

A Applications to Latent Variable Models

In the sequel, we introduce two latent variable models as examples. To apply the high dimensional EM algorithm in §2.1 and the methods for asymptotic inference in §2.2, we only need to specify the forms of $Q_n(\cdot; \cdot)$ defined in (2.1), $M_n(\cdot)$ in Algorithms 2 and 3, and $T_n(\cdot)$ in (2.4) for each model.

Gaussian Mixture Model: Let $\mathbf{y}_1, \dots, \mathbf{y}_n$ be the n i.i.d. realizations of $\mathbf{Y} \in \mathbb{R}^d$ and

$$\mathbf{Y} = Z \cdot \boldsymbol{\beta}^* + \mathbf{V}. \quad (\text{A.1})$$

Here Z is a Rademacher random variable, i.e., $\mathbb{P}(Z = +1) = \mathbb{P}(Z = -1) = 1/2$, and $\mathbf{V} \sim N(\mathbf{0}, \sigma^2 \cdot \mathbf{I}_d)$ is independent of Z , where σ is the standard deviation. We suppose σ is known. In high dimensional settings, we assume that $\boldsymbol{\beta}^* \in \mathbb{R}^d$ is sparse. To avoid the degenerate case in which the two Gaussians in the mixture are identical, here we suppose that $\boldsymbol{\beta}^* \neq \mathbf{0}$.

For the E-step (line 4 of Algorithm 1), we have

$$Q_n(\boldsymbol{\beta}'; \boldsymbol{\beta}) = -\frac{1}{2n} \sum_{i=1}^n \omega_{\boldsymbol{\beta}}(\mathbf{y}_i) \cdot \|\mathbf{y}_i - \boldsymbol{\beta}'\|_2^2 + [1 - \omega_{\boldsymbol{\beta}}(\mathbf{y}_i)] \cdot \|\mathbf{y}_i + \boldsymbol{\beta}'\|_2^2, \quad (\text{A.2})$$

$$\text{where } \omega_{\boldsymbol{\beta}}(\mathbf{y}) = \frac{1}{1 + \exp(-\langle \boldsymbol{\beta}, \mathbf{y} \rangle / \sigma^2)}.$$

The maximization implementation (Algorithm 2) of the M-step takes the form

$$M_n(\boldsymbol{\beta}) = \frac{2}{n} \sum_{i=1}^n \omega_{\boldsymbol{\beta}}(\mathbf{y}_i) \cdot \mathbf{y}_i - \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i. \quad (\text{A.3})$$

Meanwhile, for the gradient ascent implementation (Algorithm 3) of the M-step, we have

$$M_n(\boldsymbol{\beta}) = \boldsymbol{\beta} + \eta \cdot \nabla_1 Q_n(\boldsymbol{\beta}; \boldsymbol{\beta}), \quad \text{where } \nabla_1 Q_n(\boldsymbol{\beta}; \boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n [2 \cdot \omega_{\boldsymbol{\beta}}(\mathbf{y}_i) - 1] \cdot \mathbf{y}_i - \boldsymbol{\beta}.$$

Here $\eta > 0$ is the stepsize. For asymptotic inference, $T_n(\cdot)$ in (2.4) takes the form

$$T_n(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \nu_{\boldsymbol{\beta}}(\mathbf{y}_i) \cdot \mathbf{y}_i \cdot \mathbf{y}_i^\top - \mathbf{I}_d,$$

$$\text{where } \nu_{\boldsymbol{\beta}}(\mathbf{y}) = \frac{4/\sigma^2}{[1 + \exp(-2 \cdot \langle \boldsymbol{\beta}, \mathbf{y} \rangle / \sigma^2)] \cdot [1 + \exp(2 \cdot \langle \boldsymbol{\beta}, \mathbf{y} \rangle / \sigma^2)]}.$$

Mixture of Regression Model: We assume that $Y \in \mathbb{R}$ and $\mathbf{X} \in \mathbb{R}^d$ satisfy

$$Y = Z \cdot \mathbf{X}^\top \boldsymbol{\beta}^* + V, \quad (\text{A.4})$$

where $\mathbf{X} \sim N(\mathbf{0}, \mathbf{I}_d)$, $V \sim N(0, \sigma^2)$ and Z is a Rademacher random variable. Here \mathbf{X} , V and Z are independent. In the high dimensional regime, we assume $\boldsymbol{\beta}^* \in \mathbb{R}^d$ is sparse. To avoid the degenerate case, we suppose $\boldsymbol{\beta}^* \neq \mathbf{0}$. In addition, we assume that σ is known. For the E-step (line 4 of Algorithm 1), we have

$$Q_n(\boldsymbol{\beta}'; \boldsymbol{\beta}) = -\frac{1}{2n} \sum_{i=1}^n \omega_{\boldsymbol{\beta}}(\mathbf{x}_i, y_i) \cdot (y_i - \langle \mathbf{x}_i, \boldsymbol{\beta}' \rangle)^2 + [1 - \omega_{\boldsymbol{\beta}}(\mathbf{x}_i, y_i)] \cdot (y_i + \langle \mathbf{x}_i, \boldsymbol{\beta}' \rangle)^2, \quad (\text{A.5})$$

$$\text{where } \omega_{\boldsymbol{\beta}}(\mathbf{x}, y) = \frac{1}{1 + \exp(-y \cdot \langle \boldsymbol{\beta}, \mathbf{x} \rangle / \sigma^2)}.$$

For the maximization implementation (Algorithm 2) of the M-step (line 5 of Algorithm 1), we have that $M_n(\boldsymbol{\beta}) = \operatorname{argmax}_{\boldsymbol{\beta}'} Q_n(\boldsymbol{\beta}'; \boldsymbol{\beta})$ satisfies

$$\widehat{\boldsymbol{\Sigma}} \cdot M_n(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n [2 \cdot \omega_{\boldsymbol{\beta}}(\mathbf{x}_i, y_i) - 1] \cdot y_i \cdot \mathbf{x}_i, \quad \text{where } \widehat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{x}_i^\top. \quad (\text{A.6})$$

However, in high dimensional regimes, the sample covariance matrix $\widehat{\boldsymbol{\Sigma}}$ is not invertible. To estimate the inverse covariance matrix of \mathbf{X} , we use the CLIME estimator proposed by [7], i.e.,

$$\widehat{\boldsymbol{\Theta}} = \operatorname{argmin}_{\boldsymbol{\Theta} \in \mathbb{R}^{d \times d}} \|\boldsymbol{\Theta}\|_{1,1}, \quad \text{subject to } \|\widehat{\boldsymbol{\Sigma}} \cdot \boldsymbol{\Theta} - \mathbf{I}_d\|_{\infty, \infty} \leq \lambda^{\text{CLIME}}, \quad (\text{A.7})$$

where $\|\cdot\|_{1,1}$ and $\|\cdot\|_{\infty, \infty}$ are the sum and maximum of the absolute values of all entries respectively, and $\lambda^{\text{CLIME}} > 0$ is a tuning parameter. Based on (A.6), we modify the maximization implementation

of the M-step to be

$$M_n(\boldsymbol{\beta}) = \widehat{\boldsymbol{\Theta}} \cdot \frac{1}{n} \sum_{i=1}^n [2 \cdot \omega_{\boldsymbol{\beta}}(\mathbf{x}_i, y_i) - 1] \cdot y_i \cdot \mathbf{x}_i. \quad (\text{A.8})$$

For the gradient ascent implementation (Algorithm 3) of the M-step, we have

$$M_n(\boldsymbol{\beta}) = \boldsymbol{\beta} + \eta \cdot \nabla_1 Q_n(\boldsymbol{\beta}; \boldsymbol{\beta}), \quad (\text{A.9})$$

$$\text{where } \nabla_1 Q_n(\boldsymbol{\beta}, \boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n [2 \cdot \omega_{\boldsymbol{\beta}}(\mathbf{x}_i, y_i) \cdot y_i \cdot \mathbf{x}_i - \mathbf{x}_i \cdot \mathbf{x}_i^\top \cdot \boldsymbol{\beta}].$$

Here $\eta > 0$ is the stepsize. For asymptotic inference, $T_n(\cdot)$ in (2.4) takes the form

$$T_n(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \nu_{\boldsymbol{\beta}}(\mathbf{x}_i, y_i) \cdot \mathbf{x}_i \cdot \mathbf{x}_i^\top \cdot y_i^2 - \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{x}_i^\top,$$

$$\text{where } \nu_{\boldsymbol{\beta}}(\mathbf{x}, y) = \frac{4/\sigma^2}{[1 + \exp(-2 \cdot y \cdot \langle \boldsymbol{\beta}, \mathbf{x} \rangle / \sigma^2)] \cdot [1 + \exp(2 \cdot y \cdot \langle \boldsymbol{\beta}, \mathbf{x} \rangle / \sigma^2)]}.$$

It is worth noting that, for the maximization implementation of the M-step, the CLIME estimator in (A.7) requires that $\boldsymbol{\Sigma}^{-1}$ is sparse, where $\boldsymbol{\Sigma}$ is the population covariance of \mathbf{X} . Since we assume $\mathbf{X} \sim N(\mathbf{0}, \mathbf{I}_d)$, this requirement is satisfied. Nevertheless, for more general settings where $\boldsymbol{\Sigma}$ does not possess such a structure, the gradient ascent implementation of the M-step is a better choice, since it does not require inverse covariance estimation and is also more efficient in computation.

B Derivation of the EM Algorithm

Recall that in §2.1, we assume that $h_{\boldsymbol{\beta}}(\mathbf{y})$ is obtained by marginalizing over an unobserved latent variable $\mathbf{Z} \in \mathcal{Z}$, i.e.,

$$h_{\boldsymbol{\beta}}(\mathbf{y}) = \int_{\mathcal{Z}} f_{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{z}) \, d\mathbf{z}. \quad (\text{B.1})$$

Let $k_{\boldsymbol{\beta}}(\mathbf{z} \mid \mathbf{y})$ be the density of \mathbf{Z} conditioning on the observed variable $\mathbf{Y} = \mathbf{y}$, i.e.,

$$k_{\boldsymbol{\beta}}(\mathbf{z} \mid \mathbf{y}) = f_{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{z}) / h_{\boldsymbol{\beta}}(\mathbf{y}). \quad (\text{B.2})$$

Given the n observations $\mathbf{y}_1, \dots, \mathbf{y}_n$ of \mathbf{Y} , the EM algorithm aims at maximizing the log-likelihood

$$\ell_n(\boldsymbol{\beta}) = \sum_{i=1}^n \log h_{\boldsymbol{\beta}}(\mathbf{y}_i). \quad (\text{B.3})$$

Due to the unobserved latent variable \mathbf{Z} , it is difficult to directly evaluate $\ell_n(\boldsymbol{\beta})$. Instead, we turn to consider the difference between $\ell_n(\boldsymbol{\beta})$ and $\ell_n(\boldsymbol{\beta}')$. Let $k_{\boldsymbol{\beta}}(\mathbf{z} \mid \mathbf{y})$ be the density of \mathbf{Z} conditioning on the observed variable $\mathbf{Y} = \mathbf{y}$, i.e.,

$$k_{\boldsymbol{\beta}}(\mathbf{z} \mid \mathbf{y}) = f_{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{z}) / h_{\boldsymbol{\beta}}(\mathbf{y}). \quad (\text{B.4})$$

According to (B.1) and (B.3), we have

$$\begin{aligned} \frac{1}{n} \cdot [\ell_n(\boldsymbol{\beta}) - \ell_n(\boldsymbol{\beta}')] &= \frac{1}{n} \sum_{i=1}^n \log [h_{\boldsymbol{\beta}}(\mathbf{y}_i) / h_{\boldsymbol{\beta}'}(\mathbf{y}_i)] = \frac{1}{n} \sum_{i=1}^n \log \left[\int_{\mathcal{Z}} \frac{f_{\boldsymbol{\beta}}(\mathbf{y}_i, \mathbf{z})}{h_{\boldsymbol{\beta}'}(\mathbf{y}_i)} \, d\mathbf{z} \right] \\ &= \frac{1}{n} \sum_{i=1}^n \log \left[\int_{\mathcal{Z}} k_{\boldsymbol{\beta}'}(\mathbf{z} \mid \mathbf{y}_i) \cdot \frac{f_{\boldsymbol{\beta}}(\mathbf{y}_i, \mathbf{z})}{f_{\boldsymbol{\beta}'}(\mathbf{y}_i, \mathbf{z})} \, d\mathbf{z} \right] \geq \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{Z}} k_{\boldsymbol{\beta}'}(\mathbf{z} \mid \mathbf{y}_i) \cdot \log \left[\frac{f_{\boldsymbol{\beta}}(\mathbf{y}_i, \mathbf{z})}{f_{\boldsymbol{\beta}'}(\mathbf{y}_i, \mathbf{z})} \right] \, d\mathbf{z}, \end{aligned} \quad (\text{B.5})$$

where the third equality follows from (B.4) and the inequality is obtained from Jensen's inequality. On the right-hand side of (B.5) we have

$$\begin{aligned} &\frac{1}{n} \sum_{i=1}^n \int_{\mathcal{Z}} k_{\boldsymbol{\beta}'}(\mathbf{z} \mid \mathbf{y}_i) \cdot \log \left[\frac{f_{\boldsymbol{\beta}}(\mathbf{y}_i, \mathbf{z})}{f_{\boldsymbol{\beta}'}(\mathbf{y}_i, \mathbf{z})} \right] \, d\mathbf{z} \\ &= \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{Z}} k_{\boldsymbol{\beta}'}(\mathbf{z} \mid \mathbf{y}_i) \cdot \log f_{\boldsymbol{\beta}}(\mathbf{y}_i, \mathbf{z}) \, d\mathbf{z} - \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{Z}} k_{\boldsymbol{\beta}'}(\mathbf{z} \mid \mathbf{y}_i) \cdot \log f_{\boldsymbol{\beta}'}(\mathbf{y}_i, \mathbf{z}) \, d\mathbf{z}. \end{aligned} \quad (\text{B.6})$$

$\underbrace{\hspace{15em}}_{Q_n(\boldsymbol{\beta}; \boldsymbol{\beta}')}$

We define the first term on the right-hand side of (B.6) to be $Q_n(\boldsymbol{\beta}; \boldsymbol{\beta}')$. Correspondingly, we define its expectation to be $Q(\boldsymbol{\beta}; \boldsymbol{\beta}')$. Note the second term on the right-hand side of (B.6) does not depend on $\boldsymbol{\beta}$. Hence, given some fixed $\boldsymbol{\beta}'$, we can maximize the lower bound function $Q_n(\boldsymbol{\beta}; \boldsymbol{\beta}')$ over $\boldsymbol{\beta}$ to obtain a sufficiently large $\ell_n(\boldsymbol{\beta}) - \ell_n(\boldsymbol{\beta}')$. Based on such an observation, at the t -th iteration of the classical EM algorithm, we evaluate $Q_n(\boldsymbol{\beta}; \boldsymbol{\beta}^{(t)})$ at the E-step and then perform $\max_{\boldsymbol{\beta}} Q_n(\boldsymbol{\beta}; \boldsymbol{\beta}^{(t)})$ at the M-step. See [18] for more details.

C High Dimensional EM Algorithm with Resampling

To simplify the technical analysis of the high dimensional algorithm, here we introduce its resampling version (Algorithm 4).

Algorithm 4 High Dimensional EM Algorithm with Resampling.

- 1: **Parameter:** Sparsity Parameter \widehat{s} , Maximum Number of Iterations T
 - 2: **Initialization:** $\widehat{\mathcal{S}}^{\text{init}} \leftarrow \text{supp}(\boldsymbol{\beta}^{\text{init}}, \widehat{s})$, $\boldsymbol{\beta}^{(0)} \leftarrow \text{trunc}(\boldsymbol{\beta}^{\text{init}}, \widehat{\mathcal{S}}^{\text{init}})$,
 $\{\text{supp}(\cdot, \cdot) \text{ and } \text{trunc}(\cdot, \cdot) \text{ are defined in (2.2) and (2.3)}\}$
Split the Dataset into T Subsets of Size n/T
 $\{\text{Without loss of generality, we assume } n/T \text{ is an integer}\}$
 - 3: **For** $t = 0$ to $T - 1$
 - 4: **E-step:** Evaluate $Q_{n/T}(\boldsymbol{\beta}; \boldsymbol{\beta}^{(t)})$ with the t -th Data Subset
 - 5: **M-step:** $\boldsymbol{\beta}^{(t+0.5)} \leftarrow M_{n/T}(\boldsymbol{\beta}^{(t)})$
 $\{M_{n/T}(\cdot)$ is implemented as in Algorithm 2 or 3 with $Q_n(\cdot; \cdot)$ replaced by $Q_{n/T}(\cdot; \cdot)\}$
 - 6: **T-step:** $\widehat{\mathcal{S}}^{(t+0.5)} \leftarrow \text{supp}(\boldsymbol{\beta}^{(t+0.5)}, \widehat{s})$, $\boldsymbol{\beta}^{(t+1)} \leftarrow \text{trunc}(\boldsymbol{\beta}^{(t+0.5)}, \widehat{\mathcal{S}}^{(t+0.5)})$
 - 7: **End For**
 - 8: **Output:** $\widehat{\boldsymbol{\beta}} \leftarrow \boldsymbol{\beta}^{(T)}$
-

D Decorrelated Score Statistic: An Intuitive Explanation

The intuition for the decorrelated score statistic in (2.7) can be understood as follows. Since $\ell_n(\boldsymbol{\beta})$ is the log-likelihood, its score function is $\nabla \ell_n(\boldsymbol{\beta})$ and the Fisher information at $\boldsymbol{\beta}^*$ is $I(\boldsymbol{\beta}^*) = -\mathbb{E}_{\boldsymbol{\beta}^*} [\nabla^2 \ell_n(\boldsymbol{\beta}^*)] / n$, where $\mathbb{E}_{\boldsymbol{\beta}^*}(\cdot)$ means the expectation is taken under the model with parameter $\boldsymbol{\beta}^*$. The following key theorem, which restates Theorem 2.1, reveals the connection of $\nabla_1 Q_n(\cdot; \cdot)$ in (2.5) and $T_n(\cdot)$ in (2.7) with the score function and Fisher information, which forms the foundation of our inferential method.

Theorem D.1. For the true parameter $\boldsymbol{\beta}^*$ and any $\boldsymbol{\beta} \in \mathbb{R}^d$, it holds that

$$\nabla_1 Q_n(\boldsymbol{\beta}; \boldsymbol{\beta}) = \nabla \ell_n(\boldsymbol{\beta}) / n, \quad \text{and} \quad \mathbb{E}_{\boldsymbol{\beta}^*} [T_n(\boldsymbol{\beta}^*)] = -I(\boldsymbol{\beta}^*) = \mathbb{E}_{\boldsymbol{\beta}^*} [\nabla^2 \ell_n(\boldsymbol{\beta}^*)] / n. \quad (\text{D.1})$$

Proof. See §I.1 for details. □

Recall that the log-likelihood $\ell_n(\boldsymbol{\beta})$ defined in (B.3) is difficult to evaluate due to the unobserved latent variable. Theorem D.1 provides a feasible way to calculate or estimate the corresponding score function and Fisher information, since $Q_n(\cdot; \cdot)$ and $T_n(\cdot)$ have closed forms. The geometric intuition behind Theorem D.1 can be understood as follows. By (B.5) and (B.6) we have

$$\ell_n(\boldsymbol{\beta}) \geq \ell_n(\boldsymbol{\beta}') + n \cdot Q_n(\boldsymbol{\beta}; \boldsymbol{\beta}') - \sum_{i=1}^n \int_{\mathcal{Z}} k_{\boldsymbol{\beta}'}(\mathbf{z} \mid \mathbf{y}_i) \cdot \log f_{\boldsymbol{\beta}'}(\mathbf{y}_i, \mathbf{z}) \, d\mathbf{z}. \quad (\text{D.2})$$

By (D.1), both sides of (D.2) have the same gradient with respect to $\boldsymbol{\beta}$ at $\boldsymbol{\beta}' = \boldsymbol{\beta}$. Furthermore, by (B.6), (D.2) becomes an equality for $\boldsymbol{\beta}' = \boldsymbol{\beta}$. Therefore, the lower bound function on the right-hand side of (D.2) is tangent to $\ell_n(\boldsymbol{\beta})$ at $\boldsymbol{\beta}' = \boldsymbol{\beta}$. Meanwhile, according to (2.4), $T_n(\boldsymbol{\beta})$ defines a modified curvature for the right-hand side of (D.2), which is obtained by taking derivative with respect to $\boldsymbol{\beta}$, then setting $\boldsymbol{\beta}' = \boldsymbol{\beta}$ and taking the second order derivative with respect to $\boldsymbol{\beta}$. The second equation in (D.1) shows that the obtained curvature equals the curvature of $\ell_n(\boldsymbol{\beta})$ at $\boldsymbol{\beta} = \boldsymbol{\beta}^*$ in expectation (up to a renormalization factor of n). Therefore, $\nabla_1 Q_n(\boldsymbol{\beta}; \boldsymbol{\beta})$ gives the score function and $T_n(\boldsymbol{\beta}^*)$ gives a good estimate of the Fisher information $I(\boldsymbol{\beta}^*)$. Since $\boldsymbol{\beta}^*$ is unknown in practice, later we will use $T_n(\widehat{\boldsymbol{\beta}})$ or $T_n(\widehat{\boldsymbol{\beta}}_0)$ to approximate $T_n(\boldsymbol{\beta}^*)$.

In the presence of the high dimensional nuisance parameter $\boldsymbol{\gamma}^* \in \mathbb{R}^{d-1}$, the classical score test is no longer applicable. In detail, the score test for $H_0 : \boldsymbol{\alpha}^* = 0$ relies on the following Taylor expansion

of the score function $\partial \ell_n(\cdot)/\partial \alpha$

$$\frac{1}{\sqrt{n}} \cdot \frac{\partial \ell_n(\bar{\beta}_0)}{\partial \alpha} = \frac{1}{\sqrt{n}} \cdot \frac{\partial \ell_n(\beta^*)}{\partial \alpha} + \frac{1}{\sqrt{n}} \cdot \frac{\partial^2 \ell_n(\beta^*)}{\partial \alpha \partial \gamma} \cdot (\bar{\gamma} - \gamma^*) + \bar{R}. \quad (\text{D.3})$$

Here $\beta^* = [0, (\gamma^*)^\top]^\top$, \bar{R} denotes the remainder term and $\bar{\beta}_0 = (0, \bar{\gamma}^\top)^\top$, where $\bar{\gamma}$ is an estimator of the nuisance parameter γ^* . The asymptotic normality of $1/\sqrt{n} \cdot \partial \ell_n(\bar{\beta}_0)/\partial \alpha$ in (D.3) relies on the fact that $1/\sqrt{n} \cdot \partial \ell_n(\beta_0^*)/\partial \alpha$ and $\sqrt{n} \cdot (\bar{\gamma} - \gamma^*)$ are jointly normal asymptotically and \bar{R} is $o_{\mathbb{P}}(1)$. In low dimensional settings, such a necessary condition holds for $\bar{\gamma}$ being the maximum likelihood estimator. However, in high dimensional settings, the maximum likelihood estimator cannot guarantee that \bar{R} is $o_{\mathbb{P}}(1)$, since $\|\bar{\gamma} - \gamma^*\|_2$ can be large due to the curse of dimensionality. Meanwhile, for $\bar{\gamma}$ being sparsity-type estimators, in general the asymptotic normality of $\sqrt{n} \cdot (\bar{\gamma} - \gamma^*)$ does not hold. For example, let $\bar{\gamma}$ be $\hat{\gamma}$, where $\hat{\gamma} \in \mathbb{R}^{d-1}$ is the subvector of $\hat{\beta}$, i.e., the estimator attained by the proposed high dimensional EM algorithm. Note that $\hat{\gamma}$ has many zero entries due to the truncation step. As $n \rightarrow \infty$, some entries of $\sqrt{n} \cdot (\hat{\gamma} - \gamma^*)$ have limiting distributions with point mass at zero. Clearly, this limiting distribution is not Gaussian with nonzero variance. In fact, for a similar setting of high dimensional linear regression, [15] illustrate that for γ^\sharp being a subvector of the Lasso estimator and γ^* being the corresponding subvector of the true parameter, the limiting distribution of $\sqrt{n} \cdot (\gamma^\sharp - \gamma^*)$ is not Gaussian.

The decorrelated score function defined in (2.5) successfully addresses the above issues. In detail, according to (D.1) in Theorem D.1 we have

$$\sqrt{n} \cdot S_n(\hat{\beta}_0, \lambda) = \frac{1}{\sqrt{n}} \cdot \frac{\partial \ell_n(\hat{\beta}_0)}{\partial \alpha} - \frac{1}{\sqrt{n}} \cdot w(\hat{\beta}_0, \lambda)^\top \cdot \frac{\partial \ell_n(\hat{\beta}_0)}{\partial \gamma}. \quad (\text{D.4})$$

Intuitively, if we replace $w(\hat{\beta}_0, \lambda)$ with $\mathbf{w} \in \mathbb{R}^{d-1}$ that satisfies

$$\mathbf{w}^\top \cdot \frac{\partial^2 \ell_n(\beta^*)}{\partial^2 \gamma} = \frac{\partial^2 \ell_n(\beta^*)}{\partial \alpha \partial \gamma}, \quad (\text{D.5})$$

we have the following Taylor expansion of the decorrelated score function

$$\begin{aligned} \frac{1}{\sqrt{n}} \cdot \frac{\partial \ell_n(\hat{\beta}_0)}{\partial \alpha} - \frac{\mathbf{w}^\top}{\sqrt{n}} \cdot \frac{\partial \ell_n(\hat{\beta}_0)}{\partial \gamma} &= \overbrace{\frac{1}{\sqrt{n}} \cdot \frac{\partial \ell_n(\beta^*)}{\partial \alpha} - \frac{\mathbf{w}^\top}{\sqrt{n}} \cdot \frac{\partial \ell_n(\beta^*)}{\partial \gamma}}^{(i)} \\ &+ \underbrace{\frac{1}{\sqrt{n}} \cdot \frac{\partial^2 \ell_n(\beta^*)}{\partial \alpha \partial \gamma} \cdot (\hat{\gamma} - \gamma^*) - \frac{\mathbf{w}^\top}{\sqrt{n}} \cdot \frac{\partial^2 \ell_n(\beta^*)}{\partial^2 \gamma} \cdot (\hat{\gamma} - \gamma^*)}_{(ii)} + \tilde{R}, \end{aligned} \quad (\text{D.6})$$

where term (ii) is zero by (D.5). Therefore, we no longer require the asymptotic normality of $\hat{\gamma} - \gamma^*$. Also, we will prove that the new remainder term \tilde{R} in (D.6) is $o_{\mathbb{P}}(1)$, since $\hat{\gamma}$ has a fast statistical rate of convergence. Now we only need to find the \mathbf{w} that satisfies (D.5). Nevertheless, it is difficult to calculate the second order derivatives in (D.5), because it is hard to evaluate $\ell_n(\cdot)$. According to (D.1), we use the submatrices of $T_n(\cdot)$ to approximate the derivatives in (D.5). Since $[T_n(\beta)]_{\gamma, \gamma}$ is not invertible in high dimensions, we use the Dantzig selector in (2.6) to approximately solve the linear system in (D.5). Based on this intuition, we can expect that $\sqrt{n} \cdot S_n(\hat{\beta}_0, \lambda)$ is asymptotically normal, since term (i) in (D.6) is a (rescaled) average of n i.i.d. random variables for which we can apply the central limit theorem. Besides, we will prove that $-[T_n(\hat{\beta}_0)]_{\alpha|\gamma}$ in (2.7) is a consistent estimator of $\sqrt{n} \cdot S_n(\hat{\beta}_0, \lambda)$'s asymptotic variance. Hence, we can expect that the decorrelated score statistic in (2.7) is asymptotically $N(0, 1)$.

