

Multitask Quantile Regression Under the Transnormal Model

Jianqing Fan^a, Lingzhou Xue^b, and Hui Zou^c

^aDepartment of Operations Research and Financial Engineering, Princeton University, Princeton, NJ, USA; ^bDepartment of Statistics, Penn State University, University Park, PA, USA; ^cSchool of Statistics, University of Minnesota, Minneapolis, MN, USA

ABSTRACT

We consider estimating multitask quantile regression under the transnormal model, with focus on high-dimensional setting. We derive a surprisingly simple closed-form solution through rank-based covariance regularization. In particular, we propose the rank-based ℓ_1 penalization with positive-definite constraints for estimating sparse covariance matrices, and the rank-based banded Cholesky decomposition regularization for estimating banded precision matrices. By taking advantage of the alternating direction method of multipliers, nearest correlation matrix projection is introduced that inherits sampling properties of the unprojected one. Our work combines strengths of quantile regression and rank-based covariance regularization to simultaneously deal with nonlinearity and nonnormality for high-dimensional regression. Furthermore, the proposed method strikes a good balance between robustness and efficiency, achieves the “oracle”-like convergence rate, and provides the provable prediction interval under the high-dimensional setting. The finite-sample performance of the proposed method is also examined. The performance of our proposed rank-based method is demonstrated in a real application to analyze the protein mass spectroscopy data. Supplementary materials for this article are available online.

ARTICLE HISTORY

Received December 2013
Revised July 2015

KEYWORDS

Copula model; Optimal transformation; Rank correlation; Cholesky decomposition; Quantile regression; Prediction interval; Alternating direction method of multipliers

1. Introduction

Consider a multitask high-dimensional learning paradigm that independent variables $\mathbf{z} = (z_1, \dots, z_{p_0})'$ are used to simultaneously predict multiple response variables $\mathbf{y} = (y_1, \dots, y_{q_0})'$. Here, dimensions p_0 and q_0 can be of larger order of magnitude than sample size n . This work is motivated by obtaining the optimal prediction $\mathbf{t}(\mathbf{z}) = (t_1(\mathbf{z}), \dots, t_{q_0}(\mathbf{z}))'$ for \mathbf{y} . To this end, we shall solve optimal transformations \mathbf{t} from

$$\min_{\mathbf{t}: \mathbb{R}^{p_0} \mapsto \mathbb{R}^{q_0}} \sum_{j=1}^{q_0} \mathbb{E}[L(y_j - t_j(\mathbf{z}))],$$

where $L(\cdot)$ is some convex loss function. If \mathbf{z} and \mathbf{y} have a joint normal distribution, nice properties such as linearity and homoscedasticity hold for \mathbf{t} . Under normality, it is appropriate to specify $L(\cdot)$ as the squared loss, that is, $L(u) = u^2$. Thus, normality makes the neat connection between optimal prediction and ordinary least squares, and optimal transformations can be easily obtained by solving ordinary least squares.

However, observed data are often skewed or heavy-tailed, and rarely normally distributed in real-world applications. Transformations are usually used to achieve normality in regression analysis. Under the traditional low-dimensional setting, estimating transformations for regression has received considerable attention. On one hand, parametric methods focus on the parametric families of transformations, such as the celebrated Box–Cox transformation (Box and Cox 1964). On the other hand, nonparametric estimation of regression transformations is also studied, for instance, projection pursuit

regression (Friedman and Stuetzle 1981), alternating conditional expectation (Breiman and Friedman 1985), additivity and variance stabilization (Stone 1985; Tibshirani 1988), among others. However, it is nontrivial to extend them to the high-dimensional setting, since these methods would suffer from the curse of dimensionality. Without requiring the restricted normality, there are significant demands for estimating optimal transformations for high-dimensional regression.

As a nice combination of flexibility and interpretability, the Gaussian copulas have generated a lot of interests in statistics and econometrics, and are deemed as a favorable alternative to the Gaussian model in high-dimensional statistical problems, including linear or quadratic discriminant analysis (Lin and Jeon 2003; Mai and Zou 2015; Fan et al. 2013), graphical modeling (Liu, Lafferty, and Wasserman 2009; Liu et al. 2012; Xue and Zou 2012), covariance matrix estimation (Xue and Zou 2014; Wegkamp and Zhao 2016), and principal component analysis (Han and Liu 2012). The semiparametric Gaussian copula provides a semiparametric generalization of the Gaussian model by assuming the existence of univariate monotone transformations f for $\mathbf{x} \in \mathbb{R}^p$ such that $f(\mathbf{x}) \sim N_p(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$. Throughout this article, we follow Lin and Jeon (2003) to call this copula model as the transnormal model, which is also called the nonparanormal model in Liu, Lafferty, and Wasserman (2009). We will show the power of the transnormal model in estimating optimal transformations for high-dimensional regression. We suppose that $\mathbf{x} = (\mathbf{z}', \mathbf{y}')'$ consists of $\mathbf{z} = (x_1, \dots, x_{p_0})'$ and $\mathbf{y} = (x_{p_0+1}, \dots, x_{p_0+q_0})'$. Denote by F_j the marginal cumulative distribution function (CDF) of x_j and $f_j = \Phi^{-1} \cdot F_j$ for $j = 1, \dots, p_0 + q_0$. Let $\mathbf{g} =$

$(f_1, \dots, f_{p_0})'$ and $\mathbf{h} = (f_{p_0+1}, \dots, f_{p_0+q_0})'$. The transnormal model entails the existence of monotone transformations $\mathbf{f} = (\mathbf{g}', \mathbf{h}')$ such that

$$\mathbf{f}(\mathbf{x}) = \begin{pmatrix} \mathbf{g}(\mathbf{z}) \\ \mathbf{h}(\mathbf{y}) \end{pmatrix} \sim N_p \left(\boldsymbol{\mu}^* = \begin{pmatrix} \boldsymbol{\mu}_z^* \\ \boldsymbol{\mu}_y^* \end{pmatrix}, \boldsymbol{\Sigma}^* = \begin{pmatrix} \boldsymbol{\Sigma}_{zz}^* & \boldsymbol{\Sigma}_{zy}^* \\ \boldsymbol{\Sigma}_{yz}^* & \boldsymbol{\Sigma}_{yy}^* \end{pmatrix} \right), \quad (1)$$

where we may assume that $\boldsymbol{\mu}^* = \mathbf{0}$ and $\boldsymbol{\Sigma}^*$ is a correlation matrix. The transnormal model essentially assumes the joint normality of marginally normal-transformed variables, and it strikes a balance between model robustness and interpretability. This work aims to estimate transformations of predictors $\mathbf{t}(\mathbf{z})$ to optimally predict \mathbf{y} under the transnormal model (1).

Given transformations \mathbf{g} and \mathbf{h} , we derive the conditional distribution of $\mathbf{h}(\mathbf{y})$ given $\mathbf{g}(\mathbf{z})$ as

$$\mathbf{h}(\mathbf{y})|\mathbf{g}(\mathbf{z}) \sim N_{q_0}(\boldsymbol{\Sigma}_{yz}^*(\boldsymbol{\Sigma}_{zz}^*)^{-1} \cdot \mathbf{g}(\mathbf{z}), \boldsymbol{\Sigma}_{yy}^* - \boldsymbol{\Sigma}_{yz}^*(\boldsymbol{\Sigma}_{zz}^*)^{-1} \boldsymbol{\Sigma}_{zy}^*). \quad (2)$$

But, since both \mathbf{g} and \mathbf{h} are unknown, it is challenging to obtain the explicit conditional distribution of \mathbf{y} given \mathbf{z} . Moreover, linearity and homoscedasticity do not hold for the optimal transformations under the transnormal model. In the presence of nonlinear \mathbf{g} and \mathbf{h} , we should take into account the coordinate-wise heterogeneity among the conditional distributions of \mathbf{y} given \mathbf{z} . To this end, we employ quantile regression to handle such heterogeneity.

Quantile regression is introduced by the seminal article of Koenker and Bassett (1978), and it has received much attention in various topics such as survival analysis (Koenker and Geling 2001), time series analysis (Koenker and Xiao 2006), microarray analysis (Wang and He 2007), and variable selection (Zou and Yuan 2008; Bradic, Fan and Wang 2011). Let $\rho_\tau(u) = u \cdot (\tau - I_{[u \leq 0]})$ be the check loss function. Then we consider the multitask quantile regression:

$$\min_{\mathbf{t}: \mathbb{R}^{p_0} \mapsto \mathbb{R}^{q_0}} \sum_{j=1}^{q_0} E[\rho_\tau(y_j - t_j(\mathbf{z}))]. \quad (3)$$

When $\tau = \frac{1}{2}$, $\rho_{\tau=\frac{1}{2}}(u) = |u|$, and (3) reduces to multitask median regression, namely,

$$\min_{\mathbf{t}: \mathbb{R}^{p_0} \mapsto \mathbb{R}^{q_0}} \sum_{j=1}^{q_0} E[\rho_{\frac{1}{2}}(y_j - t_j(\mathbf{z}))] = \frac{1}{2} \min_{\mathbf{t}: \mathbb{R}^{p_0} \mapsto \mathbb{R}^{q_0}} E\|\mathbf{y} - \mathbf{t}(\mathbf{z})\|_{\ell_1}.$$

The use of the ℓ_1 loss in prediction was recommended in Friedman (2001, 2002). Our work shares the similar philosophy and includes multitask median regression as a special case.

