{1}{2}})(\mathbf{W} \otimes \mathbf{I}^{\frac{1}{2}})(\mathbf{y} - \mathbf{L}\boldsymbol{\beta})\} + P_{\mathbf{W}}(\mathbf{W}; \lambda_{\mathbf{W}}) \\ &\stackrel{(i)}{=} \min_{\mathbf{B}, \mathbf{W}} \frac{1}{2} \text{vec}((\mathbf{Y} - \mathbf{X}\mathbf{B})\mathbf{W}^{\frac{1}{2}})^\top \text{vec}((\mathbf{Y} - \mathbf{X}\mathbf{B})\mathbf{W}^{\frac{1}{2}}) + P_{\mathbf{W}}(\mathbf{W}; \lambda_{\mathbf{W}}) \\ &= \min_{\mathbf{B}, \mathbf{W}} \text{Tr}\{(\mathbf{Y} - \mathbf{X}\mathbf{B})\mathbf{W}(\mathbf{Y} - \mathbf{X}\mathbf{B})^\top\} + P_{\mathbf{W}}(\mathbf{W}; \lambda_{\mathbf{W}}), \end{aligned}$$

where in the first equality we use the fact that  $\mathbf{I} - \phi \mathbf{W} \otimes \mathbf{I}$  and  $\mathbf{W} \otimes \mathbf{I}$  are exchangeable, and in (i) we apply the identity  $\text{vec}((\mathbf{Y} - \mathbf{X}\mathbf{B})\mathbf{W}^{1/2}) = (\mathbf{W}^{1/2} \otimes \mathbf{I}) \text{vec}(\mathbf{Y} - \mathbf{X}\mathbf{B})$ . Then the optimization problem in equation (1.17) is equivalent with

$$\underset{\mathbf{B}, \mathbf{W}}{\text{argmin}} \frac{1}{2}\{(\mathbf{Y} - \mathbf{X}\mathbf{B})\mathbf{W}(\mathbf{Y} - \mathbf{X}\mathbf{B})^\top\} - \frac{n}{2} \log \det \mathbf{W} + P_{\mathbf{W}}(\mathbf{W}; \lambda_{\mathbf{W}}). \quad (1.20)$$

**Remark 1.4.1.** We can see the advantages of applying the GIDL framework (1.17) to non-Gaussian losses:

1. The loss function  $l$  can be customized directly while the structure information still plays a role in estimating the first order effect;
2. With the explicit appearance of loss functions, we can directly add regularizations on the loss to deal with high dimensional estimation problems, which is more nature comparing with penalization on the estimating equations (e.g. [27]).

For the non-Gaussian cases, we expand the  $\nabla l(\boldsymbol{\theta})$  to the second order term:  $\boldsymbol{\delta}(\boldsymbol{\beta}; \mathbf{W}; \mathbf{c}) = \nabla l(\boldsymbol{\theta}) - \nabla l(\mathbf{L}\boldsymbol{\beta}) - \nabla^2 l(\mathbf{L}\boldsymbol{\beta})(\boldsymbol{\theta} - \mathbf{L}\boldsymbol{\beta})$ . Denote  $\mathbf{Z} = \mathbf{I} - \phi \mathbf{W} \otimes \mathbf{I}$ ,  $\mathbf{D}(\boldsymbol{\beta}) = \nabla^2 l(\mathbf{L}\boldsymbol{\beta})$ , we now try to get rid of  $\mathbf{c}$  to obtain the estimating equations only for  $\boldsymbol{\beta}$  given  $\mathbf{W}$ . Based on the equation (1.18), we have:

$$\nabla l(\mathbf{L}\boldsymbol{\beta}) + \boldsymbol{\delta}(\boldsymbol{\beta}; \mathbf{W}; \mathbf{c}) + \mathbf{D}(\boldsymbol{\beta})\mathbf{Z}^{1/2}\mathbf{c} = -\phi \mathbf{Z}^{-1/2}(\mathbf{W} \otimes \mathbf{I})\mathbf{c}, \quad (1.21)$$

then rearranging the terms, we can obtain:

$$\begin{aligned}
\mathbf{Z}^{-1/2}\mathbf{c} &= -\mathbf{Z}^{-1/2}(\mathbf{Z}^{1/2}\mathbf{D}(\boldsymbol{\beta})\mathbf{Z}^{1/2} + \phi\mathbf{W} \otimes \mathbf{I})^{-1}\mathbf{Z}^{1/2}(\nabla l(\mathbf{L}\boldsymbol{\beta}) + \boldsymbol{\delta}(\boldsymbol{\beta}; \mathbf{W}; \mathbf{c})) \\
&= -[\mathbf{D}(\boldsymbol{\beta})\mathbf{Z} + \phi\mathbf{Z}^{-1/2}(\mathbf{W} \otimes \mathbf{I})\mathbf{Z}^{1/2}]^{-1}(\nabla l(\mathbf{L}\boldsymbol{\beta}) + \boldsymbol{\delta}(\boldsymbol{\beta}; \mathbf{W}; \mathbf{c})) \\
&= -[\mathbf{D}(\boldsymbol{\beta})(\mathbf{I} - \phi\mathbf{W} \otimes \mathbf{I}) + \phi\mathbf{W} \otimes \mathbf{I}]^{-1}(\nabla l(\mathbf{L}\boldsymbol{\beta}) + \boldsymbol{\delta}(\boldsymbol{\beta}; \mathbf{W}; \mathbf{c})),
\end{aligned} \tag{1.22}$$

where in the second equality we use the fact that  $\mathbf{Z}$  and  $\mathbf{W} \otimes \mathbf{I}$  are exchangeable. Plugging the obtained identity to the equation (2.40), we have:

$$\mathbf{L}^T(\mathbf{W} \otimes \mathbf{I})[\mathbf{D}(\boldsymbol{\beta})(\mathbf{I} - \phi\mathbf{W} \otimes \mathbf{I}) + \phi\mathbf{W} \otimes \mathbf{I}]^{-1}(\nabla l(\mathbf{L}\boldsymbol{\beta}) + \boldsymbol{\delta}(\boldsymbol{\beta}; \mathbf{W}; \mathbf{c})) = \mathbf{0}, \tag{1.23}$$

Now we incorporate equation (1.18) into the above estimating equation. With a little abuse of notation, denote

$$\boldsymbol{\delta}(\boldsymbol{\beta}; \mathbf{W}) = \boldsymbol{\delta}(\boldsymbol{\beta}; \mathbf{W}; \mathbf{c}^*(\boldsymbol{\beta}; \mathbf{W})), \tag{1.24}$$

where  $\mathbf{c}^*(\boldsymbol{\beta})$  is a root of  $\mathbf{c}$  for equation (1.18) given  $\boldsymbol{\beta}$  and  $\mathbf{W}$ . Therefore, the final estimating equation becomes:

$$\mathbf{L}^T(\mathbf{W} \otimes \mathbf{I})[\mathbf{D}(\boldsymbol{\beta}) + \phi(\mathbf{W} \otimes \mathbf{I})(\mathbf{I} - \mathbf{D}(\boldsymbol{\beta}))]^{-1}(\nabla l(\mathbf{L}\boldsymbol{\beta}) + \boldsymbol{\delta}(\boldsymbol{\beta}; \mathbf{W})) = \mathbf{0}. \tag{1.25}$$

Notice that when  $l$  is quadratic loss,  $\mathbf{D}(\boldsymbol{\beta}) = \mathbf{I}$ ,  $\boldsymbol{\delta} = \mathbf{0}$ , then the estimating equation reduces to weighted least squares given  $\mathbf{W}$  as shown in equation (1.16). When  $\boldsymbol{\delta} = \mathbf{0}$  and  $\phi = 0$ , it seems we insert a covariance matrix into the traditional GLM estimating equations, just like GEE.

### 1.4.1 Asymptotic Analysis for Coefficients

Consider the case when  $n \rightarrow \infty$  while  $p, m$  are fixed. Define the loss function as

$$\bar{l}(\mathbf{L}\boldsymbol{\beta}, \mathbf{y}) = \sum_{k=1}^m l_k(\mathbf{X}\mathbf{B}[k], \mathbf{Y}[k]) = \sum_{i=1}^n \sum_{k=1}^m l_k(\mathbf{X}[i,] \mathbf{B}[k], \mathbf{Y}[i, k]) \tag{1.26}$$

$$= \sum_{i=1}^n \bar{l}(\mathbf{L}[i,] \boldsymbol{\beta}, \mathbf{y}[i,]), \tag{1.27}$$

where  $m(i) = \{i + jn\}_{j=0,1,\dots,m-1}$  and  $l_k$  is the prespecified loss for the  $k$ th variable of  $\mathbf{Y}$ . Define

$$\bar{l}_i(\mathbf{L}\boldsymbol{\beta}, \mathbf{y}) = \bar{l}(\mathbf{L}[i,] \boldsymbol{\beta}, \mathbf{y}[i,]), \tag{1.28}$$

then for  $i = 1, 2, \dots, n$ ,  $\bar{l}_i$  are independent and follow the same distribution, since they are the losses for each independent observations, with each observation including all response variables. Then we have the asymptotic property of the estimator  $\hat{\boldsymbol{\beta}}$ :

**Theorem 1.4.1.** *Suppose the all loss functions  $l_k$ ,  $k = 1, \dots, m$  have bounded second order derivatives and denote  $\beta^*$  as the true vectorized coefficients. Assume*

1. *Given  $\beta, \mathbf{W}$ , the solution for axillary variable  $\mathbf{c}$  exists for equation (1.18);*
2. *Given  $\beta$ , the dependency matrix  $\mathbf{W}$  can be estimated  $O_p(1/\sqrt{n})$  consistently converging to a constant positive-definite matrix  $\overline{\mathbf{W}}$ , which can be different from the truth;*
3. *The high order term  $\delta(\beta^*; \mathbf{W})$  satisfies  $\sqrt{n}\delta(\beta^*; \mathbf{W}) \rightarrow \mathbf{0}$  and  $\|\frac{\partial \delta}{\partial \beta}(\beta^*; \delta(\beta^*, \mathbf{W}))\|_2/n \rightarrow 0$  for any  $\mathbf{W}$  if  $n \rightarrow \infty$ .*

Then let  $\hat{\beta}$  be the solution of the optimization problem

$$\min_{\beta, \mathbf{W}, \mathbf{c}} \phi^{-1} \bar{l}(\mathbf{L}\beta + (\mathbf{I} - \phi \mathbf{W} \otimes \mathbf{I})^{1/2} \mathbf{c}; \mathbf{y}) + \frac{1}{2} \mathbf{c}^T (\mathbf{W} \otimes \mathbf{I}) \mathbf{c} + P_{\mathbf{W}}(\mathbf{W}; \lambda_{\mathbf{W}}), \quad (1.29)$$

where  $\bar{l}$  are defined in equation (1.26),  $\sqrt{n}(\hat{\beta} - \beta^*)$  is asymptotically multivariate Gaussian with zero mean and covariance matrix:

$$\begin{aligned} \mathbf{V} = & \lim_{n \rightarrow \infty} n [\mathbf{L}^T \overline{\mathbf{W}}_T [\mathbf{D}(\beta^*) + \phi \overline{\mathbf{W}}_T (\mathbf{I} - \mathbf{D}(\beta^*))]^{-1} \mathbf{D}(\beta^*) \mathbf{L}]^{-1} \\ & [\mathbf{L}^T \overline{\mathbf{W}}_T [\mathbf{D}(\beta^*) + \phi \overline{\mathbf{W}}_T (\mathbf{I} - \mathbf{D}(\beta^*))]^{-1} \text{Cov}(\mathbf{y}) [\mathbf{D}(\beta^*) + \phi \overline{\mathbf{W}}_T (\mathbf{I} - \mathbf{D}(\beta^*))]^{-1} \overline{\mathbf{W}}_T \mathbf{L}] \\ & [\mathbf{L}^T \overline{\mathbf{W}}_T [\mathbf{D}(\beta^*) + \phi \overline{\mathbf{W}}_T (\mathbf{I} - \mathbf{D}(\beta^*))]^{-1} \mathbf{D}(\beta^*) \mathbf{L}]^{-1}, \end{aligned}$$

where  $\overline{\mathbf{W}}_T = \overline{\mathbf{W}} \otimes \mathbf{I}$ .

**Remark 1.4.2.** *For the Gaussian case,  $\mathbf{D}(\beta^*) = \mathbf{I}$ , then the covariance matrix is reduced to  $\mathbf{V} = \lim_{n \rightarrow \infty} n [\mathbf{L}^T \overline{\mathbf{W}}_T \mathbf{L}]^{-1} [\mathbf{L}^T \overline{\mathbf{W}}_T \text{Cov}(\mathbf{y}) \overline{\mathbf{W}}_T \mathbf{L}] [\mathbf{L}^T \overline{\mathbf{W}}_T \mathbf{L}]^{-1}$ . In addition, when  $\overline{\mathbf{W}}$  is correctly specified as the truth of inverse covariance matrix, then  $\mathbf{V} = \lim_{n \rightarrow \infty} n [\mathbf{L}^T \overline{\mathbf{W}}_T \mathbf{L}]^{-1}$ .*

**Remark 1.4.3.** *The conditions required in the above theorem are similar with clustered-specific model for GEE where there is also a random effect in GEE [31]. Both the random effect  $\mathbf{c}$  and the high order terms  $\delta$  need to be regularized properly. However, for GEE with a random effect, the distribution of the random effect needs to be known as a priori, while for our model, the random effect  $\mathbf{c}$  is defined in an optimization way. Our model is computational more efficient since optimization is usually more feasible than integral over a distribution.*

**Remark 1.4.4.** *When  $\phi = 0$ , the asymptotic covariance matrix has the similar form with GEE. However, there is a main difference between our estimator and GEE: for finite sample estimation,*

our method incorporates the high order information  $\delta(\beta; \mathbf{W})$  into estimation. This could be more desirable just like in literature about bootstrap, bootstrapping the entire estimator could produce better results than only the linearized part of the estimator (e.g. [3]).

## 1.5 Algorithm Design

Construct a surrogate function by linearizing the first term  $\phi^{-1}l_0$  only at  $\Theta = \mathbf{M} + \mathbf{C}(\mathbf{I} - \phi\mathbf{W})^{1/2}$  as a whole:

$$g(\mathbf{M}, \mathbf{W}, \mathbf{C}; \mathbf{M}^{[k]}, \mathbf{W}^{[k]}, \mathbf{C}^{[k]}) = \phi^{-1}l_0(\Theta^{[k]}; \mathbf{Y}) + \phi^{-1}\langle \nabla l_0(\Theta^{[k]}), \Theta - \Theta^{[k]} \rangle + \frac{\rho}{2\phi} \|\Theta - \Theta^{[k]}\|_F^2 + \frac{1}{2}Tr\{\mathbf{C}\mathbf{W}\mathbf{C}^T\} - \frac{n}{2} \log \det \mathbf{W} + P_M(\mathbf{M}) + P_W(\mathbf{W})$$

Define the next iterate as  $(\mathbf{M}^{[k+1]}, \mathbf{W}^{[k+1]}, \mathbf{C}^{[k+1]}) = \operatorname{argmin}_{\mathbf{M} \in \mathcal{X}, (\mathbf{W}, \mathbf{C}) \in \mathcal{Z}} g(\mathbf{M}, \mathbf{W}, \mathbf{C}; \mathbf{M}^{[k]}, \mathbf{W}^{[k]}, \mathbf{C}^{[k]})$ . Assume  $\nabla l_0$  is  $L$ -Lipschitz continuous. Then, as long as  $\rho \geq L$ , the sequence of iterates guarantees the objective function values to be non-increasing and so convergent. Rewrite the  $g$ -optimization

$$\min_{\mathbf{M}, \mathbf{W}, \mathbf{C}} \frac{\rho}{2\phi} \|\mathbf{M} + \mathbf{C}(\mathbf{I} - \phi\mathbf{W})^{1/2} - \Xi^{[k+1]}\|_F^2 + \frac{1}{2}Tr\{\mathbf{C}\mathbf{W}\mathbf{C}^T\} - \frac{n}{2} \log \det(\mathbf{W}) + P_M(\mathbf{M}) + P_W(\mathbf{W}),$$

with  $\Xi^{[k+1]} = \Theta^{[k]} - \nabla_{\Theta} l_0(\Theta^{[k]})/\rho$ ,

and  $1/\rho$  amounts to the step size. The problem can be further simplified—we can scale  $l_0$  by  $L$  beforehand and set  $\rho = 1$ . Then we are back to the plain multivariate Gaussian estimation.

Another appealing fact is that the auxiliary matrix  $\mathbf{C}$  does *not* have to be explicitly computed at all. Forming  $\Xi^{[k]}$  avoids the need of computing SVD or matrix square-roots. Concretely, the key quantity  $\Theta^{[k]}$  can be conveniently written as a weighted average of  $\mathbf{M}$  and  $\Xi^{[k]}$ :

$$\Theta^{[k]} = \mathbf{M} + \mathbf{C}^{[k]}(\mathbf{I} - \phi\mathbf{W}^{[k]})^{1/2} = \mathbf{M} + (\Xi^{[k]} - \mathbf{M})(\mathbf{I} - \phi\mathbf{W}^{[k]})^{1/2}(\mathbf{I} - \phi\mathbf{W}^{[k]})^{1/2} = \mathbf{M}(\mathbf{W}^{[k]}\phi) + \Xi^{[k]}(\mathbf{I} - \phi\mathbf{W}^{[k]}).$$

Hence if the scaled  $l_0$  is 1-strongly smooth, the main loop of the algorithm consists of only 3 steps:

- 1)  $\Xi^{[k]} \leftarrow \Theta^{[k-1]} - \nabla_{\Theta} l_0(\Theta^{[k-1]})$ ;
- 2)  $(\mathbf{M}^{[k]}, \mathbf{W}^{[k]}) \leftarrow \operatorname{argmin}_{\mathbf{M}, \mathbf{W}} \frac{1}{2}Tr\{(\Xi^{[k]} - \mathbf{M})\mathbf{W}(\Xi^{[k]} - \mathbf{M})^T\} - \frac{n}{2} \log \det(\mathbf{W}) + P_M(\mathbf{M}) + P_W(\mathbf{W})$
- 3)  $\Theta^{[k]} \leftarrow \Xi^{[k]} + \phi(\mathbf{M}^{[k]} - \Xi^{[k]})\mathbf{W}^{[k]}$



The optimization algorithm has great implementation ease: 1) and 3) involve some basic matrix operations to form the key matrices  $\Xi$  and  $\Theta$ , and 2) just carries out (ordinary) Gaussian learning. Correspondingly, GIDL calls multivariate Gaussian estimation *iteratively* to harness a general loss driven dependent learning with guaranteed convergence. (Special case: for  $l(\theta, y) = (\theta - y)^2/2$ ,  $\Xi^{[k]} = \Theta^{[k-1]} - (\Theta^{[k-1]} - \mathbf{Y}) = \mathbf{Y}$  stays fixed, indicating no need to iterate.) The iterative/indirect dependent learning mechanism is innovative to the best of our knowledge, and relies on quadratic distance based learning—any such optimization algorithms can be seamlessly applied in Step 2).

Under high dimensional settings with predictors where  $p$  is much larger than  $n$ . It is natural and important to assume joint row sparsity on the true coefficient matrix  $\mathbf{B}^*$  and add penalization on the objective function, which means that all the observations for some variables are zero simultaneously. Let  $\mathbf{B} = [\beta_1, \dots, \beta_p]^\top$  and for simplicity denote  $P_g(\beta; \lambda) := \sum_{i=1}^p P(\|\beta_i\|_2; \lambda)$ , where  $P$  is a sparse inducing penalty. Then  $P_g(\beta; \lambda)$  can be seen as a group penalty for the vectorized coefficient  $\beta$ . We consider the following optimization problem:

$$\min_{\beta, \mathbf{c}, \mathbf{W}} \phi^{-1} \bar{l}(\mathbf{L}\beta + (\mathbf{I} - \phi \mathbf{W} \otimes \mathbf{I})^{1/2} \mathbf{c}; \mathbf{y}) + \frac{1}{2} \mathbf{c}^\top (\mathbf{W} \otimes \mathbf{I}) \mathbf{c} - \frac{n}{2} \log \det \mathbf{W} + P_g(\beta; \lambda) + P_W(\mathbf{W}; \lambda_W) \quad (1.30)$$

With a bit abuse of notation, we denote the objective function as  $F(\beta, \mathbf{W}, \mathbf{c})$ , which can be non-convex. It is not possible to obtain global optimal through computation. Instead, we provide a computational framework and establish the statistical analysis on the fixed points obtained through the algorithm we designed in the next subsection. Following the spirits of traditional iterative computing methods like weighted least squares, we first convert the optimization problem into a Gaussian one through linearization and then apply a weighted type of algorithm to obtain the fixed-point solutions. Given the  $t$ th iterate  $(\beta^{[k]}, \mathbf{W}^{[k]}, \mathbf{c}^{[k]})$ , let  $\theta^{[k]} = \mathbf{L}\beta^{[k]} + (\mathbf{I} - \phi \mathbf{W}^{[k]} \otimes \mathbf{I})^\top \mathbf{c}^{[k]}$ , we construct the following surrogate function

$$g(\beta, \mathbf{W}, \mathbf{c}; \beta^{[k]}, \mathbf{W}^{[k]}, \mathbf{c}^{[k]}) = \phi^{-1} \bar{l}(\theta^{[k]}; \mathbf{y}) + \phi^{-1} \langle \nabla_{\theta} \bar{l}(\theta^{[k]}), \theta - \theta^{[k]} \rangle + \frac{\rho}{2\phi} \|\theta - \theta^{[k]}\|_F^2 - \frac{n}{2} \log \det \mathbf{W} + \frac{1}{2} \mathbf{c}^\top (\mathbf{W} \otimes \mathbf{I}) \mathbf{c} + P_g(\beta; \lambda) + P_W(\mathbf{W}; \lambda_W), \quad (1.31)$$

where  $\rho$  is a step size which can be tuned. Then at time  $t + 1$  we obtain

$$(\beta^{[k+1]}, \mathbf{W}^{[k+1]}, \mathbf{c}^{[k+1]}) = \underset{(\beta \in \mathbb{R}^{mp}; \mathbf{W} \in \mathcal{S}_{++}^m; \mathbf{c} \in \mathbb{R}^{mp})}{\operatorname{argmin}} g(\beta, \mathbf{W}, \mathbf{c}; \beta^{[k]}, \mathbf{W}^{[k]}, \mathbf{c}^{[k]}). \quad (1.32)$$

Convergence of our algorithm under the Lipschitz conditions on the gradients of all loss functions are natural: suppose the following Lipschitz conditions hold for all loss  $l_k$ ,

$$\|\nabla l_k(\boldsymbol{\theta}_1) - \nabla l_k(\boldsymbol{\theta}_2)\|_2 \leq L\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2, \quad \forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2. \quad (1.33)$$

Then as long as we set  $\rho \geq L$  in the above algorithm, the values of objective function  $F$  satisfies:

$$F(\boldsymbol{\beta}^{[k+1]}, \mathbf{W}^{[k+1]}, \mathbf{c}^{[k+1]}) \leq F(\boldsymbol{\beta}^{[k]}, \mathbf{W}^{[k]}, \mathbf{c}^{[k]}) \quad (1.34)$$

by standard convergence analysis of gradient descent (e.g. see [23]). Instead of directly solving the optimization problem (1.31), we can first apply linearization on  $\boldsymbol{\theta}$  to turn the problem into Gaussian and then apply the weighting scheme. Define  $\boldsymbol{\xi}^{[k+1]} = \boldsymbol{\theta}^{[k]} - \nabla \bar{l}(\boldsymbol{\theta}^{[k]}; \mathbf{y})/\rho$ , then the optimization problem (1.31) can be rewritten as

$$\min_{\boldsymbol{\beta}, \mathbf{c}, \mathbf{W}} \frac{\rho}{2\phi} \|\mathbf{L}\boldsymbol{\beta} + (\mathbf{I} - \phi\mathbf{W} \otimes \mathbf{I})^{1/2} \mathbf{c} - \boldsymbol{\xi}^{[k+1]}\|_F^2 - \frac{n}{2} \log \det \mathbf{W} + \frac{1}{2} \mathbf{c}^\top (\mathbf{W} \otimes \mathbf{I}) \mathbf{c} + P_g(\boldsymbol{\beta}; \lambda) + P_W(\mathbf{W}; \lambda_W). \quad (1.35)$$

In order to obtain the closed form of  $\mathbf{c}^{[k+1]}$  for Gaussian case, we need to set  $\rho = 1$ . This can be done by scaling the loss function  $\bar{l}$  by the Lipschitz constant  $L$ . Given  $\boldsymbol{\beta}$  and  $\mathbf{W}$ , we have the explicit form  $\mathbf{c}^{[k+1]} = (\mathbf{I} - \phi\mathbf{W} \otimes \mathbf{I})^{1/2} (\boldsymbol{\xi}^{[k+1]} - \mathbf{L}\boldsymbol{\beta})$ . Then plugging the form into the above problem (1.35), we obtain the following optimization problem for  $\boldsymbol{\beta}^{[k+1]}$  and  $\mathbf{W}^{[k+1]}$

$$\min_{\boldsymbol{\beta}, \mathbf{W}} \|\mathbf{W}^{\frac{1}{2}} \text{vec}^{-1}(\boldsymbol{\xi}^{[k+1]} - \mathbf{L}\boldsymbol{\beta})\|_F^2 - \frac{n}{2} \log \det \mathbf{W} + P_g(\boldsymbol{\beta}; \lambda) + P_W(\mathbf{W}; \lambda_W). \quad (1.36)$$

Now the optimization becomes a natural extension of traditional weighted least squares applied in high dimensional problems when given  $\mathbf{W}$ . In this case, efficient gradient descent type algorithm or coordinate descent algorithm can be applied to solve the above optimization for  $\boldsymbol{\beta}$  just like the group lasso. While when  $\boldsymbol{\beta}$  is given, we are back to a Gaussian graph learning problem, like graphical lasso when the  $\ell_1$  penalty is used for  $\mathbf{W}$ . Thus a block coordinate descent algorithm can be applied for solving  $\boldsymbol{\beta}^{[k+1]}$  and  $\mathbf{W}^{[k+1]}$ . In the inner loop of the BCD algorithm, it can be terminated as long as the function value of the surrogate function decreases. Then  $\boldsymbol{\theta}^{[k+1]}$  can also be obtained without explicitly calculating  $\mathbf{c}^{[k+1]}$ :

$$\boldsymbol{\theta}^{[k+1]} = \mathbf{L}\boldsymbol{\beta}^{[k+1]} + (\mathbf{I} - \phi\mathbf{W} \otimes \mathbf{I}^{[k+1]})^{\frac{1}{2}} \mathbf{c}^{[k+1]} = \phi\mathbf{L}\boldsymbol{\beta}^{[k+1]} + (\mathbf{I} - \phi\mathbf{W} \otimes \mathbf{I}^{[k+1]}) \boldsymbol{\xi}^{[k+1]}. \quad (1.37)$$

To conclude, the full algorithm is summarized as follows:

<p><b>Algorithm 1:</b> Multivariate Learning with Dependency</p> <p><b>Data:</b> <math>\mathbf{Y} \in \mathbb{R}^{n \times m}</math>, <math>\mathbf{X} \in \mathbb{R}^{n \times p}</math>, <math>\mathbf{y} = \text{vec}(\mathbf{Y})</math>, <math>\mathbf{L} = \mathbf{I}_{m \times m} \otimes \mathbf{X}</math>, <math>\bar{l}</math> such that <math>\nabla \bar{l}</math> is Lipschitz-1 continuous, initial value <math>\boldsymbol{\xi}^{[0]}</math>, <math>\boldsymbol{\beta}^{[0]}</math> and <math>\mathbf{W}^{[0]} \in \mathcal{S}_{++}^m</math>, <math>\phi \leq \ \mathbf{W}^{[0]}\ _2^{-1}</math>.</p> <p>Initialize iteration number <math>k \rightarrow 0</math>;</p> <p><b>while</b> <i>not convergence</i> <b>do</b></p> <p style="padding-left: 2em;"><math>k \rightarrow k + 1</math>;</p> <p style="padding-left: 2em;"><math>\boldsymbol{\theta}^{[k-1]} = \phi \mathbf{L} \boldsymbol{\beta}^{[k-1]} + (\mathbf{I} - \phi \mathbf{W} \otimes \mathbf{I}) \boldsymbol{\xi}^{[k-1]}</math>, <math>\boldsymbol{\xi}^{[k]} = \boldsymbol{\theta}^{[k-1]} - \nabla \bar{l}(\boldsymbol{\theta}^{[k]}; \mathbf{y})</math>;</p> <p style="padding-left: 2em;"><math>t \rightarrow 0</math>, <math>\boldsymbol{\beta}^{(0)} = \boldsymbol{\beta}^{[k-1]}</math>, <math>\mathbf{W}^{(0)} = \mathbf{W}^{[k-1]}</math>;</p> <p style="padding-left: 2em;"><b>while</b> <math>g(\boldsymbol{\beta}^{(t)}, \mathbf{W}^{(t)}, \mathbf{c}^{(t)}; \boldsymbol{\beta}^{k-1}, \mathbf{W}^{k-1}) &gt; g(\boldsymbol{\beta}^{(t-1)}, \mathbf{W}^{(t-1)}, \mathbf{c}^{(t-1)})</math> <i>if existing</i> <b>do</b></p> <p style="padding-left: 4em;"><math>t \rightarrow t + 1</math>;</p> <p style="padding-left: 4em;"><math>\mathbf{W}^{(t)} \rightarrow \underset{\mathbf{W}}{\text{argmin}} \ \mathbf{W}^{\frac{1}{2}} \text{vec}^{-1}(\boldsymbol{\xi}^{[k+1]} - \mathbf{L} \boldsymbol{\beta}^{(t-1)})\ _F^2 - \frac{n}{2} \log \det \mathbf{W} + P_W(\mathbf{W}; \lambda_W)</math>;</p> <p style="padding-left: 4em;"><math>\boldsymbol{\beta}^{(t)} \rightarrow \underset{\boldsymbol{\beta}}{\text{argmin}} \ (\mathbf{W} \otimes \mathbf{I}^{(t)\frac{1}{2}})(\boldsymbol{\xi}^{[k+1]} - \mathbf{L} \boldsymbol{\beta})\ _F^2 + P_g(\boldsymbol{\beta}; \lambda_W)</math></p> <p style="padding-left: 2em;"><b>end</b></p> <p style="padding-left: 2em;"><math>\boldsymbol{\beta}^{[k]} \rightarrow \boldsymbol{\beta}^{(t)}</math>, <math>\mathbf{W}^{[k]} \rightarrow \mathbf{W}^{(t)}</math>;</p> <p><b>end</b></p> <p><b>Result:</b> Output the final estimate <math>\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}^{[k]}</math>, <math>\hat{\mathbf{W}} = \mathbf{W}^{[k]}</math>.</p>
---

Note that for clustered data case, there are two main differences: first, since  $\mathbf{W}$  capture the dependency among rows, we need change  $W$ - optimization into

$$\mathbf{W}^{(t)} \rightarrow \underset{\mathbf{W}}{\text{argmin}} \|\text{vec}^{-1}(\boldsymbol{\xi}^{[k+1]} - \mathbf{L} \boldsymbol{\beta}^{(t-1)}) \mathbf{W}^{\frac{1}{2}}\|_F^2 - \frac{n}{2} \log \det \mathbf{W} + P_W(\mathbf{W}; \lambda_W). \quad (1.38)$$

Besides, there should be a constraint on  $\boldsymbol{\beta}$  optimization such that  $\boldsymbol{\beta}^{(t)}$  should be projected in to the space  $\{\boldsymbol{\beta} : \exists \boldsymbol{\beta}_0 \in \mathbb{R}^p, \boldsymbol{\beta} = \mathbf{I}_{m \times 1} \otimes \boldsymbol{\beta}_0\}$ .