From a high-level perspective, we can view $w(\hat{\beta}_0, \lambda)^\top \cdot \partial \ell_n(\hat{\beta}_0)/\partial \gamma$ in (D.4) as the projection of $\partial \ell_n(\hat{\beta}_0)/\partial \alpha$ onto the span of $\partial \ell_n(\hat{\beta}_0)/\partial \gamma$, where $w(\hat{\beta}_0, \lambda)$ is the projection coefficient. Intuitively, such a projection guarantees that in (D.4), $S_n(\hat{\beta}_0, \lambda)$ is orthogonal (uncorrelated) with $\partial \ell_n(\hat{\beta}_0)/\partial \gamma$, i.e., the score function with respect to the nuisance parameter γ . In this way, the projection corrects the effects of the high dimensional nuisance parameter. According to this intuition of decorrelation, we name $S_n(\hat{\beta}_0, \lambda)$ as the decorrelated score function.

E Implications for Specific Models: Computation and Estimation

To establish the corresponding results for specific models under the unified framework, we only need to establish Conditions 3.1-3.3 and determine the key quantities R , γ_1 , γ_2 , ν , μ , κ and ϵ . Recall that Conditions 3.1 and 3.2 and the models analyzed in our paper are identical to those in [2]. Meanwhile, note that Conditions 3.1 and 3.2 only involve the population version lower bound function $Q(\cdot; \cdot)$ and M-step $M(\cdot)$. Since [2] prove that the quantities in Conditions 3.1 and 3.2 are independent of the dimension d and sample size n , their corresponding results can be directly adapted. To establish the final results, it still remains to verify Condition 3.3 for each high dimensional latent variable model.

Gaussian Mixture Model: The following lemma, which is proved by [2], verifies Conditions 3.1 and 3.2 for Gaussian mixture model. Recall that σ is the standard deviation of each individual Gaussian distribution within the mixture.

Lemma E.1. Suppose that we have $\|\beta^*\|_2/\sigma \geq r$, where $r > 0$ is a sufficiently large constant that denotes the minimum signal-to-noise ratio. There exists some constant $C > 0$ such that Conditions *Lipschitz-Gradient-1*(γ_1, \mathcal{B}) and *Concavity-Smoothness*(μ, ν, \mathcal{B}) hold with

$$\gamma_1 = \exp(-C \cdot r^2), \quad \mu = \nu = 1, \quad \mathcal{B} = \{\beta : \|\beta - \beta^*\|_2 \leq R\} \quad \text{with } R = \kappa \cdot \|\beta^*\|_2, \quad \kappa = 1/4. \quad (\text{E.1})$$

Proof. See the proof of Corollary 1 in [2] for details. \square

Now we verify Condition 3.3 for the maximization implementation of the M-step (Algorithm 2).

Lemma E.2. For the maximization implementation of the M-step (Algorithm 2), we have that for a sufficiently large n and \mathcal{B} specified in (E.1), Condition *Statistical-Error*($\epsilon, \delta, s, n, \mathcal{B}$) holds with

$$\epsilon = C \cdot (\|\beta^*\|_\infty + \sigma) \cdot \sqrt{\frac{\log d + \log(2/\delta)}{n}}. \quad (\text{E.2})$$

Proof. See §H.4 for a detailed proof. \square

The next theorem establishes the implication of Theorem 3.4 for Gaussian mixture model.

Theorem E.3. We consider the maximization implementation of M-step (Algorithm 2). We assume $\|\beta^*\|_2/\sigma \geq r$ holds with a sufficiently large $r > 0$, and \mathcal{B} and R are as defined in (E.1). We suppose the initialization β^{init} of Algorithm 4 satisfies $\|\beta^{\text{init}} - \beta^*\|_2 \leq R/2$. Let the sparsity parameter \hat{s} be

$$\hat{s} = \left\lceil C' \cdot \max \left\{ 16 \cdot [\exp(C \cdot r^2) - 1]^{-2}, 100/9 \right\} \cdot s^* \right\rceil \quad (\text{E.3})$$

with C specified in (E.1) and $C' \geq 1$. Let the total number of iterations T in Algorithm 4 be

$$T = \left\lceil \frac{\log \left\{ C' \cdot R / [\Delta^{\text{GMM}}(s^*) \cdot \sqrt{\log d/n}] \right\}}{C \cdot r^2/2} \right\rceil, \quad (\text{E.4})$$

$$\text{where } \Delta^{\text{GMM}}(s^*) = (\sqrt{\hat{s}} + C'' \cdot \sqrt{s^*}) \cdot (\|\beta^*\|_\infty + \sigma).$$

Meanwhile, we suppose the dimension d is sufficiently large such that T in (E.4) is upper bounded by \sqrt{d} , and the sample size n is sufficiently large such that

$$C' \cdot \Delta^{\text{GMM}}(s^*) \cdot \sqrt{\frac{\log d \cdot T}{n}} \leq \min \left\{ [1 - \exp(-C \cdot r^2/2)]^2 \cdot R, 9/40 \cdot \|\beta^*\|_2 \right\}. \quad (\text{E.5})$$

We have that, with probability at least $1 - 2 \cdot d^{-1/2}$, the final estimator $\hat{\beta} = \beta^{(T)}$ satisfies

$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{C' \cdot \Delta^{\text{GMM}}(s^*)}{1 - \exp(-C \cdot r^2/2)} \cdot \sqrt{\frac{\log d \cdot T}{n}}. \quad (\text{E.6})$$

Proof. First we plug the quantities in (E.1) and (E.2) into Theorem 3.4. Recall that Theorem 3.4 requires Condition *Statistical-Error*($\epsilon, \delta/T, \hat{s}, n/T, \mathcal{B}$). Thus we need to replace δ and n with δ/T and n/T in the definition of ϵ in (E.2). Then we set $\delta = 2 \cdot d^{-1/2}$. Since T is specified in (E.4) and the dimension d is sufficiently large such that $T \leq \sqrt{d}$, we have $\log[2/(\delta/T)] \leq \log d$ in the definition of ϵ . By (E.3) and (E.5), we can then verify the assumptions in (3.8) and (3.9). Finally, by plugging in T in (E.4) into (3.10) and taking $t = T$, we can verify that in (3.9) the optimization error term is smaller than the statistical error term up to a constant factor. Therefore, we obtain (E.6). \square

To see the statistical rate of convergence with respect to s^* , d and n , for the moment we assume that R , r , $\|\beta^*\|_\infty$, $\|\beta^*\|_2$ and σ are constants. From (E.3) and (E.4), we obtain $\hat{s} = C \cdot s^*$ and therefore $\Delta^{\text{GMM}}(s^*) = C' \cdot \sqrt{s^*}$, which implies $T = C''' \cdot \log[C'' \cdot \sqrt{n/(s^* \cdot \log d)}]$. Hence, by (E.6) we have that, with high probability,

$$\|\hat{\beta} - \beta^*\|_2 \leq C \cdot \sqrt{\frac{s^* \cdot \log d \cdot \log n}{n}}.$$

Because the minimax lower bound for estimating an s^* -sparse d -dimensional vector is $\sqrt{s^* \cdot \log d/n}$, the rate of convergence in (E.6) is optimal up to a factor of $\log n$. Such a logarithmic factor results from the resampling scheme in Algorithm 4, since we only utilize n/T samples within each iteration. We expect that by directly analyzing Algorithm 1 we can eliminate such a logarithmic factor, which however incurs extra technical complexity for the analysis.

Mixture of Regression Model: The next lemma, proved by [2], verifies Conditions 3.1 and 3.2 for mixture of regression model. Recall that β^* is the regression coefficient and σ is the standard deviation of the Gaussian noise.

Lemma E.4. Suppose $\|\beta^*\|_2/\sigma \geq r$, where $r > 0$ is a sufficiently large constant that denotes the required minimum signal-to-noise ratio. Conditions *Lipschitz-Gradient-1*(γ_1, \mathcal{B}), *Lipschitz-Gradient-2*(γ_2, \mathcal{B}) and *Concavity-Smoothness*(μ, ν, \mathcal{B}) hold with

$$\begin{aligned} \gamma_1 \in (0, 1/2), \quad \gamma_2 \in (0, 1/4), \quad \mu = \nu = 1, \\ \mathcal{B} = \{\beta : \|\beta - \beta^*\|_2 \leq R\} \text{ with } R = \kappa \cdot \|\beta^*\|_2, \quad \kappa = 1/32. \end{aligned} \quad (\text{E.7})$$

Proof. See the proof of Corollary 3 in [2] for details. \square

The following lemma establishes Condition 3.3 for the two implementations of the M-step.

Lemma E.5. For \mathcal{B} specified in (E.7), we have the following results.

- For the maximization implementation of the M-step (line 5 of Algorithm 4), we have that Condition *Statistical-Error*($\epsilon, \delta, s, n, \mathcal{B}$) holds with

$$\epsilon = C \cdot \left[\max\{\|\beta^*\|_2^2 + \sigma^2, 1\} + \|\beta^*\|_2 \right] \cdot \sqrt{\frac{\log d + \log(4/\delta)}{n}} \quad (\text{E.8})$$

for sufficiently large sample size n and constant $C > 0$.

- For the gradient ascent implementation, Condition *Statistical-Error*($\epsilon, \delta, s, n, \mathcal{B}$) holds with

$$\epsilon = C \cdot \eta \cdot \max\{\|\beta^*\|_2^2 + \sigma^2, 1, \sqrt{s} \cdot \|\beta^*\|_2\} \cdot \sqrt{\frac{\log d + \log(4/\delta)}{n}} \quad (\text{E.9})$$

for sufficiently large sample size n and $C > 0$, where η denotes the stepsize in Algorithm 3.

Proof. See §H.5 for a detailed proof. \square

The next theorem establishes the implication of Theorem 3.4 for mixture of regression model.

Theorem E.6. Let $\|\beta^*\|_2/\sigma \geq r$ with $r > 0$ sufficiently large. Assuming that \mathcal{B} and R are specified in (E.7) and the initialization β^{init} satisfies $\|\beta^{\text{init}} - \beta^*\|_2 \leq R/2$, we have the following results.

- For the maximization implementation of the M-step (Algorithm 2), let \hat{s} and T be

$$\hat{s} = \lceil C \cdot \max\{16, 132/31\} \cdot s^* \rceil, \quad T = \left\lceil \frac{\log\{C' \cdot R / [\Delta_1^{\text{MR}}(s^*) \cdot \sqrt{\log d/n}]\}}{\log \sqrt{2}} \right\rceil,$$

$$\text{where } \Delta_1^{\text{MR}}(s^*) = (\sqrt{\hat{s}} + C'' \cdot \sqrt{s^*}) \cdot \left[\max\{\|\beta^*\|_2^2 + \sigma^2, 1\} + \|\beta^*\|_2 \right], \quad \text{and } C \geq 1.$$

We suppose d and n are sufficiently large such that $T \leq \sqrt{d}$ and

$$C \cdot \Delta_1^{\text{MR}}(s^*) \cdot \sqrt{\frac{\log d \cdot T}{n}} \leq \min\left\{(1 - 1/\sqrt{2})^2 \cdot R, 3/8 \cdot \|\beta^*\|_2\right\}.$$

Then with probability at least $1 - 4 \cdot d^{-1/2}$, the final estimator $\hat{\beta} = \beta^{(T)}$ satisfies

$$\|\hat{\beta} - \beta^*\|_2 \leq C' \cdot \Delta_1^{\text{MR}}(s^*) \cdot \sqrt{\frac{\log d \cdot T}{n}}. \quad (\text{E.10})$$

- For the gradient ascent implementation of the M-step (Algorithm 3) with stepsize set to $\eta = 1$, let \hat{s} and T be

$$\hat{s} = \lceil C \cdot \max\{16/9, 132/31\} \cdot s^* \rceil, \quad T = \left\lceil \frac{\log\{C' \cdot R / [\Delta_2^{\text{MR}}(s^*) \cdot \sqrt{\log d/n}]\}}{\log 2} \right\rceil,$$

$$\text{where } \Delta_2^{\text{MR}}(s^*) = (\sqrt{\hat{s}} + C'' \cdot \sqrt{s^*}) \cdot \max\{\|\beta^*\|_2^2 + \sigma^2, 1, \sqrt{\hat{s}} \cdot \|\beta^*\|_2\}, \quad \text{and } C \geq 1.$$

We suppose d and n are sufficiently large such that $T \leq \sqrt{d}$ and

$$C \cdot \Delta_2^{\text{MR}}(s^*) \cdot \sqrt{\frac{\log d \cdot T}{n}} \leq \min\{R/4, 3/8 \cdot \|\beta^*\|_2\}.$$

Then with probability at least $1 - 4 \cdot d^{-1/2}$, the final estimator $\hat{\beta} = \beta^{(T)}$ satisfies

$$\|\hat{\beta} - \beta^*\|_2 \leq C' \cdot \Delta_2^{\text{MR}}(s^*) \cdot \sqrt{\frac{\log d \cdot T}{n}}. \quad (\text{E.11})$$

Proof. The proof is similar to Theorem E.3. \square

To understand the intuition of Theorem E.6, we suppose that $\|\beta^*\|_2$, σ , R and r are constants, which implies $\hat{s} = C \cdot s^*$ and $\Delta_1^{\text{MR}}(s^*) = C' \cdot \sqrt{s^*}$, $\Delta_2^{\text{MR}}(s^*) = C'' \cdot s^*$. Therefore, for the maximization and gradient ascent implementations of the M-step, we have $T = C' \cdot \log[n/(s^* \cdot \log d)]$ and $T = C'' \cdot \log\{n/[(s^*)^2 \cdot \log d]\}$ correspondingly. Hence, by (E.10) in Theorem E.6 we have that, for the maximization implementation of the M-step,

$$\|\hat{\beta} - \beta^*\|_2 \leq C \cdot \sqrt{\frac{s^* \cdot \log d \cdot \log n}{n}} \quad (\text{E.12})$$

holds with high probability. Meanwhile, from (E.11) in Theorem E.6 we have that, for the gradient ascent implementation of the M-step,

$$\|\hat{\beta} - \beta^*\|_2 \leq C' \cdot s^* \cdot \sqrt{\frac{\log d \cdot \log n}{n}} \quad (\text{E.13})$$

holds with high probability. The statistical rates in (E.12) and (E.13) attain the $\sqrt{s^* \cdot \log d/n}$ minimax lower bound up to factors of $\sqrt{\log n}$ and $\sqrt{s^* \cdot \log n}$ respectively and are therefore near-optimal. Note that the statistical rate of convergence attained by the gradient ascent implementation of the M-step is slower by a $\sqrt{s^*}$ factor than the rate of the maximization implementation. However, our discussion in §A illustrates that, for mixture of regression model, the gradient ascent implementation does not involve estimating the inverse covariance of \mathbf{X} in (A.4). Hence, the gradient ascent implementation is more computationally efficient, and is also applicable to the settings in which \mathbf{X} has more general covariance structures.

F Implications for Specific Models: Inference

To establish the high dimensional inference results for each model, we only need to verify Conditions 4.1-4.4 and determine the key quantities ζ^{EM} , ζ^{G} , ζ^{T} and ζ^{L} . In the following, we focus on Gaussian mixture and mixture of regression models.

Gaussian Mixture Model: The following lemma verifies Conditions 4.1 and 4.2.

Lemma F.1. We have that Conditions 4.1 and 4.2 hold with

$$\zeta^{\text{EM}} = \frac{\sqrt{\hat{s}} \cdot \Delta^{\text{GMM}}(s^*)}{1 - \exp(-C \cdot r^2/2)} \cdot \sqrt{\frac{\log d \cdot T}{n}}, \quad \text{and } \zeta^{\text{G}} = (\|\beta^*\|_\infty + \sigma) \cdot \sqrt{\frac{\log d}{n}},$$

where \hat{s} , $\Delta^{\text{GMM}}(s^*)$, r and T are as defined in Theorem E.3.

Proof. See §I.5 for a detailed proof. \square

By our discussion that follows Theorem E.3, we have that \hat{s} and $\Delta^{\text{GMM}}(s^*)$ are of the same order as s^* and $\sqrt{s^*}$ respectively, and T is roughly of the order $\sqrt{\log n}$. Therefore, ζ^{EM} is roughly of the order $s^* \cdot \sqrt{\log d/n \cdot \log n}$. The following lemma verifies Condition 4.3 for Gaussian mixture model.

Lemma F.2. We have that Condition 4.3 holds with

$$\zeta^T = (\|\beta^*\|_\infty^2 + \sigma^2)/\sigma^2 \cdot \sqrt{\frac{\log d}{n}}.$$

Proof. See §I.6 for a detailed proof. \square

The following lemma establishes Condition 4.4 for Gaussian mixture model.

Lemma F.3. We have that Condition 4.4 holds with

$$\zeta^L = (\|\beta^*\|_\infty^2 + \sigma^2)^{3/2}/\sigma^4 \cdot (\log d + \log n)^{3/2}.$$

Proof. See §I.7 for a detailed proof. \square

Equipped with Lemmas F.1-F.3, we establish the inference results for Gaussian mixture model.

Theorem F.4. Under Assumption 4.5, we have that for $n \rightarrow \infty$, (4.7) holds for Gaussian mixture model.

In fact, for Gaussian mixture model we can make (4.6) in Assumption 4.5 more transparent by plugging in ζ^{EM} , ζ^{G} , ζ^T and ζ^L specified above. Particularly, for simplicity we assume all quantities except $s_{\mathbf{w}}^*$, s^* , d and n are constants. Then we can verify that (4.6) holds if

$$\max\{s_{\mathbf{w}}^*, s^*\}^2 \cdot (s^*)^2 \cdot (\log d)^5 = o[n/(\log n)^2]. \quad (\text{F.1})$$

According to the discussion following Theorem E.3, we require $s^* \cdot \log d = o(n/\log n)$ for the estimator $\hat{\beta}$ to be consistent. In comparison, (F.1) illustrates that high dimensional inference requires a higher sample complexity than parameter estimation. In the context of high dimensional generalized linear models, [26, 32] also observe the same phenomenon.

Mixture of Regression Model: The following lemma verifies Conditions 4.1 and 4.2. Recall that \hat{s} , T , $\Delta_1^{\text{MR}}(s^*)$ and $\Delta_2^{\text{MR}}(s^*)$ are defined in Theorem E.6, while σ denotes the standard deviation of the Gaussian noise in mixture of regression model.

Lemma F.5. We have that Conditions 4.1 and 4.2 hold with

$$\zeta^{\text{EM}} = \sqrt{\hat{s}} \cdot \Delta^{\text{MR}}(s^*) \cdot \sqrt{\frac{\log d \cdot T}{n}}, \quad \text{and} \quad \zeta^{\text{G}} = \max\{\|\beta^*\|_2^2 + \sigma^2, 1, \sqrt{s^*} \cdot \|\beta^*\|_2\} \cdot \sqrt{\frac{\log d}{n}},$$

where we have $\Delta^{\text{MR}}(s^*) = \Delta_1^{\text{MR}}(s^*)$ for the maximization implementation of the M-step (Algorithm 2), and $\Delta^{\text{MR}}(s^*) = \Delta_2^{\text{MR}}(s^*)$ for the gradient ascent implementation of the M-step (Algorithm 3).

Proof. See §I.8 for a detailed proof. \square

By our discussion that follows Theorem E.6, we have that \hat{s} is of the same order as s^* . For the maximization implementation of the M-step (Algorithm 2), we have that $\Delta^{\text{MR}}(s^*) = \Delta_1^{\text{MR}}(s^*)$ is of the same order as $\sqrt{s^*}$. Meanwhile, for the gradient ascent implementation in Algorithm 3, we have that $\Delta^{\text{MR}}(s^*) = \Delta_2^{\text{MR}}(s^*)$ is of the same order as s^* . Hence, ζ^{EM} is of the order $s^* \cdot \sqrt{\log d/n \cdot \log n}$ or $(s^*)^{3/2} \cdot \sqrt{\log d/n \cdot \log n}$ correspondingly, since T is roughly of the order $\sqrt{\log n}$. The next lemma establishes Condition 4.3 for mixture of regression model.

Lemma F.6. We have that Condition 4.3 holds with

$$\zeta^T = (\log n + \log d) \cdot [(\log n + \log d) \cdot \|\beta^*\|_1^2 + \sigma^2]/\sigma^2 \cdot \sqrt{\frac{\log d}{n}}.$$

Proof. See §I.9 for a detailed proof. \square

The following lemma establishes Condition 4.4 for mixture of regression model.

Lemma F.7. We have that Condition 4.4 holds with

$$\zeta^L = (\|\beta^*\|_1 + \sigma)^3 \cdot (\log n + \log d)^3/\sigma^4.$$

Proof. See §I.10 for a detailed proof. \square

Equipped with Lemmas F.5-F.7, we are now ready to establish the high dimensional inference results for mixture of regression model.

Theorem F.8. For mixture of regression model, under Assumption 4.5, (4.7) holds as $n \rightarrow \infty$.

Similar to the discussion that follows Theorem F.4, we can make (4.6) in Assumption 4.5 more explicit by plugging in $\zeta^{\text{EM}}, \zeta^{\text{G}}, \zeta^{\text{T}}$ and ζ^{L} specified in Lemmas F.5-F.7. Assuming all quantities except $s_{\mathbf{w}}^*, s^*, d$ and n are constants, we have that (4.6) holds if

$$\max\{s_{\mathbf{w}}^*, s^*\}^2 \cdot (s^*)^4 \cdot (\log d)^8 = o[n/(\log n)^2].$$

In contrast, for high dimensional estimation, we only require $(s^*)^2 \cdot \log d = o(n/\log n)$ to ensure the consistency of $\hat{\beta}$ by our discussion following Theorem E.6.

G Proof of Main Results

We lay out a proof sketch of the main theory. First we prove the results in Theorem 3.4 for parameter estimation and computation. Then we establish the results in Theorem 4.6 for inference.

G.1 Proof of Results for Computation and Estimation

Proof of Theorem 3.4: First we introduce some notations. Recall that the $\text{trunc}(\cdot, \cdot)$ function is defined in (2.3). We define $\bar{\beta}^{(t+0.5)}, \bar{\beta}^{(t+1)} \in \mathbb{R}^d$ as

$$\bar{\beta}^{(t+0.5)} = M(\beta^{(t)}), \quad \bar{\beta}^{(t+1)} = \text{trunc}(\bar{\beta}^{(t+0.5)}, \hat{\mathcal{S}}^{(t+0.5)}). \quad (\text{G.1})$$

As defined in (3.1) or (3.2), $M(\cdot)$ is the population version M-step with the maximization or gradient ascent implementation. Here $\hat{\mathcal{S}}^{(t+0.5)}$ denotes the set of index j 's with the top \hat{s} largest $|\beta_j^{(t+0.5)}|$'s. It is worth noting $\hat{\mathcal{S}}^{(t+0.5)}$ is calculated based on $\beta^{(t+0.5)}$ in the truncation step (line 6 of Algorithm 4), rather than based on $\bar{\beta}^{(t+0.5)}$ defined in (G.1).

Our goal is to characterize the relationship between $\|\beta^{(t+1)} - \beta^*\|_2$ and $\|\beta^{(t)} - \beta^*\|_2$. According to the definition of the truncation step (line 6 of Algorithm 4) and triangle inequality, we have

$$\begin{aligned} \|\beta^{(t+1)} - \beta^*\|_2 &= \left\| \text{trunc}(\beta^{(t+0.5)}, \hat{\mathcal{S}}^{(t+0.5)}) - \beta^* \right\|_2 \\ &\leq \left\| \text{trunc}(\beta^{(t+0.5)}, \hat{\mathcal{S}}^{(t+0.5)}) - \text{trunc}(\bar{\beta}^{(t+0.5)}, \hat{\mathcal{S}}^{(t+0.5)}) \right\|_2 + \left\| \text{trunc}(\bar{\beta}^{(t+0.5)}, \hat{\mathcal{S}}^{(t+0.5)}) - \beta^* \right\|_2 \\ &= \underbrace{\left\| \text{trunc}(\beta^{(t+0.5)}, \hat{\mathcal{S}}^{(t+0.5)}) - \text{trunc}(\bar{\beta}^{(t+0.5)}, \hat{\mathcal{S}}^{(t+0.5)}) \right\|_2}_{(i)} + \underbrace{\left\| \bar{\beta}^{(t+1)} - \beta^* \right\|_2}_{(ii)}, \end{aligned} \quad (\text{G.2})$$

where the last equality is obtained from (G.1). According to the definition of the $\text{trunc}(\cdot, \cdot)$ function in (2.3), the two terms within term (i) both have support $\hat{\mathcal{S}}^{(t+0.5)}$ with cardinality \hat{s} . Thus, we have

$$\begin{aligned} \left\| \text{trunc}(\beta^{(t+0.5)}, \hat{\mathcal{S}}^{(t+0.5)}) - \text{trunc}(\bar{\beta}^{(t+0.5)}, \hat{\mathcal{S}}^{(t+0.5)}) \right\|_2 &= \left\| (\beta^{(t+0.5)} - \bar{\beta}^{(t+0.5)})_{\hat{\mathcal{S}}^{(t+0.5)}} \right\|_2 \\ &\leq \sqrt{\hat{s}} \cdot \left\| (\beta^{(t+0.5)} - \bar{\beta}^{(t+0.5)})_{\hat{\mathcal{S}}^{(t+0.5)}} \right\|_{\infty} \\ &\leq \sqrt{\hat{s}} \cdot \|\beta^{(t+0.5)} - \bar{\beta}^{(t+0.5)}\|_{\infty}. \end{aligned} \quad (\text{G.3})$$

Since we have $\beta^{(t+0.5)} = M_n(\beta^{(t)})$ and $\bar{\beta}^{(t+0.5)} = M(\beta^{(t)})$, we can further establish an upper bound for the right-hand side by invoking Condition 3.3.