In this article, we will show that optimal transformations in multitask quantile regression can be efficiently estimated under the transnormal model. Although estimating optimal transformation is difficult, we derive the surprisingly closed-form optimal prediction without using any smoothing techniques as in Friedman and Stuetzle (1981), Breiman and Friedman (1985), Stone (1985), or Tibshirani (1988). The key ingredient of our proposed method is the positive-definite regularized estimation of large covariance matrices under the transnormal model. We propose two novel rank-based covariance regularization methods: the *rank-based positive definite ℓ_1 penalization* for estimating sparse covariance matrices, and the *rank-based banded Cholesky decomposition regularization*

for estimating banded inverse covariance matrices. Our proposed rank-based covariance regularization critically depends on a positive definite correlation matrix that retains the desired sampling properties of the adjusted Spearman's or Kendall's rank correlation matrix (Kendall 1948; Liu et al. 2012; Xue and Zou 2012), which can be solved by an efficient alternating direction method of multipliers. By combining both strengths of quantile regression and rank-based covariance regularization, we simultaneously address issues of nonlinearity, nonnormality and high dimensionality in estimating optimal transformations for high-dimensional regression. Especially, our proposed method achieves the "oracle"-like convergence rate and provides a provable prediction interval in high dimensions, where dimension is on the nearly exponential order of sample size.

The rest of this article is organized as follows. We first present the methodological details of optimal prediction in multitask quantile regression in Section 2. Section 3 establishes the theoretical properties of our proposed method under the transnormal model. Section 4 contains simulation results and a real application to analyze the protein mass spectroscopy data. Technical proofs are presented in the online Appendices.

2. Multitask Quantile Regression

This section presents the optimal prediction in multitask quantile regression under the transnormal model. The transnormal family of distributions retains the nice interpretation of the normal model, and enables us to make a good use of normal model and theory. In particular, it allows us to obtain the closed-form optimal prediction, which is very appealing to deal with nonnormality and high dimensionality. Moreover, the monotone transformation is easy to handle for quantile estimation.

2.1. The Closed-Form Solution

Let $Q_\tau(y_j|\mathbf{z})$ be the τ th conditional quantile of y_j given \mathbf{z} . Then, it immediately implies that $Q_\tau(y_j|\mathbf{z}) = \min_{t_j} E[\rho_\tau(y_j - t_j(\mathbf{z}))]$. Define $\mathbf{Q}_\tau(\mathbf{y}|\mathbf{z}) = (Q_\tau(y_1|\mathbf{z}), \dots, Q_\tau(y_{p_0}|\mathbf{z}))'$. The τ th equicorrelation conditional quantile $\mathbf{Q}_\tau(\mathbf{y}|\mathbf{z})$ is the analytical solution to (3), that is,

$$\mathbf{Q}_\tau(\mathbf{y}|\mathbf{z}) = \arg \min_{\mathbf{t}: \mathbb{R}^{p_0} \mapsto \mathbb{R}^{q_0}} \sum_{j=1}^{q_0} E[\rho_\tau(y_j - t_j(\mathbf{z}))].$$

By using the fact that \mathbf{h} and \mathbf{g} are monotone under the transnormal model, we have

$$\mathbf{Q}_\tau(\mathbf{y}|\mathbf{z}) = \mathbf{Q}_\tau(\mathbf{y}|\mathbf{g}(\mathbf{z})) = \mathbf{h}^{-1}(\mathbf{Q}_\tau(\mathbf{h}(\mathbf{y})|\mathbf{g}(\mathbf{z}))).$$

Since \mathbf{h} is monotonically nondecreasing, it now follows from (2) that

$$\mathbf{Q}_\tau(\mathbf{y}|\mathbf{z}) = \mathbf{h}^{-1}(\boldsymbol{\Sigma}_{yz}^*(\boldsymbol{\Sigma}_{zz}^*)^{-1} \cdot \mathbf{g}(\mathbf{z}) + \Phi^{-1}(\tau) \text{vdiag}(\boldsymbol{\Sigma}_{yy}^* - \boldsymbol{\Sigma}_{yz}^*(\boldsymbol{\Sigma}_{zz}^*)^{-1} \boldsymbol{\Sigma}_{zy}^*)), \quad (4)$$

where $\Phi(\cdot)$ is the CDF of $N(0, 1)$, and $\text{vdiag}(\mathbf{A})$ denotes the vector formed by the diagonal element of \mathbf{A} . Therefore, the transnormal model enables us to obtain the closed-form solution to multitask quantile regression.

Moreover, the closed-form solution (4) can be used to construct prediction intervals for predicting \mathbf{y} given \mathbf{z} in high

dimensions. To be more specific, we can obtain the closed-form $100(1 - \tau)\%$ prediction interval as

$$[Q_{\frac{\tau}{2}}(y|z), Q_{1-\frac{\tau}{2}}(y|z)].$$

For different values of τ , we only need to adjust the value $\Phi(\tau)$. This is also an appealing feature of the transnormal model, especially under the high-dimensional setting.

The solution $Q_{\tau}(y|z)$ uses true covariance matrices and transformations, and it is not a feasible estimator. To use (4), we need to estimate the covariance matrix Σ^* and transformation f . It turns out that these two tasks are relatively easy under the transnormal model. Section 2.2 describes how to estimate transformation f , and then Section 3 shows how to estimate the large structured covariance matrix Σ^* . With estimated transformation $\hat{f} = (\hat{g}, \hat{h})$ and covariance $\hat{\Sigma}$, it is straightforward to obtain the following plug-in estimator

$$\hat{Q}_{\tau}(y|z) = \hat{h}^{-1}(\hat{\Sigma}_{yz}(\hat{\Sigma}_{zz})^{-1} \cdot \hat{g}(z) + \Phi^{-1}(\tau) \text{vdiag}(\hat{\Sigma}_{yy} - \hat{\Sigma}_{yz}(\hat{\Sigma}_{zz})^{-1} \hat{\Sigma}_{zy})). \quad (5)$$

Given the plug-in estimator (5), we further estimate the $100(1 - \tau)\%$ prediction interval as

$$[\hat{Q}_{\frac{\tau}{2}}(y|z), \hat{Q}_{1-\frac{\tau}{2}}(y|z)].$$

Remark 1. When $\tau = \frac{1}{2}$, it reduces to optimal prediction in multitask median regression. Since $\Phi^{-1}(\frac{1}{2}) = 0$, the optimal prediction can be further simplified as

$$Q_{\frac{1}{2}}(y|z) = h^{-1}(\Sigma_{yz}^*(\Sigma_{zz}^*)^{-1} \cdot g(z)),$$

and

$$\hat{Q}_{\frac{1}{2}}(y|z) = \hat{h}^{-1}(\hat{\Sigma}_{yz}^*(\hat{\Sigma}_{zz}^*)^{-1} \cdot \hat{g}(z)).$$

Remark 2. Compared to multitask mean regression, multitask quantile regression (including median regression) is more robust against outliers in measurements, in addition to delivering the closed-form solution (4). In contrast, by solving ordinary least squares, multitask mean regression is

$$E(y|z) = E(y|g(z)) = E(h^{-1}(h(y))|g(z)).$$

But unlike the quantile regression, this cannot be simplified further unless h is linear.

Remark 3. Often $h(\cdot)$ in (4) is nonlinear. The nonlinearity of $h(\cdot)$ makes the difference of conditional quantiles at different values of τ depend on z and thus models the effect of heteroscedasticity. In contrast, Wu, Yu and Yu (2010), Zhu, Huang and Li (2012), and Fan et al. (2013) used the single-index or semiparametric quantile regression to model the heteroscedastic effect, by imposing the model $y = \mu(z'\beta) + \sigma(z'\beta) \cdot \varepsilon$, where $\sigma(z'\beta) \cdot \varepsilon$ is a heteroscedastic error and ε is independent of z . Unlike the closed-form expression (4), it would be difficult to employ semiparametric quantile regression to handle nonlinearity and high dimensionality.

2.2. Estimation of Transformation

Note that $f_j(x_j) \sim N(0, 1)$ under the transnormal model for any j . Hence, the cumulative distribution function of X_j admits

the form

$$F_j(x_j) = \Phi(f_j(x_j)), \quad \text{or} \quad f_j(x_j) = \Phi^{-1}(F_j(x_j)). \quad (6)$$

Equation (6) provides an intuitive way to estimate the transformation function f_j . Define $\tilde{F}_j(u) = \frac{1}{n} \sum_{i=1}^n I_{\{x_{ij} \leq u\}}$ as the empirical distribution function. We consider the Winsorized empirical distribution function $\hat{F}_j(\cdot)$ as follows:

$$\begin{aligned} \hat{F}_j(u) &= \delta_n \cdot I_{\{\tilde{F}_j(u) < \delta_n\}} + \tilde{F}_j(u) \cdot I_{\{\delta_n \leq \tilde{F}_j(u) \leq 1 - \delta_n\}} \\ &\quad + (1 - \delta_n) \cdot I_{\{\tilde{F}_j(u) > 1 - \delta_n\}}, \end{aligned} \quad (7)$$

where δ_n is the Winsorization parameter to avoid infinity values and to achieve better bias-variance tradeoff. Following Liu, Lafferty, and Wasserman (2009) and Mai and Zou (2015), we specify the Winsorization parameter as $\delta_n = \frac{1}{n}$, which facilitates both theoretical analysis and practical performance. Next, we estimate the transformation f in the transnormal model as follows:

$$\hat{f} = (\hat{f}_1, \dots, \hat{f}_p) = (\Phi^{-1} \circ \hat{F}_1, \dots, \Phi^{-1} \circ \hat{F}_p).$$

Remark that $\hat{f}_1, \dots, \hat{f}_p$ are nondecreasing, and they can be substituted into (5) to estimate optimal transformations.

2.3. Estimation of Correlation Matrix

This section presents two covariance regularization methods based on the *i.i.d.* transnormal data x_1, \dots, x_n : the *rank-based positive definite ℓ_1 penalization* for estimating sparse correlation matrices, and the *rank-based banded Cholesky decomposition regularization* for estimating banded precision matrices. Both estimates achieve the critical positive definiteness, and they can be used in (5) to obtain the optimal prediction.

2.3.1. Positive Definite Sparse Correlation Matrix

Sparse correlation matrices are widely used in many applications, where variables are permutation invariant. By truncating small entries to zero, thresholding (Bickel and Levina 2008b; Fan, Liao, and Mincheva 2013) is a powerful approach to encourage (conditional) sparsity when estimating large correlation matrices. However, the resulting sparse estimator may not be positive definite in practice. Xue, Ma, and Zou (2012) proposed an efficient positive-definite ℓ_1 -penalized estimator to address the indefiniteness issue of thresholding. Now, we extend the positive-definite ℓ_1 -penalized estimator to the transnormal model.