## 1.6 Experiments

### 1.6.1 Simulation Studies

In this section, we provide some simulation results to support our main theorems. Consider the Gaussian graph model  $\mathbf{Y} = \mathbf{M} + \mathbf{E}$  where we have  $\mathbf{M} = \mathbf{0}$  and the rows of  $\mathbf{E}$  are independent distributed as  $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}^*)$  with  $\mathbf{W}^* = \boldsymbol{\Sigma}^{*-1}$ . Suppose  $\mathbf{W}^*$  is sparse with index  $\mathcal{J}^*$  on the upper diagonal parts. The main goal is to perform statistical inference on the nonzero part of true dependency  $\mathbf{W}_{\mathcal{J}^*}^*$ . Given the generated data  $\mathbf{Y}$ , the steps to perform statistical inference of the true dependency matrix  $\mathbf{W}^*$  are as follows: First, an estimator  $\hat{\mathbf{W}}$  with nonzero index  $\hat{\mathcal{J}}$  can be obtained through our GIDL framework with an adaptive lasso penalty, where we assume  $\mathcal{J}^* \subset \hat{\mathcal{J}}$  by choosing proper  $\lambda$ . Then, the sample covariance of  $\text{vec}[-\sqrt{n}\{\hat{\boldsymbol{\Sigma}}_{\hat{\mathcal{J}}^*}^n - \boldsymbol{\Sigma}_{\hat{\mathcal{J}}^*}^*\}]$  can be calculated

Table 1.1: Converge probability of the confidence interval constructed based on the Hotelling’s  $t$ -Squared statistics derived from GIDL

%	Case 1	Case 2	Case 3	Case 4
$C = 0.90$	89.00	92.00	89.00	88.67
$C = 0.95$	94.67	96.00	93.00	93.67
$C = 0.99$	97.33	98.67	97.67	97.00

as a plug-in estimate of its population correspondence. Finally, based on the asymptotic normal distribution on the index  $\mathcal{J}^*$ , we can construct a Hotelling’s  $t$ -Squared type of statistics to perform the statistical inference. We perform simulations to demonstrate that the above steps can deliver a reasonable result for statistical inference. To generate the true sparse dependency matrix, we set  $\rho$  as the percentage of non-zero off-diagonal elements, randomly generate nonzero index for the off diagonal part based on probability  $\rho$  and then let  $W_{ij} \sim \mathcal{N}(0, 1)$  for  $(i, j) \in \mathcal{J}^*$ . We show the coverage probability of the constructed confidence interval under the following cases of simulation settings at confidence level  $C = 0.99, 0.95, 0.90$  by performing Monte Carlo simulations for 300 times:

1. (Case 1:)  $n = 100, m = 10, \rho = 0.1$ ;
2. (Case 2:)  $n = 500, m = 10, \rho = 0.1$ ;
3. (Case 3:)  $n = 500, m = 10, \rho = 0.5$ ;
4. (Case 4:)  $n = 500, m = 20, \rho = 0.1$ .

The simulation results are shown in table 1.1. From the table, we can see that as long as the sample size  $n$  is large enough, the coverage probabilities for all confidence intervals are almost match the corresponding confidence levels. This result implies that our method can perform a valid statistical inference for the nonzero part of the true dependency  $\mathbf{W}_{\mathcal{J}^*}^*$ .

We also perform a comparison between squared loss and Tukey’s loss when applying main theorem to calculate confidence intervals and the presence of outliers is a problem. The Tukey’s loss is a well-known robust loss which is defined by a  $\psi$ -function:  $l(\theta, y) = \int_0^{|\theta-y|} \psi(t)dt$  [5], where Tukey’s bisquare is

$$\psi(t) = \begin{cases} t \left[1 - \left(\frac{t}{c}\right)^2\right]^2, & |t| \leq c \\ 0, & |t| > c \end{cases} \quad (1.39)$$

Table 1.2: Converge probability and length of the confidence interval comparisons between squared and Tukey’s loss by GIDL

	Case 1		Case 2		Case 3		$C$
	Coverage %	Length	Coverage %	Length	Coverage %	Length	
Gaussian	100	14.99	100	17.35	100	20.06	0.90
Tukey	82.67	0.88	88	0.43	81	0.44	
Gaussian	100	16.47	100	19.65	100	22.42	0.95
Tukey	89.33	1.09	93.67	0.52	89.6667	0.53	
Gaussian	100	19.31	100	23.07	100	25.82	0.99
Tukey	97	1.44	98.33	0.67	97	0.670	

and  $c = 4.685\sigma$  is recommended. To generate synthetic data, we use the model  $\mathbf{Y} = \mathbf{M} + \mathbf{\Psi} + \mathbf{E}$  where the settings of  $\mathbf{M}$  and  $\mathbf{E}$  are still the same with the above simulation. In addition,  $\mathbf{\Psi} \in \mathbb{R}^{n \times m}$  is a sparse matrix that incorporates the outlying effects in  $\mathbf{Y}$  as its components have 0.95 probability to be zero and 0.05 probability to be 20. In order to compare the length of calculated confidence intervals, we only set one dependent edge on the off-diagonal parts of  $\mathbf{W}^*$ :  $|\mathcal{J}^*| = 1$  where its location is random assigned as before. We use the same approach as the above simulations to calculate both the confidence intervals and length of the confidence interval of  $\mathbf{W}_{\mathcal{J}^*}$  under both squared loss and Tukey’s loss. Three simulation cases are considered: Case 1:  $n = 100, m = 10$ ; Case 2:  $n = 500, m = 10$ ; Case 3:  $n = 500, m = 20$  and the average coverage probabilities and lengths at confidence level  $C = 0.90, 0.95, 0.99$  are reported for 300 Monte Carlo simulations. The simulation result is listed in table 1.2. From the table, some phenomenon can be observed: first, the confidence intervals obtained via squared loss are not informative, since the length of the interval is too large so that the covering probabilities are always 100%; moreover, the coverage probabilities of the intervals obtained via Tukey’s loss is less than what we desire, which is reasonable since our model is misspecified for the true data generating process; finally, when  $n$  is large enough, we can see the coverage probabilities for the Tukey’s loss are approaching to the desired ones. The above conclusions strike us to apply robust loss functions to make statistical inference when outlying effect must be taken into account.

### 1.6.2 Real Data

We explore the NASDAQ-100 stock market index data from 01/03/2011 to 12/31/2011. The differences between closed price of two consecutive days is the main target, where we want to detect

and make inference on the dependency structures of these differences between the top 98 largest non-financial companies listed on the Nasdaq stock market. Since the differences could be extremely large in some days, it is more reasonable to adopt some robust loss functions when apply our GIDL framework than traditional squared loss. Here we provide some exploration to demonstrate that robust loss is a more reasonable choice by making inference on the first detected dependent edges through both squared loss and Tukey’s loss. By tuning the parameter  $\lambda$  of the adaptive lasso penalty, we make sure that only one edge can be detected and then calculate the confidence interval of the weigh of the edge through our main theorem under both kinds of losses. For Gaussian loss, the detected edge is the link between ‘FAST’ and ‘ALXN’ with weight  $-5.47$ . However, the 95% confidence interval of this edge we calculate is  $[-64.8, 53.9]$ , which has a unreasonable magnitude and incorrectly contains zero. The two problems imply that assuming Gaussian distribution is even useless in making inference on the dependent edge. Instead, when we apply the Tukey’s loss, the detected edge is the link between ‘LLTC’ and ‘TXN’ with weight  $-5.48$ . The corresponding 95% confidence interval is  $[-6.2, -0.6]$ , which doesn’t contain zero and enjoy a more reasonable length. Therefore, We have reasons to believe that for GIDL, robust loss functions are essential for reasonable inference through our main theorems when the presence of outliers is a problem.

## 1.7 Summary

In this section, we explore the statistical inference problem for GIDL where we estimate the dependency matrix or incorporate the structure information into estimating the first order effect with the allowance of customizing the marginal losses. Some asymptotic theories are provided to assistant the inference. However, there are still some problems to be answered in the future. It is a great of interest to make inference for high dimensional estimated dependencies and coefficients obtained through our framework. Whether incorporating the structure information can indeed lead to more efficient estimation as the traditional asymptotic cases should be explained in the future.

## 1.8 Outlines of Proofs

### 1.8.1 Proof of Theorem 1.3.1

The first Step: Solving  $\mathcal{C}$  based on expansion of gradient of the loss:

By taking derivatives of the objective function with respect to  $\mathbf{C}$ , we have  $\nabla \bar{l}(\boldsymbol{\Theta}) \mathbf{Z}^{1/2} + \phi \mathbf{C} \mathbf{W} = \mathbf{0}$ . Note that we also have  $\nabla \bar{l}(\boldsymbol{\Theta}) = \nabla \bar{l}(\mathbf{M}) + \mathbf{C} \mathbf{Z}^{1/2} + \boldsymbol{\Delta}$ . Then  $\mathbf{C}$  can be solved as

$$\mathbf{C} = -(\nabla \bar{l}(\mathbf{M}) + \boldsymbol{\Delta}) \mathbf{Z}^{\frac{1}{2}} (\mathbf{Z} + \phi \mathbf{W})^{-1} = -(\nabla \bar{l}(\mathbf{M}) + \boldsymbol{\Delta}) \mathbf{Z}^{\frac{1}{2}}. \quad (1.40)$$

The second Step: Plugging  $\mathbf{C}$  back into the loss, obtaining an approximation Gaussian representation of the loss based through direct expansion:

By plugging the expression of  $\mathbf{C}$  into the original optimization problem, we have

$$\begin{aligned} & \bar{l}(\boldsymbol{\Theta}) + \frac{\phi}{2} \text{Tr} \{ \mathbf{C} \mathbf{W} \mathbf{C}^T \} - \bar{l}(\mathbf{M}) \\ &= \langle \nabla \bar{l}(\mathbf{M}) + \frac{1}{2} \mathbf{C} \mathbf{Z}^{\frac{1}{2}}, \mathbf{C} \mathbf{Z}^{\frac{1}{2}} \rangle + \frac{\phi}{2} \text{Tr} \{ \mathbf{C} \mathbf{W} \mathbf{C}^T \} + \Gamma(\boldsymbol{\Theta}) \\ &= \langle \nabla \bar{l}(\mathbf{M}) + \boldsymbol{\Delta}, \mathbf{C} \mathbf{Z}^{\frac{1}{2}} \rangle + \frac{1}{2} \langle \mathbf{C} \mathbf{Z}^{\frac{1}{2}}, \mathbf{C} \mathbf{Z}^{\frac{1}{2}} \rangle + \langle -\boldsymbol{\Delta}, \mathbf{C} \mathbf{Z}^{\frac{1}{2}} \rangle + \frac{\phi}{2} \text{Tr} \{ \mathbf{C} \mathbf{W} \mathbf{C}^T \} + \delta(\boldsymbol{\Theta}) \\ &= \langle \nabla \bar{l}(\mathbf{M}) + \boldsymbol{\Delta}, \mathbf{C} \mathbf{Z}^{\frac{1}{2}} \rangle + \frac{1}{2} \text{Tr} \{ \mathbf{C} (\mathbf{Z} + \phi \mathbf{W}) \mathbf{C}^T \} + \langle -\boldsymbol{\Delta}, \mathbf{C} \mathbf{Z}^{\frac{1}{2}} \rangle + \delta(\boldsymbol{\Theta}) \\ &= -\frac{1}{2} \text{Tr} \{ (\nabla \bar{l}(\mathbf{M}) + \boldsymbol{\Delta}) \mathbf{Z} (\nabla \bar{l}(\mathbf{M}) + \boldsymbol{\Delta})^T \} + \langle \boldsymbol{\Delta}, (\nabla \bar{l}(\mathbf{M}) + \boldsymbol{\Delta}) \mathbf{Z} \rangle + \delta(\boldsymbol{\Theta}) \\ &= -\frac{1}{2} \text{Tr} \{ \nabla \bar{l}(\mathbf{M}) \mathbf{Z} \nabla \bar{l}(\mathbf{M})^T \} + \delta(\boldsymbol{\Theta}). \end{aligned}$$

This results in

$$\begin{aligned} & \phi^{-1} \bar{l}(\boldsymbol{\Theta}) + \frac{1}{2} \text{Tr} \{ \mathbf{C} \mathbf{W} \mathbf{C}^T \} - \phi^{-1} \bar{l}(\mathbf{M}) \\ &= \frac{1}{2} \text{Tr} \{ \nabla \bar{l}(\mathbf{M}) \mathbf{W} \nabla \bar{l}(\mathbf{M})^T \} - \frac{1}{2\phi} \|\nabla \bar{l}(\mathbf{M})\|_F^2 + \phi^{-1} \delta(\boldsymbol{\Theta}). \end{aligned}$$

Denote the objective function as  $F(\mathbf{W}; \mathbf{C})$  and  $F(\mathbf{W})$  as the abbreviation for  $F(\mathbf{W}; \mathbf{C}(\mathbf{W}))$ . Then we have an approximation Gaussian representation of the loss:

$$F(\mathbf{W}) = \frac{1}{2} \text{Tr} \{ \nabla \bar{l}(\mathbf{M}) \mathbf{W} \nabla \bar{l}(\mathbf{M})^T \} - \frac{1}{2\phi} \|\nabla \bar{l}(\mathbf{M})\|_F^2 - \frac{n}{2} \log \det \mathbf{W} + P_{\mathbf{W}}(\mathbf{W}; \lambda_{\mathbf{W}}) + \phi^{-1} \delta(\boldsymbol{\Theta}).$$

We need to investigate the terms containing  $\boldsymbol{\Delta}(\boldsymbol{\Theta})$  and  $\delta(\boldsymbol{\Theta})$  such that the effect of them are always higher order terms in the following analysis. Based on condition 5, we have two useful results:

$$\|\boldsymbol{\Delta}^*\|_F = o_p(1)$$

and

$$\left\| \frac{d\delta(\boldsymbol{\Theta})}{d\mathbf{W}} \right\|_F = o_p(n^{1/2}).$$

The third Step: showing rate consistency of  $\widehat{\mathbf{W}}$  by drawing a circle: Next we consider to show the rate consistency of  $\widehat{\mathbf{W}}$ . In particular, we want to show that for any given  $\epsilon > 0$ , there exists a large

constant  $C_1$  such that

$$P\left\{\inf_{\|\mathbf{U}\|_F=C_1} F(\mathbf{W}^* + d_n \mathbf{U}) - F(\mathbf{W}^*) > 0\right\} \geq 1 - \epsilon,$$

where  $\mathbf{U}$  is invertible and  $d_n = n^{-1/2} + a_n$ . This implies that with probability at least  $1 - \epsilon$ , there exists a local minimum in the ball  $\{\mathbf{W} : \|\mathbf{W} - \mathbf{W}^*\|_F \leq d_n C_1\}$ . By Taylor's expansion, there exists a constant  $0 < t < 1$ , such that

$$\begin{aligned} & F(\mathbf{W}^* + d_n \mathbf{U}) - F(\mathbf{W}^*) \\ &= \frac{1}{2} \langle d_n \mathbf{U}, \nabla \bar{l}(\mathbf{M})^T \nabla \bar{l}(\mathbf{M}) \rangle - \frac{n}{2} (\log \det(\mathbf{W}^* + d_n \mathbf{U}) - \log \det \mathbf{W}^*) + \phi^{-1}[\delta(\boldsymbol{\Theta}) - \delta(\boldsymbol{\Theta}^*)] \\ &+ P_W(\mathbf{W}^* + d_n \mathbf{U}; \lambda_W) - P_W(\mathbf{W}^*; \lambda_W) \\ &\geq \frac{n}{2} \langle d_n \mathbf{U}, \hat{\boldsymbol{\Sigma}}^n - \mathbf{W}^{*-1} \rangle + \frac{n}{2} \langle d_n \mathbf{U}, \frac{1}{2} d_n \mathbf{W}^{*-1} \mathbf{U} \mathbf{W}^{*-1} \rangle - \frac{n}{2} \mathbf{R}(d_n \mathbf{U}) + \phi^{-1}[\delta(\boldsymbol{\Theta}) - \delta(\boldsymbol{\Theta}^*)] \\ &+ \sum_{(i,j) \in \mathcal{J}_W^*} d_N U_{ij} P'_W(|W_{ij}^*|; \lambda_W) \text{sgn}(|W_{ij}^*|) + \sum_{(i,j) \in \mathcal{J}_W^*} \frac{1}{2} d_N^2 P''_W(W_{ij}^*; \lambda_W) U_{ij}^2 (1 + o(1)) \\ &- \frac{1}{2} \langle d_n \mathbf{U}, 2\boldsymbol{\Delta}^{*T} \nabla \bar{l}(\mathbf{M}) + \boldsymbol{\Delta}^{*T} \boldsymbol{\Delta}^* \rangle \\ &\equiv T_1 + T_2 + T_3 + T_4 + T_5 + T_6 + T_7, \end{aligned}$$

where  $\mathbf{R}(d_n \mathbf{U})$  is the reminder for the log det function with  $\mathbf{R}(\tilde{\mathbf{U}}) = \text{Tr}(\sum_{j=3}^{\infty} (-1)^{j-1} (\tilde{\mathbf{U}} \mathbf{W}^{*-1})^j / j)$  for any  $\tilde{\mathbf{U}}$  is based on the identity

$$\log \det(\mathbf{W}^* + d_n \mathbf{U}) - \log \det \mathbf{W}^* = \log \det(\mathbf{I} + d_n \mathbf{U} \mathbf{W}^{*-1}) = \text{Tr} \log(\mathbf{I} + d_n \mathbf{U} \mathbf{W}^{*-1})$$

when  $\mathbf{I} + d_n \mathbf{U} \mathbf{W}^{*-1}$  is positive-definite.

For the term  $T_1$ ,  $\mathbb{E}(\mathbf{W}^{*-1}) = (\nabla \bar{l}(\mathbf{M}) + \boldsymbol{\Delta}^*)^T (\nabla \bar{l}(\mathbf{M}) + \boldsymbol{\Delta}^*) / n$ , it holds that  $|T_1| = O_p(C_1 d_n \sqrt{n})$  by the central limit theorem. For the term  $T_2$ , by the positive-definiteness of  $\mathbf{W}^{*-1} \otimes \mathbf{W}^{*-1}$ , we have  $T_2 \geq C_1^2 n d_n^2 / (4\gamma_{max}^2)$ ; moreover, since we have

$$|R(d_n \mathbf{U})| \leq \sum_{j=3}^{\infty} (-1)^{j-1} d_n^j \|\mathbf{U}\|_F^j \|\mathbf{W}^{*-1}\|_F^j / j \leq d_n^3 \|\mathbf{U}\|_F^3 \|\mathbf{W}^{*-1}\|_F^3 / 3,$$

this implies that  $|T_3| = O_p(C_1^3 n d_n^3)$ . By the mean value theorem, there exists a  $\tilde{\boldsymbol{\Theta}}$  between  $\boldsymbol{\Theta}$  and  $\boldsymbol{\Theta}^*$  such that  $|\delta(\boldsymbol{\Theta}) - \delta(\boldsymbol{\Theta}^*)| = \langle d\delta(\tilde{\boldsymbol{\Theta}}) / d\mathbf{W}, d_n \mathbf{U} \rangle$ , so we have  $|T_4| = o_p(C_1 \sqrt{n} d_n)$  based on  $\|d\delta(\boldsymbol{\Theta}) / d\mathbf{W}\|_F = o_p(1)$ . For the terms with respect to the penalties, we have  $T_5 \geq -\sqrt{J_W^*} n d_n a_n C_1$  and the term  $T_6$  is bounded below by  $-n d_N^2 C_1^2 \max_{(j,k) \in \mathcal{J}^*} P''_W(W_{jk}^*; \lambda_W)$ . Finally, the term  $T_7$  is at the order of  $o_p(C_1 d_n \sqrt{n})$  based on  $\|\nabla \bar{l}(\mathbf{M})\|_F = O_p(\sqrt{n})$ ,  $\|\boldsymbol{\Delta}^*\|_F = o_p(1)$ . Therefore, by choosing



a large enough  $C_1$ , the summation from  $T_1$  to  $T_7$  is dominated by the positiveness of  $T_2$ . When on the boundary, all values are greater than  $F(\mathbf{W}^*)$ , then there must be a local minimizer within the ball, denoted as  $\widehat{\mathbf{W}}$ , so  $\|\widehat{\mathbf{W}} - \mathbf{W}^*\|_F = O_p(n^{-1/2} + a_n)$  holds by choosing  $\widehat{\mathbf{W}}$  such that  $\|\widehat{\mathbf{W}} - \mathbf{W}^*\|_F$  is the smallest among all the roots within the region.

To prove the second part, for any  $\mathbf{W}$  such that  $\|\mathbf{W} - \mathbf{W}^*\|_F = O_p(n^{-1/2})$  and thus  $\max_{j \in \mathcal{J}_{\mathcal{W}}^{*c}} |W_j| = O_p(n^{-1/2})$ . Then we have

$$\begin{aligned}
\frac{\partial F(\mathbf{W})}{\partial \mathbf{W}} &= \frac{\nabla \bar{l}(\mathbf{M})^T \nabla \bar{l}(\mathbf{M})}{2} - \frac{n}{2} \mathbf{W}^{-1} + P'_W(|\mathbf{W}|; \lambda_W) \text{sgn}(\mathbf{W}) + \phi^{-1} \frac{d\delta(\Theta)}{d\mathbf{W}} \\
&= \frac{n}{2} \{\hat{\Sigma}^n - \mathbf{W}^{*-1}\} + \left(-\frac{n}{2} \mathbf{W}^{-1} + \frac{n}{2} \mathbf{W}^{*-1}\right) + \phi^{-1} \frac{d\delta(\Theta)}{d\mathbf{W}} - \frac{1}{2} (2\Delta^{*T} \nabla \bar{l}(\mathbf{M}) + \Delta^{*T} \Delta^*) \\
&\quad + P'_W(|\mathbf{W}|; \lambda_W) \text{sgn}(\mathbf{W}) \\
&= T'_1 + T'_2 + T'_3 + T'_4 + T'_5.
\end{aligned} \tag{1.41}$$

Since  $\|\mathbf{W} - \mathbf{W}^*\|_F = O_p(n^{-1/2})$ , the Frobenius norm of the first and second terms can be bounded by  $O_p(\sqrt{n})$ . For the term  $T'_3$ , we have  $\|d\delta(\tilde{\Theta})/d\mathbf{W}\|_F = o_p(1)$ . Based on  $\|\nabla \bar{l}(\mathbf{M})\|_F = O_p(\sqrt{n})$ ,  $\|\Delta^*\|_F = o_p(1)$ , we have  $\|T'_4\|_F = o_p(\sqrt{n})$ . Finally, for the last term, since  $\lim_{n \rightarrow \infty} \lim_{t \rightarrow 0+} P'_W(t; \lambda_W)/\lambda_W > 0$  and  $\lim_{n \rightarrow \infty} \lambda_W/\sqrt{n} \rightarrow \infty$ , all its elements can dominate the Frobenius norm of the other terms. Then for  $j \in \mathcal{J}_{\mathcal{W}}^{*c}$ , we have  $\frac{\partial F(\mathbf{W})}{\partial W_j} < 0$  for  $0 < W_j$  and  $\frac{\partial F(\mathbf{W})}{\partial W_j} > 0$  for  $W_j > 0$ . Note that when  $n \rightarrow \infty$ ,  $a_n \rightarrow 0$ . Since our defined  $\widehat{\mathbf{W}}$  is a root  $n$  consistent estimator, it satisfies  $\widehat{\mathbf{W}}_{\mathcal{J}_{\mathcal{W}}^{*c}} = 0$ .

For the third part, by the above Taylor expansion (2.80), the  $\sqrt{n}$  consistency of  $\widehat{\mathbf{W}}$ , we have

$$\begin{aligned}
\mathbf{0} &= \frac{n}{2} \text{vec} \{\hat{\Sigma}^n - \mathbf{W}^{*-1}\} + \frac{n}{2} \mathbf{W}^{*-1} \otimes \mathbf{W}^{*-1} \text{vec}(\widehat{\mathbf{W}} - \mathbf{W}^*) (1 + o_p(1)) \\
&\quad + \text{vec} \{P'_W(|\mathbf{W}^*|; \lambda_W) \text{sgn}(\mathbf{W}^*)\} + \tilde{\mathbf{R}},
\end{aligned}$$

where  $\tilde{\mathbf{R}}$  is a reminder term satisfying  $\|\tilde{\mathbf{R}}\|_F = o_p(\sqrt{n})$ . Therefore, by only considering the submatrix on  $\mathcal{J}_{\mathcal{W}^*}$  and based on the conditions for penalties, we have

$$\sqrt{n}(\text{vec} \widehat{\mathbf{W}}_{\mathcal{J}_{\mathcal{W}}^*} - \text{vec} \mathbf{W}_{\mathcal{J}_{\mathcal{W}}^*}^*) = \{\Sigma_{\mathcal{J}_{\mathcal{W}}^*}^* \otimes \Sigma_{\mathcal{J}_{\mathcal{W}}^*}^*\}^{-1} \text{vec} \left[ -\sqrt{n} \{\hat{\Sigma}_{\mathcal{J}_{\mathcal{W}}^*} - \Sigma_{\mathcal{J}_{\mathcal{W}}^*}^*\} \right] (1 + o_p(1)),$$

where recall  $\Sigma^* = \mathbf{W}^{*-1}$ . Finally, the conclusion follows by Slutsky's theorem.

