# Transfer Learning for High-dimensional Quantile Regression with Statistical Guarantee

Anonymous authors Paper under double-blind review

# Abstract

The task of transfer learning is to improve estimation/inference of a target model by migrating data from closely related source populations. In this article, we propose transfer learning algorithms for high-dimensional Quantile Regression (QR) models with the technique of convolution-type smoothing. Given the transferable source populations, we derive  $\ell_1/\ell_2$ -estimation error bounds for the estimators of the target regression coefficients under mild conditions. Theoretical analysis shows that the upper bounds are improved over those of the classical penalized QR estimator with only the target data, as long as the target and the sources are sufficiently similar to each other. When the set of informative sources is unknown, a transferable source detection algorithm is proposed to detect informative sources from all available sources. Thorough simulation studies justify our theoretical analysis.

# 1 Introduction

Transfer learning (Torrey & Shavlik, 2010) has been growing popular and drawing increasing attention in machine learning, which achieves great success in a wide range of real applications with limited available training data. Transfer learning aims to transfer knowledge from related source tasks/domains to enhance the learning or performance of the target task/domain, which typically involves two main subproblems. First, some criteria should be come up with to quantify the relatedness/similarity among target and source tasks. Intuitively, a high similarity would enhance the performance, while a low similarity would be harmful for the target task, which is known as "negative transfer" in the literature. Second, a transfer procedure should be carefully designed to transfer the "critical" knowledge from source domains, just like the human intelligence of leveraging prior experiences to tackle novel problems. A well designed transfer algorithm should not only identify the positive transfer sources thereby enlarging their impact, but also avoid the negative transfer in any case. All in all, transfer learning has become an active and promising research area, and substantial contributions has also been made recently to the theoretical guarantee for transfer learning in both supervised, semi-supervised, and unsupervised settings, see for example the context of classification by Cai & Wei (2021); Reeve et al. (2021), high-dimensional (generalized) linear regression by Li et al. (2022b); Tian & Feng (2022); Lin & Li (2022), graphical model by Li et al. (2022a); He et al. (2022). As far as we know, there exist no work on transfer learning for quantile regression and we aim to fill this gap in this paper.

## Comparison with the existing work and our contribution

A few works explore transfer learning under the high-dimensional setting. Bastani (2021) studied the transfer learning problem under a high-dimensional generalized linear models (GLM) with one single known transferable source data and the dimensionality p is assume to be larger than the sample size of the target dataset  $n_{\text{target}}$  while smaller than that of the source dataset  $n_{\text{source}}$ . A two-step transfer learning algorithm was developed and the  $\ell_1$ -estimation error bound was derived when the contrast  $\delta^*$  between target and source coefficients is  $\ell_0$ -sparse. More specifically, their estimator requires  $n_{\text{target}} = \mathcal{O}(s^2 \log^2(p/\xi)/\xi^2)$  as long as  $n_{\text{source}} \gtrsim \mathcal{O}(s^2 p^2 \log^2(p/\xi)/\xi^2)$ , where  $\xi > ||\delta^*||_1$  and  $s = ||\delta^*||_0$ . Li et al. (2022b) studied the highdimensional linear regression problem under some weaker assumptions, where both target and source samples are high-dimensional. Multiple source datasets are available and the transferable set may even be unknown in their paper. With  $\ell_q$ -sparse contrasts for  $q \in [0, 1)$  and  $\ell_0$ -sparse target parameter, the  $\ell_2$ -estimation error bound was derived and proved to be minimax optimal under some conditions. In the setting where the transferable set is unknown, a source detection algorithm was proposed to consistently select the informative sources. Tian & Feng (2022) further investigated multi-source transfer learning on high-dimensional generalized linear models (GLM). They assumed both target and source data to be high-dimensional and the contrast to be  $\ell_1$ -sparse. Given the informative sources to transfer, the  $\ell_1/\ell_2$ -estimation error was derived and proved to be minimax optimal under mild conditions. Tian & Feng (2022) also established a transferable source detection algorithm to identify the informative sources. In addition, they constructed the corresponding confidence interval for individual regression parameter. Li et al. (2021) proposed a federated transfer learning approach to consolidate data from different populations and from multiple medical associations. The target and source data are both high-dimensional in their discussion and they characterized the contrasts to be  $\ell_0$ -sparse. Compared with Tian & Feng (2022), their approach achieves a faster convergence rate under some conditions and has weaker requirements on the level of heterogeneity for data from diverse populations.

Inspired by the two-step algorithm in Bastani (2021), Li et al. (2022b) and Tian & Feng (2022), we propose a multi-source transfer learning method under high-dimensional quantile regression. To overcome the nonsmoothness and non-convexity of the quantile loss, motivated by He et al. (2021) and Tan et al. (2022), we employ the convolution-type smoothed quantile regression. Assuming the contrasts between the target and each source coefficients to be  $\ell_1$ -sparse, we establish the  $\ell_1/\ell_2$ -estimation error bounds that are proved to be sharper than the bounds of the classical  $\ell_1$ -penalized quantile regression (Belloni & Chernozhukov, 2011) under some conditions. Notably, our results need the sample size of the target data to be  $\mathcal{O}(s^3 \log p)$ . However, the results in Li et al. (2022b) and Tian & Feng (2022) only require the size of target sample to be  $\mathcal{O}(s^2 \log p)$ . The difference is caused by the bandwidth involved in the smoothing method. The estimation error bounds clearly relies on the smoothing bandwidth, which leads to a more restricted sample size.

In this paper, we propose transfer learning algorithms for quantile regression with high-dimensional data and we assume the contrast between target and source coefficients to be  $\ell_0$ -sparse or  $\ell_1$ -sparse. In the setting where the sources are sufficiently close to the target, our theoretical analysis and simulation results show that the estimation error bound of the target coefficients is improved compared to the classical  $\ell_1$ -penalized quantile regression model (Belloni & Chernozhukov, 2011) using only the target data under mild conditions. To overcome the lack of smoothness and convexity of the check loss, we employed the convolution-type smoothed quantile regression and analyzed the (local) restricted strong convexity of the empirical smoothed quantile loss functions in the transferring and debiasing steps. We also extended the source detection algorithm in Tian & Feng (2022) to the quantile regression setting. Simulation results show that the algorithm works well in discovering useful sources. In contrast to the case with  $\ell_1$ -sparse contrasts, the algorithm with  $\ell_0$ -sparse contrasts learns the source coefficients independently, which greatly reduces the communications cost across different sources. Furthermore, the algorithm with  $\ell_0$ -sparse contrasts has fewer assumptions on the level of heterogeneity for data from different sources.

The most related work is a concurrent paper by Zhang & Zhu (2022), which also considered the smoothed quantile regression models under transfer learning framework. They proposed a smoothed two-step transfer learning algorithm as well as a new source detection method based on the K-means clustering algorithm, which does not need the input of a threshold in contrast to the source detection algorithm in Tian & Feng (2022). In addition, they further extended their work to the distributed quantile regression and model averaging setup. However, compared with Zhang & Zhu (2022), our work doesn't require the restrictive conditions on the kernels that  $\sup_{|h| \leq 1} K(u/h)/h < M_k$  almost everywhere in u. In addition, given that the contrasts are characterized in  $\ell_0$ -norm instead of  $\ell_1$ -norm, we introduce an algorithm which is motivated from Li et al. (2022b) and Li et al. (2021). The  $\ell_1/\ell_2$ -estimation error bounds are also established and proved to be sharper than the bounds of the classical  $\ell_1$ -penalized quantile regression (Belloni & Chernozhukov, 2011) under some conditions.

Before ending this section, we introduce the notations used throughout the paper. For any symmetric, positive semidefinite matrix  $\mathbf{A} \in \mathbb{R}^{k \times k}$ , if its vector of eigenvalues is denoted by  $\gamma(\mathbf{A})$  and ordered as  $\gamma_1(\mathbf{A}) \geq \ldots \geq \gamma_p(\mathbf{A}) \geq 0$ , the operator norm of  $\mathbf{A}$  is  $||\mathbf{A}||_2 = \gamma_1(\mathbf{A})$ . Moreover, the vector norm induced

by  $\boldsymbol{A}$  is  $||\boldsymbol{u}||_A = ||\boldsymbol{A}^{1/2}\boldsymbol{u}||_2$  for any  $\boldsymbol{u} \in \mathbb{R}^k$ . For any real numbers s and  $t, s \vee t$  denotes  $\max(s, t)$  and  $s \wedge t$  denotes  $\min(s, t)$ . For two sequences  $\{a_n\}_{n\geq 1}$  and  $\{b_n\}_{n\geq 1}$ , which consist of non-negative numbers,  $a_n \leq b_n$  means that there exists a constant C > 0 such that  $a_n \leq Cb_n$ .  $a_n \approx b_n$  is equivalent to  $a_n \leq b_n$  and  $b_n \leq a_n$ . For r, l > 0, define the  $\ell_2$ -ball and  $\ell_1$ -cone as

$$\mathbb{B}_{\Sigma}(r) = \{ \boldsymbol{\delta} \in \mathbb{R}^p : ||\boldsymbol{\delta}||_{\Sigma} \leq r \} \text{ and } \mathbb{C}_{\Sigma}(l) = \{ \boldsymbol{\delta} \in \mathbb{R}^p : ||\boldsymbol{\delta}||_1 \leq l ||\boldsymbol{\delta}||_{\Sigma} \}.$$

# 2 Methodology

#### 2.1 Problem Setup

Given the predictors  $x \in \mathbb{R}^p$  and a scalar response variable  $y \in \mathbb{R}$ , the  $\tau$ -th conditional quantile functions of y given x is written as

$$F_{y|x}^{-1}(\tau) = \inf\{y : F_{y|x}(y) \ge \tau\},\$$

where  $F_{y|x}(\cdot)$  is the conditional distribution function of y given x. Consider the following linear quantile regression model at a given  $\tau \in (0, 1)$ :

$$F_{y|x}^{-1}(\tau) = \langle \boldsymbol{x}, \boldsymbol{\beta}^*(\tau) \rangle = \boldsymbol{x}^{\mathrm{T}} \boldsymbol{\beta}^*(\tau),$$

where  $\boldsymbol{\beta}^*(\tau) = (\beta_1^*(\tau), \dots, \beta_p^*(\tau))^{\mathrm{T}} \in \mathbb{R}^p$  is the true quantile regression coefficient.

Let  $\{(y_i, \boldsymbol{x}_i)\}_{i=1}^n$  be a random sample from  $(y, \boldsymbol{x})$ . The preceding model assumption is equivalent to the following model

$$y_i = \boldsymbol{x}_i^{\mathrm{T}} \boldsymbol{\beta}^* + \epsilon_i \text{ and } \mathbb{P}(\epsilon_i \leq 0 | \boldsymbol{x}_i) = \tau.$$

The  $\ell_1$ -penalized quantile regression estimator (Belloni & Chernozhukov, 2011) is generally defined as one of the solution to the optimization problem

$$\min_{\boldsymbol{\beta}=(\beta_1,\ldots,\beta_p)^{\mathrm{T}}\in\mathbb{R}^p} \left\{ \underbrace{\frac{1}{n} \sum_{i=1}^n \rho_{\tau}(y_i - \boldsymbol{x}_i^{\mathrm{T}}\boldsymbol{\beta})}_{=:\hat{Q}(\boldsymbol{\beta})} + \lambda ||\boldsymbol{\beta}||_1 \right\},\tag{1}$$

where  $\rho_{\tau}(u)$  is defined as  $\rho_{\tau}(u) = u\{\tau - \mathbb{1}(u < 0)\}$ , also referred to as the  $\tau$ -quantile check loss function. Let  $\hat{F}(\cdot; \beta)$  be the empirical cumulative distribution function of the residuals  $\{r_i(\beta) := y_i - \mathbf{x}^{\mathsf{T}}\beta\}_{i=1}^n$ , i.e.,  $\hat{F}(u; \beta) = (1/n) \sum_{i=1}^n \mathbb{1}\{r_i(\beta) \le u\}$  for any  $u \in \mathbb{R}$ . Then the empirical quantile loss  $\hat{Q}(\beta)$  in (1) can be written as

$$\hat{Q}(\boldsymbol{\beta}) = \int_{-\infty}^{\infty} \rho_{\tau}(u) d\hat{F}(u; \boldsymbol{\beta}).$$
<sup>(2)</sup>

As the empirical cumulative distribution function  $\hat{F}(\cdot; \beta)$  is discontinuous, the empirical quantile loss is nondifferentiable, which brings great challenges to both computation and statistical theory establishment. The kernel smoothing method (Horowitz, 1998) is commonly utilized to tackle this issue. However, the smoothed loss is still non-convex, thereby we further consider the convolution-type smoothed quantile loss function, which is not only convex but also differentiable and brings great convenience in terms of both computation and theoretical analysis. In the following, we briefly introduce the convolution-type smoothed quantile loss function, which was firstly introduced by Tan et al. (2022).

Let  $K(\cdot)$  be a non-negative kernel function that is symmetric around 0 and integrates to 1, and h > 0 be a bandwidth. That is

$$K_h(u) = (1/h)K(u/h), \ \bar{K}(u) = \int_{-\infty}^u K(v)dv \text{ and } \bar{K}_h(u) = \bar{K}(u/h), \ u \in \mathbb{R}.$$

The empirical smoothed loss function can be defined as

$$\hat{Q}_{h}(\beta) = \frac{1}{n} \sum_{i=1}^{n} l_{h}(y_{i} - \boldsymbol{x}_{i}^{\mathrm{T}}\beta) \text{ with } l_{h}(u) = (\rho_{\tau} * K_{h})(u) = \int_{-\infty}^{\infty} \rho_{\tau}(v) K_{h}(v - u) dv,$$

where \* denotes the convolution operator. Therefore, the  $\ell_1$ -penalized convolution smoothed estimator is given by

$$\hat{\boldsymbol{\beta}} \in \operatorname*{arg\,min}_{\boldsymbol{\beta} \in \mathbb{R}^p} \Big\{ \hat{Q}_h(\boldsymbol{\beta}) + \lambda ||\boldsymbol{\beta}||_1 \Big\},$$

where the smoothing bandwidth h adapts to the sample size n and the dimension p while  $\hat{\beta}$  depends on the quantile index  $\tau$ , bandwidth h, and penalty level  $\lambda$ .

In the following, we consider the multi-source transfer learning scenario, where we have a target data set  $(\mathbf{X}^{(0)}, \mathbf{y}^{(0)})$  and K source data sets with the k-th source denoted as  $(\mathbf{X}^{(k)}, \mathbf{y}^{(k)})$ , where  $\mathbf{X}^{(k)} \in \mathbb{R}^{n_k \times p}$ ,  $\mathbf{y}^{(k)} \in \mathbb{R}^{n_k}$  for  $k = 0, \ldots, K$ . The *i*-th row of  $\mathbf{X}^{(k)}$  and the *i*-th element of  $\mathbf{y}^{(k)}$  are denoted as  $\mathbf{x}_i^{(k)}$  and  $\mathbf{y}_i^{(k)}$ , respectively. The goal is to transfer useful information from the source datasets to improve the estimation accuracy of the target parameters. Denote the true target parameter as  $\boldsymbol{\beta}^* = \boldsymbol{\omega}^{(0)}$ . We assume the responses in the target and source data all follow the linear quantile regression model, that is,

$$y_i^{(k)} = \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega}^{(k)} \rangle + \epsilon_i^{(k)} \text{ and } \mathbb{P}(\epsilon_i \le 0 | \boldsymbol{x}_i^{(k)}) = \tau, \ k = 0, \dots, K.$$

We build our quantile regression transfer learning procedure in the high-dimensional regime with a sparsity assumption. In other words, we assume the dimension p is much larger than the sample size  $n_k$  for all k while the target model is s-sparse, which satisfies  $||\boldsymbol{\beta}^*||_0 = s$ . Define the k-th contrast as  $\boldsymbol{\delta}^{(k)} = \boldsymbol{\beta}^* - \boldsymbol{\omega}^{(k)}$  and  $||\boldsymbol{\delta}^{(k)}||_q$  is referred to as the transferring level of source k in the literature, where  $q \in \{0, 1\}$ . Define the level-m transferring set  $\mathcal{A}_m = \{k : ||\boldsymbol{\delta}^{(k)}||_q \leq m\}$  as the set of sources which has transferring level lower than m. Denote  $n_{\mathcal{A}_m} = \sum_{k \in \mathcal{A}_m} n_k$ ,  $\alpha_k = n_k/(n_{\mathcal{A}_m} + n_0)$  for  $k \in \{0\} \cup \mathcal{A}_m$  and  $K_{\mathcal{A}_m} = |\mathcal{A}_m|$ .

As stated in the introduction, we will consider two types of transferring level, corresponding to  $q \in \{0, 1\}$  respectively. In the case of q = 0, the transferring set corresponds to the source data whose contrast vectors have at most m nonzero elements. In the case of q = 1, all the coefficients of the contrast vectors can be nonzero, but their absolute magnitude decays at a relatively rapid rate. It will be seen later that as long as m is relatively small, the source data in  $\mathcal{A}_m$  can be useful in improving the estimation accuracy of  $\beta^*$ . In addition, the logistic of the algorithm with  $\ell_1$ -normed  $\mathcal{A}_m$  and the algorithm with  $\ell_2$ -normed  $\mathcal{A}_m$  are quite different and we will elaborate on these two different algorithms in the following sections.

#### 2.2 Algorithm with $\ell_1$ -norm constrained transferring set

In this section, we propose the transfer learning algorithm with  $\ell_1$ -norm constrained transferring set, which is motivated by Tian & Feng (2022). This algorithm involves two steps. The first step of our algorithm is to transfer the information from useful sources by pooling all the data in transferable set  $\mathcal{A}_m$  and target set  $\mathcal{A}_0$  to obtain a primal estimator. We also call it the transferring step. To be more precise, we define a total smoothed loss function for the target and source datasets in the transferable set  $\mathcal{A}_m$ , i.e.,

$$\hat{Q}_h(\boldsymbol{\omega}) = \frac{1}{n_{\mathcal{A}_m} + n_0} \sum_{k \in \mathcal{A}_m} \sum_{i=1}^{n_k} l_h(y_i^{(k)} - \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega} \rangle),$$

where

$$l_h(u) = (\rho_\tau * K_h)(u) = \int_{-\infty}^{\infty} \rho_\tau(v) K_h(v-u) dv.$$

Then for the transferring step, we aim to find the minimizer to the following optimization problem with respect to  $w \in \mathbb{R}^p$ :

minimize 
$$\{\hat{Q}_h(\boldsymbol{\omega}) + \lambda_{\boldsymbol{\omega}} ||\boldsymbol{\omega}||_1\}.$$

We denote the minimizer as  $\hat{\omega}^{\mathcal{A}_m}$ , i.e.,  $\hat{\omega}^{\mathcal{A}_m} = \arg \min_{\omega} \left\{ \hat{Q}_h(\omega) + \lambda_{\omega} ||\omega||_1 \right\}$ . By selecting an appropriate bandwidth h, the iteratively reweighted  $\ell_1$ -penalized SQR estimator proposed by Tan et al. (2022) shares the same upper bounds for both  $\ell_1$  and  $\ell_2$  errors as the  $\ell_1$ -QR estimator, as indicated by Belloni & Chernozhukov (2011). Furthermore, they introduced coordinate descent and ADMM-based algorithms for solving  $\ell_1$ -penalized quantile regression, which are computationally efficient especially for large-scale problems.

Denote the true parameter in the first step as  $\omega^{\mathcal{A}_m}$ , and  $\omega^{\mathcal{A}_m}$  has the following explicit form:

$$\omega^{\mathcal{A}_m} = oldsymbol{eta} + \delta^{\mathcal{A}_m},$$

where  $\delta^{\mathcal{A}_m} = \sum_{k \in \mathcal{A}_m} \alpha_k \delta^{(k)}$  and  $\alpha_k = n_k / (n_{\mathcal{A}_m} + n_0)$ . For the second step (the debiasing step), we correct the bias,  $\delta^{\mathcal{A}_m}$ , based on the estimator  $\hat{\omega}^{\mathcal{A}_m}$  acquired in the transferring step. The smoothed loss function for the target data with respect to  $\delta$  is defined as

$$\hat{Q}_g^{(0)}(\hat{\boldsymbol{\omega}}^{\mathcal{A}_m} + \boldsymbol{\delta}) = rac{1}{n_0}\sum_{i=1}^{n_0} l_g(y_i^{(0)} - \langle \boldsymbol{x}_i^{(0)}, \hat{\boldsymbol{\omega}}^{\mathcal{A}_m} + \boldsymbol{\delta} 
angle).$$

The error of the debiasing step is under control for a relatively small m, since  $\delta^{\mathcal{A}_m}$  is a  $\ell_1$ -sparse highdimensional vector.

We call this algorithm Oracle  $\ell_1$ -Trans-SQR as we first assume that all useful sources are known as a priori. Algorithm 1 formally presents the Oracle  $\ell_1$ -Trans-SQR algorithm.

# Algorithm 1: Oracle $\ell_1$ -Trans-SQR

**Input:** Target data  $(\boldsymbol{X}^{(0)}, \boldsymbol{y}^{(0)})$ , source data  $\{(\boldsymbol{X}^{(k)}, \boldsymbol{y}^{(k)})\}_{k=1}^{K}$ , penalty parameters  $\lambda_{\omega}$  and  $\lambda_{\delta}$ , transferring set  $\mathcal{A}_{m}$ .

**Output:** The estimator  $\beta$ .

1 Transferring step:

$$\hat{\boldsymbol{\omega}}^{\mathcal{A}_m} \leftarrow \operatorname*{arg\,min}_{\boldsymbol{\omega}} \Big\{ \hat{Q}_h(\boldsymbol{\omega}) + \lambda_{\boldsymbol{\omega}} ||\boldsymbol{\omega}||_1 \Big\},$$

2 Debiasing step:

$$\hat{oldsymbol{eta}}^{\mathcal{A}_m} \leftarrow rgmin_{\delta} \left\{ \hat{Q}_g^{(0)}(\hat{oldsymbol{\omega}}^{\mathcal{A}_m} + oldsymbol{\delta}) + \lambda_{\delta} ||oldsymbol{\delta}||_1 
ight\},$$

3 return  $\hat{eta} = \hat{oldsymbol{\omega}}^{\mathcal{A}_m} + \hat{oldsymbol{\delta}}^{\mathcal{A}_m}.$ 

If  $\mathcal{A}_m$  is unknown, then we need a detection algorithm to find useful transferable sets in practice. We propose a transferrable source detection algorithm which is inspired from the Algorithm 2 in Tian & Feng (2022). Firstly, partition the target data into q subsets. Secondly, fit the penalized smoothed quantile regression on each combination of (q-1) target subsets and calculate the loss on the remaining target subset. In the following, consider the average cross-validation loss  $\hat{L}_0^{(0)}$ . Run the transferring step on each combination of (q-1) target subsets and each source data, and evaluate the loss function on the remaining target subset. Similarly compute the average cross-validation loss  $\hat{L}_0^{(k)}$  for each source. Thirdly, calculate the difference between  $\hat{L}_0^{(0)}$  and  $\hat{L}_0^{(k)}$  for each k and compare it with a predefined threshold. Finally select the sources whose difference is less than the threshold and include them in the set  $\hat{\mathcal{A}}$ . The detailed transferable source detection procedure is summarized in Algorithm 2.

With the transferrable source detection algorithm, we propose a feasible Algorithm 3 in practice, in which we first detect useful source datasets  $\hat{\mathcal{A}}$  by Algorithm 2 and then run Algorithm 1 using datasets  $\{(\mathbf{X}^{(k)}, \mathbf{y}^{(k)})\}_{k \in \{0\} \cup \hat{\mathcal{A}}}$ .

## 2.3 The proposed algorithm with $\ell_0$ -norm constrained transferring set

In this section we consider a more strict transferable set  $\mathcal{A}'_m = \{k : ||\delta^{(k)}||_0 \leq m\}$ , where the  $\ell_1$ -norm discussed in Section 2.2 is replaced by  $\ell_0$ -norm. Compared with the  $\ell_1$ -norm, the theoretical analysis of the transfer learning procedure under  $\ell_0$ -norm is free of the restrictive Assumption 3.4 below, which requires "sufficient" similarity between the target covariance matrix and transferable source covariance matrices. However, as  $\ell_0$ -norm is not additive, it is not easy to combine target and source data to estimate a primary

Algorithm	<b>2</b> :	Transferable	Source	Detection
-----------	------------	--------------	--------	-----------

- **Input:** Target data  $(\boldsymbol{X}^{(0)}, \boldsymbol{y}^{(0)})$ , all source data  $\{(\boldsymbol{X}^{(k)}, \boldsymbol{y}^{(k)})\}_{k=1}^{K}$ , a threshold  $C_0$ , penalty parameters  $\{\{\lambda^{(k)[r]}\}_{k=0}^{K}\}_{r=1}^{q}$ , where q is the number of folds chosen.
- **Output:** The set of transferable sources  $\hat{\mathcal{A}}$ .

1 Randomly divide  $(\boldsymbol{X}^{(0)}, \boldsymbol{y}^{(0)})$  into q equal-sized sets  $\{(\boldsymbol{X}^{(0)[i]}, \boldsymbol{y}^{(0)[i]})\}_{i=1}^{q}$ .

- 2 for r = 1 to q do
- $\begin{array}{c|c} \mathbf{3} & \hat{\boldsymbol{\beta}}^{(0)[r]} \leftarrow \text{fit the penalized quantile regression on } \{(\boldsymbol{X}^{(0)[i]}, \boldsymbol{y}^{(0)[i]})\}_{i=1}^{q} \setminus (\boldsymbol{X}^{(0)[r]}, \boldsymbol{y}^{(0)[r]}) \text{ with penalty parameter } \lambda^{(0)[r]}. \end{array}$
- 4  $\hat{\boldsymbol{\beta}}^{(k)[r]} \leftarrow \text{run the transferring step in Algorithm 1 with} \\ \{ (\boldsymbol{X}^{(0)[i]}, \boldsymbol{y}^{(0)[i]}) \}_{i=1}^{q} \setminus (\boldsymbol{X}^{(0)[r]}, \boldsymbol{y}^{(0)[r]}) \cup (\boldsymbol{X}^{(k)}, \boldsymbol{y}^{(k)}) \text{ and penalty parameter } \lambda^{(k)[r]} \text{ for all } k \neq 0.$
- 5 Calculate the loss function  $\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]})$  on  $(\boldsymbol{X}^{(0)[r]}, \boldsymbol{y}^{(0)[r]})$  for  $k = 1, \dots, K$ .
- $\begin{array}{l} \mathbf{6} \ \ \hat{L}_{0}^{(k)} \leftarrow \sum_{r=1}^{q} \hat{L}_{0}^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]})/q, \quad \hat{L}_{0}^{(0)} \leftarrow \sum_{r=1}^{q} \hat{L}_{0}^{[r]}(\hat{\boldsymbol{\beta}}^{(0)[r]})/q, \ \hat{\sigma} = \sqrt{\sum_{r=1}^{q} (\hat{L}_{0}^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]}) \hat{L}_{0}^{(0)})^{2}/(q-1)}.\\ \mathbf{7} \ \ \hat{\mathcal{A}} \leftarrow \{k \neq 0: \hat{L}_{0}^{(k)} \hat{L}_{0}^{(0)} \leq C_{0}(\hat{\sigma} \lor 0.01)\}.\\ \mathbf{8} \ \ \mathbf{return} \ \hat{\mathcal{A}}. \end{array}$

## Algorithm 3: Trans-SQR

**Input:** Target data  $(\boldsymbol{X}^{(0)}, \boldsymbol{y}^{(0)})$ , all source data  $\{(\boldsymbol{X}^{(k)}, \boldsymbol{y}^{(k)})\}_{k=1}^{K}$ , a threshold  $C_0$  and penalty parameters  $\{\{\lambda^{(k)[r]}\}_{k=0}^{K}\}_{r=1}^{q}$ .