First, we introduce the “oracle” positive-definite ℓ_1 -penalized estimator using the “oracle” transformations f , that is,

$$\begin{aligned} \hat{\Sigma}_{\ell_1}^o &= \arg \min_{\Sigma} \frac{1}{2} \|\Sigma - \hat{R}^o\|_F^2 + \lambda \|\Sigma\|_{1,\text{off}} \\ \text{subject to} \quad &\text{diag}(\Sigma) = 1; \Sigma \succeq \epsilon I \end{aligned}$$

where \hat{R}^o is the sample correlation matrix of “oracle” data $f(x_1), \dots, f(x_n)$, $\|\Sigma\|_{1,\text{off}}$ is the ℓ_1 norm of off-diagonal elements in Σ , and ϵ is some arbitrarily small constant (say, $\epsilon = 10^{-4}$) satisfying that $\lambda_{\min}(\Sigma^*) \geq \epsilon$.

Motivated by $\hat{\Sigma}_{\ell_1}^o$, we only need to derive a comparable alternative to \hat{R}^o . For space consideration, we focus on

the adjusted Spearman's rank correlation throughout this article, since the same analysis can be easily adapted to the adjusted Kendall's rank correlation. Let $\mathbf{s}_j = (s_{1j}, s_{2j}, \dots, s_{nj})'$ be the ranks for $(x_{1j}, x_{2j}, \dots, x_{nj})'$. Denote by $\hat{r}_{jl} = \text{corr}(\mathbf{s}_j, \mathbf{s}_l)$ the Spearman's rank correlation, and by $\tilde{r}_{jl}^s = 2 \sin(\frac{\pi}{6} \hat{r}_{jl})$ the adjusted Spearman's rank correlation. It is well-known that \tilde{r}_{jl}^s corrects the bias of \hat{r}_{jl} (Kendall 1948). Now we consider the rank correlation matrix $\tilde{\mathbf{R}}^s = (\tilde{r}_{jl}^s)_{p \times p}$, which does not require estimating any transformation function. Thus, the rank-based positive-definite ℓ_1 penalization is as follows:

$$\hat{\Sigma}_{\ell_1}^s = \arg \min_{\Sigma} \frac{1}{2} \|\Sigma - \tilde{\mathbf{R}}^s\|_F^2 + \lambda \|\Sigma\|_{1, \text{off}} \quad \text{subject to} \quad \text{diag}(\Sigma) = \mathbf{1}; \Sigma \succeq \epsilon \mathbf{I}. \quad (8)$$

Remark that $\hat{\Sigma}_{\ell_1}^s$ will guarantee the critical positive definiteness, and it can be efficiently solved by the alternating direction method of multiplier in Xue, Ma, and Zou (2012).

2.3.2. Banded Cholesky Decomposition Regularization

When an ordering exists among variables, the bandable structure is used in estimating large covariance matrices. With guaranteed positive definiteness, banded Cholesky decomposition regularization receives much attention; for example, Wu and Pourahmadi (2003); Huang et al. (2006); Bickel and Levina (2008a); and Levina, Rothman, and Zhu (2008).

To motivate our proposal, we first present the “oracle” estimator under the transnormal model. Based on “oracle” data $\mathbf{f}(\mathbf{x}_1), \dots, \mathbf{f}(\mathbf{x}_n)$, the “oracle” banded Cholesky decomposition regularization estimates Σ^* through banding the Cholesky factor of its inverse Θ^* . Suppose Θ^* has the decomposition $\Theta^* = (\mathbf{I} - \mathbf{A})' \mathbf{D}^{-1} (\mathbf{I} - \mathbf{A})$, where \mathbf{D} is a diagonal matrix and $\mathbf{A} = (a_{jl})_{p \times p}$ is a lower triangular matrix with zero diagonal. Let $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_p)'$. Since $\mathbf{f}(\mathbf{x}) \sim N_p(\mathbf{0}, \Sigma^*)$, we have $(\mathbf{I} - \mathbf{A}) \cdot \mathbf{f}(\mathbf{x}) \sim N_p(\mathbf{0}, \mathbf{D})$. As in Bickel and Levina (2008a), the “oracle” estimator is derived by regressing $f_j(x_j)$ on its closest $\min\{k, j-1\}$ predecessors, that is, $\hat{\mathbf{a}}_1^o = \mathbf{0}$, and for $j = 2, \dots, p$, $\hat{\mathbf{a}}_j^o = \arg \min_{\mathbf{a}_j \in \mathcal{A}_j(k)} \frac{1}{n} \sum_{i=1}^n (f_j(x_{ij}) - \mathbf{a}_j' \mathbf{f}(\mathbf{x}_i))^2$, where $\mathcal{A}_j(k) = \{(\alpha_1, \dots, \alpha_p) : \alpha_l = 0 \text{ if } l < j-k \text{ or } l \geq j\}$. Then, \mathbf{A} is estimated by the k -banded lower triangular matrix $\hat{\mathbf{A}}^o = (\hat{\mathbf{a}}_1^o, \dots, \hat{\mathbf{a}}_p^o)'$, and \mathbf{D} is estimated by $\hat{\mathbf{D}}^o = \text{diag}(\hat{d}_1^o, \dots, \hat{d}_p^o)$ with the residual sum of squares $\hat{d}_j^o = \frac{1}{n} \sum_{i=1}^n (f_j(x_{ij}) - (\hat{\mathbf{a}}_j^o)' \mathbf{f}(\mathbf{x}_i))^2$. Hence, the “oracle” estimator ends up with a positive-definite estimator $\hat{\Sigma}_{\text{chol}}^o = (\mathbf{I} - \hat{\mathbf{A}}^o)^{-1} (\hat{\mathbf{D}}^o)^{-1} [(\mathbf{I} - \hat{\mathbf{A}}^o)']^{-1}$, which has the k -banded precision matrix $\hat{\Theta}_{\text{chol}}^o = (\mathbf{I} - \hat{\mathbf{A}}^o)' (\hat{\mathbf{D}}^o)^{-1} (\mathbf{I} - \hat{\mathbf{A}}^o)$.

It is important to observe that $\hat{\mathbf{R}}^o = \frac{1}{n} \sum_{i=1}^n \mathbf{f}(\mathbf{x}_i) (\mathbf{f}(\mathbf{x}_i))'$ plays the central role there. To see this point, we notice that estimating $\hat{\mathbf{a}}_j^o$ and \hat{d}_j^o only depends on the quadratic term

$$\frac{1}{n} \sum_{i=1}^n (f_j(x_{ij}) - \mathbf{a}_j' \mathbf{f}(\mathbf{x}_i))^2 = \mathbf{a}_j' \hat{\mathbf{R}}^o \mathbf{a}_j - 2 \mathbf{a}_j' \hat{\mathbf{r}}_j^o + \hat{r}_{jj}^o, \quad (9)$$

where $\hat{\mathbf{R}}^o = (\hat{r}_{jl}^o)_{p \times p}$ and $\hat{\mathbf{r}}_j^o = (\hat{r}_{j1}^o, \dots, \hat{r}_{jp}^o)'$ is its j th row. Thus, we only need a positive-definite comparable alternative to $\hat{\mathbf{R}}^o$.

The adjusted rank correlation matrix $\tilde{\mathbf{R}}^s$ achieves the “oracle”-like rate of convergence, but it is not guaranteed to be positive definite (Devlin, Gnanadesikan and Kettenring 1975). We employ the nearest correlation matrix projection to inherit the desired sampling properties of $\tilde{\mathbf{R}}^s$ by solving the following optimization problem:

$$\hat{\mathbf{R}}^s = \arg \min_{\mathbf{R}} \|\mathbf{R} - \tilde{\mathbf{R}}^s\|_{\max} \quad \text{subject to} \quad \mathbf{R} \succeq \epsilon \mathbf{I}; \text{diag}(\mathbf{R}) = \mathbf{1}, \quad (10)$$

where $\|\cdot\|_{\max}$ is the entrywise ℓ_∞ norm and $\epsilon > 0$ is some arbitrarily small constant satisfying that $\lambda_{\min}(\Sigma^*) \geq \epsilon$. Zhao, Roeder, and Liu (2012) considered a related projection, but used a smooth surrogate function $\|\mathbf{R} - \tilde{\mathbf{R}}^s\|_{\max}^v = \max_{\|\mathbf{U}\|_1 \leq 1} \langle \mathbf{U}, \mathbf{R} - \tilde{\mathbf{R}}^s \rangle - \frac{v}{2} \|\mathbf{U}\|_F^2$, which would introduce approximation error. Qi and Sun (2006) solved a nearest correlation matrix projection under the Frobenius norm. We use the entrywise ℓ_∞ -norm for theoretical considerations, as we now demonstrate. Since Σ^* is feasible to (10), we have $\|\hat{\mathbf{R}}^s - \tilde{\mathbf{R}}^s\|_{\max} \leq \|\Sigma^* - \tilde{\mathbf{R}}^s\|_{\max}$. By using the triangular inequality, $\hat{\mathbf{R}}^s$ retains the desired sampling properties as $\tilde{\mathbf{R}}^s$, that is,

$$\begin{aligned} \|\hat{\mathbf{R}}^s - \Sigma^*\|_{\max} &\leq \|\hat{\mathbf{R}}^s - \tilde{\mathbf{R}}^s\|_{\max} + \|\tilde{\mathbf{R}}^s - \Sigma^*\|_{\max} \\ &\leq 2 \|\tilde{\mathbf{R}}^s - \Sigma^*\|_{\max}. \end{aligned}$$

The details of nearest correlation matrix projection are presented in the online Appendix.