### 1.8.2 Proof of Theorem 1.3.2

For simplicity, when the mean matrix  $\mathbf{M}$  is known, we can perform a location transformation such that  $\mathbf{Y} \leftarrow \mathbf{Y} - \mathbf{M}$ . Then based on the assumption that  $\mathbf{M} = \mathbf{1}\boldsymbol{\alpha}^T$ , we can scale each loss

function  $l_k(\cdot)$  by  $l_k''(\alpha_k) = 1/4$ . Based on the scaling, we have  $\mathbf{D} = \mathbf{I}$ . Apply the Theorem 2 in [23], there is a  $\mathbf{W}_{gidl}$ , such that

$$\mathbf{W}^{gidl^{-1}} = \frac{(\nabla l(\mathbf{M}) + \mathbf{\Delta}(\mathbf{M}, \mathbf{W}, \mathbf{C}))^T (\nabla l(\mathbf{M}) + \mathbf{\Delta}(\mathbf{M}, \mathbf{W}, \mathbf{C}))}{n}. \quad (1.42)$$

Note that due to relocation and scaling,  $-\nabla l(\mathbf{M}) = 4\mathbf{Y}$ , thus as the regularity condition for high order term satisfied, we have

$$\mathbf{W}^{gidl^{-1}} \rightarrow 16 \frac{\mathbf{Y}^T \mathbf{Y}}{n}. \quad (1.43)$$

Denote  $\hat{\mathbf{W}} = (\mathbf{Y}^T \mathbf{Y} / n)^{-1}$ , we now prove as  $n \rightarrow \infty$ , the stationary condition for pseudo-likelihood can be satisfied when applying the graph structure decided by  $\hat{\mathbf{W}}$ . Note that  $\hat{\mathbf{W}}_{ii} = 1$  for  $i = 1, \dots, m$ . Based on inverse of block matrix (Schur component), we have

$$\hat{\mathbf{W}}[\{-s\}, s] \hat{\mathbf{W}}[s, s]^{-1} = -(\mathbf{Y}[\{-s\}]^T \mathbf{Y}[\{-s\}])^{-1} \mathbf{Y}[\{-s\}]^T \mathbf{Y}[s]. \quad (1.44)$$

Then if we consider the graph structure defined by  $-\hat{\mathbf{W}}$  that is the same with  $\hat{\mathbf{W}}$ , we have the property

$$\begin{aligned} -\frac{1}{n} \sum_{i=1}^n Y_{ik} \sum_{t \neq s} \hat{\mathbf{W}}_{st} Y_{it} &= -\frac{1}{n} \sum_{i=1}^n Y_{ik} \sum_{t \neq s} \hat{\mathbf{W}}_{st} Y_{it} \\ &= \frac{1}{n} \mathbf{Y}[k]^T \mathbf{Y}[\{-s\}] (\mathbf{Y}[\{-s\}]^T \mathbf{Y}[\{-s\}])^{-1} \mathbf{Y}[\{-s\}]^T \mathbf{Y}[s] \\ &\stackrel{(i)}{=} \frac{1}{n} \mathbf{Y}[k]^T \mathbf{Y}[s], \end{aligned} \quad (1.45)$$

where (i) is because  $\mathbf{Y}[\{-s\}] (\mathbf{Y}[\{-s\}]^T \mathbf{Y}[\{-s\}])^{-1} \mathbf{Y}[\{-s\}]^T$  is a projection matrix that can map  $\mathbf{Y}[k]$  to itself. Moreover, if we apply the graph structure for  $-\hat{\mathbf{W}}$  into the gradient for pseudo-likelihood, we have

$$\begin{aligned} -\frac{\partial l}{\partial W_{sk}} &= -\frac{1}{n} \sum_{i=1}^n \left[ Y_{is} Y_{ik} - Y_{ik} \tanh\left(\sum_{t \neq s} \hat{\mathbf{W}}_{st} Y_{it}\right) \right] \\ &\stackrel{(i)}{=} -\frac{1}{n} \sum_{i=1}^n \left[ Y_{is} Y_{ik} - \tanh\left(\sum_{t \neq s} Y_{ik} \hat{\mathbf{W}}_{st} Y_{it}\right) \right] \\ &= -\frac{1}{n} \mathbf{Y}[s]^T \mathbf{Y}[k, ] + \frac{1}{n} \sum_{i=1}^n \tanh(Y_{ik} \mathbf{Y}_{i\{-s\}} (\mathbf{Y}_{\{-s\}}^T \mathbf{Y}_{\{-s\}})^{-1} \mathbf{Y}_{\{-s\}}^T \mathbf{Y}_s) \\ &\stackrel{(ii)}{=} \frac{1}{n} \sum_{i=1}^n (t Y_{ik} \mathbf{Y}_{i\{-s\}} (\mathbf{Y}_{\{-s\}}^T \mathbf{Y}_{\{-s\}})^{-1} \mathbf{Y}_{\{-s\}}^T \mathbf{Y}_s)^3 \\ &\rightarrow 0. \end{aligned} \quad (1.46)$$

where (i) is because  $y_{ik}$  can only take +1 or -1 and in (ii) we apply the above identity.

### 1.8.3 Proof of Theorem 1.4.1

Let  $\mathbf{W}_T = \mathbf{W} \otimes \mathbf{I}$  and  $\overline{\mathbf{W}}_T = \overline{\mathbf{W}} \otimes \mathbf{I}$ , then  $\hat{\beta}$  is the solution of the estimation problem  $U(\beta; \delta(\beta; \mathbf{W})) = \mathbf{L}^T \mathbf{W}_T [\mathbf{D}(\beta) + \phi \mathbf{W}_T (\mathbf{I} - \mathbf{D}(\beta))]^{-1} (\nabla l(\mathbf{L}\beta) + \delta(\beta; \mathbf{W})) = \mathbf{0}$ . By Taylor expansion:

$$0 = U(\hat{\beta}; \overline{\mathbf{W}}; \delta(\hat{\beta}; \overline{\mathbf{W}})) = U(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}})) + \frac{dU(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}}))}{d\beta^*} (\hat{\beta} - \beta^*) (1 + o(1)),$$

thus

$$\sqrt{n}(\hat{\beta} - \beta^*) \approx - \left[ \frac{dU(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}}))}{n d\beta^*} \right]^{-1} \cdot \frac{U(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}}))}{\sqrt{n}}$$

Let  $U_i(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}})) := \mathbf{L}^T \overline{\mathbf{W}}_T [\mathbf{D}(\beta^*) + \phi \overline{\mathbf{W}}_T (\mathbf{I} - \mathbf{D}(\beta^*))]^{-1} (\nabla \bar{l}_i(\mathbf{L}\beta^*) + \delta(\beta^*; \overline{\mathbf{W}}))$ . Based on definition (1.28),  $U_i(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}}))$  follows i.i.d distribution for  $i = 1, \dots, n$ . Note that as  $n \rightarrow \infty$ , by central limit theorem:

$$\frac{U(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}}))}{\sqrt{n}} = \sqrt{n} \left[ \frac{\sum_{i=1}^n U_i(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}}))}{n} - 0 \right] \rightarrow \mathbb{N}(\mathbf{0}, \mathbf{V}_1),$$

where

$$\mathbf{V}_1 = \lim_{n \rightarrow \infty} \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \mathbf{L}^T \overline{\mathbf{W}}_T [\mathbf{D}(\beta^*) + \phi \overline{\mathbf{W}}_T (\mathbf{I} - \mathbf{D}(\beta^*))]^{-1} \nabla \bar{l}_i((\mathbf{L}\beta^*), \mathbf{y}) \cdot \nabla \bar{l}_i((\mathbf{L}\beta^*), \mathbf{y})^T [\mathbf{D}(\beta^*) + \phi \overline{\mathbf{W}}_T (\mathbf{I} - \mathbf{D}(\beta^*))]^{-1} \overline{\mathbf{W}}_T \mathbf{L} \right]$$

$$\frac{dU(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}}))}{n d\beta^*} = \frac{\partial U(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}}))}{n \partial \beta} + \frac{\partial U(\beta^*; \delta(\beta^*; \overline{\mathbf{W}}))}{\partial \delta} \cdot \frac{\partial \delta}{\partial \beta}(\beta^*; \delta(\beta^*; \overline{\mathbf{W}})) \cdot \frac{1}{n}. \quad (1.47)$$

Since we have

$$\begin{aligned} \frac{\partial U(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}}))}{\partial \delta} &= \mathbf{L}^T \overline{\mathbf{W}}_T [\mathbf{D}(\beta^*) + \phi \overline{\mathbf{W}}_T (\mathbf{I} - \mathbf{D}(\beta^*))]^{-1}, \\ \text{and } \left\| \frac{\partial \delta}{\partial \beta}(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}})) \right\|_2 / n &\rightarrow 0, \end{aligned} \quad (1.48)$$

then by law of large numbers:

$$\frac{\partial U(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}}))}{n \partial \beta} = \frac{\sum_{i=1}^n \partial U_i(\beta^*; \overline{\mathbf{W}}; \delta(\beta^*; \overline{\mathbf{W}}))}{n \partial \beta} \rightarrow \mathbf{V}_2, \quad (1.49)$$

where

$$\begin{aligned} \mathbf{V}_2 &= \lim_{n \rightarrow \infty} \left[ \frac{1}{n} \sum_{i=1}^n \mathbf{L}^T \overline{\mathbf{W}}_T [\mathbf{D}(\beta^*) + \phi \overline{\mathbf{W}}_T (\mathbf{I} - \mathbf{D}(\beta^*))]^{-1} \nabla^2 \bar{l}_i(\mathbf{L}\beta^*) \mathbf{L} \right] \\ &= \lim_{n \rightarrow \infty} \left[ \mathbf{L}^T \overline{\mathbf{W}}_T [\mathbf{D}(\beta^*) + \phi \overline{\mathbf{W}}_T (\mathbf{I} - \mathbf{D}(\beta^*))]^{-1} \mathbf{D}(\beta^*) \mathbf{L} \right] \end{aligned} \quad (1.50)$$

To conclude, as  $n$  goes to infinity,  $\sqrt{n}(\hat{\beta} - \beta^*) \rightarrow \mathbb{N}(\mathbf{0}, \mathbf{V}_2^{-1} \mathbf{V}_1 (\mathbf{V}_2^{-1})^T)$ .

# CHAPTER 2

## SLACK EMPIRICAL LIKELIHOOD (SEL)

### 2.1 Introduction

Inference is one of the core topics in statistics, and the goal is to quantify the uncertainty of population parameter through confidence intervals and hypothesis tests. Estimating equations play an important role in hypothesis test problems. Given  $n$  observations  $X_1, \dots, X_n$  of a random variable  $X$  distributed from  $\mathbb{P}$  and a parameter of interest  $\boldsymbol{\theta}$ , many hypothesis test problems can be formulated into the following form with the assistance of estimating equations

$$H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0, \quad \text{such that} \quad \mathbb{E}(\mathbf{g}(X; \boldsymbol{\theta})) = \mathbf{0}, \quad (2.1)$$

where the expectation is taken with respect to the distribution  $\mathbb{P}$ . However, the underlying distribution  $\mathbb{P}$  could be unknown in practice, and we may not know how to adaptly choose proper distribution families for real dataset. To handle the issue, a popular way is to use the empirical distribution  $\mathbb{P}_n : P(X = X_i) = 1/n$  as a plug-in estimation of the unknown distribution. In order to gain more flexibility and adaptivity on hypothesis tests to handle data imperfection, [13] proposed the empirical likelihood method by replacing the empirical distribution with  $\mathbb{P}_n^{(w)} : P(X = X_i) = w_i$ ,  $\sum_{i=1}^n w_i = 1$ ,  $w_i \geq 0$ .

Hypothesis test problem in form (2.1) is widely assumed because it is beyond the traditional likelihood and even loss function setups. However, many modern statistical problems arising in the last several years go beyond that form. For example, consider the hypothesis test for the lasso estimator, which is a solution of minimizing sum square of residuals with  $\ell_1$  norm regularization. It is known that the estimation equation of the lasso can be characterized in the following way:

$$\boldsymbol{\beta} = \Theta_S(\boldsymbol{\beta} - \mathbf{X}^T(\mathbf{X}\boldsymbol{\beta} - \mathbf{y}); \lambda), \quad \text{where} \quad \Theta_S(t; \lambda) = \begin{cases} 0, & \text{if } |t| < \lambda; \\ (|t| - \lambda)\text{sgn}(t), & \text{if } |t| \geq \lambda \end{cases}$$

is the soft-thresholding function,  $\text{sgn}(t)$  returns the sign of  $t$  and  $\Theta_S(\cdot; \lambda)$  makes effect on all components of a vector. Due to the *many-to-one* mapping of the soft-thresholding function near

zero, there is no direct way to transform estimating equation of the lasso into the form (2.1). Hence, traditional nonparametrization for the lasso problem through estimating equations will fail.

In this paper, we consider an even more general type of formulation of hypothesis test problems than (2.1), which includes hypothesis tests for the lasso. In particular, the estimating equation for the lasso can not be written as sample additive form because of the nonlinearity of  $\Theta_S$ . To achieve the purpose, we use  $p$  additional variables to restate the optimality conditions, which are slack variables. With the help of these slack variables, the estimating equation can be written as sample additive form hence nonparametrization inference can be performed through the framework of empirical likelihood.

In this article, we propose a slack empirical likelihood framework for statistical inference. By applying our framework into problems like regressions with affine inequality constraints and high dimensional inference for regressions, we show that SEL offers a mechanism of nonparametrization and robustification for the parametric assumptions. In addition, our SEL can be applied to a large variety of penalized and constrained form of optimization problem, which could reflect elementwise sparsity, row-wise sparsity, and low-rankness of the inference targeted matrix. Both asymptotic and nonasymptotic analysis are provided for SEL. Finally, some simulations are performed to demonstrate the advantages of our approach.

The organization of this article is as follows: in section 2.2, we have a brief review the empirical likelihood methods and introduce our general framework for slack empirical likelihood; we show how to construct SEL from general optimization problem in section 2.3; in section 2.4, we apply our framework for slack empirical likelihood into some specific examples, also to show the details how our methods can be applied for traditional asymptotic inference; some explorations on how to apply SEL for high dimensional problems is offered in section 2.5; numerical evidence are provided in section 2.6.

## 2.2 Slack Empirical Likelihood

### 2.2.1 Statistic-driven EL

We first review the some classical EL. Given  $n$  observations  $\mathbf{z}_i \in \mathbb{R}^m$  and  $\boldsymbol{\mu}^\circ$ , a location of interest, [14] proposed to test  $H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}^\circ$  by  $L(\hat{w}_1, \dots, \hat{w}_n) - L(1/n, \dots, 1/n)$ , where  $L(\hat{w}_1, \dots, \hat{w}_n)$

is the optimal function value obtained by

$$\min_{\{w_i\} \in \mathbb{R}^n} -\sum \log w_i \equiv L(w_1, \dots, w_n) \text{ s.t. } \sum w_i = 1, w_i \geq 0, \sum w_i \mathbf{z}_i = \boldsymbol{\mu}^\circ.$$

From the last constraint,  $w_i$  has a probability meaning in an associated multinomial distribution  $\mathbf{z} = \mathbf{z}_i$  ( $1 \leq i \leq n$ ), thereby called the empirical likelihood method.

EL is actually statistic driven although it is widely used as a non-parametric tool for data analysis. Throughout the paper, we assume a supervised setup unless otherwise specified, which has  $n$  samples of  $p$  predictors and  $m$  responses available:  $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{S} \subset \mathbb{R}^p \times \mathbb{R}^m$ . We call  $\mathcal{S}$  the sample space and construct  $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]^T \in \mathbb{R}^{n \times p}$  and  $\mathbf{Y} = [\mathbf{y}_1 \dots \mathbf{y}_n]^T \in \mathbb{R}^{n \times m}$ . Let  $\mathbf{B}$  be determined by a set of estimating equations:

$$\sum_{i=1}^n \mathbf{T}(\mathbf{B}; \mathbf{x}_i, \mathbf{y}_i) = \mathbf{0} \tag{2.2}$$

which is assumed to be *additive in samples*. In this paper, the sample additivity is a key in defining EL in various statistics contexts. For location EL,  $\mathbf{B}$  is the location parameter matrix, but can be much more general. In statistical estimation, (2.2) can often be obtained by minimizing a certain objective—the optimal  $\mathbf{B} \in \mathbb{R}^{p \times m}$  that minimizes  $\min_{\mathbf{B}} f(\mathbf{B}; \mathbf{X}, \mathbf{Y}) \triangleq \sum_i l(\mathbf{B}; \mathbf{x}_i, \mathbf{y}_i)$  can be obtained by solving  $\nabla_{\mathbf{B}} f = 0$  or (2.2) with  $\mathbf{T}(\mathbf{B}; \mathbf{x}_i, \mathbf{y}_i) = \nabla_{\mathbf{B}} l(\mathbf{B}; \mathbf{x}_i, \mathbf{y}_i)$ , assuming  $l$  is differentiable in  $\mathbb{R}^{p \times m}$ . We call  $\mathbf{T}(\mathbf{B}; \mathbf{x}_i, \mathbf{y}_i)$  the residuals, and assume  $\mathbf{T}_i(\mathbf{B}) \in \mathcal{G} \subset \mathbb{R}^{p \times m}$ . The *residual space*  $\mathcal{G}$  is often taken to be the full Euclidian space but not necessarily so which may impose more structural properties on  $\{w_i\}$ .

With the desired residuals provided by (2.2), we can define the EL statistic for any given point  $\mathbf{B}^\circ$  in the parameter space  $\Omega \subset \mathbb{R}^{p \times m}$  by solving the following optimization problem (with  $\mathbf{w} = [w_i] \in \mathbb{R}^n$ ):

$$\min_{\mathbf{w}} -\sum_{i=1}^n \log(nw_i) \quad \text{s.t.} \quad \sum_{i=1}^n w_i = 1, w_i \geq 0 \tag{2.3}$$

$$\sum_{i=1}^n w_i \mathbf{T}(\mathbf{B}^\circ; \mathbf{x}_i, \mathbf{y}_i) = \mathbf{0}.$$

An essential component of (2.3) is the last constraint, the so-called *structural constraint*. We will denote  $\mathbf{T}(\mathbf{B}^\circ; \mathbf{x}_i, \mathbf{y}_i)$  by  $\mathbf{T}_i^\circ$  for simplicity. A deep result is that under the null hypothesis  $H_0 : \mathbf{B} = \mathbf{B}^\circ$ , twice the optimal function value obtained in (2.3) follows an asymptotic  $\chi_p^2$  distribution as  $n \rightarrow +\infty$  in a nonparametric sense under some regularity conditions [16].

The residuals  $\mathbf{T}_i^\circ$  are obtained on the  $i$ th observation using the  $\mathbf{B}^\circ$ , such as  $\mathbf{x}_i(\mathbf{x}_i^T \mathbf{B}^\circ - \mathbf{y}_i^T)$  in regression. For example, when  $\mathbf{X}$  is a column of ones and  $l(\boldsymbol{\mu}^\circ; \mathbf{z}_i) = \|\boldsymbol{\mu}^\circ - \mathbf{z}_i\|_2^2/2$ ,  $\mathbf{T}(\boldsymbol{\mu}^\circ; \mathbf{x}_i, \mathbf{z}_i) = \boldsymbol{\mu}^\circ - \mathbf{z}_i$ , and so (2.3) degenerates to the location empirical likelihood, and if  $l(\boldsymbol{\beta}^\circ; \mathbf{x}_i, y_i) = (\mathbf{x}_i^T \boldsymbol{\beta}^\circ - y_i)^2/2$ ,  $\mathbf{T}(\boldsymbol{\beta}^\circ; \mathbf{x}_i, y_i) = (\boldsymbol{\beta}^{\circ T} \mathbf{x}_i - y_i) \mathbf{x}_i^T$ , corresponding to the celebrated regression empirical likelihood [15]. In a sense, the residuals and structural constraint are rooted in a parametric model, but EL offers a mechanism of nonparametrization and robustification.

To gain more insights, we rewrite (2.3) for supervised learning with an estimation criterion of  $f(\mathbf{B}) = \bar{l}(\boldsymbol{\Theta}; \mathbf{Y}) = \sum_i l_0(\mathbf{B}^T \mathbf{x}_i; \mathbf{y}_i)$ , where a loss  $l_0 \in \mathcal{C}^1$  is placed on the *systematic component*  $\boldsymbol{\Theta} = \mathbf{X}\mathbf{B}$ . Clearly,

$$\nabla_{\mathbf{B}} f(\mathbf{B}) = \mathbf{X}^T \nabla_{\boldsymbol{\Theta}} \bar{l}(\boldsymbol{\Theta}).$$

Given  $\mathbf{B}^\circ \in \Omega$ , let  $\mathbf{R} = \nabla_{\boldsymbol{\Theta}} \bar{l}|_{\boldsymbol{\Theta}=\mathbf{X}\mathbf{B}^\circ}$ . We can also write the matrix  $\mathbf{R} = [\mathbf{r}_1, \dots, \mathbf{r}_n]^T$  with  $\mathbf{r}_i = \nabla l_0(\text{vec}(\boldsymbol{\Theta}[i, \cdot]); \mathbf{y}_i)$ . Then  $\mathbf{T}_i^\circ = \mathbf{x}_i \mathbf{r}_i^T$ . An instance is given by a vector generalized linear model (GLM) with  $m$  responses. Let  $b$  be the cumulant function and adopt the canonical link  $g = (b')^{-1}$ . Then  $l(\mathbf{B}; \mathbf{x}_i, \mathbf{y}_i)$  or  $l_0(\mathbf{B}^T \mathbf{x}_i; \mathbf{y}_i)$  is  $-\langle \mathbf{B}^T \mathbf{x}_i, \mathbf{y}_i \rangle + \langle \mathbf{1}, b(\mathbf{B}^T \mathbf{x}_i) \rangle$ . Direct calculation shows

$$\mathbf{R} = b'(\mathbf{X}\mathbf{B}^\circ) - \mathbf{Y},$$

where  $b'(\cdot)$  is applied element-wise. When  $b'(\cdot)$  is the identity function, it gives the multivariate regression empirical likelihood.

Let  $\mathcal{C} = \{w_1, \dots, w_n : \sum_{i=1}^n w_i = 1, w_i \geq 0, 1 \leq i \leq n\}$  and  $\mathbf{W} = \text{diag}\{w_i\}$ . The associated EL problem can restated in the following form

$$\min_{\mathbf{w} \in \mathcal{C}} -\langle \mathbf{1}, \log(n\mathbf{w}) \rangle \quad \text{s.t.} \quad \mathbf{X}^T \mathbf{W} \mathbf{R} = \mathbf{0}. \quad (2.4)$$

To generalize the idea of EL, the structural constraint  $\mathbf{X}^T \mathbf{W} \mathbf{R} = \mathbf{0}$  is the key.

### 2.2.2 Introducing Slack Variables

The definitions and extensions in the previous section apply well when the objective function is differentiable in the full parameter space. Modern statistical applications however pose new challenges. In particular, (a) when the parameter space  $\Omega$  is restricted, evaluating the gradient in  $\mathbb{R}^{p \times m}$  may not deliver reasonable residuals; (b) the objective function for high dimensional sparse learning is typically nonsmooth. As we will see, introducing auxiliary slack variables helps. We demonstrate our techniques with some examples.

**Example 1. (Nonnegative regression EL)** As an extension of the celebrated regression EL, we consider EL in a setting where all coefficients are nonnegative. This corresponds to the nonnegative least squares problem:  $\min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$  s.t.  $\beta_j \geq 0$  ( $1 \leq j \leq p$ ), or  $\beta \in \mathbb{R}_+^p$  where  $\mathbb{R}_+ = [0, \infty]$ . The parameter space is restricted and *closed*. Regression EL can be applied if  $\beta^\circ$  is an interior point, but if  $\beta^\circ$  lies on the boundary, i.e.,  $\beta_j^\circ = 0$  for some  $j$ , which is of practical interest in significance tests, regular EL does not apply, and the normal-equation-based residuals must be corrected.

**Example 2. (EL with affine inequality constraints)** Consider test with constraints  $\mathbf{A}\beta \preceq \alpha$  with  $\mathbf{A} \in \mathbb{R}^{k \times p}$  ( $k$  can be smaller than or greater than  $p$ ). If the constraints contain equalities  $\mathbf{A}_1\beta = \alpha_1$ , they can be converted to inequality constraints  $\mathbf{A}_1\beta \preceq \alpha_1$  and  $-\mathbf{A}_1\beta \preceq \alpha_1$ , so we can just consider a general case  $\mathbf{A}\beta \preceq \alpha$ . The constraints should be compatible such that the true value  $\beta^\circ$  is indeed in the feasible region of these constraints. This corresponds to the least squares problem with constraints:  $\min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$  s.t.  $\mathbf{A}\beta \preceq \alpha$ . Still The parameter space is restricted and closed. And if  $\beta^\circ$  is an interior point, regression EL can still be applied.

**Example 3. (Large- $p$  EL)** Consider a sparse learning problem  $\min_{\beta} \sum_{i=1}^n l_0(\mathbf{x}_i^T \beta) + \sum_{j=1}^p P(|\beta_j|; \lambda)$ , where  $l_0$  is a loss function defined on the systematic component  $\mathbf{x}_i^T \beta$  and  $P$  is a penalty function to promote the sparsity in  $\beta$ . Examples of  $P$  include the  $\ell_1$  penalty,  $\ell_0$ -penalty and SCAD [4], which are all popular for building a high-dimensional linear model. We assume that the regularization parameter  $\lambda$  is *given* (either by theory—see, e.g., [2], or by tuning such as cross-validation) so that the criterion is fully specified. A new class of empirical likelihoods say  $L_1$ -EL or  $L_0$ -EL would be useful for high-dimensional inference, but the non differentiability along with possible nonconvexity makes it difficult to obtain sample-additive EEs.

Example 1 is tricky because the closed restricted parameter space  $\mathbb{R}_+^p$  makes the gradient-derived residuals improper for any target point on its boundary. Example 3 shows a nondifferentiability issue caused by the objective function. The good news is that we can employ the same trick of slack variables to rewrite the optimality conditions to extend EL, which we call slack empirical likelihood (**SEL**).