Our subsequent proof will establish an upper bound for term (ii) in (G.2) in two steps. We first characterize the relationship between $\|\bar{\beta}^{(t+1)} - \beta^*\|_2$ and $\|\bar{\beta}^{(t+0.5)} - \beta^*\|_2$ and then the relationship between $\|\bar{\beta}^{(t+0.5)} - \beta^*\|_2$ and $\|\beta^{(t)} - \beta^*\|_2$. The next lemma accomplishes the first step. Recall that \hat{s} is the sparsity parameter in Algorithm 4, while s^* is the sparsity level of the true parameter β^* .

Lemma G.1. Suppose that we have

$$\|\bar{\beta}^{(t+0.5)} - \beta^*\|_2 \leq \kappa \cdot \|\beta^*\|_2 \quad (\text{G.4})$$

for some $\kappa \in (0, 1)$. Assuming that we have

$$\hat{s} \geq \frac{4 \cdot (1 + \kappa)^2}{(1 - \kappa)^2} \cdot s^*, \quad \text{and} \quad \sqrt{\hat{s}} \cdot \|\beta^{(t+0.5)} - \bar{\beta}^{(t+0.5)}\|_{\infty} \leq \frac{(1 - \kappa)^2}{2 \cdot (1 + \kappa)} \cdot \|\beta^*\|_2, \quad (\text{G.5})$$

then it holds that

$$\|\bar{\beta}^{(t+1)} - \beta^*\|_2 \leq \frac{C \cdot \sqrt{s^*}}{\sqrt{1-\kappa}} \cdot \|\beta^{(t+0.5)} - \bar{\beta}^{(t+0.5)}\|_\infty + (1 + 4 \cdot \sqrt{s^*/\hat{s}})^{1/2} \cdot \|\bar{\beta}^{(t+0.5)} - \beta^*\|_2. \quad (\text{G.6})$$

Proof. The proof is based on fine-grained analysis of the relationship between $\widehat{\mathcal{S}}^{(t+0.5)}$ and the true support \mathcal{S}^* . In particular, we focus on three index sets, namely, $\mathcal{I}_1 = \mathcal{S}^* \setminus \widehat{\mathcal{S}}^{(t+0.5)}$, $\mathcal{I}_2 = \mathcal{S}^* \cap \widehat{\mathcal{S}}^{(t+0.5)}$ and $\mathcal{I}_3 = \widehat{\mathcal{S}}^{(t+0.5)} \setminus \mathcal{S}^*$. Among them, \mathcal{I}_2 characterizes the similarity between $\widehat{\mathcal{S}}^{(t+0.5)}$ and \mathcal{S}^* , while \mathcal{I}_1 and \mathcal{I}_3 characterize their difference. The key proof strategy is to represent the three distances in (G.6) with the ℓ_2 -norms of the restrictions of $\bar{\beta}^{(t+0.5)}$ and β^* on the three index sets. In particular, we focus on $\|\bar{\beta}_{\mathcal{I}_1}^{(t+0.5)}\|_2$ and $\|\beta_{\mathcal{I}_1}^*\|_2$. In order to quantify these ℓ_2 -norms, we establish a fine-grained characterization for the absolute values of $\bar{\beta}^{(t+0.5)}$'s entries that are selected and eliminated within the truncation operation $\bar{\beta}^{(t+1)} \leftarrow \text{trunc}(\bar{\beta}^{(t+0.5)}, \widehat{\mathcal{S}}^{(t+0.5)})$. See §H.1 for a detailed proof. \square

Lemma G.1 is central to the proof of Theorem 3.4. In detail, the assumption in (G.4) guarantees $\bar{\beta}^{(t+0.5)}$ is within the basin of attraction. In (G.5), the first assumption ensures the sparsity parameter \hat{s} in Algorithm 4 is set to be sufficiently large, while second ensures the statistical error is sufficiently small. These assumptions will be verified in the proof of Theorem 3.4. The intuition behind (G.6) is as follows. Let $\bar{\mathcal{S}}^{(t+0.5)} = \text{supp}(\bar{\beta}^{(t+0.5)}, \hat{s})$, where $\text{supp}(\cdot, \cdot)$ is defined in (2.2). By triangle inequality, the left-hand side of (G.6) satisfies

$$\|\bar{\beta}^{(t+1)} - \beta^*\|_2 \leq \underbrace{\|\bar{\beta}^{(t+1)} - \text{trunc}(\bar{\beta}^{(t+0.5)}, \bar{\mathcal{S}}^{(t+0.5)})\|_2}_{(i)} + \underbrace{\|\text{trunc}(\bar{\beta}^{(t+0.5)}, \bar{\mathcal{S}}^{(t+0.5)}) - \beta^*\|_2}_{(ii)}. \quad (\text{G.7})$$

Intuitively, the two terms on right-hand side of (G.6) reflect terms (i) and (ii) in (G.7) correspondingly. In detail, for term (i) in (G.7), recall that according to (G.1) and line 6 of Algorithm 4 we have

$$\bar{\beta}^{(t+1)} = \text{trunc}(\bar{\beta}^{(t+0.5)}, \widehat{\mathcal{S}}^{(t+0.5)}), \quad \text{where } \widehat{\mathcal{S}}^{(t+0.5)} = \text{supp}(\beta^{(t+0.5)}, \hat{s}).$$

As the sample size n is sufficiently large, $\bar{\beta}^{(t+0.5)}$ and $\beta^{(t+0.5)}$ are close, since they are attained by the population version and sample version M-steps correspondingly. Hence, $\bar{\mathcal{S}}^{(t+0.5)} = \text{supp}(\bar{\beta}^{(t+0.5)}, \hat{s})$ and $\widehat{\mathcal{S}}^{(t+0.5)} = \text{supp}(\beta^{(t+0.5)}, \hat{s})$ should be similar. Thus, in term (i), $\bar{\beta}^{(t+1)} = \text{trunc}(\bar{\beta}^{(t+0.5)}, \widehat{\mathcal{S}}^{(t+0.5)})$ should be close to $\text{trunc}(\bar{\beta}^{(t+0.5)}, \bar{\mathcal{S}}^{(t+0.5)})$ up to some statistical error, which is reflected by the first term on the right-hand side of (G.6).

Also, we turn to quantify the relationship between $\|\bar{\beta}^{(t+0.5)} - \beta^*\|_2$ in (G.6) and term (ii) in (G.7). The truncation in term (ii) preserves the top \hat{s} coordinates of $\bar{\beta}^{(t+0.5)}$ with the largest magnitudes while setting others to zero. Intuitively speaking, the truncation incurs additional error to $\bar{\beta}^{(t+0.5)}$'s distance to β^* . Meanwhile, note that when $\bar{\beta}^{(t+0.5)}$ is close to β^* , $\bar{\mathcal{S}}^{(t+0.5)}$ is similar to \mathcal{S}^* . Therefore, the incurred error can be controlled, because the truncation does not eliminate many relevant entries. In particular, as shown in the second term on the right-hand side of (G.6), such incurred error decays as \hat{s} increases, since in this case $\widehat{\mathcal{S}}^{(t+0.5)}$ includes more entries. According to the discussion for term (i) in (G.7), $\bar{\mathcal{S}}^{(t+0.5)}$ is similar to $\widehat{\mathcal{S}}^{(t+0.5)}$, which implies that $\bar{\mathcal{S}}^{(t+0.5)}$ should also cover more entries. Thus, fewer relevant entries are wrongly eliminated by the truncation and hence the incurred error is smaller. The extreme case is that, when $\hat{s} \rightarrow \infty$, term (ii) in (G.7) becomes $\|\bar{\beta}^{(t+0.5)} - \beta^*\|_2$, which indicates that no additional error is incurred by the truncation. Correspondingly, the second term on the right-hand side of (G.6) also becomes $\|\bar{\beta}^{(t+0.5)} - \beta^*\|_2$.

Next, we turn to characterize the relationship between $\|\bar{\beta}^{(t+0.5)} - \beta^*\|_2$ and $\|\beta^{(t)} - \beta^*\|_2$. Recall $\bar{\beta}^{(t+0.5)} = M(\beta^{(t)})$ is defined in (G.1). The next lemma, which is adapted from Theorems 1 and 3 in [2], characterizes the contraction property of the population version M-step defined in (3.1) or (3.2).

Lemma G.2. Under the assumptions of Theorem 3.4, the following results hold for $\beta^{(t)} \in \mathcal{B}$.

- For the maximization implementation of the M-step (Algorithm 2), we have

$$\|\bar{\beta}^{(t+0.5)} - \beta^*\|_2 \leq (\gamma_1/\nu) \cdot \|\beta^{(t)} - \beta^*\|_2. \quad (\text{G.8})$$

- For the gradient ascent implementation of the M-step (Algorithm 3), we have

$$\|\bar{\beta}^{(t+0.5)} - \beta^*\|_2 \leq \left(1 - 2 \cdot \frac{\nu - \gamma_2}{\nu + \mu}\right) \cdot \|\beta^{(t)} - \beta^*\|_2. \quad (\text{G.9})$$

Here γ_1, γ_2, μ and ν are defined in Conditions 3.1 and 3.2.

Proof. The proof strategy is to first characterize the M-step using $Q(\cdot; \beta^*)$. According to Condition *Concavity-Smoothness*(μ, ν, \mathcal{B}), $-Q(\cdot; \beta^*)$ is ν -strongly convex and μ -smooth, and thus enjoys desired optimization guarantees. Then Condition *Lipschitz-Gradient-1*(γ_1, \mathcal{B}) or *Lipschitz-Gradient-2*(γ_2, \mathcal{B}) is invoked to characterize the difference between $Q(\cdot; \beta^*)$ and $Q(\cdot; \beta^{(t)})$. We provide the proof in §H.2 for the sake of completeness. \square

Equipped with Lemmas G.1 and G.2, we are now ready to prove Theorem 3.4.

Proof. To unify the subsequent proof for the maximization and gradient implementations of the M-step, we employ $\rho \in (0, 1)$ to denote $\rho_1 := \gamma_1/\nu$ in (G.8) or $\rho_2 := 1 - 2 \cdot (\nu - \gamma_2)/(\nu + \mu)$ in (G.9). In the following we stick to the former one to avoid confusion. The proof for the latter one is exactly the same. By the definitions of $\bar{\beta}^{(t+0.5)}$ and $\beta^{(t+0.5)}$ in (G.1) and Algorithm 4, Condition *Statistical-Error*($\epsilon, \delta/T, \hat{s}, n/T, \mathcal{B}$) implies

$$\|\beta^{(t+0.5)} - \bar{\beta}^{(t+0.5)}\|_\infty = \|M_{n/T}(\beta^{(t)}) - M(\beta^{(t)})\|_\infty \leq \epsilon$$

holds with probability at least $1 - \delta/T$. Then by taking union bound we have that, the event

$$\mathcal{E} = \left\{ \|\beta^{(t+0.5)} - \bar{\beta}^{(t+0.5)}\|_\infty \leq \epsilon, \text{ for all } t \in \{0, \dots, T-1\} \right\} \quad (\text{G.10})$$

occurs with probability at least $1 - \delta$. Conditioning on \mathcal{E} , in the following we prove that

$$\|\beta^{(t)} - \beta^*\|_2 \leq \frac{(\sqrt{\hat{s}} + C/\sqrt{1-\kappa} \cdot \sqrt{s^*}) \cdot \epsilon}{1 - \sqrt{\rho}} + \rho^{t/2} \cdot \|\beta^{(0)} - \beta^*\|_2, \quad \text{for all } t \in \{1, \dots, T\} \quad (\text{G.11})$$

by mathematical induction.

Before we lay out the proof, we first prove $\beta^{(0)} \in \mathcal{B}$. Recall β^{init} is the initialization of Algorithm 4 and R is the radius of the basin of attraction \mathcal{B} . By the assumption of Theorem 3.4, we have

$$\|\beta^{\text{init}} - \beta^*\|_2 \leq R/2. \quad (\text{G.12})$$

Therefore, (G.12) implies $\|\beta^{\text{init}} - \beta^*\|_2 < \kappa \cdot \|\beta^*\|_2$ since $R = \kappa \cdot \|\beta^*\|_2$. Invoking the auxiliary result in Lemma H.1, we obtain

$$\|\beta^{(0)} - \beta^*\|_2 \leq (1 + 4 \cdot \sqrt{s^*/\hat{s}})^{1/2} \cdot \|\beta^{\text{init}} - \beta^*\|_2 \leq (1 + 4 \cdot \sqrt{1/4})^{1/2} \cdot R/2 < R. \quad (\text{G.13})$$

Here the second inequality is from (G.12) as well as the assumption in (3.8), which implies $s^*/\hat{s} \leq (1 - \kappa)^2 / [4 \cdot (1 + \kappa)^2] \leq 1/4$. Thus, (G.13) implies $\beta^{(0)} \in \mathcal{B}$. In the sequel, we prove that (G.11) holds for $t = 1$. By invoking Lemma G.2 and setting $t = 0$ in (G.8), we obtain

$$\|\bar{\beta}^{(0.5)} - \beta^*\|_2 \leq \rho \cdot \|\beta^{(0)} - \beta^*\|_2 \leq \rho \cdot R < R = \kappa \cdot \|\beta^*\|_2,$$

where the second inequality is from (G.13). Hence, the assumption in (G.4) of Lemma G.1 holds for $\bar{\beta}^{(0.5)}$. Furthermore, by the assumptions in (3.8) and (3.9) of Theorem 3.4, we can also verify that the assumptions in (G.5) of Lemma G.1 hold conditioning on the event \mathcal{E} defined in (G.10). By invoking Lemma G.1 we have that (G.6) holds for $t = 0$. Further plugging $\|\beta^{(t+0.5)} - \bar{\beta}^{(t+0.5)}\|_\infty \leq \epsilon$ in (G.10) into (G.6) with $t = 0$, we obtain

$$\|\bar{\beta}^{(1)} - \beta^*\|_2 \leq \frac{C \cdot \sqrt{s^*}}{\sqrt{1-\kappa}} \cdot \epsilon + (1 + 4 \cdot \sqrt{s^*/\hat{s}})^{1/2} \cdot \|\bar{\beta}^{(0.5)} - \beta^*\|_2. \quad (\text{G.14})$$

Setting $t = 0$ in (G.8) of Lemma G.2 and then plugging (G.8) into (G.14), we obtain

$$\|\bar{\beta}^{(1)} - \beta^*\|_2 \leq \frac{C \cdot \sqrt{s^*}}{\sqrt{1-\kappa}} \cdot \epsilon + (1 + 4 \cdot \sqrt{s^*/\hat{s}})^{1/2} \cdot \rho \cdot \|\beta^{(0)} - \beta^*\|_2. \quad (\text{G.15})$$

For $t = 0$, plugging (G.3) into term (i) in (G.2), and (G.15) into term (ii) in (G.2), and then applying $\|\beta^{(t+0.5)} - \bar{\beta}^{(t+0.5)}\|_\infty \leq \epsilon$ with $t = 0$ in (G.10), we obtain

$$\begin{aligned} \|\beta^{(1)} - \beta^*\|_2 &\leq \sqrt{\hat{s}} \cdot \|\beta^{(0.5)} - \bar{\beta}^{(0.5)}\|_\infty + \frac{C \cdot \sqrt{s^*}}{\sqrt{1-\kappa}} \cdot \epsilon + (1 + 4 \cdot \sqrt{s^*/\hat{s}})^{1/2} \cdot \rho \cdot \|\beta^{(0)} - \beta^*\|_2 \\ &\leq (\sqrt{\hat{s}} + C/\sqrt{1-\kappa} \cdot \sqrt{s^*}) \cdot \epsilon + (1 + 4 \cdot \sqrt{s^*/\hat{s}})^{1/2} \cdot \rho \cdot \|\beta^{(0)} - \beta^*\|_2. \end{aligned} \quad (\text{G.16})$$

By our assumption that $\hat{s} \geq 16 \cdot (1/\rho - 1)^{-2} \cdot s^*$ in (3.8), we have $(1 + 4 \cdot \sqrt{s^*/\hat{s}})^{1/2} \leq 1/\sqrt{\rho}$ in (G.16). Hence, from (G.16) we obtain

$$\|\beta^{(1)} - \beta^*\|_2 \leq (\sqrt{\hat{s}} + C/\sqrt{1-\kappa} \cdot \sqrt{s^*}) \cdot \epsilon + \sqrt{\rho} \cdot \|\beta^{(0)} - \beta^*\|_2, \quad (\text{G.17})$$

which implies that (G.11) holds for $t = 1$, since we have $1 - \sqrt{\rho} < 1$ in (G.11).

Suppose we have that (G.11) holds for some $t \geq 1$. By (G.11) we have

$$\begin{aligned} \|\beta^{(t)} - \beta^*\|_2 &\leq \frac{(\sqrt{\hat{s}} + C/\sqrt{1-\kappa} \cdot \sqrt{s^*}) \cdot \epsilon}{1 - \sqrt{\rho}} + \rho^{t/2} \cdot \|\beta^{(0)} - \beta^*\|_2 \\ &\leq (1 - \sqrt{\rho}) \cdot R + \sqrt{\rho} \cdot R = R, \end{aligned} \quad (\text{G.18})$$

where the second inequality is from (G.13) and our assumption $(\sqrt{\hat{s}} + C/\sqrt{1-\kappa} \cdot \sqrt{s^*}) \cdot \epsilon \leq (1 - \sqrt{\rho})^2 \cdot R$ in (3.9). Therefore, by (G.18) we have $\beta^{(t)} \in \mathcal{B}$. Then (G.8) in Lemma G.2 implies

$$\|\bar{\beta}^{(t+0.5)} - \beta^*\|_2 \leq \rho \cdot \|\beta^{(t)} - \beta^*\|_2 \leq \rho \cdot R < R = \kappa \cdot \|\beta^*\|_2,$$

where the third inequality is from $\rho \in (0, 1)$. Following the same proof for (G.17), we obtain

$$\begin{aligned} \|\beta^{(t+1)} - \beta^*\|_2 &\leq (\sqrt{\hat{s}} + C/\sqrt{1-\kappa} \cdot \sqrt{s^*}) \cdot \epsilon + \sqrt{\rho} \cdot \|\beta^{(t)} - \beta^*\|_2 \\ &\leq \left(1 + \frac{\sqrt{\rho}}{1 - \sqrt{\rho}}\right) \cdot (\sqrt{\hat{s}} + C/\sqrt{1-\kappa} \cdot \sqrt{s^*}) \cdot \epsilon + \sqrt{\rho} \cdot \rho^{t/2} \cdot \|\beta^{(0)} - \beta^*\|_2 \\ &= \frac{(\sqrt{\hat{s}} + C/\sqrt{1-\kappa} \cdot \sqrt{s^*}) \cdot \epsilon}{1 - \sqrt{\rho}} + \rho^{(t+1)/2} \cdot \|\beta^{(0)} - \beta^*\|_2. \end{aligned}$$

Here the second inequality is obtained by plugging in (G.11) for t . Hence we have that (G.11) holds for $t + 1$. By induction, we conclude that (G.11) holds conditioning on the event \mathcal{E} defined in (G.10), which occurs with probability at least $1 - \delta$. By plugging the specific definitions of ρ into (G.11), and applying $\|\beta^{(0)} - \beta^*\|_2 \leq R$ in (G.13) to the right-hand side of (G.11), we obtain (3.10). The results of the gradient descent implementation follows from the same proof with $\rho = \rho_2$. \square

G.2 Proof of Results for Inference

Proof of Theorem 4.6: We establish the asymptotic normality of the decorrelated score statistic defined in (2.7) in two steps. We first prove the asymptotic normality of $\sqrt{n} \cdot S_n(\hat{\beta}_0, \lambda)$, where $\hat{\beta}_0$ is defined in (2.7) and $S_n(\cdot, \cdot)$ is defined in (2.5). Then we prove that $-[T_n(\hat{\beta}_0)]_{\alpha|\gamma}$ defined in (2.7) is a consistent estimator of $\sqrt{n} \cdot S_n(\hat{\beta}_0, \lambda)$'s asymptotic variance. The next lemma accomplishes the first step. Recall $I(\beta^*) = -\mathbb{E}_{\beta^*}[\nabla^2 \ell_n(\beta^*)]/n$ is the Fisher information for $\ell_n(\beta^*)$ defined in (B.3).

Lemma G.3. Under the assumptions of Theorem 4.6, we have that for $n \rightarrow \infty$,

$$\sqrt{n} \cdot S_n(\hat{\beta}_0, \lambda) \xrightarrow{D} N(0, [I(\beta^*)]_{\alpha|\gamma}),$$

where $[I(\beta^*)]_{\alpha|\gamma}$ is defined in (4.2).

Proof. Our proof consists of two steps. Note that by the definition in (2.5) we have

$$\sqrt{n} \cdot S_n(\hat{\beta}_0, \lambda) = \sqrt{n} \cdot [\nabla_1 Q_n(\hat{\beta}_0; \hat{\beta}_0)]_\alpha - \sqrt{n} \cdot w(\hat{\beta}_0, \lambda)^\top \cdot [\nabla_1 Q_n(\hat{\beta}_0; \hat{\beta}_0)]_\gamma. \quad (\text{G.19})$$

Recall that $\mathbf{w}^* = [I(\beta^*)]_{\gamma, \gamma}^{-1} \cdot [I(\beta^*)]_{\gamma, \alpha}$ is defined in (4.1). At the first step, we prove

$$\sqrt{n} \cdot S_n(\hat{\beta}_0, \lambda) = \sqrt{n} \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_\alpha - \sqrt{n} \cdot (\mathbf{w}^*)^\top \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_\gamma + o_{\mathbb{P}}(1). \quad (\text{G.20})$$

In other words, replacing $\widehat{\beta}_0$ and $w(\widehat{\beta}_0, \lambda)$ in (G.19) with the corresponding population quantities β^* and \mathbf{w}^* only introduces an $o_{\mathbb{P}}(1)$ error term. Meanwhile, by Theorem D.1 we have $\nabla_1 Q_n(\beta^*; \beta^*) = \nabla \ell_n(\beta^*)/n$. Recall that $\ell_n(\cdot)$ is the log-likelihood defined in (B.3), which implies that in (G.20)

$$\sqrt{n} \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_{\alpha} - \sqrt{n} \cdot (\mathbf{w}^*)^{\top} \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_{\gamma} = \sqrt{n} \cdot [1, -(\mathbf{w}^*)^{\top}] \cdot \nabla \ell_n(\beta^*)/n$$

is a (rescaled) average of n i.i.d. random variables. At the second step, we calculate the mean and variance of each term within this average and invoke the central limit theorem. Finally we combine these two steps by invoking Slutsky's theorem. See §I.3 for a detailed proof. \square

The next lemma establishes the consistency of $-[T_n(\widehat{\beta}_0)]_{\alpha|\gamma}$ for estimating $[I(\beta^*)]_{\alpha|\gamma}$. Recall that $[T_n(\widehat{\beta}_0)]_{\alpha|\gamma} \in \mathbb{R}$ and $[I(\beta^*)]_{\alpha|\gamma} \in \mathbb{R}$ are defined in (2.7) and (4.2) respectively.

Lemma G.4. Under the assumptions of Theorem 4.6, we have

$$[T_n(\widehat{\beta}_0)]_{\alpha|\gamma} + [I(\beta^*)]_{\alpha|\gamma} = o_{\mathbb{P}}(1). \quad (\text{G.21})$$

Proof. For notational simplicity, we abbreviate $w(\widehat{\beta}_0, \lambda)$ in the definition of $[T_n(\widehat{\beta}_0)]_{\alpha|\gamma}$ as $\widehat{\mathbf{w}}_0$. By (2.7) and (4.3), we have

$$[T_n(\widehat{\beta}_0)]_{\alpha|\gamma} = (1, -\widehat{\mathbf{w}}_0^{\top}) \cdot T_n(\widehat{\beta}_0) \cdot (1, -\widehat{\mathbf{w}}_0^{\top})^{\top}, \quad [I(\beta^*)]_{\alpha|\gamma} = [1, -(\mathbf{w}^*)^{\top}] \cdot I(\beta^*) \cdot [1, -(\mathbf{w}^*)^{\top}]^{\top}.$$

First, we establish the relationship between $\widehat{\mathbf{w}}_0$ and \mathbf{w}^* by analyzing the Dantzig selector in (2.6). Meanwhile, by Theorem D.1 we have $\mathbb{E}_{\beta^*} [T_n(\beta^*)] = -I(\beta^*)$. Then by triangle inequality we have

$$\left| [T_n(\widehat{\beta}_0)]_{\alpha|\gamma} + [I(\beta^*)]_{\alpha|\gamma} \right| \leq \underbrace{\left| [T_n(\widehat{\beta}_0)]_{\alpha|\gamma} - [T_n(\beta^*)]_{\alpha|\gamma} \right|}_{(i)} + \underbrace{\left| [T_n(\beta^*)]_{\alpha|\gamma} - \mathbb{E}_{\beta^*} [T_n(\beta^*)]_{\alpha|\gamma} \right|}_{(ii)}.$$

We prove term (i) is $o_{\mathbb{P}}(1)$ by quantifying the Lipschitz continuity of $T_n(\cdot)$ using Condition 4.4. We then prove term (ii) is $o_{\mathbb{P}}(1)$ by concentration analysis. Together with the result on the relationship between $\widehat{\mathbf{w}}_0$ and \mathbf{w}^* we establish (G.21). See §I.4 for a detailed proof. \square

Combining Lemmas G.3 and G.4 using Slutsky's theorem, we obtain Theorem 4.6.

H Proof of Results for Computation and Estimation

We provide the detailed proof of the main results in §3 for computation and parameter estimation. We first lay out the proof for the general framework, and then the proof for specific models.