Now, we propose the rank-based estimator on the basis of $\hat{\mathbf{R}}^s$. Let $\hat{\mathbf{R}}^s = (\hat{r}_{jl}^s)_{p \times p}$ and $\hat{\mathbf{r}}_j^s = (\hat{r}_{j1}^s, \dots, \hat{r}_{jp}^s)'$. We substitute $\hat{\mathbf{R}}^s$ into the quadratic term (9) as

$$\mathbf{a}_j' \hat{\mathbf{R}}^s \mathbf{a}_j - 2 \mathbf{a}_j' \hat{\mathbf{r}}_j^s + \hat{r}_{jj}^s = \mathbf{a}_j' \hat{\mathbf{R}}^s \mathbf{a}_j - 2 \mathbf{a}_j' \hat{\mathbf{r}}_j^s + 1.$$

Accordingly, we estimate \mathbf{A} by

$$\hat{\mathbf{A}}^s = (\hat{\mathbf{a}}_1^s, \dots, \hat{\mathbf{a}}_p^s)',$$

with $\hat{\mathbf{a}}_1^s = \mathbf{0}$, and for $j = 2, \dots, p$,

$$\hat{\mathbf{a}}_j^s = \arg \min_{\mathbf{a}_j \in \mathcal{A}_j(k)} \mathbf{a}_j' \hat{\mathbf{R}}^s \mathbf{a}_j - 2 \mathbf{a}_j' \hat{\mathbf{r}}_j^s + 1,$$

where $\mathcal{A}_j(k) = \{(\alpha_1, \dots, \alpha_p) : \alpha_i = 0 \text{ if } i < j-k \text{ or } i \geq j\}$. In addition, we estimate \mathbf{D} by

$$\hat{\mathbf{D}}^s = \text{diag}(\hat{d}_1^s, \dots, \hat{d}_p^s),$$

with

$$\hat{d}_j^s = (\hat{\mathbf{a}}_j^s)' \hat{\mathbf{R}}^s \hat{\mathbf{a}}_j^s - 2 (\hat{\mathbf{a}}_j^s)' \hat{\mathbf{r}}_j^s + 1.$$

Thus, the rank-based banded Cholesky decomposition regularization yields

$$\hat{\Sigma}_{\text{chol}}^s = (\mathbf{I} - \hat{\mathbf{A}}^s)^{-1} \hat{\mathbf{D}}^s [(\mathbf{I} - \hat{\mathbf{A}}^s)']^{-1},$$

which has the k -banded precision matrix $\hat{\Theta}_{\text{chol}}^s = (\mathbf{I} - \hat{\mathbf{A}}^s)' (\hat{\mathbf{D}}^s)^{-1} (\mathbf{I} - \hat{\mathbf{A}}^s)$.

3. Theoretical Properties

This section presents theoretical properties of our proposed methods. We use several matrix norms: the ℓ_1 norm $\|\mathbf{U}\|_{\ell_1} = \max_j \sum_i |u_{ij}|$, the spectral norm $\|\mathbf{U}\|_{\ell_2} = \lambda_{\max}^{1/2}(\mathbf{U}'\mathbf{U})$, and the ℓ_∞ norm $\|\mathbf{U}\|_{\ell_\infty} = \max_i \sum_j |u_{ij}|$. For a symmetric matrix, its matrix ℓ_1 norm coincides with its matrix ℓ_∞ norm. We use c or C to denote constants that do not depend on n or p .

Throughout this section, we follow Bickel and Levina (2008a, b) to assume that

$$\varepsilon_0 \leq \lambda_{\min}(\mathbf{\Sigma}^*) \leq \lambda_{\max}(\mathbf{\Sigma}^*) \leq \frac{1}{\varepsilon_0}.$$

3.1. A General Theory

First of all, we consider any regularized estimator $\hat{\mathbf{\Sigma}}$ satisfying the following condition:

(C1) There exists a regularized estimator $\hat{\mathbf{\Sigma}} = \begin{pmatrix} \hat{\mathbf{\Sigma}}_{zz} & \hat{\mathbf{\Sigma}}_{zy} \\ \hat{\mathbf{\Sigma}}_{yz} & \hat{\mathbf{\Sigma}}_{yy} \end{pmatrix}$ satisfying that

$$\|\hat{\mathbf{\Sigma}} - \mathbf{\Sigma}^*\|_{\ell_1} = O_p(\xi_{n,p}), \quad \text{with } \xi_{n,p} = o((\log n)^{-1/2}).$$

In the following theorem, we show that the conditional distribution of $\mathbf{h}(\mathbf{y})$ given $\mathbf{g}(\mathbf{z})$ can be well estimated under Condition (C1).

Theorem 1. Assume the data follow the transnormal model. Suppose that there exists $\kappa \in (0, 1)$ such that $n^\kappa \gg \log n \log p_0$. Given Condition (C1), we have the error bounds concerning the estimated conditional distribution of $\mathbf{h}(\mathbf{y})$ given $\mathbf{g}(\mathbf{z})$ in (2) as follows:

(i) On the conditional mean vector:

$$\begin{aligned} & \|\hat{\mathbf{\Sigma}}_{yz} \hat{\mathbf{\Sigma}}_{zz}^{-1} \cdot \hat{\mathbf{g}}(\mathbf{z}) - \mathbf{\Sigma}_{yz}^* (\mathbf{\Sigma}_{zz}^*)^{-1} \cdot \mathbf{g}(\mathbf{z})\|_{\max} \\ &= O_p \left(\sqrt{\frac{\log n \log p_0}{n^\kappa}} + \xi_{n,p} \sqrt{\log n} \right). \end{aligned}$$

(ii) On the conditional variance matrix:

$$\begin{aligned} & \|(\hat{\mathbf{\Sigma}}_{yy} - \hat{\mathbf{\Sigma}}_{yz} \hat{\mathbf{\Sigma}}_{zz}^{-1} \hat{\mathbf{\Sigma}}_{zy}) - (\mathbf{\Sigma}_{yy}^* - \mathbf{\Sigma}_{yz}^* (\mathbf{\Sigma}_{zz}^*)^{-1} \mathbf{\Sigma}_{zy}^*)\|_{\ell_2} \\ &= O_p(\xi_{n,p}). \end{aligned}$$

Next, we show that the plug-in estimator $\hat{\mathbf{Q}}_\tau(\mathbf{y}|\mathbf{z})$ is asymptotically as good as $\mathbf{Q}_\tau(\mathbf{y}|\mathbf{z})$ under mild regularity conditions. Let $\psi_j(\cdot)$ be the probability density function of X_j . Define $\mathbf{L} = \mathbf{\Sigma}_{yz}^* (\mathbf{\Sigma}_{zz}^*)^{-1} \cdot \mathbf{g}(\mathbf{z}) + \Phi^{-1}(\tau) \text{vdiag}(\mathbf{\Sigma}_{yy}^* - \mathbf{\Sigma}_{yz}^* (\mathbf{\Sigma}_{zz}^*)^{-1} \mathbf{\Sigma}_{zy}^*)$, and let $\mathbf{L} = (L_1, \dots, L_{q_0})'$.

Theorem 2. Under the conditions of Theorem 1, we have the error bound concerning the plug-in estimator in (5) as follows:

$$\begin{aligned} & \|\hat{\mathbf{Q}}_\tau(\mathbf{y}|\mathbf{z}) - \mathbf{Q}_\tau(\mathbf{y}|\mathbf{z})\|_{\max} \\ &= O_p \left(\frac{1}{M} \sqrt{\frac{\log q_0}{n}} + \frac{1}{M} \sqrt{\frac{\log n \log p_0}{n^\kappa}} + \frac{\xi_{n,p}}{M} \sqrt{\log n} \right), \end{aligned} \quad (11)$$

where M is the minimum of $\psi_{p_0+j}(x)$ over $x \in \mathcal{I}_j = [-|2L_j|, |2L_j|] \cup [-f_{p_0+j}^{-1}(|2L_j|), f_{p_0+j}^{-1}(|2L_j|)]$ for any $j = 1, \dots, q_0$, that is, $M = \min_{j=1, \dots, q_0} \min_{x \in \mathcal{I}_j} \psi_{p_0+j}(x)$.

Theorem 2 immediately implies that the plug-in estimator $\hat{\mathbf{Q}}_\tau(\mathbf{y}|\mathbf{z})$ can be used to construct provable prediction intervals in high dimensions. Therefore, we use

$$[\hat{\mathbf{Q}}_{\frac{\tau}{2}}(\mathbf{y}|\mathbf{z}), \hat{\mathbf{Q}}_{1-\frac{\tau}{2}}(\mathbf{y}|\mathbf{z})]$$

to construct the $100(1 - \tau)\%$ prediction interval for predicting \mathbf{y} given \mathbf{z} .

In what follows, we show how Condition (C1) is satisfied. Specifically, we consider two parameter spaces for $\mathbf{\Sigma}^*$:

$$\begin{aligned} \mathcal{G}_q &= \left\{ \mathbf{\Sigma} : \max_j \sum_{j \neq i} |a_{ij}|^q \leq s_0 \right\} \\ \mathcal{H}_\alpha &= \left\{ \mathbf{\Sigma} : \max_j \sum_{j < i-k} |a_{ij}| \leq c_0 k^{-\alpha}, \forall k \right\}. \end{aligned}$$

These two parameter spaces were previously studied by Bickel and Levina (2008a,b) and Cai and Zhou (2012). We will show that the proposed estimators achieve the “oracle”-like convergence rate over \mathcal{G}_q and \mathcal{H}_α , respectively. As a result, we obtain the optimal prediction result.

3.2. Positive-Definite Sparse Correlation Matrix

First, we derive the estimation bound for the rank-based positive-definite ℓ_1 penalization and show its application to estimate the conditional distribution of $\mathbf{h}(\mathbf{y})$ given $\mathbf{g}(\mathbf{z})$.