Let  $\bar{l}(\mathbf{X}\beta; \mathbf{y}) = \sum_i l_0(\mathbf{x}_i^T \beta; \mathbf{y}_i)$  and recall the sparse estimation problem

$$\min_{\beta} f(\beta) \triangleq \bar{l}(\mathbf{X}\beta; \mathbf{y}) + \sum_{j=1}^p P(|\beta_j|; \lambda). \quad (2.5)$$



For simplicity, we assume  $\bar{l}$  to be differentiable, but the technique applies to any objective function that is only directionally differentiable. Apart from convex penalties, nonconvex penalties have recently received much attention, due to their capability of accommodating coherent designs and reducing estimation bias. All of the sparsity-promoting penalties are associated with **thresholding** rules. To make the connection more clear, we use the framework of  $\Theta$ -estimators [21] to study EL driven by an *arbitrary* thresholding function  $\Theta$ . A thresholding function  $\Theta(t; \lambda)$  is defined for  $-\infty < t < \infty$  and  $0 \leq \lambda < \infty$  such that (i)  $\Theta(-t; \lambda) = -\Theta(t; \lambda)$ ; (ii)  $\Theta(t; \lambda) \leq \Theta(t'; \lambda)$  for  $t \leq t'$ ; (iii)  $\lim_{t \rightarrow \infty} \Theta(t; \lambda) = \infty$ ; (iv)  $0 \leq \Theta(t; \lambda) \leq \infty$  for  $0 \leq t < \infty$ . Given any  $u \geq 0$ ,  $\Theta^{-1}(u; \lambda) \triangleq \sup\{t : \Theta(t; \lambda) \leq u\}$ , and so the threshold is  $\tau(\lambda) \triangleq \Theta^{-1}(0; \lambda)$ . For simplicity, we always assume that  $\lambda$  is the threshold parameter, i.e.,  $\tau(\lambda) = \lambda$ . A vector version of  $\Theta$  is defined component-wise if either  $t$  or  $\lambda$  is replaced by a vector. Given any thresholding  $\Theta$ , we say  $P$  is induced by  $\Theta$  if

$$P(t; \lambda) = P_{\Theta}(t; \lambda) + q(t; \lambda) = \int_0^{|t|} (\Theta^{-1}(u; \lambda) - u) du + q(t; \lambda), \quad (2.6)$$

where  $q$  is an arbitrary nonnegative function satisfying  $q(t; \lambda) = 0$  if  $t = \Theta(s; \lambda)$  for some  $s \in \mathbb{R}$ . The mapping from  $P$  to  $\Theta$  is *many-to-one* if  $\Theta$  has discontinuities. The  $\Theta$ - $P$  framework covers all aforementioned penalties: ridge scaling  $\Theta_R(t; \eta) = t/(1 + \eta)$  corresponds to the  $\ell_2$ -penalty  $P(t; \eta) = \eta t^2/2$ , soft-thresholding  $\Theta_S(t; \lambda) = \text{sgn}(t)(|t| - \lambda)_+$  gives the  $\ell_1$ -penalty  $P(t; \lambda) = \lambda|t|$ , and hard-thresholding  $\Theta_H(t; \lambda) = t1_{|t| \geq \lambda}$  induces infinitely many  $\ell_0$ -type penalties including the discrete  $\ell_0$  penalty and the capped  $\ell_1$  penalty; see [21] for other examples of elastic net, SCAD, MCP,  $\ell_r$  ( $0 \leq r < 1$ ), and so on. Given (2.6), a class of iterative thresholding procedures can be used to solve (2.5), the resulting estimates called  $\Theta$ -estimates.

The key to defining EL lies in the characterization a ‘fit’ or  $\Theta$ -estimate. Let  $f_{\Theta}(\boldsymbol{\beta}) = \bar{l}(\mathbf{X}\boldsymbol{\beta}; \mathbf{y}) + \sum_{j=1}^p P_{\Theta}(|\beta_j|; \lambda)$ . From [21], any fixed point or locally optimal solution  $\hat{\boldsymbol{\beta}}$  of the iterative thresholding procedure for solving (2.5) satisfies  $D_{\pm \mathbf{e}_j} f_{\Theta}(\boldsymbol{\beta}) \geq 0$  for all  $1 \leq j \leq p$ , where  $\mathbf{e}_j$  is the vector with the  $j$ th component 1 and the other components 0, and  $D_{\mathbf{u}} f_{\Theta}(\boldsymbol{\beta})$  denotes the directional derivative of  $f_{\Theta}$  at  $\boldsymbol{\beta}$  with increment  $\mathbf{u}$ , i.e.,  $D_{\mathbf{u}} f_{\Theta}(\boldsymbol{\beta}) = \lim_{\epsilon \rightarrow 0^+} [f_{\Theta}(\boldsymbol{\beta} + \epsilon \mathbf{u}) - f_{\Theta}(\boldsymbol{\beta})]/\epsilon$ . These conditions do lead to a set of  $\Theta$ -equations [21, Theorem 1]:

$$\hat{\boldsymbol{\beta}} = \Theta(\hat{\boldsymbol{\beta}} - \mathbf{X}^T \nabla \bar{l}(\mathbf{X}\hat{\boldsymbol{\beta}}); \lambda), \quad (2.7)$$

under the mild assumption that  $\Theta(\cdot; \lambda)$  is continuous at  $\hat{\beta} - \mathbf{X}^T \nabla \bar{l}(\mathbf{X} \hat{\beta})$ . Moreover, all  $\Theta$ -estimates enjoy provable statistical guarantees.

However, unless  $\Theta$  is the ridge scaling, (2.7) is **not** sample additive due to the nonlinearity of  $\Theta$ . To achieve the purpose, we use  $p$  additional variables (and inequalities) to restate the optimality conditions. Given  $\beta^\circ \in \mathbb{R}^p$ , define  $\mathcal{J} = \{j : \beta_j^\circ \neq 0\}$  and  $\mathcal{J}^c = \{j : \beta_j^\circ = 0\}$ . For any  $j \in \mathcal{J}$ , a simple derivative calculation shows that  $\Theta^{-1}(|\beta_j^\circ|; \lambda) \text{sgn}(\beta_j^\circ) = \beta_j^\circ - \mathbf{X}[j]^T \nabla \bar{l}(\mathbf{X} \beta^\circ)$ , and for  $j \in \mathcal{J}^c$ ,  $\beta_j^\circ = 0$ , and so  $-\lambda \leq \mathbf{X}[j]^T \nabla \bar{l}(\mathbf{X} \beta^\circ) \leq \lambda$  by definition. Let

$$\gamma_j^\circ = \begin{cases} \Theta^{-1}(|\beta_j^\circ|; \lambda) \text{sgn}(\beta_j^\circ) - \beta_j^\circ, & \text{if } j \in \mathcal{J} \\ 0, & \text{if } j \in \mathcal{J}^c. \end{cases} \quad (2.8)$$

(2.7) is now replaced by  $\mathbf{X}^T \nabla \bar{l}(\mathbf{X} \beta^\circ) + \gamma^\circ + \mathbf{s} = \mathbf{0}$ , which is sample additive, where  $\mathbf{s} \in \mathbb{R}^p$ ,  $\mathbf{s}_{\mathcal{J}} = \mathbf{0}$ , and  $|s_j| \leq \lambda$ ,  $j \in \mathcal{J}^c$ . Correspondingly, given any  $\Theta(\cdot; \lambda)$  and  $\beta^\circ \in \mathbb{R}^p$ , we cast SEL as a **joint** optimization with respect to the weights  $\mathbf{w}$  and the slack variables  $\mathbf{s} = [s_j]$ :

$$\begin{aligned} \min_{\mathbf{w} \in \mathcal{C}, \mathbf{s} \in \mathbb{R}^p} -\langle \mathbf{1}, \log(\mathbf{w}) \rangle \quad \text{s.t.} \quad & \mathbf{X}^T \mathbf{W} \nabla \bar{l}(\mathbf{X} \beta^\circ) + \gamma^\circ + \mathbf{s} = \mathbf{0} \\ & \mathbf{s} \circ \beta^\circ = \mathbf{0}, \|\mathbf{s}\|_\infty \leq \lambda, \end{aligned} \quad (2.9)$$

where  $\circ$  denotes the elementwise product and  $\|\mathbf{s}\|_\infty = \max |s_j|$ . If we define  $\mathcal{Z} = \mathcal{J}^c$ , then  $\mathbf{s}$  satisfies  $-\lambda \leq s_j \leq \lambda$  for  $j \in \mathcal{Z}$  and  $s_j = 0$  for  $j \in \mathcal{J}$ .

A useful version is the sparsity EL under an  $\ell_0$ -constraint  $\|\beta\|_0 \leq q$ . From [19], all locally optimal solutions to  $\min_{\|\beta\|_0 \leq q} \bar{l}(\mathbf{X} \beta; \mathbf{y})$  satisfy  $\hat{\beta} = \Theta^\#(\hat{\beta} - \mathbf{X}^T \nabla \bar{l}(\mathbf{X} \hat{\beta}); q)$ , assuming no ties occur and  $\|\mathbf{X}\|_2 \leq 1/L$  with  $L$  the Lipschitz constant of  $\nabla \bar{l}$ . Here, the quantile thresholding  $\Theta^\#(\boldsymbol{\alpha}; q)$  for any  $\boldsymbol{\alpha} \in \mathbb{R}^p$  is defined to be a vector  $\boldsymbol{\gamma} \in \mathbb{R}^p$  satisfying  $\gamma_{(j)} = \alpha_{(j)}$  if  $1 \leq j \leq q$ , and 0 otherwise, where  $\alpha_{(1)}, \dots, \alpha_{(p)}$  are the order statistics of  $\alpha_1, \dots, \alpha_p$  satisfying  $|\alpha_{(1)}| \geq \dots \geq |\alpha_{(p)}|$ . Rewriting the  $\Theta^\#$ -equation by use of slack variables, we get the SEL problem for any given  $\|\beta^\circ\|_0 \leq q$ :

$$\begin{aligned} \min_{\mathbf{w} \in \mathcal{C}, \mathbf{s}} -\langle \mathbf{1}, \log(\mathbf{w}) \rangle \quad \text{s.t.} \quad & \mathbf{X}^T \mathbf{W} \nabla \bar{l}(\mathbf{X} \beta^\circ) + \mathbf{s} = \mathbf{0} \\ & \mathbf{s} \circ \beta^\circ = \mathbf{0}, \|\mathbf{s}\|_\infty < |\beta_{(q)}^\circ|, \end{aligned} \quad (2.10)$$

where  $|\beta_{(q)}^\circ|$  is the  $q$ th largest element in  $|\beta_j^\circ|$ ,  $1 \leq j \leq p$ . Note that all  $\mathbf{s}$ -constraints in SEL, either in equalities or inequalities, are **linear**.

Nonnegative regression EL can be defined similarly. Consider  $\min_{\beta \succeq \mathbf{0}} \sum_i l_0(\beta; \mathbf{x}_i, y_i) = f(\beta)$ , where  $l$  is differentiable in the augmented parameter space  $\mathbb{R}^p$  but not necessarily convex. Any

optimal solution  $\beta$  obeys  $D_{\mathbf{u}}f(\beta) \geq 0$  for all feasible  $\mathbf{u}$ . Specifically, taking  $\mathbf{u} = \pm \mathbf{e}_j$  for  $j \in \mathcal{J}(\beta) = \{j : \beta_j \neq 0\}$  and  $\mathbf{u} = \mathbf{e}_j$  for  $j \in \mathcal{J}^c(\beta)$  leads to the estimating equations with slack variables:

$$\sum_i \{\nabla l(\beta; \mathbf{x}_i, y_i) + (1/n)\mathbf{s}\} = \mathbf{0},$$

where  $\mathbf{s}_{\mathcal{J}^c} \succeq \mathbf{0}$  and  $\mathbf{s}_{\mathcal{J}} = \mathbf{0}$ . Given  $\beta^\circ$ , adding  $w_i$  in the first term and replacing  $1/n$  by  $1/w_i$  gives the SEL test statistic:  $\min_{\mathbf{w} \in \mathcal{C}, \mathbf{s}} -\langle \mathbf{1}, \log(n\mathbf{w}) \rangle$  s.t.  $\mathbf{X}^T \mathbf{W} \nabla \bar{l}(\mathbf{X}\beta^\circ) + \mathbf{s} = \mathbf{0}$ ,  $\mathbf{s} \circ \beta^\circ = \mathbf{0}$ ,  $\mathbf{s} \succeq \mathbf{0}$ . When  $l(\beta; \mathbf{x}_i, y_i) = (\mathbf{x}_i^T \beta - y_i)^2/2$ , we get

$$\begin{aligned} \min_{\mathbf{w} \in \mathcal{C}, \mathbf{s}} -\langle \mathbf{1}, \log(n\mathbf{w}) \rangle \quad \text{s.t.} \quad & \mathbf{X}^T \mathbf{W} \nabla \bar{l}(\mathbf{X}\beta^\circ) + \mathbf{s} = \mathbf{0} \\ & \mathbf{s} \circ \beta^\circ = \mathbf{0}, \mathbf{s} \succeq \mathbf{0}. \end{aligned} \quad (2.11)$$

If  $\beta^\circ \succ \mathbf{0}$ , then  $\mathbf{s} = \mathbf{0}$  and (2.11) reduces to the ordinary regression EL.

For EL with affine inequality constraints, Consider  $\min_{\mathbf{A}\beta \preceq \mathbf{a}} \sum_i l(\beta; \mathbf{x}_i, y_i) = f(\beta)$ , now when we consider the optimal solution  $\beta$  for all feasible  $\mathbf{u}$ : when  $\mathbf{A}[i, ]\beta < a_i$ ,  $\mathbf{u} = \pm \mathbf{A}[i, ]/\|\mathbf{A}[i, ]\|_{\ell_2}$ , and if  $\mathbf{A}[i, ]\beta = a_i$ ,  $\mathbf{u} = -\mathbf{A}[i, ]/\|\mathbf{A}[i, ]\|_{\ell_2}$ . Now let  $\mathcal{J} = \{j = 1, 2, \dots, k \mid \mathbf{A}[j, ]\beta < a_j\}$  and  $\mathcal{J}^c = \{j = 1, 2, \dots, k \mid \mathbf{A}[j, ]\beta = a_j\}$ . Then given  $\beta^\circ$  in the feasible region, suppose the loss is  $\ell_2$  type, we can obtain the SEL statistics:

$$\begin{aligned} \min_{\mathbf{w} \in \mathcal{C}, \mathbf{s}} -\langle \mathbf{1}, \log(n\mathbf{w}) \rangle \quad \text{s.t.} \quad & \mathbf{X}^T \mathbf{W} \nabla \bar{l}(\mathbf{X}\beta^\circ) + \mathbf{s} = \mathbf{0} \\ & \mathbf{A}[\mathcal{J}, ]\mathbf{s} = \mathbf{0}, \mathbf{A}[\mathcal{J}^c, ]\mathbf{s} \preceq \mathbf{0}. \end{aligned} \quad (2.12)$$

## 2.3 Constructing SEL from General Optimizations

Consider the following optimization problem:

$$\min_{\mathbf{B}} f(\mathbf{B}) := \min_{\mathbf{B}} \bar{l}(\mathbf{X}\mathbf{B}; \mathbf{Y}) + P(\mathbf{B}; \lambda), \quad (2.13)$$

where  $\bar{l}(\mathbf{X}\mathbf{B}; \mathbf{Y}) = \sum_i l_0(\mathbf{x}_i^T \mathbf{B}; \mathbf{y}_i)$  with a well-defined differentiable loss function  $l_0$ , which is assumed to be  $L$ -Lipschitz continuous.  $\mathbf{Y} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{X} \in \mathbb{R}^{n \times p}$  and  $\mathbf{B} \in \mathbb{R}^{p \times m}$ . Here  $P(\mathbf{B}; \lambda)$  is general, which can be element-wise sparsity, group-wise sparsity and low-rankness inducing penalties. For ease of introduction, we first study SEL construction with a special family of penalties in section, then we extend our analysis to SEL with general penalties or constraints.

### 2.3.1 SEL for $\Theta$ Estimators

To start with, we first study SEL constructions from optimizations with a family of regularizations  $P_\Theta$ , which are induced by some thresholding functions  $\Theta(t; \lambda)$  defined as follows:

*Definition 4.* (Thresholding function). A thresholding function is a real valued function  $\Theta(t; \lambda)$  defined for  $-\infty < t < \infty$  and  $0 \leq \lambda < \infty$  such that (i)  $\Theta(-t; \lambda) = -\Theta(t; \lambda)$ ; (ii)  $\Theta(t; \lambda) \leq \Theta(t'; \lambda)$  for  $t \leq t'$ ; (iii)  $\lim_{t \rightarrow \infty} \Theta(t; \lambda) = \infty$ ; (iv)  $0 \leq \Theta(t; \lambda) \leq t$  for  $0 \leq t < \infty$ .

We always assume  $\lambda$  as the threshold parameter such that  $\lambda = \Theta^{-1}(0; \lambda)$ . This family of penalties  $P_\Theta$  can be constructed from  $\Theta$  as follows:

$$P_\Theta(t; \lambda) = \int_0^{|t|} (\Theta^{-1}(u; \lambda) - u) du = \int_0^{|t|} (\sup\{s : \Theta(s; \lambda) \leq u\} - u) du. \quad (2.14)$$

As shown in Theorem 1 in [21], the local minimum points  $\hat{\mathbf{B}}$  for optimization (2.13) can be characterized by the following  $\Theta$  equation when penalty  $P$  belongs to the family of  $P_\Theta$ :

$$\hat{\mathbf{B}} = \Theta(\hat{\mathbf{B}} - \mathbf{X}^T \nabla \bar{l}(\mathbf{X} \hat{\mathbf{B}}); \lambda), \quad (2.15)$$

under mild condition such that  $\Theta(\cdot; \lambda)$  is continuous at  $\hat{\mathbf{B}} - \mathbf{X}^T \nabla \bar{l}(\mathbf{X} \hat{\mathbf{B}})$ .

**(Element-wise thresholdings)** Consider the case when the penalty function  $P_\Theta(\cdot; \lambda)$  for a matrix  $\mathbf{B}$  is define as follows:

$$P_\Theta(\mathbf{B}; \lambda) = \sum_{i=1}^p \sum_{j=1}^m P(B[i, j]; \lambda), \quad (2.16)$$

then the corresponding  $\Theta(\mathbf{B}; \lambda)$  is a function on all components of  $\mathbf{B}$ . Given  $\mathbf{B}^\circ \in \mathbb{R}^{p \times m}$ , define  $\mathcal{J} = \{(i, j) : B_{i,j}^\circ \neq 0\}$  and  $\mathcal{J}^c = \{(i, j) : B_{i,j}^\circ = 0\}$ . If  $(i, j) \in \mathcal{J}$ , by inverting the  $\Theta$  function, we have  $\Theta^{-1}(|B_{i,j}^\circ|; \lambda) \text{sgn}(B_{i,j}^\circ) = B_{i,j}^\circ - [\mathbf{X}^T \nabla \bar{l}(\mathbf{X} \mathbf{B}^\circ)]_{i,j}$ ; If  $(i, j) \in \mathcal{J}^c$ , we have  $B_{i,j}^\circ = 0$ , which implies  $|\mathbf{X}^T \nabla \bar{l}(\mathbf{X} \mathbf{B}^\circ)|_{i,j} \leq \lambda$ . Define a matrix  $\mathbf{\Gamma}^\circ$  such that

$$\Gamma_{i,j}^\circ = \begin{cases} \Theta^{-1}(|B_{i,j}^\circ|; \lambda) \text{sgn}(B_{i,j}^\circ) - B_{i,j}^\circ & \text{if } (i, j) \in \mathcal{J}, \\ 0 & \text{if } (i, j) \in \mathcal{J}^c. \end{cases} \quad (2.17)$$

Then by introducing the slack variable  $\mathbf{S}$  such that  $\mathbf{X}^T \nabla \bar{l}(\mathbf{X} \mathbf{B}^\circ) + \mathbf{\Gamma}^\circ + \mathbf{S} = \mathbf{0}$ , we have constraints  $\mathbf{S} \circ \mathbf{B}^\circ = \mathbf{0}$  and  $\|\mathbf{S}\|_{\max} \leq \lambda$ . Finally, the corresponding SEL is:

$$\begin{aligned} \min_{w \in \mathcal{C}, \mathbf{S} \in \mathbb{R}^{p \times m}} -\langle \mathbf{1}, \log(\mathbf{n} \mathbf{w}) \rangle \quad \text{s.t.} \quad & \mathbf{X}^T \mathbf{W} \nabla \bar{l}(\mathbf{X} \mathbf{B}^\circ) + \mathbf{\Gamma}^\circ + \mathbf{S} = \mathbf{0} \\ & \mathbf{S} \circ \mathbf{B}^\circ = \mathbf{0}, \quad \|\mathbf{S}\|_{\max} \leq \lambda, \end{aligned} \quad (2.18)$$

where  $\mathcal{C} = \{w_1, \dots, w_n : \sum_{i=1}^n w_i = 1, w_i \geq 0, 1 \leq i \leq n\}$  and  $\mathbf{W} = \text{diag}\{w_i\}$ .

**(Group-wise thresholdings)** Sometimes it is more interesting to consider group-wise sparsity structure on coefficient  $\mathbf{B}$ . Let  $P_{2,\Theta}$  be a penalty function for a matrix that promotes sparsity in the following way

$$P_{2,\Theta}(\mathbf{B}; \lambda) = \sum_{i=1}^p P_{\Theta}(\|\mathbf{B}[i,]\|_2; \lambda). \quad (2.19)$$

As discussed in [20], we define multivariate threshold function as follows:

*Definition 5.* (Multivariate Threshold function). Given any  $\Theta, \vec{\Theta}$  is defined for any vector  $\boldsymbol{\alpha} \in \mathbb{R}^m$  such that  $\vec{\Theta}(\boldsymbol{\alpha}; \lambda) = \boldsymbol{\alpha}\Theta(\|\boldsymbol{\alpha}\|_2; \lambda) / \|\boldsymbol{\alpha}\|_2$  for  $\boldsymbol{\alpha} \neq \mathbf{0}$  and  $\mathbf{0}$  otherwise. For any matrix  $\mathbf{A} = (\boldsymbol{\alpha}_1 \dots \boldsymbol{\alpha}_n)^T \in \mathbb{R}^{n \times m}$ ,  $\vec{\Theta}(\mathbf{A}; \lambda) = \{\vec{\Theta}(\boldsymbol{\alpha}_1; \lambda) \dots \vec{\Theta}(\boldsymbol{\alpha}_n; \lambda)\}^T$ .

Then the corresponding  $\Theta$  equation can be written as  $\hat{\mathbf{B}} = \vec{\Theta}(\hat{\mathbf{B}} - \mathbf{X}^T \nabla \bar{l}(\mathbf{X}\hat{\mathbf{B}}); \lambda)$ . Given  $\mathbf{B}^\circ \in \mathbb{R}^{p \times m}$ , define  $\mathcal{J} = \{i : \mathbf{B}^\circ[i,] \neq \mathbf{0}\}$  and  $\mathcal{J}^c = \{i : \mathbf{B}^\circ[i,] = \mathbf{0}\}$ . If  $i \in \mathcal{J}$ , by inverting the  $\vec{\Theta}$  function, we have  $\mathbf{B}^\circ[i,] \Theta^{-1}(\|\mathbf{B}^\circ[i,]\|_2; \lambda) / \|\mathbf{B}^\circ[i,]\|_2 = \mathbf{B}^\circ[i,] - [\mathbf{X}^T \nabla \bar{l}(\mathbf{X}\mathbf{B}^\circ)][i,]$ ; If  $i \in \mathcal{J}^c$ , we have  $\mathbf{B}^\circ[i,] = \mathbf{0}$ , which implies  $\|[\mathbf{X}^T \nabla \bar{l}(\mathbf{X}\mathbf{B}^\circ)][i,]\|_2 \leq \lambda$ . Define a matrix  $\Gamma^\circ$  such that

$$\Gamma^\circ[i,] = \begin{cases} \mathbf{B}^\circ[i,] \Theta^{-1}(\|\mathbf{B}^\circ[i,]\|_2; \lambda) / \|\mathbf{B}^\circ[i,]\|_2 - \mathbf{B}^\circ[i,] & \text{if } i \in \mathcal{J}, \\ \mathbf{0} & \text{if } i \in \mathcal{J}^c. \end{cases} \quad (2.20)$$

Then by introducing a slack variable  $\mathbf{S}$  such that  $\mathbf{X}^T \nabla \bar{l}(\mathbf{X}\mathbf{B}^\circ) + \Gamma^\circ + \mathbf{S} = \mathbf{0}$ , we have constraints  $\mathbf{S}[\mathcal{J},] = \mathbf{0}$  and  $\|\mathbf{S}\|_{2,\infty} \leq \lambda$ . Finally, the corresponding SEL is:

$$\begin{aligned} \min_{w \in \mathcal{C}, \mathbf{S} \in \mathbb{R}^{p \times m}} -\langle \mathbf{1}, \log(\mathbf{n}w) \rangle \quad \text{s.t.} \quad & \mathbf{X}^T \mathbf{W} \nabla \bar{l}(\mathbf{X}\mathbf{B}^\circ) + \Gamma^\circ + \mathbf{S} = \mathbf{0} \\ & \mathbf{S}[\mathcal{J},] = \mathbf{0}, \quad \|\mathbf{S}\|_{2,\infty} \leq \lambda. \end{aligned} \quad (2.21)$$

**(Thresholdings on singular values)** Consider the following low-rankness promoting penalty induced by  $P_{\Theta}$ :

$$P_{\sigma,\Theta}(\mathbf{B}; \lambda) = \sum_{i=1}^{p \wedge m} P_{\Theta}(\sigma_i^{\mathbf{B}}; \lambda), \quad (2.22)$$

where  $\sigma_i^{\mathbf{B}}, i = 1, \dots, p \wedge m$  are the singular values of  $\mathbf{B}$  ordered in the descending order. As discussed in [22], the corresponding  $\Theta$  function defined on a matrix is:

*Definition 6.* (Matrix threshold function). Given any threshold function  $\Theta(\cdot; \lambda)$ , its matrix version  $\Theta^\sigma$  is defined for  $\mathbf{B} \in \mathbb{R}^{n \times m}$  as follows

$$\Theta^\sigma(\mathbf{B}; \lambda) = \mathbf{U} \text{diag}\{\Theta(\sigma_i^{\mathbf{B}}; \lambda)\} \mathbf{V}^T$$

where  $\mathbf{U}, \mathbf{V}$ , and  $\sigma_i^{\mathbf{B}}$  are obtained from the complete SVD of  $\mathbf{B} : \mathbf{B} = \mathbf{U} \text{diag}(\sigma_i^{\mathbf{B}}) \mathbf{V}^T$ .

Based on the analysis in [22], the corresponding  $\Theta$ -equation can be written as  $\mathbf{B} = \Theta^\sigma(\mathbf{B} - \mathbf{X}^T \nabla \bar{l}(\mathbf{X}\mathbf{B}; \mathbf{Y}); \lambda)$ . Given  $\mathbf{B}^\circ \in \mathbb{R}^{p \times m}$ , let  $\mathbf{B}^\circ = \mathbf{U}_{p \times p}^\circ \mathbf{D}_{p \times m}^\circ \mathbf{V}_{m \times m}^{\circ T}$  is the complete form SVD, where the singular values of  $\mathbf{B}^\circ$  are ordered in the descending order on the diagonal part of  $\mathbf{D}^\circ$ . Let  $\mathbf{U}^\circ = [\mathbf{U}^{\circ\circ}, \mathbf{U}_\perp^{\circ\circ}]$  and  $\mathbf{V}^\circ = [\mathbf{V}^{\circ\circ}, \mathbf{V}_\perp^{\circ\circ}]$  such that  $\mathbf{B}^\circ = \mathbf{U}_{p \times r^\circ}^{\circ\circ} \mathbf{D}[1 : r^\circ; 1 : r^\circ]^\circ \mathbf{V}_{m \times r^\circ}^{\circ\circ T}$  is the compact SVD, where  $r^\circ$  be the rank of  $\mathbf{B}^\circ$ , then we have  $\mathbf{D}^\circ[i, i] \neq \mathbf{0}$  if  $1 \leq i \leq r^\circ$  while  $\mathbf{D}^\circ[i, i] = \mathbf{0}$  if  $r^\circ + 1 \leq i \leq p \wedge m$ . An important equation can be obtained by plugging  $\mathbf{B}^\circ$  into the  $\Theta$ -equation:

$$\begin{aligned}
\mathbf{D}^\circ &= \mathbf{U}^{\circ T} \Theta^\sigma(\mathbf{B}^\circ - \mathbf{X}^T \nabla \bar{l}(\mathbf{X}\mathbf{B}^\circ); \lambda) \mathbf{V}^\circ \\
&\stackrel{(i)}{=} \mathbf{U}^{\circ T} \mathbf{U}' \Theta^\sigma(\mathbf{D}'; \lambda) \mathbf{V}'^T \mathbf{V}^\circ \\
&\stackrel{(ii)}{=} \mathbf{U}^{\circ\circ T} \mathbf{U}'' \Theta^\sigma(\mathbf{D}'; \lambda) \mathbf{V}''^T \mathbf{V}^{\circ\circ} \\
&\stackrel{(iii)}{=} \Theta^\sigma(\mathbf{D}^\circ - \mathbf{U}^{\circ\circ T} \mathbf{X}^T \nabla \bar{l}(\mathbf{X}\mathbf{B}^\circ)[1 : r^\circ, 1 : r^\circ] \mathbf{V}^{\circ\circ}; \lambda),
\end{aligned} \tag{2.23}$$

where  $\mathbf{B}^\circ - \mathbf{X}^T \nabla \bar{l}(\mathbf{X}\mathbf{B}^\circ) = \mathbf{U}' \mathbf{D}' \mathbf{V}'^T$  is the complete SVD and  $\mathbf{U}' = [\mathbf{U}'', \mathbf{U}'_\perp]'$  and  $\mathbf{V}' = [\mathbf{V}'', \mathbf{V}'_\perp]'$  are defined similarly. (i) is based on the definition of  $\Theta^\sigma$  function, (ii) holds because the lower diagonal part of  $\Theta^\sigma(\mathbf{D}'; \lambda)$  is zero matrix, and (iii) is based on the uniqueness of compact SVD. In fact, the corresponding columns of  $\mathbf{U}''$  with  $\mathbf{U}^{\circ\circ}$  and  $\mathbf{V}''$  with  $\mathbf{V}^{\circ\circ}$  must be the same up to any rotation/reflection applied to both sets of columns for any repeated singular values.