H.1 Proof of Lemma G.1

Proof. Recall $\bar{\beta}^{(t+0.5)}$ and $\widehat{\beta}^{(t+1)}$ are defined in (G.1). Note that in (G.4) of Lemma G.1 we assume

$$\|\bar{\beta}^{(t+0.5)} - \beta^*\|_2 \leq \kappa \cdot \|\beta^*\|_2, \quad (\text{H.1})$$

which implies

$$(1 - \kappa) \cdot \|\beta^*\|_2 \leq \|\bar{\beta}^{(t+0.5)}\|_2 \leq (1 + \kappa) \cdot \|\beta^*\|_2. \quad (\text{H.2})$$

For notational simplicity, we define

$$\bar{\theta} = \bar{\beta}^{(t+0.5)} / \|\bar{\beta}^{(t+0.5)}\|_2, \quad \theta = \beta^{(t+0.5)} / \|\beta^{(t+0.5)}\|_2, \quad \text{and} \quad \theta^* = \beta^* / \|\beta^*\|_2. \quad (\text{H.3})$$

Note that $\bar{\theta}$ and θ^* are unit vectors, while θ is not, since it is obtained by normalizing $\beta^{(t+0.5)}$ with $\|\bar{\beta}^{(t+0.5)}\|_2$. Recall that the $\text{supp}(\cdot, \cdot)$ function is defined in (2.2). Hence we have

$$\text{supp}(\theta^*) = \text{supp}(\beta^*) = \mathcal{S}^*, \quad \text{and} \quad \text{supp}(\theta, \widehat{s}) = \text{supp}(\beta^{(t+0.5)}, \widehat{s}) = \widehat{\mathcal{S}}^{(t+0.5)}, \quad (\text{H.4})$$

where the last equality follows from line 6 of Algorithm 4. To ease the notation, we define

$$\mathcal{I}_1 = \mathcal{S}^* \setminus \widehat{\mathcal{S}}^{(t+0.5)}, \quad \mathcal{I}_2 = \mathcal{S}^* \cap \widehat{\mathcal{S}}^{(t+0.5)}, \quad \text{and} \quad \mathcal{I}_3 = \widehat{\mathcal{S}}^{(t+0.5)} \setminus \mathcal{S}^*. \quad (\text{H.5})$$

Let $s_1 = |\mathcal{I}_1|$, $s_2 = |\mathcal{I}_2|$ and $s_3 = |\mathcal{I}_3|$ correspondingly. Also, we define $\Delta = \langle \bar{\theta}, \theta^* \rangle$. Note that

$$\Delta = \langle \bar{\theta}, \theta^* \rangle = \sum_{j \in \mathcal{S}^*} \bar{\theta}_j \cdot \theta_j^* = \sum_{j \in \mathcal{I}_1} \bar{\theta}_j \cdot \theta_j^* + \sum_{j \in \mathcal{I}_2} \bar{\theta}_j \cdot \theta_j^* \leq \|\bar{\theta}_{\mathcal{I}_1}\|_2 \cdot \|\theta_{\mathcal{I}_1}^*\|_2 + \|\bar{\theta}_{\mathcal{I}_2}\|_2 \cdot \|\theta_{\mathcal{I}_2}^*\|_2. \quad (\text{H.6})$$

Here the first equality is from $\text{supp}(\boldsymbol{\theta}^*) = \mathcal{S}^*$, the second equality is from (H.5) and the last inequality is from Cauchy-Schwarz inequality. Furthermore, from (H.6) we have

$$\begin{aligned}\Delta^2 &\leq \left(\|\bar{\boldsymbol{\theta}}_{\mathcal{I}_1}\|_2 \cdot \|\boldsymbol{\theta}_{\mathcal{I}_1}^*\|_2 + \|\bar{\boldsymbol{\theta}}_{\mathcal{I}_2}\|_2 \cdot \|\boldsymbol{\theta}_{\mathcal{I}_2}^*\|_2 \right)^2 \\ &\leq \|\bar{\boldsymbol{\theta}}_{\mathcal{I}_1}\|_2^2 \cdot \left(\|\boldsymbol{\theta}_{\mathcal{I}_1}^*\|_2^2 + \|\boldsymbol{\theta}_{\mathcal{I}_2}^*\|_2^2 \right) + \|\bar{\boldsymbol{\theta}}_{\mathcal{I}_2}\|_2^2 \cdot \left(\|\boldsymbol{\theta}_{\mathcal{I}_1}^*\|_2^2 + \|\boldsymbol{\theta}_{\mathcal{I}_2}^*\|_2^2 \right) \\ &= \|\bar{\boldsymbol{\theta}}_{\mathcal{I}_1}\|_2^2 + \|\bar{\boldsymbol{\theta}}_{\mathcal{I}_2}\|_2^2 \\ &\leq 1 - \|\bar{\boldsymbol{\theta}}_{\mathcal{I}_3}\|_2^2.\end{aligned}\tag{H.7}$$

To obtain the second inequality, we expand the square and apply $2ab \leq a^2 + b^2$. In the equality and the last inequality of (H.7), we use the fact that $\boldsymbol{\theta}^*$ and $\bar{\boldsymbol{\theta}}$ are both unit vectors.

By (2.2) and (H.4), $\widehat{\mathcal{S}}^{(t+0.5)}$ contains the index j 's with the top \widehat{s} largest $|\beta_j^{(t+0.5)}|$'s. Therefore, we have

$$\frac{\|\boldsymbol{\beta}_{\mathcal{I}_3}^{(t+0.5)}\|_2^2}{s_3} = \frac{\sum_{j \in \mathcal{I}_3} (\beta_j^{(t+0.5)})^2}{s_3} \geq \frac{\sum_{j \in \mathcal{I}_1} (\beta_j^{(t+0.5)})^2}{s_1} = \frac{\|\boldsymbol{\beta}_{\mathcal{I}_1}^{(t+0.5)}\|_2^2}{s_1},\tag{H.8}$$

because from (H.5) we have $\mathcal{I}_3 \subseteq \widehat{\mathcal{S}}^{(t+0.5)}$ and $\mathcal{I}_1 \cap \widehat{\mathcal{S}}^{(t+0.5)} = \emptyset$. Taking square roots of both sides of (H.8) and then dividing them by $\|\bar{\boldsymbol{\beta}}^{(t+0.5)}\|_2$ (which is nonzero according to (H.2)), by the definition of $\boldsymbol{\theta}$ in (H.3) we obtain

$$\frac{\|\boldsymbol{\theta}_{\mathcal{I}_3}\|_2}{\sqrt{s_3}} \geq \frac{\|\boldsymbol{\theta}_{\mathcal{I}_1}\|_2}{\sqrt{s_1}}.\tag{H.9}$$

Equipped with (H.9), we now quantify the relationship between $\|\bar{\boldsymbol{\theta}}_{\mathcal{I}_3}\|_2$ and $\|\bar{\boldsymbol{\theta}}_{\mathcal{I}_1}\|_2$. For notational simplicity, let

$$\tilde{\epsilon} = 2 \cdot \|\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_\infty = 2 \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)} - \boldsymbol{\beta}^{(t+0.5)}\|_\infty / \|\bar{\boldsymbol{\beta}}^{(t+0.5)}\|_2.\tag{H.10}$$

Note that we have

$$\max \left\{ \frac{\|\boldsymbol{\theta}_{\mathcal{I}_3} - \bar{\boldsymbol{\theta}}_{\mathcal{I}_3}\|_2}{\sqrt{s_3}}, \frac{\|\boldsymbol{\theta}_{\mathcal{I}_1} - \bar{\boldsymbol{\theta}}_{\mathcal{I}_1}\|_2}{\sqrt{s_1}} \right\} \leq \max \left\{ \|\boldsymbol{\theta}_{\mathcal{I}_3} - \bar{\boldsymbol{\theta}}_{\mathcal{I}_3}\|_\infty, \|\boldsymbol{\theta}_{\mathcal{I}_1} - \bar{\boldsymbol{\theta}}_{\mathcal{I}_1}\|_\infty \right\} \leq \|\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_\infty = \tilde{\epsilon}/2,$$

which implies

$$\begin{aligned}\frac{\|\bar{\boldsymbol{\theta}}_{\mathcal{I}_3}\|_2}{\sqrt{s_3}} &\geq \frac{\|\boldsymbol{\theta}_{\mathcal{I}_3}\|_2}{\sqrt{s_3}} - \frac{\|\boldsymbol{\theta}_{\mathcal{I}_3} - \bar{\boldsymbol{\theta}}_{\mathcal{I}_3}\|_2}{\sqrt{s_3}} \geq \frac{\|\boldsymbol{\theta}_{\mathcal{I}_1}\|_2}{\sqrt{s_1}} - \frac{\|\boldsymbol{\theta}_{\mathcal{I}_3} - \bar{\boldsymbol{\theta}}_{\mathcal{I}_3}\|_2}{\sqrt{s_3}} \\ &\geq \frac{\|\bar{\boldsymbol{\theta}}_{\mathcal{I}_1}\|_2}{\sqrt{s_1}} - \frac{\|\bar{\boldsymbol{\theta}}_{\mathcal{I}_1} - \boldsymbol{\theta}_{\mathcal{I}_1}\|_2}{\sqrt{s_1}} - \frac{\|\boldsymbol{\theta}_{\mathcal{I}_3} - \bar{\boldsymbol{\theta}}_{\mathcal{I}_3}\|_2}{\sqrt{s_3}} \geq \frac{\|\bar{\boldsymbol{\theta}}_{\mathcal{I}_1}\|_2}{\sqrt{s_1}} - \tilde{\epsilon},\end{aligned}\tag{H.11}$$

where second inequality is obtained from (H.9), while the first and third are from triangle inequality. Plugging (H.11) into (H.7), we obtain

$$\Delta^2 \leq 1 - \|\bar{\boldsymbol{\theta}}_{\mathcal{I}_3}\|_2^2 \leq 1 - \left(\sqrt{s_3/s_1} \cdot \|\bar{\boldsymbol{\theta}}_{\mathcal{I}_1}\|_2 - \sqrt{s_3} \cdot \tilde{\epsilon} \right)^2.$$

Since by definition we have $\Delta = \langle \bar{\boldsymbol{\theta}}, \boldsymbol{\theta}^* \rangle \in [-1, 1]$, solving for $\|\bar{\boldsymbol{\theta}}_{\mathcal{I}_1}\|_2$ in the above inequality yields

$$\|\bar{\boldsymbol{\theta}}_{\mathcal{I}_1}\|_2 \leq \sqrt{s_1/s_3} \cdot \sqrt{1 - \Delta^2} + \sqrt{s_1} \cdot \tilde{\epsilon} \leq \sqrt{s^*/\widehat{s}} \cdot \sqrt{1 - \Delta^2} + \sqrt{s^*} \cdot \tilde{\epsilon}.\tag{H.12}$$

Here we employ the fact that $s_1 \leq s^*$ and $s_1/s_3 \leq (s_1 + s_2)/(s_3 + s_2) = s^*/\widehat{s}$, which follows from (H.5) and our assumption in (G.5) that $s^*/\widehat{s} \leq (1 - \kappa)^2/[4 \cdot (1 + \kappa)^2] < 1$.

In the following, we prove that the right-hand side of (H.12) is upper bounded by Δ , i.e.,

$$\sqrt{s^*/\widehat{s}} \cdot \sqrt{1 - \Delta^2} + \sqrt{s^*} \cdot \tilde{\epsilon} \leq \Delta.\tag{H.13}$$

We can verify that a sufficient condition for (H.13) to hold is that

$$\begin{aligned}\Delta &\geq \frac{\sqrt{s^*} \cdot \tilde{\epsilon} + [s^* \cdot \tilde{\epsilon}^2 - (s^*/\widehat{s} + 1) \cdot (s^* \cdot \tilde{\epsilon}^2 - s^*/\widehat{s})]^{1/2}}{s^*/\widehat{s} + 1} \\ &= \frac{\sqrt{s^*} \cdot \tilde{\epsilon} + [-(s^* \cdot \tilde{\epsilon})^2/\widehat{s} + (s^*/\widehat{s} + 1) \cdot (s^*/\widehat{s})]^{1/2}}{s^*/\widehat{s} + 1},\end{aligned}\tag{H.14}$$

which is obtained by solving for Δ in (H.13). When we are solving for Δ in (H.13), we use the fact that $\sqrt{s^*} \cdot \tilde{\epsilon} \leq \Delta$, which holds because

$$\sqrt{s^*} \cdot \tilde{\epsilon} \leq \sqrt{\hat{s}} \cdot \tilde{\epsilon} = 2 \cdot \frac{\sqrt{\hat{s}} \cdot \|\bar{\beta}^{(t+0.5)} - \beta^{(t+0.5)}\|_\infty}{\|\bar{\beta}^{(t+0.5)}\|_2} \leq \frac{1 - \kappa}{1 + \kappa} \leq \Delta. \quad (\text{H.15})$$

The first inequality is from our assumption in (G.5) that $s^*/\hat{s} \leq (1 - \kappa)^2/[4 \cdot (1 + \kappa)^2] < 1$. The equality is from the definition of $\tilde{\epsilon}$ in (H.10). The second inequality follows from our assumption in (G.5) that

$$\sqrt{\hat{s}} \cdot \|\beta^{(t+0.5)} - \bar{\beta}^{(t+0.5)}\|_\infty \leq \frac{(1 - \kappa)^2}{2 \cdot (1 + \kappa)} \cdot \|\beta^*\|_2$$

and the first inequality in (H.2). To prove the last inequality in (H.15), we note that (H.1) implies

$$\|\bar{\beta}^{(t+0.5)}\|_2^2 + \|\beta^*\|_2^2 - 2 \cdot \langle \bar{\beta}^{(t+0.5)}, \beta^* \rangle = \|\bar{\beta}^{(t+0.5)} - \beta^*\|_2^2 \leq \kappa^2 \cdot \|\beta^*\|_2^2.$$

This together with (H.3) implies

$$\begin{aligned} \Delta = \langle \bar{\theta}, \theta^* \rangle &= \frac{\langle \bar{\beta}^{(t+0.5)}, \beta^* \rangle}{\|\bar{\beta}^{(t+0.5)}\|_2 \cdot \|\beta^*\|_2} \geq \frac{\|\bar{\beta}^{(t+0.5)}\|_2^2 + \|\beta^*\|_2^2 - \kappa^2 \cdot \|\beta^*\|_2^2}{2 \cdot \|\bar{\beta}^{(t+0.5)}\|_2 \cdot \|\beta^*\|_2} \\ &\geq \frac{(1 - \kappa)^2 + 1 - \kappa^2}{2 \cdot (1 + \kappa)} = \frac{1 - \kappa}{1 + \kappa}, \end{aligned} \quad (\text{H.16})$$

where in the second inequality we use both sides of (H.2). In summary, we have that (H.15) holds. Now we verify that (H.14) holds. By (H.15) we have

$$\sqrt{\hat{s}} \cdot \tilde{\epsilon} \leq \frac{1 - \kappa}{1 + \kappa} < 1 < \sqrt{(s^* + \hat{s})/\hat{s}},$$

which implies $\tilde{\epsilon} \leq \sqrt{s^* + \hat{s}}/\hat{s}$. For the right-hand side of (H.14) we have

$$\begin{aligned} \frac{\sqrt{s^*} \cdot \tilde{\epsilon} + [-(s^* \cdot \tilde{\epsilon})^2/\hat{s} + (s^*/\hat{s} + 1) \cdot (s^*/\hat{s})]^{1/2}}{s^*/\hat{s} + 1} &\leq \frac{\sqrt{s^*} \cdot \tilde{\epsilon} + [(s^*/\hat{s} + 1) \cdot (s^*/\hat{s})]^{1/2}}{s^*/\hat{s} + 1} \\ &\leq 2 \cdot \sqrt{s^*/(s^* + \hat{s})}, \end{aligned} \quad (\text{H.17})$$

where the last inequality is obtained by plugging in $\tilde{\epsilon} \leq \sqrt{s^* + \hat{s}}/\hat{s}$. Meanwhile, note that we have

$$2 \cdot \sqrt{s^*/(s^* + \hat{s})} \leq 2 \cdot \sqrt{1/[1 + 4 \cdot (1 + \kappa)^2/(1 - \kappa)^2]} \leq (1 - \kappa)/(1 + \kappa) \leq \Delta, \quad (\text{H.18})$$

where the first inequality is from our assumption in (G.5) that $s^*/\hat{s} \leq (1 - \kappa)^2/[4 \cdot (1 + \kappa)^2]$, while the last inequality is from (H.16). Combining (H.17) and (H.18), we then obtain (H.14). By (H.14) we further establish (H.13), i.e., the right-hand side of (H.12) is upper bounded by Δ , which implies

$$\|\bar{\theta}_{\mathcal{I}_1}\|_2 \leq \Delta. \quad (\text{H.19})$$

Furthermore, according to (H.6) we have

$$\Delta \leq \|\bar{\theta}_{\mathcal{I}_1}\|_2 \cdot \|\theta_{\mathcal{I}_1}^*\|_2 + \|\bar{\theta}_{\mathcal{I}_2}\|_2 \cdot \|\theta_{\mathcal{I}_2}^*\|_2 \leq \|\bar{\theta}_{\mathcal{I}_1}\|_2 \cdot \|\theta_{\mathcal{I}_1}^*\|_2 + \sqrt{1 - \|\bar{\theta}_{\mathcal{I}_1}\|_2^2} \cdot \sqrt{1 - \|\theta_{\mathcal{I}_1}^*\|_2^2}, \quad (\text{H.20})$$

where in the last inequality we use the fact θ^* and $\bar{\theta}$ are unit vectors. Now we solve for $\|\theta_{\mathcal{I}_1}^*\|_2$ in (H.20). According to (H.19) and the fact that $\|\theta_{\mathcal{I}_1}^*\|_2 \leq \|\theta^*\|_2 = 1$, on the right-hand side of (H.20) we have $\|\bar{\theta}_{\mathcal{I}_1}\|_2 \cdot \|\theta_{\mathcal{I}_1}^*\|_2 \leq \|\bar{\theta}_{\mathcal{I}_1}\|_2 \leq \Delta$. Thus, we have

$$\left(\Delta - \|\bar{\theta}_{\mathcal{I}_1}\|_2 \cdot \|\theta_{\mathcal{I}_1}^*\|_2\right)^2 \leq \left(1 - \|\bar{\theta}_{\mathcal{I}_1}\|_2^2\right) \cdot \left(1 - \|\theta_{\mathcal{I}_1}^*\|_2^2\right).$$

Further by solving for $\|\theta_{\mathcal{I}_1}^*\|_2$ in the above inequality, we obtain

$$\begin{aligned} \|\theta_{\mathcal{I}_1}^*\|_2 &\leq \|\bar{\theta}_{\mathcal{I}_1}\|_2 \cdot \Delta + \sqrt{1 - \|\bar{\theta}_{\mathcal{I}_1}\|_2^2} \cdot \sqrt{1 - \Delta^2} \leq \|\bar{\theta}_{\mathcal{I}_1}\|_2 + \sqrt{1 - \Delta^2} \\ &\leq (1 + \sqrt{s^*/\hat{s}}) \cdot \sqrt{1 - \Delta^2} + \sqrt{s^*} \cdot \tilde{\epsilon}, \end{aligned} \quad (\text{H.21})$$

where in the second inequality we use the fact that $\Delta \leq 1$, which follows from its definition, while in the last inequality we plug in (H.12). Then combining (H.12) and (H.21), we obtain

$$\|\bar{\boldsymbol{\theta}}_{\mathcal{I}_1}\|_2 \cdot \|\boldsymbol{\theta}_{\mathcal{I}_1}^*\|_2 \leq \left[\sqrt{s^*/\hat{s}} \cdot \sqrt{1-\Delta^2} + \sqrt{s^*} \cdot \tilde{\epsilon} \right] \cdot \left[(1 + \sqrt{s^*/\hat{s}}) \cdot \sqrt{1-\Delta^2} + \sqrt{s^*} \cdot \tilde{\epsilon} \right]. \quad (\text{H.22})$$

Note that by (G.1) and the definition of $\bar{\boldsymbol{\theta}}$ in (H.3), we have

$$\bar{\boldsymbol{\beta}}^{(t+1)} = \text{trunc}(\bar{\boldsymbol{\beta}}^{(t+0.5)}, \hat{\mathcal{S}}^{(t+0.5)}) = \text{trunc}(\bar{\boldsymbol{\theta}}, \hat{\mathcal{S}}^{(t+0.5)}) \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)}\|_2.$$

Therefore, we have

$$\begin{aligned} \left\langle \bar{\boldsymbol{\beta}}^{(t+1)} / \|\bar{\boldsymbol{\beta}}^{(t+0.5)}\|_2, \boldsymbol{\beta}^* / \|\boldsymbol{\beta}^*\|_2 \right\rangle &= \left\langle \text{trunc}(\bar{\boldsymbol{\theta}}, \hat{\mathcal{S}}^{(t+0.5)}), \boldsymbol{\theta}^* \right\rangle = \langle \bar{\boldsymbol{\theta}}_{\mathcal{I}_2}, \boldsymbol{\theta}_{\mathcal{I}_2}^* \rangle = \langle \bar{\boldsymbol{\theta}}, \boldsymbol{\theta}^* \rangle - \langle \bar{\boldsymbol{\theta}}_{\mathcal{I}_1}, \boldsymbol{\theta}_{\mathcal{I}_1}^* \rangle \\ &\geq \langle \bar{\boldsymbol{\theta}}, \boldsymbol{\theta}^* \rangle - \|\bar{\boldsymbol{\theta}}_{\mathcal{I}_1}\|_2 \cdot \|\boldsymbol{\theta}_{\mathcal{I}_1}^*\|_2, \end{aligned}$$

where the second and third equalities follow from (H.5). Let $\bar{\chi} = \|\bar{\boldsymbol{\beta}}^{(t+0.5)}\|_2 \cdot \|\boldsymbol{\beta}^*\|_2$. Plugging (H.22) into the right-hand side of the above inequality and then multiplying $\bar{\chi}$ on both sides, we obtain

$$\begin{aligned} &\langle \bar{\boldsymbol{\beta}}^{(t+1)}, \boldsymbol{\beta}^* \rangle \quad (\text{H.23}) \\ &\geq \langle \bar{\boldsymbol{\beta}}^{(t+0.5)}, \boldsymbol{\beta}^* \rangle \\ &\quad - \left[\sqrt{s^*/\hat{s}} \cdot \sqrt{\bar{\chi} \cdot (1-\Delta^2)} + \sqrt{s^*} \cdot \sqrt{\bar{\chi}} \cdot \tilde{\epsilon} \right] \cdot \left[(1 + \sqrt{s^*/\hat{s}}) \cdot \sqrt{\bar{\chi} \cdot (1-\Delta^2)} + \sqrt{s^*} \cdot \sqrt{\bar{\chi}} \cdot \tilde{\epsilon} \right] \\ &= \langle \bar{\boldsymbol{\beta}}^{(t+0.5)}, \boldsymbol{\beta}^* \rangle - (\sqrt{s^*/\hat{s}} + s^*/\hat{s}) \cdot \bar{\chi} \cdot (1-\Delta^2) \\ &\quad - (1 + 2 \cdot \sqrt{s^*/\hat{s}}) \cdot \underbrace{\sqrt{\bar{\chi} \cdot (1-\Delta^2)}}_{(i)} \cdot \underbrace{\sqrt{s^*} \cdot \sqrt{\bar{\chi}} \cdot \tilde{\epsilon}}_{(ii)} - (\sqrt{s^*} \cdot \sqrt{\bar{\chi}} \cdot \tilde{\epsilon})^2. \end{aligned}$$

For term (i) in (H.23), note that $\sqrt{1-\Delta^2} \leq \sqrt{2} \cdot (1-\Delta)$. By (H.3) and the definition that $\Delta = \langle \bar{\boldsymbol{\theta}}, \boldsymbol{\theta}^* \rangle$, for term (i) we have

$$\begin{aligned} \sqrt{\bar{\chi} \cdot (1-\Delta^2)} &\leq \sqrt{2 \cdot \bar{\chi} \cdot (1-\Delta)} \leq \sqrt{2 \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)}\|_2 \cdot \|\boldsymbol{\beta}^*\|_2 - 2 \cdot \langle \bar{\boldsymbol{\beta}}^{(t+0.5)}, \boldsymbol{\beta}^* \rangle} \quad (\text{H.24}) \\ &\leq \sqrt{\|\bar{\boldsymbol{\beta}}^{(t+0.5)}\|_2^2 + \|\boldsymbol{\beta}^*\|_2^2 - 2 \cdot \langle \bar{\boldsymbol{\beta}}^{(t+0.5)}, \boldsymbol{\beta}^* \rangle} = \|\bar{\boldsymbol{\beta}}^{(t+0.5)} - \boldsymbol{\beta}^*\|_2. \end{aligned}$$

For term (ii) in (H.23), by the definition of $\tilde{\epsilon}$ in (H.10) we have

$$\begin{aligned} \sqrt{\bar{\chi}} \cdot \tilde{\epsilon} &= \sqrt{\|\bar{\boldsymbol{\beta}}^{(t+0.5)}\|_2 \cdot \|\boldsymbol{\beta}^*\|_2} \cdot 2 \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)} - \boldsymbol{\beta}^{(t+0.5)}\|_\infty / \|\bar{\boldsymbol{\beta}}^{(t+0.5)}\|_2 \\ &= 2 \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)} - \boldsymbol{\beta}^{(t+0.5)}\|_\infty \cdot \sqrt{\|\boldsymbol{\beta}^*\|_2 / \|\bar{\boldsymbol{\beta}}^{(t+0.5)}\|_2} \leq \frac{2}{\sqrt{1-\kappa}} \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)} - \boldsymbol{\beta}^{(t+0.5)}\|_\infty, \quad (\text{H.25}) \end{aligned}$$

where the last inequality is obtained from (H.2). Plugging (H.24) and (H.25) into (H.23), we obtain

$$\begin{aligned} \langle \bar{\boldsymbol{\beta}}^{(t+1)}, \boldsymbol{\beta}^* \rangle &\geq \langle \bar{\boldsymbol{\beta}}^{(t+0.5)}, \boldsymbol{\beta}^* \rangle - (\sqrt{s^*/\hat{s}} + s^*/\hat{s}) \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)} - \boldsymbol{\beta}^*\|_2^2 \quad (\text{H.26}) \\ &\quad - (1 + 2 \cdot \sqrt{s^*/\hat{s}}) \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)} - \boldsymbol{\beta}^*\|_2 \cdot \frac{2 \cdot \sqrt{s^*}}{\sqrt{1-\kappa}} \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)} - \boldsymbol{\beta}^{(t+0.5)}\|_\infty \\ &\quad - \frac{4 \cdot s^*}{1-\kappa} \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)} - \boldsymbol{\beta}^{(t+0.5)}\|_\infty^2. \end{aligned}$$

Meanwhile, according to (G.1) we have that $\bar{\boldsymbol{\beta}}^{(t+1)}$ is obtained by truncating $\bar{\boldsymbol{\beta}}^{(t+0.5)}$, which implies

$$\|\bar{\boldsymbol{\beta}}^{(t+1)}\|_2^2 + \|\boldsymbol{\beta}^*\|_2^2 \leq \|\bar{\boldsymbol{\beta}}^{(t+0.5)}\|_2^2 + \|\boldsymbol{\beta}^*\|_2^2. \quad (\text{H.27})$$

Subtracting two times both sides of (H.26) from (H.27), we obtain

$$\begin{aligned} \|\bar{\boldsymbol{\beta}}^{(t+1)} - \boldsymbol{\beta}^*\|_2^2 &\leq (1 + 2 \cdot \sqrt{s^*/\hat{s}} + 2 \cdot s^*/\hat{s}) \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)} - \boldsymbol{\beta}^*\|_2^2 \\ &\quad + (1 + 2 \cdot \sqrt{s^*/\hat{s}}) \cdot \frac{4 \cdot \sqrt{s^*}}{\sqrt{1-\kappa}} \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)} - \boldsymbol{\beta}^{(t+0.5)}\|_\infty \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)} - \boldsymbol{\beta}^*\|_2 \\ &\quad + \frac{8 \cdot s^*}{1-\kappa} \cdot \|\bar{\boldsymbol{\beta}}^{(t+0.5)} - \boldsymbol{\beta}^{(t+0.5)}\|_\infty^2. \end{aligned}$$

We can easily verify that the above inequality implies

$$\begin{aligned} \|\bar{\beta}^{(t+1)} - \beta^*\|_2^2 &\leq (1 + 2 \cdot \sqrt{s^*/\hat{s}} + 2 \cdot s^*/\hat{s}) \cdot \left[\|\bar{\beta}^{(t+0.5)} - \beta^*\|_2 + \frac{2 \cdot \sqrt{s^*}}{\sqrt{1-\kappa}} \cdot \|\bar{\beta}^{(t+0.5)} - \beta^{(t+0.5)}\|_\infty \right]^2 \\ &\quad + \frac{8 \cdot s^*}{1-\kappa} \cdot \|\bar{\beta}^{(t+0.5)} - \beta^{(t+0.5)}\|_\infty^2. \end{aligned}$$

Taking square roots of both sides and utilizing the fact that $\sqrt{a^2 + b^2} \leq a + b$ ($a, b > 0$), we obtain

$$\begin{aligned} \|\bar{\beta}^{(t+1)} - \beta^*\|_2 &\leq (1 + 4 \cdot \sqrt{s^*/\hat{s}})^{1/2} \cdot \|\bar{\beta}^{(t+0.5)} - \beta^*\|_2 \\ &\quad + \frac{C \cdot \sqrt{s^*}}{\sqrt{1-\kappa}} \cdot \|\bar{\beta}^{(t+0.5)} - \beta^{(t+0.5)}\|_\infty, \end{aligned} \quad (\text{H.28})$$

where $C > 0$ is a constant. Here we utilize the fact that $s^*/\hat{s} \leq \sqrt{s^*/\hat{s}}$ and

$$1 + 2 \cdot \sqrt{s^*/\hat{s}} + 2 \cdot s^*/\hat{s} \leq 5,$$

both of which follow from our assumption that $s^*/\hat{s} \leq (1-\kappa)^2/[4 \cdot (1+\kappa)^2] < 1$ in (G.5). By (H.28) we conclude the proof of Lemma G.1. \square

H.2 Proof of Lemma G.2

In the following, we prove (G.8) and (G.9) for the maximization and gradient ascent implementation of the M-step correspondingly.