Theorem 3. Assume the data follow the transnormal model with $\mathbf{\Sigma}^* \in \mathcal{G}_q$, and also assume that $n \gg \log p$. Suppose that $\lambda_{\min}(\mathbf{\Sigma}^*) \geq \varepsilon_0 \gg s_0 (\log p/n)^{\frac{1-q}{2}} + \epsilon$. With probability tending to 1, the rank-based positive-definite ℓ_1 penalization with $\lambda = c(\log p/n)^{1/2}$ achieves the following upper bound under the matrix ℓ_1 norm,

$$\sup_{\mathbf{\Sigma}^* \in \mathcal{G}_q} \|\hat{\mathbf{\Sigma}}_{\ell_1}^s - \mathbf{\Sigma}^*\|_{\ell_1} \leq C \cdot s_0 \left(\frac{\log p}{n} \right)^{\frac{1-q}{2}}.$$

Remark 4. Both the rank-based estimator $\hat{\mathbf{\Sigma}}_{\ell_1}^s$ and its “oracle” counterpart achieve the minimax optimal convergence rate under the matrix ℓ_1 norm (Cai and Zhou 2012).

Corollary 1. Under the conditions of Theorems 1, 2, and 3, with $\hat{\mathbf{\Sigma}} = \hat{\mathbf{\Sigma}}_{\ell_1}^s$ we have

$$\begin{aligned} & \sup_{\mathbf{\Sigma}^* \in \mathcal{G}_q} \left\| \hat{\mathbf{\Sigma}}_{yz} \hat{\mathbf{\Sigma}}_{zz}^{-1} \cdot \hat{\mathbf{g}}(\mathbf{z}) - \mathbf{\Sigma}_{yz}^* (\mathbf{\Sigma}_{zz}^*)^{-1} \cdot \mathbf{g}(\mathbf{z}) \right\|_{\max} \\ &= O_p \left(\sqrt{\frac{\log n \log p_0}{n^\kappa}} + s_0 \left(\frac{\log p}{n} \right)^{\frac{1-q}{2}} \sqrt{\log n} \right), \end{aligned}$$

and

$$\sup_{\mathbf{\Sigma}^* \in \mathcal{G}_q} \left\| (\hat{\mathbf{\Sigma}}_{yy} - \hat{\mathbf{\Sigma}}_{yz} \hat{\mathbf{\Sigma}}_{zz}^{-1} \hat{\mathbf{\Sigma}}_{zy}) - (\mathbf{\Sigma}_{yy}^* - \mathbf{\Sigma}_{yz}^* (\mathbf{\Sigma}_{zz}^*)^{-1} \mathbf{\Sigma}_{zy}^*) \right\|_{\ell_2}$$

$$= O_P \left(s_0 \left(\frac{\log p}{n} \right)^{\frac{1-q}{2}} \right).$$

In light of [Theorem 3](#) and [Corollary 1](#), we can derive the convergence rate for the proposed optimal prediction in the following theorem.

Theorem 4. Under same conditions of [Theorems 1, 2, and 3](#), with $\hat{\Sigma} = \hat{\Sigma}_{\ell_1}^s$ we have

$$\sup_{\Sigma^* \in \mathcal{H}_q} \|\hat{Q}_\tau(y|z) - Q_\tau(y|z)\|_{\max} = O_P \left(\frac{1}{M} \sqrt{\frac{\log q_0}{n}} + \frac{1}{M} \sqrt{\frac{\log n \log p_0}{n^\kappa}} + \frac{s_0}{M} \left(\frac{\log p}{n} \right)^{\frac{1-q}{2}} \sqrt{\log n} \right).$$

3.3. Banded Cholesky Decomposition Regularization

Next, we derive the estimation bound for the rank-based banded Cholesky decomposition regularization with applications to estimate the conditional distribution of $\mathbf{h}(\mathbf{y})$ given $\mathbf{g}(\mathbf{z})$.

Theorem 5. Assume the data follow the transnormal model with $\Sigma^* \in \mathcal{H}_\alpha$, and also assume that $n \gg k^2 \log p$. With probability tending to 1, the rank-based banded Cholesky decomposition regularization achieves the following upper bound under the matrix ℓ_1 norm,

$$\sup_{\Sigma^* \in \mathcal{H}_\alpha} \|\hat{\Sigma}_{\text{chol}}^s - \Sigma^*\|_{\ell_1} \leq Ck \left(\frac{\log p}{n} \right)^{1/2} + Ck^{-\alpha}.$$

If $k = c \cdot \left(\frac{\log p}{n} \right)^{\frac{-1}{2(\alpha+1)}}$ in the rank-based banded Cholesky decomposition regularization, we have

$$\sup_{\Sigma^* \in \mathcal{H}_\alpha} \|\hat{\Sigma}_{\text{chol}}^s - \Sigma^*\|_{\ell_1} = O_P \left(\left(\frac{\log p}{n} \right)^{\frac{\alpha}{2(\alpha+1)}} \right).$$

Remark 5. Bickel and Levina (2008a) studied the banded Cholesky decomposition regularization under the normal model. Their analysis directly applies to the “oracle” banded Cholesky decomposition regularization. By Theorem 3 of Bickel and Levina (2008a), we have

$$\sup_{\Sigma^* \in \mathcal{H}_\alpha} \|\hat{\Sigma}_{\text{chol}}^o - \Sigma^*\|_{\ell_1} = O_P \left(\left(\frac{\log p}{n} \right)^{\frac{\alpha}{2(\alpha+1)}} \right).$$

Therefore, the rank-based banded Cholesky decomposition regularization achieves the same convergence rate as the “oracle” counterpart.

Corollary 2. Under the conditions of [Theorems 1, 2, and 5](#), with $\hat{\Sigma} = \hat{\Sigma}_{\text{chol}}^s$ we have

$$\sup_{\Sigma^* \in \mathcal{H}_\alpha} \left\| \hat{\Sigma}_{yz} \hat{\Sigma}_{zz}^{-1} \cdot \hat{\mathbf{g}}(\mathbf{z}) - \Sigma_{yz}^* (\Sigma_{zz}^*)^{-1} \cdot \mathbf{g}(\mathbf{z}) \right\|_{\max} = O_P \left(\sqrt{\frac{\log n \log p_0}{n^\kappa}} + \left(\frac{\log p}{n} \right)^{\frac{\alpha}{2(\alpha+1)}} \sqrt{\log n} \right),$$

and

$$\sup_{\Sigma^* \in \mathcal{H}_\alpha} \left\| \left(\hat{\Sigma}_{yy} - \hat{\Sigma}_{yz} \hat{\Sigma}_{zz}^{-1} \hat{\Sigma}_{zy} \right) - \left(\Sigma_{yy}^* - \Sigma_{yz}^* (\Sigma_{zz}^*)^{-1} \Sigma_{zy}^* \right) \right\|_{\ell_2} = O_P \left(\left(\frac{\log p}{n} \right)^{\frac{\alpha}{2(\alpha+1)}} \right).$$

In light of [Theorem 5](#) and [Corollary 2](#), we can derive the convergence rate for the proposed optimal prediction in the following theorem.

Theorem 6. Under the conditions of [Theorems 1, 2 and 5](#), with $\hat{\Sigma} = \hat{\Sigma}_{\text{chol}}^s$ we have

$$\sup_{\Sigma^* \in \mathcal{H}_q} \left\| \hat{Q}_\tau(y|z) - Q_\tau(y|z) \right\|_{\max} = O_P \left(\frac{1}{M} \sqrt{\frac{\log q_0}{n}} + \frac{1}{M} \sqrt{\frac{\log n \log p_0}{n^\kappa}} + \frac{\sqrt{\log n}}{M} \left(\frac{\log p}{n} \right)^{\frac{\alpha}{2(\alpha+1)}} \right).$$

4. Numerical Properties

This section demonstrates the numerical performance of the proposed methods. For space consideration, we only focus on multitask median regression and rank-based banded Cholesky decomposition regularization.

4.1. Simulation Studies

Simulation studies investigate the “oracle” estimator, the proposed rank-based estimator and the “naïve” estimator. The “oracle” estimator serves as a benchmark in the numerical comparison, and the “naïve” estimator directly regularizes on the covariance of the original data. We summarize notation and details of three regularized estimator in [Table 1](#).

In Models 1–3, we consider three different designs for the inverse covariance matrix Ω^* :

Model 1: $\Omega^* = (\mathbf{I} - \mathbf{A})' \mathbf{D}^{-1} (\mathbf{I} - \mathbf{A})$: $d_{ii} = 0.01$, $a_{i+1,i} = 0.8$ and $d_{ij} = a_{ij} = 0$ otherwise;

Model 2: Ω^* : $\omega_{ii}^* = 1$, $\omega_{i,i+1}^* = 0.5$, $\omega_{i,i+2}^* = 0.25$ and $\omega_{ij}^* = 0$ otherwise;

Model 3: Ω^* : $\omega_{ii}^* = 1$, $\omega_{i,i+1}^* = 0.4$, $\omega_{i,i+2}^* = \omega_{i,i+3}^* = 0.2$, $\omega_{i,i+4}^* = 0.1$ and $\omega_{ij}^* = 0$ otherwise.

Models 1–3 are autoregressive models considered by Huang et al. (2006) and Xue and Zou (2012). Based on $\Gamma^* = (\Omega^*)^{-1} = (\gamma_{ij}^*)_{p \times p}$, we calculate the true correlation matrix $\Sigma^* = (\sigma_{ij}^*)_{p \times p} = \left(\frac{\gamma_{ij}^*}{\sqrt{\gamma_{ii}^* \gamma_{jj}^*}} \right)_{p \times p}$ and also its inverse $\Theta^* = (\Sigma^*)^{-1}$.

Table 1. List of three regularized estimators in the simulation study.

Methods	Details
$(\hat{\Sigma}_{\text{chol}}^o, \hat{\Theta}_{\text{chol}}^o, \hat{t}_1^o(\mathbf{z}))$	regularizing the “oracle” sample correlation matrix
$(\hat{\Sigma}_{\text{chol}}^s, \hat{\Theta}_{\text{chol}}^s, \hat{t}_1^s(\mathbf{z}))$	regularizing the adjusted Spearman’s rank correlation matrix
$(\hat{\Sigma}_{\text{chol}}^n, \hat{\Theta}_{\text{chol}}^n, \hat{t}_1^n(\mathbf{z}))$	regularizing the usual sample correlation matrix

Table 2. Performance of estimating the correlation matrix with $\hat{\Sigma}_{\text{chol}}^o$, $\hat{\Sigma}_{\text{chol}}^s$ and $\hat{\Sigma}_{\text{chol}}^n$.