Based on equation (2.23), if we have  $1 \leq i \leq r^\circ$ , which implies  $\mathbf{D}^\circ[i, i] \neq \mathbf{0}$ , we can inverse the  $\Theta^\sigma$  function as before to obtain  $\Theta^{\sigma^{-1}}(\mathbf{D}^\circ[1 : r^\circ, 1 : r^\circ]; \lambda) = \mathbf{D}^\circ[1 : r^\circ, 1 : r^\circ] - \text{diag}(\sigma_i^{\mathbf{X}^T \nabla \bar{l}(\mathbf{X}\mathbf{B}^\circ)})[1 : r^\circ, 1 : r^\circ]$ . In addition, when we have  $r^\circ + 1 \leq i \leq p \wedge m$ , which implies  $\mathbf{D}^\circ[i, i] = \mathbf{0}$ , the inequality constraint can be obtained:

$$\sigma_{r^\circ+1}(\mathbf{X}^T \nabla \bar{l}(\mathbf{X}\mathbf{B}^\circ)) \leq \lambda. \tag{2.24}$$

Define a matrix  $\mathbf{\Gamma}^\circ$  such that

$$\mathbf{\Gamma}^\circ = \mathbf{U}^{\circ\circ} (\Theta^{\sigma^{-1}}(\mathbf{D}^\circ[1 : r^\circ, 1 : r^\circ]; \lambda) - \mathbf{D}^\circ[1 : r^\circ, 1 : r^\circ]) \mathbf{V}^{\circ\circ T}. \tag{2.25}$$

We can construct a slack variable  $\mathbf{S}$  such that  $\mathbf{X}^T \nabla \bar{l}(\mathbf{X}\mathbf{B}^\circ) + \mathbf{\Gamma}^\circ + \mathbf{S} = \mathbf{0}$ , then the corresponding SEL is:

$$\begin{aligned}
&\min_{w \in \mathcal{C}, \mathbf{S} \in \mathbb{R}^{p \times m}} -\langle \mathbf{1}, \log(\mathbf{n}w) \rangle \quad \text{s.t.} \quad \mathbf{X}^T \mathbf{W} \nabla \bar{l}(\mathbf{X}\mathbf{B}^\circ) + \mathbf{\Gamma}^\circ + \mathbf{S} = \mathbf{0} \\
&\mathbf{U}^{\circ\circ T} \mathbf{S} = \mathbf{0}, \quad \mathbf{S} \mathbf{V}^{\circ\circ} = \mathbf{0}, \quad \|\mathbf{S}\|_2 \leq \lambda.
\end{aligned} \tag{2.26}$$

**Remark 2.3.1.** When the penalty  $P_\Theta$  is  $\ell_1$  norm, then  $\mathbf{\Gamma}^\circ = \lambda \mathbf{U}^{\circ\circ} \mathbf{V}^{\circ\circ T}$ ; while when the penalty is  $P_H$ , then  $\mathbf{\Gamma}^\circ = \mathbf{0}$ .

### 2.3.2 SEL for General Penalties or Constraints

This section we generalize the above methodology into optimization with general penalties and even constraints. First, as discussed in [21] and [22], when a thresholding function  $\Theta(t; \lambda)$  is given, then the corresponding penalties can be constructed as:

$$P(t; \lambda) = P_{\Theta}(t; \lambda) + q(t; \lambda),$$

where  $q$  is an arbitrary function satisfying  $q(t, \lambda) \geq 0, \forall t \in \mathbb{R}$  and  $q(t; \lambda) = 0$  if  $t = \Theta(s; \lambda)$  for some  $s \in \mathbb{R}$ . Due to the nonconvex nature of the optimization problem, it is useful to construct the following surrogate function:

$$g(\mathbf{B}; \mathbf{B}^-) = \bar{l}(\mathbf{X}\mathbf{B}^-) + \langle \nabla \bar{l}(\mathbf{X}\mathbf{B}^-), \mathbf{B} - \mathbf{B}^- \rangle + \rho \|\mathbf{B} - \mathbf{B}^-\|_F^2 + P(\mathbf{B}; \lambda). \quad (2.27)$$

In addition, as long as a step size is proper chosen  $\rho > \sqrt{L} \|\mathbf{X}\|_2$ , if we set  $\mathbf{B}^{(t+1)} \in \operatorname{argmin}_{\mathbf{B}} g(\mathbf{B}; \mathbf{B}^{(t)})$ , the sequence of iterates satisfies

$$f(\mathbf{B}^{(t+1)}) \leq g(\mathbf{B}^{(t+1)}; \mathbf{B}^{(t)}) \leq g(\mathbf{B}^{(t)}; \mathbf{B}^{(t)}) = f(\mathbf{B}^{(t)}). \quad (2.28)$$

Then the optimal solutions of problem (2.50) can be characterized as:

$$\mathbf{B} = \Theta(\mathbf{B} - \mathbf{X}^T \nabla \bar{l}(\mathbf{X}\mathbf{B}) / \rho; \lambda), \quad (2.29)$$

which correspond to the set of fixed points of the original optimization problem (2.13). Then similar techniques can be applied to construct SEL from this  $\Theta$ -equation. For the constrained-type problems like

$$\min_{\mathbf{B}} \bar{l}(\mathbf{X}\mathbf{B}; \mathbf{Y}) \quad \text{s.t.} \quad \|\mathbf{B}\|_0 \leq J, \quad (2.30)$$

the quartile thresholding  $\Theta^{\#}(\cdot; q)$  is helpful. Given any  $\mathbf{s} \in \mathbb{R}^{n \times m}$ , we define  $\Theta^{\#}(\mathbf{s}; q)$  to be a vector  $\mathbf{t}$  satisfying  $t_{(j)} = s_{(j)}$  if  $1 \leq j \leq q$  and 0 otherwise, where  $s_{(1)}, \dots, s_{(n)}$  are the order statistics of  $s_1, \dots, s_n$  satisfying  $|s_1| \geq \dots \geq |s_n|$ , and  $t_{(1)}, \dots, t_{(n)}$  are defined similarly. To avoid ambiguity, we assume no ties occur in performing  $\Theta^{\#}(\mathbf{s}; q)$ . Then the optimality condition for problem (2.30) is

$$\operatorname{vec} \mathbf{B} = \Theta^{\#}(\operatorname{vec}(\mathbf{B} - \mathbf{X}^T \nabla \bar{l}(\mathbf{X}\mathbf{B}) / \rho); J)$$

based on constructing a similar surrogate function:

$$g(\mathbf{B}; \mathbf{B}^-) = \bar{l}(\mathbf{X}\mathbf{B}^-) + \langle \nabla \bar{l}(\mathbf{X}\mathbf{B}^-), \mathbf{B} - \mathbf{B}^- \rangle + \rho \|\mathbf{B} - \mathbf{B}^-\|_F^2 \quad \text{s.t.} \quad \|\mathbf{B}\|_0 \leq J. \quad (2.31)$$

### 2.3.3 SEL for Sparse Reduced Rank Regression

In this section, we consider the case where the target matrix  $\mathbf{B}$  is simultaneous sparse and low-rank. To impose low-rankness, we in addition to the penalized optimization, we add another low-rank constraint  $\text{rank}(\mathbf{B}) \leq r$ . We first consider the simplest case when the loss  $l$  is induced by the Frobenius norm. Luckily, this low-rank constraint doesn't induce the need of an above type of surrogate function when a reparametrization is applied: let  $\mathbf{B} = \mathbf{S}\mathbf{V}^T$ , where  $\mathbf{S} \in \mathbb{R}^{p \times r}$  carries the sparsity structure of  $\mathbf{B}$  and  $\mathbf{V} \in \mathbb{O}^{m \times r}$  is an orthogonal matrix. Then the original optimization problem (2.13) can be reparametrized to:

$$\min_{\mathbf{S} \in \mathbb{R}^{p \times r}, \mathbf{V} \in \mathbb{O}^{m \times r}} \frac{\|\mathbf{X}\mathbf{S}\mathbf{V}^T - \mathbf{Y}\|_F^2}{2} + P(\mathbf{S}; \lambda). \quad (2.32)$$

A simple block coordinate descent algorithm for  $\mathbf{S}$  given  $\mathbf{V}$  and  $\mathbf{V}$  given  $\mathbf{S}$  can be applied to solve problem (2.32). To develop SEL for this type of parametrization, we need sample-additive estimation equations for optimizations of the two blocks. First, given the orthogonal matrix  $\mathbf{V}$ , since we have  $\|\mathbf{X}\mathbf{S}\mathbf{V}^T - \mathbf{Y}\|_F^2 = \|\mathbf{X}\mathbf{S} - \mathbf{Y}\mathbf{V}\|_F^2$ , the sample additive estimating equations can be derived based on the aforementioned sections with the transformed response matrix  $\mathbf{Y}\mathbf{V}$ . It is more tricky to derive an sample additive estimating equation for  $\mathbf{V}$  given  $\mathbf{S}$  with the constraint  $\mathbf{V}^T\mathbf{V} = \mathbf{I}$ . We apply the gradient on Riemannian manifold of  $\mathbf{V}$  to achieve this goal. Denote a function  $F(\mathbf{V})$ , the Euclidean gradient of  $\mathbf{V}$  as  $\nabla F(\mathbf{V})$  and the Stiefel manifold gradient as  $\nabla F^{(s)}(\mathbf{V})$ . Then given any matrix  $\mathbf{\Delta}_1, \mathbf{\Delta}_2 \in \mathbb{R}^{m \times r}$ , the canonical metric for Stiefel manifold is defined as  $g_V(\mathbf{\Delta}_1, \mathbf{\Delta}_2) = \text{Tr}(\mathbf{\Delta}_1^T(\mathbf{I} - \frac{1}{2}\mathbf{V}\mathbf{V}^T)\mathbf{\Delta}_2)$ . The identity induced by the metric is  $\langle \nabla F(\mathbf{V}), \mathbf{\Delta} \rangle = g_V(\nabla F^{(s)}(\mathbf{V}), \mathbf{\Delta})$ . Hence, we have  $\nabla F(\mathbf{V}) = (\mathbf{I} - \frac{1}{2}\mathbf{V}\mathbf{V}^T)\nabla F^{(s)}$ . Since  $\mathbf{V}^T\mathbf{V} = \mathbf{I}$ , we have  $\mathbf{V}^T\nabla F(\mathbf{V}) + \nabla F(\mathbf{V})^T\mathbf{V} = \mathbf{0}$ . Then we can see

$$\begin{aligned} \nabla F(\mathbf{V}) &= \nabla F(\mathbf{V}) - \frac{1}{2}\mathbf{V}(\mathbf{V}^T\nabla F(\mathbf{V}) + \nabla F(\mathbf{V})^T\mathbf{V}) \\ &= \nabla F(\mathbf{V}) - \mathbf{V}\mathbf{V}^T\nabla F(\mathbf{V}) + \frac{1}{2}\mathbf{V}\mathbf{V}^T\nabla F(\mathbf{V}) - \frac{1}{2}\mathbf{V}\nabla F(\mathbf{V})^T\mathbf{V} \\ &= (\mathbf{I} - \frac{1}{2}\mathbf{V}\mathbf{V}^T)(\nabla F(\mathbf{V}) - \mathbf{V}\nabla F(\mathbf{V})^T\mathbf{V}), \end{aligned} \quad (2.33)$$

which implies that  $\nabla F(\mathbf{V}) = \nabla F(\mathbf{V}) - \mathbf{V}\nabla F(\mathbf{V})^T\mathbf{V}$ . Let  $F_i(\mathbf{V}) = \|\mathbf{y}_i - \mathbf{x}_i^T\mathbf{S}\mathbf{V}^T\|_2^2/2$  such that  $\|\mathbf{X}\mathbf{S}\mathbf{V}^T - \mathbf{Y}\|_F^2 = \sum_{i=1}^n F_i(\mathbf{V})$ . Then we have  $\nabla F_i(\mathbf{V}) = \mathbf{x}_i(\mathbf{x}_i^T\mathbf{S}\mathbf{V}^T - \mathbf{y}_i)$  and the sample additive



estimating equation for  $\mathbf{V}$  given  $\mathbf{S}$ :

$$\mathbf{x}_i(\mathbf{x}_i^T \mathbf{S} \mathbf{V}^T - \mathbf{y}_i) - \mathbf{V} \mathbf{x}_i(\mathbf{x}_i^T \mathbf{S} \mathbf{V}^T - \mathbf{y}_i) \mathbf{V} = \mathbf{0}.$$

## 2.4 Asymptotic Analysis for SEL

### 2.4.1 SEL with Affine Inequality Constraints

Consider test with affine constraints  $\mathbf{A}\boldsymbol{\beta} \preceq \boldsymbol{\alpha}$  with  $\mathbf{A} \in \mathbb{R}^{k \times p}$  with rank  $r$  ( $k$  can be smaller than or greater than  $p$ ). The constraints should be compatible such that the true value  $\boldsymbol{\beta}^\circ$  is indeed in the feasible region of these constraints in equation (2.12). By introducing the slack variable, the null hypothesis is  $H_0 : \boldsymbol{\beta} = \boldsymbol{\beta}^\circ$ , such that: there is a vector  $\mathbf{s} \in \mathbb{R}^p$ , so that:

$$\begin{aligned} \mathbb{E}[\mathbf{X}^T \nabla \bar{l}(\mathbf{X}\boldsymbol{\beta}) + \mathbf{s}] &= \mathbf{0}, \\ \mathbf{A}[\mathcal{J}, ]\mathbf{s} &= \mathbf{0}, \quad \mathbf{A}[\mathcal{J}^c, ]\mathbf{s} \preceq \mathbf{0}. \end{aligned} \tag{2.34}$$

However, the constraints are not identifiable so we need to do some operations to revise the constrained matrix  $\mathbf{A}[\mathcal{J}, ]$  and  $\mathbf{A}[\mathcal{J}^c, ]$ :

1. If  $\mathbf{A}[\mathcal{J}^c, ]\mathbf{s} \preceq \mathbf{0}$  implies some equality constraints from the solution region of  $\mathbf{s}$  defined by only inequalities, then we put them into  $\mathbf{A}[\mathcal{J}, ]\mathbf{s} = \mathbf{0}$ ;
2. Pick up the linear independent constrains in both equality and inequality ones, such that  $\mathbf{A}[\mathcal{J}, ] \in \mathbb{R}^{r_1 \times p}$  and  $\text{rank}(\mathbf{A}[\mathcal{J}, ]) = r_1 \leq r \leq \min(k, p)$ ;  $\mathbf{A}[\mathcal{J}^c, ] \in \mathbb{R}^{k_2 \times p}$  and  $\text{rank}(\mathbf{A}[\mathcal{J}^c, ]) = r_2 \leq r \leq \min(k, p)$ .  $r_1 + r_2 = r$  and  $k_2 \leq 2r_2$ .

If  $\mathbf{A} = -\mathbf{I}$  and  $\boldsymbol{\alpha} = \mathbf{0}$ , then (2.12) reduces to the non-negative regression EL; if  $\mathcal{J} = \{1, 2, \dots, k\}$  and  $\mathcal{J}^c = \emptyset$ , (2.12) reduces to the traditional regression EL. The following conditions are required for us to obtain an asymptotic analysis of large  $n$  SEL:

1. Let  $\boldsymbol{\epsilon} = -\nabla \bar{l}(\mathbf{X}\boldsymbol{\beta}^\circ)$ . Suppose  $\mathbf{z}_i = -\mathbf{x}_i^T \boldsymbol{\epsilon}_i$ ,  $i = 1, 2, \dots, n$  follows a distribution with mean  $\mathbf{0}$  and covariance matrix  $\mathbf{V}$ ;
2. We assume that the interior of convex hull of  $\mathbf{z}_1, \dots, \mathbf{z}_n$  contains zero with probability 1 as  $n \rightarrow \infty$ ;
3. The minimum and maximum eigenvalues of the matrix  $\mathbf{V}$  is bounded below and above by constants;

Then we have the following theorem to characterize the asymptotic distribution of large  $n$  SEL under affine inequality constraints:

**Theorem 2.4.1.** *Suppose the above conditions 1-3 hold. Then under the inequality constraints  $\mathbf{A}\boldsymbol{\beta} \preceq \boldsymbol{\alpha}$  with  $\mathbf{A} \in \mathbb{R}^{k \times p}$  and  $\text{rank}(\mathbf{A}) = r$ , the SEL statistics for test (2.34) follows a chi-bar-square distribution as follows:*

$$\begin{aligned}
P(c_n(\boldsymbol{\beta}^\circ) \geq t) &= \sum_{j=|\mathcal{J}|}^p P(\chi_j^2 \geq t) \delta_j, \\
\delta_j &= \sum_{|\mathcal{I}|=j} P(\boldsymbol{\nu}_{\mathbf{A}_{\mathcal{I}^c-\mathcal{J}}} \preceq \mathbf{0}) P(\boldsymbol{\nu}_{\mathbf{A}_{\mathcal{I}^c}} | \boldsymbol{\nu}_{\mathbf{A}_{\mathcal{I}}} \preceq \mathbf{0}), \text{ with } p > |\mathcal{I}|, \mathcal{J} \subsetneq \mathcal{I}, \\
\delta_{|\mathcal{J}|} &= 1 - \sum_{j=|\mathcal{J}|+1}^p \delta_j,
\end{aligned} \tag{2.35}$$

where we denote  $\sqrt{n} \sum_{i=1}^n \frac{1}{n} \mathbf{z}_{\mathcal{H}i} \xrightarrow{d} \boldsymbol{\nu}_{\mathbf{A}_{\mathcal{H}}}$  for any index set  $\mathcal{H}$  and  $\mathbf{z}_{\mathcal{H}i}$  is the subvector of  $\mathbf{z}_i$  indexed by  $\mathcal{H}$ .

Since  $|\mathcal{J}| = r_1$  which is the rank of all the independent equality constraints, for  $\mathbf{A} = -\mathbf{I}$  and  $\mathcal{J} \subset \{1, 2, \dots, p\}$  and  $\mathcal{J}^c \neq \emptyset$ , then  $c_n(\boldsymbol{\beta}^\circ)$  reduces to the nonnegative regression EL; while if  $\mathcal{J} = \{1, 2, \dots, p\}$  and  $\mathcal{J}^c = \emptyset$ , then  $c_n(\boldsymbol{\beta}^\circ)$  reduces to the traditional regression EL with  $\chi_p^2$  distribution. We provide some numeric evidence that the distribution of SEL for large  $n$  cases with affine inequality constraints is indeed in a chi-bar-square form. In particular, we show that if only a single chi-square distribution is applied to generate p-values SEL, then the corresponding tests will be either too optimistic or conservative.

Two specific examples are considered: nonnegative regressions and regressions with variable sections. Suppose the data generating model is  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^\circ + \boldsymbol{\epsilon}$ , where  $\boldsymbol{\epsilon}$  follow a normal distribution with covariance  $\sigma^2 \mathbf{I}$ . Let  $n = 100$ ,  $p = 2$ ,  $\sigma = 1$ ,  $\boldsymbol{\beta}^\circ = (5, 0)^T$ ,  $X_{ij} \sim \mathcal{N}(0, 1)$  for all  $1 \leq i \leq n$  and  $1 \leq j \leq p$ . The simulation settings of the two cases are as follows:

1. Nonnegative regression: The  $p$ -value of SEL test is generated from a mixture of chi-square distribution with degrees of freedom one and two, where the weights for  $\chi_2^2$  and  $\chi_1^2$  are both 0.5;
2. Regression with variable selection: Let  $\lambda = \sigma \sqrt{\log p/n}$ . The  $p$ -value of SEL test is generated from a mixture of chi-square distribution with degrees of freedom one and two, where the weight for  $\chi_2^2$  is the probability that the maximum of 100 standard Gaussian random variables is greater than  $\lambda$ , while the weight for  $\chi_1^2$  can be obtained by one minus the above weight for  $\chi_2^2$ ;

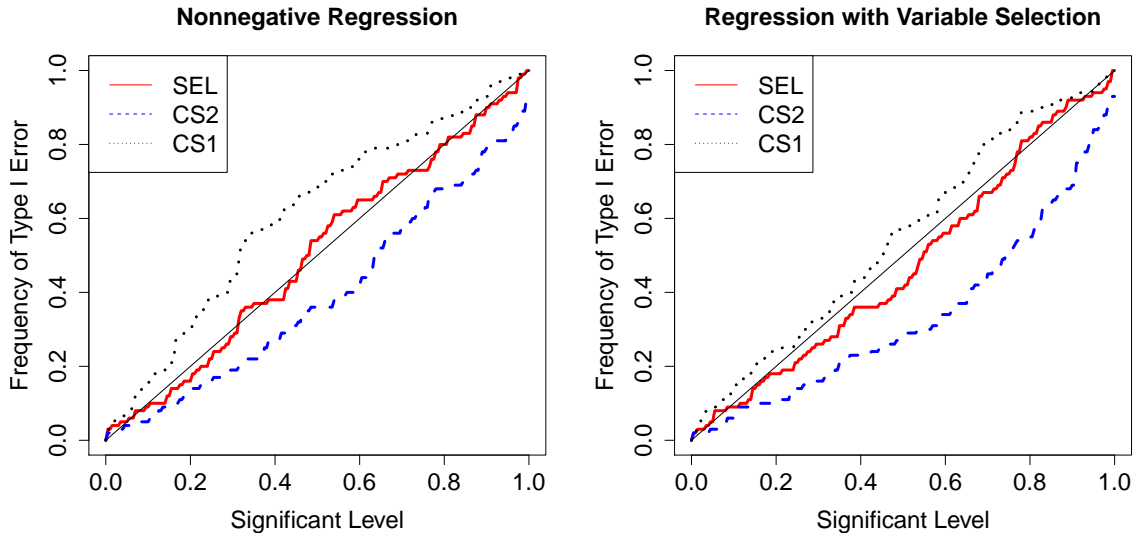


Figure 2.1: Percentage of the type I error made by 100 hypothesis tests on for  $H_0 : \beta = \beta^\circ$  v.s. the respective significant level by large  $n$  SEL.

The simulations are repeated for 100 times where we test the hypothesis  $\beta = \beta^\circ$  in both cases and record all the  $p$ -values of SEL tests. For comparison, we also record all the  $p$ -values when we ‘assume’ the distribution of SEL test is a single  $\chi_2^2$  or  $\chi_1^2$ . Finally, by ranging the significant level  $\alpha$  from zero to one, we show how the frequencies of the type I error change with respect to the significant level. The simulation results are shown in the following figure 2.1, where ‘CS1’ and ‘CS2’ stand for the testing results when the  $p$ - values are generated from single  $\chi_1^2$  and  $\chi_2^2$  distributions. From the figure, we can see the validity of our proposed theorem for large  $n$  SEL.

## 2.5 Analysis of Large $p$ SEL

Recall that our defined large  $p$  SEL contains a bias term  $\gamma^\circ$ , which can be obtained directly from the null hypothesis  $\beta^\circ$ . Therefore, the null hypothesis can be refined and the bias term can be eliminated. Finally the making inference for large  $p$  SEL corresponds the following estimating equation:

$$\beta = \Theta_H(\beta - \mathbf{X}^T \nabla \bar{l}(\mathbf{X} \beta^\circ)/n; \lambda), \tag{2.36}$$

where  $\Theta_H$  is the hard-thresholding operator. Unfortunately, when  $p > n$ , it is impossible to directly learn the marginal limiting distribution of the empirical likelihood ratio as in the lower dimensional

case. However, we can still make practical inference conditional on the selection event of the given data  $\mathbf{X}, \mathbf{y}$ . Some notations are worth to mention: let  $\boldsymbol{\alpha} \succeq \boldsymbol{\beta}$  mean all components of  $\boldsymbol{\alpha}$  is greater than the respective component of  $\boldsymbol{\beta}$ . For a matrix  $\mathbf{A}$  and an index set  $\mathcal{J}$ ,  $\mathbf{A}_{\mathcal{J}}$  means the the submatrix of  $\mathbf{A}$  with columns indexed by  $\mathcal{J}$  and  $\mathbf{A}_{\mathcal{J}, \mathcal{I}}$  is denoted as the submatrix of  $\mathbf{A}$  with column indexed by  $\mathcal{I}$  and row indexed by  $\mathcal{J}$ .

### 2.5.1 A Limiting Result for Large $p$ SEL for Regression

In this section, we consider the example of sparse regressions to illustrate our analysis of SEL. For high dimensional regression with  $P_H$  penalty, (2.43) corresponds to the null hypothesis  $H_0 : \boldsymbol{\beta} = \boldsymbol{\beta}^\circ$ , and there is a vector  $\mathbf{s} \in \mathbb{R}^p$ , so that:

$$\begin{aligned} \mathbb{E}[\mathbf{X}^T(\mathbf{X}\boldsymbol{\beta}^\circ - \mathbf{y}) + \mathbf{s}] &= \mathbf{0}, \\ \mathbf{s}_{\mathcal{J}} &= \mathbf{0}, \|\mathbf{s}\|_\infty \leq \lambda. \end{aligned} \tag{2.37}$$

The corresponding regression large  $p$  SEL is:

$$\begin{aligned} \min_{\mathbf{w} \in \mathcal{C}, \mathbf{s}} -\langle \mathbf{1}, \log(n\mathbf{w}) \rangle \quad \text{s.t.} \quad & \mathbf{X}^T \mathbf{W} \nabla \bar{l}(\mathbf{X}\boldsymbol{\beta}^\circ) + \mathbf{s} = \mathbf{0} \\ & \mathbf{s} \circ \boldsymbol{\beta}^\circ = \mathbf{0}, \|\mathbf{s}\|_\infty \leq \lambda. \end{aligned} \tag{2.38}$$

We consider the case when the selection is consistent, where we suppose given  $\mathbf{X}, \mathbf{y}$  and  $\lambda$ , there exists a  $\hat{\boldsymbol{\beta}}$  satisfying estimating equation (2.36), with nonzero index set  $\mathcal{J} = \{j : \hat{\beta}_j \neq 0\}$  the same with the truth  $\boldsymbol{\beta}^\circ$ . We first provide some equivalent characterizations of the selection event, where we can obtain a sparse solution  $\hat{\boldsymbol{\beta}}$  of estimating equation (2.36) given  $\mathbf{X}, \mathbf{y}$  and  $\lambda$  with nonzero index set  $\mathcal{J}$ . Clearly, by directly inverting the estimating equation, we have  $\hat{\boldsymbol{\beta}}_{\mathcal{J}} = (\mathbf{X}_{\mathcal{J}}^T \mathbf{X}_{\mathcal{J}})^{-1} \mathbf{X}_{\mathcal{J}}^T \mathbf{y}$  while  $\hat{\boldsymbol{\beta}}_{\mathcal{J}^c} = \mathbf{0}$ . In addition, the following inequality constraints are required:

$$\begin{cases} |\hat{\boldsymbol{\beta}}_{\mathcal{J}}| \succeq \lambda \mathbf{1} \\ \|\mathbf{X}_{\mathcal{J}^c}^T (\mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{y})\|_\infty < n\lambda, \end{cases} \iff \begin{cases} |(\mathbf{X}_{\mathcal{J}}^T \mathbf{X}_{\mathcal{J}})^{-1} \mathbf{X}_{\mathcal{J}}^T \mathbf{y}| \succeq \lambda \mathbf{1} \\ \|\mathbf{X}_{\mathcal{J}^c}^T (\mathbf{I} - \mathbf{P}_{\mathbf{X}_{\mathcal{J}}}) \mathbf{y}\|_\infty < n\lambda, \end{cases} \tag{2.39}$$

where  $\mathbf{P}_{\mathbf{X}_{\mathcal{J}}} = \mathbf{X}_{\mathcal{J}} (\mathbf{X}_{\mathcal{J}}^T \mathbf{X}_{\mathcal{J}})^{-1} \mathbf{X}_{\mathcal{J}}^T$  is a projection matrix induced by  $\mathbf{X}_{\mathcal{J}}$  and  $\mathbf{1}$  is a vector with all ones. Furthermore, under the knowledge of null hypothesis, we'd like to represent the above constraints into constraints of  $\boldsymbol{\epsilon}$ , which is

$$\begin{cases} |\boldsymbol{\beta}_{\mathcal{J}}^\circ + (\mathbf{X}_{\mathcal{J}}^T \mathbf{X}_{\mathcal{J}})^{-1} \mathbf{X}_{\mathcal{J}}^T \boldsymbol{\epsilon}| \succeq \lambda \mathbf{1} \\ \|\mathbf{X}_{\mathcal{J}^c}^T (\mathbf{I} - \mathbf{P}_{\mathbf{X}_{\mathcal{J}}}) \boldsymbol{\epsilon}\|_\infty < n\lambda. \end{cases} \tag{2.40}$$

Equation (2.40) is helpful to to show the equivalence between prime and dual problem of SEL.