Proof of (G.8): To prove (G.8) for the maximization implementation of the M-step (Algorithm 2), note that by the self-consistency property [18] we have

$$\beta^* = \operatorname{argmax}_{\beta} Q(\beta; \beta^*). \quad (\text{H.29})$$

Hence, β^* satisfies the following first-order optimality condition

$$\langle \beta - \beta^*, \nabla_1 Q(\beta^*; \beta^*) \rangle \leq 0, \quad \text{for all } \beta,$$

where $\nabla_1 Q(\cdot, \cdot)$ denotes the gradient taken with respect to the first variable. In particular, it implies

$$\langle \bar{\beta}^{(t+0.5)} - \beta^*, \nabla_1 Q(\beta^*; \beta^*) \rangle \leq 0. \quad (\text{H.30})$$

Meanwhile, by (G.1) and the definition of $M(\cdot)$ in (3.1), we have

$$\bar{\beta}^{(t+0.5)} = M(\beta^{(t)}) = \operatorname{argmax}_{\beta} Q(\beta; \beta^{(t)}).$$

Hence we have the following first-order optimality condition

$$\langle \beta - \bar{\beta}^{(t+0.5)}, \nabla_1 Q(\bar{\beta}^{(t+0.5)}; \beta^{(t)}) \rangle \leq 0, \quad \text{for all } \beta,$$

which implies

$$\langle \beta^* - \bar{\beta}^{(t+0.5)}, \nabla_1 Q(\bar{\beta}^{(t+0.5)}; \beta^{(t)}) \rangle \leq 0. \quad (\text{H.31})$$

Combining (H.30) and (H.31), we then obtain

$$\langle \beta^* - \bar{\beta}^{(t+0.5)}, -\nabla_1 Q(\beta^*; \beta^*) \rangle \leq \langle \beta^* - \bar{\beta}^{(t+0.5)}, -\nabla_1 Q(\bar{\beta}^{(t+0.5)}; \beta^{(t)}) \rangle,$$

which implies

$$\begin{aligned} &\langle \beta^* - \bar{\beta}^{(t+0.5)}, \nabla_1 Q(\bar{\beta}^{(t+0.5)}; \beta^*) - \nabla_1 Q(\beta^*; \beta^*) \rangle \\ &\leq \langle \beta^* - \bar{\beta}^{(t+0.5)}, \nabla_1 Q(\bar{\beta}^{(t+0.5)}; \beta^*) - \nabla_1 Q(\bar{\beta}^{(t+0.5)}; \beta^{(t)}) \rangle. \end{aligned} \quad (\text{H.32})$$

In the following, we establish upper and lower bounds for both sides of (H.32) correspondingly. By applying Condition *Lipschitz-Gradient-1*(γ_1, \mathcal{B}), for the right-hand side of (H.32) we have

$$\begin{aligned} &\langle \beta^* - \bar{\beta}^{(t+0.5)}, \nabla_1 Q(\bar{\beta}^{(t+0.5)}; \beta^*) - \nabla_1 Q(\bar{\beta}^{(t+0.5)}; \beta^{(t)}) \rangle \\ &\leq \|\beta^* - \bar{\beta}^{(t+0.5)}\|_2 \cdot \left\| \nabla_1 Q(\bar{\beta}^{(t+0.5)}; \beta^*) - \nabla_1 Q(\bar{\beta}^{(t+0.5)}; \beta^{(t)}) \right\|_2 \\ &\leq \gamma_1 \cdot \|\beta^* - \bar{\beta}^{(t+0.5)}\|_2 \cdot \|\beta^* - \beta^{(t)}\|_2, \end{aligned} \quad (\text{H.33})$$

where the last inequality is from (3.3). Meanwhile, for the left-hand side of (H.32), we have

$$Q(\bar{\beta}^{(t+0.5)}; \beta^*) \leq Q(\beta^*; \beta^*) + \langle \nabla_1 Q(\beta^*; \beta^*), \bar{\beta}^{(t+0.5)} - \beta^* \rangle - \nu/2 \cdot \|\bar{\beta}^{(t+0.5)} - \beta^*\|_2^2, \quad (\text{H.34})$$

$$Q(\beta^*; \beta^*) \leq Q(\bar{\beta}^{(t+0.5)}; \beta^*) + \langle \nabla_1 Q(\bar{\beta}^{(t+0.5)}; \beta^*), \beta^* - \bar{\beta}^{(t+0.5)} \rangle - \nu/2 \cdot \|\bar{\beta}^{(t+0.5)} - \beta^*\|_2^2 \quad (\text{H.35})$$

by (3.6) in Condition *Concavity-Smoothness*(μ, ν, \mathcal{B}). By adding (H.34) and (H.35), we obtain

$$\nu \cdot \|\bar{\beta}^{(t+0.5)} - \beta^*\|_2^2 \leq \langle \beta^* - \bar{\beta}^{(t+0.5)}, \nabla_1 Q(\bar{\beta}^{(t+0.5)}; \beta^*) - \nabla_1 Q(\beta^*; \beta^*) \rangle. \quad (\text{H.36})$$

Plugging (H.33) and (H.36) into (H.32), we obtain

$$\nu \cdot \|\bar{\beta}^{(t+0.5)} - \beta^*\|_2^2 \leq \gamma_1 \cdot \|\beta^* - \bar{\beta}^{(t+0.5)}\|_2 \cdot \|\beta^* - \beta^{(t)}\|_2,$$

which implies (G.8) in Lemma G.2.

Proof of (G.9): We turn to prove (G.9). The self-consistency property in (H.29) implies that β^* is the maximizer of $Q(\cdot; \beta^*)$. Furthermore, (3.5) and (3.6) in Condition *Concavity-Smoothness*(μ, ν, \mathcal{B}) ensure that $-Q(\cdot; \beta^*)$ is μ -smooth and ν -strongly convex. By invoking standard optimization results for minimizing strongly convex and smooth objective functions, e.g., in [21], for stepsize $\eta = 2/(\nu + \mu)$, we have

$$\left\| \beta^{(t)} + \eta \cdot \nabla_1 Q(\beta^{(t)}; \beta^*) - \beta^* \right\|_2 \leq \left(\frac{\mu - \nu}{\mu + \nu} \right) \cdot \|\beta^{(t)} - \beta^*\|_2, \quad (\text{H.37})$$

i.e., the gradient ascent step decreases the distance to β^* by a multiplicative factor. Hence, for the gradient ascent implementation of the M-step, i.e., $M(\cdot)$ defined in (3.2), we have

$$\begin{aligned} \|\bar{\beta}^{(t+0.5)} - \beta^*\|_2 &= \left\| M(\beta^{(t)}) - \beta^* \right\|_2 \\ &= \left\| \beta^{(t)} + \eta \cdot \nabla_1 Q(\beta^{(t)}; \beta^{(t)}) - \beta^* \right\|_2 \\ &\leq \left\| \beta^{(t)} + \eta \cdot \nabla_1 Q(\beta^{(t)}; \beta^*) - \beta^* \right\|_2 + \eta \cdot \left\| \nabla_1 Q(\beta^{(t)}; \beta^*) - \nabla_1 Q(\beta^{(t)}; \beta^{(t)}) \right\|_2 \\ &\leq \left(\frac{\mu - \nu}{\mu + \nu} \right) \cdot \|\beta^{(t)} - \beta^*\|_2 + \eta \cdot \gamma_2 \cdot \|\beta^{(t)} - \beta^*\|_2, \end{aligned} \quad (\text{H.38})$$

where the last inequality is from (H.37) and (3.4) in Condition *Lipschitz-Gradient-2*(γ_2, \mathcal{B}). Plugging $\eta = 2/(\nu + \mu)$ into (H.38), we obtain

$$\|\bar{\beta}^{(t+0.5)} - \beta^*\|_2 \leq \left(\frac{\mu - \nu + 2 \cdot \gamma_2}{\mu + \nu} \right) \cdot \|\beta^{(t)} - \beta^*\|_2,$$

which implies (G.9). Thus, we conclude the proof of Lemma G.2.

H.3 Auxiliary Lemma for Proving Theorem 3.4

The following lemma characterizes the initialization step in line 2 of Algorithm 4.

Lemma H.1. Suppose that we have $\|\beta^{\text{init}} - \beta^*\|_2 \leq \kappa \cdot \|\beta^*\|_2$ for some $\kappa \in (0, 1)$. Assuming that $\hat{s} \geq 4 \cdot (1 + \kappa)^2 / (1 - \kappa)^2 \cdot s^*$, we have $\|\beta^{(0)} - \beta^*\|_2 \leq (1 + 4 \cdot \sqrt{s^* / \hat{s}})^{1/2} \cdot \|\beta^{\text{init}} - \beta^*\|_2$.

Proof. Following the same proof of Lemma G.1 with both $\bar{\beta}^{(t+0.5)}$ and $\beta^{(t+0.5)}$ replaced with β^{init} , $\bar{\beta}^{(t+1)}$ replaced with $\beta^{(0)}$ and $\hat{S}^{(t+0.5)}$ replaced with \hat{S}^{init} , we reach the conclusion. \square

H.4 Proof of Lemma E.2

Proof. Recall that $Q(\cdot; \cdot)$ is the expectation of $Q_n(\cdot; \cdot)$. According to (A.2) and (3.1), we have

$$M(\beta) = \mathbb{E}[2 \cdot \omega_\beta(\mathbf{Y}) \cdot \mathbf{Y} - \mathbf{Y}]$$

with $\omega_\beta(\cdot)$ being the weight function defined in (A.2), which together with (A.3) implies

$$M_n(\beta) - M(\beta) = \frac{1}{n} \sum_{i=1}^n [2 \cdot \omega_\beta(\mathbf{y}_i) - 1] \cdot \mathbf{y}_i - \mathbb{E}\{[2 \cdot \omega_\beta(\mathbf{Y}) - 1] \cdot \mathbf{Y}\}. \quad (\text{H.39})$$

Recall \mathbf{y}_i is the i -th realization of \mathbf{Y} , which follows the mixture distribution. For any $u > 0$, we have

$$\begin{aligned} \mathbb{E} \left\{ \exp \left[u \cdot \|M_n(\boldsymbol{\beta}) - M(\boldsymbol{\beta})\|_\infty \right] \right\} &= \mathbb{E} \left\{ \max_{j \in \{1, \dots, d\}} \exp \left[u \cdot |[M_n(\boldsymbol{\beta}) - M(\boldsymbol{\beta})]_j| \right] \right\} \\ &\leq \sum_{j=1}^d \mathbb{E} \left\{ \exp \left[u \cdot |[M_n(\boldsymbol{\beta}) - M(\boldsymbol{\beta})]_j| \right] \right\}. \end{aligned} \quad (\text{H.40})$$

Based on (H.39), we apply the symmetrization result in Lemma J.4 to the right-hand side of (H.40). Then we have

$$\mathbb{E} \left\{ \exp \left[u \cdot \|M_n(\boldsymbol{\beta}) - M(\boldsymbol{\beta})\|_\infty \right] \right\} \leq \sum_{j=1}^d \mathbb{E} \left\{ \exp \left[u \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot [2 \cdot \omega_{\boldsymbol{\beta}}(\mathbf{y}_i) - 1] \cdot y_{i,j} \right| \right] \right\}, \quad (\text{H.41})$$

where ξ_1, \dots, ξ_n are i.i.d. Rademacher random variables that are independent of $\mathbf{y}_1, \dots, \mathbf{y}_n$. Then we invoke the contraction result in Lemma J.5 by setting

$$f(y_{i,j}) = y_{i,j}, \quad \mathcal{F} = \{f\}, \quad \psi_i(v) = [2 \cdot \omega_{\boldsymbol{\beta}}(\mathbf{y}_i) - 1] \cdot v, \quad \text{and} \quad \phi(v) = \exp(u \cdot v),$$

where u is the variable of the moment generating function in (H.40). From the definition of $\omega_{\boldsymbol{\beta}}(\cdot)$ in (A.2) we have $|2 \cdot \omega_{\boldsymbol{\beta}}(\mathbf{y}_i) - 1| \leq 1$, which implies

$$|\psi_i(v) - \psi_i(v')| \leq |[2 \cdot \omega_{\boldsymbol{\beta}}(\mathbf{y}_i) - 1] \cdot (v - v')| \leq |v - v'|, \quad \text{for all } v, v' \in \mathbb{R}.$$

Therefore, by Lemma J.5 we obtain

$$\mathbb{E} \left\{ \exp \left[u \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot [2 \cdot \omega_{\boldsymbol{\beta}}(\mathbf{y}_i) - 1] \cdot y_{i,j} \right| \right] \right\} \leq \mathbb{E} \left\{ \exp \left[2 \cdot u \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_{i,j} \right| \right] \right\} \quad (\text{H.42})$$

for the right-hand side of (H.41), where $j \in \{1, \dots, d\}$. Here note that in Gaussian mixture model we have $y_{i,j} = z_i \cdot \beta_j^* + v_{i,j}$, where z_i is a Rademacher random variable and $v_{i,j} \sim N(0, \sigma^2)$. Therefore, according to Example 5.8 in [28] we have $\|z_i \cdot \beta_j^*\|_{\psi_2} \leq |\beta_j^*|$ and $\|v_{i,j}\|_{\psi_2} \leq C \cdot \sigma$. Hence by Lemma J.1 we have

$$\|y_{i,j}\|_{\psi_2} = \|z_i \cdot \beta_j^* + v_{i,j}\|_{\psi_2} \leq C \cdot \sqrt{|\beta_j^*|^2 + C' \cdot \sigma^2} \leq C \cdot \sqrt{\|\boldsymbol{\beta}^*\|_\infty^2 + C' \cdot \sigma^2}.$$

Since $|\xi_i \cdot y_{i,j}| = |y_{i,j}|$, $\xi_i \cdot y_{i,j}$ and $y_{i,j}$ have the same ψ_2 -norm. Because ξ_i is a Rademacher random variable independent of $y_{i,j}$, we have $\mathbb{E}(\xi_i \cdot y_{i,j}) = 0$. By Lemma 5.5 in [28], we obtain

$$\mathbb{E}[\exp(u' \cdot \xi_i \cdot y_{i,j})] \leq \exp\left[(u')^2 \cdot C \cdot (\|\boldsymbol{\beta}^*\|_\infty^2 + C' \cdot \sigma^2)\right], \quad \text{for all } u' \in \mathbb{R}. \quad (\text{H.43})$$

Hence, for the right-hand side of (H.42) we have

$$\begin{aligned} \mathbb{E} \left\{ \exp \left[2 \cdot u \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_{i,j} \right| \right] \right\} &\leq \mathbb{E} \left(\max \left\{ \exp \left[2 \cdot u \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_{i,j} \right], \exp \left[-2 \cdot u \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_{i,j} \right] \right\} \right) \\ &\leq \mathbb{E} \left\{ \exp \left[2 \cdot u \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_{i,j} \right] \right\} + \mathbb{E} \left\{ \exp \left[-2 \cdot u \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_{i,j} \right] \right\} \\ &\leq 2 \cdot \exp \left[C \cdot u^2 \cdot (\|\boldsymbol{\beta}^*\|_\infty^2 + C' \cdot \sigma^2) / n \right]. \end{aligned} \quad (\text{H.44})$$

Here the last inequality is obtained by plugging (H.43) with $u' = 2 \cdot u/n$ and $u' = -2 \cdot u/n$ respectively into the two terms. Plugging (H.44) into (H.42) and then into (H.41), we obtain

$$\mathbb{E} \left\{ \exp \left[u \cdot \|M_n(\boldsymbol{\beta}) - M(\boldsymbol{\beta})\|_\infty \right] \right\} \leq 2 \cdot d \cdot \exp \left[C \cdot u^2 \cdot (\|\boldsymbol{\beta}^*\|_\infty^2 + C' \cdot \sigma^2) / n \right].$$

By Chernoff bound we have that, for all $u > 0$ and $v > 0$,

$$\begin{aligned} \mathbb{P} \left[\|M_n(\boldsymbol{\beta}) - M(\boldsymbol{\beta})\|_\infty > v \right] &\leq \mathbb{E} \left\{ \exp \left[u \cdot \|M_n(\boldsymbol{\beta}) - M(\boldsymbol{\beta})\|_\infty \right] \right\} / \exp(u \cdot v) \\ &\leq 2 \cdot \exp \left[C \cdot u^2 \cdot (\|\boldsymbol{\beta}^*\|_\infty^2 + C' \cdot \sigma^2) / n - u \cdot v + \log d \right]. \end{aligned}$$

Minimizing the right-hand side over u we obtain

$$\mathbb{P} \left[\|M_n(\boldsymbol{\beta}) - M(\boldsymbol{\beta})\|_\infty > v \right] \leq 2 \cdot \exp \left\{ -n \cdot v^2 / \left[4 \cdot C \cdot (\|\boldsymbol{\beta}^*\|_\infty^2 + C' \cdot \sigma^2) \right] + \log d \right\}.$$

Setting the right-hand side to be δ , we have that

$$v = C \cdot (\|\beta^*\|_\infty^2 + C' \cdot \sigma^2)^{1/2} \cdot \sqrt{\frac{\log d + \log(2/\delta)}{n}} \leq C'' \cdot (\|\beta^*\|_\infty + \sigma) \cdot \sqrt{\frac{\log d + \log(2/\delta)}{n}}$$

holds for some constants C , C' and C'' , which completes the proof of Lemma E.2. \square

H.5 Proof of Lemma E.5

In the sequel, we first establish the result for the maximization implementation of the M-step and then for the gradient ascent implementation.

Maximization Implementation: For the maximization implementation we need to estimate the inverse covariance matrix $\Theta^* = \Sigma^{-1}$ with the CLIME estimator $\hat{\Theta}$ defined in (A.7). The following lemma from [7] quantifies the statistical rate of convergence of $\hat{\Theta}$. Recall that $\|\cdot\|_{1,\infty}$ is defined as the maximum of the row ℓ_1 -norms of a matrix.

Lemma H.2. For $\Sigma = \mathbf{I}_d$ and $\lambda^{\text{CLIME}} = C \cdot \sqrt{\log d/n}$ in (A.7), we have that

$$\|\hat{\Theta} - \Theta^*\|_{1,\infty} \leq C' \cdot \sqrt{\frac{\log d + \log(1/\delta)}{n}}$$

holds with probability at least $1 - \delta$, where C and C' are positive constants.

Proof. See the proof of Theorem 6 in [7] for details. \square

Now we are ready to prove (E.8) of Lemma E.5.

Proof. Recall that $Q(\cdot; \cdot)$ is the expectation of $Q_n(\cdot; \cdot)$. According to (A.5) and (3.1), we have

$$M(\beta) = \mathbb{E} \left\{ [2 \cdot \omega_\beta(\mathbf{X}, Y) - 1] \cdot Y \cdot \mathbf{X} \right\} \quad (\text{H.45})$$

with $\omega_\beta(\cdot, \cdot)$ being the weight function defined in (A.5), which together with (A.8) implies

$$M_n(\beta) - M(\beta) = \hat{\Theta} \cdot \frac{1}{n} \sum_{i=1}^n [2 \cdot \omega_\beta(\mathbf{x}_i, y_i) - 1] \cdot y_i \cdot \mathbf{x}_i - \mathbb{E} \left\{ [2 \cdot \omega_\beta(\mathbf{X}, Y) - 1] \cdot Y \cdot \mathbf{X} \right\}.$$

Here $\hat{\Theta}$ is the CLIME estimator defined in (A.7). For notational simplicity, we denote

$$\bar{\omega}_i = 2 \cdot \omega_\beta(\mathbf{x}_i, y_i) - 1, \quad \text{and} \quad \bar{\omega} = 2 \cdot \omega_\beta(\mathbf{X}, Y) - 1. \quad (\text{H.46})$$

It is worth noting that both $\bar{\omega}_i$ and $\bar{\omega}$ depend on β . Note that we have

$$\begin{aligned} & \|M_n(\beta) - M(\beta)\|_\infty \quad (\text{H.47}) \\ & \leq \left\| \hat{\Theta} \cdot \left[\frac{1}{n} \sum_{i=1}^n \bar{\omega}_i \cdot y_i \cdot \mathbf{x}_i - \mathbb{E}(\bar{\omega} \cdot Y \cdot \mathbf{X}) \right] \right\|_\infty + \left\| (\hat{\Theta} - \mathbf{I}_d) \cdot \mathbb{E}(\bar{\omega} \cdot Y \cdot \mathbf{X}) \right\|_\infty \\ & \leq \underbrace{\|\hat{\Theta}\|_{1,\infty}}_{(i)} \cdot \underbrace{\left\| \frac{1}{n} \sum_{i=1}^n \bar{\omega}_i \cdot y_i \cdot \mathbf{x}_i - \mathbb{E}(\bar{\omega} \cdot Y \cdot \mathbf{X}) \right\|_\infty}_{(ii)} + \underbrace{\|\hat{\Theta} - \mathbf{I}_d\|_{1,\infty}}_{(iii)} \cdot \underbrace{\|\mathbb{E}(\bar{\omega} \cdot Y \cdot \mathbf{X})\|_\infty}_{(iv)}. \end{aligned}$$

Analysis of Term (i): For term (i) in (H.47), recall that by our model assumption we have $\Sigma = \mathbf{I}_d$, which implies $\Theta^* = \Sigma^{-1} = \mathbf{I}_d$. By Lemma H.2, for a sufficiently large sample size n , we have that

$$\|\hat{\Theta}\|_{1,\infty} \leq \|\hat{\Theta} - \mathbf{I}_d\|_{1,\infty} + \|\mathbf{I}_d\|_{1,\infty} \leq 1/2 + 1 = 3/2 \quad (\text{H.48})$$

holds with probability at least $1 - \delta/4$.