**Output:** The estimator  $\hat{\boldsymbol{\beta}}$ .

- 1 Run Algorithm 2 (Transferable Source Detection Algorithm) and output  $\hat{\mathcal{A}}$ .
- 2 Run Algorithm 1 (Oracle Trans-SQR) using data  $\{(X^{(k)}, y^{(k)})\}_{k \in \{0\} \cup \hat{\mathcal{A}}}$ .

3 return  $\hat{\beta}$ .

estimator for the true target parameter. Instead, we correct each source data independently and incorporate the corrected source and target data to make predictions. Certain adjustments need to be made on the proposed transfer learning procedure in Algorithm 1.

This  $\ell_0$ -norm constrained transfer algorithm is inspired by the idea in Li et al. (2021). Unlike the transferring step in Algorithm 1, the first step of the algorithm in this section is to train each source separately to get primal estimators of  $\omega^{(k)}$ ,  $k \in \{1, \ldots, K\}$ , where the smoothed loss function for each source k is

$$\hat{Q}_h(\boldsymbol{\omega}) = rac{1}{n_k} \sum_{i=1}^{n_k} l_h(y_i^{(k)} - \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega} \rangle).$$

In the second step, as the debiasing step in Algorithm 1, we adjust for the differences  $\hat{\delta}^{(k)}$  for all k using the target data, which is obtained via

$$\hat{\boldsymbol{\delta}}^{(k)} = \arg\min_{\boldsymbol{\delta}} \left\{ \hat{Q}_g^{(0)}(\hat{\boldsymbol{\omega}}^{(k)} + \boldsymbol{\delta}) + \lambda_{\boldsymbol{\delta}} ||\boldsymbol{\delta}||_1 \right\},\$$

where the smoothed loss function with respect to  $\boldsymbol{\delta}$  is defined as

$$\hat{Q}_{g}^{(0)}(\hat{\boldsymbol{\omega}}^{(k)} + \boldsymbol{\delta}) = rac{1}{n_{0}}\sum_{i=1}^{n_{0}}l_{g}(y_{i}^{(0)} - \langle \boldsymbol{x}_{i}^{(0)}, \hat{\boldsymbol{\omega}}^{(k)} + \boldsymbol{\delta} 
angle).$$

Then a threshold for each  $\hat{\delta}^{(k)}$  is computed by only keeping the largest  $\sqrt{n_0/\log p}$  elements of  $\hat{\delta}^{(k)}$  and letting all the other elements be zero. In the third step, with the estimated "bias" from the second step, the corrected source data has the following form:

$$\left\{\left(oldsymbol{X}^{(k)},oldsymbol{y}^{(k)}+oldsymbol{X}^{(k)} ilde{oldsymbol{\delta}}^{(k)}
ight)
ight\}_{k=1}^{K}$$

Algorithm 4: Oracle Trans-SQR with  $\ell_0$ -norm constrained transferring set

**Input:** Target data  $(\boldsymbol{X}^{(0)}, \boldsymbol{y}^{(0)})$ , source data  $\{(\boldsymbol{X}^{(k)}, \boldsymbol{y}^{(k)})\}_{k=1}^{K}$ , penalty parameters  $\lambda_{\omega}, \lambda_{\delta}$  and  $\lambda_{\beta}$ , transferring set  $\mathcal{A}'_{m}$ . Let  $n = n_0 + n_{\mathcal{A}'_{m}}$ .

**Output:** The estimator  $\hat{\boldsymbol{\beta}}$ .

1 For each  $k \in \mathcal{A}'_m$ ,

$$\hat{\boldsymbol{\omega}}^{(k)} \leftarrow \operatorname*{arg\,min}_{\boldsymbol{\omega}} \left\{ \hat{Q}_h(\boldsymbol{\omega}) + \lambda_{\omega}^{(k)} ||\boldsymbol{\omega}||_1 
ight\}.$$

**2** For each  $k \in \mathcal{A}'_m$ ,

$$\hat{\boldsymbol{\delta}}^{(k)} \leftarrow rgmin_{\delta} \left\{ \hat{Q}_{g}^{(0)}(\hat{\boldsymbol{\omega}}^{(k)} + \boldsymbol{\delta}) + \lambda_{\delta} || \boldsymbol{\delta} ||_{1} 
ight\}$$

Threshold  $\hat{\delta}^{(k)}$  via  $\tilde{\delta}^{(k)} = \mathcal{H}_{\sqrt{n_0/\log p}}(\hat{\delta}^{(k)})$ , where  $\mathcal{H}_k(\boldsymbol{b})$  is formed by setting all but the largest k elements of  $\boldsymbol{b}$  to zero.

**3** Joint estimation using source and target data:

$$\hat{\boldsymbol{\beta}} \leftarrow \operatorname*{arg\,min}_{\boldsymbol{\beta}} \left\{ \frac{1}{n} \sum_{i=1}^{n_0} l_r(y_i^{(0)} - \langle \boldsymbol{x}_i^{(0)}, \boldsymbol{\beta} \rangle) \right. \\ \left. + \frac{1}{n} \sum_{k \in \mathcal{A}'_m} \sum_{i=1}^{n_k} l_r(y_i^{(k)} - \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\beta} - \tilde{\boldsymbol{\delta}}^{(k)} \rangle) + \lambda_{\boldsymbol{\beta}} ||\boldsymbol{\beta}||_1 \right\}.$$

4 return  $\hat{\beta}$ .

Then, we combine all the corrected sources and target data to estimate the parameter  $\beta$  which is of our interest. The above algorithm estimate the source parameters and the contrast vectors individually, while in the Oracle  $\ell_1$ -Trans-SQR proposed in Section 2.2, a pooled analysis is conducted with data from target and sources, which relies on the homogeneous designs of the covariance matrices among target and source data.

## **3** Statistical theory

In this section, we establish theoretical guarantees on the algorithms in the above section.

Assumption 3.1. The conditional density of  $\epsilon$  given  $\boldsymbol{x}$  satisfies  $f_{\epsilon|\boldsymbol{x}}(u) \leq f_u$  almost surely for some  $f_u \geq f_l > 0$ . Moreover, there exists  $l_0 > 0$  such that  $|f_{\epsilon|\boldsymbol{x}}(u) - f_{\epsilon|\boldsymbol{x}}(v)| \leq l_0|u-v|$  for all  $u, v \in \mathbb{R}$  almost surely, and

$$\inf_{t\in[0,1],\boldsymbol{v}\in\mathbb{S}^{p-1}}\mathbb{E}[f_{\epsilon|\boldsymbol{x}}(t\langle\boldsymbol{z},\boldsymbol{v}\rangle)\langle\boldsymbol{z},\boldsymbol{v}\rangle^2]\geq f_l.$$

Assumption 3.2. The kernel function  $K : \mathbb{R} \to [0,\infty)$  is symmetric around zero, and satisfies  $\int_{-\infty}^{\infty} K(u) du = 1$  and  $\int_{-\infty}^{\infty} u^2 K(u) du < \infty$ . For  $l = 1, 2, ..., \text{ let } \kappa_l = \int_{-\infty}^{\infty} |u|^l K(u) du$  be the *l*-th absolute moment of  $K(\cdot)$ . Assume  $\sup_{u \in \mathbb{R}} K(u) \leq \kappa_u$  for some  $\kappa_u \in (0, 1]$ .

Assumption 3.3. The covariate vector  $\boldsymbol{x}$  is compactly supported with

$$\zeta_p := \sup_{\boldsymbol{x} \in \mathbb{R}^p} ||\Sigma^{-1/2} \boldsymbol{x}||_2 < \infty$$

and  $||\boldsymbol{x}||_{\infty} \leq B$  almost surely for some  $B \geq 1$ , where  $\Sigma = \mathbb{E}(\boldsymbol{x}\boldsymbol{x}^{\mathrm{T}})$  is positive definite. Without loss of generality, assume B = 1. In addition to  $\zeta_p$ ,  $\sup_{\boldsymbol{u} \in \mathbb{S}^{p-1}} \mathbb{E}(|\boldsymbol{u}^{\mathrm{T}}\Sigma^{-1/2}\boldsymbol{x}|^4) \leq \zeta_p^2$ .

Under Assumption 3.3, for every  $\delta \in (0, 1]$ , define

$$\eta_{\delta} = \inf \left\{ \eta > 0 : \mathbb{E} \left[ (\boldsymbol{z}^{\mathrm{T}} \boldsymbol{v})^{2} \mathbb{1} (|\boldsymbol{z}^{\mathrm{T}} \boldsymbol{v}| > \eta) \right] \le \delta \text{ for all } \boldsymbol{v} \in \mathbb{S}^{p-1} \right\},$$
(3)

where  $\boldsymbol{z} = \Sigma^{-1/2} \boldsymbol{x}$ . Since  $\mathbb{E}(\boldsymbol{z}^{\mathrm{T}} \boldsymbol{v})^2 = 1$  for any  $\boldsymbol{v} \in \mathbb{S}^{p-1}$ ,  $\eta_{\delta}$  is well-defined for each  $\delta$ , and depends implicitly on the underlying distribution of  $\boldsymbol{z}$ .

Assumption 3.4. Denote

$$\tilde{\Sigma} = \sum_{k=0}^{K} \alpha_k \int_0^1 \nabla^2 Q^{(k)} ((1-t)\boldsymbol{\beta}^* + t\boldsymbol{\omega}^*) dt$$
$$\tilde{\Sigma}^{(k)} = \int_0^1 \nabla^2 Q^{(k)} ((1-t)\boldsymbol{\beta}^* + t\boldsymbol{\omega}^{(k)}) dt,$$

where  $\nabla^2 Q^{(k)}((1-t)\boldsymbol{\beta}^* + t\boldsymbol{\omega}) = \mathbb{E}\{f_{\epsilon|\boldsymbol{x}}(t\boldsymbol{\omega} - t\boldsymbol{\beta}^*) \cdot \boldsymbol{x}^{(k)}(\boldsymbol{x}^{(k)})^{\mathrm{T}}\}$ . Define

$$C_1 = \sup_{0 \le k \le K} ||\tilde{\Sigma}^{-1}\tilde{\Sigma}^{(k)}||_1.$$

Let  $C_1$  be bounded, that is  $C_1 < \infty$ .

Assumption 3.1 imposes the Lipschitz continuity on the conditional density  $f_{\epsilon|\mathbf{x}}(\cdot)$ . Assumption 3.2 holds for most commonly used kernel functions, for instance, uniform kernel, Gaussian kernel, etc.

Compared with Tian & Feng (2022) and Li et al. (2022b), Assumption 3.3 is different. Note that quantile regression has Hessian matrix  $\nabla^2 \hat{Q}_h(\beta) = (1/n) \sum_{i=1}^n K_h(\boldsymbol{x}_i^{\mathrm{T}}\beta - y_i)\boldsymbol{x}_i\boldsymbol{x}_i^{\mathrm{T}}$ , where  $\hat{Q}_h(\beta)$  is the smoothed empirical quantile loss and  $K_h(u) = (1/h)K(u/h)$ . Unlike the generalized linear regression, there is a smoothing bandwidth h in the denominator. We import Assumption 3.3 to provide convenience for bounding the difference  $\nabla \hat{Q}_h(\beta) - \nabla \hat{Q}_h(\beta^*)$  in the debiasing step.

Assumption 3.4 restricts the difference between the target covariance matrix and transferable source covariance matrix in some sense, which guarantees the estimator at the transferring step is close to the true parameter  $\beta^*$ . This assumption is commonly used in other transfer learning works, Tian & Feng (2022); Li et al. (2022b); Zhang & Zhu (2022); Huang et al. (2022).

Formally, we consider the parameter space

$$\Theta(s,m) = \left\{ \boldsymbol{\beta}^*, \{\boldsymbol{\omega}^{(k)}\}_{k \in \mathcal{A}_m} : ||\boldsymbol{\beta}^*||_0 \le s, \sup_{k \in \mathcal{A}_m} ||\boldsymbol{\omega}^{(k)} - \boldsymbol{\beta}^*||_1 \le m \right\}.$$

#### 3.1 Estimation with $\ell_1$ -norm constrained transferring set

**Proposition 3.1.** (Local Restricted Strong Convexity) Assume Assumptions 3.1 - 3.3 hold. Let  $\Delta = \omega - \omega^*$ ,  $n = n_{\mathcal{A}_m} + n_0$  and  $\kappa_l = \min_{|u| < 1} K(u) > 0$ . If

 $\max\{r/(4\eta_{1/4}), 16m\zeta_p\} \le h \le f_l/l_0 \text{ and } nh \gtrsim f_u f_l^{-2} \log p \max\{l^2/(\eta_{1/4}^2), \eta_{1/4}^4/\zeta_p^2\},$ 

then for any  $\boldsymbol{\omega} \in \boldsymbol{\omega}^* + \mathbb{B}_{\Sigma}(r) \cap \mathbb{C}_{\Sigma}(l)$ ,

$$\hat{Q}_{h}(\boldsymbol{\omega}) - \hat{Q}_{h}(\boldsymbol{\omega}^{*}) - \langle \nabla \hat{Q}_{h}(\boldsymbol{\omega}^{*}), \boldsymbol{\omega} - \boldsymbol{\omega}^{*} \rangle \ge \phi_{1} ||\Delta||_{\Sigma}^{2} - \phi_{2} \sqrt{\frac{\log p + \log n}{nh}} ||\Delta||_{1} ||\Delta||_{\Sigma},$$
(4)

with probability at least  $1 - (pn)^{-1}$ , where  $\phi_1 = \frac{\kappa_l f_l}{10}$ ,  $\phi_2 = C' \kappa_l$ .

**Proposition 3.2.** Assume Assumptions 3.1 - 3.3 hold. Let  $\boldsymbol{v} = \beta_1 - \beta_2$ ,  $\alpha_1 = \kappa_l f_l/10$  and  $\alpha_2 = C^2 \kappa_l^2/(2\alpha_1)$ . If  $4r\eta_{1/4} \leq g \leq f_u/l_0$  with  $\eta_{1/4}$  defined in (3) and  $n_0g \gtrsim f_u f_l^{-2} \max\{s, l^2 \log p\}$ , then for any  $\boldsymbol{v} \in \mathbb{B}_{\Sigma}(r) \cap \mathbb{C}_{\Sigma}(l)$ ,

$$\hat{Q}_{g}^{(0)}(\boldsymbol{\beta}_{1}) - \hat{Q}_{g}^{(0)}(\boldsymbol{\beta}_{2}) - \langle \nabla \hat{Q}_{g}^{(0)}(\boldsymbol{\beta}_{2}), \boldsymbol{\beta}_{1} - \boldsymbol{\beta}_{2} \rangle \geq \alpha_{1} ||\boldsymbol{v}||_{\Sigma}^{2} - C\kappa_{l} \sqrt{\frac{\log p + \log n_{0}}{n_{0}g}} ||\boldsymbol{v}||_{1} ||\boldsymbol{v}||_{\Sigma}$$

with probability at least  $1 - (pn_0)^{-1}$ . By the arithmetic mean-geometric mean

By the arithmetic mean-geometric mean inequality

$$C\kappa_l \sqrt{\frac{\log p + \log n_0}{n_0 g}} ||\boldsymbol{v}||_1 ||\boldsymbol{v}||_{\Sigma} \le \frac{\alpha_1}{2} ||\boldsymbol{v}||_{\Sigma}^2 + \frac{C^2 \kappa_l^2}{2\alpha_1} \frac{\log p + \log n_0}{n_0 g} ||\boldsymbol{v}||_1^2,$$

we have

$$\hat{Q}_{g}^{(0)}(\boldsymbol{\beta}_{1}) - \hat{Q}_{g}^{(0)}(\boldsymbol{\beta}_{2}) - \langle \nabla \hat{Q}_{g}^{(0)}(\boldsymbol{\beta}_{2}), \boldsymbol{\beta}_{1} - \boldsymbol{\beta}_{2} \rangle \geq \frac{\alpha_{1}}{2} ||\boldsymbol{v}||_{\Sigma}^{2} - \alpha_{2} \frac{\log p + \log n_{0}}{n_{0}g} ||\boldsymbol{v}||_{1}^{2},$$

with probability at least  $1 - (pn_0)^{-1}$ .

In the debiasing step, we need another restricted strong convexity condition with both  $|| \cdot ||_1$  and  $|| \cdot ||_{\Sigma}$  in the lower bound. Proposition 3.2 provides that kind condition.

Finally, with the above establishments of restricted strong convexity, we are able to obtain the main result for the two-step transfer learning algorithm on quantile regression.

**Theorem 3.1.** Assume Assumptions 3.1 - 3.4 hold. Suppose  $n_0 \ge Cs^2 \log p$  and  $n_{\mathcal{A}_m} \gtrsim n_0$ , where C > 0 is a constant. Also let

$$\log(p)/(n_{\mathcal{A}_m} + n_0) \lesssim h \leq \min\{f_l/(2l_0\kappa_1), (s^{1/2}\lambda_{\omega})^{1/2}\} \\ s \log(p)/n_0 \lesssim g \leq (\log(p)/n_0)^{1/4}.$$

We take  $\lambda_{\omega} = C_{\omega} \sqrt{\log(p)/(n_{\mathcal{A}_m} + n_0)}$ ,  $\lambda_{\delta} = C_{\delta} \sqrt{\log(p)/n_0}$ , where  $C_{\omega}$  and  $C_{\delta}$  are sufficiently large constants, then

$$||\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*||_{\Sigma} \lesssim \sqrt{m} \left(\frac{\log p}{n_0}\right)^{1/4} + \sqrt{s} \left(\frac{\log p}{n_0}\right)^{1/4} \left(\frac{\log p}{n_{\mathcal{A}_m} + n_0}\right)^{1/4},\tag{5}$$

$$||\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*||_1 \lesssim s \sqrt{\frac{\log p}{n_{\mathcal{A}_m} + n_0}} + \left(\frac{\log p}{n_{\mathcal{A}_m} + n_0}\right)^{\frac{1}{4}} \sqrt{sm} + m,\tag{6}$$

with probability at least  $1 - p^{-1}$ .

**Remark 3.1.** In the trivial case where  $\mathcal{A}_m$  is an empty set, the upper bound in (5) is  $\mathcal{O}_P(\sqrt{s\log(p)/n_0})$ . When  $\mathcal{A}_m$  is non-empty, the upper bound in (5) is sharper than  $\sqrt{s\log(p)/n_0}$  and the upper bound in (6) is sharper than  $s\sqrt{\log(p)/n_0}$ , if  $n_{\mathcal{A}_m} \gtrsim n_0$  and  $m < s(\log(p)/n_0)^{1/2}$ .

The above theorem gives the convergence rate of the Trans-SQR estimator under  $\ell_1/\ell_2$ -errors. As the above remarks stated, if the total sample size of the transferable sources is significantly larger than the target sample size, the Trans-SQR estimator could even achieve a sharper convergence rate with some proper choices of the transferable level of the contrasts and the smoothing bandwidth in the debiasing step. As some previous works show, our theorem shares similar estimation error bounds as the results in Tian & Feng (2022) and Li et al. (2022b).

# 3.2 Estimation with $\ell_0$ -norm constrained transferring set

**Proposition 3.3.** (RSC in Step 2) Assume Assumptions 3.1 - 3.3 hold. Let  $\boldsymbol{v} = \beta_1 - \beta_2$ . If  $4r\eta_{1/4} \leq g \leq f_u/l_0$  with  $\eta_{1/4}$  defined in (3) and  $n_0g \gtrsim f_u f_l^{-2} \max\{s, l^2 \log p\}$ , then for any  $\boldsymbol{v} \in \mathbb{B}_{\Sigma}(r) \cap \mathbb{C}_{\Sigma}(l)$ ,

$$\hat{Q}_{g}^{(0)}(m{eta}_{1}) - \hat{Q}_{g}^{(0)}(m{eta}_{2}) - \langle 
abla \hat{Q}_{g}^{(0)}(m{eta}_{2}), m{eta}_{1} - m{eta}_{2} 
angle \geq 0.2 f_{l} \kappa_{l} ||m{v}||_{\Sigma}^{2},$$

with probability at least  $1 - (pn_0)^{-1}$ .

Theorem 3.2. Assume Assumptions 3.1 - 3.3 hold. Let

$$\log(p)/n_0 \lesssim h \le \min\{f_l/(2l_0\kappa_1), (s^{1/2}\lambda_{\omega})^{1/2}\}$$

$$(s+m)\log(p)/n_0 \lesssim g \le ((s+m)\log(p)/n_0)^{1/4}$$

$$m\log(p)/n \lesssim r \le (m\log(p)/n)^{1/4},$$

where  $n = n_0 + n_{\mathcal{A}'_m}$ . Meanwhile, suppose  $m \leq s, n_k \geq n_0$  and  $n_0 \geq Cs^2 \log p$ , where C > 0 is a constant. We take  $\lambda_{\omega}^{(k)} = C_{\omega} \sqrt{\log(p)/n_k}$ ,  $\lambda_{\delta} = C_{\delta} \sqrt{\log(p)/n_0}$  and  $\lambda_{\beta} = C_{\beta} \sqrt{\log(p)/n}$ , where  $C_{\omega}$ ,  $C_{\delta}$  and  $C_{\beta}$  are sufficiently large constants, then

$$||\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*||_{\Sigma} \lesssim \sqrt{\frac{s\log p}{n}} + \sqrt{\frac{sm\log p}{n_0}},\tag{7}$$

$$||\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*||_1 \lesssim s \sqrt{\frac{\log p}{n}} + s \sqrt{\frac{m \log p}{n_0}},\tag{8}$$

with probability at least  $1 - p^{-1}$ .

The above theorem gives the convergence rate of the Trans-SQR estimator under  $\ell_1/\ell_2$ -errors, where the contrast vectors are characterized in terms of the  $\ell_0$ -norm. If the sample size of the target data is large enough and the total sample size of the transferable sources is significantly larger than the target sample size, the Trans-SQR estimator could achieve a sharp convergence rate with some proper choices of the transferable level of the contrasts.

**Remark 3.2.** As mentioned above, Assumption 3.4 is to make sure that the estimation error in the transferring step is small enough. However, Theorem 3.2 does not require Assumption 3.4 because Algorithm 4 learns the parameter  $\boldsymbol{w}^{(k)}$  independently in Step 1 and reduces the bias in Step 2. For Step 1, the upper bound of the difference between the estimator  $\hat{\boldsymbol{w}}^{(k)}$  and true parameter  $\boldsymbol{w}^{(k)}$  can be controlled by the sample size of the each source data and the  $\ell_0$  transferable level m. For Step 2, the estimated  $\hat{\boldsymbol{\delta}}$  could also be closed enough to the true difference between the target and source parameter by having an appropriate target sample size. Therefore, if both the target and source sample sizes are large enough, the error of Algorithm 4 would be well controlled without Assumption 3.4.

# 4 Numerical Studies

In this section, we evaluate the performance of our proposed algorithms via numerical experiments. The methods in the following section include Smoothed Quantile Regression (SQR) on target data, the Oracle-Trans-SQR,  $\mathcal{A}_m$ -Trans-SQR and the Naive-Trans-SQR, which naïvely assumes  $\mathcal{A}_m = 1, \ldots, K$  in the Oracle

Trans-SQR. The purpose of including the Naive-Trans-Lasso is to understand the overall informative level of the auxiliary samples.

# 4.1 Transfer learning on $\ell_1$ -normed $\mathcal{A}_m$

We consider p = 500,  $n_0 = 200$ , and  $n_1, \ldots, n_K = 150$  for K = 10. The covariates from target  $\boldsymbol{x}_i^{(0)}$  are i.i.d. Gaussian with mean zero and covariance matrix  $\boldsymbol{\Sigma}$  with  $\boldsymbol{\Sigma}_{jj'} = 0.5^{|j-j'|}$  for all  $i = 1, \ldots, n_0$  and  $\epsilon_i^{(0)}$ are i.i.d. Gaussian with mean zero and variance one for all *i*. For  $k \in \mathcal{A}_m$ ,  $\boldsymbol{x}_i^{(k)} \sim \mathcal{N}(\mathbf{0}_p, \boldsymbol{\Sigma} + \boldsymbol{\epsilon}\boldsymbol{\epsilon}^{\mathrm{T}})$ , where  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}_p, 0.3^2 \boldsymbol{I}_p)$ . For the target, the true parameter  $\boldsymbol{\beta}^*$ , we set s = 5,  $\beta_j = 0.5$  for  $j \in \{1, \ldots, s\}$ , and  $\beta_j = 0$  otherwise. Denote  $\mathcal{R}_p^{(k)}$  as *p* independent Rademacher variables.  $\mathcal{R}_p^{(k)}$  is independent with  $\mathcal{R}_p^{(k')}$ for any  $k \neq k'$ . For any source data *k* in  $\mathcal{A}_m$ , we let the true parameter  $\boldsymbol{\omega}^{(k)} = \boldsymbol{\beta}^* + (m/p)\mathcal{R}_p^{(k)}$ , where  $m \in \{5, 10\}$ . For any source data *k'* not in  $\mathcal{A}_m$ , the true parameter  $\boldsymbol{\omega}^{(k')} = \boldsymbol{\beta}^* + (2m/p)\mathcal{R}_p^{(k')}$ . We train the four methods with 100 reproductions and record their average  $\ell_2$ -estimation errors under different settings of  $\tau$ . Figure 1 shows the changes of the estimation errors along with the amount of the transferable sources.

We observe from Figure 1 that the Oracle-Trans-SQR has the best performance among all the methods and  $\mathcal{A}_m$ -Trans-SQR has almost the same performance as the Oracle-Trans-SQR, which indicates that the transferable source detection algorithm still works under the smoothed quantile regression models. Meanwhile, compared with SQR on target, the estimation errors of the Oracle-Trans-SQR and  $\mathcal{A}_m$ -Trans-SQR are always smaller, which means that the source data which share some similarities in  $\ell_1$ -norm with the target data could improve the estimation. Another observation is that the performance of  $\mathcal{A}_m$ -Trans-SQR consistently improves as more and more source data are transferable. This matches the theoretical  $\ell_2$ -estimation error bounds which become sharper as  $n_{\mathcal{A}_m}$  grows.