Method	Model 1			Model 2			Model 3		
	$p = 100$	$p = 200$	$p = 500$	$p = 100$	$p = 200$	$p = 500$	$p = 100$	$p = 200$	$p = 500$
matrix ℓ_1 norm									
$\hat{\Sigma}_{\text{chol}}^o$	1.34 (0.05)	1.60 (0.13)	1.95 (0.27)	0.89 (0.03)	0.91 (0.04)	0.99 (0.08)	1.00 (0.02)	1.04 (0.05)	1.16 (0.12)
$\hat{\Sigma}_{\text{chol}}^s$	1.49 (0.07)	1.68 (0.11)	2.08 (0.21)	0.92 (0.03)	0.99 (0.05)	1.10 (0.14)	1.03 (0.02)	1.06 (0.04)	1.21 (0.12)
$\hat{\Sigma}_{\text{chol}}^n$	3.21 (0.06)	3.49 (0.09)	3.78 (0.24)	1.45 (0.02)	1.48 (0.04)	1.63 (0.11)	1.11 (0.03)	1.22 (0.05)	1.28 (0.11)
matrix ℓ_2 norm									
$\hat{\Sigma}_{\text{chol}}^o$	0.78 (0.03)	0.95 (0.06)	1.10 (0.13)	0.48 (0.01)	0.49 (0.02)	0.54 (0.04)	0.52 (0.01)	0.54 (0.02)	0.59 (0.05)
$\hat{\Sigma}_{\text{chol}}^s$	0.84 (0.03)	1.00 (0.06)	1.26 (0.11)	0.49 (0.01)	0.53 (0.02)	0.58 (0.05)	0.53 (0.01)	0.55 (0.02)	0.63 (0.06)
$\hat{\Sigma}_{\text{chol}}^n$	1.67 (0.05)	1.73 (0.08)	1.89 (0.14)	0.77 (0.01)	0.80 (0.02)	0.84 (0.04)	0.55 (0.01)	0.60 (0.02)	0.68 (0.05)

NOTE: Each metric is averaged over 100 independent replications with standard errors in the bracket.

Then, we generate n independent normal data from $N_p(\mathbf{0}, \Sigma^*)$ and transfer the normal data to the desired transnormal data $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ using

$$[f_1^{-1}, f_2^{-1}, f_3^{-1}, f_4^{-1}, f_1^{-1}, f_2^{-1}, f_3^{-1}, f_4^{-1}, \dots], \quad (12)$$

where $f_1(x) = x^{\frac{1}{3}}$ (power transformation), $f_2(x) = \log(\frac{x}{1-x})$ (logit transformation), $f_3(x) = \log(x)$ (log transformation) and $f_4(x) = f_1(x)I_{\{x < -1\}} + f_2(x)I_{\{-1 \leq x \leq 1\}} + (f_3(x-1) + 1)I_{\{x > 1\}}$.

In all cases, we let $n = 200$ and $p = 100, 200$, and 500 . For each estimator, tuning parameter is chosen by cross-validation. Estimation accuracy is measured by the matrix ℓ_1 -norm and ℓ_2 -norm averaged over 100 independent replications.

Tables 2 and 3 summarize the numerical performance of estimating the correlation matrix and its inverse for three estimators. We see that the “naïve” estimator performs the worst in the presence of nonnormality. The rank-based estimator effectively deals with the transnormal data, and performs comparably with the “oracle” estimator. The numeric evidence is consistent with theoretical results presented in Section 3.

Next, we compare the performance of estimating optimal transformations for multitask median regression. In each of the 100 replications, we simulate another $4 \times n$ independent transnormal data $(\mathbf{x}_{n+1}, \dots, \mathbf{x}_{5n})$ as the testing data. Now, we take the first $p_0 = p/2$ variables in \mathbf{x} as response \mathbf{y} and take the second $q_0 = p/2$ variables in \mathbf{x} as predictor \mathbf{z} . Now, given three estimators $\hat{\Sigma}_{\text{chol}}^o$, $\hat{\Sigma}_{\text{chol}}^s$ and $\hat{\Sigma}_{\text{chol}}^n$, we examine their performance in predicting \mathbf{y}_i based on \mathbf{z}_i for $i = n+1, \dots, 5n$. Recall that $\hat{\mathbf{t}}_1(\mathbf{z}) = \hat{\mathbf{h}}^{-1}(\hat{\Sigma}_{yz}(\hat{\Sigma}_{zz})^{-1} \cdot \hat{\mathbf{g}}(\mathbf{z}))$ is the closed-form solution for multitask median regression. In what follows, we consider three optimal prediction methods using $\hat{\Sigma}_{\text{chol}}^o$, $\hat{\Sigma}_{\text{chol}}^s$ and $\hat{\Sigma}_{\text{chol}}^n$ respectively:

- the “oracle” estimator: $\hat{\mathbf{t}}_1^o(\mathbf{z}) = \mathbf{h}^{-1}(\hat{\Sigma}_{yz}^o(\hat{\Sigma}_{zz}^o)^{-1} \cdot \mathbf{g}(\mathbf{z}))$
- the rank-based estimator: $\hat{\mathbf{t}}_1^s(\mathbf{z}) = \hat{\mathbf{h}}^{-1}(\hat{\Sigma}_{yz}^s(\hat{\Sigma}_{zz}^s)^{-1} \cdot \hat{\mathbf{g}}(\mathbf{z}))$
- the “naïve” estimator: $\hat{\mathbf{t}}_1^n(\mathbf{z}) = \hat{\mu}_y + \hat{\Sigma}_{yz}^n(\hat{\Sigma}_{zz}^n)^{-1} \cdot (\mathbf{z} - \hat{\mu}_z)$

Table 4 shows their performances. Prediction accuracy is measured by the difference of prediction errors, which is defined

Table 3. Performance of estimating the inverse correlation matrix with $\hat{\Theta}_{\text{chol}}^o$, $\hat{\Theta}_{\text{chol}}^s$ and $\hat{\Theta}_{\text{chol}}^n$.

Method	Model 1			Model 2			Model 3		
	$p = 100$	$p = 200$	$p = 500$	$p = 100$	$p = 200$	$p = 500$	$p = 100$	$p = 200$	$p = 500$
matrix ℓ_1 norm									
$\hat{\Theta}_{\text{chol}}^o$	4.07 (0.21)	4.72 (0.37)	5.14 (0.57)	1.83 (0.08)	1.95 (0.12)	1.99 (0.25)	2.21 (0.04)	2.29 (0.09)	2.37 (0.18)
$\hat{\Theta}_{\text{chol}}^s$	4.21 (0.24)	4.80 (0.32)	5.28 (0.55)	1.91 (0.07)	1.97 (0.14)	2.02 (0.25)	2.27 (0.04)	2.33 (0.07)	2.41 (0.16)
$\hat{\Theta}_{\text{chol}}^n$	7.18 (0.22)	7.55 (0.35)	8.02 (0.54)	3.29 (0.06)	3.48 (0.14)	3.70 (0.29)	2.81 (0.03)	2.93 (0.10)	3.08 (0.19)
matrix ℓ_2 norm									
$\hat{\Theta}_{\text{chol}}^o$	3.18 (0.10)	3.41 (0.27)	4.07 (0.35)	1.26 (0.04)	1.36 (0.07)	1.41 (0.13)	1.68 (0.04)	1.74 (0.08)	1.75 (0.10)
$\hat{\Theta}_{\text{chol}}^s$	3.29 (0.16)	3.59 (0.25)	4.10 (0.30)	1.35 (0.04)	1.39 (0.07)	1.45 (0.08)	1.69 (0.04)	1.75 (0.08)	1.77 (0.11)
$\hat{\Theta}_{\text{chol}}^n$	4.98 (0.18)	5.15 (0.25)	5.33 (0.32)	2.64 (0.04)	2.81 (0.11)	2.89 (0.18)	2.29 (0.03)	2.32 (0.07)	2.39 (0.18)

NOTE: Each metric is averaged over 100 independent replications with standard errors in the bracket.

Table 4. Performance of three optimal prediction methods $\hat{t}_1^o(z)$, $\hat{t}_1^s(z)$, and $\hat{t}_1^n(z)$.

Method	Model 1			Model 2			Model 3		
	$p = 100$	$p = 200$	$p = 500$	$p = 100$	$p = 200$	$p = 500$	$p = 100$	$p = 200$	$p = 500$
$\hat{t}_1^s(z)$ vs. $\hat{t}_1^o(z)$	0.02 (0.00)	0.03 (0.00)	0.08 (0.01)	0.16 (0.02)	0.26 (0.03)	0.66 (0.04)	0.12 (0.01)	0.22 (0.01)	0.56 (0.04)
$\hat{t}_1^n(z)$ vs. $\hat{t}_1^o(z)$	0.11 (0.00)	0.19 (0.01)	0.21 (0.03)	13.50 (0.22)	26.07 (0.58)	63.80 (2.16)	7.38 (0.14)	14.70 (0.37)	35.68 (1.28)

NOTE: Prediction accuracy are measured by the relative prediction error and averaged over 100 independent replications, with the standard errors shown in the bracket.

by subtracting the “oracle” prediction error, that is,

$$\text{DPE}(\hat{t}_1(z)) = \frac{1}{4n} \sum_{i=n+1}^{5n} \|\hat{t}_1(z_i) - y_i\|_{\ell_1} - \frac{1}{4n} \sum_{i=n+1}^{5n} \|\hat{t}_1^o(z) - y_i\|_{\ell_1}.$$

The relative prediction error $\text{DPE}(\hat{t}_1(z))$ is averaged over 100 independent replications.

As shown in Table 4, we see that the rank-based estimator performs similarly to its oracle counterpart, and outperforms the “naïve” estimator. The numeric evidence is also consistent with theoretical results presented in Section 3.