Analysis of Term (ii): For term (ii) in (H.47), we have that for $u > 0$,

$$\begin{aligned}
& \mathbb{E} \left\{ \exp \left[u \cdot \left\| \frac{1}{n} \sum_{i=1}^n \bar{\omega}_i \cdot y_i \cdot \mathbf{x}_i - \mathbb{E}(\bar{\omega} \cdot Y \cdot \mathbf{X}) \right\|_{\infty} \right] \right\} \\
&= \mathbb{E} \left\{ \max_{j \in \{1, \dots, d\}} \exp \left[u \cdot \left| \frac{1}{n} \sum_{i=1}^n \bar{\omega}_i \cdot y_i \cdot x_{i,j} - \mathbb{E}(\bar{\omega} \cdot Y \cdot X_j) \right| \right] \right\} \\
&\leq \sum_{j=1}^d \mathbb{E} \left\{ \exp \left[u \cdot \left| \frac{1}{n} \sum_{i=1}^n \bar{\omega}_i \cdot y_i \cdot x_{i,j} - \mathbb{E}(\bar{\omega} \cdot Y \cdot X_j) \right| \right] \right\} \\
&\leq \sum_{j=1}^d \mathbb{E} \left\{ \exp \left[u \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot \bar{\omega}_i \cdot y_i \cdot x_{i,j} \right| \right] \right\}, \tag{H.49}
\end{aligned}$$

where ξ_1, \dots, ξ_n are i.i.d. Rademacher random variables. The last inequality follows from Lemma J.4. Furthermore, for the right-hand side of (H.49), we invoke the contraction result in Lemma J.5 by setting

$$f(y_i \cdot x_{i,j}) = y_i \cdot x_{i,j}, \quad \mathcal{F} = \{f\}, \quad \psi_i(v) = \bar{\omega}_i \cdot v, \quad \text{and} \quad \phi(v) = \exp(u \cdot v),$$

where u is the variable of the moment generating function in (H.49). From the definitions in (A.5) and (H.46) we have $|\bar{\omega}_i| = |2 \cdot \omega_{\beta}(\mathbf{x}_i, y_i) - 1| \leq 1$, which implies

$$|\psi_i(v) - \psi_i(v')| \leq \left| [2 \cdot \omega_{\beta}(\mathbf{x}_i, y_i) - 1] \cdot (v - v') \right| \leq |v - v'|, \quad \text{for all } v, v' \in \mathbb{R}.$$

By Lemma J.5, we obtain

$$\mathbb{E} \left\{ \exp \left[u \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot \bar{\omega}_i \cdot y_i \cdot x_{i,j} \right| \right] \right\} \leq \mathbb{E} \left\{ \exp \left[2 \cdot u \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_i \cdot x_{i,j} \right| \right] \right\} \tag{H.50}$$

for $j \in \{1, \dots, d\}$ on the right-hand side of (H.49). Recall that in mixture of regression model we have $y_i = z_i \cdot \langle \beta^*, \mathbf{x}_i \rangle + v_i$, where z_i is a Rademacher random variable, $v_i \sim N(0, \sigma^2)$, and $\mathbf{x}_i \sim N(\mathbf{0}, \mathbf{I}_d)$. Then by Example 5.8 in [28] we have $\|z_i \cdot \langle \beta^*, \mathbf{x}_i \rangle\|_{\psi_2} = \|\langle \beta^*, \mathbf{x}_i \rangle\|_{\psi_2} \leq C \cdot \|\beta^*\|_2$ and $\|v_{i,j}\|_{\psi_2} \leq C' \cdot \sigma$. By Lemma J.1 we further have

$$\|y_i\|_{\psi_2} = \|z_i \cdot \langle \beta^*, \mathbf{x}_i \rangle + v_i\|_{\psi_2} \leq \sqrt{C \cdot \|\beta^*\|_2^2 + C' \cdot \sigma^2}.$$

Note that we have $\|x_{i,j}\|_{\psi_2} \leq C''$ since $x_{i,j} \sim N(0, 1)$. Therefore, by Lemma J.2 we have

$$\|\xi_i \cdot y_i \cdot x_{i,j}\|_{\psi_1} = \|y_i \cdot x_{i,j}\|_{\psi_1} \leq \max \{C \cdot \|\beta^*\|_2^2 + C' \cdot \sigma^2, C''\} \leq C''' \cdot \max \{\|\beta^*\|_2^2 + \sigma^2, 1\}.$$

Since ξ_i is a Rademacher random variable independent of $y_i \cdot x_{i,j}$, we have $\mathbb{E}(\xi_i \cdot y_i \cdot x_{i,j}) = 0$. Hence, by Lemma 5.15 in [28], we obtain

$$\mathbb{E}[\exp(u' \cdot \xi_i \cdot y_i \cdot x_{i,j})] \leq \exp \left[C \cdot (u')^2 \cdot \max \{\|\beta^*\|_2^2 + \sigma^2, 1\}^2 \right] \tag{H.51}$$

for all $|u'| \leq C' / \max \{\|\beta^*\|_2^2 + \sigma^2, 1\}$. Hence we have

$$\begin{aligned}
& \mathbb{E} \left\{ \exp \left[2 \cdot u \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_i \cdot x_{i,j} \right| \right] \right\} \\
&\leq \mathbb{E} \left(\max \left\{ \exp \left[2 \cdot u \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_i \cdot x_{i,j} \right], \exp \left[-2 \cdot u \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_i \cdot x_{i,j} \right] \right\} \right) \\
&\leq \mathbb{E} \left\{ \exp \left[2 \cdot u \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_i \cdot x_{i,j} \right] \right\} + \mathbb{E} \left\{ \exp \left[-2 \cdot u \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_i \cdot x_{i,j} \right] \right\} \\
&\leq 2 \cdot \exp \left[C \cdot u^2 \cdot \max \{\|\beta^*\|_2^2 + \sigma^2, 1\}^2 / n \right]. \tag{H.52}
\end{aligned}$$

The last inequality is obtained by plugging (H.51) with $u' = 2 \cdot u/n$ and $u' = -2 \cdot u/n$ correspondingly into the two terms. Here $|u| \leq C' \cdot n / \max \{\|\beta^*\|_2^2 + \sigma^2, 1\}$. Plugging (H.52) into (H.50) and further into (H.49), we obtain

$$\mathbb{E} \left\{ \exp \left[u \cdot \left\| \frac{1}{n} \sum_{i=1}^n \bar{\omega}_i \cdot y_i \cdot \mathbf{x}_i - \mathbb{E}(\bar{\omega} \cdot Y \cdot \mathbf{X}) \right\|_{\infty} \right] \right\} \leq 2 \cdot d \cdot \exp \left[C \cdot u^2 \cdot \max \{\|\beta^*\|_2^2 + \sigma^2, 1\}^2 / n \right].$$

By Chernoff bound we have that, for all $v > 0$ and $|u| \leq C' \cdot n / \max \{ \|\beta^*\|_2^2 + \sigma^2, 1 \}$,

$$\begin{aligned} & \mathbb{P} \left[\left\| \frac{1}{n} \sum_{i=1}^n \bar{\omega}_i \cdot y_i \cdot \mathbf{x}_i - \mathbb{E}(\bar{\omega} \cdot Y \cdot \mathbf{X}) \right\|_{\infty} > v \right] \\ & \leq \mathbb{E} \left\{ \exp \left[u \cdot \left\| \frac{1}{n} \sum_{i=1}^n \bar{\omega}_i \cdot y_i \cdot \mathbf{x}_i - \mathbb{E}(\bar{\omega} \cdot Y \cdot \mathbf{X}) \right\|_{\infty} \right] \right\} / \exp(u \cdot v) \\ & \leq 2 \cdot \exp \left[C \cdot u^2 \cdot \max \{ \|\beta^*\|_2^2 + \sigma^2, 1 \}^2 / n - u \cdot v + \log d \right]. \end{aligned}$$

Minimizing over u on its right-hand side we have that, for $0 < v \leq C'' \cdot \max \{ \|\beta^*\|_2^2 + \sigma^2, 1 \}$,

$$\begin{aligned} & \mathbb{P} \left[\left\| \frac{1}{n} \sum_{i=1}^n \bar{\omega}_i \cdot y_i \cdot \mathbf{x}_i - \mathbb{E}(\bar{\omega} \cdot Y \cdot \mathbf{X}) \right\|_{\infty} > v \right] \\ & \leq 2 \cdot \exp \left\{ -n \cdot v^2 / \left[4 \cdot C \cdot \max \{ \|\beta^*\|_2^2 + \sigma^2, 1 \}^2 \right] + \log d \right\}. \end{aligned}$$

Setting the right-hand side of the above inequality to be $\delta/2$, we have that

$$\left\| \frac{1}{n} \sum_{i=1}^n \bar{\omega}_i \cdot y_i \cdot \mathbf{x}_i - \mathbb{E}(\bar{\omega} \cdot Y \cdot \mathbf{X}) \right\|_{\infty} \leq v = C \cdot \max \{ \|\beta^*\|_2^2 + \sigma^2, 1 \} \cdot \sqrt{\frac{\log d + \log(4/\delta)}{n}} \quad (\text{H.53})$$

holds with probability at least $1 - \delta/2$ for a sufficiently large n .

Analysis of Term (iii): For term (iii) in (H.47), by Lemma H.2 we have

$$\|\hat{\Theta} - \mathbf{I}_d\|_{1,\infty} \leq C \cdot \sqrt{\frac{\log d + \log(4/\delta)}{n}} \quad (\text{H.54})$$

with probability at least $1 - \delta/4$ for a sufficiently large n .

Analysis of Term (iv): For term (iv) in (H.47), recall that by (H.45) and (H.46) we have

$$M(\beta) = \mathbb{E} \left\{ [2 \cdot \omega_{\beta}(\mathbf{X}, Y) - 1] \cdot Y \cdot \mathbf{X} \right\} = \mathbb{E}(\bar{\omega} \cdot Y \cdot \mathbf{X}),$$

which implies

$$\begin{aligned} \|\mathbb{E}(\bar{\omega} \cdot Y \cdot \mathbf{X})\|_{\infty} &= \|M(\beta)\|_{\infty} \leq \|M(\beta) - \beta^*\|_2 + \|\beta^*\|_2 \\ &\leq \|\beta - \beta^*\|_2 + \|\beta^*\|_2 \leq (1 + 1/32) \cdot \|\beta^*\|_2, \end{aligned} \quad (\text{H.55})$$

where the first inequality follows from triangle inequality and $\|\cdot\|_{\infty} \leq \|\cdot\|_2$, the second inequality is from the proof of (G.8) in Lemma G.2 with $\bar{\beta}^{(t+0.5)}$ replaced with β and the fact that $\gamma_1/\nu < 1$ in (G.8), and the third inequality holds since in Condition *Statistical-Error*($\epsilon, \delta, s, n, \mathcal{B}$) we suppose that $\beta \in \mathcal{B}$, and for mixture of regression model \mathcal{B} is specified in (E.7).

Plugging (H.48), (H.53), (H.54) and (H.55) into (H.47), by union bound we have that

$$\begin{aligned} & \|M_n(\beta) - M(\beta)\|_{\infty} \\ & \leq C \cdot \max \{ \|\beta^*\|_2^2 + \sigma^2, 1 \} \cdot \sqrt{\frac{\log d + \log(4/\delta)}{n}} + C' \cdot \|\beta^*\|_2 \cdot \sqrt{\frac{\log d + \log(4/\delta)}{n}} \\ & \leq C'' \cdot \left[\max \{ \|\beta^*\|_2^2 + \sigma^2, 1 \} + \|\beta^*\|_2 \right] \cdot \sqrt{\frac{\log d + \log(4/\delta)}{n}} \end{aligned}$$

holds with probability at least $1 - \delta$. Therefore, we conclude the proof of (E.8) in Lemma E.5. \square

Gradient Ascent Implementation: In the following, we prove (E.9) in Lemma E.5.

Proof. Recall that $Q(\cdot; \cdot)$ is the expectation of $Q_n(\cdot; \cdot)$. According to (A.5) and (3.2), we have

$$M(\beta) = \beta + \eta \cdot \mathbb{E} [2 \cdot \omega_{\beta}(\mathbf{X}, Y) \cdot Y \cdot \mathbf{X} - \beta]$$

with $\omega_\beta(\cdot, \cdot)$ being the weight function defined in (A.5), which together with (A.9) implies

$$\begin{aligned} & \|M_n(\boldsymbol{\beta}) - M(\boldsymbol{\beta})\|_\infty \tag{H.56} \\ &= \left\| \eta \cdot \frac{1}{n} \sum_{i=1}^n [2 \cdot \omega_\beta(\mathbf{x}_i, y_i) \cdot y_i \cdot \mathbf{x}_i - \mathbf{x}_i \cdot \mathbf{x}_i^\top \cdot \boldsymbol{\beta}] - \eta \cdot \mathbb{E}[2 \cdot \omega_\beta(\mathbf{X}, Y) \cdot Y \cdot \mathbf{X} - \boldsymbol{\beta}] \right\|_\infty \\ &\leq \eta \cdot \underbrace{\left\| \frac{1}{n} \sum_{i=1}^n [2 \cdot \omega_\beta(\mathbf{x}_i, y_i) \cdot y_i \cdot \mathbf{x}_i] - \mathbb{E}[2 \cdot \omega_\beta(\mathbf{X}, Y) \cdot Y \cdot \mathbf{X}] \right\|_\infty}_{(i)} + \eta \cdot \underbrace{\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{x}_i^\top \cdot \boldsymbol{\beta} - \boldsymbol{\beta} \right\|_\infty}_{(ii)}. \end{aligned}$$

Here $\eta > 0$ denotes the stepsize in Algorithm 3.

Analysis of Term (i): For term (i) in (H.56), we redefine $\bar{\omega}_i$ and $\bar{\omega}$ in (H.46) as

$$\bar{\omega}_i = 2 \cdot \omega_\beta(\mathbf{x}_i, y_i), \quad \text{and} \quad \bar{\omega} = 2 \cdot \omega_\beta(\mathbf{X}, Y). \tag{H.57}$$

Note that $|\bar{\omega}_i| = |2 \cdot \omega_\beta(\mathbf{x}_i, y_i)| \leq 2$. Following the same way we establish the upper bound of term (ii) in (H.47), under the new definitions of $\bar{\omega}_i$ and $\bar{\omega}$ in (H.57) we have that

$$\begin{aligned} & \left\| \frac{1}{n} \sum_{i=1}^n [2 \cdot \omega_\beta(\mathbf{x}_i, y_i) \cdot y_i \cdot \mathbf{x}_i] - \mathbb{E}[2 \cdot \omega_\beta(\mathbf{X}, Y) \cdot Y \cdot \mathbf{X}] \right\|_\infty \\ & \leq C \cdot \max \{ \|\boldsymbol{\beta}^*\|_2^2 + \sigma^2, 1 \} \cdot \sqrt{\frac{\log d + \log(4/\delta)}{n}} \end{aligned}$$

holds with probability at least $1 - \delta/2$.

Analysis of Term (ii): For term (ii) in (H.56), we have

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{x}_i^\top \cdot \boldsymbol{\beta} - \boldsymbol{\beta} \right\|_\infty \leq \underbrace{\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{x}_i^\top - \mathbf{I}_d \right\|_{\infty, \infty}}_{(ii).a} \cdot \underbrace{\|\boldsymbol{\beta}\|_1}_{(ii).b}.$$

For term (ii).a, recall by our model assumption we have $\mathbb{E}(\mathbf{X} \cdot \mathbf{X}^\top) = \mathbf{I}_d$ and \mathbf{x}_i 's are the independent realizations of \mathbf{X} . Hence we have

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{x}_i^\top - \mathbf{I}_d \right\|_{\infty, \infty} \leq \max_{j \in \{1, \dots, d\}} \max_{k \in \{1, \dots, d\}} \left| \frac{1}{n} \sum_{i=1}^n x_{i,j} \cdot x_{i,k} - \mathbb{E}(X_j \cdot X_k) \right|.$$

Since $X_j, X_k \sim N(0, 1)$, according to Example 5.8 in [28] we have $\|X_j\|_{\psi_2} = \|X_k\|_{\psi_2} \leq C$. By Lemma J.2, $X_j \cdot X_k$ is a sub-exponential random variable with $\|X_j \cdot X_k\|_{\psi_1} \leq C'$. Moreover, we have $\|X_j \cdot X_k - \mathbb{E}(X_j \cdot X_k)\|_{\psi_1} \leq C''$ by Lemma J.3. Then by Bernstein's inequality (Proposition 5.16 in [28]) and union bound, we have

$$\mathbb{P} \left[\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{x}_i^\top - \mathbf{I}_d \right\|_{\infty, \infty} > v \right] \leq 2 \cdot d^2 \cdot \exp(-C \cdot n \cdot v^2)$$

for $0 < v \leq C'$ and a sufficiently large sample size n . Setting its right-hand side to be $\delta/2$, we have

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{x}_i^\top - \mathbf{I}_d \right\|_{\infty, \infty} \leq C'' \cdot \sqrt{\frac{2 \cdot \log d + \log(4/\delta)}{n}}$$

holds with probability at least $1 - \delta/2$. For term (ii).b we have $\|\boldsymbol{\beta}\|_1 \leq \sqrt{s} \cdot \|\boldsymbol{\beta}\|_2$, since in Condition *Statistical-Error*($\epsilon, \delta, s, n, \mathcal{B}$) we assume $\|\boldsymbol{\beta}\|_0 \leq s$. Furthermore, we have $\|\boldsymbol{\beta}\|_2 \leq \|\boldsymbol{\beta}^*\|_2 + \|\boldsymbol{\beta}^* - \boldsymbol{\beta}\|_2 \leq (1 + 1/32) \cdot \|\boldsymbol{\beta}^*\|_2$, because in Condition *Statistical-Error*($\epsilon, \delta, s, n, \mathcal{B}$) we assume that $\boldsymbol{\beta} \in \mathcal{B}$, and for mixture of regression model \mathcal{B} is specified in (E.7).

Plugging the above results into the right-hand side of (H.56), by union bound we have that

$$\begin{aligned} \|M_n(\boldsymbol{\beta}) - M(\boldsymbol{\beta})\|_\infty &\leq \eta \cdot C \cdot \max \{ \|\boldsymbol{\beta}^*\|_2^2 + \sigma^2, 1 \} \cdot \sqrt{\frac{\log d + \log(4/\delta)}{n}} \\ &\quad + \eta \cdot C' \cdot \sqrt{\frac{2 \cdot \log d + \log(4/\delta)}{n}} \cdot \sqrt{s} \cdot \|\boldsymbol{\beta}^*\|_2 \\ &\leq \eta \cdot C'' \cdot \max \{ \|\boldsymbol{\beta}^*\|_2^2 + \sigma^2, 1, \sqrt{s} \cdot \|\boldsymbol{\beta}^*\|_2 \} \cdot \sqrt{\frac{\log d + \log(4/\delta)}{n}} \end{aligned}$$

holds with probability at least $1 - \delta$. Therefore, we conclude the proof of (E.9) in Lemma E.5. \square

I Proof of Results for Inference

In the following, we provide the detailed proof of the theoretical results for asymptotic inference in §4. We first present the proof of the general results, and then the proof for specific models.

I.1 Proof of Theorem 2.1

Proof. In the sequel we establish the two equations in (2.8) respectively.

Proof of the First Equation: According to the definition of the lower bound function $Q_n(\cdot; \cdot)$ in (2.1), we have

$$Q_n(\beta'; \beta) = \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{Z}} k_{\beta}(\mathbf{z} | \mathbf{y}_i) \cdot \log f_{\beta'}(\mathbf{y}_i, \mathbf{z}) \, d\mathbf{z}. \quad (\text{I.1})$$

Here $k_{\beta}(\mathbf{z} | \mathbf{y}_i)$ is the density of the latent variable \mathbf{Z} conditioning on the observed variable $\mathbf{Y} = \mathbf{y}_i$ under the model with parameter β . Hence we obtain

$$\nabla_1 Q_n(\beta; \beta) = \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{Z}} k_{\beta}(\mathbf{z} | \mathbf{y}_i) \cdot \frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta}{f_{\beta}(\mathbf{y}_i, \mathbf{z})} \, d\mathbf{z} = \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{Z}} \frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta}{h_{\beta}(\mathbf{y}_i)} \, d\mathbf{z}, \quad (\text{I.2})$$

where $h_{\beta}(\mathbf{y}_i)$ is the marginal density function of \mathbf{Y} evaluated at \mathbf{y}_i , and the second equality follows from the fact that

$$k_{\beta}(\mathbf{z} | \mathbf{y}_i) = f_{\beta}(\mathbf{y}_i, \mathbf{z}) / h_{\beta}(\mathbf{y}_i), \quad (\text{I.3})$$

since $k_{\beta}(\mathbf{z} | \mathbf{y}_i)$ is the conditional density. According to the definition in (B.3), we have

$$\nabla \ell_n(\beta) = \sum_{i=1}^n \frac{\partial \log h_{\beta}(\mathbf{y}_i)}{\partial \beta} = \sum_{i=1}^n \frac{\partial h_{\beta}(\mathbf{y}_i) / \partial \beta}{h_{\beta}(\mathbf{y}_i)} = \sum_{i=1}^n \int_{\mathcal{Z}} \frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta}{h_{\beta}(\mathbf{y}_i)} \, d\mathbf{z}, \quad (\text{I.4})$$

where the last equality is from (B.1). Comparing (I.2) and (I.4), we obtain $\nabla_1 Q_n(\beta; \beta) = \nabla \ell_n(\beta) / n$.

Proof of the Second Equation: For the second equation in (D.1), by (I.1) and (I.3) we have

$$Q_n(\beta'; \beta) = \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{Z}} \frac{f_{\beta}(\mathbf{y}_i, \mathbf{z})}{h_{\beta}(\mathbf{y}_i)} \cdot \log f_{\beta'}(\mathbf{y}_i, \mathbf{z}) \, d\mathbf{z}.$$

By calculation we obtain

$$\begin{aligned} & \nabla_{1,2}^2 Q_n(\beta; \beta) \\ &= \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{Z}} \frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta}{f_{\beta}(\mathbf{y}_i, \mathbf{z})} \otimes \left\{ \frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta \cdot h_{\beta}(\mathbf{y}_i)}{[h_{\beta}(\mathbf{y}_i)]^2} - \frac{f_{\beta}(\mathbf{y}_i, \mathbf{z}) \cdot \partial h_{\beta}(\mathbf{y}_i) / \partial \beta}{[h_{\beta}(\mathbf{y}_i)]^2} \right\} \, d\mathbf{z}. \end{aligned} \quad (\text{I.5})$$

Here \otimes denotes the vector outer product. Note that in (I.5) we have

$$\begin{aligned} & \int_{\mathcal{Z}} \frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta}{f_{\beta}(\mathbf{y}_i, \mathbf{z})} \otimes \frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta}{h_{\beta}(\mathbf{y}_i)} \, d\mathbf{z} = \int_{\mathcal{Z}} \left[\frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta}{f_{\beta}(\mathbf{y}_i, \mathbf{z})} \right]^{\otimes 2} \cdot \frac{f_{\beta}(\mathbf{y}_i, \mathbf{z})}{h_{\beta}(\mathbf{y}_i)} \, d\mathbf{z} \\ &= \int_{\mathcal{Z}} \left[\frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta}{f_{\beta}(\mathbf{y}_i, \mathbf{z})} \right]^{\otimes 2} \cdot k_{\beta}(\mathbf{z} | \mathbf{y}_i) \, d\mathbf{z} \\ &= \mathbb{E}_{\beta} \left[\tilde{S}_{\beta}(\mathbf{Y}, \mathbf{Z})^{\otimes 2} \mid \mathbf{Y} = \mathbf{y}_i \right], \end{aligned} \quad (\text{I.6})$$

where $\mathbf{v}^{\otimes 2}$ denotes $\mathbf{v} \otimes \mathbf{v}$. Here $\tilde{S}_{\beta}(\cdot, \cdot)$ is defined as

$$\tilde{S}_{\beta}(\mathbf{y}, \mathbf{z}) = \frac{\partial \log f_{\beta}(\mathbf{y}, \mathbf{z})}{\partial \beta} = \frac{\partial f_{\beta}(\mathbf{y}, \mathbf{z}) / \partial \beta}{f_{\beta}(\mathbf{y}, \mathbf{z})} \in \mathbb{R}^d, \quad (\text{I.7})$$

i.e., the score function for the complete likelihood, which involves both the observed variable \mathbf{Y} and the latent variable \mathbf{Z} . Meanwhile, in (I.5) we have

$$\begin{aligned} \int_{\mathcal{Z}} \frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta}{f_{\beta}(\mathbf{y}_i, \mathbf{z})} \otimes \frac{f_{\beta}(\mathbf{y}_i, \mathbf{z}) \cdot \partial h_{\beta}(\mathbf{y}_i) / \partial \beta}{[h_{\beta}(\mathbf{y}_i)]^2} d\mathbf{z} &= \left[\int_{\mathcal{Z}} \frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta}{h_{\beta}(\mathbf{y}_i)} d\mathbf{z} \right] \otimes \frac{\partial h_{\beta}(\mathbf{y}_i) / \partial \beta}{h_{\beta}(\mathbf{y}_i)} \\ &= \left[\int_{\mathcal{Z}} \frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta}{h_{\beta}(\mathbf{y}_i)} d\mathbf{z} \right]^{\otimes 2}, \end{aligned} \quad (\text{I.8})$$

where in the last equality we utilize the fact that

$$\int_{\mathcal{Z}} f_{\beta}(\mathbf{y}_i, \mathbf{z}) d\mathbf{z} = h_{\beta}(\mathbf{y}_i), \quad (\text{I.9})$$

because $h_{\beta}(\cdot)$ is the marginal density function of \mathbf{Y} . By (I.3) and (I.7), for the right-hand side of (I.8) we have

$$\mathbb{E}_{\beta} \left[\tilde{S}_{\beta}(\mathbf{Y}, \mathbf{Z}) \mid \mathbf{Y} = \mathbf{y}_i \right] = \int_{\mathcal{Z}} \frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta}{f_{\beta}(\mathbf{y}_i, \mathbf{z})} \cdot k_{\beta}(\mathbf{z} \mid \mathbf{y}_i) d\mathbf{z} = \int_{\mathcal{Z}} \frac{\partial f_{\beta}(\mathbf{y}_i, \mathbf{z}) / \partial \beta}{h_{\beta}(\mathbf{y}_i)} d\mathbf{z}. \quad (\text{I.10})$$

Plugging (I.10) into (I.8) and then plugging (I.6) and (I.8) into (I.5) we obtain

$$\nabla_{1,2}^2 Q_n(\beta; \beta) = \frac{1}{n} \sum_{i=1}^n \left(\mathbb{E}_{\beta} \left[\tilde{S}_{\beta}(\mathbf{Y}, \mathbf{Z})^{\otimes 2} \mid \mathbf{Y} = \mathbf{y}_i \right] - \left\{ \mathbb{E}_{\beta} \left[\tilde{S}_{\beta}(\mathbf{Y}, \mathbf{Z}) \mid \mathbf{Y} = \mathbf{y}_i \right] \right\}^{\otimes 2} \right).$$

Setting $\beta = \beta^*$ in the above equality, we obtain

$$\mathbb{E}_{\beta^*} \left[\nabla_{1,2}^2 Q_n(\beta^*; \beta^*) \right] = \mathbb{E}_{\beta^*} \left\{ \text{Cov}_{\beta^*} \left[\tilde{S}_{\beta^*}(\mathbf{Y}, \mathbf{Z}) \mid \mathbf{Y} \right] \right\}. \quad (\text{I.11})$$