## 4.2 Transfer learning on $\ell_0$ -normed $\mathcal{A}_m$

We consider p = 500,  $n_0 = 200$ , and assume that there are 2, 4, 6, 8, 10 transferable sources with the sample sizes 400. The covariates from target  $\boldsymbol{x}_i^{(0)}$  are i.i.d. Gaussian with mean zero and covariance matrix  $\boldsymbol{\Sigma}$  with  $\Sigma_{jj'} = 0.5^{|j-j'|}$ . The covariates from source  $\boldsymbol{x}_i^{(k)}$  are also i.i.d. Gaussian with mean zero, but with covariance matrix  $\boldsymbol{\Sigma} + \boldsymbol{\epsilon}\boldsymbol{\epsilon}^{\mathrm{T}}$ , where  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}_p, 0.3^2 \boldsymbol{I}_p)$ . For the target, the true parameter  $\boldsymbol{\beta}^*$ , we set s = 5,  $\beta_j = 1$  for  $j \in \{1, \ldots, s\}$ , and  $\beta_j = 0$  otherwise. For the source, their true parameter  $\boldsymbol{w}^{(k)}$  is generated from  $w_j^{(k)} = \beta_j^* + \Delta \mathbb{1}(j \in M)$  where M is a random subset of [p] with |M| = m. We take  $m \in \{2, 4\}$ , and  $\Delta = 2$ . Figure 3 and 4 show the  $\ell_2$ -estimation errors in different settings of m.

From the results, Trans-SQR with  $\ell_0$ -norm constrained transferring set has better performances than SQR only on target and SQR on all sources and target. Meanwhile, when the target data sample size  $n_0$  becomes larger, the performance of Trans-SQR increases quickly, which accords with our results that the estimation error is depend on the target sample size. There are considerable decreases in estimation errors of Trans-SQR when the transferable level increases or  $\Delta$  increases, which corresponds to the difference on components between target and source populations.

# 5 Conclusion

This paper studies transfer learning for high-dimensional quantile regression models, employing convolutiontype smoothing techniques. The proposed algorithms focus on leveraging  $\ell_1/\ell_0$ -normed transferable source populations to improve estimation accuracy of the target regression coefficients. We derive error bounds for the estimators in terms of  $\ell_1/\ell_2$ -norms for the algorithms. Theoretical analysis reveals that these error bounds surpass those of the classical penalized quantile regression estimator, which only utilizes the target data, provided that the target and source populations exhibit sufficient similarity. Furthermore, we propose a transferable source detection algorithm to identify informative sources from the available sources when the set of informative sources is unknown. Numerical experiments validate our theoretical results.



Figure 1:  $\ell_2$  estimation errors of several methods under quantile levels  $\tau = 0.25, 0.5, 0.75$ , over 100 repetitions, where Oracle-Trans-SQR is Algorithm 1



Figure 2:  $\ell_2$  estimation errors of several methods for Gaussian and  $t_{1.5}$  errors, over 100 repetitions.



Figure 3:  $\ell_2$  estimation errors of several methods with  $\ell_0$  constraints for  $t_{1.5}$  errors, over 100 repetitions.



Figure 4:  $\ell_2$  estimation errors of several methods with  $\ell_0$  constraints for Gaussian errors, over 100 repetitions.

# References

- Hamsa Bastani. Predicting with proxies: Transfer learning in high dimension. *Management Science*, 67(5): 2964–2984, 2021.
- Alexandre Belloni and Victor Chernozhukov. l<sub>1</sub>-penalized quantile regression in high-dimensional sparse models. The Annals of Statistics, 39(1):82–130, 2011.
- Olivier Bousquet. Concentration inequalities for sub-additive functions using the entropy method. In Stochastic inequalities and applications, pp. 213–247. Springer, 2003.
- T Tony Cai and Hongji Wei. Transfer learning for nonparametric classification: Minimax rate and adaptive classifier. *The Annals of Statistics*, 49(1):100–128, 2021.
- Xuming He, Xiaoou Pan, Kean Ming Tan, and Wen-Xin Zhou. Smoothed quantile regression with large-scale inference. *Journal of Econometrics*, 2021.
- Yong He, Qiushi Li, Qinqin Hu, and Lei Liu. Transfer learning in high-dimensional semiparametric graphical models with application to brain connectivity analysis. *Statistics in medicine*, 41(21):4112–4129, 2022.
- Joel L Horowitz. Bootstrap methods for median regression models. *Econometrica*, pp. 1327–1351, 1998.
- Jiayu Huang, Mingqiu Wang, and Yuanshan Wu. Transfer learning with high-dimensional quantile regression. arXiv preprint arXiv: 2211.14578, 2022.
- Michel Ledoux and Michel Talagrand. Probability in Banach Spaces: isoperimetry and processes, volume 23. Springer Science & Business Media, 1991.
- Sai Li, Tianxi Cai, and Rui Duan. Targeting underrepresented populations in precision medicine a federated transfer learning approach. arXiv preprint arXiv:2108.12112, 2021.
- Sai Li, T Tony Cai, and Hongzhe Li. Transfer learning in large-scale gaussian graphical models with false discovery rate control. Journal of the American Statistical Association, (just-accepted):1–13, 2022a.
- Sai Li, T Tony Cai, and Hongzhe Li. Transfer learning for high-dimensional linear regression: Prediction, estimation and minimax optimality. Journal of the Royal Statistical Society. Series B, Statistical Methodology, 84(1):149–173, 2022b.
- Lu Lin and Weiyu Li. A correlation-ratio transfer learning and variational stein's paradox. arXiv preprint arXiv:2206.06086, 2022.
- Po-Ling Loh and Martin J Wainwright. Regularized m-estimators with nonconvexity: Statistical and algorithmic theory for local optima. *Journal of Machine Learning Research*, 16:559–616, 2015.
- Henry WJ Reeve, Timothy I Cannings, and Richard J Samworth. Adaptive transfer learning. The Annals of Statistics, 49(6):3618–3649, 2021.
- Kean Ming Tan, Lan Wang, and Wen-Xin Zhou. High-dimensional quantile regression: Convolution smoothing and concave regularization. *Journal of the Royal Statistical Society Series B*, 84(1):205–233, 2022.
- Ye Tian and Yang Feng. Transfer learning under high-dimensional generalized linear models. *Journal of the American Statistical Association*, (just-accepted):1–30, 2022.
- Lisa Torrey and Jude Shavlik. In handbook of research on machine learning applications and trends: algorithms, methods, and techniques: algorithms, methods, and techniques. IGI global, pp. 242–264, 2010.
- Roman Vershynin. High-dimensional probability: An introduction with applications in data science, volume 47. Cambridge university press, 2018.
- Xiao-Tong Yuan, Ping Li, and Tong Zhang. Gradient hard thresholding pursuit. Journal of Machine Learning Research, 18(166):1–43, 2018.
- Yijiao Zhang and Zhongyi Zhu. Transfer learning for high-dimensional quantile regression via convolution smoothing. arXiv preprint arXiv:2212.00428, 2022.

# A Appendix: Proofs of the main results

#### A.1 Technical Lemmas

For  $\boldsymbol{\omega} \in \mathbb{R}^p$ , suppose  $\Delta = \boldsymbol{\omega} - \boldsymbol{\omega}^*$ . Define

$$\hat{R}_{h}(\Delta) = \hat{Q}_{h}(\boldsymbol{\omega}) - \hat{Q}_{h}(\boldsymbol{\omega}^{*}) - \langle \nabla \hat{Q}_{h}(\boldsymbol{\omega}^{*}), \boldsymbol{\omega} - \boldsymbol{\omega}^{*} \rangle, \\ \hat{D}_{h}(\Delta) = \hat{Q}_{h}(\boldsymbol{\omega}) - \hat{Q}_{h}(\boldsymbol{\omega}^{*}),$$

and their population counterparts  $R_h(\Delta) = \mathbb{E}\{\hat{R}_h(\Delta)\}$  and  $D_h(\Delta) = \mathbb{E}\{\hat{D}_h(\Delta)\}$ , where  $\boldsymbol{\omega}^*$  is the true parameter of the transferring step in the algorithm.

Lemma A.1. Let  $\beta^*$  be the true target parameter, then  $||\omega^* - \beta^*||_1 \leq C_1 m$ , where  $C_1 = \sup_k ||\tilde{\Sigma}^{-1}\tilde{\Sigma}^{(k)}||_1$ and  $\tilde{\Sigma}^{-1}, \tilde{\Sigma}^{(k)}$  are given in Assumption 3.4.

Note that  $\boldsymbol{w}^*$  has the explicit form,  $\boldsymbol{w}^* = \boldsymbol{\beta}^* + \boldsymbol{\delta}^*$ . Lemma A.1 gives an upper bound of the distance between the true  $\boldsymbol{\beta}^*$  and the true estimate in transferring step. In other words, the  $\ell_1$ -norm of  $\boldsymbol{\delta}^*$  is controlled by  $\boldsymbol{m}$ . Lemma A.2. Define  $\boldsymbol{\pi}_h^* = \boldsymbol{\pi}_h(\boldsymbol{\beta}^*) \in \mathbb{R}^p$ , where  $\boldsymbol{\pi}_h(\boldsymbol{\beta}) = \nabla \hat{Q}_h(\boldsymbol{\beta}) - \nabla Q_h(\boldsymbol{\beta})$ . Assumptions 4.1 - 4.3 ensure that for any t > 0,

$$||\boldsymbol{\pi}_{h}^{*}||_{\infty} \leq \sigma \sqrt{\{\tau(1-\tau) + Ch^{2}\}\frac{2t}{n_{\mathcal{A}_{m}} + n_{0}}} + \max(1-\tau,\tau)\frac{t}{n_{\mathcal{A}_{m}} + n_{0}},$$

with probability at least  $1 - 2pe^{-t}$ , where  $C = (\tau + 1)l_0\kappa_2$  and  $\sigma = \max_{1 \le j \le p} \sigma_{jj}$ .

In both transferring and debiased steps, we need to restrict the regularization parameters  $\lambda_{\omega}$  (or  $\lambda_{\delta}$ ) to be no smaller than  $2||\boldsymbol{\pi}_{h}^{*}||_{\infty}$  (or  $2||\boldsymbol{\pi}_{g}^{*}||_{\infty}$ ). This Lemma helps to specify the choice of the parameters.

**Lemma A.3.** Define  $b_h^* = ||\Sigma^{-1/2} \nabla Q_h(\boldsymbol{\omega}^*)||_2$ , which quantifies the smoothing bias, then for some  $\kappa_2 > 0$ 

$$b_h^* \le l_0 \kappa_2 \frac{h^2}{2},$$

where  $l_0$  is the Lipschitz constant of the density  $f_{\epsilon | \boldsymbol{x}}(\cdot)$ .

Lemma A.4. For r, l > 0, define

$$\psi(r,l) = \sup_{\boldsymbol{\beta} \in \boldsymbol{\beta}^* + \mathbb{B}_{\Sigma}(r) \cap \mathbb{C}_1(l)} \left\| \frac{1}{n} \sum_{i=1}^n (1-\mathbb{E}) \left\{ \bar{K}_h(\boldsymbol{x}_i^{\mathrm{T}} \boldsymbol{\beta} - y_i) - \bar{K}_h(\boldsymbol{x}_i^{\mathrm{T}} \boldsymbol{\beta}^* - y_i) \right\} \boldsymbol{x}_i \right\|_{\infty}.$$

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

$$\psi(r,l) \lesssim \frac{l}{h} \sqrt{\frac{\log p}{n}} + f_u^{1/2} r \sqrt{\frac{t+\log p}{nh}} + \frac{t+\log p}{n}.$$

#### A.2 Proof of Proposition 3.1

Define the Taylor error

$$\mathcal{T}(\boldsymbol{\omega},\boldsymbol{\omega}^*) = \hat{Q}_h(\boldsymbol{\omega}) - \hat{Q}_h(\boldsymbol{\omega}^*) - \langle \nabla \hat{Q}_h(\boldsymbol{\omega}^*), \boldsymbol{\omega} - \boldsymbol{\omega}^* \rangle.$$

If

$$\frac{||\Delta||_1}{||\Delta||_{\Sigma}} > \frac{\phi_1}{\phi_2} \sqrt{\frac{nh}{\log p + \log n}},$$

then the lower bound in (4) is negative. Since  $\hat{Q}_h(\cdot)$  is convex, we have  $\mathcal{T}(\boldsymbol{\omega}, \boldsymbol{\omega}^*) \geq 0$ . Thus (4) holds trivially in that case. In the following proofs,  $\boldsymbol{\omega} \in \boldsymbol{\omega}^* + \mathbb{B}_{\Sigma}(r) \cap \mathbb{C}_{\Sigma}(l)$ , where  $l = \phi_1 \sqrt{nh/(\log p + \log n)}/\phi_2$ .

It follows from a second-order Taylor expansion that

$$\begin{aligned} \mathcal{T}(\boldsymbol{\omega},\boldsymbol{\omega}^*) \\ &= \frac{1}{2} (\boldsymbol{\omega} - \boldsymbol{\omega}^*)^{\mathrm{T}} \nabla^2 \hat{Q}_h \big( t \boldsymbol{\omega} + (1-t) \boldsymbol{\omega}^* \big) (\boldsymbol{\omega} - \boldsymbol{\omega}^*) \\ &= \frac{1}{2(n_{\mathcal{A}_m} + n_0)} \sum_{k \in \mathcal{A}_m \cup \{0\}} \sum_{i=1}^{n_k} K_h \big\{ y_i^{(k)} - \langle \boldsymbol{x}_i^{(k)}, t_i^{(k)} \boldsymbol{\omega} + (1-t_i^{(k)}) \boldsymbol{\omega}^* \rangle \big\} \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega} - \boldsymbol{\omega}^* \rangle^2 \\ &= \frac{1}{2(n_{\mathcal{A}_m} + n_0)} \sum_{k \in \mathcal{A}_m \cup \{0\}} \sum_{i=1}^{n_k} K_h \big\{ \epsilon_i - t_i^{(k)} \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega} - \boldsymbol{\omega}^* \rangle - \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega}^* - \boldsymbol{\omega}^{(k)} \rangle \big\} \\ &\cdot \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega} - \boldsymbol{\omega}^* \rangle^2, \end{aligned}$$

for some  $t_i^{(k)} \in [0, 1]$ . For each *i* and *k*, define the event  $\mathcal{F}_{i,k}$ ,

$$\mathcal{F}_{i,k} = \{ |\epsilon_i| \le h/4 \} \cap \{ |\langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega} - \boldsymbol{\omega}^* \rangle| \le ||\boldsymbol{\omega} - \boldsymbol{\omega}^*||_{\Sigma} \cdot h/(2r) \} \cap \{ |\langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega}^* - \boldsymbol{\omega}^{(k)} \rangle| \le h/4 \},$$

for all  $\boldsymbol{\omega} - \boldsymbol{\omega}^* \in \mathbb{B}_{\Sigma}(r)$ . Thus

$$\mathcal{T}(\boldsymbol{\omega},\boldsymbol{\omega}^*) \geq \frac{\kappa_l}{2(n_{\mathcal{A}_m} + n_0)h} \sum_{k \in \mathcal{A}_m \cup \{0\}} \sum_{i=1}^{n_k} \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega} - \boldsymbol{\omega}^* \rangle^2 \mathbb{1}_{\mathcal{F}_{i,k}},$$
(9)

where  $\kappa_l = \min_{|u| \leq 1} K(u)$ . For a truncation level R > 0, define functions

$$\varphi_R(u) = \begin{cases} u^2 & |u| \le \frac{R}{2}, \\ (R - |u|)^2 & \frac{R}{2} < |u| \le R, \\ 0 & |u| > R. \end{cases}$$

By this construction,  $\varphi_R(u) \leq u^2 \cdot \mathbb{1}\{|u| \leq R\}$ ,  $\varphi_{cR}(cu) = c^2 \varphi_R(u)$  and  $\varphi_R$  is R-Lipschitz. In addition, we define the trapezoidal function

$$\psi_R(u) = \begin{cases} 1 & |u| \le \frac{R}{2}, \\ 2 - \frac{2}{R}|u| & \frac{R}{2} < |u| \le R, \\ 0 & |u| > R, \end{cases}$$

and note that  $\psi_R$  is (2/R)-Lipschitz and  $\psi_R(u) \leq \mathbb{1}\{|u| \leq R\}$ .

With these two new-defined function and the notation  $\Delta = \omega - \omega^*$ ,  $n = n_{\mathcal{A}_m} + n_0$ , we have established the lower bound of (9)

$$\mathcal{T}(\boldsymbol{\omega},\boldsymbol{\omega}^{*}) \geq \frac{\kappa_{l}}{2nh} ||\Delta||_{\Sigma}^{2} \sum_{k \in \mathcal{A}_{m} \cup \{0\}} \sum_{i=1}^{n_{k}} \mathbb{1}\{|\epsilon_{i}| \leq h/4\} \varphi_{||\Delta||_{\Sigma} \cdot h/(2r)} \left(\langle \boldsymbol{x}_{i}^{(k)}, \Delta \rangle\right) \psi_{h/4} \left(\langle \boldsymbol{x}_{i}^{(k)}, \boldsymbol{\omega}^{*} - \boldsymbol{\omega}^{(k)} \rangle\right) \\ \geq \frac{\kappa_{l}}{2} ||\Delta||_{\Sigma}^{2} \cdot \underbrace{\frac{1}{nh} \sum_{k,i} \mathbb{1}\{|\epsilon_{i}| \leq h/4\} \varphi_{h/(2r)} \left(\langle \boldsymbol{x}_{i}^{(k)}, \Delta \rangle/||\Delta||_{\Sigma}\right) \psi_{h/4} \left(\langle \boldsymbol{x}_{i}^{(k)}, \boldsymbol{\omega}^{*} - \boldsymbol{\omega}^{(k)} \rangle\right)}_{D_{0}(\boldsymbol{\omega}, \boldsymbol{\omega}^{*})} \tag{10}$$

In the following proofs, we bound  $\mathbb{E}D_0(\omega, \omega^*)$  and  $D_0(\omega, \omega^*) - \mathbb{E}D_0(\omega, \omega^*)$ , respectively. First, we show that

$$\mathbb{E}D_0(\boldsymbol{\omega}, \boldsymbol{\omega}^*) \ge 0.218 f_l. \tag{11}$$

Note that

$$\begin{aligned} \left|\frac{h}{2}f_{\epsilon|\boldsymbol{x}}(0)\right| - \left|\mathbb{E}\left[\mathbbm{1}\left\{|\epsilon_{i}| \leq h/4\right\}|\boldsymbol{x}_{i}^{(k)}\right]\right| &\leq \left|\mathbb{E}\left[\mathbbm{1}\left\{|\epsilon_{i}| \leq h/4\right\}|\boldsymbol{x}_{i}^{(k)}\right] - \frac{h}{2}f_{\epsilon|\boldsymbol{x}}(0)\right] \\ &\leq \int_{-h/4}^{h/4} |f_{\epsilon|\boldsymbol{x}}(t) - f_{\epsilon|\boldsymbol{x}}(0)|dt \\ &\leq \frac{l_{0}h^{2}}{16}. \end{aligned}$$

Hence we obtain

$$\left| \mathbb{E}[\mathbb{1}\{|\epsilon_i| \le h/4\} | \boldsymbol{x}_i^{(k)}] \right| \ge \frac{h}{2} f_l - \frac{l_0 h^2}{16}.$$

Provided  $h \leq f_l/l_0 \leq f_u/l_0$ , we have

$$\left| \mathbb{E}[\mathbb{1}\{|\epsilon_i| \le h/4\} | \boldsymbol{x}_i^{(k)}] \right| \ge \frac{7f_l h}{16}.$$

Meanwhile

$$\left| \mathbb{E}[\mathbb{1}\{|\epsilon_i| \le h/4\} | \boldsymbol{x}_i^{(k)}] \right| - \left| \frac{h}{2} f_{\epsilon|\boldsymbol{x}}(0) \right| \le \int_{-h/4}^{h/4} |f_{\epsilon|\boldsymbol{x}}(t) - f_{\epsilon|\boldsymbol{x}}(0)| dt$$

implies

$$\left|\mathbb{E}[\mathbb{1}\{|\epsilon_i| \le h/4\} | \boldsymbol{x}_i^{(k)}]\right| \le \frac{9f_u h}{16}.$$

Then

$$\begin{split} \mathbb{E}D_{0}(\boldsymbol{\omega},\boldsymbol{\omega}^{*}) \\ &= \frac{1}{nh}\sum_{k,i}\mathbb{E}\Big[\mathbb{E}\big[\mathbbm{1}\{|\epsilon_{i}| \leq h/4\}\big|\boldsymbol{x}_{i}^{(k)}\big]\varphi_{h/(2r)}(\langle\boldsymbol{x}_{i}^{(k)},\Delta\rangle/||\Delta||_{\Sigma})\psi_{h/4}(\langle\boldsymbol{x}_{i}^{(k)},\boldsymbol{\omega}^{*}-\boldsymbol{\omega}^{(k)}\rangle)\Big] \\ &\geq \frac{7f_{l}}{16}\mathbb{E}\Big[\varphi_{h/(2r)}(\langle\boldsymbol{x},\Delta\rangle/||\Delta||_{\Sigma})\psi_{h/4}(\langle\boldsymbol{x},\boldsymbol{\omega}^{*}-\boldsymbol{\omega}^{(k)}\rangle)\Big] \\ &\geq \frac{7f_{l}}{16}\bigg(1-\mathbb{E}\big[\langle\boldsymbol{x},\Delta/||\Delta||_{\Sigma}\rangle^{2}\mathbbm{1}\{|\langle\boldsymbol{x},\Delta/||\Delta||_{\Sigma}\rangle|>h/(4r)\}\big] \\ &\quad -\mathbb{E}\Big[\langle\boldsymbol{x},\Delta/||\Delta||_{\Sigma}\rangle^{2}\mathbbm{1}\{|\langle\boldsymbol{x},\boldsymbol{\omega}^{*}-\boldsymbol{\omega}^{(k)}\rangle|>h/8\}\Big]\bigg). \end{split}$$

By the definition of  $\eta_{\delta}$ , as long as  $0 < r \le h/(4\eta_{1/4})$ ,

$$\sup_{\Delta \in \mathbb{B}_{\Sigma}(r)} \mathbb{E}\Big[ \langle \boldsymbol{x}, \Delta / ||\Delta||_{\Sigma} \rangle^{2} \mathbb{1}\{ |\langle \boldsymbol{x}, \Delta / ||\Delta||_{\Sigma} \rangle| > h/(4r) \} \Big] \leq \frac{1}{4}$$

Moreover,  $\boldsymbol{\omega}^* \in \boldsymbol{\omega}^{(k)} + \mathbb{B}_1(m)$ .

$$\mathbb{E}\Big[\langle \boldsymbol{x}, \Delta/||\Delta||_{\Sigma}\rangle^{2} \mathbb{1}\{|\langle \boldsymbol{x}, \boldsymbol{\omega}^{*} - \boldsymbol{\omega}^{(k)}\rangle| > h/8\}\Big]$$

$$\leq \left(\frac{8}{h}\right)^{2} \mathbb{E}\Big[\langle \boldsymbol{x}, \Delta/||\Delta||_{\Sigma}\rangle^{2} \langle \boldsymbol{x}, \boldsymbol{\omega}^{*} - \boldsymbol{\omega}^{(k)}\rangle^{2}\Big]$$

$$\leq \left(\frac{8}{h}\right)^{2} \Big(\mathbb{E}\langle \boldsymbol{x}, \Delta/||\Delta||_{\Sigma}\rangle^{4}\Big)^{1/2} \Big(\mathbb{E}\langle \boldsymbol{x}, \boldsymbol{\omega}^{*} - \boldsymbol{\omega}^{(k)}\rangle^{4}\Big)^{1/2}$$

$$\leq \left(\frac{8}{h}\right)^{2} \zeta_{p} \cdot \left[m^{4} \mathbb{E}\langle \boldsymbol{x}, (\boldsymbol{\omega}^{*} - \boldsymbol{\omega}^{(k)})/||\boldsymbol{\omega}^{*} - \boldsymbol{\omega}^{(k)}||_{\Sigma}\rangle^{4}\right]^{1/2}$$

$$\leq \left(\frac{8}{h}\right)^{2} \zeta_{p} \cdot m^{2} \zeta_{p} = \left(\frac{8m}{h}\right)^{2} \zeta_{p}^{2}.$$

The inequality comes from  $\mathbb{E}\langle \boldsymbol{x}, \boldsymbol{u}/||\boldsymbol{u}||_{\Sigma}\rangle^4 \leq \zeta_p^2$ .