Based on the above analysis, we propose the following assumptions for our main theorem:

1. Suppose  $\beta^\circ$  is sparse with nonzero index  $\mathcal{J} = \{j : \beta_j^\circ \neq 0\}$ ,  $\min_{j \in \mathcal{J}} |\beta_j^\circ| > (1 + \kappa)\lambda$  for a constant  $\kappa > 0$ , and  $|\mathcal{J}| < n$ . Assume  $\epsilon = \mathbf{y} - \mathbf{X}\beta^\circ$ . Let  $\mathbf{z}_i = -\mathbf{x}_i^T \epsilon_i$ , and  $\mathbf{z}_1, \dots, \mathbf{z}_n$  follows a distribution with mean  $\mathbf{0}$  and covariance matrix  $\mathbf{V}$ . In addition, the minimum and maximum eigenvalues of the matrix  $\mathbf{V}$  is bounded below and above by constants;
2. Given  $\mathbf{X}, \mathbf{y}, \lambda$ , there exists vector  $\hat{\beta}$ , such that it satisfies the estimating equation (2.36), and the nonzero index set is also  $\mathcal{J}$ .
3. Let  $\mathbf{z}_{\mathcal{J}i}$  be the subvector of  $\mathbf{z}_i$  indexed by  $\mathcal{J}$ . We assume that the interior of the convex hull of  $\mathbf{z}_{\mathcal{J}1}, \dots, \mathbf{z}_{\mathcal{J}n}$  contains zero with probability 1 as  $n \rightarrow \infty$ ;
4. The minimum and maximum eigenvalues of the matrix  $\mathbf{X}_{\mathcal{J}}^T \mathbf{X}_{\mathcal{J}}/n$  is bounded below and above by constants;
5. There exists  $\delta > 0$ , such that  $\|\mathbf{X}_{\mathcal{J}c}^T \mathbf{X}_{\mathcal{J}} (\mathbf{X}_{\mathcal{J}}^T \mathbf{X}_{\mathcal{J}})^{-1}\|_\infty \leq \delta$ , and  $\delta n^{1/2} = o(\lambda)$ .

To make practical inference for high dimensional regression SEL, we propose the following theorem:

**Theorem 2.5.1.** *Under the above assumptions 1-5, denote the variable  $\mathbf{z}_\lambda^{\mathbf{X}, \beta^\circ}$  as*

$$\left( \sum_{i=1}^n \mathbf{z}_{\mathcal{J}i} \right)^T \left( \sum_{i=1}^n \mathbf{z}_{\mathcal{J}i} \mathbf{z}_{\mathcal{J}i}^T \right)^{-1} \left( \sum_{i=1}^n \mathbf{z}_{\mathcal{J}i} \right). \quad (2.41)$$

*Then if as  $n$  and  $p \rightarrow \infty$ ,  $\lambda \rightarrow 0$ , there exist constants  $c_1, c_2$  such that with probability converging to 1, the optimal SEL ratio  $c(\beta^\circ) = -2\langle \mathbf{1}, \log(n\hat{\mathbf{w}}) \rangle$  satisfies*

$$P(c(\beta^\circ) \geq t) = P(\mathbf{z}_\lambda^{\mathbf{X}, \beta^\circ} \geq t) + o_p(1), \quad (2.42)$$

*given any positive value  $t$ , where  $\hat{\mathbf{w}}$  is the optimizer of weights in problem (2.43).*

Our theory renders a limiting behavior of the distribution of  $c_n(\beta^\circ)$  when section consistency can be obtained. However, signal strength assumptions is unnecessary since the information is already incorporated into the null hypothesis. We demonstrate the validity of theorem 2.5.1, the limiting theory of large  $p$  SEL. We set  $\lambda$  to be large enough such that the assumptions of the main theorem can be satisfied. The details of simulation settings are as follows:

1. Case 1: let  $n = 100$ ,  $p = 200$ ,  $\sigma = 1$ ,  $\beta^\circ = (5, 0, \dots, 0)^T$ ,  $X_{ij} \sim \mathcal{N}(0, 1)$  for all  $1 \leq i \leq n$  and  $1 \leq j \leq p$ ,  $\lambda = \sigma\sqrt{2\log p}$ ;
2. Case 2: let  $n = 100$ ,  $p = 200$ ,  $\sigma = 2.5$ ,  $\beta^\circ = (5, 0, \dots, 0)^T$ ,  $X_{ij} \sim \mathcal{N}(0, 1)$  for all  $1 \leq i \leq n$  and  $1 \leq j \leq p$ ,  $\lambda = \sigma\sqrt{2\log p}$ ;

3. Case 3: let  $n = 200$ ,  $p = 500$ ,  $\sigma = 1$ ,  $\beta^\circ = (5, 0, \dots, 0)^T$ ,  $X_{ij} \sim \mathcal{N}(0, 1)$  for all  $1 \leq i \leq n$  and  $1 \leq j \leq p$ ,  $\lambda = \sigma\sqrt{2\log p}$ ;
4. Case 4: let  $n = 100$ ,  $p = 200$ ,  $\sigma = 1$ ,  $\beta^\circ = (5, 5, 0, \dots, 0)^T$ ,  $X_{ij} \sim \mathcal{N}(0, 1)$  for all  $1 \leq i \leq n$  and  $1 \leq j \leq p$ ,  $\lambda = \sigma\sqrt{2\log p}$ .

The simulation is repeated for 100 times where we still test the hypothesis  $\beta = \beta^\circ$  and record all the  $P$ -values. Finally, by ranging the significant level  $\alpha$  from zero to one, we show how the frequencies of the type I error change with respect to the significant level in the following Figure 2.2. From the figure, we can see that the type I error can be successfully controlled by the respective significant levels when the selection consistency can happen, which implies that our main theorem is effective in a large range of settings in high dimensional inference problem.

To make practical inference in real applications, we consider the case of SEL where the  $\lambda$  is fixed in the following section.

### 2.5.2 Nonasymptotic Analysis of Large $p$ SEL for a Fixed $\lambda$

In this section, we consider large  $p$  SEL with a general loss function and a fixed  $\lambda$ :

$$\begin{aligned} \min_{\mathbf{w} \in \mathcal{C}, \mathbf{s}} -\langle \mathbf{1}, \log(n\mathbf{w}) \rangle \quad \text{s.t.} \quad & \mathbf{X}^T \mathbf{W} \nabla \bar{l}(\mathbf{X}\beta^\circ) + \mathbf{s} = \mathbf{0} \\ & \mathbf{s} \circ \beta^\circ = \mathbf{0}, \|\mathbf{s}\|_\infty \leq \lambda. \end{aligned} \quad (2.43)$$

We show that the optimal slack empirical likelihood ratio follows a specific distribution given the data when  $\lambda$  is fixed. For notation ease, denote  $\mathbf{z}_i = [\mathbf{z}_{\mathcal{J}_i}^T, \mathbf{z}_{\mathcal{J}^c_i}^T]^T = -\mathbf{x}^T \epsilon_i$  where  $\mathbf{z}_{\mathcal{J}_i} = -\mathbf{x}_{\mathcal{J}_i}^T \epsilon_i$  and  $\mathbf{z}_{\mathcal{J}^c_i} = -\mathbf{x}_{\mathcal{J}^c_i}^T \epsilon_i$ . We assume that the following condition hold:

1. The optimization (2.43) is feasible: there exists a  $\mathbf{w}^0 = \{w_1^0, \dots, w_n^0\}$  such that  $\sum_{i=1}^n w_i^0 \mathbf{z}_{\mathcal{J}_i} = \mathbf{0}$ ,  $-\lambda \mathbf{1} \preceq \sum_{i=1}^n w_i^0 \mathbf{z}_{\mathcal{J}^c_i} \preceq \lambda \mathbf{1}$ , and  $0 < w_i^0 < 1$  for  $i = 1, 2, \dots, n$ .

The above assumption guarantees that the Slater's conditions hold, therefore, dual optimization can render an equivalent characterization of the prime problem.

To obtain the Lagrange dual function, first, by introducing Lagrange multipliers  $\phi \in \mathbb{R}$ ,  $\mathbf{u}_0 \in \mathbb{R}^{|\mathcal{J}|}$ ,  $\mathbf{h} \in \mathbb{R}^p$  and  $\mathbf{v}^-, \mathbf{v}^+ \in \mathbb{R}^{|\mathcal{J}^c|}$ , the optimization problem can be formulated to the following Lagrangian form:

$$\begin{aligned} T(\mathbf{w}, \mathbf{s}, \phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}^-, \mathbf{v}^+) = & -\langle \mathbf{1}, \log(n\mathbf{w}) \rangle + \phi \left( \sum_{i=1}^n w_i - 1 \right) + \mathbf{h}^T (-\mathbf{X}^T \mathbf{W} \epsilon + \mathbf{s}) \\ & + \mathbf{u}_0^T \mathbf{s}_{\mathcal{J}} + \mathbf{v}^{-T} (-\mathbf{s}_{\mathcal{J}^c} - \lambda \mathbf{1}) + \mathbf{v}^{+T} (\mathbf{s}_{\mathcal{J}^c} - \lambda \mathbf{1}), \end{aligned} \quad (2.44)$$

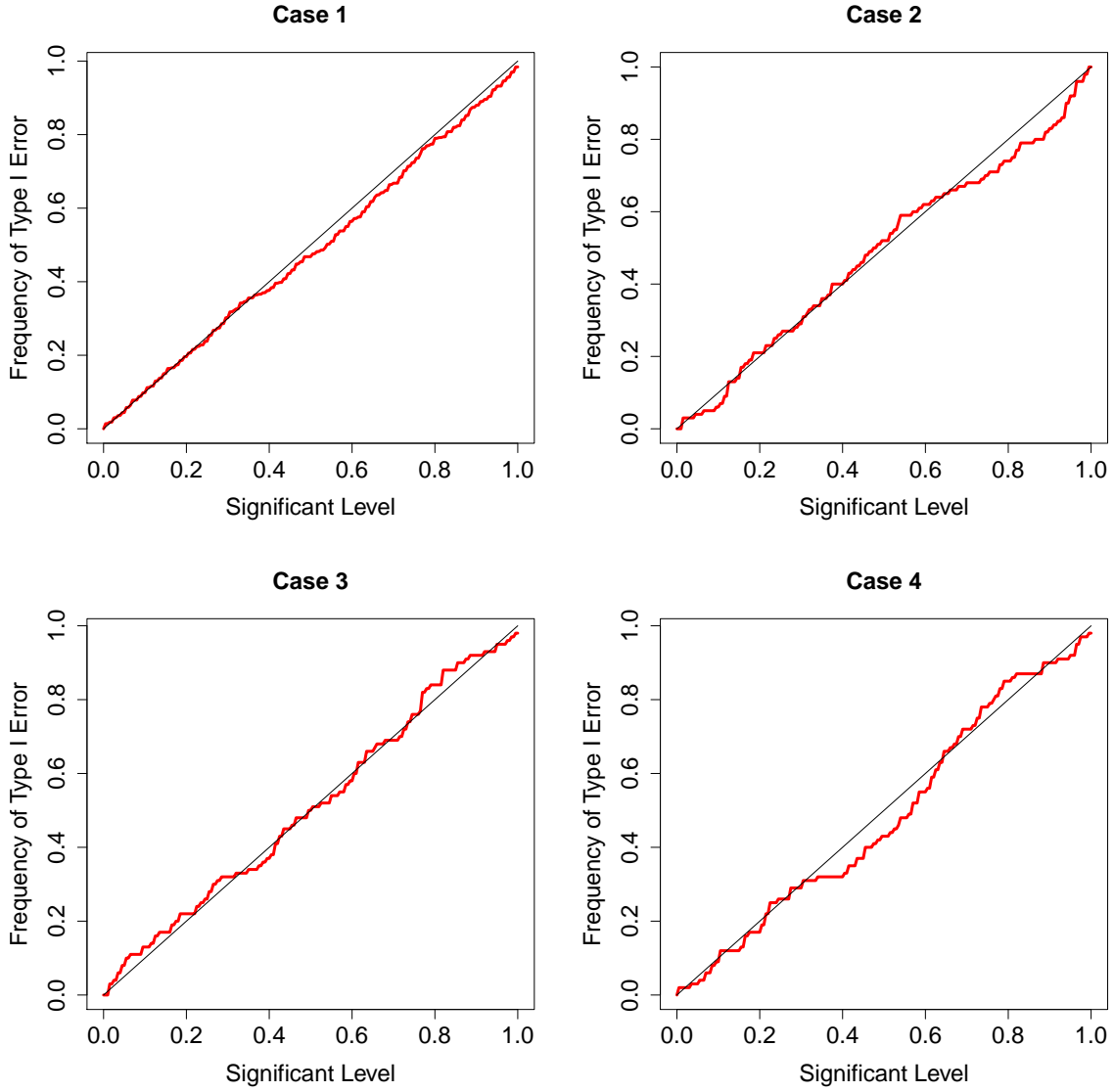


Figure 2.2: Percentage of the type I error made by 100 hypothesis tests on for  $H_0 : \beta = \beta^\circ$  v.s. the respective significant level by conditional inference of SEL.

and the Lagrange dual function is:

$$d(\phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}^-, \mathbf{v}^+) = \min_{\mathbf{w}, \mathbf{s}} T(\mathbf{w}, \mathbf{s}, \phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}^-, \mathbf{v}^+), \quad (2.45)$$

which is expected to be maximized. Taking derivative of the Lagrangian  $T$  with respect to  $\mathbf{w}$  and  $\mathbf{s}$ , we have

$$\begin{cases} -\frac{1}{\hat{w}_i} + \mathbf{h}^T \mathbf{z}_i + \phi = 0, \\ \mathbf{h}_{\mathcal{J}} = -\mathbf{u}_0, \\ \mathbf{h}_{\mathcal{J}^c} = \mathbf{v}^- - \mathbf{v}^+. \end{cases} \quad (2.46)$$

Plugging equation (2.72) into the Lagrangian to get the Lagrange dual function, we have:

$$\begin{aligned} d(\phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}^-, \mathbf{v}^+) &= -\sum_{i=1}^n \log\left(\frac{n}{\phi + \mathbf{h}^T \mathbf{z}_i}\right) + \phi \left(\sum_{i=1}^n \frac{1}{\phi + \mathbf{h}^T \mathbf{z}_i} - 1\right) \\ &\quad + \sum_{i=1}^n \frac{\mathbf{h}^T \mathbf{z}_i}{\phi + \mathbf{h}^T \mathbf{z}_i} - \mathbf{v}^{+T} \lambda \mathbf{1} - \mathbf{v}^{-T} \lambda \mathbf{1}, \\ &= -\sum_{i=1}^n \log\left(\frac{n}{\phi + \mathbf{h}^T \mathbf{z}_i}\right) + n - \phi - \mathbf{v}^{+T} \lambda \mathbf{1} - \mathbf{v}^{-T} \lambda \mathbf{1}. \end{aligned} \quad (2.47)$$

Note that this dual problem implicitly require that the constraints  $\phi + \mathbf{h}^T \mathbf{z}_i > 0$  for  $i = 1, \dots, n$ , this can be guarantee by the prespecified condition and the first equation of (2.72). To maximize equation (2.73), taking derivatives with respect to  $\mathbf{u}_0$  and  $\phi$ , we have:

$$\begin{cases} \nabla_{\mathbf{u}_0} d(\phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}^-, \mathbf{v}^+) = \sum_{i=1}^n \frac{-\mathbf{z}_{\mathcal{J}i}}{\phi + \mathbf{h}^T \mathbf{z}_i} = \mathbf{0}, \\ \nabla_{\phi} d(\phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}^-, \mathbf{v}^+) = \sum_{i=1}^n \frac{1}{\phi + \mathbf{h}^T \mathbf{z}_i} - 1 = 0. \end{cases} \quad (2.48)$$

In addition, taking derivatives with respect to  $\mathbf{v}^-$  and  $\mathbf{v}^+$ , we have:

$$\begin{cases} \nabla_{\mathbf{v}^-} d(\phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}^-, \mathbf{v}^+) = \sum_{i=1}^n \frac{\mathbf{z}_{\mathcal{J}^c i}}{\phi + \mathbf{h}^T \mathbf{z}_i} - \lambda \mathbf{1}, \\ \nabla_{\mathbf{v}^+} d(\phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}^-, \mathbf{v}^+) = \sum_{i=1}^n \frac{-\mathbf{z}_{\mathcal{J}^c i}}{\phi + \mathbf{h}^T \mathbf{z}_i} - \lambda \mathbf{1}. \end{cases} \quad (2.49)$$

Not all components of equation (2.49) can be set to zeros since the dual optimization is limited to  $\mathbf{v}^- \succeq \mathbf{0}$  and  $\mathbf{v}^+ \succeq \mathbf{0}$ . However, we always have

$$\begin{cases} \mathbf{v}^{-T} \left( \sum_{i=1}^n \frac{\mathbf{z}_{\mathcal{J}^c i}}{\phi + \mathbf{h}^T \mathbf{z}_i} - \lambda \mathbf{1} \right) = 0, \\ \mathbf{v}^{+T} \left( \sum_{i=1}^n \frac{-\mathbf{z}_{\mathcal{J}^c i}}{\phi + \mathbf{h}^T \mathbf{z}_i} - \lambda \mathbf{1} \right) = 0, \end{cases} \quad (2.50)$$



since either the components of  $\mathbf{v}^-$  and  $\mathbf{v}^+$  or the corresponding components of gradients  $\nabla_{\mathbf{v}^-}d$  and  $\nabla_{\mathbf{v}^+}d$  are zeros when maximizing the dual function. Then by combining equation (2.72), (2.48) and (2.50), we have:

$$\phi + \mathbf{v}^{+T}\lambda\mathbf{1} + \mathbf{v}^{-T}\lambda\mathbf{1} - n = 0, \quad (2.51)$$

which gives  $\phi = n - \mathbf{v}^{+T}\lambda\mathbf{1} - \mathbf{v}^{-T}\lambda\mathbf{1}$  and implies that  $\hat{w}_i = 1/(n - \mathbf{u}_0^T \mathbf{z}_{\mathcal{J}_i} + \mathbf{v}^{+T}(-\mathbf{z}_{\mathcal{J}^c_i} - \lambda\mathbf{1}) + \mathbf{v}^{-T}(\mathbf{z}_{\mathcal{J}^c_i} - \lambda\mathbf{1}))$ .

Next, based on the form of  $\hat{\mathbf{w}}$ , we show that the obtained optimal  $\hat{\mathbf{w}}$  by solving SEL can indeed define a unique distribution. Note that the optimization (2.43) can be written as the following equivalent form:

$$\min_{\mathbf{w} \in \mathcal{C}} -\langle \mathbf{1}, \log(n\mathbf{w}) \rangle \quad \text{s.t.} \quad \sum_{i=1}^n w_i \mathbf{z}_{\mathcal{J}_i} = \mathbf{0}, \quad -\lambda\mathbf{1} \leq \sum_{i=1}^n w_i \mathbf{z}_{\mathcal{J}^c_i} \leq \lambda\mathbf{1}. \quad (2.52)$$

Then the following theorem indicates that through dual optimization, the induced distribution of  $\hat{\mathbf{w}}$  is unique.

**Theorem 2.5.2.** *Denote  $\hat{\mathbf{w}} = \{\hat{w}_1, \dots, \hat{w}_n\}$  such that*

$$\hat{w}_i = \frac{1}{n - \mathbf{u}_0^T \mathbf{z}_{\mathcal{J}_i} + \mathbf{v}^{-T}(\mathbf{z}_{\mathcal{J}^c_i} - \lambda) + \mathbf{v}^{+T}(-\mathbf{z}_{\mathcal{J}^c_i} - \lambda)}$$

for some  $\mathbf{u}_0 \in \mathbb{R}^{|\mathcal{J}|}$ , and  $\mathbf{v}^-, \mathbf{v}^+ \in \mathbb{R}^{|\mathcal{J}^c|}$  with  $\mathbf{v}^- \succeq \mathbf{0}$  and  $\mathbf{v}^+ \succeq \mathbf{0}$ , then for any other  $\mathbf{w}' = \{w'_1, \dots, w'_n\}$  that is different from  $\hat{\mathbf{w}}$ , such that  $\sum_{i=1}^n w'_i \mathbf{z}_{\mathcal{J}_i} = \mathbf{0}$ ,  $-\lambda\mathbf{1} \preceq \sum_{i=1}^n w'_i \mathbf{z}_{\mathcal{J}^c_i} \preceq \lambda\mathbf{1}$ , and  $0 < w'_i < 1$  for  $i = 1, 2, \dots, n$ , we have

$$-\langle \mathbf{1}, \log(n\mathbf{w}') \rangle < -\langle \mathbf{1}, \log(n\hat{\mathbf{w}}) \rangle.$$

Since the optimizer  $\hat{\mathbf{w}}$  is defined as a unique distribution according to the data, the optimal value  $c_0(\beta^\circ) = -2\langle \mathbf{1}, \log(n\hat{\mathbf{w}}) \rangle$  follows a specific distribution, which can be applied to make practical nonparametric inference. In order to calculate  $p$ -value, we can apply bootstrap to generate pseudo-samples of  $\mathbf{x}$  and  $\mathbf{y}$ , and then apply SEL optimization to calculate the corresponding  $c_0(\beta^\circ)$  for the samples. Therefore, the distribution of  $c_0(\beta^\circ)$  can be obtained and the nonparametric inference can be applied. In the following section, we show some numeric evidence to support our main theorems and demonstrate that our SEL approach can indeed provide robustification for some statistical inference problems.

## 2.6 Experiments

### 2.6.1 Simulation Studies

In this part, we compare our nonparametric inference method with other approaches that rely on the correct model and distribution specification. Throughout the simulations, all SEL optimizations are solved by *CVXR* package in R. In particular, two circumstances are considered:

**(Estimating sparse multivariate mean)** Let the data generating model be  $\mathbf{y} = \boldsymbol{\beta}^\circ + \boldsymbol{\epsilon}$  with  $\boldsymbol{\epsilon}$  follows a specific distribution. Let  $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^p$  are i.i.d observations of responses with respect to  $\boldsymbol{\epsilon}_1, \dots, \boldsymbol{\epsilon}_n$ . Suppose that  $\boldsymbol{\beta}^\circ$  is sparse with nonzero index set  $\mathcal{J}$ . Based on the estimating equation, given any  $\lambda$ , we first calculate the hard thresholding estimator with:

$$\begin{cases} \hat{\beta}_i = \frac{1}{n} \sum_{k=1}^n y_{ki}, & \text{if } |\frac{1}{n} \sum_{k=1}^n y_{ki}| \geq \lambda, \\ \hat{\beta}_i = 0, & \text{if } |\frac{1}{n} \sum_{k=1}^n y_{ki}| < \lambda, \end{cases} \quad (2.53)$$

where  $y_{ki}$  is the  $i$ th component of  $\mathbf{y}_k$ . Suppose this hard thresholding estimator has nonzero index set  $\mathcal{I}$  with  $\mathcal{J} \subset \mathcal{I}$  and  $|\mathcal{I}| < n$ , then clearly, the selection event induced by  $\hat{\boldsymbol{\beta}}$  is

$$\min_{i \in \mathcal{I}} |\beta_i^\circ + \frac{1}{n} \sum_{k=1}^n \epsilon_{ki}| \geq \lambda, \quad \text{and} \quad \max_{i \in \mathcal{I}^c} |\frac{1}{n} \sum_{k=1}^n \epsilon_{ki}| < \lambda. \quad (2.54)$$

When the distribution of the noise  $\boldsymbol{\epsilon}$  is assumed to be normal with mean zero and variance  $\sigma^2 \mathbf{I}$ , the exact inference for components  $\mathcal{I} - \mathcal{J}$  can be made by drawing samples from Gaussian distributions  $\mathcal{N}(0, \sigma^2)$  lower truncated at  $\lambda$  with 0.5 probability and  $\mathcal{N}(0, \sigma^2)$  upper truncated at  $-\lambda$  with 0.5 probability. Therefore, we can calculate the converge probability of the truth by constructing confidence intervals of components indexed by  $\mathcal{I} - \mathcal{J}$  through the above exact distribution.

On the other hand, the respective SEL to test  $\boldsymbol{\beta}^\circ$  is defined as follows:

$$\begin{aligned} \min_{\mathbf{w} \in \mathcal{C}, \mathbf{s}} -\langle \mathbf{1}, \log(n\mathbf{w}) \rangle \quad \text{s.t.} \quad & \sum_{i=1}^n w_i (\boldsymbol{\beta}^\circ - \mathbf{y}_i) + \mathbf{s} = \mathbf{0} \\ & \mathbf{s} \circ \boldsymbol{\beta}^\circ = \mathbf{0}, \quad \|\mathbf{s}\|_\infty \leq \lambda. \end{aligned} \quad (2.55)$$

In addition, the  $P$ -value of the optimal SEL ratio  $c(\boldsymbol{\beta}^\circ) = -2\langle \mathbf{1}, \log(n\hat{\mathbf{w}}) \rangle$  can be constructed from bootstrap by resampling  $\boldsymbol{\epsilon} = \mathbf{y} - \boldsymbol{\beta}^\circ$ . We consider the cases where  $n = 100$ ,  $p = 200$  and  $n = 200$ ,  $p = 500$ , the truth  $\boldsymbol{\beta}^\circ = [5, 0, \dots, 0]^T$ ,  $\lambda = \sqrt{\log p/n}$  and three types of distribution of  $\boldsymbol{\epsilon}$ :

1. Case 1:  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, 3)$ ;
2. Case 2:  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, 3)$  with probability 0.9 while  $\boldsymbol{\epsilon} \sim T$  distribution with  $df = 3$  (which also has a variance of 3) with probability 0.1;

Table 2.1: Comparisons between SEL and exact Gaussian inference

Unit %	$n = 100, p = 200$			$n = 200, p = 500$			
	Case 1	Case 2	Case 3	Case 1	Case 2	Case 3	
$C = 0.90$	89.60	89.31	89.55	85.77	88.28	86.97	SEL
	89.96	80.50	68.25	87.78	87.37	54.51	GEI
$C = 0.95$	95.30	96.33	94.18	93.99	93.79	92.99	SEL
	94.98	87.42	80.50	93.79	92.38	61.52	GEI
$C = 0.99$	99.33	98.11	99.47	98.60	98.20	97.59	SEL
	99.19	95.60	68.25	98.60	96.59	54.51	GEI

3. Case 3:  $\epsilon \sim T$  distribution with  $df = 3$ .

Therefore, for Gaussian exact inference, the noise specifications vary from correct to contaminated and finally totally misspecified. We show the coverage probability of the constructed confidence intervals at different confidence levels  $C = 0.90, 0.95$  and  $0.99$  in Table 2.1 for both our SEL method and the Gaussian exact inference (GEI) approach. From the tables, we can see that the exact inference by assuming Gaussian noise will easily fall even if there is only 10% contamination of the noises, however, our SEL approach can keep the coverage level even if the distribution is totally misspecified. The simulation results demonstrates the robustness of our SEL method with respect to the distribution assumptions of the effective noise.