4.2. Application to the Protein Mass Spectroscopy Data

Recent advances in high-throughput mass spectroscopy technology has enabled biomedical researchers to simultaneously analyze thousands of proteins. This subsection illustrates the power of the proposed rank-based method in an application to study the prostate cancer using the protein mass spectroscopy data (Adam et al. 2002), which was previously analyzed by Levina, Rothman, and Zhu (2008). This dataset consists of the protein mass spectroscopy measurements for the blood serum samples of 157 healthy people and 167 prostate cancer patients. In each blood serum sample, the protein mass spectroscopy measures the intensity for the ordered time-of-flight values, which are related to the mass over charge ratio of proteins. In our analysis, we exclude the measurements with mass over charge ratio below 2000 to avoid chemical artifacts, and perform the same preprocessing as in Levina, Rothman, and Zhu (2008) to smooth the intensity profile. This gives a total of $p = 218$ ordered mass over charge ratio indices for each sample. Then, we have the control data $\mathbf{x}_i^{\text{co}} = (x_{i,1}^{\text{co}}, \dots, x_{i,218}^{\text{co}})$ for $i = 1, \dots, 157$, and the cancer data $\mathbf{x}_j^{\text{ca}} = (x_{j,1}^{\text{ca}}, \dots, x_{j,218}^{\text{ca}})$ for $j = 1, \dots, 167$. We refer the readers to Adam et al. (2002) and Levina, Rothman, and Zhu (2008) for more details about this dataset and also the preprocessing procedure. We use the more stable intensity measurements in the latter 168 mass over charge ratio indices to

predict the more volatile intensity measurements in the first 50 indices.

The analysis of Levina, Rothman, and Zhu (2008) is based on the multivariate normal assumption. We perform normality tests on both control and cancer data. Table 5 shows the number of rejecting null hypotheses among 218 normality tests. At least 50% of mass spectroscopy measurements appear to be non-normal. Even after Bonferroni correction, there are at least 30 indices in the control data and 116 indices in the cancer data rejecting all null hypotheses. Figure 1 illustrates nonnormality (e.g., heavy tails, skewness) in two indices (109, 218).

The nonnormality of measurements suggests us to try the transnormal model. Let \mathbf{z}^{co} (\mathbf{z}^{ca}) and \mathbf{y}^{co} (\mathbf{y}^{ca}) denote the intensity measurements in the latter 168 and first 50 indices of the control (cancer) data. We divide this dataset into training sets $(\mathbf{x}_1^{\text{co}}, \dots, \mathbf{x}_{120}^{\text{co}})$, $(\mathbf{x}_1^{\text{ca}}, \dots, \mathbf{x}_{120}^{\text{ca}})$ and testing sets $(\mathbf{x}_{121}^{\text{co}}, \dots, \mathbf{x}_{157}^{\text{co}})$, $(\mathbf{x}_{121}^{\text{ca}}, \dots, \mathbf{x}_{167}^{\text{ca}})$. Note that the sample estimator $\hat{t}_1^{\text{sample}}(z) = \hat{\mu}_y + \hat{\Sigma}_{yz}^{\text{sample}} (\hat{\Sigma}_{zz}^{\text{sample}})^{-1} \cdot (z - \hat{\mu}_z)$ is infeasible since the usual sample correlation matrix $\hat{\Sigma}_{zz}^{\text{sample}}$ is not invertible. We consider two different methods to predict y : the proposed rank-based estimator $\hat{t}_1^s(z) = \hat{\mathbf{h}}^{-1} (\hat{\Sigma}_{yz}^s (\hat{\Sigma}_{zz}^s)^{-1} \cdot \hat{\mathbf{g}}(z))$ by using the proposed rank-based banded Cholesky decomposition regularization and the “naïve” estimator $\hat{t}_1^n(z) = \hat{\mu}_y + \hat{\Sigma}_{yz}^n (\hat{\Sigma}_{zz}^n)^{-1} \cdot (z - \hat{\mu}_z)$ by using the banded Cholesky decomposition regularization (Bickel and Levina 2008a). Following Levina, Rothman, and Zhu (2008), tuning parameters are chosen via cross-validation. The difference of prediction error at j th index is computed as

$$\text{DPE}_j(\hat{t}_1(z)) = \text{PE}_j(\hat{t}_1^s(z)) - \text{PE}_j(\hat{t}_1^n(z)),$$

where $\text{PE}_j(\hat{t}_1(z))$ is prediction error at the j th index computed by averaging the absolute prediction error $|\hat{t}_1(z_i)_j - y_{ij}|$ over $i = 121, \dots, 157$ for the control data and over $i = 121, \dots, 167$ for the cancer data. The differences of prediction errors are shown in Figure 2. As demonstrated in Figure 2, the rank-based method outperforms the “naïve” estimator in 38 out of 50 indices for the control data and 46 out of 50 indices for the cancer data.

Table 5. Testing for normality.

	Significance level	Anderson–Darling	Cramer–von Mises	Kolmogorov–Smirnov
Control data	0.05	158	138	119
	0.05/218	71	52	30
Cancer data	0.05	196	177	156
	0.05/218	136	134	116

NOTE: This table shows the number of rejecting the normality hypothesis at different significance levels among 218 mass over charge ratio indices.

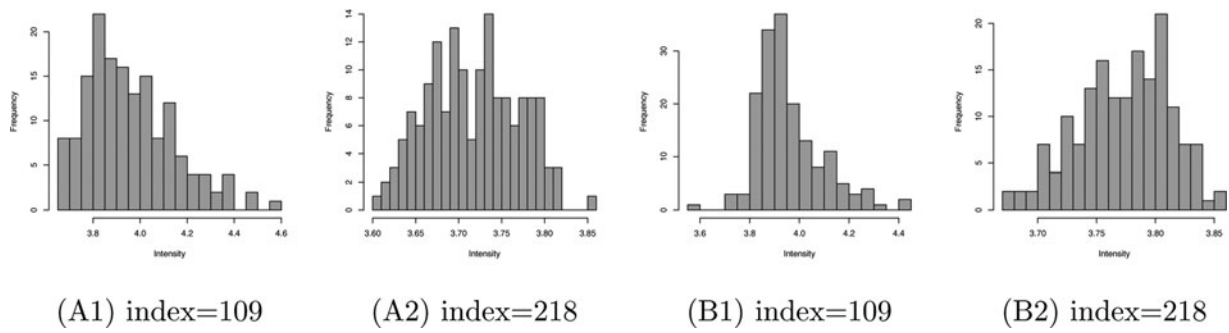
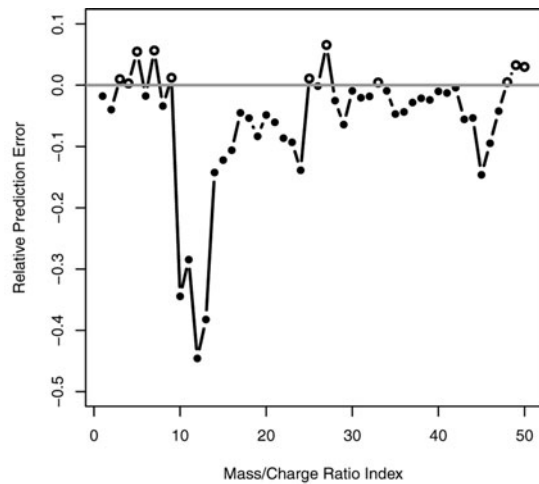


Figure 1. Illustration of the nonnormality issue of the protein mass spectroscopy data: (A1) and (A2) for the control data; (B1) and (B2) for the cancer data.

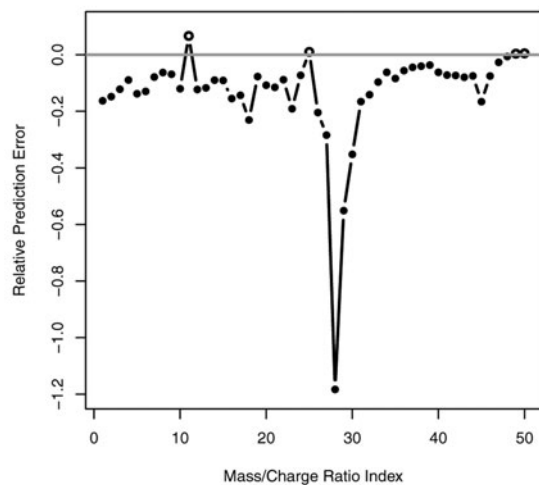
Moreover, we further use $[\hat{Q}_{\frac{\tau}{2}}(y|z), \hat{Q}_{1-\frac{\tau}{2}}(y|z)]$ to construct the $100(1 - \tau)\%$ prediction interval. The predict target and prediction intervals are averaged over the testing set. Specifically, we estimate the $100(1 - \tau)\%$ prediction interval as

follows,

$$100(1 - \tau)\% \text{ PI} = \left[\frac{1}{34} \sum_{i=206}^{239} \hat{Q}_{\frac{\tau}{2}}(y_i|z_i), \frac{1}{34} \sum_{i=206}^{239} \hat{Q}_{1-\frac{\tau}{2}}(y_i|z_i) \right].$$

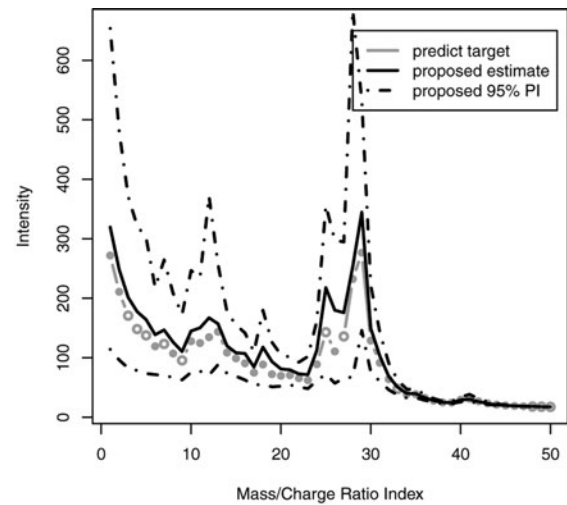


(A)

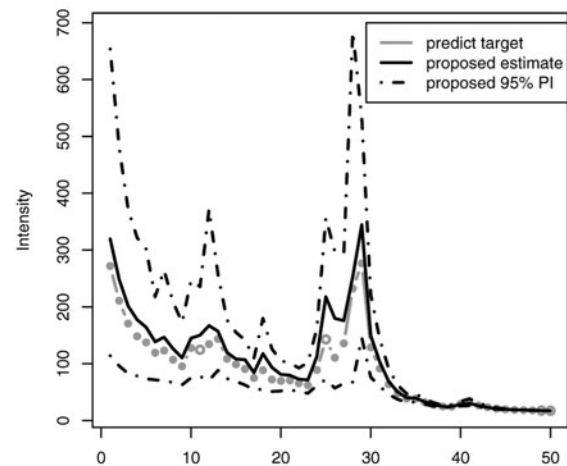


(B)

Figure 2. The differences of prediction errors between the existing Cholesky-based method (Bickel and Levina 2008a) and the proposed rank-based estimate with $\tau = 0.5$. The solid dots below the horizontal line represent the mass/charge ratio indices, where our proposed method outperforms the BL method. (A) Relative prediction errors for the control data and (B) relative prediction errors for the cancer data.