Provided  $16m\zeta_p \leq h$ ,

$$\mathbb{E}D_0(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) > 0.218 f_l.$$

Next we find a lower bound of  $D_0(\boldsymbol{\omega}, \boldsymbol{\omega}^*) - \mathbb{E}D_0(\boldsymbol{\omega}, \boldsymbol{\omega}^*)$  over subset  $\Lambda(r, l) := \{(\boldsymbol{\omega}, \boldsymbol{\omega}^*) : \boldsymbol{\omega}^* \in \boldsymbol{\omega}^{(k)} + \mathbb{B}_1(m), \ \boldsymbol{\omega} \in \boldsymbol{\omega}^* + \mathbb{B}_{\Sigma}(r) \cap \mathbb{C}_{\Sigma}(l)\}$ . Define

$$\Omega(r,l) = \sup_{(\boldsymbol{\omega},\boldsymbol{\omega}^*) \in \Lambda(r,l)} \{-D_0(\boldsymbol{\omega},\boldsymbol{\omega}^*) + \mathbb{E}D_0(\boldsymbol{\omega},\boldsymbol{\omega}^*)\}.$$

Write  $D_0(\boldsymbol{\omega}, \boldsymbol{\omega}^*) = n^{-1} \sum_{k \in \mathcal{A}_m \cup \{0\}} \sum_{i=1}^{n_k} w_{i,k}(\boldsymbol{\omega}, \boldsymbol{\omega}^*)$ , where

$$w_{i,k}(\boldsymbol{\omega},\boldsymbol{\omega}^*) = h^{-1} \mathbb{1}\{|\epsilon_i| \le h/4\} \varphi_{h/(2r)}(\langle \boldsymbol{x}_i^{(k)}, \Delta \rangle / ||\Delta||_{\Sigma}) \psi_{h/4}(\langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega}^* - \boldsymbol{\omega}^{(k)} \rangle)$$

satisfies  $0 \le w_{i,k}(\boldsymbol{\omega}, \boldsymbol{\omega}^*) \le h/(4r)^2$ , since  $\varphi_R(u) \le (R/2)^2$  and  $\psi_R(u) \in [0, 1]$ . Moreover,

$$\begin{split} \mathbb{E}w_{i,k}^2(\boldsymbol{\omega}, \boldsymbol{\omega}^*) &= \mathbb{E}\Big[h^{-2}\mathbbm{1}\{|\epsilon_i| \le h/4\}\varphi_{h/(2r)}^2(\langle \boldsymbol{x}_i^{(k)}, \Delta \rangle / ||\Delta||_{\Sigma})\psi_{h/4}^2(\langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega}^* - \boldsymbol{\omega}^{(k)} \rangle)\Big] \\ &\le \frac{9f_u}{16h} \cdot \mathbb{E}\varphi_{h/(2r)}^2(\langle \boldsymbol{x}_i^{(0)}, \Delta \rangle / ||\Delta||_{\Sigma}) \\ &\le \frac{9f_u}{16h} \cdot \mathbb{E}(\langle \boldsymbol{x}_i^{(0)}, \Delta \rangle / ||\Delta||_{\Sigma})^4 = \frac{9f_u \zeta_p^2}{16h}. \end{split}$$

Using Bousquet's version of Talagrand's inequality yields that, for any z > 0,

$$\begin{aligned} \Omega(r,l) &\leq \mathbb{E}\Omega(r,l) + \{\mathbb{E}\Omega(r,l)\}^{1/2} \frac{1}{2r} \sqrt{\frac{hz}{n}} + \frac{3\sqrt{2}\zeta_p}{4} \sqrt{\frac{f_u z}{nh}} + \frac{h}{(4r)^2} \frac{z}{3n} \\ &\leq \mathbb{E}\Omega(r,l) + \frac{1}{4} \mathbb{E}\Omega(r,l) + \frac{1}{4r^2} \frac{hz}{n} + \frac{3\sqrt{2}\zeta_p}{4} \sqrt{\frac{f_u z}{nh}} + \frac{h}{(4r)^2} \frac{z}{3n} \\ &\leq \frac{5}{4} \mathbb{E}\Omega(r,l) + \frac{3\sqrt{2}\zeta_p}{4} \sqrt{\frac{f_u z}{nh}} + \frac{13}{3} \frac{hz}{(4r)^2 n} \end{aligned}$$
(12)

holds with probability at least  $1-e^{-z}$ . To bound the expectation  $\mathbb{E}\Omega(r, l)$ , using Rademacher symmetrization and the connection between Gaussian and Rademacher complexities, Lemma 5.5 in Ledoux & Talagrand (1991), we have

$$\mathbb{E}\Omega(r,l) \le 2\sqrt{\frac{\pi}{2}} \mathbb{E}\left[\sup_{(\boldsymbol{\omega},\boldsymbol{\omega}^*)\in\Lambda(r,l)} \mathbb{G}_{\boldsymbol{\omega},\boldsymbol{\omega}^*}\right],\tag{13}$$

where

 $\sim$ 

$$:= \frac{1}{nh} \sum_{k \in \mathcal{A}_m \cup \{0\}} \sum_{i=1}^{n_k} e_i \mathbb{1}\{|\epsilon_i| \le h/4\} \varphi_{h/(2r)}(\langle \boldsymbol{x}_i^{(k)}, \Delta/||\Delta||_{\Sigma}\rangle) \psi_{h/4}(\langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega}^* - \boldsymbol{\omega}^{(k)}\rangle),$$

and  $e_i$  are independent standard normal variables. Note that  $\mathbb{G}_{\boldsymbol{\omega},\boldsymbol{\omega}^*}$  is a Gaussian process conditioned on  $\{(y_i^{(k)}, \boldsymbol{x}_i^{(k)})\}_{i=1}^{n_k}$  for  $k \in \mathcal{A}_m \cup \{0\}$ . For  $(\boldsymbol{\omega}, \boldsymbol{\omega}^*)$  and  $(\boldsymbol{\omega}', \boldsymbol{\omega}'^*)$ , write  $\Delta = \boldsymbol{\omega} - \boldsymbol{\omega}^*$ ,  $\Delta' = \boldsymbol{\omega}' - \boldsymbol{\omega}'^*$  and

 $\chi_i = \mathbb{1}\{|\epsilon_i| \le h/4\}, \text{ then }$ 

$$\begin{split} & \mathbb{G}_{\boldsymbol{\omega},\boldsymbol{\omega}^{*}} - \mathbb{G}_{\boldsymbol{\omega}',\boldsymbol{\omega}'^{*}} \\ &= \mathbb{G}_{\boldsymbol{\omega},\boldsymbol{\omega}^{*}} - \mathbb{G}_{\boldsymbol{\omega}',\boldsymbol{\omega}'+\Delta} + \mathbb{G}_{\boldsymbol{\omega}',\boldsymbol{\omega}'+\Delta} - \mathbb{G}_{\boldsymbol{\omega}',\boldsymbol{\omega}'^{*}} \\ &= \frac{1}{nh} \sum_{k,i} e_{i} \chi_{i} \varphi_{h/(2r)} (\langle \boldsymbol{x}_{i}^{(k)}, \Delta/||\Delta||_{\Sigma} \rangle) \big\{ \psi_{h/4} (\langle \boldsymbol{x}_{i}^{(k)}, \boldsymbol{\omega}^{*} - \boldsymbol{\omega}^{(k)} \rangle) \\ &- \psi_{h/4} (\langle \boldsymbol{x}_{i}^{(k)}, \boldsymbol{\omega}'^{*} - \boldsymbol{\omega}^{(k)} \rangle) \big\} \\ &+ \frac{1}{nh} \sum_{k,i} e_{i} \chi_{i} \psi_{h/4} (\langle \boldsymbol{x}_{i}^{(k)}, \boldsymbol{\omega}'^{*} - \boldsymbol{\omega}^{(k)} \rangle) \big\{ \varphi_{h/(2r)} \big( (\boldsymbol{x}_{i}^{(k)})^{\mathrm{T}} \Delta/||\Delta||_{\Sigma} \big) \\ &- \varphi_{h/(2r)} \big( (\boldsymbol{x}_{i}^{(k)})^{\mathrm{T}} \Delta'/||\Delta'||_{\Sigma} \big) \big\}. \end{split}$$

Note that  $\varphi_R$  and  $\psi_R$  are Lipschitz continuous, and  $\varphi_R(u) \leq (R/2)^2$ . Let  $\mathbb{E}^*$  be the conditional expectation given  $\{(y_i^{(k)}, \boldsymbol{x}_i^{(k)})\}_{i=1}^{n_k}$ . Consequently,

$$\mathbb{E}^* \left( \mathbb{G}_{\boldsymbol{\omega},\boldsymbol{\omega}^*} - \mathbb{G}_{\boldsymbol{\omega}',\boldsymbol{\omega}'+\Delta} \right)^2 \leq \frac{1}{(nh)^2} \left( \frac{8}{h} \right)^2 \left( \frac{h}{4r} \right)^4 \sum_{k \in \mathcal{A}_m \cup \{0\}} \sum_{i=1}^{n_k} \chi_i \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega} - \boldsymbol{\omega}' \rangle^2$$
$$= \frac{1}{4r^4 n^2} \sum_{k \in \mathcal{A}_m \cup \{0\}} \sum_{i=1}^{n_k} \chi_i \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega} - \boldsymbol{\omega}' \rangle^2 \tag{14}$$

and

$$\mathbb{E}^{*} \left( \mathbb{G}_{\boldsymbol{\omega}',\boldsymbol{\omega}'+\Delta} - \mathbb{G}_{\boldsymbol{\omega}',\boldsymbol{\omega}'^{*}} \right)^{2} \tag{15}$$

$$\leq \frac{1}{(nh)^{2}} \left( \frac{h}{2r} \right)^{2} \sum_{k \in \mathcal{A}_{m} \cup \{0\}} \sum_{i=1}^{n_{k}} \chi_{i} \left( \langle \boldsymbol{x}_{i}^{(k)}, \Delta / ||\Delta||_{\Sigma} \rangle - \langle \boldsymbol{x}_{i}^{(k)}, \Delta / ||\Delta'||_{\Sigma} \rangle \right)^{2}$$

$$= \frac{1}{4r^{2}n^{2}} \sum_{k \in \mathcal{A}_{m} \cup \{0\}} \sum_{i=1}^{n_{k}} \chi_{i} \langle \boldsymbol{x}_{i}^{(k)}, \Delta / ||\Delta||_{\Sigma} - \Delta' / ||\Delta'||_{\Sigma} \rangle^{2}. \tag{16}$$

Motivated by the last two inequalities, we have

$$\mathbb{E}^* \left( \mathbb{G}_{\boldsymbol{\omega},\boldsymbol{\omega}^*} - \mathbb{G}_{\boldsymbol{\omega}',\boldsymbol{\omega}'^*} \right)^2 \leq \frac{1}{2r^4 n^2} \sum_{k \in \mathcal{A}_m \cup \{0\}} \sum_{i=1}^{n_k} \chi_i \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega} - \boldsymbol{\omega}' \rangle^2 \\ + \frac{1}{2r^2 n^2} \sum_{k \in \mathcal{A}_m \cup \{0\}} \sum_{i=1}^{n_k} \chi_i \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\Delta}/||\boldsymbol{\Delta}||_{\Sigma} - \boldsymbol{\Delta}'/||\boldsymbol{\Delta}'||_{\Sigma} \rangle^2.$$

Define another Gaussian process  $\mathbb{Z}_{\boldsymbol{\omega},\boldsymbol{\omega}^*}$  as

$$\mathbb{Z}_{\boldsymbol{\omega},\boldsymbol{\omega}^*} = \frac{1}{2^{1/2}r^2n} \sum_{k \in \mathcal{A}_m \cup \{0\}} \sum_{i=1}^{n_k} e'_i \chi_i \langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega}^* - \boldsymbol{\omega}^{(k)} \rangle \\ + \frac{1}{2^{1/2}rn} \sum_{k \in \mathcal{A}_m \cup \{0\}} \sum_{i=1}^{n_k} e''_i \chi_i \langle \boldsymbol{x}_i^{(k)}, \Delta \rangle / ||\Delta||_{\Sigma}$$

such that  $\mathbb{E}^*(\mathbb{G}_{\boldsymbol{\omega},\boldsymbol{\omega}^*} - \mathbb{G}_{\boldsymbol{\omega}',\boldsymbol{\omega}'^*})^2 \leq \mathbb{E}^*(\mathbb{Z}_{\boldsymbol{\omega},\boldsymbol{\omega}^*} - \mathbb{Z}_{\boldsymbol{\omega}',\boldsymbol{\omega}'^*})^2$ , where  $\{e'_i\}$  and  $\{e''_i\}$  are two dependent copies of  $\{e_i\}$ . Applying Theorem 7.2.11 in Vershynin (2018), we obtain

$$\mathbb{E}^* \Big( \sup_{\boldsymbol{\omega}, \boldsymbol{\omega}^*} \mathbb{G}_{\boldsymbol{\omega}, \boldsymbol{\omega}^*} \Big) \le \mathbb{E}^* \Big( \sup_{\boldsymbol{\omega}, \boldsymbol{\omega}^*} \mathbb{Z}_{\boldsymbol{\omega}, \boldsymbol{\omega}^*} \Big).$$
(17)

To bound the supremum of  $\mathbb{Z}_{\omega,\omega^*}$ , using the cone constraint and  $\omega^* \in \omega^{(k)} + \mathbb{B}_1(m)$ , we have

$$\mathbb{E}^{*}\left(\sup_{\boldsymbol{\omega},\boldsymbol{\omega}^{*}} \mathbb{Z}_{\boldsymbol{\omega},\boldsymbol{\omega}^{*}}\right) \leq \frac{\sqrt{2m}}{2r^{2}} \mathbb{E} \left\| \frac{1}{n} \sum_{k \in \mathcal{A}_{m} \cup \{0\}} \sum_{i=1}^{n_{k}} e_{i}' \chi_{i} \boldsymbol{x}_{i}^{(k)} \right\|_{\infty} + \frac{\sqrt{2l}}{2r} \mathbb{E} \left\| \frac{1}{n} \sum_{k \in \mathcal{A}_{m} \cup \{0\}} \sum_{i=1}^{n_{k}} e_{i}'' \chi_{i} \boldsymbol{x}_{i}^{(k)} \right\|_{\infty}.$$
(18)

Thus, by (13) (17) and (18), we have

$$\mathbb{E}\Omega(r,l) \leq \sqrt{\pi} \bigg\{ \frac{m}{r^2} \mathbb{E} \bigg\| \frac{1}{n} \sum_{k \in \mathcal{A}_m \cup \{0\}} \sum_{i=1}^{n_k} e'_i \chi_i \boldsymbol{x}_i^{(k)} \bigg\|_{\infty} + \frac{l}{r} \mathbb{E} \bigg\| \frac{1}{n} \sum_{k \in \mathcal{A}_m \cup \{0\}} \sum_{i=1}^{n_k} e''_i \chi_i \boldsymbol{x}_i^{(k)} \bigg\|_{\infty} \bigg\}.$$
(19)

It remains to find the bound of the two  $l_{\infty}$ -norm terms on the right-hand side of (19). Note that the variable  $|n^{-1}\sum_{k\in\mathcal{A}_m\cup\{0\}}\sum_{i=1}^{n_k}e'_i\chi_i x_{ij}^{(k)}|$  is zero-mean for  $j=1,\ldots,p$ .

$$\mathbb{E}\left[\exp\left(\lambda\left|\frac{1}{n}\sum_{k}\sum_{i=1}^{n_{k}}e_{i}'\chi_{i}x_{ij}^{(k)}\right|\right)\right] \leq \prod_{k}\prod_{i=1}^{n_{k}}\mathbb{E}\left[\exp\left(\left|\frac{\lambda e_{i}'\chi_{i}x_{ij}^{(k)}}{n}\right|\right)\right]$$
$$\leq \prod_{k}\prod_{i=1}^{n_{k}}\exp\left\{\left|\frac{\lambda^{2}\chi_{i}^{2}(x_{ij}^{(k)})^{2}}{2n^{2}}\right|\right\}$$
$$=\exp\left\{\frac{\lambda^{2}}{2n^{2}}\sum_{k}\sum_{i=1}^{n_{k}}\chi_{i}^{2}(x_{ij}^{(k)})^{2}\right\}.$$

Thus,  $|n^{-1}\sum_{k\in\mathcal{A}_m\cup\{0\}}\sum_{i=1}^{n_k} e'_i\chi_i x_{ij}^{(k)}|$  is sub-Gaussian with parameter  $n^{-1}\sqrt{\sum_k\sum_{i=1}^{n_k}\chi_i^2(x_{ij}^{(k)})^2}$ . Applying Lemma 15 in Loh & Wainwright (2015), we have

$$\mathbb{E}\left[\left\|\frac{1}{n}\sum_{k}\sum_{i=1}^{n_{k}}e_{i}'\chi_{i}\boldsymbol{x}_{i}^{(k)}\right\|_{\infty}\left|\boldsymbol{x}_{i}^{(k)}\right| \leq \frac{c}{n}\cdot\max_{j=1,\dots,p}\sqrt{\sum_{k}\sum_{i=1}^{n_{k}}\chi_{i}^{2}(\boldsymbol{x}_{ij}^{(k)})^{2}\cdot\sqrt{\log p}}\right.$$
$$\leq \frac{c\sqrt{\log p}}{n}\sqrt{\sum_{k}\sum_{i=1}^{n_{k}}\chi_{i}^{2}},$$

implying that

$$\mathbb{E} \left\| \frac{1}{n} \sum_{k} \sum_{i=1}^{n_{k}} e_{i}' \chi_{i} \boldsymbol{x}_{i}^{(k)} \right\|_{\infty} \leq c \sqrt{\frac{\log p}{n}} \cdot \mathbb{E} \left[ \sqrt{\frac{\sum_{k} \sum_{i=1}^{n_{k}} \chi_{i}^{2}}{n}} \right] \\
\leq c \sqrt{\frac{\log p}{n}} \cdot \sqrt{\mathbb{E} \left[ \frac{\sum_{k} \sum_{i=1}^{n_{k}} \chi_{i}^{2}}{n} \right]} \\
\leq c \sqrt{\frac{\log p}{n}} \cdot \sqrt{\frac{9f_{u}h}{16}}$$
(20)

Similarly,

$$\mathbb{E} \left\| \frac{1}{n} \sum_{k} \sum_{i=1}^{n_k} e_i'' \chi_i \boldsymbol{x}_i^{(k)} \right\|_{\infty} \le c \sqrt{\frac{\log p}{n}} \cdot \sqrt{\frac{9f_u h}{16}}.$$
(21)

Finally, if we take  $r = h/(4\eta_{1/4})$ ,  $m = h/(16\zeta_p)$  and  $z = t + \log p$ , combining (19), (20), (21) with Bousquet's version inequality (12), we conclude that

$$\Omega(r,l) \le 0.018 f_l + c' l \sqrt{\left(t + \log p\right)/(nh)}$$

with probability at least  $1 - p^{-1}e^{-t}$  for any t > 0, as long as

$$nh \gtrsim f_u f_l^{-2} \log p \max\{l^2/(\eta_{1/4}^2), \eta_{1/4}^4/\zeta_p^2\}$$

This, together with (9), (10) and (11), we have

$$\mathcal{T}(\boldsymbol{\omega}, \boldsymbol{\omega}^*) \geq \frac{\kappa_l}{2} ||\Delta||_{\Sigma}^2 \left[ 0.218 f_l - \left( 0.018 f_l + c' l \sqrt{\frac{t + \log p}{nh}} \right) \right]$$
$$\geq \frac{\kappa_l}{2} ||\Delta||_{\Sigma}^2 \left( 0.2 f_l - c' l \sqrt{\frac{t + \log p}{nh}} \right)$$
(22)

with probability at least  $1 - p^{-1}e^{-t}$ .

It remains to extend this bound to one that is uniform in the ratio  $||\Delta||_1/||\Delta||_{\Sigma}$ , which we do via a peeling argument. Consider the inequality

$$\frac{1}{||\Delta||_{\Sigma}^{2}}\mathcal{T}(\boldsymbol{\omega},\boldsymbol{\omega}^{*}) \geq \frac{\kappa_{l}f_{l}}{10} - c'\gamma\kappa_{l}\frac{||\Delta||_{1}}{||\Delta||_{\Sigma}}\sqrt{\frac{t+\log p}{nh}},$$
(23)

as well as the event

$$A := \left\{ (23) \text{ holds and } \frac{||\Delta||_1}{||\Delta||_{\Sigma}} \le \frac{f_l}{10c'\gamma} \sqrt{\frac{nh}{t + \log p}} \right\}$$

where  $\gamma > 1$  need to be determined. Over A, we have

$$1 \le \frac{||\Delta||_1}{||\Delta||_{\Sigma}} \le \frac{f_l}{10c'\gamma} \sqrt{\frac{nh}{t + \log p}}.$$

Define the function

$$g(l) := c' l \kappa_l \sqrt{\frac{t + \log p}{nh}}, \text{ and } i(\Delta) := ||\Delta||_1/||\Delta||_{\Sigma},$$

as well as the set

$$\Theta_k = \{ \Delta \in \mathbb{R}^p : \gamma^{k-1} \mu \le g(i(\Delta)) \le \gamma^k \mu \},\$$

for each  $k = 1, ..., N := \lceil \log(\sqrt{(nh)/(t + \log p)}) / \log(\gamma) \rceil$ , where  $\mu = c' \kappa_l \sqrt{(t + \log p)/(nh)}$ . By the union bound, we then have

$$\mathbb{P}(A^{c}) \leq \sum_{k=1}^{N} \mathbb{P}\left\{ \exists \Delta \in \Theta_{k} \text{ s.t. } \frac{\kappa_{l}f_{l}}{10} - \frac{1}{||\Delta||_{\Sigma}^{2}}\mathcal{T}(\boldsymbol{\omega}, \boldsymbol{\omega}^{*}) > \gamma^{k}\mu \right\}$$
$$\leq \sum_{k=1}^{N} \mathbb{P}\left\{ \sup_{||\Delta||_{1}/||\Delta||_{\Sigma} \leq g^{-1}(\gamma^{k}\mu)} \frac{\kappa_{l}f_{l}}{10} - \frac{1}{||\Delta||_{\Sigma}^{2}}\mathcal{T}(\boldsymbol{\omega}, \boldsymbol{\omega}^{*}) > \gamma^{k}\mu \right\}$$
$$\leq \sum_{k=1}^{N} p^{-1}e^{-t} \leq N \cdot p^{-1}e^{-t}.$$

Taking  $\gamma = e^{1/e}$  and  $t = \log\{e \log(l/r_l)\} + u$  yields that with probability at least  $1 - p^{-1}e^{-u}$ ,

$$\frac{\kappa_l f_l}{10} - \frac{1}{||\Delta||_{\Sigma}^2} \mathcal{T}(\boldsymbol{\omega}, \boldsymbol{\omega}^*) \leq \frac{c'' ||\Delta||_1}{||\Delta||_{\Sigma}} \sqrt{\frac{\log p + \log\{e \log(l/r_l)\} + u}{nh}}$$

Multiplying by  $||\Delta||_{\Sigma}^2$  on both sides yields

$$\mathcal{T}(\boldsymbol{\omega}, \boldsymbol{\omega}^*) \geq \frac{\kappa_l f_l}{10} ||\Delta||_{\Sigma}^2 - c'' ||\Delta||_1 ||\Delta||_{\Sigma} \sqrt{\frac{\log p + \log\{e \log(l/r_l)\} + u}{nh}}$$

# A.3 Proof of Proposition 3.2

The Taylor error around  $\beta_2$  in the direction  $\beta_1 - \beta_2$  is given by

$$\mathcal{T}(\beta_1,\beta_2) = \hat{Q}_g^{(0)}(\beta_1) - \hat{Q}_g^{(0)}(\beta_2) - \langle \nabla \hat{Q}_g^{(0)}(\beta_2), \beta_1 - \beta_2 \rangle.$$

For a given kernel function  $K(\cdot)$  and bandwidth g > 0, the smoothed quantile loss  $\hat{Q}_g^{(0)}$  can be written as  $(n_0g)^{-1}\sum_{i=1}^{n_0}\int_{-\infty}^{\infty}\rho_{\tau}(u)K\{(u+\langle \boldsymbol{x}_i^{(0)},\boldsymbol{\beta}\rangle-y_i^{(0)})/g\}du$ . Therefore

$$\begin{aligned} \mathcal{T}(\boldsymbol{\beta}_{1},\boldsymbol{\beta}_{2}) &= \frac{1}{2n_{0}} \sum_{i=1}^{n_{0}} K_{g} \big\{ \epsilon_{i} - \langle \boldsymbol{x}_{i}^{(0)}, (1-t)\boldsymbol{\beta}_{1} + t\boldsymbol{\beta}_{2} - \boldsymbol{\beta}^{*} \rangle \big\} \langle \boldsymbol{x}_{i}^{(0)}, \boldsymbol{\beta}_{1} - \boldsymbol{\beta}_{2} \rangle^{2} \\ &= \frac{1}{2n_{0}} \sum_{i=1}^{n_{0}} K_{g} \big\{ \epsilon_{i} - t \langle \boldsymbol{x}_{i}^{(0)}, \boldsymbol{\beta}_{2} - \boldsymbol{\beta}_{1} \rangle - \langle \boldsymbol{x}_{i}^{(0)}, \boldsymbol{\beta}_{1} - \boldsymbol{\beta}^{*} \rangle \big\} \langle \boldsymbol{x}_{i}^{(0)}, \boldsymbol{\beta}_{1} - \boldsymbol{\beta}_{2} \rangle^{2} \end{aligned}$$

for some  $t \in [0, 1]$ , For each *i*, define the event  $\mathcal{E}_i$ ,

$$\mathcal{E}_i = \{ |\epsilon_i| \le g/4 \} \cap \{ |\langle \boldsymbol{x}_i^{(0)}, \boldsymbol{\beta}_1 - \boldsymbol{\beta}_2 \rangle| \le g ||\boldsymbol{\beta}_1 - \boldsymbol{\beta}_2||_{\Sigma}/(2r) \} \cap \{ |\langle \boldsymbol{x}_i^{(0)}, \boldsymbol{\beta}_1 - \boldsymbol{\beta}^* \rangle| \le g/4 \},$$

for all  $\beta_1 - \beta_2 \in \mathbb{B}_{\Sigma}(r)$ . Thus

$$\mathcal{T}(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) \ge \frac{\kappa_l}{2n_0 g} \sum_{i=1}^{n_0} \langle \boldsymbol{x}_i^{(0)}, \boldsymbol{\delta} \rangle^2 \mathbb{1}_{\mathcal{E}_i},$$
(24)

where  $\kappa_l = \min_{|u| \leq 1} K(u)$  and  $\boldsymbol{\delta} = \boldsymbol{\beta}_1 - \boldsymbol{\beta}_2$ . For a truncation level R > 0, define functions

$$\varphi_R(u) = \begin{cases} u^2 & |u| \le \frac{R}{2}, \\ (R - |u|)^2 & \frac{R}{2} < |u| \le R, \\ 0 & |u| > R. \end{cases}$$

By this construction,  $\varphi_R(u) \leq u^2 \cdot \mathbb{1}\{|u| \leq R\}$ ,  $\varphi_{cR}(cu) = c^2 \varphi_R(u)$  and  $\varphi_R$  is R-Lipschitz. In addition, we define the trapezoidal function

$$\psi_R(u) = \begin{cases} 1 & |u| \le \frac{R}{2}, \\ 2 - \frac{2}{R}|u| & \frac{R}{2} < |u| \le R, \\ 0 & |u| > R, \end{cases}$$

and note that  $\psi_R$  is (2/R)-Lipschitz and  $\psi_R(u) \leq \mathbb{1}\{|u| \leq R\}$ . From these two new-defined function, (24) implies

$$\mathcal{T}(\boldsymbol{\beta}_{1},\boldsymbol{\beta}_{2}) \geq \frac{\kappa_{l}}{2n_{0}g} ||\boldsymbol{\delta}||_{\Sigma}^{2} \sum_{i=1}^{n_{0}} \mathbb{1}\{|\epsilon_{i}| \leq g/4\} \varphi_{g||\boldsymbol{\delta}||_{\Sigma}/(2r)}(\langle \boldsymbol{x}_{i}^{(0)},\boldsymbol{\delta} \rangle) \psi_{g/4}(\langle \boldsymbol{x}_{i}^{(0)},\boldsymbol{\beta}_{1}-\boldsymbol{\beta}^{*} \rangle)$$

$$\geq \frac{\kappa_{l}}{2} ||\boldsymbol{\delta}||_{\Sigma}^{2} \cdot \underbrace{\frac{1}{n_{0}g} \sum_{i=1}^{n_{0}} \mathbb{1}\{|\epsilon_{i}| \leq g/4\} \varphi_{g/(2r)}(\langle \boldsymbol{x}_{i}^{(0)},\boldsymbol{\delta} \rangle/||\boldsymbol{\delta}||_{\Sigma}) \psi_{g/4}(\langle \boldsymbol{x}_{i}^{(0)},\boldsymbol{\beta}_{1}-\boldsymbol{\beta}^{*} \rangle)}_{D_{0}(\boldsymbol{\beta}_{1},\boldsymbol{\beta}_{2})}$$

$$(25)$$

In the following proofs, we bound  $\mathbb{E}D_0(\beta_1, \beta_2)$  and  $D_0(\beta_1, \beta_2) - \mathbb{E}D_0(\beta_1, \beta_2)$ , respectively. Write  $\boldsymbol{\nu} = \Sigma^{1/2} \boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma} \in \mathbb{S}^{p-1}$ . Note that

$$\begin{aligned} \left|\frac{g}{2}f_{\epsilon|\boldsymbol{x}}(0)\right| - \left|\mathbb{E}[\mathbbm{1}\{|\epsilon_i| \le g/4\} | \boldsymbol{x}_i^{(0)}]\right| \le \left|\mathbb{E}[\mathbbm{1}\{|\epsilon_i| \le g/4\} | \boldsymbol{x}_i^{(0)}] - \frac{g}{2}f_{\epsilon|\boldsymbol{x}}(0)\right| \\ \le \int_{-g/4}^{g/4} |f_{\epsilon|\boldsymbol{x}}(t) - f_{\epsilon|\boldsymbol{x}}(0)| dt. \end{aligned}$$