Meanwhile, for $\beta = \beta^*$, by the property of Fisher information we have

$$I(\beta^*) = \text{Cov}_{\beta^*} \left[\frac{\partial \log h_{\beta^*}(\mathbf{Y})}{\partial \beta} \right] = \text{Cov}_{\beta^*} \left\{ \mathbb{E}_{\beta^*} \left[\tilde{S}_{\beta^*}(\mathbf{Y}, \mathbf{Z}) \mid \mathbf{Y} \right] \right\}. \quad (\text{I.12})$$

Here the last equality is obtained by taking $\beta = \beta^*$ in

$$\begin{aligned} \frac{\partial \log h_{\beta}(\mathbf{Y})}{\partial \beta} &= \frac{\partial h_{\beta}(\mathbf{Y})}{\partial \beta} \cdot \frac{1}{h_{\beta}(\mathbf{Y})} = \int_{\mathcal{Z}} \frac{\partial f_{\beta}(\mathbf{Y}, \mathbf{z}) / \partial \beta}{h_{\beta}(\mathbf{Y})} d\mathbf{z} = \int_{\mathcal{Z}} \frac{\partial f_{\beta}(\mathbf{Y}, \mathbf{z}) / \partial \beta}{f_{\beta}(\mathbf{Y}, \mathbf{z})} \cdot k_{\beta}(\mathbf{z} \mid \mathbf{Y}) d\mathbf{z} \\ &= \int_{\mathcal{Z}} \tilde{S}_{\beta}(\mathbf{Y}, \mathbf{z}) \cdot k_{\beta}(\mathbf{z} \mid \mathbf{Y}) d\mathbf{z} \\ &= \mathbb{E}_{\beta} \left[\tilde{S}_{\beta}(\mathbf{Y}, \mathbf{Z}) \mid \mathbf{Y} \right], \end{aligned}$$

where the second equality follows from (I.9), the third equality follows from (I.3), while the second last equality follows from (I.7). Combining (I.11) and (I.12), by the law of total variance we have

$$\begin{aligned} I(\beta^*) + \mathbb{E}_{\beta^*} \left[\nabla_{1,2}^2 Q_n(\beta^*; \beta^*) \right] &= \text{Cov}_{\beta^*} \left\{ \mathbb{E}_{\beta^*} \left[\tilde{S}_{\beta^*}(\mathbf{Y}, \mathbf{Z}) \mid \mathbf{Y} \right] \right\} + \mathbb{E}_{\beta^*} \left\{ \text{Cov}_{\beta^*} \left[\tilde{S}_{\beta^*}(\mathbf{Y}, \mathbf{Z}) \mid \mathbf{Y} \right] \right\} \\ &= \text{Cov}_{\beta^*} \left[\tilde{S}_{\beta^*}(\mathbf{Y}, \mathbf{Z}) \right]. \end{aligned} \quad (\text{I.13})$$

In the following, we prove

$$\mathbb{E}_{\beta^*} \left[\nabla_{1,1}^2 Q_n(\beta^*; \beta^*) \right] = -\text{Cov}_{\beta^*} \left[\tilde{S}_{\beta^*}(\mathbf{Y}, \mathbf{Z}) \right]. \quad (\text{I.14})$$

According to (I.1) we have

$$\nabla_{1,1}^2 Q_n(\beta; \beta) = \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{Z}} k_{\beta}(\mathbf{z} \mid \mathbf{y}_i) \cdot \frac{\partial^2 \log f_{\beta}(\mathbf{y}_i, \mathbf{z})}{\partial^2 \beta} d\mathbf{z} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\beta} \left[\frac{\partial^2 \log f_{\beta}(\mathbf{Y}, \mathbf{Z})}{\partial^2 \beta} \mid \mathbf{Y} = \mathbf{y}_i \right]. \quad (\text{I.15})$$

Let $\tilde{\ell}(\beta) = \log f_{\beta}(\mathbf{Y}, \mathbf{Z})$ be the complete log-likelihood, which involves both the observed variable \mathbf{Y} and the latent variable \mathbf{Z} , and $\tilde{I}(\beta)$ be the corresponding Fisher information. By setting $\beta = \beta^*$ in (I.15) and taking expectation, we obtain

$$\mathbb{E}_{\beta^*} \left[\nabla_{1,1}^2 Q_n(\beta^*; \beta^*) \right] = \mathbb{E}_{\beta^*} \left\{ \mathbb{E}_{\beta^*} \left[\frac{\partial^2 \log f_{\beta^*}(\mathbf{Y}, \mathbf{Z})}{\partial^2 \beta} \mid \mathbf{Y} \right] \right\} = \mathbb{E}_{\beta^*} \left[\nabla^2 \tilde{\ell}(\beta^*) \right] = -\tilde{I}(\beta^*). \quad (\text{I.16})$$

Since $\tilde{S}_\beta(\mathbf{Y}, \mathbf{Z})$ defined in (I.7) is the score function for the complete log-likelihood $\tilde{\ell}(\beta)$, according to the relationship between the score function and Fisher information, we have

$$\tilde{I}(\beta^*) = \text{Cov}_{\beta^*} \left[\tilde{S}_{\beta^*}(\mathbf{Y}, \mathbf{Z}) \right],$$

which together with (I.16) implies (I.14). By further plugging (I.14) into (I.13), we obtain

$$\mathbb{E}_{\beta^*} \left[\nabla_{1,1}^2 Q_n(\beta^*; \beta^*) + \nabla_{1,2}^2 Q_n(\beta^*; \beta^*) \right] = -I(\beta^*),$$

which establishes the first equality of the second equation in (D.1). In addition, the second equality of the second equation in (D.1) follows from the property of Fisher information. Thus, we conclude the proof of Theorem D.1. \square

I.2 Auxiliary Lemmas for Proving Theorem 4.6

In this section, we lay out several lemmas on the Dantzig selector defined in (2.6). The first lemma, which is from [5], characterizes the cone condition for the Dantzig selector.

Lemma I.1. Any feasible solution \mathbf{w} in (2.6) satisfies

$$\left\| [w(\beta, \lambda) - \mathbf{w}]_{\overline{\mathcal{S}_w}} \right\|_1 \leq \left\| [w(\beta, \lambda) - \mathbf{w}]_{\mathcal{S}_w} \right\|_1,$$

where $w(\beta, \lambda)$ is the minimizer of (2.6), \mathcal{S}_w is the support of \mathbf{w} and $\overline{\mathcal{S}_w}$ is its complement.

Proof. See Lemma B.3 in [5] for a detailed proof. \square

In the sequel, we focus on analyzing $w(\hat{\beta}, \lambda)$. The results for $w(\hat{\beta}_0, \lambda)$ can be obtained similarly. The next lemma characterizes the restricted eigenvalue of $T_n(\hat{\beta})$, which is defined as

$$\hat{\rho}_{\min} = \inf_{\mathbf{v} \in \mathcal{C}} \frac{\mathbf{v}^\top \cdot [-T_n(\hat{\beta})] \cdot \mathbf{v}}{\|\mathbf{v}\|_2^2}, \quad \text{where } \mathcal{C} = \left\{ \mathbf{v} : \|\mathbf{v}_{\overline{\mathcal{S}_w^*}}\|_1 \leq \|\mathbf{v}_{\mathcal{S}_w^*}\|_1, \mathbf{v} \neq \mathbf{0} \right\}. \quad (\text{I.17})$$

Here \mathcal{S}_w^* is the support of \mathbf{w}^* defined in (4.1).

Lemma I.2. Under Assumption 4.5 and Conditions 4.1, 4.3 and 4.4, for a sufficiently large sample size n , we have $\hat{\rho}_{\min} \geq \rho_{\min}/2 > 0$ with high probability, where ρ_{\min} is specified in (4.4).

Proof. By triangle inequality we have

$$\begin{aligned} \hat{\rho}_{\min} &\geq \inf_{\mathbf{v} \in \mathcal{C}} \frac{\mathbf{v}^\top \cdot [-T_n(\hat{\beta})] \cdot \mathbf{v}}{\|\mathbf{v}\|_2^2} \geq \inf_{\mathbf{v} \in \mathcal{C}} \frac{\mathbf{v}^\top \cdot I(\beta^*) \cdot \mathbf{v} - \left| \mathbf{v}^\top \cdot [I(\beta^*) + T_n(\hat{\beta})] \cdot \mathbf{v} \right|}{\|\mathbf{v}\|_2^2} \\ &\geq \underbrace{\inf_{\mathbf{v} \in \mathcal{C}} \frac{\mathbf{v}^\top \cdot I(\beta^*) \cdot \mathbf{v}}{\|\mathbf{v}\|_2^2}}_{(i)} - \underbrace{\sup_{\mathbf{v} \in \mathcal{C}} \frac{\left| \mathbf{v}^\top \cdot [I(\beta^*) + T_n(\hat{\beta})] \cdot \mathbf{v} \right|}{\|\mathbf{v}\|_2^2}}_{(ii)}, \end{aligned} \quad (\text{I.18})$$

where \mathcal{C} is defined in (I.17).

Analysis of Term (i): For term (i) in (I.18), by (4.4) in Assumption 4.5 we have

$$\inf_{\mathbf{v} \in \mathcal{C}} \frac{\mathbf{v}^\top \cdot I(\beta^*) \cdot \mathbf{v}}{\|\mathbf{v}\|_2^2} \geq \inf_{\mathbf{v} \neq \mathbf{0}} \frac{\mathbf{v}^\top \cdot I(\beta^*) \cdot \mathbf{v}}{\|\mathbf{v}\|_2^2} = \lambda_d[I(\beta^*)] \geq \rho_{\min}. \quad (\text{I.19})$$

Analysis of Term (ii): For term (ii) in (I.18) we have

$$\sup_{\mathbf{v} \in \mathcal{C}} \frac{\left| \mathbf{v}^\top \cdot [I(\beta^*) + T_n(\hat{\beta})] \cdot \mathbf{v} \right|}{\|\mathbf{v}\|_2^2} \leq \sup_{\mathbf{v} \in \mathcal{C}} \frac{\|\mathbf{v}\|_1^2 \cdot \left\| [I(\beta^*) + T_n(\hat{\beta})] \right\|_{\infty, \infty}}{\|\mathbf{v}\|_2^2}. \quad (\text{I.20})$$

By the definition of \mathcal{C} in (I.17), for any $\mathbf{v} \in \mathcal{C}$ we have

$$\|\mathbf{v}\|_1 = \|\mathbf{v}_{\mathcal{S}_w^*}\|_1 + \|\mathbf{v}_{\overline{\mathcal{S}_w^*}}\|_1 \leq 2 \cdot \|\mathbf{v}_{\mathcal{S}_w^*}\|_1 \leq 2 \cdot \sqrt{s_w^*} \cdot \|\mathbf{v}_{\mathcal{S}_w^*}\|_2 \leq 2 \cdot \sqrt{s_w^*} \cdot \|\mathbf{v}\|_2.$$

Therefore, the right-hand side of (I.20) is upper bounded by

$$4 \cdot s_w^* \cdot \left\| [I(\beta^*) + T_n(\hat{\beta})] \right\|_{\infty, \infty} \leq \underbrace{4 \cdot s_w^* \cdot \left\| [I(\beta^*) + T_n(\beta^*)] \right\|_{\infty, \infty}}_{(ii).a} + \underbrace{4 \cdot s_w^* \cdot \left\| [T_n(\hat{\beta}) - T_n(\beta^*)] \right\|_{\infty, \infty}}_{(ii).b}.$$

For term (ii).a, by Theorem D.1 and Condition 4.3 we have

$$4 \cdot s_{\mathbf{w}}^* \cdot \|I(\boldsymbol{\beta}^*) + T_n(\boldsymbol{\beta}^*)\|_{\infty, \infty} = 4 \cdot s_{\mathbf{w}}^* \cdot \|T_n(\boldsymbol{\beta}^*) - \mathbb{E}_{\boldsymbol{\beta}^*} [T_n(\boldsymbol{\beta}^*)]\|_{\infty, \infty} = O_{\mathbb{P}}(s_{\mathbf{w}}^* \cdot \zeta^{\text{T}}) = o_{\mathbb{P}}(1),$$

where the last equality is from (4.6) in Assumption 4.5, since for λ specified in (4.5) we have

$$s_{\mathbf{w}}^* \cdot \zeta^{\text{T}} \leq s_{\mathbf{w}}^* \cdot \lambda \leq \max\{\|\mathbf{w}^*\|_1, 1\} \cdot s_{\mathbf{w}}^* \cdot \lambda = o(1).$$

For term (ii).b, by Conditions 4.1 and 4.4 we have

$$4 \cdot s_{\mathbf{w}}^* \cdot \|T(\widehat{\boldsymbol{\beta}}) - T_n(\boldsymbol{\beta}^*)\|_{\infty, \infty} = 4 \cdot s_{\mathbf{w}}^* \cdot O_{\mathbb{P}}(\zeta^{\text{L}}) \cdot \|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_1 = O_{\mathbb{P}}(s_{\mathbf{w}}^* \cdot \zeta^{\text{L}} \cdot \zeta^{\text{EM}}) = o_{\mathbb{P}}(1),$$

where the last equality is also from (4.6) in Assumption 4.5, since for λ specified in (4.5) we have

$$s_{\mathbf{w}}^* \cdot \zeta^{\text{L}} \cdot \zeta^{\text{EM}} \leq s_{\mathbf{w}}^* \cdot \lambda \leq \max\{\|\mathbf{w}^*\|_1, 1\} \cdot s_{\mathbf{w}}^* \cdot \lambda = o(1).$$

Hence, term (ii) in (I.18) is $o_{\mathbb{P}}(1)$. Since ρ_{\min} is a constant, for a sufficiently large n we have that term (ii) is upper bounded by $\rho_{\min}/2$ with high probability. Further by plugging this and (I.19) into (I.18), we conclude that $\widehat{\rho}_{\min} \geq \rho_{\min}/2$ holds with high probability. \square

The next lemma quantifies the statistical accuracy of $w(\widehat{\boldsymbol{\beta}}, \lambda)$, where $w(\cdot, \cdot)$ is defined in (2.6).

Lemma I.3. Under Assumption 4.5 and Conditions 4.1-4.4, for λ specified in (4.5) we have that

$$\max\left\{\|w(\widehat{\boldsymbol{\beta}}, \lambda) - \mathbf{w}^*\|_1, \|w(\widehat{\boldsymbol{\beta}}_0, \lambda) - \mathbf{w}^*\|_1\right\} \leq \frac{16 \cdot s_{\mathbf{w}}^* \cdot \lambda}{\rho_{\min}}$$

holds with high probability. Here ρ_{\min} is specified in (4.4), while \mathbf{w}^* and $s_{\mathbf{w}}^*$ are defined (4.1).

Proof. For λ specified in (4.5), we verify that \mathbf{w}^* is a feasible solution in (2.6) with high probability. For notational simplicity, we define the following event

$$\mathcal{E} = \left\{ \left\| [T_n(\widehat{\boldsymbol{\beta}})]_{\gamma, \alpha} - [T_n(\widehat{\boldsymbol{\beta}})]_{\gamma, \gamma} \cdot \mathbf{w}^* \right\|_{\infty} \leq \lambda \right\}. \quad (\text{I.21})$$

By the definition of \mathbf{w}^* in (4.1), we have $[I(\boldsymbol{\beta}^*)]_{\gamma, \alpha} - [I(\boldsymbol{\beta}^*)]_{\gamma, \gamma} \cdot \mathbf{w}^* = 0$. Hence, we have

$$\begin{aligned} \left\| [T_n(\widehat{\boldsymbol{\beta}})]_{\gamma, \alpha} - [T_n(\widehat{\boldsymbol{\beta}})]_{\gamma, \gamma} \cdot \mathbf{w}^* \right\|_{\infty} &= \left\| [T_n(\widehat{\boldsymbol{\beta}}) + I(\boldsymbol{\beta}^*)]_{\gamma, \alpha} - [T_n(\widehat{\boldsymbol{\beta}}) + I(\boldsymbol{\beta}^*)]_{\gamma, \gamma} \cdot \mathbf{w}^* \right\|_{\infty} \\ &\leq \left\| T_n(\widehat{\boldsymbol{\beta}}) + I(\boldsymbol{\beta}^*) \right\|_{\infty, \infty} + \left\| T_n(\widehat{\boldsymbol{\beta}}) + I(\boldsymbol{\beta}^*) \right\|_{\infty, \infty} \cdot \|\mathbf{w}^*\|_1, \end{aligned} \quad (\text{I.22})$$

where the last inequality is from triangle inequality and Hölder's inequality. Note that we have

$$\left\| T_n(\widehat{\boldsymbol{\beta}}) + I(\boldsymbol{\beta}^*) \right\|_{\infty, \infty} \leq \left\| T_n(\boldsymbol{\beta}^*) + I(\boldsymbol{\beta}^*) \right\|_{\infty, \infty} + \left\| T(\widehat{\boldsymbol{\beta}}) - T_n(\boldsymbol{\beta}^*) \right\|_{\infty, \infty}. \quad (\text{I.23})$$

On the right-hand side, by Theorem D.1 and Condition 4.3 we have

$$\left\| T_n(\boldsymbol{\beta}^*) + I(\boldsymbol{\beta}^*) \right\|_{\infty, \infty} = \left\| T_n(\boldsymbol{\beta}^*) - \mathbb{E}_{\boldsymbol{\beta}^*} [T_n(\boldsymbol{\beta}^*)] \right\|_{\infty, \infty} = O_{\mathbb{P}}(\zeta^{\text{T}}),$$

while by Conditions 4.1 and 4.4 we have

$$\left\| T(\widehat{\boldsymbol{\beta}}) - T_n(\boldsymbol{\beta}^*) \right\|_{\infty, \infty} = O_{\mathbb{P}}(\zeta^{\text{L}}) \cdot \|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_1 = O_{\mathbb{P}}(\zeta^{\text{L}} \cdot \zeta^{\text{EM}}).$$

Plugging the above equations into (I.23) and further plugging (I.23) into (I.22), by (4.5) we have

$$\left\| [T_n(\widehat{\boldsymbol{\beta}})]_{\gamma, \alpha} - [T_n(\widehat{\boldsymbol{\beta}})]_{\gamma, \gamma} \cdot \mathbf{w}^* \right\|_{\infty} \leq C \cdot (\zeta^{\text{T}} + \zeta^{\text{L}} \cdot \zeta^{\text{EM}}) \cdot (1 + \|\mathbf{w}^*\|_1) = \lambda$$

holds with high probability for a sufficiently large constant $C \geq 1$. In other words, \mathcal{E} occurs with high probability. The subsequent proof will be conditioning on \mathcal{E} and the following event

$$\mathcal{E}' = \{\widehat{\rho}_{\min} \geq \rho_{\min}/2 > 0\}, \quad (\text{I.24})$$

which also occurs with high probability according to Lemma I.2. Here $\widehat{\rho}_{\min}$ is defined in (I.17).

For notational simplicity, we denote $w(\widehat{\boldsymbol{\beta}}, \lambda) = \widehat{\mathbf{w}}$. By triangle inequality we have

$$\begin{aligned} &\left\| [T_n(\widehat{\boldsymbol{\beta}})]_{\gamma, \gamma} \cdot (\widehat{\mathbf{w}} - \mathbf{w}^*) \right\|_{\infty} \\ &\leq \left\| [T_n(\widehat{\boldsymbol{\beta}})]_{\gamma, \alpha} - [T_n(\widehat{\boldsymbol{\beta}})]_{\gamma, \gamma} \cdot \mathbf{w}^* \right\|_{\infty} + \left\| [T_n(\widehat{\boldsymbol{\beta}})]_{\gamma, \gamma} \cdot \widehat{\mathbf{w}} - [T_n(\widehat{\boldsymbol{\beta}})]_{\gamma, \alpha} \right\|_{\infty} \\ &\leq 2 \cdot \lambda, \end{aligned} \quad (\text{I.25})$$

where the last inequality follows from (2.6) and (I.21). Moreover, by (I.17) and (I.24) we have

$$(\widehat{\mathbf{w}} - \mathbf{w}^*)^\top \cdot [-T_n(\widehat{\beta})]_{\gamma, \gamma} \cdot (\widehat{\mathbf{w}} - \mathbf{w}^*) \geq \widehat{\rho}_{\min} \cdot \|\widehat{\mathbf{w}} - \mathbf{w}^*\|_2^2 \geq \rho_{\min}/2 \cdot \|\widehat{\mathbf{w}} - \mathbf{w}^*\|_2^2. \quad (\text{I.26})$$

Meanwhile, by Lemma I.1 we have

$$\|\widehat{\mathbf{w}} - \mathbf{w}^*\|_1 = \|(\widehat{\mathbf{w}} - \mathbf{w}^*)_{S_{\mathbf{w}}^*}\|_1 + \|(\widehat{\mathbf{w}} - \mathbf{w}^*)_{\overline{S_{\mathbf{w}}^*}}\|_1 \leq 2 \cdot \|(\widehat{\mathbf{w}} - \mathbf{w}^*)_{S_{\mathbf{w}}^*}\|_1 \leq 2 \cdot \sqrt{s_{\mathbf{w}}^*} \cdot \|\widehat{\mathbf{w}} - \mathbf{w}^*\|_2.$$

Plugging the above inequality into (I.26), we obtain

$$(\widehat{\mathbf{w}} - \mathbf{w}^*)^\top \cdot [-T_n(\widehat{\beta})]_{\gamma, \gamma} \cdot (\widehat{\mathbf{w}} - \mathbf{w}^*) \geq \rho_{\min}/(8 \cdot s_{\mathbf{w}}^*) \cdot \|\widehat{\mathbf{w}} - \mathbf{w}^*\|_1^2. \quad (\text{I.27})$$

Note that by (I.25), the left-hand side of (I.27) is upper bounded by

$$\|\widehat{\mathbf{w}} - \mathbf{w}^*\|_1 \cdot \left\| [T_n(\widehat{\beta})]_{\gamma, \gamma} \cdot (\widehat{\mathbf{w}} - \mathbf{w}^*) \right\|_\infty \leq \|\widehat{\mathbf{w}} - \mathbf{w}^*\|_1 \cdot 2 \cdot \lambda. \quad (\text{I.28})$$

By (I.27) and (I.28), we then obtain $\|\widehat{\mathbf{w}} - \mathbf{w}^*\|_1 \leq 16 \cdot s_{\mathbf{w}}^* \cdot \lambda / \rho_{\min}$ conditioning on \mathcal{E} and \mathcal{E}' , both of which hold with high probability. Note that the proof for $w(\widehat{\beta}_0, \lambda)$ follows similarly. Therefore, we conclude the proof of Lemma I.3. \square

I.3 Proof of Lemma G.3

Proof. Our proof strategy is as follows. First we prove that

$$\sqrt{n} \cdot S_n(\widehat{\beta}_0, \lambda) = \sqrt{n} \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_\alpha - \sqrt{n} \cdot (\mathbf{w}^*)^\top \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_\gamma + o_{\mathbb{P}}(1), \quad (\text{I.29})$$

where β^* is the true parameter and \mathbf{w}^* is defined in (4.1). We then prove

$$\sqrt{n} \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_\alpha - \sqrt{n} \cdot (\mathbf{w}^*)^\top \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_\gamma \xrightarrow{D} N(0, [I(\beta^*)]_{\alpha|\gamma}), \quad (\text{I.30})$$

where $[I(\beta^*)]_{\alpha|\gamma}$ is defined in (4.2). Throughout the proof, we abbreviate $w(\widehat{\beta}_0, \lambda)$ as $\widehat{\mathbf{w}}_0$. Also, it is worth noting that our analysis is under the null hypothesis where $\beta^* = [\alpha^*, (\gamma^*)^\top]^\top$ with $\alpha^* = 0$.

Proof of (I.29): For (I.29), by the definition of the decorrelated score function in (2.5) we have

$$S_n(\widehat{\beta}_0, \lambda) = [\nabla_1 Q_n(\widehat{\beta}_0; \widehat{\beta}_0)]_\alpha - \widehat{\mathbf{w}}_0^\top \cdot [\nabla_1 Q_n(\widehat{\beta}_0; \widehat{\beta}_0)]_\gamma.$$

By the mean-value theorem, we obtain

$$\begin{aligned} S_n(\widehat{\beta}_0, \lambda) &= \overbrace{[\nabla_1 Q_n(\beta^*; \beta^*)]_\alpha - \widehat{\mathbf{w}}_0^\top \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_\gamma}^{(i)} \\ &\quad + \underbrace{[T_n(\beta^\sharp)]_{\gamma, \alpha}^\top \cdot (\widehat{\beta}_0 - \beta^*) - \widehat{\mathbf{w}}_0^\top \cdot [T_n(\beta^\sharp)]_{\gamma, \gamma} \cdot (\widehat{\beta}_0 - \beta^*)}_{(ii)}, \end{aligned} \quad (\text{I.31})$$

where we have $T_n(\beta) = \nabla_{1,1}^2 Q_n(\beta; \beta) + \nabla_{1,2}^2 Q_n(\beta; \beta)$ as defined in (2.4), and β^\sharp is an intermediate value between β^* and $\widehat{\beta}_0$.