**(High dimensional regression)** Consider the regression model  $\mathbf{y} = \mathbf{X}\beta^\circ + \epsilon$  where  $\beta^\circ$  is sparse with nonzero index  $\mathcal{J}$  and the target is to make inference on the coefficient  $\beta^\circ$ . Post-selection approach [8] is a well-known exact inference method for the above type of sparse regressions which relies on the regression model and Gaussian assumption of the noises. First, an estimator  $\hat{\beta}$  with nonzero index  $\mathcal{I}$  is obtained by the formulation of the lasso problem given  $\lambda$ , where  $\mathcal{J} \subset \mathcal{I}$  is necessary. Then, via the KKT condition of the lasso problem, we can obtain some valuable information of the data  $\mathbf{X}, \mathbf{y}$ , which is in the form of affine constraints of  $\mathbf{y}$ . Finally, based on the Gaussian assumption, the exact distribution of the target statistics, like the OLS estimator on  $\mathcal{I}$ , conditional on the above obtained affine constraints can be calculated, which is often in the truncated Gaussian form. Therefore, the statistical inference on the target statistics can be carried out. We compare our large  $p$  SEL (2.43) with the above truncated Gaussian inference for the lasso (TGL). To obtain the  $P$ -value of the optimal SEL ratio, we generate bootstrap sample  $\mathbf{X}', \mathbf{y}'$ . Then in order to avoid duplicate effective noise defined by the estimating equations, we still use

Table 2.2: Comparisons between SEL and truncated Gaussian inference of the lasso

Unit %	$n = 100, p = 200$			$n = 200, p = 500$			
	Case 1	Case 2	Case 3	Case 1	Case 2	Case 3	
$C = 0.90$	88.78	90.20	86.50	88.84	89.41	84.25	SEL
	88.16	84.08	70.14	90.26	87.06	63.19	TGL
$C = 0.95$	93.67	95.39	92.02	94.30	96.47	90.80	SEL
	93.88	90.00	75.46	95.49	92.94	68.71	TGL
$C = 0.99$	98.98	99.59	96.73	98.81	97.65	97.75	SEL
	98.78	94.08	80.37	99.29	98.82	73.01	TGL

original  $X$  to be the predictor, while the effective noise is generate according to  $\mathbf{X}^T(\mathbf{X}'^T)^+\mathbf{X}'^T(\mathbf{y}' - \mathbf{X}'\beta^\circ)$ , where  $(\mathbf{X}'^T)^+$  is the Moore-Penrose inverse of  $\mathbf{X}'^T$ . We consider the cases where  $n = 100$ ,  $p = 200$  and  $n = 200$ ,  $p = 500$ , the truth  $\beta^\circ = [5, 0, \dots, 0]^T$ ,  $\lambda = 1.5 * \sqrt{\log p/n}$  and the same three types of distribution of  $\epsilon$  as above case of estimating sparse multivariate mean. The simulation results are shown in Table 2.2. From the simulation results, we can see that when the distribution assumption is correct, our method can perform almost as good as the exact inference approach. However, when the distribution is totally misspecified, the exact approach will fall completely, while our large  $p$  SEL can still some reasonable results. The simulations results support our opinion that our SEL approach can provide nonparametrization and robustification for statistical inference.

### 2.6.2 Real Data

Riboflavin data denotes how the log-transformed riboflavin production rate varies according to variables measuring the logarithm of the expression level of 4088 genes of *Bacillus subtilis*. Many research has been done to discover the association between response and predictors, like [1]. Two tasks are of particular interest: one is to select particular meaningful genes that could contribute to the response, the other is to quantify the uncertainty of the selected genes. Our proposed methodology focus on the second task. The specific data set can be founded in *hdi* R package, which contains  $n = 71$  observations and  $p = 4088$  predictors. To validate our methodology, we wish to test whether the gene ‘YXLD at’ is significant as it is shown in [1], this gene is significant at 0.05 level through multi sample-splitting inference method but not significant when tested via projected-based method [32];[26]. We apply the following procedure to construct a valid hypothesis test based on our SEL methodology:

1. Choosing proper  $\lambda$  to test the null hypothesis  $\beta^\circ = \mathbf{0}$  and obtaining the first  $p$ -value  $p_1$ ;

2. Denoting  $j$  as the index of ‘YXLD at’ and Calculating the OLS estimated value at index  $j$ :  
 $(\mathbf{x}_j^T \mathbf{x}_j)^{-1} \mathbf{x}_j^T \mathbf{y} = -0.55$ ;
3. Choosing another proper  $\lambda$  such that the SEL is feasible if  $\beta^\circ$  is changed on the component of ‘YXLD at’ around  $-0.55$ , then computing the 95% confidence interval of ‘YXLD at’.

After performing the above steps, we can get  $p_1 = 0.01$  and the confidence interval is  $[-0.75, -0.33]$ , which implies that at 5% significant level, we should reject the null hypothesis that the coefficients are zeros while we don’t have enough evidence to reject the statement that the gene ‘YXLD at’ is significant. The result delivers similar information as in [1], which supports the efficacy of our proposed methodology.

## 2.7 Summary

In this section, we extend traditional empirical likelihood approach by incorporating a slack variable to characterize some inequality constraints of the inference target. Both asymptotic and nonasymptotic theories are provided to support the inference. In addition, some simulations show that our inference approach offers a mechanism of nonparametrization and robustification for the parametric assumptions. In the future, it is important to explore composite slack empirical likelihood where the inference target is solved through a joint optimization instead of prespecified. This composite approach provide an efficient way to learn the hypothesis, which is extremely useful in high dimensional inference problem.

## 2.8 Outlines of Proofs

### 2.8.1 Proof of Theorem 2.4.1

To obtain the Lagrange dual function, first, by introducing Lagrange multipliers  $\phi \in \mathbb{R}$ ,  $\mathbf{u}_0 \in \mathbb{R}^{|\mathcal{J}|}$ ,  $\mathbf{h} \in \mathbb{R}^p$  and  $\mathbf{v} \in \mathbb{R}^{|\mathcal{J}^c|}$ , the optimization problem can be formulated to the following Lagrangian form:

$$T(\mathbf{w}, \mathbf{s}, \phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}) = -\langle \mathbf{1}, \log(n\mathbf{w}) \rangle + \phi \left( \sum_{i=1}^n w_i - 1 \right) + \mathbf{h}^T (-\mathbf{X}^T \mathbf{W} \boldsymbol{\epsilon} + \mathbf{s}) + \mathbf{u}_0^T \mathbf{A}[\mathcal{J}, ] \mathbf{s} + \mathbf{v}^T (\mathbf{A}[\mathcal{J}^c, ] \mathbf{s}), \quad (2.56)$$

and the Lagrange dual function is:

$$d(\phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}) = \min_{\mathbf{w}, \mathbf{s}} T(\mathbf{w}, \mathbf{s}, \phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}), \quad (2.57)$$

which is expected to be maximized. To see that the Slater's condition is hold, note that the equality constraints for the primal problem are affine and we will show the conditions on inequality constraints later. Taking derivative of the Lagrangian  $T$  with respect to  $\mathbf{w}$  and  $\mathbf{s}$ , we have

$$\begin{cases} -\frac{1}{\hat{w}_i} + \mathbf{h}^T \mathbf{z}_i + \phi = 0, \\ \mathbf{h}_{\mathcal{J}} = -\mathbf{A}[\mathcal{J},]^T \mathbf{u}_0, \\ \mathbf{h}_{\mathcal{J}^c} = -\mathbf{A}[\mathcal{J}^c,]^T \mathbf{v}. \end{cases} \quad (2.58)$$

Plugging equation (2.58) into the Lagrangian to get the Lagrange dual function, we have:

$$\begin{aligned} d(\phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}) &= -\sum_{i=1}^n \log\left(\frac{n}{\phi + \mathbf{h}^T \mathbf{z}_i}\right) + \phi \left(\sum_{i=1}^n \frac{1}{\phi + \mathbf{h}^T \mathbf{z}_i} - 1\right) + \sum_{i=1}^n \frac{\mathbf{h}^T \mathbf{z}_i}{\phi + \mathbf{h}^T \mathbf{z}_i}, \\ &= -\sum_{i=1}^n \log\left(\frac{n}{\phi + \mathbf{h}^T \mathbf{z}_i}\right) + n - \phi. \end{aligned} \quad (2.59)$$

Note that this dual problem implicitly require that the constraints  $\phi + \mathbf{h}^T \mathbf{z}_i > 0$  for  $i = 1, \dots, n$ , this can be guarantee by the prespecified condition and the first equation of (2.58). To maximize equation (2.59), taking derivatives with respect to  $\mathbf{u}_0$  and  $\phi$ , we have:

$$\begin{cases} \nabla_{\mathbf{u}_0} d(\phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}) = \sum_{i=1}^n \frac{-\mathbf{z}_{\mathcal{J}i}}{\phi + \mathbf{h}^T \mathbf{z}_i} = \mathbf{0}, \\ \nabla_{\phi} d(\phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}) = \sum_{i=1}^n \frac{1}{\phi + \mathbf{h}^T \mathbf{z}_i} - 1 = 0, \end{cases} \quad (2.60)$$

where  $\mathbf{z}_{\mathcal{J}i}$  is the subvector of  $\mathbf{z}_i$  indexed by  $\mathcal{J}$  and  $\mathbf{z}_{\mathcal{J}^c i}$  can be defined similarly. In addition, taking derivatives with respect to  $\mathbf{v}$ , we have:

$$\nabla_{\mathbf{v}} d(\phi, \mathbf{h}, \mathbf{u}_0, \mathbf{v}) = \sum_{i=1}^n \frac{\mathbf{z}_{\mathcal{J}^c i}}{\phi + \mathbf{h}^T \mathbf{z}_i}. \quad (2.61)$$

Not all components of equation (2.61) can be set to zeros since the dual optimization is limited to  $\mathbf{v} \succeq \mathbf{0}$ . However, we always have

$$\mathbf{v}^T \left( \sum_{i=1}^n \frac{\mathbf{z}_{\mathcal{J}^c i}}{\phi + \mathbf{h}^T \mathbf{z}_i} \right) = 0, \quad (2.62)$$

since the components of  $\mathbf{v}$  or the corresponding components of gradients  $\nabla_{\mathbf{v}} d$  are zeros when maximizing the dual function. Then by combining equation (2.58), (2.60) and (2.62), we have  $\hat{\phi} = n$  which gives that  $\hat{w}_i = 1/(n - \mathbf{u}_0^T \mathbf{z}_{\mathcal{J}i} - \mathbf{v}^T \mathbf{z}_{\mathcal{J}^c i})$ . In addition, based on the condition 2, we have  $0 < \hat{w}_i < 1$ , so the requirement of Slater's condition for the inequality of  $\mathbf{w}$  is satisfied.

Denote  $\mathbf{u} = -[\mathbf{u}_0^T, \mathbf{v}^T]^T/n$ . To better characterize the dual variable  $\mathbf{u}$ , let's assume for one realization,  $i \in \mathcal{I}$ ,  $\mathbf{A}[\mathcal{I}, ]\mathbf{s} = 0$  or  $i \in \mathcal{I}$  and for  $i \in \mathcal{I}^c$ ,  $\mathbf{A}[i, ]\mathbf{s} < 0$ , where  $\mathcal{I}$  is a subset for  $\{1, \dots, r\}$  and  $\mathcal{J}$  is a subset for  $\mathcal{I}$ ,  $\mathcal{I}^c$  is the complement of  $\mathcal{I}$  with respect to  $\{1, \dots, k\}$ , which implies  $\mathbf{u}_{\mathcal{I}^c} = \mathbf{0}$  and  $\mathbf{u}_{\mathcal{I}-\mathcal{J}} \preceq \mathbf{0}$ . The difference between  $\mathcal{I}$  and  $\mathcal{J}$  is that for  $\mathbf{A}[\mathcal{J}, ]\mathbf{s} = 0$ , the constraints are active because  $\beta^\circ$ ; while for  $\mathbf{A}[\mathcal{I} - \mathcal{J}, ]\mathbf{s} = 0$  is just because the realization of the random error. Note that the rank for  $\mathbf{A}[\mathcal{I}, ]$  is always less than or equal to the rank of the matrix  $\mathbf{A}$ , which is  $r \leq \min(k, p)$ . Thus  $\mathcal{I}$  is random and we still obtain the final distribution of the likelihood ratio statistics conditioning on each component corresponding to one case for  $\mathcal{A}$ .

By plugging  $\hat{\mathbf{u}}$  to the empirical likelihood, we can get:  $\hat{\mathbf{u}} = \operatorname{argmax}_{\mathbf{u}} \ell(\mathbf{u}) = \sum_{i=1}^n \log(1 + \mathbf{u}^T \mathbf{z}_{\mathcal{I}i})$ , under the constraints  $\mathbf{u}_{\mathcal{I}-\mathcal{J}} \preceq \mathbf{0}$ . Therefore, we want to obtain the distribution of  $c_n(\beta^\circ) = 2 \sum_{i=1}^n \log(1 + \hat{\mathbf{u}}^T \mathbf{z}_{\mathcal{I}i})$  under the cases of all possible  $\mathcal{I}$ . If there are no constraints on  $\mathbf{u}$ , we have  $\sqrt{n} \sum_{i=1}^n \frac{1}{n} \mathbf{z}_{\mathcal{I}i} \xrightarrow{d} \mathbf{A}[\mathcal{I}, ]\boldsymbol{\nu}$  and  $\sum_{i=1}^n \frac{1}{n} \mathbf{z}_{\mathcal{I}i} \mathbf{z}_{\mathcal{I}i}^T \xrightarrow{d} \mathbf{A}[\mathcal{I}, ]\mathbf{V}\mathbf{A}[\mathcal{I}, ]^T$  where  $\boldsymbol{\nu} \sim \mathcal{N}(\mathbf{0}, \mathbf{V})$ . Let  $\boldsymbol{\nu}_{\mathcal{A}\mathcal{I}} = \mathbf{A}[\mathcal{I}, ]\boldsymbol{\nu}$  and  $\mathbf{V}_{\mathcal{A}\mathcal{I}} = \boldsymbol{\nu}^T \mathbf{A}[\mathcal{I}, ]^T (\mathbf{A}[\mathcal{I}, ]\mathbf{V}\mathbf{A}[\mathcal{I}, ]^T)^{-1} \mathbf{A}[\mathcal{I}, ]$ . Then based on Lemma A.0.1 we have  $\sqrt{n}\mathbf{u} \xrightarrow{d} \mathbf{V}_{\mathcal{A}\mathcal{I}}^{-1} \boldsymbol{\nu}_{\mathcal{A}\mathcal{I}}$  and  $c_n(\beta^\circ) \xrightarrow{d} \boldsymbol{\nu}_{\mathcal{A}\mathcal{I}}^T \mathbf{V}_{\mathcal{A}\mathcal{I}}^{-1} \boldsymbol{\nu}_{\mathcal{A}\mathcal{I}}$  under this realization.

For the inequality constraints, apply Taylor expansion:

$$\begin{aligned} \mathbf{A}[\mathcal{I}^c, ]\mathbf{s} &= \mathbf{A}[\mathcal{I}^c, ]\mathbf{X}^T \mathbf{W} (\mathbf{X}\beta^\circ - \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n \frac{\mathbf{A}[\mathcal{I}^c, ]\mathbf{X}[i, ]^T (\mathbf{x}_i^T \beta^\circ - \mathbf{y}_i)}{1 + \mathbf{u}^T \mathbf{z}_{\mathcal{I}i}} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{A}[\mathcal{I}^c, ]\mathbf{X}[i, ]^T (\mathbf{x}_i^T \beta^\circ - \mathbf{y}_i) - \frac{1}{n} \sum_{i=1}^n \mathbf{A}[\mathcal{I}^c, ]\mathbf{X}[i, ]^T (\mathbf{x}_i^T \beta^\circ - \mathbf{y}_i) \mathbf{z}_{\mathcal{I}i}^T \mathbf{u} \\ &\quad + o\left(\frac{1}{n} \sum_{i=1}^n \mathbf{A}[\mathcal{I}^c, ]\mathbf{X}[i, ]^T (\mathbf{x}_i^T \beta^\circ - \mathbf{y}_i) \mathbf{z}_{\mathcal{I}i}^T \mathbf{u}\right), \end{aligned} \quad (2.63)$$

we notice that

$$\begin{aligned} &\frac{1}{n} \sum_{i=1}^n \mathbf{A}[\mathcal{I}^c, ]\mathbf{X}[i, ]^T \mathbf{z}_{\mathcal{I}i}^T \mathbf{u} (\mathbf{x}_i^T \beta^\circ - \mathbf{y}_i) \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{A}[\mathcal{I}^c, ]\mathbf{X}[i, ]^T (\mathbf{x}_i^T \beta^\circ - \mathbf{y}_i)^2 \mathbf{X}[i, ]\mathbf{A}[\mathcal{I}, ]\mathbf{u} \\ &= \hat{\mathbf{V}}_{\mathcal{I}^c\mathcal{I}} (\hat{\mathbf{V}}_{\mathcal{I}\mathcal{I}}^{-1}) \frac{1}{n} \sum_{i=1}^n \mathbf{z}_{\mathcal{I}i}, \end{aligned} \quad (2.64)$$

where  $\hat{\mathbf{V}}_{\mathcal{I}^c\mathcal{I}} = \frac{1}{n} \sum_{i=1}^n \mathbf{A}[\mathcal{I}^c, ]\mathbf{X}[i, ]^T (\mathbf{x}_i^T \beta^\circ - \mathbf{y}_i)^2 \mathbf{X}[i, ]\mathbf{A}[\mathcal{I}^c, ]^T$  and  $\hat{\mathbf{V}}_{\mathcal{I}\mathcal{I}}^{-1} = (\frac{1}{n} \sum_{i=1}^n \mathbf{z}_{\mathcal{I}i} \mathbf{z}_{\mathcal{I}i}^T)^{-1}$ . Let  $\boldsymbol{\nu}_{\mathcal{A}\mathcal{I}^c} = \mathbf{A}[\mathcal{I}^c, ]\boldsymbol{\nu}$ , by law of the large numbers, the above decomposition implies

$$\sqrt{n}\mathbf{A}[\mathcal{I}^c, ]\mathbf{s} \rightarrow \boldsymbol{\nu}_{\mathcal{A}\mathcal{I}^c} \mid \boldsymbol{\nu}_{\mathcal{A}\mathcal{I}}, \quad (2.65)$$

which is independent of  $\boldsymbol{v}_{\mathcal{I}}$ . For exactly one case of  $\mathcal{I}$ , under the following realizations, the Slater's condition is satisfied asymptotically:

$$\begin{aligned}\boldsymbol{v}_{A_{\mathcal{I}^c}} | \boldsymbol{v}_{A_{\mathcal{I}}} &\preceq \mathbf{0}, \\ \sqrt{n}\boldsymbol{u}_{\mathcal{I}-\mathcal{J}} &\preceq \mathbf{0},\end{aligned}\tag{2.66}$$

since they imply  $\mathbf{A}[\mathcal{I}^c, ]\boldsymbol{s} \preceq \mathbf{0}$  for large enough  $n$ , which satisfies the Slater's condition for the inequality constraints of  $\boldsymbol{s}$ . Based on Corollary A.0.4, we have:

$$P(\boldsymbol{v}_{A_{\mathcal{I}}}^T \mathbf{V}_{A_{\mathcal{I}}}^{-1} \boldsymbol{v}_{A_{\mathcal{I}}} \geq t, -\boldsymbol{v}_{A_{\mathcal{I}}} \succeq \mathbf{0}) = P(\boldsymbol{v}_{A_{\mathcal{I}}}^T \mathbf{V}_{A_{\mathcal{I}}}^{-1} \boldsymbol{v}_{A_{\mathcal{I}}} \geq t)P(-\boldsymbol{v}_{A_{\mathcal{I}}} \succeq \mathbf{0})\tag{2.67}$$

Consider all possible  $\mathcal{I}$ , finally we can obtain:

$$\begin{aligned}P(c_n(\boldsymbol{\beta}^\circ) \geq t) &= \sum_A P(c_n(\boldsymbol{\beta}^\circ) \geq t, \sqrt{n}\boldsymbol{u}_{\mathcal{I}-\mathcal{J}} \preceq \mathbf{0}, \sqrt{n}\mathbf{A}[\mathcal{I}^c, ]\boldsymbol{s} \preceq \mathbf{0}), \\ &\stackrel{d}{\rightarrow} \sum_{\mathcal{I}} P(\boldsymbol{v}_{A_{\mathcal{I}}}^T \mathbf{V}_{A_{\mathcal{I}}}^{-1} \boldsymbol{v}_{A_{\mathcal{I}}} \geq t, \boldsymbol{v}_{A_{\mathcal{I}-\mathcal{J}}} \preceq \mathbf{0}, \boldsymbol{v}_{A_{\mathcal{I}^c}} | \boldsymbol{v}_{A_{\mathcal{I}}} \preceq \mathbf{0}) \\ &= \sum_{\mathcal{I}} P(\boldsymbol{v}_{A_{\mathcal{I}}}^T \mathbf{V}_{A_{\mathcal{I}}}^{-1} \boldsymbol{v}_{A_{\mathcal{I}}} \geq t)P(\boldsymbol{v}_{A_{\mathcal{I}-\mathcal{J}}} \preceq \mathbf{0})P(\boldsymbol{v}_{A_{\mathcal{I}^c}} | \boldsymbol{v}_{A_{\mathcal{I}}} \preceq \mathbf{0}),\end{aligned}\tag{2.68}$$

the probability of  $\mathcal{I} = \mathcal{J}$  is 1 minus the sum of all other probabilities. Thus  $c_n(\boldsymbol{\beta}^\circ)$  follows a chi-bar-square distribution as follows:

$$\begin{aligned}P(c_n(\boldsymbol{\beta}^\circ) \geq t) &= \sum_{j=|\mathcal{J}|}^p P(\chi_j^2 \geq t)\delta_j, \\ \delta_j &= \sum_{|\mathcal{I}|=j} P(\boldsymbol{v}_{A_{\mathcal{I}-\mathcal{J}}} \preceq \mathbf{0})P(\boldsymbol{v}_{A_{\mathcal{I}^c}} | \boldsymbol{v}_{A_{\mathcal{I}}} \preceq \mathbf{0}), \text{ with } p > |\mathcal{I}|, \mathcal{J} \subsetneq \mathcal{I}, \\ \delta_{|\mathcal{J}|} &= 1 - \sum_{j=|\mathcal{J}|+1}^p \delta_j.\end{aligned}\tag{2.69}$$

## 2.8.2 Proof of Theorem 2.5.1

In this section, we prove the above main theorem. By introducing Lagrange multipliers  $\phi \in \mathbb{R}$ ,  $\boldsymbol{u}_0 \in \mathbb{R}^{|\mathcal{J}|}$  and  $\boldsymbol{h}_0 \in \mathbb{R}^p$ , the optimization problem can be formulated to the following Lagrangian form:

$$T(\boldsymbol{w}, \boldsymbol{s}, \phi, \boldsymbol{h}, \boldsymbol{u}_0) = -\langle \mathbf{1}, \log(n\boldsymbol{w}) \rangle + \phi \left( \sum_{i=1}^n w_i - 1 \right) + \boldsymbol{h}^T (-\mathbf{X}^T \mathbf{W} \boldsymbol{\epsilon} + \boldsymbol{s}) + \boldsymbol{u}_0^T \boldsymbol{s}_{\mathcal{J}},\tag{2.70}$$

and the Lagrange dual function is:

$$d(\phi, \boldsymbol{h}, \boldsymbol{u}_0) = \min_{\boldsymbol{w}, \boldsymbol{s}} T(\boldsymbol{w}, \boldsymbol{s}, \phi, \boldsymbol{h}, \boldsymbol{u}_0),\tag{2.71}$$



which is expected to be maximized. Note that the equality constraints for the primal problem are affine and we will show that the Slater's conditions for the inequality constraints hold later. Taking derivative of the Lagrangian  $T(\mathbf{w}, \phi, \mathbf{h}, \mathbf{s}, \mathbf{u}_0)$  with respect to  $\mathbf{w}$  and  $\mathbf{s}$ , we have

$$\begin{cases} -\frac{1}{\hat{w}_i} + \mathbf{h}^T \mathbf{x}_i \epsilon_i + \phi = 0, \\ \mathbf{h}_{\mathcal{J}} = \mathbf{u}_0, \\ \mathbf{h}_{\mathcal{J}^c} = \mathbf{0}. \end{cases} \quad (2.72)$$

Plugging equation (2.72) into the Lagrangian to get the Lagrange dual function, we have:

$$\begin{aligned} d(\phi, \mathbf{u}_0) &= -\sum_{i=1}^n \log\left(\frac{n}{\phi + \mathbf{u}_0^T \mathbf{z}_i}\right) + \phi\left(\sum_{i=1}^n \frac{1}{\phi + \mathbf{z}_i} - 1\right) + \sum_{i=1}^n \frac{\mathbf{u}_0^T \mathbf{z}_i}{\phi + \mathbf{u}_0^T \mathbf{z}_i} \\ &= -\sum_{i=1}^n \log\left(\frac{n}{\phi + \mathbf{u}_0^T \mathbf{z}_i}\right) + n - \phi. \end{aligned} \quad (2.73)$$

Note that this dual problem implicitly require that the constraints  $\phi + \mathbf{u}_0^T \mathbf{z}_i > 0$  for  $i = 1, \dots, n$ , but we will show our prespecified conditions could guarantee them later. To maximize equation (2.73), taking derivatives with respect to  $\mathbf{u}_0$  and  $\phi$ , we have:

$$\begin{cases} \nabla_{\mathbf{u}_0} d(\phi, \mathbf{u}_0) = \sum_{i=1}^n \frac{\mathbf{z}_i}{\phi + \mathbf{u}_0^T \mathbf{z}_i} = 0, \\ \nabla_{\phi} d(\phi, \mathbf{h}) = \sum_{i=1}^n \frac{1}{\phi + \mathbf{u}_0^T \mathbf{z}_i} - 1 = 0, \end{cases} \quad (2.74)$$

which gives  $\phi = n$  and implies that  $\hat{w}_i = 1/(n + \mathbf{u}_0^T \mathbf{z}_i)$ . Let  $\mathbf{u} = n\mathbf{u}_0$ , the dual optimization becomes

$$\max_{\mathbf{u}} \sum_{i=1}^n \log(1 + \mathbf{u}^T \mathbf{z}_i) \quad \text{s.t.} \quad 0 \leq \hat{w}_i(\mathbf{u}) \leq 1, i = 1, \dots, n; \quad \|\hat{\mathbf{s}}(\mathbf{u})\|_{\infty} \leq \lambda. \quad (2.75)$$

Now we denote the unconstrained maximizer of  $\max_{\mathbf{u}} \sum_{i=1}^n \log(1 + \mathbf{u}^T \mathbf{z}_i)$  as  $\tilde{\mathbf{u}}$ , where  $\tilde{\mathbf{u}}$  is a traditional  $M$ -estimator for the large  $n$  case. Then based on the condition 3, the interior of the convex hull of  $\mathbf{z}_1, \dots, \mathbf{z}_n$  contains zero, which implies that  $0 < \hat{w}_i(\tilde{\mathbf{u}}) < 1, i = 1, \dots, n$ . In addition, denote  $\mathbf{r}(\tilde{\mathbf{u}}) = \hat{\mathbf{s}}(\tilde{\mathbf{u}})_{\mathcal{J}^c} - \mathbf{X}_{\mathcal{J}^c}^T \boldsymbol{\epsilon}/n$ , then we have

$$\begin{aligned} \hat{\mathbf{s}}(\tilde{\mathbf{u}})_{\mathcal{J}^c} &= \mathbf{X}_{\mathcal{J}^c}^T \hat{\mathbf{W}}(\tilde{\mathbf{u}}) \boldsymbol{\epsilon} = \frac{1}{n} \sum_{i=1}^n \frac{\mathbf{X}_{i, \mathcal{J}^c}^T \boldsymbol{\epsilon}_i}{1 + \tilde{\mathbf{u}}^T \mathbf{z}_i} \\ &= \frac{1}{n} \mathbf{X}_{\mathcal{J}^c}^T (\mathbf{I} - \mathbf{P}_{\mathbf{X}_{\mathcal{J}}}) \boldsymbol{\epsilon} + \frac{1}{n} \mathbf{X}_{\mathcal{J}^c}^T \mathbf{P}_{\mathbf{X}_{\mathcal{J}}} \boldsymbol{\epsilon} + \mathbf{r}(\tilde{\mathbf{u}}). \end{aligned} \quad (2.76)$$