Then,

$$\begin{aligned} \left|\frac{g}{2}f_{\epsilon|\boldsymbol{x}}(0)\right| &- \left|\mathbb{E}[\mathbbm{1}\{|\epsilon_i| \le g/4\} | \boldsymbol{x}_i^{(0)}]\right| \le \frac{l_0 g^2}{16} \\ &\left|\mathbb{E}[\mathbbm{1}\{|\epsilon_i| \le g/4\} | \boldsymbol{x}_i^{(0)}]\right| \ge \frac{g}{2}f_l - \frac{l_0 g^2}{16} \end{aligned}$$

Provided  $g \leq f_l/l_0 \leq f_u/l_0$ , we have

$$\left| \mathbb{E}[\mathbb{1}\{|\epsilon_i| \le g/4\} | \boldsymbol{x}_i^{(0)}] \right| \ge \frac{7f_l g}{16}.$$

Meanwhile

$$\left| \mathbb{E}[\mathbb{1}\{|\epsilon_i| \le g/4\} | \boldsymbol{x}_i^{(0)}] \right| - \left| \frac{g}{2} f_{\epsilon|\boldsymbol{x}}(0) \right| \le \int_{-g/4}^{g/4} |f_{\epsilon|\boldsymbol{x}}(t) - f_{\epsilon|\boldsymbol{x}}(0)| dt$$

implies

$$\left| \mathbb{E}[\mathbb{1}\{|\epsilon_i| \le g/4\} | \boldsymbol{x}_i^{(0)}] \right| \le \frac{9f_u g}{16}.$$

Then

$$\begin{split} \mathbb{E}D_{0}(\boldsymbol{\beta}_{1},\boldsymbol{\beta}_{2}) \\ &= \frac{1}{n_{0}g}\sum_{i=1}^{n_{0}}\mathbb{E}\Big[\mathbb{E}[\mathbbm{1}\{|\boldsymbol{\epsilon}_{i}|\leq g/4\}|\boldsymbol{x}_{i}^{(0)}]\varphi_{g/(2r)}(\langle\boldsymbol{x}_{i}^{(0)},\boldsymbol{\delta}\rangle/||\boldsymbol{\delta}||_{\Sigma})\psi_{g/4}(\langle\boldsymbol{x}_{i}^{(0)},\boldsymbol{\beta}_{1}-\boldsymbol{\beta}^{*}\rangle)\Big] \\ &\geq \frac{7f_{l}}{16}\mathbb{E}\Big[\varphi_{g/(2r)}(\langle\boldsymbol{x},\boldsymbol{\delta}\rangle/||\boldsymbol{\delta}||_{\Sigma})\psi_{g/4}(\langle\boldsymbol{x},\boldsymbol{\beta}_{1}-\boldsymbol{\beta}^{*}\rangle)\Big] \\ &\geq \frac{7f_{l}}{16}\Big(1-\mathbb{E}\Big[\langle\boldsymbol{x},\boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma}\rangle^{2}\mathbbm{1}\{|\langle\boldsymbol{x},\boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma}\rangle|>g/(4r)\}\Big] \\ &\quad -\mathbb{E}\Big[\langle\boldsymbol{x},\boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma}\rangle^{2}\mathbbm{1}\{|\langle\boldsymbol{x},\boldsymbol{\beta}_{1}-\boldsymbol{\beta}^{*}\rangle|>g/8\}\Big]\Big). \end{split}$$

By the definition of  $\eta_{\delta}$ , as long as  $0 < r \le g/(4\eta_{1/4})$ ,

$$\sup_{\boldsymbol{\delta}\in\mathbb{B}_{\Sigma}(r)}\mathbb{E}\Big[\langle \boldsymbol{x},\boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma}\rangle^{2}\mathbb{1}\{|\langle \boldsymbol{x},\boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma}\rangle|>g/(4r)\}\Big]\leq\frac{1}{4}.$$

Moreover,  $\beta_1 \in \beta^* + \mathbb{B}_{\Sigma}(r/2)$ .

$$\mathbb{E}\left[\langle \boldsymbol{x}, \boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma}\rangle^{2}\mathbb{1}\{|\langle \boldsymbol{x}, \boldsymbol{\beta}_{1} - \boldsymbol{\beta}^{*}\rangle| > g/8\}\right]$$

$$\leq \left(\frac{8}{g}\right)^{2} \mathbb{E}\left[\langle \boldsymbol{x}, \boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma}\rangle^{2}\langle \boldsymbol{x}, \boldsymbol{\beta}_{1} - \boldsymbol{\beta}^{*}\rangle^{2}\right]$$

$$\leq \left(\frac{8}{g}\right)^{2} \left(\mathbb{E}\langle \boldsymbol{x}, \boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma}\rangle^{4}\right)^{1/2} \left(\mathbb{E}\langle \boldsymbol{x}, \boldsymbol{\beta}_{1} - \boldsymbol{\beta}^{*}\rangle^{4}\right)^{1/2}$$

$$\leq \left(\frac{8}{g}\right)^{2} \zeta_{p} \cdot \left[\left(\frac{r}{2}\right)^{4} \mathbb{E}\langle \boldsymbol{x}, (\boldsymbol{\beta}_{1} - \boldsymbol{\beta}^{*})/||\boldsymbol{\beta}_{1} - \boldsymbol{\beta}^{*}||_{\Sigma}\rangle^{4}\right]^{1/2}$$

$$\leq \left(\frac{8}{g}\right)^{2} \zeta_{p} \cdot \left(\frac{r}{2}\right)^{2} \zeta_{p} = \left(\frac{4r}{g}\right)^{2} \zeta_{p}^{2}.$$

The inequality comes from  $\mathbb{E}\langle \boldsymbol{x}, \boldsymbol{u}/||\boldsymbol{u}||_{\Sigma}\rangle^4 \leq \zeta_p^2$ .

Provided  $8r\zeta_p \leq g$ ,

$$\mathbb{E}D_0(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) > 0.218 f_l.$$
<sup>(26)</sup>

Next we find a lower bound of  $D_0(\beta_1, \beta_2) - \mathbb{E}D_0(\beta_1, \beta_2)$  over  $\Lambda(r, l) := \{(\beta_1, \beta_2) : \beta_1 \in \beta^* + \mathbb{B}_{\Sigma}(r/2), \beta_2 \in \beta_1 + \mathbb{B}_{\Sigma}(r) \cap \mathbb{C}_{\Sigma}(l), \operatorname{supp}(\beta_1) \subseteq S\}$ . Define

$$\Omega(r,l) = \sup_{(\boldsymbol{\beta}_1,\boldsymbol{\beta}_2)\in\Lambda(r,l)} \{-D_0(\boldsymbol{\beta}_1,\boldsymbol{\beta}_2) + \mathbb{E}D_0(\boldsymbol{\beta}_1,\boldsymbol{\beta}_2)\}.$$

Write  $D_0(\beta_1, \beta_2) = (1/n_0) \sum_{i=1}^{n_0} w_i(\beta_1, \beta_2)$ , where

$$w_i(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) = g^{-1} \mathbb{1}\{|\boldsymbol{\epsilon}_i| \le g/4\} \varphi_{g/(2r)}(\langle \boldsymbol{x}_i^{(0)}, \boldsymbol{\delta} \rangle / ||\boldsymbol{\delta}||_{\Sigma}) \psi_{g/4}(\langle \boldsymbol{x}_i^{(0)}, \boldsymbol{\beta}_1 - \boldsymbol{\beta}^* \rangle)$$

satisfies  $0 \le w_i(\beta_1, \beta_2) \le g/(4r)^2$ , since  $\varphi_R(u) \le (R/2)^2$  and  $\psi_R(u) \in [0, 1]$ . Moreover,

$$\begin{split} \mathbb{E}w_i^2(\boldsymbol{\beta}_1,\boldsymbol{\beta}_2) &= \mathbb{E}\Big[g^{-2}\mathbbm{1}\{|\epsilon_i| \leq g/4\}\varphi_{g/(2r)}^2(\langle \boldsymbol{x}_i^{(0)},\boldsymbol{\delta}\rangle/||\boldsymbol{\delta}||_{\Sigma})\psi_{g/4}^2(\langle \boldsymbol{x}_i^{(0)},\boldsymbol{\beta}_1-\boldsymbol{\beta}^*\rangle)\Big] \\ &\leq \frac{9f_u}{16g} \cdot \mathbb{E}\varphi_{g/(2r)}^2(\langle \boldsymbol{x}_i^{(0)},\boldsymbol{\delta}\rangle/||\boldsymbol{\delta}||_{\Sigma}) \\ &\leq \frac{9f_u}{16g} \cdot \mathbb{E}(\langle \boldsymbol{x}_i^{(0)},\boldsymbol{\delta}\rangle/||\boldsymbol{\delta}||_{\Sigma})^4 = \frac{9f_u\zeta_p^2}{16g}. \end{split}$$

Using Bousquet's version of Talagrand's inequality yields that, for any z > 0,

$$\Omega(r,l) \leq \mathbb{E}\Omega(r,l) + \{\mathbb{E}\Omega(r,l)\}^{1/2} \frac{1}{2r} \sqrt{\frac{gz}{n_0}} + \frac{3\sqrt{2}\zeta_p}{4} \sqrt{\frac{f_u z}{n_0 g}} + \frac{g}{(4r)^2} \frac{z}{3n_0} \\
\leq \mathbb{E}\Omega(r,l) + \frac{1}{4} \mathbb{E}\Omega(r,l) + \frac{1}{4r^2} \frac{gz}{n_0} + \frac{3\sqrt{2}\zeta_p}{4} \sqrt{\frac{f_u z}{n_0 g}} + \frac{g}{(4r)^2} \frac{z}{3n_0} \\
\leq \frac{5}{4} \mathbb{E}\Omega(r,l) + \frac{3\sqrt{2}\zeta_p}{4} \sqrt{\frac{f_u z}{n_0 g}} + \frac{13}{3} \frac{gz}{(4r)^2 n_0}$$
(27)

holds with probability at least  $1-e^{-z}$ . To bound the expectation  $\mathbb{E}\Omega(r, l)$ , using Rademacher symmetrization and the connection between Gaussian and Rademacher complexities, Lemma 5.5 in Ledoux & Talagrand (1991), we have

$$\mathbb{E}\Omega(r,l) \le 2\sqrt{\frac{\pi}{2}} \mathbb{E}\left[\sup_{(\beta_1,\beta_2)\in\Lambda(r,l)} \mathbb{G}_{\beta_1,\beta_2}\right],\tag{28}$$

where  $\mathbb{G}_{\beta_1,\beta_2} := (n_0 g)^{-1} \sum_{i=1}^{n_0} e_i \mathbb{1}\{|\epsilon_i| \leq g/4\} \varphi_{g/(2r)}(\langle \boldsymbol{x}_i^{(0)}, \boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma}\rangle) \psi_{g/4}(\langle \boldsymbol{x}_i^{(0)}, \boldsymbol{\beta}_1 - \boldsymbol{\beta}^*\rangle)$  and  $e_i$  are independent standard normal variables. Note that  $\mathbb{G}_{\beta_1,\beta_2}$  is a Gaussian process conditioned on  $\{(y_i^{(0)}, \boldsymbol{x}_i^{(0)})\}_{i=1}^{n_0}$  and  $\mathbb{G}_{\beta^*,\beta^*} = 0$ . For  $(\beta_1,\beta_2)$  and  $(\beta'_1,\beta'_2)$ , write  $\boldsymbol{\delta} = \beta_1 - \beta_2$ ,  $\boldsymbol{\delta}' = \beta'_1 - \beta'_2$  and  $\chi_i = \mathbb{1}\{|\epsilon_i| \leq g/4\}$ , then

$$\begin{split} &\mathbb{G}_{\beta_{1},\beta_{2}} - \mathbb{G}_{\beta_{1}',\beta_{2}'} \\ &= \mathbb{G}_{\beta_{1},\beta_{2}} - \mathbb{G}_{\beta_{1}',\beta_{1}'+\delta} + \mathbb{G}_{\beta_{1}',\beta_{1}'+\delta} - \mathbb{G}_{\beta_{1}',\beta_{2}'} \\ &= \frac{1}{n_{0}g} \sum_{i=1}^{n_{0}} e_{i}\chi_{i}\varphi_{g/(2r)}(\langle \boldsymbol{x}_{i}^{(0)}, \boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma}\rangle)\{\psi_{g/4}(\langle \boldsymbol{x}_{i}^{(0)}, \beta_{1} - \boldsymbol{\beta}^{*}\rangle) - \psi_{g/4}(\langle \boldsymbol{x}_{i}^{(0)}, \beta_{1}' - \boldsymbol{\beta}^{*}\rangle)\} \\ &+ \frac{1}{n_{0}g} \sum_{i=1}^{n_{0}} e_{i}\chi_{i}\psi_{g/4}(\langle \boldsymbol{x}_{i}^{(0)}, \beta_{1}' - \boldsymbol{\beta}^{*}\rangle)\{\varphi_{g/(2r)}(\boldsymbol{x}_{i}^{(0)T}\boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma}) - \varphi_{g/(2r)}(\boldsymbol{x}_{i}^{(0)T}\boldsymbol{\delta}'/||\boldsymbol{\delta}'||_{\Sigma})\}. \end{split}$$

Note that  $\varphi_R$  and  $\psi_R$  are Lipschitz continuous, and  $\varphi_R(u) \leq (R/2)^2$ . Let  $\mathbb{E}^*$  be the conditional expectation given  $\{(y_i^{(0)}, \boldsymbol{x}_i^{(0)})\}_{i=1}^{n_0}$ . Consequently,

$$\mathbb{E}^{*} \left( \mathbb{G}_{\beta_{1},\beta_{2}} - \mathbb{G}_{\beta_{1}',\beta_{1}'+\delta} \right)^{2} \leq \frac{1}{(n_{0}g)^{2}} \left( \frac{8}{g} \right)^{2} \left( \frac{g}{4r} \right)^{4} \sum_{i=1}^{n_{0}} \chi_{i} \langle \boldsymbol{x}_{i}^{(0)}, \beta_{1} - \beta_{1}' \rangle^{2}$$
$$= \frac{1}{4r^{4}n_{0}^{2}} \sum_{i=1}^{n_{0}} \chi_{i} \langle \boldsymbol{x}_{i}^{(0)}, \beta_{1} - \beta_{1}' \rangle^{2}$$
(29)

and

$$\mathbb{E}^{*} \left( \mathbb{G}_{\boldsymbol{\beta}_{1}^{\prime},\boldsymbol{\beta}_{1}^{\prime}+\boldsymbol{\delta}} - \mathbb{G}_{\boldsymbol{\beta}_{1}^{\prime},\boldsymbol{\beta}_{2}^{\prime}} \right)^{2} \leq \frac{1}{(n_{0}g)^{2}} \left( \frac{g}{2r} \right)^{2} \sum_{i=1}^{n_{0}} \chi_{i} \left( \langle \boldsymbol{x}_{i}^{(0)}, \boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma} \right) - \langle \boldsymbol{x}_{i}^{(0)}, \boldsymbol{\delta}^{\prime}/||\boldsymbol{\delta}^{\prime}||_{\Sigma} \rangle^{2}$$
$$= \frac{1}{4r^{2}n_{0}^{2}} \sum_{i=1}^{n_{0}} \chi_{i} \langle \boldsymbol{x}_{i}^{(0)}, \boldsymbol{\delta}/||\boldsymbol{\delta}||_{\Sigma} - \boldsymbol{\delta}^{\prime}/||\boldsymbol{\delta}^{\prime}||_{\Sigma} \rangle^{2}.$$
(30)

Motivated by the last two inequalities, we have

$$\mathbb{E}^* \left( \mathbb{G}_{oldsymbol{eta}_1,oldsymbol{eta}_2} - \mathbb{G}_{oldsymbol{eta}_1',oldsymbol{eta}_2'} 
ight)^2 \ \leq rac{1}{2r^4n_0^2} \sum_{i=1}^{n_0} \chi_i \langle oldsymbol{x}_i^{(0)}, oldsymbol{eta}_1 - oldsymbol{eta}_1' 
angle^2 + rac{1}{2r^2n_0^2} \sum_{i=1}^{n_0} \chi_i \langle oldsymbol{x}_i^{(0)}, oldsymbol{\delta}/||oldsymbol{\delta}||_{\Sigma} - oldsymbol{\delta}'/||oldsymbol{\delta}'||_{\Sigma} 
angle^2.$$

Define another Gaussian process  $\mathbb{Z}_{\pmb{\beta}_1,\pmb{\beta}_2}$  as

$$\begin{aligned} \mathbb{Z}_{\boldsymbol{\beta}_{1},\boldsymbol{\beta}_{2}} &= \frac{1}{2^{1/2}r^{2}n_{0}}\sum_{i=1}^{n_{0}}e_{i}'\chi_{i}\langle\boldsymbol{x}_{i}^{(0)},\boldsymbol{\beta}_{1}-\boldsymbol{\beta}^{*}\rangle + \frac{1}{2^{1/2}rn_{0}}\sum_{i=1}^{n_{0}}e_{i}''\chi_{i}\langle\boldsymbol{x}_{i}^{(0)},\boldsymbol{\beta}_{2}-\boldsymbol{\beta}_{1}\rangle/||\boldsymbol{\delta}||_{\Sigma} \\ &= \frac{1}{2^{1/2}r^{2}n_{0}}\sum_{i=1}^{n_{0}}e_{i}'\chi_{i}\langle\boldsymbol{x}_{i,\mathcal{S}}^{(0)},(\boldsymbol{\beta}_{1}-\boldsymbol{\beta}^{*})_{\mathcal{S}}\rangle + \frac{1}{2^{1/2}rn_{0}}\sum_{i=1}^{n_{0}}e_{i}''\chi_{i}\langle\boldsymbol{x}_{i}^{(0)},\boldsymbol{\beta}_{2}-\boldsymbol{\beta}_{1}\rangle/||\boldsymbol{\delta}||_{\Sigma} \end{aligned}$$

such that  $\mathbb{E}^*(\mathbb{G}_{\beta_1,\beta_2} - \mathbb{G}_{\beta'_1,\beta'_2})^2 \leq \mathbb{E}^*(\mathbb{Z}_{\beta_1,\beta_2} - \mathbb{Z}_{\beta'_1,\beta'_2})^2$ , where  $\{e'_i\}$  and  $\{e''_i\}$  are two dependent copies of  $\{e_i\}$ . The second equity holds since  $\operatorname{supp}(\beta_1)$ ,  $\operatorname{supp}(\beta^*) \subseteq S$ . Applying Theorem 7.2.11 in Vershynin (2018), we obtain

$$\mathbb{E}^* \Big( \sup_{\beta_1,\beta_2} \mathbb{G}_{\beta_1,\beta_2} \Big) \le \mathbb{E}^* \Big( \sup_{\beta_1,\beta_2} \mathbb{Z}_{\beta_1,\beta_2} \Big).$$
(31)

To bound the supremum of  $\mathbb{Z}_{\beta_1,\beta_2}$ , using the cone constraint and  $\beta_1 \in \beta^* + \mathbb{B}_{\Sigma}(r/2)$ , we have

$$\mathbb{E}^{*}\left(\sup_{\beta_{1},\beta_{2}} \mathbb{Z}_{\beta_{1},\beta_{2}}\right) \leq \frac{\sqrt{2}}{4r} \mathbb{E} \left\| \frac{1}{n_{0}} \sum_{i=1}^{n_{0}} e_{i}' \chi_{i} \mathbf{S}^{-1/2} \mathbf{x}_{i,\mathcal{S}}^{(0)} \right\|_{2} + \frac{\sqrt{2}l}{2r} \mathbb{E} \left\| \frac{1}{n_{0}} \sum_{i=1}^{n_{0}} e_{i}'' \chi_{i} \mathbf{x}_{i}^{(0)} \right\|_{\infty} \\
\leq \frac{\sqrt{2}}{4r} \sqrt{\frac{9f_{u}g}{16} \frac{s}{n_{0}}} + \frac{\sqrt{2}l}{2r} \mathbb{E} \left\| \frac{1}{n_{0}} \sum_{i=1}^{n_{0}} e_{i}'' \chi_{i} \mathbf{x}_{i}^{(0)} \right\|_{\infty},$$
(32)

where  $\boldsymbol{S} = \Sigma_{SS} = \mathbb{E}(\boldsymbol{x}_{S}\boldsymbol{x}_{S}^{\mathrm{T}})$ . Thus, by (28) (31) and (32), we have

$$\mathbb{E}\Omega(r,l) \le \sqrt{\pi} \left\{ \frac{3}{8r} \sqrt{\frac{f_u gs}{n_0}} + \frac{l}{r} \mathbb{E} \left\| \frac{1}{n_0} \sum_{i=1}^{n_0} e_i'' \chi_i \boldsymbol{x}_i^{(0)} \right\|_{\infty} \right\}.$$
(33)

It remains to find the bound of the second term on the right-hand side of (33). Note that the variable  $|n_0^{-1} \sum_{i=1}^{n_0} e''_i \chi_i x_{ij}^{(0)}|$  is zero-mean for  $j = 1, \ldots, p$ .

$$\mathbb{E}\left[\exp\left(\lambda \left|\frac{1}{n_{0}} \sum_{i=1}^{n_{0}} e_{i}''\chi_{i}x_{ij}^{(0)}\right|\right)\right] \leq \prod_{i=1}^{n_{0}} \mathbb{E}\left[\exp\left(\left|\frac{\lambda e_{i}''\chi_{i}x_{ij}^{(0)}}{n_{0}}\right|\right)\right]$$
$$\leq \prod_{i=1}^{n_{0}} \exp\left\{\left|\frac{\lambda^{2}\chi_{i}^{2}(x_{ij}^{(0)})^{2}}{2n_{0}^{2}}\right|\right\}$$
$$= \exp\left\{\frac{\lambda^{2}}{2n_{0}^{2}} \sum_{i=1}^{n_{0}} \chi_{i}^{2}(x_{ij}^{(0)})^{2}\right\}.$$

Thus,  $|n_0^{-1} \sum_{i=1}^{n_0} e_i'' \chi_i x_{ij}^{(0)}|$  is sub-Gaussian with parameter  $n_0^{-1} \sqrt{\sum_{i=1}^{n_0} \chi_i^2 (x_{ij}^{(0)})^2}$ . Applying Lemma 15 in Loh & Wainwright (2015), we have

$$\mathbb{E}\left[\left\|\frac{1}{n_0}\sum_{i=1}^{n_0} e_i''\chi_i \boldsymbol{x}_i^{(0)}\right\|_{\infty} \middle| \boldsymbol{x}_i^{(0)}\right] \le \frac{c_0}{n_0} \cdot \max_{j=1,\dots,p} \sqrt{\sum_{i=1}^{n_0} \chi_i^2(\boldsymbol{x}_{ij}^{(0)})^2} \cdot \sqrt{\log p} \\ \le \frac{c_0 \sqrt{\log p}}{n_0} \sqrt{\sum_{i=1}^{n_0} \chi_i^2},$$

implying that

$$\mathbb{E} \left\| \frac{1}{n_0} \sum_{i=1}^{n_0} e_i'' \chi_i \boldsymbol{x}_i^{(0)} \right\|_{\infty} \le c_0 \sqrt{\frac{\log p}{n_0}} \cdot \mathbb{E} \left[ \sqrt{\frac{\sum_{i=1}^{n_0} \chi_i^2}{n_0}} \right]$$
$$\le c_0 \sqrt{\frac{\log p}{n_0}} \cdot \sqrt{\mathbb{E} \left[ \frac{\sum_{i=1}^{n_0} \chi_i^2}{n_0} \right]}$$
$$\le c_0 \sqrt{\frac{\log p}{n_0}} \cdot \sqrt{\frac{9f_u g}{16}}$$

Therefore,

$$\mathbb{E} \left\| \frac{1}{n_0} \sum_{i=1}^{n_0} e_i'' \chi_i \boldsymbol{x}_i^{(0)} \right\|_{\infty} \le \frac{3c_0}{4} \sqrt{\frac{f_u g \log p}{n_0}}.$$
(34)

Plug this bound to (33), we obtain

$$\mathbb{E}\Omega(r,l) \le \sqrt{\pi} \left\{ \frac{3}{8r} \sqrt{\frac{f_u gs}{n_0}} + \frac{3c_0 l}{4r} \sqrt{\frac{f_u g \log p}{n_0}} \right\}.$$
(35)

Finally, if we take  $r = g/(48c_0)$  and  $z = t + \log p$ , combining (35) with Bousquet's version inequality (27), we conclude that

$$\Omega(r,l) \le \sqrt{\pi} \left\{ \frac{15}{32r} \sqrt{\frac{f_u gs}{n_0}} + \frac{15c_0 l}{16r} \sqrt{\frac{f_u g \log p}{n_0}} \right\} + \frac{3\sqrt{2}\zeta_p}{4} \sqrt{\frac{f_u z}{n_0 g}} + \frac{13}{3} \frac{gz}{(4r)^2 n_0}$$
$$\Omega(r,l) \le 0.018 f_l + c' l \sqrt{(t+\log p)/(n_0 g)}$$

with probability at least  $1 - p^{-1}e^{-t}$  as long as  $n_0g \gtrsim f_u f_l^{-2} \max\{s, l^2 \log p\}$ . This, together with (24), (25) and (26), we have

$$\mathcal{T}(\boldsymbol{\beta}_{1},\boldsymbol{\beta}_{2}) \geq \frac{\kappa_{l}}{2} ||\boldsymbol{\delta}||_{\Sigma}^{2} \left[ 0.218f_{l} - \left( 0.018f_{l} + c'l\sqrt{\frac{t + \log p}{n_{0}g}} \right) \right]$$
$$\geq \frac{\kappa_{l}}{2} ||\boldsymbol{\delta}||_{\Sigma}^{2} \left( 0.2f_{l} - c'l\sqrt{\frac{t + \log p}{n_{0}g}} \right)$$
(36)

with probability at least  $1 - p^{-1}e^{-t}$ .