(A)



(B)

Figure 3. Prediction intervals for predicting protein mass spectroscopy intensities using our proposed method. The predict target, proposed point estimate ($\tau = 0.5$), and 95% prediction intervals are averaged over the testing set. The solid dots are the mass/charge ratio indices in which the proposed method outperforms the BL method. (A) 95% Prediction interval for the control data and (B) 95% prediction interval for the cancer data.

We plot the 95% prediction intervals for both the control data and the cancer data in Figure 3. Figure 3 shows that the estimated prediction intervals cover most of the predict target. This suggests that the transnormal model is a good fit to the data.

Supplementary Materials

The online supplement contains the appendices for the article.

Acknowledgement

The authors sincerely thank the associate editor and referees for their help comments and suggestions.

Funding

Jianqing Fan's research is supported in part by R01GM100474-04 and National Science Foundation grants DMS-1206464 and DMS-1406266. Lingzhou Xue's research is supported by the National Institutes of Health grant R01-GM072611-09 and National Science Foundation grant DMS-1505256. Hui Zou's research is supported by NSF grants DMS-0846068 and DMS-1505111.

References

- Adam, B., Qu, Y., Davis, J., Ward, M., Clements, M., Cazares, L., Semmes, O., Schellhammer, P., Yasui, Y., Feng, Z., and Wright, G. (2002), "Serum Protein Fingerprinting Coupled With a Pattern-Matching Algorithm Distinguishes Prostate Cancer From Benign Prostate Hyperplasia And Healthy Men," *Cancer Research*, 62, 3609–3614. [1733]
- Bickel, P., and Levina, E. (2008a), "Regularized Estimation of Large Covariance Matrices," *The Annals of Statistics*, 36, 199–227. [1729,1730,1731,1733]
- (2008b), "Covariance Regularization by Thresholding," *The Annals of Statistics*, 36, 2577–2604. [1728,1730]
- Box, G. E. P., and Cox, D. R. (1964), "An Analysis of Transformations" (with discussion), *Journal of the Royal Statistical Society, Series B*, 26, 211–252. [1726]
- Bradic, J., Fan, J., and Wang, W. (2011), "Penalized Composite Quasilielihood for Ultrahigh Dimensional Variable Selection," *Journal of the Royal Statistical Society, Series B*, 73, 325–349. [1727]
- Breiman, L., and Friedman, J. (1985), "Estimating Optimal Transformations for Multiple Regression and Correlation" (with discussion), *Journal of the American Statistical Association*, 80, 580–598. [1726,1727]
- Cai, T., and Zhou, H. (2012), "Minimax Estimation of Large Covariance Matrices Under ℓ_1 Norm" (with discussion), *Statistica Sinica*, 22, 1319–1378. [1730]
- Devlin, S., Gnanadesikan, R., and Kettenring, J. (1975), "Robust Estimation and Outlier Detection With Correlation Coefficients," *Biometrika*, 62, 531–545. [1729]
- Fan, Y., Härdle, W., Wang, W., and Zhu, L. (2013), "Composite Quantile Regression for the Single-Index Model," *SFB 649 Discussion Papers*, No.2013-010. [1728]
- Fan, J., Ke, T., Liu, H., and Xia, L. (2013), "QUADRO: A Supervised Dimension Reduction Method via Rayleigh Quotient Optimization," *The Annals of Statistics*, to appear. [1726]
- Fan, J., Liao, Y., and Mincheva, M. (2013), "Large Covariance Estimation by Thresholding Principal Orthogonal Complements" (with discussion), *Journal of the Royal Statistical Society, Series B*, 75, 603–680. [1728]
- Friedman, J. H. (2001), "Greedy Function Approximation: A Gradient Boosting Machine," *The Annals of Statistics*, 29, 1189–1232. [1727]
- (2002), "Stochastic Gradient Boosting," *Computational Statistics & Data Analysis*, 38, 367–378. [1727]
- Friedman, J. H., and Stuetzle, W. (1981), "Projection Pursuit Regression," *Journal of the American Statistical Association*, 76, 817–823. [1726,1727]
- Han, F., and Liu, H. (2012), "Semiparametric Principal Component Analysis," *Advances in Neural Information Processing Systems*, 171–179. [1726]
- Huang, J., Liu, N., Pourahmadi, M., and Liu, L. (2006), "Covariance Matrix Selection and Estimation via Penalised Normal Likelihood," *Biometrika*, 93, 85–98. [1729,1731]
- Kendall, M. (1948), *Rank Correlation Methods*, London: Charles Griffin and Co. Ltd. [1727,1729]
- Koenker, R., and Bassett, G. J. (1978), "Regression Quantiles," *Econometrica*, 46, 33–50. [1727]
- Koenker, R., and Geling, O. (2001), "Reappraising Medfly Longevity: A Quantile Regression Survival Analysis," *Journal of the American Statistical Association*, 96, 458–468. [1727]
- Koenker, R., and Xiao, Z. (2006), "Quantile Autoregression" (with discussion), *Journal of the American Statistical Association*, 101, 980–990. [1727]
- Levina, E., Rothman, A., and Zhu, J. (2008), "Sparse Estimation of Large Covariance Matrices via a Nested Lasso Penalty," *The Annals of Applied Statistics*, 2, 245–263. [1729,1733]
- Lin, Y., and Jeon, Y. (2003), "Discriminant Analysis Through a Semiparametric Model," *Biometrika*, 90, 379–392. [1726]
- Liu, H., Han, F., Yuan, M., Lafferty, J., and Wasserman, L. (2012), "High-Dimensional Semiparametric Gaussian Copula Graphical Models," *The Annals of Statistics*, 40, 2293–2326. [1726,1727]
- Liu, H., Lafferty, J., and Wasserman, L. (2009), "The Nonparanormal: Semiparametric Estimation of High-Dimensional Undirected Graphs," *Journal of Machine Learning Research*, 10, 1–37. [1726,1728]
- Mai, Q., and Zou, H. (2015), "Semiparametric Sparse Discriminant Analysis in Ultra-high Dimensions," *Journal of Multivariate Analysis*, 135, 175–188. [1726,1728]
- Qi, H., and Sun, D. (2006), "A Quadratically Convergent Newton Method for Computing the Nearest Correlation Matrix," *SIAM Journal on Matrix Analysis and Applications*, 28, 360–385. [1729]
- Stone, C. J. (1985), "Additive Regression and Other Nonparametric Models," *The Annals of Statistics*, 13, 689–705. [1726,1727]
- Tibshirani, R. (1988), "Estimating Transformations for Regression via Additivity and Variance Stabilization," *Journal of the American Statistical Association*, 83, 394–405. [1726,1727]
- Wang, H., and He, X. (2007), "Detecting Differential Expressions in GeneChip Microarray Studies: A Quantile Approach," *Journal of the American Statistical Association*, 102, 104–112. [1727]
- Wegkamp, M., and Zhao, Y. (2016), "Adaptive Estimation of the Copula Correlation Matrix for Semiparametric Elliptical Copulas," *Bernoulli*, 22, 1184–1226. [1726]
- Wu, W., and Pourahmadi, M. (2003), "Nonparametric Estimation of Large Covariance Matrices of Longitudinal Data," *Biometrika*, 90, 831–844. [1729]
- Wu, T., Yu, K., and Yu, Y. (2010), "Single-Index Quantile Regression," *Journal of Multivariate Analysis*, 101, 1607–1621. [1728]
- Xue, L., Ma, S., and Zou, H. (2012), "Positive Definite ℓ_1 Penalized Estimation of Large Covariance Matrices," *Journal of the American Statistical Association*, 107, 1480–1491. [1728,1729]
- Xue, L., and Zou, H. (2012), "Regularized Rank Estimation of High-dimensional Nonparanormal Graphical Models," *The Annals of Statistics*, 40, 2541–2571. [1726,1727,1731]
- (2014), "Rank-Based Tapering Estimation of Bandable Correlation Matrices," *Statistica Sinica*, 24, 83–100. [1726]
- Zhao, T., Roeder, K., and Liu, H. (2012), "Smooth-projected Neighborhood Pursuit for High-dimensional Nonparanormal Graph Estimation," *Advances in Neural Information Processing Systems*, 162–170. [1729]
- Zhu, L., Huang, M., and Li, R. (2012), "Semiparametric Quantile Regression With High-Dimensional Covariates," *Statistica Sinica*, 22, 1379–1401. [1728]
- Zou, H., and Yuan, M. (2008), "Composite Quantile Regression and the Oracle Model Selection Theory," *The Annals of Statistics*, 36, 1108–1126. [1727]