Based on the asymptotic property of  $\tilde{\mathbf{u}}$  (Lemma A.0.1) such that  $\max_i(\tilde{\mathbf{u}}^T \mathbf{z}_i) = o_p(1)$  we have  $\|\mathbf{r}(\tilde{\mathbf{u}})\|_{\infty} = o_p(\|\mathbf{X}_{\mathcal{J}^c}^T \boldsymbol{\epsilon}/n\|_{\infty})$ . Moreover, based on the selection event induced condition 2, we have

$$\left\| \frac{1}{n} \mathbf{X}_{\mathcal{J}^c}^T (\mathbf{I} - \mathbf{P}_{\mathbf{X}_{\mathcal{J}}}) \boldsymbol{\epsilon} \right\|_{\infty} < \lambda. \quad (2.77)$$

For the last term  $\frac{1}{n} \mathbf{X}_{\mathcal{J}^c}^T \mathbf{P}_{\mathbf{X}_{\mathcal{J}}} \boldsymbol{\epsilon}$ , since  $\left\| \frac{1}{n} \mathbf{X}_{\mathcal{J}^c}^T \mathbf{P}_{\mathbf{X}_{\mathcal{J}}} \boldsymbol{\epsilon} \right\|_{\infty} \leq \left\| \mathbf{X}_{\mathcal{J}^c}^T \mathbf{X}_{\mathcal{J}} (\mathbf{X}_{\mathcal{J}}^T \mathbf{X}_{\mathcal{J}})^{-1} \right\|_{\infty} \left\| \mathbf{X}_{\mathcal{J}}^T \boldsymbol{\epsilon} \right\|_{\infty}$ , based on Lemma A.0.2, since  $\mathbf{X}_{\mathcal{J}} \boldsymbol{\epsilon}$  is a  $|\mathcal{J}|$ -dimensional vector with finite variance, with probability converging to 1, there is a constant  $C$  such that

$$\left\| \frac{1}{n} \mathbf{X}_{\mathcal{J}^c}^T \mathbf{P}_{\mathbf{X}_{\mathcal{J}}} \boldsymbol{\epsilon} \right\|_{\infty} \leq C\delta/n^{1/2} = o(\lambda). \quad (2.78)$$

Therefore with probability one, we have,

$$\left\| \hat{\mathbf{s}}(\tilde{\mathbf{u}})_{\mathcal{J}^c} \right\|_{\infty} < \lambda + o(\lambda) \leq \lambda. \quad (2.79)$$

Since we also have  $\hat{\mathbf{s}}_{\mathcal{J}} = \mathbf{0}$ , the above analysis implies that the inequality constraints  $\left\| \hat{\mathbf{s}}(\tilde{\mathbf{u}}) \right\|_{\infty} \leq \lambda$  holds automatically for the unconstrained optimizer  $\tilde{\mathbf{u}}$  when  $n$  is sufficiently large. Note that both  $0 \leq w_i \leq 1$  for  $i = 1, 2, \dots, n$  and  $\left\| \mathbf{s} \right\|_{\infty} \leq \lambda$  are all affine inequality constraints, so the existence of  $\tilde{\mathbf{u}}$  such that  $0 \leq \hat{w}_i(\tilde{\mathbf{u}}) \leq 1, i = 1, \dots, n$  and  $\left\| \hat{\mathbf{s}}(\tilde{\mathbf{u}}) \right\|_{\infty} \leq \lambda$  hold guarantees the requirement for inequality constraints of the Slater's condition. Hence, strict feasibility holds and the duality gap is zero, which means the transformation from prime to dual is equivalent. Moreover, when  $\min_{j \in \mathcal{J}} |\beta_j^{\circ}| > (1 + \kappa)\lambda$ , the event  $\{|\beta_{\mathcal{J}}^{\circ} + (\mathbf{X}_{\mathcal{J}}^T \mathbf{X}_{\mathcal{J}})^{-1} \mathbf{X}_{\mathcal{J}}^T \boldsymbol{\epsilon}| \succeq \lambda \mathbf{1}\}$  holds with probability approaching one since  $\left\| \mathbf{X}_{\mathcal{J}}^T \boldsymbol{\epsilon} / n \right\|_{\infty} = O(1/n) = o(\lambda)$  and  $(\mathbf{X}_{\mathcal{J}}^T \mathbf{X}_{\mathcal{J}})^{-1} / n$  is positive definite and finite. Finally, given any  $t > 0$ , based on asymptotic analysis of traditional EL, we have  $c_n(\beta^{\circ}) = \boldsymbol{\epsilon}^T \mathbf{X}_{\mathcal{J}} (\mathbf{X}_{\mathcal{J}} \mathbf{D}(\boldsymbol{\epsilon})^2 \mathbf{X}_{\mathcal{J}})^{-1} \mathbf{X}_{\mathcal{J}}^T \boldsymbol{\epsilon} + o_p(1)$ , so the obtained conditional critical probability can be obtained

$$\begin{aligned} P(c_n(\beta^{\circ}) \geq t \mid \{|\beta_{\mathcal{J}}^{\circ} + (\mathbf{X}_{\mathcal{J}}^T \mathbf{X}_{\mathcal{J}})^{-1} \mathbf{X}_{\mathcal{J}}^T \boldsymbol{\epsilon}| \succeq \lambda \mathbf{1}, \text{ and } \left\| \mathbf{X}_{\mathcal{J}^c}^T (\mathbf{I} - \mathbf{P}_{\mathbf{X}_{\mathcal{J}}}) \boldsymbol{\epsilon} \right\|_{\infty} < n\lambda\}) \\ = P(\mathbf{z}_{\lambda}^{\mathbf{X}, \beta^{\circ}} \geq t + o_p(1)), \end{aligned} \quad (2.80)$$

and the finally conclusion can be obtained if the distribution of  $\mathbf{z}_{\lambda}^{\mathbf{X}, \beta^{\circ}}$  is continuous at  $t$ .

### 2.8.3 Proof of Theorem 2.5.2

$$\begin{aligned}
-\langle \mathbf{1}, \log(n\mathbf{w}') \rangle &= -\langle \mathbf{1}, \log(n\hat{\mathbf{w}}) \rangle - \sum_{i=1}^n \log\left(\frac{w'_i}{\hat{w}_i}\right) \\
&< -\langle \mathbf{1}, \log(n\hat{\mathbf{w}}) \rangle - \sum_{i=1}^n \left(\frac{w'_i}{\hat{w}_i} - 1\right) \\
&= -\langle \mathbf{1}, \log(n\hat{\mathbf{w}}) \rangle - \sum_{i=1}^n (w'_i(n - \mathbf{u}_0^T \mathbf{z}_{\mathcal{J}i} \\
&\quad + \mathbf{v}^{-T}(\mathbf{z}_{\mathcal{J}c_i} - \lambda) + \mathbf{v}^{-T}(-\mathbf{z}_{\mathcal{J}c_i} - \lambda)) - 1) \\
&= -\langle \mathbf{1}, \log(n\hat{\mathbf{w}}) \rangle - \sum_{i=1}^n (nw'_i - 1) - \mathbf{u}_0^T \left(\sum_{i=1}^n w'_i \mathbf{z}_{\mathcal{J}i}\right) \\
&\quad - \mathbf{v}^{-T} \left(\sum_{i=1}^n w'_i (\mathbf{z}_{\mathcal{J}c_i} - \lambda)\right) - \mathbf{v}^{-T} \left(\sum_{i=1}^n w'_i (-\mathbf{z}_{\mathcal{J}c_i} - \lambda)\right) \\
&\leq -\langle \mathbf{1}, \log(n\hat{\mathbf{w}}) \rangle,
\end{aligned}$$

where the first inequality follows by the fact that  $\hat{\mathbf{w}}$  and  $\mathbf{w}'$  are different and the last inequality holds based on  $\mathbf{v}^- \succeq \mathbf{0}$  and  $\mathbf{v}^+ \succeq \mathbf{0}$ .

# APPENDIX A

## TECHNICAL LEMMAS

**Lemma A.0.1.** *Suppose  $\mathbf{z}_i \in \mathbb{R}^p$ ,  $i = 1, \dots, n$ ,  $p < n$  are independent and mean 0 and finite and positive definite variance  $\mathbf{V}$ , denote*

$$\hat{\mathbf{h}} = \operatorname{argmax}_{\mathbf{h}} \sum_{i=1}^n \log(1 + \mathbf{h}^T \mathbf{z}_i) \quad (\text{A.1})$$

$$c_n = 2 \sum_{i=1}^n \log(1 + \hat{\mathbf{h}}^T \mathbf{z}_i).$$

Then we have  $\hat{\mathbf{h}} \rightarrow (\sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^T)^{-1} \sum_{i=1}^n \mathbf{z}_i$  and  $c_n \rightarrow (\sum_{i=1}^n \mathbf{z}_i^T) (\sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^T)^{-1} (\sum_{i=1}^n \mathbf{z}_i)$ .

*Proof.* To solve optimization (A.1), we take the derivative and find that  $\hat{\mathbf{h}}$  should satisfy

$$\sum_{i=1}^n \frac{\mathbf{z}_i}{1 + \mathbf{h}^T \mathbf{z}_i} = \mathbf{0}. \quad (\text{A.2})$$

Then we need to show  $\hat{\mathbf{h}}$  is root- $n$  consistency with respect to zero  $\|\hat{\mathbf{h}}\|_2 = O_p(n^{-1/2})$ . Let  $\mathbf{g}_i(\mathbf{h}) = \mathbf{z}_i / (1 + \mathbf{h}^T \mathbf{z}_i)$ , then  $\hat{\mathbf{h}}$  satisfies  $\sum_{i=1}^n \mathbf{g}_i(\hat{\mathbf{h}}) = \mathbf{0}$ .

We consider to directly show the root- $n$  consistency of  $\hat{\mathbf{h}}$ . Denote  $d(\mathbf{h}) = \sum_{i=1}^n \log(1 + \mathbf{h}^T \mathbf{z}_i)$ . Then, we want to show that for any given  $\epsilon > 0$ , there exists a large constant  $C$  such that

$$P\left\{ \sup_{\|\mathbf{u}\|_2=C} d(n^{-\frac{1}{2}}\mathbf{u}) - d(\mathbf{0}) < 0 \right\} \geq 1 - \epsilon. \quad (\text{A.3})$$

This implies that with probability at least  $1 - \epsilon$ , there exists a local maximum in the ball  $\{n^{-1/2}\mathbf{u} : \|\mathbf{u}\|_2 \leq C\}$ . Hence, there is a local maximizer of  $d(\mathbf{h})$  such that  $\|\hat{\mathbf{h}}\|_2 = O_p(n^{-1/2})$ . By Taylor's expansion (the first Taylor expansion), there is  $0 < t < 1$  such that

$$\begin{aligned} d(n^{-\frac{1}{2}}\mathbf{u}) - d(\mathbf{0}) &= \frac{1}{n^{1/2}} \sum_{i=1}^n \mathbf{z}_i^T \mathbf{u} + \left(-\frac{1}{2n} \sum_{i=1}^n \mathbf{u}^T \mathbf{z}_i \mathbf{z}_i^T \mathbf{u} + \frac{1}{2} \mathbf{u}^T \mathbf{V} \mathbf{u}\right) \\ &\quad + \left(-\frac{1}{2} \mathbf{u}^T \mathbf{V} \mathbf{u}\right) + \left(-\frac{1}{2n} \sum_{i=1}^n \mathbf{u}^T \nabla \mathbf{g}_i(tn^{-\frac{1}{2}}\mathbf{u})^T \mathbf{u} + \frac{1}{2n} \sum_{i=1}^n \mathbf{u}^T \mathbf{z}_i \mathbf{z}_i^T \mathbf{u}\right) \\ &\equiv I_1 + I_2 + I_3 + I_4. \end{aligned} \quad (\text{A.4})$$

By central limit theorem,  $\sum_{i=1}^n \|\mathbf{z}_i\|_2/\sqrt{n} = O_p(1)$ , so the term  $I_1$  is upper bounded by  $CC'$ , where  $C' > 0$  is a constant chosen based on  $\epsilon$ ; while the second term  $I_2$  can be smaller than  $C$  with probability tending to 1 by WLLN.  $I_3$  is upper bounded by  $-C^2\lambda_{\min}(\mathbf{V})$ , which is a negative constant based on the positive definite of  $\mathbf{V}$ . It remains to show  $I_4 = o_p(1)$ .

$$\begin{aligned}
|I_4| &= \left| \frac{1}{2n} \sum_{i=1}^n \mathbf{u}^T \nabla g_i(tn^{-1/2}\mathbf{u})^T \mathbf{u} - \frac{1}{2n} \sum_{i=1}^n \mathbf{u}^T \mathbf{z}_i \mathbf{z}_i^T \mathbf{u} \right| \\
&= \left| \frac{1}{2n} \sum_{i=1}^n -\frac{\mathbf{u}^T \mathbf{z}_i \mathbf{z}_i^T \mathbf{u}}{(1 + tn^{-1/2}\mathbf{u}^T \mathbf{z}_i)^2} - \frac{1}{2n} \sum_{i=1}^n \mathbf{u}^T \mathbf{z}_i \mathbf{z}_i^T \mathbf{u} \right| \\
&= \left| \frac{1}{2n} \sum_{i=1}^n \frac{tn^{-1/2}\mathbf{u}^T \mathbf{z}_i \mathbf{u}^T \mathbf{z}_i \mathbf{z}_i^T \mathbf{u}}{(1 + tn^{-1/2}\mathbf{u}^T \mathbf{z}_i)^2} \right| + \left| \frac{1}{2n} \sum_{i=1}^n \frac{tn^{-1/2}\mathbf{u}^T \mathbf{z}_i \mathbf{u}^T \mathbf{z}_i \mathbf{z}_i^T \mathbf{u}}{1 + tn^{-1/2}\mathbf{u}^T \mathbf{z}_i} \right| \\
&\leq C^3 \frac{1}{2n} \sum_{i=1}^n \frac{n^{-1/2}\|\mathbf{z}_i\|_2^3}{(1 + tn^{-1/2}\mathbf{u}^T \mathbf{z}_i)^2} + C^3 \frac{1}{2n} \sum_{i=1}^n \frac{n^{-1/2}\|\mathbf{z}_i\|_2^3}{|1 + tn^{-1/2}\mathbf{u}^T \mathbf{z}_i|}.
\end{aligned} \tag{A.5}$$

Based on the finite second moment and i.i.d. of  $\mathbf{z}_i, i = 1, 2, \dots, n$ , we have  $\frac{1}{n} \sum_{i=1}^n \|\mathbf{z}_i\|_2^3 = o_p(n^{1/2})$  and  $\max_{1 \leq i \leq n} \|\mathbf{z}_i\|_2 = o_p(n^{1/2})$  according to Lemma A.0.2. Denote  $z_{\max} = \max_{1 \leq i \leq n} \|\mathbf{z}_i\|_2$ . Then for sufficient large  $n$ , the right hand side of (A.6) can be bounded by

$$C^3 \frac{1}{2n} \sum_{i=1}^n \frac{n^{-1/2}\|\mathbf{z}_i\|_2^3}{1 - n^{-1/2}Cz_{\max}} + C^3 \frac{1}{2n} \sum_{i=1}^n \frac{n^{-1/2}\|\mathbf{z}_i\|_2^3}{(1 - n^{-1/2}Cz_{\max})^2} = o_p(1).$$

Hence the term  $I_4$  can be less than  $C$  with probability tending to 1. To summarize, we have  $I_1 + I_3 + I_3 + I_4 \leq C(C' + 2) - C^2\lambda_{\min}(\mathbf{V}) < 0$  by choosing a large enough constant  $C$  such that  $C > (C' + 2)/\lambda_{\min}(\mathbf{V})$ . Hence,  $d(n^{-1/2}\mathbf{u}) - d(\mathbf{0}) < 0$  with probability at least  $1 - \epsilon$ . This implies that there is a local maximum with in the ball  $\{n^{-1/2}\mathbf{u} : \|\mathbf{u}\|_2 \leq C\}$  with probability  $1 - \epsilon$ . Therefore, by picking up the root closed to zero as  $\hat{\mathbf{h}}$ , we have  $\|\hat{\mathbf{h}}\| = O_p(n^{-1/2})$  and  $\mathbf{z}_i^T \hat{\mathbf{h}} \leq z_{\max}\|\hat{\mathbf{h}}\| = o_p(1)$  for  $i = 1, 2, \dots, n$ .

After obtaining root- $n$  consistency of  $\hat{\mathbf{h}}$ , we can expand the estimating equation (A.2) for  $\mathbf{z}_i^T \mathbf{h}$  (the second Taylor expansion), which means that there exist constants  $0 \leq t_i \leq 1$  for  $i = 1, \dots, n$ , such that

$$\sum_{i=1}^n \mathbf{z}_i - \sum_{i=1}^n \mathbf{z}_i \hat{\mathbf{h}}^T \mathbf{z}_i + \sum_{i=1}^n \mathbf{z}_i (t_i \hat{\mathbf{h}}^T \mathbf{z}_i)^2 = \mathbf{0}. \tag{A.6}$$

Then by defining

$$\mathbf{r} = \left( \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^T \right)^{-1} \sum_{i=1}^n \mathbf{z}_i (t_i \hat{\mathbf{h}}^T \mathbf{z}_i)^2, \tag{A.7}$$

we have

$$\hat{\mathbf{h}} = \left( \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^T \right)^{-1} \sum_{i=1}^n \mathbf{z}_i + \mathbf{r}, \quad (\text{A.8})$$

with  $\|\mathbf{r}\|_2 = o(n^{-1/2})$ . We define  $\mathbf{v}$  as a normal random variable such that we have  $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{V})$ . Then by plugging the form of  $\hat{\mathbf{h}}$  to the EL statistics, there exist  $0 \leq t'_i \leq 1, i = 1, 2, \dots, n$ , such that (the third Taylor expansion):

$$\begin{aligned} -2\langle \mathbf{1}, \log(n\hat{\mathbf{w}}) \rangle &= 2 \sum_{i=1}^n \log(1 + \hat{\mathbf{h}}^T \mathbf{z}_i) \\ &= 2 \sum_{i=1}^n \left\{ \hat{\mathbf{h}}^T \mathbf{z}_i - \frac{1}{2} \hat{\mathbf{h}}^T \mathbf{z}_i \mathbf{z}_i^T \hat{\mathbf{h}} + \frac{1}{3} (t'_i \hat{\mathbf{h}}^T \mathbf{z}_i)^3 \right\} \\ &= \left( \sum_{i=1}^n \mathbf{z}_i \right)^T \left( \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^T \right)^{-1} \left( \sum_{i=1}^n \mathbf{z}_i \right) + 2\mathbf{r}^T \sum_{i=1}^n \mathbf{z}_i \\ &\quad - 2\mathbf{r}^T \sum_{i=1}^n \mathbf{z}_i - \mathbf{r}^T \left( \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^T \right) \mathbf{r} + \frac{2}{3} \sum_{i=1}^n (t'_i \hat{\mathbf{h}}^T \mathbf{z}_i)^3 \\ &\stackrel{(i)}{=} \left( \sqrt{n} \sum_{i=1}^n \frac{1}{n} \mathbf{z}_i \right)^T \left( \sum_{i=1}^n \frac{1}{n} \mathbf{z}_i \mathbf{z}_i^T \right)^{-1} \left( \sqrt{n} \sum_{i=1}^n \frac{1}{n} \mathbf{z}_i \right) + o_p(1) \end{aligned} \quad (\text{A.9})$$

where (i) is because  $\mathbf{r}^T \left( \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^T \right) \mathbf{r} = o_p(1)$  based on law of large numbers and the finiteness of  $\mathbf{V}$  and  $\| \sum_{i=1}^n (t'_i \hat{\mathbf{h}}^T \mathbf{z}_i)^3 \|_2 \leq z_{\max} \|\hat{\mathbf{h}}\|_2 \sum_{i=1}^n \hat{\mathbf{h}}^T \mathbf{z}_i \mathbf{z}_i^T \hat{\mathbf{h}} = o_p(1)$ .  $\square$

**Lemma A.0.2.** (Lemma 3, [14]) Let  $\xi_i \geq 0$  be i.i.d. random variables and define  $\xi_{\max} = \max_{1 \leq i \leq n} \xi_i$ . If  $E(\xi_1^2) < \infty$ , then

$$\xi_{\max} = o(n^{\frac{1}{2}}) \quad \text{and} \quad \sum_{i=1}^n \frac{\xi_i^3}{n} = o(n^{\frac{1}{2}}). \quad (\text{A.10})$$

**Lemma A.0.3.** ([18]) Suppose that  $\mathbf{y}$  is  $N(\mathbf{0}, \mathbf{I})$ ,  $\mathbf{Q}$  is an idempotent symmetric matrix of rank  $r$  and  $\mathcal{A}$  is a cone defined by

$$\mathcal{A} = \{ \mathbf{v} \in \mathbb{R}^q : \mathbf{d}_i^T \mathbf{v} \leq 0, i = 1, \dots, k \}, \quad (\text{A.11})$$

where  $\mathbf{d}_i (i = 1, \dots, k)$  are vectors such that either  $\mathbf{Q}\mathbf{d}_i = \mathbf{0}$  or  $\mathbf{Q}\mathbf{d}_i = \mathbf{d}_i$ . Then the conditional distribution of  $\mathbf{y}^T \mathbf{Q} \mathbf{y}$ , given  $\mathbf{y} \in \mathcal{A}$ , is that of a chi-squared random variable with  $r$  degrees of freedom.  $\square$

Based on this Lemma, we can easily get the following Corollary:

**Corollary A.0.4.** *Suppose  $\mathbf{v} \in \mathbb{R}^p$ ,  $\mathbf{v} \sim N(\mathbf{0}, \mathbf{V})$  and  $\mathbf{V}$  is positive definite, then for any  $t$ , we have:*

$$P(\mathbf{v}^T \mathbf{V}^{-1} \mathbf{v} \geq t, \mathbf{V}^{-1} \mathbf{v} \succeq 0) = P(\mathbf{v}^T \mathbf{V}^{-1} \mathbf{v} \geq t)P(\mathbf{V}^{-1} \mathbf{v} \succeq 0). \quad (\text{A.12})$$

*Proof.* Let  $\mathbf{y} = \mathbf{V}^{-\frac{1}{2}} \mathbf{v}$ , and  $\mathbf{Q} = \mathbf{I}_p$ , then directly apply Lemma A.0.3. □

# BIBLIOGRAPHY

- [1] Peter Bühlmann, Markus Kalisch, and Lukas Meier. High-dimensional statistics with a view toward applications in biology. 2014.
- [2] T. T. Cai, G. Xu, and J. Zhang. On recovery of sparse signals via  $\ell_1$  minimization. *IEEE Transactions on Information Theory*, 55(7):3388–3397, July 2009.
- [3] Ruben Dezeure, Peter Bühlmann, and Cun-Hui Zhang. High-dimensional simultaneous inference with the bootstrap. *Test*, 26(4):685–719, 2017.
- [4] Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, 96(456):1348–1360, 2001.
- [5] Frank R. Hampel, Elvezio M. Ronchetti, Peter J. Rousseeuw, and Werner A. Stahel. *Robust statistics*. John Wiley & Sons, New York, 2005.
- [6] Holger Höfling and Robert Tibshirani. Estimation of sparse binary pairwise markov networks using pseudo-likelihoods. *Journal of Machine Learning Research*, 10(Apr):883–906, 2009.
- [7] Dimitris Karlis. An EM algorithm for multivariate Poisson distribution and related models. *Journal of Applied Statistics*, 30(1):63–77, 2003.
- [8] Jason D Lee, Dennis L Sun, Yuekai Sun, Jonathan E Taylor, et al. Exact post-selection inference, with application to the lasso. *The Annals of Statistics*, 44(3):907–927, 2016.
- [9] Kung-Yee Liang and Scott L Zeger. Longitudinal data analysis using generalized linear models. *Biometrika*, 73(1):13–22, 1986.
- [10] Po-Ling Loh and Martin J Wainwright. Structure estimation for discrete graphical models: Generalized covariance matrices and their inverses. *THE ANNALS of STATISTICS*, pages 3022–3049, 2013.
- [11] Willard G Manning, Anirban Basu, and John Mullahy. Generalized modeling approaches to risk adjustment of skewed outcomes data. *Journal of health economics*, 24(3):465–488, 2005.
- [12] Hamed Masnadi-shirazi and Nuno Vasconcelos. On the design of loss functions for classification: theory, robustness to outliers, and savageboost. In *Advances in Neural Information Processing Systems 21*, pages 1049–1056. 2009.
- [13] Art B Owen. Empirical likelihood ratio confidence intervals for a single functional. *Biometrika*, 75(2):237–249, 1988.



- [14] Art B Owen. Empirical likelihood ratio confidence regions. *The Annals of Statistics*, 18(1):90–120, 1990.
- [15] Art B Owen. Empirical likelihood for linear models. *The Annals of Statistics*, 19(4):1725–1747, 1991.
- [16] Jin Qin and Jerry Lawless. Empirical likelihood and general estimating equations. *The Annals of Statistics*, 22(1):300–325, 1994.
- [17] Pradeep Ravikumar, Martin J. Wainwright, and John D. Lafferty. High-dimensional Ising model selection using  $\ell_1$ -regularized logistic regression. *Ann. Statist.*, 38(3):1287–1319, 06 2010.
- [18] Alexander Shapiro. Asymptotic distribution of test statistics in the analysis of moment structures under inequality constraints. *Biometrika*, 72(1):133–144, 1985.
- [19] Y. She and K. Chen. Robust reduced-rank regression. *Biometrika*, 104(3):633–647, 2017.
- [20] Yiyuan She. An iterative algorithm for fitting nonconvex penalized generalized linear models with grouped predictors. *Computational Statistics & Data Analysis*, 56(10):2976–2990, October 2012.
- [21] Yiyuan She. On the finite-sample analysis of  $\Theta$ -estimators. *Electronic Journal of Statistics*, 10(2):1874–1895, 2016.
- [22] Yiyuan She and Kun Chen. Robust reduced rank regression. *Biometrika*, 104(3):633–647, 2017.
- [23] Yiyuan She, Shao Tang, and Qiaoya Zhang. Indirect gaussian graph learning beyond gaussianity. *IEEE Transactions on Network Science and Engineering*, 2019.
- [24] Auke Tellegen, David Watson, and Lee Anna Clark. On the dimensional and hierarchical structure of affect. *Psychological Science*, 10(4):297–303, 1999.
- [25] Konstantinos Trohidis, Grigorios Tsoumakas, George Kalliris, and Ioannis P Vlahavas. Multi-label classification of music into emotions. In *ISMIR*, volume 8, pages 325–330, 2008.
- [26] Sara Van de Geer, Peter Bühlmann, Ya’acov Ritov, Ruben Dezeure, et al. On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42(3):1166–1202, 2014.
- [27] Lan Wang, Jianhui Zhou, and Annie Qu. Penalized generalized estimating equations for high-dimensional longitudinal data analysis. *Biometrics*, 68(2):353–360, 2012.

- [28] Lingzhou Xue, Hui Zou, and Tianxi Cai. Nonconcave penalized composite conditional likelihood estimation of sparse Ising models. *Ann. Statist.*, 40(3):1403–1429, 06 2012.
- [29] Eunho Yang, Genevera Allen, Zhandong Liu, and Pradeep K Ravikumar. Graphical models via generalized linear models. In *Advances in Neural Information Processing Systems*, pages 1358–1366, 2012.
- [30] Eunho Yang, Pradeep Ravikumar, Genevera I. Allen, and Zhandong Liu. Graphical models via univariate exponential family distributions. *Journal of Machine Learning Research*, 16(1):3813–3847, January 2015.
- [31] Scott L Zeger, Kung-Yee Liang, and Paul S Albert. Models for longitudinal data: a generalized estimating equation approach. *Biometrics*, pages 1049–1060, 1988.
- [32] Cun-Hui Zhang and Stephanie S Zhang. Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):217–242, 2014.

## BIOGRAPHICAL SKETCH

Peng Zhao was born in Henan, China. He received his B.S. in Statistics from Beijing Institute of Technology in 2015. In 2017, he attained a M.S. in Statistics from Florida State University and then pursue a doctorate in Statistics.