It remains to extend this bound to one that is uniform in the ratio  $||\delta||_1/||\delta||_{\Sigma}$ , which we do via a peeling argument. Consider the inequality

$$\mathcal{T}(\boldsymbol{\beta}_{1},\boldsymbol{\beta}_{2}) \geq \frac{\kappa_{l}}{2} ||\boldsymbol{\delta}||_{\Sigma}^{2} \left( 0.2f_{l} - c'l\sqrt{\frac{t + \log p}{n_{0}g}} \right)$$
$$\frac{1}{||\boldsymbol{\delta}||_{\Sigma}^{2}} \mathcal{T}(\boldsymbol{\beta}_{1},\boldsymbol{\beta}_{2}) \geq \frac{\kappa_{l}f_{l}}{10} - \frac{c'l\kappa_{l}}{2}\sqrt{\frac{t + \log p}{n_{0}g}}$$
$$\frac{\kappa_{l}f_{l}}{10} - \frac{1}{||\boldsymbol{\delta}||_{\Sigma}^{2}} \mathcal{T}(\boldsymbol{\beta}_{1},\boldsymbol{\beta}_{2}) \leq c''l\sqrt{\frac{t + \log p}{n_{0}g}}.$$

For some  $\gamma > 1$  to be determined and positive integers  $k = 1, \ldots, N := \lceil \log(l/r_l) / \log(\gamma) \rceil$ , define the set  $\Theta_k = \{ \boldsymbol{\delta} \in \mathbb{R}^p : \gamma^{k-1}r_l \leq ||\boldsymbol{\delta}||_1 / ||\boldsymbol{\delta}||_{\Sigma} \leq \gamma^k r_l \}$ , so that  $\{ \boldsymbol{\delta} \in \mathbb{R}^p : r_l \leq ||\boldsymbol{\delta}||_1 / ||\boldsymbol{\delta}||_{\Sigma} \leq l \} \subseteq \cup_{k=1}^N \Theta_k$ . Then

$$\begin{aligned} & \mathbb{P}\bigg\{ \exists \boldsymbol{\delta} \in \{\boldsymbol{\delta} \in \mathbb{R}^p : r_l \leq ||\boldsymbol{\delta}||_1 / ||\boldsymbol{\delta}||_{\Sigma} \leq l \} \text{ s.t.} \\ & \frac{\kappa_l f_l}{10} - \frac{1}{||\boldsymbol{\delta}||_{\Sigma}^2} \mathcal{T}(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) > \frac{c'' \gamma ||\boldsymbol{\delta}||_1}{||\boldsymbol{\delta}||_{\Sigma}} \sqrt{\frac{t + \log p}{n_0 g}} \bigg\} \\ & \leq \sum_{k=1}^N \mathbb{P}\bigg\{ \exists \boldsymbol{\delta} \in \Theta_k \text{ s.t. } \frac{\kappa_l f_l}{10} - \frac{1}{||\boldsymbol{\delta}||_{\Sigma}^2} \mathcal{T}(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) > c'' \gamma^k r_l \sqrt{\frac{t + \log p}{n_0 g}} \bigg\} \\ & \leq \sum_{k=1}^N \mathbb{P}\bigg\{ \sup_{||\boldsymbol{\delta}||_1 / ||\boldsymbol{\delta}||_{\Sigma} \leq \gamma^k r_l} \frac{\kappa_l f_l}{10} - \frac{1}{||\boldsymbol{\delta}||_{\Sigma}^2} \mathcal{T}(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) > c'' \gamma^k r_l \sqrt{\frac{t + \log p}{n_0 g}} \bigg\} \\ & \leq \sum_{k=1}^N \mathbb{P}\bigg\{ \sup_{||\boldsymbol{\delta}||_1 / ||\boldsymbol{\delta}||_{\Sigma} \leq \gamma^k r_l} \frac{\kappa_l f_l}{10} - \frac{1}{||\boldsymbol{\delta}||_{\Sigma}^2} \mathcal{T}(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) > c'' \gamma^k r_l \sqrt{\frac{t + \log p}{n_0 g}} \bigg\} \\ & \leq \sum_{k=1}^N p^{-1} e^{-t} \leq \lceil \log(l/r_l) / \log(\gamma) \rceil p^{-1} e^{-t}. \end{aligned}$$

Taking  $\gamma = e^{1/e}$  and  $t = \log\{e \log(l/r_l)\} + u$  yields that with probability at least  $1 - p^{-1}e^{-u}$ ,

$$\frac{\kappa_l f_l}{10} - \frac{1}{||\boldsymbol{\delta}||_{\Sigma}^2} \mathcal{T}(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) \leq \frac{c'' ||\boldsymbol{\delta}||_1}{||\boldsymbol{\delta}||_{\Sigma}} \sqrt{\frac{\log p + \log\{e \log(l/r_l)\} + u}{n_0 g}}$$

Multiplying by  $||\boldsymbol{\delta}||_{\Sigma}^2$  on both sides yields

$$\mathcal{T}(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) \geq \frac{\kappa_l f_l}{10} ||\boldsymbol{\delta}||_{\Sigma}^2 - c'' ||\boldsymbol{\delta}||_1 ||\boldsymbol{\delta}||_{\Sigma} \sqrt{\frac{\log p + \log\{e \log(l/r_l)\} + u}{n_0 g}}$$

## A.4 Proof of Theorem 3.1

**Transferring step:** Let  $\omega^*$  be the true parameter of the transferring step and S be the active set of  $\beta^*$  with cardinality s. The symmetric Bregman divergence between  $\hat{\omega}^{\mathcal{A}_m}$  and  $\omega^*$  is defined as

$$\langle \nabla \hat{Q}_h(\hat{\omega}^{\mathcal{A}_m}) - \nabla \hat{Q}_h(\omega^*), \hat{\omega}^{\mathcal{A}_m} - \omega^* \rangle \ge 0.$$
 (37)

Let  $\hat{\Delta} = \hat{\omega}^{\mathcal{A}_m} - \omega^*$ . By optimality, there exists a subgradient  $\hat{g} \in \partial \sum_{i=1}^n q_{\lambda_\omega}(|\omega_i|)$ , such that

$$\nabla \hat{Q}_h(\hat{\boldsymbol{\omega}}^{\mathcal{A}_m}) + \lambda_\omega \hat{g} = 0.$$

Then (37) is equivalent to

$$-\langle \nabla \hat{Q}_h(\boldsymbol{\omega}^*), \hat{\Delta} \rangle - \lambda_{\boldsymbol{\omega}} \langle \hat{g}, \hat{\Delta} \rangle \ge 0.$$
(38)

Note that

$$\begin{split} \langle \hat{g}, \boldsymbol{\omega}^{*} - \hat{\boldsymbol{\omega}}^{\mathcal{A}_{m}} \rangle &\leq ||\boldsymbol{\omega}^{*}||_{1} - ||\hat{\boldsymbol{\omega}}^{\mathcal{A}_{m}}||_{1} = ||\boldsymbol{\omega}^{*}_{\mathcal{S}}||_{1} + ||\boldsymbol{\omega}^{*}_{\mathcal{S}^{c}}||_{1} - ||\hat{\Delta} + \boldsymbol{\omega}^{*}||_{1} \\ &= ||\boldsymbol{\omega}^{*}_{\mathcal{S}}||_{1} + ||\boldsymbol{\omega}^{*}_{\mathcal{S}^{c}}||_{1} - ||\hat{\Delta}_{\mathcal{S}} + \boldsymbol{\omega}^{*}_{\mathcal{S}}||_{1} - ||\hat{\Delta}_{\mathcal{S}^{c}} + \boldsymbol{\omega}^{*}_{\mathcal{S}^{c}}||_{1} \\ &\leq ||\hat{\Delta}_{\mathcal{S}}||_{1} - ||\hat{\Delta}_{\mathcal{S}^{c}}||_{1} + 2||\boldsymbol{\omega}^{*}_{\mathcal{S}^{c}}||_{1}. \end{split}$$

Then,

$$\begin{aligned} \langle \nabla \hat{Q}_{h}(\hat{\omega}^{\mathcal{A}_{m}}) - \nabla \hat{Q}_{h}(\omega^{*}), \hat{\Delta} \rangle \\ &= \lambda_{\omega} \langle \hat{g}, \omega^{*} - \hat{\omega}^{\mathcal{A}_{m}} \rangle + \langle \nabla \hat{Q}_{h}(\omega^{*}) - \nabla Q_{h}(\omega^{*}), \omega^{*} - \hat{\omega}^{\mathcal{A}_{m}} \rangle + \langle \nabla Q_{h}(\omega^{*}), \omega^{*} - \hat{\omega}^{\mathcal{A}_{m}} \rangle \\ &\leq \lambda_{\omega} \Big( ||\hat{\Delta}_{\mathcal{S}}||_{1} - ||\hat{\Delta}_{\mathcal{S}^{c}}||_{1} + 2||\omega^{*}_{\mathcal{S}^{c}}||_{1} \Big) + \underbrace{||\nabla \hat{Q}_{h}(\omega^{*}) - \nabla Q_{h}(\omega^{*})||_{\infty}}_{||\pi^{*}_{h}||_{\infty}} ||\hat{\Delta}||_{1} \\ &+ \underbrace{||\Sigma^{-1/2} \nabla Q_{h}(\omega^{*})||_{2}}_{b^{*}_{h}} ||\hat{\Delta}||_{\Sigma}. \end{aligned} \tag{39}$$

Conditioned on the event  $\{\lambda_{\omega} \geq 2 || \boldsymbol{\pi}_{h}^{*} ||_{\infty}\}$ , (39) becomes

$$\langle \nabla \hat{Q}_h(\hat{\boldsymbol{\omega}}^{\mathcal{A}_m}) - \nabla \hat{Q}_h(\boldsymbol{\omega}^*), \hat{\Delta} \rangle \leq \lambda_{\boldsymbol{\omega}} \Big( ||\hat{\Delta}_{\mathcal{S}}||_1 - ||\hat{\Delta}_{\mathcal{S}^c}||_1 + 2||\boldsymbol{\omega}_{\mathcal{S}^c}^*||_1 \Big) + \frac{\lambda_{\boldsymbol{\omega}}}{2} ||\hat{\Delta}||_1 + b_h^* ||\hat{\Delta}||_{\Sigma}$$

Lemma A.1 implies  $||\boldsymbol{\omega}_{\mathcal{S}^c}^*||_1 \leq C_1 m$ , so we have

$$\langle \nabla \hat{Q}_h(\hat{\omega}^{\mathcal{A}_m}) - \nabla \hat{Q}_h(\omega^*), \hat{\Delta} \rangle \leq \frac{3}{2} \lambda_{\omega} ||\hat{\Delta}_{\mathcal{S}}||_1 - \frac{1}{2} \lambda_{\omega} ||\hat{\Delta}_{\mathcal{S}^c}||_1 + 2\lambda_{\omega} C_1 m + b_h^* ||\hat{\Delta}||_{\Sigma}.$$

Since  $\langle \nabla \hat{Q}_h(\hat{\omega}^{\mathcal{A}_m}) - \nabla \hat{Q}_h(\omega^*), \hat{\Delta} \rangle \ge 0$ ,  $\hat{\Delta}$  satisfies the constraint  $||\hat{\Delta}_{\mathcal{S}^c}||_1 \le 3||\hat{\Delta}_{\mathcal{S}}||_1 + 4C_1m + 2\lambda_{\omega}^{-1}b_h^*||\hat{\Delta}||_{\Sigma}$ , from which it follows that

$$|\hat{\Delta}||_{1} \le 4s^{1/2} ||\hat{\Delta}||_{2} + 4C_{1}m + 2\lambda_{\omega}^{-1}b_{h}^{*}||\hat{\Delta}||_{\Sigma}.$$
(40)

Now we claim that when  $\lambda_{\omega} \geq 2||\boldsymbol{\pi}_{h}^{*}||_{\infty}$ , with probability at least  $1 - p^{-1}$ , it holds that

$$||\hat{\Delta}||_{\Sigma} \le \frac{8\phi_2 C_1 m}{\phi_1} \sqrt{\frac{\log p + \log(n_{\mathcal{A}_m} + n_0)}{n_{\mathcal{A}_m} + n_0}} + \frac{3\lambda_{\omega}\gamma_p^{-1/2}s^{1/2} + 2b_h^*}{\phi_1} + 2\sqrt{\frac{C_1\lambda_{\omega}m}{\phi_1}}.$$
 (41)

If the claim does not hold, consider  $\mathbb{C} = \{\Delta : 1.5\lambda_{\omega} ||\Delta_{\mathcal{S}}||_1 - 0.5\lambda_{\omega} ||\Delta_{\mathcal{S}^c}||_1 + 2\lambda_{\omega}C_1m + b_h^*||\Delta||_{\Sigma} \ge 0\}$ . For any  $t \in (0, 1)$ ,

$$\begin{aligned} \frac{1}{2}\lambda_{\boldsymbol{\omega}}||t\hat{\Delta}_{\mathcal{S}^{c}}||_{1} &= t \cdot \frac{1}{2}\lambda_{\boldsymbol{\omega}}||\hat{\Delta}_{\mathcal{S}^{c}}||_{1} \leq t \cdot \left(\frac{3}{2}\lambda_{\boldsymbol{\omega}}||\hat{\Delta}_{\mathcal{S}}||_{1} + 2\lambda_{\boldsymbol{\omega}}C_{1}m + b_{h}^{*}||\hat{\Delta}||_{\Sigma}\right) \\ &\leq \frac{3}{2}\lambda_{\boldsymbol{\omega}}||t\hat{\Delta}_{\mathcal{S}}||_{1} + 2\lambda_{\boldsymbol{\omega}}C_{1}m + b_{h}^{*}||t\hat{\Delta}||_{\Sigma}, \end{aligned}$$

which implies that  $t\hat{\Delta} \in \mathbb{C}$ . We could find some t satisfying that  $||t\hat{\Delta}||_{\Sigma} \leq 1$  and

$$|t\hat{\Delta}||_{\Sigma} > \frac{8\phi_2 C_1 m}{\phi_1} \sqrt{\frac{\log p + \log(n_{\mathcal{A}_m} + n_0)}{n_{\mathcal{A}_m} + n_0}} + \frac{3\lambda_{\omega}\gamma_p^{-1/2}s^{1/2} + 2b_h^*}{\phi_1} + 2\sqrt{\frac{C_1\lambda_{\omega}m}{\phi_1}}.$$

Denote  $\tilde{\Delta} = t\hat{\Delta}$  and  $F(\Delta) = \hat{Q}_h(\omega^* + \Delta) - \hat{Q}_h(\omega^*) + \lambda_{\omega}(||\omega^* + \Delta||_1 - ||\omega^*||_1)$ . Since  $F(\mathbf{0}) = 0$  and  $F(\hat{\Delta}) \le 0$ , by convexity,

$$F(\tilde{\Delta}) = F(t\hat{\Delta} + (1-t)\mathbf{0}) \le tF(\hat{\Delta}) \le 0.$$

However,

$$\begin{split} F(\tilde{\Delta}) &= \hat{D}_{h}(\tilde{\Delta}) - \lambda_{\omega} ||\boldsymbol{\omega}^{*}||_{1} + \lambda_{\omega} ||\boldsymbol{\omega}^{*} + \tilde{\Delta}||_{1} \\ &= \hat{R}_{h}(\tilde{\Delta}) + \langle \nabla \hat{Q}_{h}(\boldsymbol{\omega}^{*}), \tilde{\Delta} \rangle - \lambda_{\omega} ||\boldsymbol{\omega}^{*}||_{1} + \lambda_{\omega} ||\boldsymbol{\omega}^{*} + \tilde{\Delta}||_{1} \\ &= \hat{R}_{h}(\tilde{\Delta}) - (\lambda_{\omega} ||\boldsymbol{\omega}^{*}||_{1} - \lambda_{\omega} ||\boldsymbol{\omega}^{*} + \tilde{\Delta}||_{1} - \langle \nabla \hat{Q}_{h}(\boldsymbol{\omega}^{*}) - \nabla Q_{h}(\boldsymbol{\omega}^{*}), \tilde{\Delta} \rangle \\ &- \langle \nabla Q_{h}(\boldsymbol{\omega}^{*}), \tilde{\Delta} \rangle ). \end{split}$$

Then by Proposition 3.1 and (39),

$$\begin{split} F(\tilde{\Delta}) &\geq \phi_1 ||\tilde{\Delta}||_{\Sigma}^2 - \phi_2 \sqrt{\frac{\log p + \log n}{nh}} ||\tilde{\Delta}||_1 ||\tilde{\Delta}||_{\Sigma} - \frac{3}{2} \lambda_{\omega}||\tilde{\Delta}_{\mathcal{S}}||_1 + \frac{1}{2} \lambda_{\omega}||\tilde{\Delta}_{\mathcal{S}^c}||_1 \\ &- 2\lambda_{\omega}C_1m - b_h^* ||\tilde{\Delta}||_{\Sigma} \\ &\geq \phi_1 ||\tilde{\Delta}||_{\Sigma}^2 - \phi_2 \sqrt{\frac{\log p + \log n}{nh}} ||\tilde{\Delta}||_1 ||\tilde{\Delta}||_{\Sigma} - \frac{3}{2} \lambda_{\omega}||\tilde{\Delta}_{\mathcal{S}}||_1 - 2\lambda_{\omega}C_1m - b_h^*||\tilde{\Delta}||_{\Sigma} \end{split}$$

Note that  $||\tilde{\Delta}_{\mathcal{S}}||_1 \leq s^{1/2} ||\tilde{\Delta}||_2 \leq \gamma_p^{-1/2} s^{1/2} ||\tilde{\Delta}||_{\Sigma}$  and (40). Therefore, when

$$(n_{\mathcal{A}_m} + n_0)h > 16\phi_1^{-2}\phi_2^2(\log p + \log n)\max\{16s\gamma_p^{-1}, 4\lambda_\omega^{-2}(b_h^*)^2\},\$$

we have  $\phi_2 \sqrt{(\log p + \log n)/(nh)} (4\sqrt{s\gamma_p^{-1/2}} + 2\lambda_{\omega}^{-1}b_h^*) \le \phi_1/2$ . Then it follows that,

$$F(\tilde{\Delta}) \geq \frac{1}{2}\phi_1 ||\tilde{\Delta}||_{\Sigma}^2 - \left(4\phi_2\sqrt{\frac{\log p + \log n}{nh}}C_1m + \frac{3}{2}\lambda_{\omega}\gamma_p^{-1/2}s^{1/2} + b_h^*\right)||\tilde{\Delta}||_{\Sigma} - 2\lambda_{\omega}C_1m$$
  
> 0,

which contradicts with  $F(\tilde{\Delta}) \leq 0$ . Therefore the claim holds.

It remains to control the probability of the event  $\{\lambda_{\omega} \geq ||\boldsymbol{\pi}_{h}^{*}||_{\infty}\}$  and the probability of the local RSC condition. By Lemma A.2, we pick

$$\lambda_{\omega} = 2 \left[ \sigma \sqrt{\{\tau(1-\tau) + Ch^2\} \frac{4\log(2p)}{n_{\mathcal{A}_m} + n_0}} + \max(1-\tau,\tau) \frac{2\log(2p)}{n_{\mathcal{A}_m} + n_0} \right],$$

so that  $\{\lambda_{\omega} \geq 2 || \pi_h^* ||_{\infty}\}$ . From Lemma A.3, we have  $b_h^* \leq Ch^2$ . Now with probability at least  $1 - (pn)^{-1}$ ,

$$||\hat{\Delta}||_{\Sigma} \lesssim m\sqrt{\frac{\log p + \log(n_{\mathcal{A}_m} + n_0)}{n_{\mathcal{A}_m} + n_0}} + \sqrt{\frac{s\log p}{n_{\mathcal{A}_m} + n_0}} + h^2 + \left(\frac{\log p}{n_{\mathcal{A}_m} + n_0}\right)^{\frac{1}{4}}\sqrt{m}.$$

We then let  $h^2 \leq s^{1/2} \lambda_{\omega}$ , so that

$$||\hat{\Delta}||_{\Sigma} \lesssim m\sqrt{\frac{\log p + \log(n_{\mathcal{A}_m} + n_0)}{n_{\mathcal{A}_m} + n_0}} + \sqrt{\frac{s\log p}{n_{\mathcal{A}_m} + n_0}} + \left(\frac{\log p}{n_{\mathcal{A}_m} + n_0}\right)^{\frac{1}{4}}\sqrt{m},\tag{42}$$

with probability at least  $1 - (pn)^{-1}$ .

Note that

$$\begin{split} ||\hat{\Delta}||_{1} &\leq 4 ||\hat{\Delta}_{\mathcal{S}}||_{1} + 4C_{1}m + 2\lambda_{\omega}^{-1}b_{h}^{*}||\hat{\Delta}||_{\Sigma} \\ &\leq 4\sqrt{s}||\hat{\Delta}||_{2} + 4C_{1}m + 2\lambda_{\omega}^{-1}b_{h}^{*}||\hat{\Delta}||_{\Sigma} \\ &\leq (4\sqrt{s}\gamma_{p}^{-1/2} + l_{0}\kappa_{2}\lambda_{\omega}^{-1}h^{2})||\hat{\Delta}||_{\Sigma} + 4C_{1}m, \end{split}$$

which encloses

$$||\hat{\Delta}||_1 \lesssim m\sqrt{\frac{s\log(p) + s\log(n_{\mathcal{A}_m} + n_0)}{n_{\mathcal{A}_m} + n_0}} + s\sqrt{\frac{\log p}{n_{\mathcal{A}_m} + n_0}} + \left(\frac{\log p}{n_{\mathcal{A}_m} + n_0}\right)^{\frac{1}{4}}\sqrt{sm} + m,$$

with probability at least  $1 - (pn)^{-1}$ . **Debiasing step:** Denote  $\delta^* = \beta^* - \omega^*$ ,  $\hat{\delta}^{\mathcal{A}_m} = \hat{\beta} - \hat{\omega}^{\mathcal{A}_m}$  and  $\hat{v}^{\mathcal{A}_m} = \hat{\delta}^{\mathcal{A}_m} - \delta^*$ . Similar to (39), we have

$$\begin{split} \langle \nabla \hat{Q}_{g}^{(0)}(\hat{\omega}^{\mathcal{A}_{m}} + \hat{\delta}^{\mathcal{A}_{m}}) - \nabla \hat{Q}_{g}^{(0)}(\beta^{*}), \hat{\beta} - \beta^{*} \rangle \\ &= \lambda_{\delta} \langle \hat{f}, \beta^{*} - \hat{\beta} \rangle + \langle \nabla \hat{Q}_{g}^{(0)}(\beta^{*}) - \nabla Q_{g}^{(0)}(\beta^{*}), \beta^{*} - \hat{\beta} \rangle + \langle \nabla Q_{g}^{(0)}(\beta^{*}), \beta^{*} - \hat{\beta} \rangle \\ &\leq \lambda_{\delta} \big( ||\beta^{*} - \hat{\omega}^{\mathcal{A}_{m}}||_{1} - ||\hat{\beta} - \hat{\omega}^{\mathcal{A}_{m}}||_{1} \big) + \underbrace{\|\nabla \hat{Q}_{g}^{(0)}(\beta^{*}) - \nabla Q_{g}^{(0)}(\beta^{*})\|_{\infty}}_{||\pi_{g}^{*}||_{\infty}} \|\hat{\beta} - \beta^{*}\|_{1} \\ &+ \underbrace{\|\sum^{-1/2} \nabla Q_{g}^{(0)}(\beta^{*})\|_{2}}_{b_{g}^{*}} \|\hat{\beta} - \beta^{*}\|_{\Sigma} \\ &\leq \frac{3}{2} \lambda_{\delta} ||\beta^{*} - \hat{\omega}^{\mathcal{A}_{m}}||_{1} - \frac{1}{2} \lambda_{\delta} ||\hat{\beta} - \hat{\omega}^{\mathcal{A}_{m}}||_{1} + b_{g}^{*}\|\hat{\beta} - \beta^{*}\|_{\Sigma} \\ &\leq \frac{3}{2} \lambda_{\delta} ||\beta^{*} - \omega^{*}||_{1} + \frac{3}{2} \lambda_{\delta} ||\hat{\Delta}||_{1} - \frac{1}{2} \lambda_{\delta} ||\hat{\beta} - \hat{\omega}^{\mathcal{A}_{m}}||_{1} + b_{g}^{*}\|\hat{\beta} - \beta^{*}\|_{\Sigma} \\ &\leq \frac{3}{2} \lambda_{\delta} Cm + \frac{3}{2} \lambda_{\delta} ||\hat{\Delta}||_{1} - \frac{1}{2} \lambda_{\delta} ||\hat{\beta} - \hat{\omega}^{\mathcal{A}_{m}}||_{1} + b_{g}^{*}\|\hat{\beta} - \beta^{*}\|_{\Sigma}. \end{split}$$

$$\tag{43}$$

On the other hand,

$$\begin{split} \langle \nabla \hat{Q}_{g}^{(0)}(\hat{\omega}^{\mathcal{A}_{m}} + \hat{\delta}^{\mathcal{A}_{m}}) - \nabla \hat{Q}_{g}^{(0)}(\beta^{*}), \hat{\beta} - \beta^{*} \rangle \\ &\leq \lambda_{\delta} (||\beta^{*} - \hat{\omega}^{\mathcal{A}_{m}}||_{1} - ||\hat{\beta} - \hat{\omega}^{\mathcal{A}_{m}}||_{1}) + \frac{\lambda_{\delta}}{2} ||\hat{\beta} - \beta^{*}||_{1} + b_{g}^{*}||\hat{\beta} - \beta^{*}||_{\Sigma} \\ &\leq \lambda_{\delta} ||\beta^{*}_{S} - \hat{\omega}^{\mathcal{A}_{m}}_{S}||_{1} + \lambda_{\delta}||\beta^{*}_{S^{c}} - \hat{\omega}^{\mathcal{A}_{m}}_{S^{c}}||_{1} - \lambda_{\delta}||\hat{\beta}_{S} - \hat{\omega}^{\mathcal{A}_{m}}_{S^{c}}||_{1} - \lambda_{\delta}||\hat{\beta}_{S^{c}} - \hat{\omega}^{\mathcal{A}_{m}}_{S^{c}}||_{1} \\ &+ \frac{\lambda_{\delta}}{2} ||\hat{\beta} - \beta^{*}||_{1} + b_{g}^{*}||\hat{\beta} - \beta^{*}||_{\Sigma} \\ &\leq \lambda_{\delta} (||\beta^{*}_{S} - \hat{\omega}^{\mathcal{A}_{m}}_{S}||_{1} - ||\hat{\beta}_{S} - \hat{\omega}^{\mathcal{A}_{m}}_{S}||_{1}) - \lambda_{\delta} (||\beta^{*}_{S^{c}} - \hat{\omega}^{\mathcal{A}_{m}}_{S^{c}}||_{1} + ||\hat{\beta}_{S^{c}} - \hat{\omega}^{\mathcal{A}_{m}}_{S^{c}}||_{1}) \\ &+ 2\lambda_{\delta} ||\beta^{*}_{S} - \hat{\omega}^{\mathcal{A}_{m}}_{S}||_{1} + \frac{\lambda_{\delta}}{2} ||\hat{\beta} - \beta^{*}||_{1} + b_{g}^{*}||\hat{\beta} - \beta^{*}||_{\Sigma} \\ &\leq \lambda_{\delta} ||\beta^{*}_{S} - \hat{\beta}_{S}||_{1} - \lambda_{\delta} ||\beta^{*}_{S^{c}} - \hat{\beta}_{S^{c}}||_{1} + 2\lambda_{\delta} ||\beta^{*}_{S^{c}} - \omega^{*}_{S^{c}}||_{1} + 2\lambda_{\delta} ||\hat{\Delta}_{S^{c}}||_{1} \\ &+ \frac{\lambda_{\delta}}{2} ||\hat{\beta} - \beta^{*}||_{1} + b_{g}^{*}||\hat{\beta} - \beta^{*}||_{\Sigma} \\ &\leq \frac{3}{2}\lambda_{\delta} ||\beta^{*}_{S} - \hat{\beta}_{S}||_{1} - \frac{1}{2}\lambda_{\delta} ||\beta^{*}_{S^{c}} - \hat{\beta}_{S^{c}}||_{1} + 2\lambda_{\delta} C_{1}m + 2\lambda_{\delta} ||\hat{\Delta}_{S^{c}}||_{1} + b_{g}^{*}||\hat{\beta} - \beta^{*}||_{\Sigma} \\ &\leq \frac{3}{2}\lambda_{\delta} ||\beta^{*}_{S} - \hat{\beta}_{S}||_{1} - \frac{1}{2}\lambda_{\delta} ||\beta^{*}_{S^{c}} - \hat{\beta}_{S^{c}}||_{1} + 2\lambda_{\delta} C_{1}m + b_{g}^{*}||\hat{\beta} - \beta^{*}||_{\Sigma} \\ &\leq \frac{3}{2}\lambda_{\delta} ||\beta^{*}_{S} - \hat{\beta}_{S}||_{1} - \frac{1}{2}\lambda_{\delta} ||\beta^{*}_{S^{c}} - \hat{\beta}_{S^{c}}||_{1} + 2\lambda_{\delta} C_{1}m + b_{g}^{*}||\hat{\beta} - \beta^{*}||_{\Sigma} \\ &+ 2\lambda_{\delta} \left(m\sqrt{\frac{\sin(p) + \sin(p(n_{\mathcal{A}_{m} + n_{0})}{n_{\mathcal{A}_{m} + n_{0}}}} + s\sqrt{\frac{\log p}{n_{\mathcal{A}_{m} + n_{0}}}} + \left(\frac{\log p}{n_{\mathcal{A}_{m} + n_{0}}}\right)^{\frac{1}{4}}\sqrt{sm} + m\right).$$

Thus

$$|\boldsymbol{\beta}^{*} - \hat{\boldsymbol{\beta}}||_{1} \leq 4\gamma_{p}^{-1/2}\sqrt{s}||\boldsymbol{\beta}^{*} - \hat{\boldsymbol{\beta}}||_{\Sigma} + 2b_{g}^{*}\lambda_{\delta}^{-1}||\boldsymbol{\beta}^{*} - \hat{\boldsymbol{\beta}}||_{\Sigma} + 4\left(m\sqrt{\frac{s\log(p) + s\log(n_{\mathcal{A}_{m}} + n_{0})}{n_{\mathcal{A}_{m}} + n_{0}}} + s\sqrt{\frac{\log p}{n_{\mathcal{A}_{m}} + n_{0}}} + \left(\frac{\log p}{n_{\mathcal{A}_{m}} + n_{0}}\right)^{\frac{1}{4}}\sqrt{sm} + m\right)$$
(45)

Set  $r = g/(48c_0)$  and  $R = (4\gamma_p^{-1/2}\sqrt{s} + 2b_g^*\lambda_{\delta}^{-1})r + 4C\sqrt{s}$ , if  $m \leq C\sqrt{s}$  for some positive constant C and  $n_{\mathcal{A}_m} + n_0 \geq s \log p$ . Denote  $\Theta(r, R) = \mathbb{B}_{\Sigma}(r) \cap \mathbb{C}_{\Sigma}(R)$  and  $\tilde{\beta} = (1 - \eta)\beta^* + \eta\hat{\beta}$ , where  $\eta = \sup\{u \in [0, 1] : \beta^* + u(\hat{\beta} - \beta^*) \in \beta^* + \Theta(r, R)\}$ . If  $\hat{\beta} \notin \Theta(r, R)$ , then  $\eta \in (0, 1)$  and  $\tilde{\beta}$  falls onto the boundary of  $\Theta(r, R)$ ; otherwise  $\tilde{\beta} = \hat{\beta}$ .

Combining (43) and Proposition 3.2, we have

$$\frac{\alpha_1}{2}||\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*||_{\Sigma}^2 - \alpha_2 \cdot \frac{\log p + \log n_0}{n_0 g}||\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*||_1^2 \le \frac{3}{2}\lambda_\delta Cm + \frac{3}{2}\lambda_\delta ||\tilde{\boldsymbol{\Delta}}||_1 + b_g^* \left\|\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\right\|_{\Sigma}.$$
(46)

Besides, (43) implies

$$||\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\omega}}^{\mathcal{A}_m}||_1 \leq 3Cm + 3||\hat{\Delta}||_1 + \frac{2b_g^*}{\lambda_{\delta}} \|\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_{\Sigma}.$$

As a result,

$$\begin{aligned} ||\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*||_1 &\leq ||\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\omega}}^{\mathcal{A}_m}||_1 + ||\hat{\boldsymbol{\omega}}^{\mathcal{A}_m} - \boldsymbol{\beta}^*||_1 \\ &\leq 4Cm + 4||\hat{\boldsymbol{\Delta}}||_1 + \frac{2b_g^*}{\lambda_{\delta}} \|\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_{\Sigma}. \end{aligned}$$

Let  $\alpha = \alpha_1 - 4b_g^{*2}\lambda_{\delta}^{-2}$ , then (46) becomes

$$\begin{split} \alpha_1 ||\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*||_{\Sigma}^2 &\leq 2\alpha_2 \cdot \frac{\log p + \log n_0}{n_0 g} (16Cm^2 + 16||\hat{\boldsymbol{\Delta}}||_1^2 + 4b_g^{*2}\lambda_{\delta}^{-2} \left\|\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\right\|_{\Sigma}^2) \\ &\quad + 3\lambda_{\delta}Cm + 3\lambda_{\delta}||\hat{\boldsymbol{\Delta}}||_1 + 2b_g^* \left\|\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\right\|_{\Sigma} \\ \alpha ||\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*||_{\Sigma}^2 - 2b_g^* \left\|\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\right\|_{\Sigma} + \frac{b_g^{*2}}{\alpha} &\leq \frac{\log p}{n_0 g} (m^2 + ||\hat{\boldsymbol{\Delta}}||_1^2) + \lambda_{\delta}m + \lambda_{\delta}||\hat{\boldsymbol{\Delta}}||_1 \\ \alpha \left(||\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*||_{\Sigma} - \frac{b_g^*}{\alpha}\right)^2 &\leq \frac{\log p}{n_0 g} \left(m^2 + ||\hat{\boldsymbol{\Delta}}||_1^2\right) + \lambda_{\delta}m + \lambda_{\delta}||\hat{\boldsymbol{\Delta}}||_1. \end{split}$$

Thus,

$$||\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*||_{\Sigma} \lesssim \sqrt{\frac{\log p}{n_0 g}} (m + ||\hat{\boldsymbol{\Delta}}||_1) + \sqrt{\lambda_\delta m} + \sqrt{\lambda_\delta} ||\hat{\boldsymbol{\Delta}}||_1 + b_g^*.$$

Let  $\lambda_{\delta} = C\sqrt{\log(p)/n_0}$  and  $g \asymp (\log(p)/n_0)^{1/4}$ , then

$$\begin{split} ||\tilde{\beta} - \beta^*||_{\Sigma} \lesssim \left(\frac{\log p}{n_0}\right)^{3/8} ||\hat{\Delta}||_1 + \sqrt{m} \left(\frac{\log p}{n_0}\right)^{1/4} + \left(\frac{\log p}{n_0}\right)^{1/4} \sqrt{||\hat{\Delta}||_1} + \left(\frac{\log p}{n_0}\right)^{1/2} \\ \lesssim m \left(\frac{\log p}{n_0}\right)^{3/8} + s \left(\frac{\log p}{n_0}\right)^{3/8} \sqrt{\frac{\log p}{n_{\mathcal{A}_m} + n_0}} \\ + \sqrt{sm} \left(\frac{\log p}{n_0}\right)^{3/8} \left(\frac{\log p}{n_{\mathcal{A}_m} + n_0}\right)^{1/4} + \sqrt{m} \left(\frac{\log p}{n_0}\right)^{1/4} \\ + \sqrt{s} \left(\frac{\log p}{n_0}\right)^{1/4} \left(\frac{\log p}{n_{\mathcal{A}_m} + n_0}\right)^{1/4} + (sm)^{1/4} \left(\frac{\log p}{n_0}\right)^{1/4} \left(\frac{\log p}{n_{\mathcal{A}_m} + n_0}\right)^{1/8}. \end{split}$$

If  $n_0 > s^2 \log p$ ,  $\hat{\beta}$  falls in the interior of  $\Theta(r, R)$ , so we must have  $\hat{\beta} \in \Theta(r, R)$ . Consequently,  $\hat{\beta} = \tilde{\beta}$  satisfies the claimed bound,

$$||\hat{\beta} - \beta^*||_{\Sigma} \lesssim \sqrt{m} \left(\frac{\log p}{n_0}\right)^{1/4} + \sqrt{s} \left(\frac{\log p}{n_0}\right)^{1/4} \left(\frac{\log p}{n_{\mathcal{A}_m} + n_0}\right)^{1/4}.$$
(47)

In addition, if  $m \leq s\sqrt{\log(p)/n_0}$ , the above upper bound is sharper than  $\sqrt{s\log(p)/n_0}$ . Then by (45), we have

$$||\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*||_1 \lesssim s \sqrt{\frac{\log p}{n_{\mathcal{A}_m} + n_0}} + \left(\frac{\log p}{n_{\mathcal{A}_m} + n_0}\right)^{\frac{1}{4}} \sqrt{sm} + m.$$

## A.5 Proof of Proposition 3.3

The method is similar to the proof of Proposition 3.2. At first, the divergence is given by

$$D(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) = \langle \nabla \hat{Q}_g^{(0)}(\boldsymbol{\beta}_1) - \nabla \hat{Q}_g^{(0)}(\boldsymbol{\beta}_2), \boldsymbol{\beta}_1 - \boldsymbol{\beta}_2 \rangle$$

For a given kernel function  $K(\cdot)$  and bandwidth g > 0, the smoothed quantile loss  $\hat{Q}_g^{(0)}$  can be written as  $(n_0g)^{-1}\sum_{i=1}^{n_0}\int_{-\infty}^{\infty}\rho_{\tau}(u)K\{(u+\langle \boldsymbol{x}_i^{(0)},\boldsymbol{\beta}\rangle-y_i^{(0)})/g\}du$ . Therefore

$$D(oldsymbol{eta}_1,oldsymbol{eta}_2) \geq rac{\kappa_l}{n_0g}\sum_{i=1}^{n_0} \langle oldsymbol{x}_i^{(0)},oldsymbol{eta}_1-oldsymbol{eta}_2
angle^2 \mathbbm{1}_{\mathcal{E}_i},$$

where the event  $\mathcal{E}_i$  is defined by,

$$\mathcal{E}_{i} = \{ |\epsilon_{i}| \leq g/4 \} \cap \{ |\langle \boldsymbol{x}_{i}^{(0)}, \boldsymbol{\beta}_{1} - \boldsymbol{\beta}_{2} \rangle | \leq g || \boldsymbol{\beta}_{1} - \boldsymbol{\beta}_{2} ||_{\Sigma} / (2r) \} \cap \{ |\langle \boldsymbol{x}_{i}^{(0)}, \boldsymbol{\beta}_{1} - \boldsymbol{\beta}^{*} \rangle | \leq g/4 \}.$$

for all  $\beta_1 - \beta_2 \in \mathbb{B}_{\Sigma}(r)$  and  $\kappa_l = \min_{|u| \leq 1} K(u)$ . For a truncation level R > 0, define functions  $\varphi_R(u)$ and  $\psi_R(u)$  as previous proof. By this construction,  $\varphi_R(u) \leq u^2 \cdot \mathbb{1}\{|u| \leq R\}, \ \varphi_{cR}(cu) = c^2 \varphi_R(u), \ \varphi_R$  is R-Lipschitz,  $\psi_R$  is (2/R)-Lipschitz and  $\psi_R(u) \leq \mathbb{1}\{|u| \leq R\}$ . From these two new-defined function, we have

$$D(\boldsymbol{\beta}_{1},\boldsymbol{\beta}_{2}) \geq \frac{\kappa_{l}}{n_{0}g} ||\boldsymbol{\delta}||_{\Sigma}^{2} \sum_{i=1}^{n_{0}} \mathbb{1}\{|\epsilon_{i}| \leq g/4\} \varphi_{g||\boldsymbol{\delta}||_{\Sigma}/(2r)}(\langle \boldsymbol{x}_{i}^{(0)},\boldsymbol{\delta} \rangle) \psi_{g/4}(\langle \boldsymbol{x}_{i}^{(0)},\boldsymbol{\beta}_{1}-\boldsymbol{\beta}^{*} \rangle)$$

$$\geq \kappa_{l} ||\boldsymbol{\delta}||_{\Sigma}^{2} \cdot \underbrace{\frac{1}{n_{0}g} \sum_{i=1}^{n_{0}} \mathbb{1}\{|\epsilon_{i}| \leq g/4\} \varphi_{g/(2r)}(\langle \boldsymbol{x}_{i}^{(0)},\boldsymbol{\delta} \rangle/||\boldsymbol{\delta}||_{\Sigma}) \psi_{g/4}(\langle \boldsymbol{x}_{i}^{(0)},\boldsymbol{\beta}_{1}-\boldsymbol{\beta}^{*} \rangle), \qquad (48)$$

where  $\delta = \beta_1 - \beta_2$ . Finally, with a similar proof as Proposition 3.2, if  $r = g/(48c_0)$  and  $n_0g \gtrsim f_u f_l^{-2} \max\{s, l^2 \log p\}$ , then

$$D_0(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) \ge 0.2f_l,$$

with probability at least  $1 - (pn_0)^{-1}$ . Therefore,

$$D(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) \geq 0.2 f_l \kappa_l || \boldsymbol{\beta}_1 - \boldsymbol{\beta}_2 ||_{\Sigma}^2.$$

## A.6 Proof of Theorem 3.2

For step 1, the parameter  $\omega^{(k)}$  is at most s + m sparse. Therefore, similarly as Theorem 1 in Tan et al. (2022), we have

$$\|\hat{\boldsymbol{\omega}}^{(k)} - \boldsymbol{\omega}^{(k)}\|_2^2 \lesssim \frac{(s+m)\log p}{n^{(k)}}, \ \|\hat{\boldsymbol{\omega}}^{(k)} - \boldsymbol{\omega}^{(k)}\|_1 \lesssim (s+m)\sqrt{\frac{\log p}{n^{(k)}}}, \ k \in \mathcal{A}'_m$$

with probability at least  $1 - p^{-1}$ , provided that the bandwidth h satisfies

$$\max\left(\frac{\sigma_x}{f_l}\sqrt{\frac{(s+m)\log p}{n^{(k)}}}, \frac{\sigma_x^2 f_u}{f_l^2}\frac{(s+m)\log p}{n^{(k)}}\right) \lesssim h \le \min\{f_l/(2l_0), (s^{1/2}\lambda_{\omega}^{(k)})\},$$

where  $\sigma_x^2 = \max_{1 \ge j \ge p} \sigma_{jj}$ ,  $\sigma_{jj}$  are the diagonal elements of  $\Sigma$ . For step 2, denote  $\boldsymbol{\delta}^{(k)} = \boldsymbol{\beta}^* - \boldsymbol{\omega}^{(k)}$ ,  $\hat{\boldsymbol{\delta}}^{(k)} = \hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\omega}}^{(k)}$  and  $\hat{\boldsymbol{v}}^{(k)} = \hat{\boldsymbol{\delta}}^{(k)} - \boldsymbol{\delta}^{(k)}$ . For each  $k \in \mathcal{A}'_m$ ,

$$\langle \nabla \hat{Q}_{g}^{(0)}(\hat{\omega}^{(k)} + \hat{\delta}^{(k)}) - \nabla \hat{Q}_{g}^{(0)}(\hat{\omega}^{(k)} + \delta^{(k)}), \hat{v}^{(k)} \rangle$$

$$= \lambda_{\delta} \langle \hat{f}, \delta^{(k)} - \hat{\delta}^{(k)} \rangle + \langle \nabla \hat{Q}_{g}^{(0)}(\hat{\omega}^{(k)} + \delta^{(k)}) - \nabla \hat{Q}_{g}^{(0)}(\beta^{*}), \delta^{(k)} - \hat{\delta}^{(k)} \rangle$$

$$+ \langle \nabla \hat{Q}_{g}^{(0)}(\beta^{*}) - \nabla Q_{g}^{(0)}(\beta^{*}), \delta^{(k)} - \hat{\delta}^{(k)} \rangle + \langle \nabla Q_{g}^{(0)}(\beta^{*}), \delta^{(k)} - \hat{\delta}^{(k)} \rangle$$

$$\leq \lambda_{\delta} (||\delta^{*}||_{1} - ||\hat{\delta}^{(k)}||_{1}) + \langle \nabla \hat{Q}_{g}^{(0)}(\hat{\omega}^{(k)} + \delta^{*}) - \nabla \hat{Q}_{g}^{(0)}(\beta^{*}), \delta^{(k)} - \hat{\delta}^{(k)} \rangle$$

$$+ \underbrace{\| \nabla \hat{Q}_{g}^{(0)}(\beta^{*}) - \nabla Q_{g}^{(0)}(\beta^{*}) \|_{\infty}}_{||\pi_{g}^{*}||_{\infty}} \| \hat{v}^{(k)} \|_{1} + \underbrace{\| \Sigma^{-1/2} \nabla Q_{g}^{(0)}(\beta^{*}) \|_{2}}_{b_{g}^{*}} \| \hat{v}^{(k)} \|_{\Sigma}.$$

$$(49)$$

Then by Lemma A.4 with  $t = 2 \log p$  and some choice of  $(r_k, l_k)$ . For each  $k \in \mathcal{A}'_m$ , we let  $r_k = \sqrt{(s+m)\log(p)/n^{(k)}}$  and  $l = (s+m)\sqrt{\log(p)/n^{(k)}}$ . When  $\lambda_\delta \ge 2||\boldsymbol{\pi}_g^*||_{\infty}$ ,

$$\begin{split} \langle \nabla \hat{Q}_{g}^{(0)}(\hat{\boldsymbol{\omega}}^{(k)} + \hat{\boldsymbol{\delta}}^{(k)}) - \nabla \hat{Q}_{g}^{(0)}(\hat{\boldsymbol{\omega}}^{(k)} + \boldsymbol{\delta}^{(k)}), \hat{\boldsymbol{v}}^{(k)} \rangle \\ &\leq \lambda_{\delta} \big( ||\boldsymbol{\delta}^{(k)}||_{1} - ||\hat{\boldsymbol{\delta}}^{(k)}||_{1} \big) + \underbrace{c_{0} \bigg( \frac{s+m}{g} \sqrt{\frac{\log p}{n^{(k)}}} \sqrt{\frac{\log p}{n_{0}}} + \frac{\log p}{n_{0}} + \sqrt{s+m} \sqrt{\frac{\log p}{n^{(k)}}} \bigg)}_{C_{v}} \| \hat{\boldsymbol{v}}^{(k)} \|_{1} \\ &+ \frac{\lambda_{\delta}}{2} \| \hat{\boldsymbol{v}}^{(k)} \|_{1} + b_{g}^{*} \| \hat{\boldsymbol{v}}^{(k)} \|_{\Sigma}. \end{split}$$

Since

$$\|\boldsymbol{\delta}_{\mathcal{S}_{k}}^{(k)}\|_{1} - \|\hat{\boldsymbol{\delta}}_{\mathcal{S}_{k}}^{(k)}\|_{1} \leq \|(\boldsymbol{\delta}^{(k)} - \hat{\boldsymbol{\delta}}^{(k)})_{\mathcal{S}_{k}}\|_{1} \text{ and } \|\boldsymbol{\delta}_{\mathcal{S}_{k}^{c}}^{(k)}\|_{1} - \|\hat{\boldsymbol{\delta}}_{\mathcal{S}_{k}^{c}}^{(k)}\|_{1} = -\|(\hat{\boldsymbol{\delta}}^{(k)} - \boldsymbol{\delta}^{(k)})_{\mathcal{S}_{k}^{c}}\|_{1},$$

when  $C_v < \lambda_{\delta}/2$ , we obtain

$$\begin{split} \langle \nabla \hat{Q}_{g}^{(0)}(\hat{\boldsymbol{\omega}}^{(k)} + \hat{\boldsymbol{\delta}}^{(k)}) - \nabla \hat{Q}_{g}^{(0)}(\hat{\boldsymbol{\omega}}^{(k)} + \boldsymbol{\delta}^{(k)}), \hat{\boldsymbol{v}}^{(k)} \rangle &\leq \left(\frac{3}{2}\lambda_{\delta} + C_{v}\right) \left\| \left(\boldsymbol{\delta}^{(k)} - \hat{\boldsymbol{\delta}}^{(k)}\right)_{\mathcal{S}_{k}} \right\|_{1} + b_{g}^{*} \left\| \hat{\boldsymbol{v}}^{(k)} \right\|_{\Sigma} \\ &\leq m^{1/2} \left(\frac{3}{2}\lambda_{\delta} + C_{v}\right) \left\| \hat{\boldsymbol{v}}^{(k)} \right\|_{2} + b_{g}^{*} \left\| \hat{\boldsymbol{v}}^{(k)} \right\|_{\Sigma}. \end{split}$$

By Proposition 3.3, the RSC of  $\langle \nabla \hat{Q}_g^{(0)}(\hat{\boldsymbol{\omega}}^{(k)} + \hat{\boldsymbol{\delta}}^{(k)}) - \nabla \hat{Q}_g^{(0)}(\hat{\boldsymbol{\omega}}^{(k)} + \boldsymbol{\delta}^{(k)}), \hat{\boldsymbol{v}}^{(k)} \rangle$ , we have

$$0.2f_{l}\kappa_{l}\|\hat{\boldsymbol{v}}^{(k)}\|_{\Sigma}^{2} \leq m^{1/2}\left(\frac{3}{2}\lambda_{\delta}+C_{v}\right)\|\hat{\boldsymbol{v}}^{(k)}\|_{2}+b_{g}^{*}\|\hat{\boldsymbol{v}}^{(k)}\|_{\Sigma},$$

with probability at least  $1 - (pn_0)^{-1}$ . Therefore, if we let  $g^2 \leq m^{1/2} \lambda_{\delta}$ ,

$$\left\|\hat{\boldsymbol{v}}^{(k)}\right\|_{\Sigma}^{2} \lesssim \frac{m\log p}{n_{0}}$$

By Lemma 17 in Yuan et al. (2018) and the condition  $m \lesssim \sqrt{n_0/\log p}$ , we have

$$\left\|\tilde{\boldsymbol{\delta}}^{(k)} - \boldsymbol{\delta}^{(k)}\right\|_{\Sigma}^{2} \lesssim \frac{m\log p}{n_{0}} \text{ and } \left\|\tilde{\boldsymbol{\delta}}^{(k)} - \boldsymbol{\delta}^{(k)}\right\|_{1} \lesssim m\sqrt{\frac{\log p}{n_{0}}}$$

For step 3, let  $\tilde{\delta}^{(0)} = \delta^{(0)} = 0$ , then the loss function in step 3 could be written as:

$$\frac{1}{n_0 + n_{\mathcal{A}'_m}} \sum_{k \in \{0\} \cup \mathcal{A}'_m} \sum_{i=1}^{n_k} l_r(y_i^{(k)} - \langle \mathbf{X}_i^{(k)}, \beta - \tilde{\delta}^{(k)} \rangle) =: \sum_{k \in \{0\} \cup \mathcal{A}'_m} \hat{Q}_r^{(k)}(\beta - \tilde{\delta}^{(k)}).$$

The symmetric Bregman divergence is defined as

$$\left\langle \sum_{k \in \{0\} \cup \mathcal{A}'_m} \left( \nabla \hat{Q}_r^{(k)}(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\delta}}^{(k)}) - \nabla \hat{Q}_r^{(k)}(\boldsymbol{\beta}^* - \tilde{\boldsymbol{\delta}}^{(k)}) \right), \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^* \right\rangle.$$

To simplify the notations, define  $\nabla \hat{R}_r(\boldsymbol{\beta}) = \sum_{k \in \{0\} \cup \mathcal{A}'_m} \nabla \hat{Q}_r^{(k)}(\boldsymbol{\beta} - \tilde{\boldsymbol{\delta}}^{(k)})$ . Similarly as above, we have an oracle inequality for  $\hat{\boldsymbol{\beta}}$ ,

$$\begin{split} &\langle \nabla \hat{R}_{r}(\hat{\beta}) - \nabla \hat{R}_{r}(\beta^{*}), \hat{\beta} - \beta^{*} \rangle \\ &\leq \lambda_{\beta} \big( \|\beta^{*}\|_{1} - \|\hat{\beta}\|_{1} \big) + \sum_{k \in \{0\} \cup \mathcal{A}'_{m}} \big\langle \nabla \hat{Q}_{r}^{(k)}(\omega^{(k)} + \delta^{(k)} - \tilde{\delta}^{(k)}) - \nabla \hat{Q}_{r}^{(k)}(\omega^{(k)}), \beta^{*} - \hat{\beta} \big\rangle \\ &+ \sum_{k \in \{0\} \cup \mathcal{A}'_{m}} \big\langle \nabla \hat{Q}_{r}^{(k)}(\omega^{(k)}) - \nabla Q_{r}^{(k)}(\omega^{(k)}), \beta^{*} - \hat{\beta} \big\rangle + \sum_{k \in \{0\} \cup \mathcal{A}'_{m}} \big\langle \nabla Q_{r}^{(k)}(\omega^{(k)}), \beta^{*} - \hat{\beta} \big\rangle \\ &\leq \lambda_{\beta} \big( \|\beta^{*}\|_{1} - \|\hat{\beta}\|_{1} \big) + \sum_{k \in \{0\} \cup \mathcal{A}'_{m}} \big\langle \nabla \hat{Q}_{r}^{(k)}(\omega^{(k)} + \delta^{(k)} - \tilde{\delta}^{(k)}) - \nabla \hat{Q}_{r}^{(k)}(\omega^{(k)}), \beta^{*} - \hat{\beta} \big\rangle \\ &+ \sum_{k \in \{0\} \cup \mathcal{A}'_{m}} \underbrace{ \|\nabla \hat{Q}_{r}^{(k)}(\omega^{(k)}) - \nabla Q_{r}^{(k)}(\omega^{(k)})\|_{\infty}}_{||\pi_{r}^{(k)}||_{\infty}} \|\beta^{*} - \hat{\beta}\|_{1} \\ &+ \sum_{k \in \{0\} \cup \mathcal{A}'_{m}} \underbrace{ \|\Sigma^{-1/2} \nabla Q_{r}^{(k)}(\omega^{(k)})\|_{2}}_{b_{r}^{*}} \|\beta^{*} - \hat{\beta}\|_{\Sigma}. \end{split}$$

For the second term above, by Lemma A.4,

$$\underbrace{ \left\langle \nabla \hat{Q}_{r}^{(k)}(\boldsymbol{\omega}^{(k)} + \boldsymbol{\delta}^{(k)} - \tilde{\boldsymbol{\delta}}^{(k)}) - \nabla \hat{Q}_{r}^{(k)}(\boldsymbol{\omega}^{(k)}), \boldsymbol{\beta}^{*} - \hat{\boldsymbol{\beta}} \right\rangle}_{C_{v}} \leq \underbrace{ c_{0} \left( \frac{m}{r} \sqrt{\frac{\log p}{n_{0}}} \sqrt{\frac{\log p}{n_{0} + n_{\mathcal{A}_{m}^{\prime}}}} + \frac{\log p}{n_{0} + n_{\mathcal{A}_{m}^{\prime}}} + \sqrt{m} \sqrt{\frac{\log p}{n_{0}}} \right)}_{C_{v}^{\prime}} \left\| \boldsymbol{\beta}^{*} - \hat{\boldsymbol{\beta}} \right\|_{1}.$$

If we set  $\lambda_{\beta} \geq 2 ||\boldsymbol{\pi}_r^*||_{\infty}$  and  $C'_v \leq \lambda_{\beta}/2$ , then

$$\begin{split} \langle \nabla \hat{R}_{r}(\hat{\beta}) - \nabla \hat{R}_{r}(\beta^{*}), \hat{\beta} - \beta^{*} \rangle \\ &\leq \left(\frac{3}{2}\lambda_{\beta} + C_{v}'\right) \left\| \left(\beta^{*} - \hat{\beta}\right)_{\mathcal{S}} \right\|_{1} + b_{r}^{*} \left\|\beta^{*} - \hat{\beta}\right\|_{\Sigma} \\ &\leq s^{1/2} \left(\frac{3}{2}\lambda_{\beta} + C_{v}'\right) \left\|\beta^{*} - \hat{\beta}\right\|_{2} + b_{r}^{*} \left\|\beta^{*} - \hat{\beta}\right\|_{\Sigma}. \end{split}$$

Under the RSC of  $\langle \nabla \hat{R}_r(\hat{\beta}) - \nabla \hat{R}_r(\beta^*), \hat{\beta} - \beta^* \rangle$ , we have

$$\langle \nabla \hat{R}_r(\hat{\boldsymbol{\beta}}) - \nabla \hat{R}_r(\boldsymbol{\beta}^*), \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^* \rangle \ge c_1 \| \boldsymbol{\beta}^* - \hat{\boldsymbol{\beta}} \|_{\Sigma}^2,$$

with probability as least  $1 - (pn)^{-1}$ , where  $n = n_0 + n_{\mathcal{A}'_m}$  and  $c_1$  is a positive constant. The proof of the RSC in step 3 is similar to Proposition 3.3. Thus,

$$\left\|\boldsymbol{\beta}^* - \hat{\boldsymbol{\beta}}\right\|_{\Sigma} \le s^{1/2} \left(\frac{3}{2}\lambda_{\beta} + C'_v\right) \gamma_p^{-1/2} + b_r^*.$$

Through a similar proof as Lemma A.2, we obtain  $\lambda_{\beta} \lesssim \sqrt{\log(p)/n}$ . If  $s \log(p)/n \leq r^2 \leq s^{1/2} \lambda_{\beta}$ , we have

$$\left\|\boldsymbol{\beta}^* - \hat{\boldsymbol{\beta}}\right\|_{\Sigma} \lesssim \sqrt{\frac{s\log p}{n}} + \sqrt{\frac{sm\log p}{n_0}} \text{ and } \left\|\boldsymbol{\beta}^* - \hat{\boldsymbol{\beta}}\right\|_1 \lesssim s\sqrt{\frac{\log p}{n}} + s\sqrt{\frac{m\log p}{n_0}},$$

with probability at least  $1 - p^{-1}$ .

# **B** Proof of Lemmas

# B.1 Proof of Lemma A.1

Define  $\boldsymbol{\omega}^{(k)}$  for all  $0 \leq k \leq K$  as the true parameters of each local source model, then note that  $\nabla Q^{(k)}(\boldsymbol{\omega}^{(k)}) = 0$  and  $\nabla Q(\boldsymbol{\omega}^*) = \sum_{k=0}^{K} \alpha_k \nabla Q^{(k)}(\boldsymbol{\omega}^*) = 0$ . So we have

$$\nabla Q(\boldsymbol{\omega}^*) - \nabla Q(\boldsymbol{\beta}^*) + \nabla Q(\boldsymbol{\beta}^*) - \sum_{k=1}^K \alpha_k \nabla Q^{(k)}(\boldsymbol{\omega}^{(k)}) = 0$$
$$\nabla Q(\boldsymbol{\omega}^*) - \nabla Q(\boldsymbol{\beta}^*) = \sum_{k=1}^K \alpha_k \nabla Q^{(k)}(\boldsymbol{\omega}^{(k)}) - \nabla Q(\boldsymbol{\beta}^*)$$

Note that  $\nabla Q^{(0)}(\boldsymbol{\omega}^{(0)}) = Q^{(0)}(\boldsymbol{\beta}^*) = 0$ , so

$$\sum_{k=0}^{K} \alpha_k (\nabla Q^{(k)}(\boldsymbol{\omega}^*) - \nabla Q^{(k)}(\boldsymbol{\beta}^*)) = \sum_{k=1}^{K} \alpha_k (\nabla Q^{(k)}(\boldsymbol{\omega}^{(k)}) - \nabla Q^{(k)}(\boldsymbol{\beta}^*))$$

By the second-order Taylor expansions and Assumption 3.4,

$$\sum_{k=0}^{K} \alpha_k \int_0^1 \nabla^2 Q^{(k)} ((1-t)\beta^* + t\omega^*) dt (\omega^* - \beta^*) = \sum_{k=1}^{K} \alpha_k \int_0^1 \nabla^2 Q^{(k)} ((1-t)\beta^* + t\omega^{(k)}) dt (\omega^{(k)} - \beta^*) dt (\omega^{(k)} - \beta^*)$$

By the definition of the parameter space

$$\Theta(s,m) = \Big\{ \boldsymbol{\beta}^*, \{\boldsymbol{\omega}^{(k)}\} : ||\boldsymbol{\beta}^*||_0 \le s, \sup_{k \in \mathcal{A}_m} ||\boldsymbol{\omega}^{(k)} - \boldsymbol{\beta}^*||_1 \le m \Big\},\$$

We have  $||\boldsymbol{\omega}^{(k)} - \boldsymbol{\beta}^*||_1 \leq m$ . Let  $C_1 = \sup_k ||\tilde{\Sigma}^{-1}\tilde{\Sigma}^{(k)}||_1$ . Then Lemma A.1 is proved.

# B.2 Proof of Lemma A.2

For the transferring steps,

$$\nabla \hat{Q}_h(\boldsymbol{\omega}) = \frac{1}{n_{\mathcal{A}_m} + n_0} \sum_{k=0}^K \sum_{i=1}^{n_k} \left\{ \bar{K} \left( \frac{\langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega} \rangle - y_i^{(k)}}{h} \right) - \tau \right\} \boldsymbol{x}_i^{(k)}$$
$$\nabla^2 \hat{Q}_h(\boldsymbol{\omega}) = \frac{1}{n_{\mathcal{A}_m} + n_0} \sum_{k=0}^K \sum_{i=1}^{n_k} K \left( \frac{\langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega} \rangle - y_i^{(k)}}{h} \right) \boldsymbol{x}_i^{(k)} (\boldsymbol{x}_i^{(k)})^{\mathrm{T}}.$$

Let  $\xi_i^{(k)} = \bar{K}\{(\langle \boldsymbol{x}_i^{(k)}, \boldsymbol{\omega} \rangle - y_i^{(k)})/h\} - \tau$ , then  $\nabla \hat{Q}_h(\boldsymbol{\omega}) = (n_{\mathcal{A}_m} + n_0)^{-1} \sum_{k=0}^K \sum_{i=1}^{n_k} \xi_i^{(k)} \boldsymbol{x}_i^{(k)}$  and

$$\|\boldsymbol{\pi}_{h}^{*}\|_{\infty} = \left\|\frac{1}{n_{\mathcal{A}_{m}} + n_{0}} \sum_{k=0}^{K} \sum_{i=1}^{n_{k}} \left\{\xi_{i}^{(k)} \boldsymbol{x}_{i}^{(k)} - \mathbb{E}(\xi_{i}^{(k)} \boldsymbol{x}_{i}^{(k)})\right\}\right\|_{\infty}$$

The upper bound of  $||\pi_h^*||_{\infty}$  involves two quantities that are related to

$$\mathbb{E}\bigg[\bar{K}^2\bigg(\frac{\langle \boldsymbol{x}^{(k)}, \boldsymbol{\omega} - \boldsymbol{\omega}^{(k)} \rangle - \epsilon}{h}\bigg) \Big| \boldsymbol{x}^{(k)}\bigg] \text{ and } \mathbb{E}\bigg[\bigg(\frac{\langle \boldsymbol{x}^{(k)}, \boldsymbol{\omega} - \boldsymbol{\omega}^{(k)} \rangle - \epsilon}{h}\bigg) \Big| \boldsymbol{x}^{(k)}\bigg].$$

For the first term, by a change of variable and integration by parts, we obtain

$$\mathbb{E}\left[\bar{K}^{2}\left(\frac{\langle \boldsymbol{x}^{(k)}, \boldsymbol{\omega} - \boldsymbol{\omega}^{(k)} \rangle - \epsilon}{h}\right) \middle| \boldsymbol{x}^{(k)} \right] = \int_{-\infty}^{\infty} \bar{K}^{2}(-u/h) f_{\epsilon|\boldsymbol{x}}(u) du$$
$$= h \int_{-\infty}^{\infty} \bar{K}^{2}(v) f_{\epsilon|\boldsymbol{x}}(-vh) dv$$
$$= 2 \int_{-\infty}^{\infty} K(v) \bar{K}(v) F_{\epsilon|\boldsymbol{x}}(-vh) dv.$$
(50)

By the fact that  $F_{\epsilon|\boldsymbol{x}}(0) = \tau$ , we have

$$F_{\epsilon|\boldsymbol{x}}(-vh) = F_{\epsilon|\boldsymbol{x}}(0) + \int_{0}^{-vh} f_{\epsilon|\boldsymbol{x}}(t)dt$$
$$= \tau - hvf_{\epsilon|\boldsymbol{x}}(0) + \int_{0}^{-vh} \{f_{\epsilon|\boldsymbol{x}}(t) - f_{\epsilon|\boldsymbol{x}}(0)\}dt.$$
(51)

Moreover, it can be shown that

$$a_K := \int_{-\infty}^{\infty} v K(v) \bar{K}(v) dv = \int_0^{\infty} K(v) \{1 - K(v)\} dv > 0 \text{ and } a_K \le \kappa_1,$$
(52)

where  $\kappa_1 = \int |u| K(u) du$ .

Substituting (51) into (50), and by (52), we obtain

$$\mathbb{E}\left[\bar{K}^{2}\left(\frac{\langle \boldsymbol{x}^{(k)}, \boldsymbol{\omega} - \boldsymbol{\omega}^{(k)} \rangle - \epsilon}{h}\right) \middle| \boldsymbol{x}^{(k)} \right] = 2\tau \int_{-\infty}^{\infty} K(v)\bar{K}(v)dv - 2hf_{\epsilon|\boldsymbol{x}}(0) \int_{-\infty}^{\infty} vK(v)\bar{K}(v)dv + 2\int_{-\infty}^{\infty} \int_{0}^{-vh} \{f_{\epsilon|\boldsymbol{x}}(t) - f_{\epsilon|\boldsymbol{x}}(0)\}K(v)\bar{K}(v)dtdv \\ \leq \tau - 2a_{K}hf_{\epsilon|\boldsymbol{x}}(0) + l_{0}h^{2} \int_{-\infty}^{\infty} v^{2}K(v)\bar{K}(v)dv \\ \leq \tau + l_{0}\kappa_{2}h^{2},$$

where the first inequality holds using the Lipschitz condition on  $f_{\epsilon|x}$  in Assumption 4.1 and the last inequality holds by Assumption 4.2. Through a similar calculation, the Lipschitz condition on  $f_{\epsilon|x}$  ensures that

$$\left|\mathbb{E}\left[\left(\frac{\langle \boldsymbol{x}^{(k)}, \boldsymbol{\omega} - \boldsymbol{\omega}^{(k)} \rangle - \epsilon}{h}\right) \middle| \boldsymbol{x}^{(k)}\right] - \tau\right| \leq \frac{l_0}{2}\kappa_2 h^2.$$

Hence

$$\begin{split} \mathbb{E}(\xi_i^{(k)} x_{ij}^{(k)})^2 &= \mathbb{E}_x \left\{ (x_{ij}^{(k)})^2 \cdot \mathbb{E}((\xi_i^{(k)})^2 | \boldsymbol{x}_i^{(k)}) \right\} \\ \mathbb{E}(\xi^2 | \boldsymbol{x}) &= \mathbb{E}\left[ \left( \bar{K} \left( \frac{\langle \boldsymbol{x}, \boldsymbol{\omega} \rangle - \boldsymbol{y}}{h} \right) - \tau \right)^2 | \boldsymbol{x} \right] \\ &= \underbrace{\mathbb{E}\left[ \bar{K}^2 \left( \frac{\langle \boldsymbol{x}, \boldsymbol{\omega} \rangle - \boldsymbol{y}}{h} \right) | \boldsymbol{x} \right]}_{\leq \tau + l_0 \kappa_2 h^2} - 2\tau \underbrace{\mathbb{E}\left[ \bar{K} \left( \frac{\langle \boldsymbol{x}, \boldsymbol{\omega} \rangle - \boldsymbol{y}}{h} \right) | \boldsymbol{x} \right]}_{\geq \tau - \frac{l_0}{2} \kappa_2 h^2} + \tau^2 \\ &\leq \tau (1 - \tau) + Ch^2, \end{split}$$

where  $C = (\tau + 1)l_0\kappa_2$ . Then, by Assumption 4.3, we have

$$\mathbb{E}(\xi_i^{(k)} x_{ij}^{(k)})^2 \le \tau (1-\tau)\sigma_{jj} + C\sigma_{jj}h^2.$$

Also by Assumption 4.3 and  $|\xi_i^{(k)}| \leq \max(1-\tau,\tau)$ , for  $s = 3, 4, \ldots$ ,

$$\mathbb{E}(|\xi_i^{(k)} x_{ij}^{(k)}|^s) \leq \max(1-\tau,\tau)^{s-2} \mathbb{E}_{\boldsymbol{x}}\{|x_{ij}^{(k)}|^s \cdot \mathbb{E}[(\xi_i^{(k)})^2 | \boldsymbol{x}_i^{(k)}]\} \\ \leq \max(1-\tau,\tau)^{s-2} \{\tau(1-\tau) + Ch^2\} \\ \leq \frac{s!}{2} \{\tau(1-\tau) + Ch^2\} \max(1-\tau,\tau)^{s-2}.$$

Thus it follows from Bernstein's inequality and union bound that for every  $t \ge 0$ ,

$$||\boldsymbol{\pi}_{h}^{*}||_{\infty} \leq \sigma \sqrt{\{\tau(1-\tau) + Ch^{2}\}\frac{2t}{n_{\mathcal{A}_{m}} + n_{0}}} + \max(1-\tau,\tau)\frac{t}{n_{\mathcal{A}_{m}} + n_{0}}$$

with probability at least  $1 - 2pe^{-t}$ .

For the debiasing step, through the similar proof we could get same results with different sample size and smoothing bandwidth.

## B.3 Proof of Lemma A.3

Note that

$$\begin{split} b_h^* &= ||\Sigma^{-1/2} \nabla Q_h(\boldsymbol{\omega}^*)||_2 \\ &= \left\| \Sigma^{-1/2} \left( \sum_{k=0}^K \alpha_k \mathbb{E} \left[ \mathbb{E} \left\{ \bar{K} \left( \frac{\langle \boldsymbol{x}^{(k)}, \boldsymbol{\omega}^* \rangle - y^{(k)}}{h} \right) - \tau \middle| \boldsymbol{x}^{(k)} \right\} \boldsymbol{x}^{(k)} \right] \right) \right\|_2 \\ &\leq \sup_{\boldsymbol{u} \in \mathbb{S}^{p-1}} \sum_{k=0}^K \mathbb{E} \left[ \bar{K} \left( \frac{\langle \boldsymbol{x}^{(k)}, \boldsymbol{\omega}^* \rangle - y^{(k)}}{h} \right) - \tau \right] \langle \Sigma^{-1/2} \boldsymbol{x}^{(k)}, \boldsymbol{u} \rangle \\ &\leq \frac{l_0}{2} \kappa_2 h^2. \end{split}$$

## B.4 Proof of Lemma A.4

For  $k = 1, \ldots, p$ , define that

$$\psi_k(r,l) = \sup_{\boldsymbol{v} \in \mathbb{B}_{\Sigma}(r) \cap \mathbb{B}_1(l)} \left| \frac{1}{n} \sum_{i=1}^n (1-\mathbb{E}) \underbrace{\left\{ \bar{K}_h(\boldsymbol{x}_i^{\mathrm{T}} \boldsymbol{v} - \epsilon_i) - \bar{K}_h(-\epsilon_i) \right\} x_{ik}}_{=:g_{\boldsymbol{v},k}(y_i, \boldsymbol{x}_i)} \right|,$$

where  $\boldsymbol{v} = \boldsymbol{\beta} - \boldsymbol{\beta}^*$ . Note that  $\psi(r, l) \leq \max_{1 \leq k \leq p} \{\psi_k(r, l) + |\mathbb{E}g_{\boldsymbol{v},k}(y_i, \boldsymbol{x}_i)|\}$ . In the following, we bound  $\psi_k(r, l)$  and  $\mathbb{E}g_{\boldsymbol{v},k}(y_i, \boldsymbol{x}_i)$ , respectively.

Let  $\sigma$  be any positive constant such that  $\sigma^2 \geq \sup_{\boldsymbol{v} \in \mathbb{B}_{\Sigma}(r) \cap \mathbb{B}_1(l)} \mathbb{E}g_{\boldsymbol{v},k}^2(y_i, \boldsymbol{x}_i)$ . By the bounded design, we note that  $\sup_{\boldsymbol{v}} |g_{\boldsymbol{v},k}(y_i, \boldsymbol{x}_i)| \leq |x_{ik}| \leq 1$ . Applying Theorem 7.3 in Bousquet (2003), Bousquet's version of Talagrand's inequality, we obtain that for any z > 0,

$$\psi_k(r,l) \le \mathbb{E}\psi_k(r,l) + \sqrt{\{\sigma^2 + 2\mathbb{E}\psi_k(r,l)\}\frac{2z}{n} + \frac{z}{3n}}$$
(53)

holds with probability at least  $1 - e^{-z}$ . For the second moment  $\mathbb{E}g_{\boldsymbol{v},k}^2(y_i, \boldsymbol{x}_i)$ , by a change of variable and Minkowski's integral inequality we derive that

$$\begin{split} \mathbb{E}g_{\boldsymbol{v},k}^{2}(y_{i},\boldsymbol{x}_{i}) &= \mathbb{E}\left[x_{ik}^{2}\int_{-\infty}^{\infty}\left\{\bar{K}_{h}(\boldsymbol{x}_{i}^{\mathrm{T}}\boldsymbol{v}-t)-\bar{K}_{h}(-t)\right\}^{2}f_{\epsilon_{i}|\boldsymbol{x}_{i}}(t)dt\right]\\ &= \mathbb{E}\left[x_{ik}^{2}\int_{-\infty}^{\infty}\left\{\bar{K}_{h}(u)-\bar{K}_{h}(u-\boldsymbol{x}_{i}^{\mathrm{T}}\boldsymbol{v})\right\}^{2}f_{\epsilon_{i}|\boldsymbol{x}_{i}}(\boldsymbol{x}_{i}^{\mathrm{T}}\boldsymbol{v}-u)du\right]\\ &= h\mathbb{E}\left[x_{ik}^{2}\int_{-\infty}^{\infty}\left\{\bar{K}(v)-\bar{K}(v-\boldsymbol{x}_{i}^{\mathrm{T}}\boldsymbol{v}/h)\right\}^{2}f_{\epsilon_{i}|\boldsymbol{x}_{i}}(\boldsymbol{x}_{i}^{\mathrm{T}}\boldsymbol{v}-vh)dv\right]\\ &\leq f_{u}h^{-1}\mathbb{E}\left[x_{ik}^{2}(\boldsymbol{x}_{i}^{\mathrm{T}}\boldsymbol{v})^{2}\int_{-\infty}^{\infty}\left\{\int_{0}^{1}K(v-w\boldsymbol{x}_{i}^{\mathrm{T}}\boldsymbol{v}/h)dw\right\}^{2}dv\right]\\ &\leq f_{u}h^{-1}\mathbb{E}\left(x_{ik}^{2}(\boldsymbol{x}_{i}^{\mathrm{T}}\boldsymbol{v})^{2}\left[\int_{0}^{1}\left\{\int_{-\infty}^{\infty}K^{2}(v-w\boldsymbol{x}_{i}^{\mathrm{T}}\boldsymbol{v}/h)dv\right\}^{1/2}dw\right]^{2}\right)\\ &\leq \kappa_{u}f_{u}h^{-1}\mathbb{E}(x_{ik}\cdot\boldsymbol{x}_{i}^{\mathrm{T}}\boldsymbol{v})^{2}\leq \kappa_{u}f_{u}h^{-1}r^{2}. \end{split}$$

It remains to bound  $\mathbb{E}\psi_k(r,l)$ . Note that  $|g_{\boldsymbol{v},k}(y_i,\boldsymbol{x}_i) - g_{\boldsymbol{v}',k}(y_i,\boldsymbol{x}_i)| \leq (\kappa_u/h)|\boldsymbol{x}_i^{\mathrm{T}}\boldsymbol{v} - \boldsymbol{x}_i^{\mathrm{T}}\boldsymbol{v}'|$ , for any  $\boldsymbol{v}, \boldsymbol{v}'$ . Hence using Rademacher symmetrization and Talagrand's contraction principle, we have

$$\mathbb{E}\psi_{k}(r,l) \leq 2\mathbb{E}\left[\sup_{\boldsymbol{v}\in\mathbb{B}_{\Sigma}(r)\cap\mathbb{B}_{1}(l)}\left|\frac{1}{n}\sum_{i=1}^{n}e_{i}g_{\boldsymbol{v},k}(y_{i},\boldsymbol{x}_{i})\right|\right] \\ \leq 4\kappa_{u}\mathbb{E}\left[\sup_{\boldsymbol{v}\in\mathbb{B}_{\Sigma}(r)\cap\mathbb{B}_{1}(l)}\left|\frac{1}{nh}\sum_{i=1}^{n}e_{i}\boldsymbol{x}_{i}^{\mathrm{T}}\boldsymbol{v}\right|\right] \leq 4\kappa_{u}\frac{l}{h}\mathbb{E}\left\|\frac{1}{n}\sum_{i=1}^{n}e_{i}\boldsymbol{x}_{i}\right\|_{\infty},$$
(54)

where  $e_1, \ldots, e_n$  are independent Rademacher variables. Applying Hoeffding's moment inequality,

$$\mathbb{E}_e \left\| \frac{1}{n} \sum_{i=1}^n e_i \boldsymbol{x}_i \right\|_{\infty} \le \max_{1 \le k \le p} \left( \sum_{i=1}^n x_{ik}^2 \right)^{1/2} \frac{\sqrt{2\log(2p)}}{n},\tag{55}$$

where  $\mathbb{E}_e$  denotes the expectation over  $\{e_i\}_{i=1}^n$ . By (54) and (55), we obtain

$$\mathbb{E}\psi_k(r,l) \le 4\kappa_u \frac{l}{h} \sqrt{\frac{2\log(2p)}{n}}$$

Taking  $z = t + \log p$  in (53), we have that

$$\psi_k(r,l) \lesssim \frac{l}{h} \sqrt{\frac{\log p}{n}} + f_u^{1/2} r \sqrt{\frac{t+\log p}{nh}} + \frac{t+\log p}{n}$$
(56)

holds with probability at least  $1 - e^{-t}$ .

Next we find an union upper bound of  $|\mathbb{E}g_{\boldsymbol{v},k}(y_i,\boldsymbol{x}_i)|$ . Similarly as the method to bound the second moment, we derive that

$$\begin{split} \mathbb{E}g_{\boldsymbol{v},k}(y_i, \boldsymbol{x}_i) &= h\mathbb{E}\left[x_{ik}\int_{-\infty}^{\infty}\left\{\bar{K}(v) - \bar{K}(v - \boldsymbol{x}_i^{\mathrm{T}}\boldsymbol{v}/h)\right\}f_{\epsilon_i|\boldsymbol{x}_i}(\boldsymbol{x}_i^{\mathrm{T}}\boldsymbol{v} - vh)dv\right] \\ &\leq f_u\mathbb{E}\left[|x_{ik}||\boldsymbol{x}_i^{\mathrm{T}}\boldsymbol{v}|\int_{-\infty}^{\infty}\left\{\int_{0}^{1}K(v - w\boldsymbol{x}_i^{\mathrm{T}}\boldsymbol{v}/h)dw\right\}dv\right] \\ &\leq f_u\mathbb{E}\left(|x_{ik}||\boldsymbol{x}_i^{\mathrm{T}}\boldsymbol{v}|\left[\int_{0}^{1}\left\{\int_{-\infty}^{\infty}K(v - w\boldsymbol{x}_i^{\mathrm{T}}\boldsymbol{v}/h)dv\right\}dw\right]\right) \\ &\leq \kappa_u f_u\mathbb{E}|x_{ik}\boldsymbol{x}_i^{\mathrm{T}}\boldsymbol{v}| \leq \kappa_u f_u r. \end{split}$$

Finally taking the union bound, we obtain that with probability at least  $1 - e^{-t}$ ,

$$\psi(r,l) \lesssim \frac{l}{h} \sqrt{\frac{\log p}{n}} + f_u^{1/2} r \sqrt{\frac{t + \log p}{nh}} + \frac{t + \log p}{n}.$$