Analysis of Term (i): For term (i) in (I.31), we have

$$\begin{aligned} &[\nabla_1 Q_n(\beta^*; \beta^*)]_\alpha - \widehat{\mathbf{w}}_0^\top \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_\gamma \\ &= [\nabla_1 Q_n(\beta^*; \beta^*)]_\alpha - (\mathbf{w}^*)^\top \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_\gamma + (\mathbf{w}^* - \widehat{\mathbf{w}}_0)^\top \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_\gamma. \end{aligned} \quad (\text{I.32})$$

For the right-hand side of (I.32), we have

$$(\mathbf{w}^* - \widehat{\mathbf{w}}_0)^\top \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_\gamma \leq \|\mathbf{w}^* - \widehat{\mathbf{w}}_0\|_1 \cdot \left\| [\nabla_1 Q_n(\beta^*; \beta^*)]_\gamma \right\|_\infty. \quad (\text{I.33})$$

By Lemma I.3, we have $\|\mathbf{w}^* - \widehat{\mathbf{w}}_0\|_1 = O_{\mathbb{P}}(s_{\mathbf{w}}^* \cdot \lambda)$, where λ is specified in (4.5). Meanwhile, we have

$$\begin{aligned} \left\| [\nabla_1 Q_n(\beta^*; \beta^*)]_\gamma \right\|_\infty &\leq \left\| \nabla_1 Q_n(\beta^*; \beta^*) \right\|_\infty = \left\| \nabla_1 Q_n(\beta^*; \beta^*) - \nabla_1 Q(\beta^*; \beta^*) \right\|_\infty \\ &= O_{\mathbb{P}}(\zeta^G), \end{aligned}$$

where the first equality follows from the self-consistency property [18] that $\beta^* = \operatorname{argmax}_{\beta} Q(\beta; \beta^*)$, which gives $\nabla_1 Q(\beta^*; \beta^*) = \mathbf{0}$. Here the last equality is from Condition 4.2. Therefore, (I.33) implies

$$(\mathbf{w}^* - \widehat{\mathbf{w}}_0)^\top \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_\gamma = O_{\mathbb{P}}(s_{\mathbf{w}}^* \cdot \lambda \cdot \zeta^G) = o_{\mathbb{P}}(1/\sqrt{n}),$$

where the second equality is from $s_{\mathbf{w}}^* \cdot \lambda \cdot \zeta^G = o(1/\sqrt{n})$ in (4.6) of Assumption 4.5. Thus, by (I.32) we conclude that term (i) in (I.31) equals

$$[\nabla_1 Q_n(\beta^*; \beta^*)]_{\alpha} - (\mathbf{w}^*)^{\top} \cdot [\nabla_1 Q_n(\beta^*; \beta^*)]_{\gamma} + o_{\mathbb{P}}(1/\sqrt{n}).$$

Analysis of Term (ii): By triangle inequality, term (ii) in (I.31) is upper bounded by

$$\begin{aligned} & \underbrace{\left| [T_n(\widehat{\beta}_0)]_{\gamma, \alpha}^{\top} \cdot (\widehat{\beta}_0 - \beta^*) - \widehat{\mathbf{w}}_0^{\top} \cdot [T_n(\widehat{\beta}_0)]_{\gamma, \gamma} \cdot (\widehat{\beta}_0 - \beta^*) \right|}_{\text{(ii).a}} \\ & + \underbrace{\left| [T_n(\beta^{\sharp})]_{\gamma, \alpha}^{\top} \cdot (\widehat{\beta}_0 - \beta^*) - [T_n(\widehat{\beta}_0)]_{\gamma, \alpha}^{\top} \cdot (\widehat{\beta}_0 - \beta^*) \right|}_{\text{(ii).b}} \\ & + \underbrace{\left| \widehat{\mathbf{w}}_0^{\top} \cdot [T_n(\widehat{\beta}_0)]_{\gamma, \gamma} \cdot (\widehat{\beta}_0 - \beta^*) - \widehat{\mathbf{w}}_0^{\top} \cdot [T_n(\beta^{\sharp})]_{\gamma, \gamma} \cdot (\widehat{\beta}_0 - \beta^*) \right|}_{\text{(ii).c}}. \end{aligned} \quad (\text{I.34})$$

By Hölder's inequality, term (ii).a in (I.34) is upper bounded by

$$\begin{aligned} \|\widehat{\beta}_0 - \beta^*\|_1 \cdot \left\| [T_n(\widehat{\beta}_0)]_{\gamma, \alpha} - \widehat{\mathbf{w}}_0^{\top} \cdot [T_n(\widehat{\beta}_0)]_{\gamma, \gamma} \right\|_{\infty} &= \|\widehat{\beta}_0 - \beta^*\|_1 \cdot \lambda \\ &\leq O_{\mathbb{P}}(\zeta^{\text{EM}}) \cdot \lambda = o_{\mathbb{P}}(1/\sqrt{n}), \end{aligned} \quad (\text{I.35})$$

where the first inequality holds because $\widehat{\mathbf{w}}_0 = w(\widehat{\beta}_0, \lambda)$ is a feasible solution in (2.6). Meanwhile, Condition 4.1 gives $\|\widehat{\beta} - \beta^*\|_1 = O_{\mathbb{P}}(\zeta^{\text{EM}})$. Also note that by definition we have $(\widehat{\beta}_0)_{\alpha} = (\beta^*)_{\alpha} = 0$, which implies $\|\widehat{\beta}_0 - \beta^*\|_1 \leq \|\widehat{\beta} - \beta^*\|_1$. Hence, we have

$$\|\widehat{\beta}_0 - \beta^*\|_1 = O_{\mathbb{P}}(\zeta^{\text{EM}}), \quad (\text{I.36})$$

which implies the first equality in (I.35). The last equality in (I.35) follows from $\zeta^{\text{EM}} \cdot \lambda = o(1/\sqrt{n})$ in (4.6) of Assumption 4.5. Note that term (ii).b in (I.34) is upper bounded by

$$\left\| [T_n(\beta^{\sharp})]_{\gamma, \alpha} - [T_n(\widehat{\beta}_0)]_{\gamma, \alpha} \right\|_{\infty} \cdot \|\widehat{\beta}_0 - \beta^*\|_1 \leq \left\| T_n(\beta^{\sharp}) - T_n(\widehat{\beta}_0) \right\|_{\infty, \infty} \cdot \|\widehat{\beta}_0 - \beta^*\|_1. \quad (\text{I.37})$$

For the first term on the right-hand side of (I.37), by triangle inequality we have

$$\left\| T_n(\beta^{\sharp}) - T_n(\widehat{\beta}_0) \right\|_{\infty, \infty} \leq \left\| T_n(\beta^{\sharp}) - T_n(\beta^*) \right\|_{\infty, \infty} + \left\| T_n(\widehat{\beta}_0) - T_n(\beta^*) \right\|_{\infty, \infty}.$$

By Condition 4.4, we have

$$\left\| T_n(\widehat{\beta}_0) - T_n(\beta^*) \right\|_{\infty, \infty} = O_{\mathbb{P}}(\zeta^{\text{L}}) \cdot \|\widehat{\beta}_0 - \beta^*\|_1, \quad (\text{I.38})$$

and

$$\left\| T_n(\beta^{\sharp}) - T_n(\beta^*) \right\|_{\infty, \infty} = O_{\mathbb{P}}(\zeta^{\text{L}}) \cdot \|\beta^{\sharp} - \beta^*\|_1 \leq O_{\mathbb{P}}(\zeta^{\text{L}}) \cdot \|\widehat{\beta}_0 - \beta^*\|_1, \quad (\text{I.39})$$

where the last inequality in (I.39) holds because β^{\sharp} is defined as an intermediate value between β^* and $\widehat{\beta}_0$. Further by plugging (I.36) into (I.38), (I.39) as well as the second term on the right-hand side of (I.37), we have that term (ii).b in (I.34) is $O_{\mathbb{P}}[\zeta^{\text{L}} \cdot (\zeta^{\text{EM}})^2]$. Moreover, by our assumption in (4.6) of Assumption 4.5 we have

$$\zeta^{\text{L}} \cdot (\zeta^{\text{EM}})^2 \leq \max\{1, \|\mathbf{w}^*\|_1\} \cdot \zeta^{\text{L}} \cdot (\zeta^{\text{EM}})^2 = o(1/\sqrt{n}).$$

Thus, we conclude that term (ii).b is $o_{\mathbb{P}}(1/\sqrt{n})$. Similarly, term (ii).c in (I.34) is upper bounded by

$$\|\widehat{\mathbf{w}}_0\|_1 \cdot \left\| T_n(\beta^{\sharp}) - T_n(\widehat{\beta}_0) \right\|_{\infty, \infty} \cdot \|\widehat{\beta}_0 - \beta^*\|_1. \quad (\text{I.40})$$

By triangle inequality and Lemma I.3, the first term in (I.40) is upper bounded by

$$\|\mathbf{w}^*\|_1 + \|\widehat{\mathbf{w}}_0 - \mathbf{w}^*\|_1 = \|\mathbf{w}^*\|_1 + O_{\mathbb{P}}(s_{\mathbf{w}}^* \cdot \lambda).$$

Meanwhile, for the second and third terms in (I.40), by the same analysis for term (ii).b in (I.34) we have

$$\left\| T_n(\beta^{\sharp}) - T_n(\widehat{\beta}_0) \right\|_{\infty, \infty} \cdot \|\widehat{\beta}_0 - \beta^*\|_1 = O_{\mathbb{P}}[\zeta^{\text{L}} \cdot (\zeta^{\text{EM}})^2].$$

By (4.6) in Assumption 4.5, since $s_{\mathbf{w}^*} \cdot \lambda = o(1)$, we have

$$\begin{aligned} (\|\mathbf{w}^*\|_1 + s_{\mathbf{w}^*} \cdot \lambda) \cdot \zeta^L \cdot (\zeta^{\text{EM}})^2 &\leq [\max\{1, \|\mathbf{w}^*\|_1\} + o(1)] \cdot \zeta^L \cdot (\zeta^{\text{EM}})^2 \\ &= o(1/\sqrt{n}). \end{aligned}$$

Therefore, term (ii).c in (I.34) is $o_{\mathbb{P}}(1/\sqrt{n})$. Hence, by (I.34) we conclude that term (ii) in (I.31) is $o_{\mathbb{P}}(1/\sqrt{n})$. Combining the analysis for terms (i) and (ii) in (I.31), we then obtain (I.29). In the sequel, we turn to prove the second part on asymptotic normality.

Proof of (I.30): Note that by Theorem D.1, we have

$$\begin{aligned} \sqrt{n} \cdot [\nabla_1 Q_n(\boldsymbol{\beta}^*; \boldsymbol{\beta}^*)]_{\alpha} - \sqrt{n} \cdot (\mathbf{w}^*)^{\top} \cdot [\nabla_1 Q_n(\boldsymbol{\beta}^*; \boldsymbol{\beta}^*)]_{\gamma} &= \sqrt{n} \cdot [1, -(\mathbf{w}^*)^{\top}] \cdot \nabla_1 Q_n(\boldsymbol{\beta}^*; \boldsymbol{\beta}^*) \\ &= \sqrt{n} \cdot [1, -(\mathbf{w}^*)^{\top}] \cdot \nabla \ell_n(\boldsymbol{\beta}^*)/n. \end{aligned} \tag{I.41}$$

Recall that $\ell_n(\boldsymbol{\beta}^*)$ is the log-likelihood function defined in (B.3). Hence, $[1, -(\mathbf{w}^*)^{\top}] \cdot \nabla \ell_n(\boldsymbol{\beta}^*)/n$ is the average of n independent random variables. Meanwhile, the score function has mean zero at $\boldsymbol{\beta}^*$, i.e., $\mathbb{E}[\nabla \ell_n(\boldsymbol{\beta}^*)] = \mathbf{0}$. For the variance of the rescaled average in (I.41), we have

$$\begin{aligned} \text{Var} \left\{ \sqrt{n} \cdot [1, -(\mathbf{w}^*)^{\top}] \cdot \nabla \ell_n(\boldsymbol{\beta}^*)/n \right\} &= [1, -(\mathbf{w}^*)^{\top}] \cdot \text{Cov}[\nabla \ell_n(\boldsymbol{\beta}^*)/\sqrt{n}] \cdot [1, -(\mathbf{w}^*)^{\top}]^{\top} \\ &= [1, -(\mathbf{w}^*)^{\top}] \cdot I(\boldsymbol{\beta}^*) \cdot [1, -(\mathbf{w}^*)^{\top}]^{\top}. \end{aligned}$$

Here the second equality is from the fact that the covariance of the score function equals the Fisher information (up to renormalization). Hence, the variance of each item in the average in (I.41) is

$$\begin{aligned} [1, -(\mathbf{w}^*)^{\top}] \cdot I(\boldsymbol{\beta}^*) \cdot [1, -(\mathbf{w}^*)^{\top}]^{\top} &= [I(\boldsymbol{\beta}^*)]_{\alpha, \alpha} - 2 \cdot (\mathbf{w}^*)^{\top} \cdot [I(\boldsymbol{\beta}^*)]_{\gamma, \alpha} + (\mathbf{w}^*)^{\top} \cdot [I(\boldsymbol{\beta}^*)]_{\gamma, \gamma} \cdot \mathbf{w}^* \\ &= [I(\boldsymbol{\beta}^*)]_{\alpha, \alpha} - [I(\boldsymbol{\beta}^*)]_{\gamma, \alpha} \cdot [I(\boldsymbol{\beta}^*)]_{\gamma, \gamma}^{-1} \cdot [I(\boldsymbol{\beta}^*)]_{\gamma, \alpha} \\ &= [I(\boldsymbol{\beta}^*)]_{\alpha|\gamma}, \end{aligned}$$

where the second and third equalities are from (4.1) and (4.2). Hence, by the central limit theorem we obtain (I.30). Finally, combining (I.29) and (I.30) by invoking Slutsky's theorem, we obtain

$$\sqrt{n} \cdot S_n(\widehat{\boldsymbol{\beta}}_0, \lambda) \xrightarrow{D} N(0, [I(\boldsymbol{\beta}^*)]_{\alpha|\gamma}),$$

which concludes the proof of Lemma G.3. \square

I.4 Proof of Lemma G.4

Proof. Throughout the proof, we abbreviate $w(\widehat{\boldsymbol{\beta}}_0, \lambda)$ as $\widehat{\mathbf{w}}_0$. Our proof is under the null hypothesis where $\boldsymbol{\beta}^* = [\alpha^*, (\boldsymbol{\gamma}^*)^{\top}]^{\top}$ with $\alpha^* = 0$. Recall that \mathbf{w}^* is defined in (4.1). Then by the definitions of $[T_n(\widehat{\boldsymbol{\beta}}_0)]_{\alpha|\gamma}$ and $[I(\boldsymbol{\beta}^*)]_{\alpha|\gamma}$ in (2.7) and (4.2), we have

$$\begin{aligned} [T_n(\widehat{\boldsymbol{\beta}}_0)]_{\alpha|\gamma} &= (1, -\widehat{\mathbf{w}}_0^{\top}) \cdot T_n(\widehat{\boldsymbol{\beta}}_0) \cdot (1, -\widehat{\mathbf{w}}_0^{\top})^{\top} \\ &= [T_n(\widehat{\boldsymbol{\beta}}_0)]_{\alpha, \alpha} - 2 \cdot \widehat{\mathbf{w}}_0^{\top} \cdot [T_n(\widehat{\boldsymbol{\beta}}_0)]_{\gamma, \alpha} + \widehat{\mathbf{w}}_0^{\top} \cdot [T_n(\widehat{\boldsymbol{\beta}}_0)]_{\gamma, \gamma} \cdot \widehat{\mathbf{w}}_0, \\ [I(\boldsymbol{\beta}^*)]_{\alpha|\gamma} &= [I(\boldsymbol{\beta}^*)]_{\alpha, \alpha} - [I(\boldsymbol{\beta}^*)]_{\gamma, \alpha}^{\top} \cdot [I(\boldsymbol{\beta}^*)]_{\gamma, \gamma}^{-1} \cdot [I(\boldsymbol{\beta}^*)]_{\gamma, \alpha} \\ &= [I(\boldsymbol{\beta}^*)]_{\alpha, \alpha} - 2 \cdot (\mathbf{w}^*)^{\top} \cdot [I(\boldsymbol{\beta}^*)]_{\gamma, \alpha} + (\mathbf{w}^*)^{\top} \cdot [I(\boldsymbol{\beta}^*)]_{\gamma, \gamma} \cdot \mathbf{w}^*. \end{aligned}$$

By triangle inequality, we have

$$\begin{aligned} &\left| [T_n(\widehat{\boldsymbol{\beta}}_0)]_{\alpha|\gamma} + [I(\boldsymbol{\beta}^*)]_{\alpha|\gamma} \right| \\ &\leq \underbrace{\left| [T_n(\widehat{\boldsymbol{\beta}}_0)]_{\alpha, \alpha} + [I(\boldsymbol{\beta}^*)]_{\alpha, \alpha} \right|}_{(i)} + 2 \cdot \underbrace{\left| \widehat{\mathbf{w}}_0^{\top} \cdot [T_n(\widehat{\boldsymbol{\beta}}_0)]_{\gamma, \alpha} + (\mathbf{w}^*)^{\top} \cdot [I(\boldsymbol{\beta}^*)]_{\gamma, \alpha} \right|}_{(ii)} \\ &\quad + \underbrace{\left| \widehat{\mathbf{w}}_0^{\top} \cdot [T_n(\widehat{\boldsymbol{\beta}}_0)]_{\gamma, \gamma} \cdot \widehat{\mathbf{w}}_0 + (\mathbf{w}^*)^{\top} \cdot [I(\boldsymbol{\beta}^*)]_{\gamma, \gamma} \cdot \mathbf{w}^* \right|}_{(iii)}. \end{aligned} \tag{I.42}$$

Analysis of Term (i): For term (i) in (I.42), by Theorem D.1 and triangle inequality we have

$$\left| [T_n(\widehat{\beta}_0)]_{\alpha,\alpha} + [I(\beta^*)]_{\alpha,\alpha} \right| \leq \underbrace{\left| [T_n(\widehat{\beta}_0)]_{\alpha,\alpha} - [T_n(\beta^*)]_{\alpha,\alpha} \right|}_{(i).a} + \underbrace{\left| [T_n(\beta^*)]_{\alpha,\alpha} - \{\mathbb{E}_{\beta^*} [T_n(\beta^*)]\}_{\alpha,\alpha} \right|}_{(i).b}. \quad (\text{I.43})$$

For term (i).a in (I.43), by Condition 4.4 we have

$$\begin{aligned} \left| [T_n(\widehat{\beta}_0)]_{\alpha,\alpha} - [T_n(\beta^*)]_{\alpha,\alpha} \right| &\leq \left\| T_n(\widehat{\beta}_0) - T_n(\beta^*) \right\|_{\infty,\infty} \\ &= O_{\mathbb{P}}(\zeta^L) \cdot \|\widehat{\beta}_0 - \beta^*\|_1. \end{aligned} \quad (\text{I.44})$$

Note that we have $(\widehat{\beta}_0)_{\alpha} = (\beta^*)_{\alpha} = 0$ by definition, which implies $\|\widehat{\beta}_0 - \beta^*\|_1 \leq \|\widehat{\beta} - \beta^*\|_1$. Hence, by Condition 4.1 we have

$$\|\widehat{\beta}_0 - \beta^*\|_1 = O_{\mathbb{P}}(\zeta^{\text{EM}}). \quad (\text{I.45})$$

Moreover, combining (I.44) and (I.45), by (4.6) in Assumption 4.5 we have

$$\zeta^L \cdot \zeta^{\text{EM}} \leq \max\{\|\mathbf{w}^*\|_1, 1\} \cdot s_{\mathbf{w}}^* \cdot \lambda = o(1)$$

for λ specified in (4.5). Hence we obtain

$$\begin{aligned} \left| [T_n(\widehat{\beta}_0)]_{\alpha,\alpha} - [T_n(\beta^*)]_{\alpha,\alpha} \right| &\leq \left\| T_n(\widehat{\beta}_0) - T_n(\beta^*) \right\|_{\infty,\infty} \\ &= O_{\mathbb{P}}(\zeta^L \cdot \zeta^{\text{EM}}) = o_{\mathbb{P}}(1). \end{aligned} \quad (\text{I.46})$$

Meanwhile, for term (i).b in (I.43) we have

$$\left| [T_n(\beta^*)]_{\alpha,\alpha} - \{\mathbb{E}_{\beta^*} [T_n(\beta^*)]\}_{\alpha,\alpha} \right| \leq \left\| T_n(\beta^*) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\|_{\infty,\infty} = O_{\mathbb{P}}(\zeta^{\text{T}}) = o_{\mathbb{P}}(1), \quad (\text{I.47})$$

where the second last equality follows from Condition 4.3, while the last equality holds because our assumption in (4.6) of Assumption 4.5 implies

$$\zeta^{\text{T}} \leq \max\{\|\mathbf{w}^*\|_1, 1\} \cdot s_{\mathbf{w}}^* \cdot \lambda = o(1)$$

for λ specified in (4.5).

Analysis of Term (ii): For term (ii) in (I.42), by Theorem D.1 and triangle inequality, we have

$$\begin{aligned} &\left| \widehat{\mathbf{w}}_0^{\text{T}} \cdot [T_n(\widehat{\beta}_0)]_{\gamma,\alpha} + (\mathbf{w}^*)^{\text{T}} \cdot [I(\beta^*)]_{\gamma,\alpha} \right| \quad (\text{I.48}) \\ &\leq \underbrace{\left| (\widehat{\mathbf{w}}_0 - \mathbf{w}^*)^{\text{T}} \cdot \left\{ T_n(\widehat{\beta}_0) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\}_{\gamma,\alpha} \right|}_{(ii).a} + \underbrace{\left| (\widehat{\mathbf{w}}_0 - \mathbf{w}^*)^{\text{T}} \cdot \left\{ \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\}_{\gamma,\alpha} \right|}_{(ii).b} \\ &\quad + \underbrace{\left| (\mathbf{w}^*)^{\text{T}} \cdot \left\{ T_n(\widehat{\beta}_0) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\}_{\gamma,\alpha} \right|}_{(ii).c}. \end{aligned}$$

By Hölder's inequality, term (ii).a in (I.48) is upper bounded by

$$\|\widehat{\mathbf{w}}_0 - \mathbf{w}^*\|_1 \cdot \left\| \left\{ T_n(\widehat{\beta}_0) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\}_{\gamma,\alpha} \right\|_{\infty} \leq \|\widehat{\mathbf{w}}_0 - \mathbf{w}^*\|_1 \cdot \left\| T_n(\widehat{\beta}_0) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\|_{\infty,\infty}.$$

By Lemma I.3, we have $\|\widehat{\mathbf{w}}_0 - \mathbf{w}^*\|_1 = O_{\mathbb{P}}(s_{\mathbf{w}}^* \cdot \lambda)$. Meanwhile, we have

$$\left\| T_n(\widehat{\beta}_0) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\|_{\infty,\infty} \leq \left\| T_n(\widehat{\beta}_0) - T_n(\beta^*) \right\|_{\infty,\infty} + \left\| T_n(\beta^*) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\|_{\infty,\infty} = o_{\mathbb{P}}(1).$$

where the second equality follows from (I.46) and (I.47). Therefore, term (ii).a is $o_{\mathbb{P}}(1)$, since (4.6) in Assumption 4.5 implies $s_{\mathbf{w}}^* \cdot \lambda = o(1)$. Meanwhile, by Hölder's inequality, term (ii).b in (I.48) is upper bounded by

$$\|\widehat{\mathbf{w}}_0 - \mathbf{w}^*\|_1 \cdot \left\| \left\{ \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\}_{\gamma,\alpha} \right\|_{\infty} \leq \|\widehat{\mathbf{w}}_0 - \mathbf{w}^*\|_1 \cdot \left\| \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\|_{\infty,\infty}. \quad (\text{I.49})$$

By Lemma I.3, we have $\|\widehat{\mathbf{w}}_0 - \mathbf{w}^*\|_1 = O_{\mathbb{P}}(s_{\mathbf{w}}^* \cdot \lambda)$. Meanwhile, we have $\mathbb{E}_{\beta^*} [T_n(\beta^*)] = -I(\beta^*)$ by Theorem D.1. Furthermore, (4.4) in Assumption 4.5 implies

$$\left\| I(\beta^*) \right\|_{\infty,\infty} \leq \left\| I(\beta^*) \right\|_2 \leq C, \quad (\text{I.50})$$

where $C > 0$ is some constant. Therefore, from (I.49) we have that term (ii).b in (I.48) is $O_{\mathbb{P}}(s_{\mathbf{w}}^* \cdot \lambda)$. By (4.6) in Assumption 4.5, we have $s_{\mathbf{w}}^* \cdot \lambda = o(1)$. Thus, we conclude that term (ii).b is $o_{\mathbb{P}}(1)$. For

term (ii).c, we have

$$\begin{aligned}
& \left| (\mathbf{w}^*)^\top \cdot \left\{ T_n(\widehat{\beta}_0) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\}_{\gamma, \alpha} \right| \\
& \leq \|\mathbf{w}^*\|_1 \cdot \left\| T_n(\widehat{\beta}_0) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\|_{\infty, \infty} \\
& \leq \|\mathbf{w}^*\|_1 \cdot \left\| T_n(\widehat{\beta}_0) - T_n(\beta^*) \right\|_{\infty, \infty} + \|\mathbf{w}^*\|_1 \cdot \left\| T_n(\beta^*) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\|_{\infty, \infty} \\
& = O_{\mathbb{P}}(\|\mathbf{w}^*\|_1 \cdot \zeta^L \cdot \zeta^{\text{EM}}) + O_{\mathbb{P}}(\|\mathbf{w}^*\|_1 \cdot \zeta^T) = o_{\mathbb{P}}(1).
\end{aligned}$$

Here the first and second inequalities are from Hölder's inequality and triangle inequality, the first equality follows from (I.46) and (I.47), and the second equality holds because (4.6) in Assumption 4.5 implies

$$\|\mathbf{w}^*\|_1 \cdot (\zeta^L \cdot \zeta^{\text{EM}} + \zeta^T) \leq \max\{\|\mathbf{w}^*\|_1, 1\} \cdot s_{\mathbf{w}}^* \cdot \lambda = o(1)$$

for λ specified in (4.5).

Analysis of Term (iii): For term (iii) in (I.42), by (D.1) in Theorem D.1 we have

$$\begin{aligned}
& \left| \widehat{\mathbf{w}}_0^\top \cdot [T_n(\widehat{\beta}_0)]_{\gamma, \gamma} \cdot \widehat{\mathbf{w}}_0 + (\mathbf{w}^*)^\top \cdot [I(\beta^*)]_{\gamma, \gamma} \cdot \mathbf{w}^* \right| \tag{I.51} \\
& \leq \underbrace{\left| \widehat{\mathbf{w}}_0^\top \cdot \left\{ T_n(\widehat{\beta}_0) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\}_{\gamma, \gamma} \cdot \widehat{\mathbf{w}}_0 \right|}_{\text{(iii).a}} + \underbrace{\left| \widehat{\mathbf{w}}_0^\top \cdot [I(\beta^*)]_{\gamma, \gamma} \cdot \widehat{\mathbf{w}}_0 - (\mathbf{w}^*)^\top \cdot [I(\beta^*)]_{\gamma, \gamma} \cdot \mathbf{w}^* \right|}_{\text{(iii).b}}.
\end{aligned}$$

For term (iii).a in (I.51), we have

$$\begin{aligned}
\left| \widehat{\mathbf{w}}_0^\top \cdot \left\{ T_n(\widehat{\beta}_0) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\}_{\gamma, \gamma} \cdot \widehat{\mathbf{w}}_0 \right| & \leq \|\widehat{\mathbf{w}}_0\|_1^2 \cdot \left\| \left\{ T_n(\widehat{\beta}_0) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\}_{\gamma, \gamma} \right\|_{\infty, \infty} \\
& \leq \|\widehat{\mathbf{w}}_0\|_1^2 \cdot \left\| T_n(\widehat{\beta}_0) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\|_{\infty, \infty}. \tag{I.52}
\end{aligned}$$

For $\|\widehat{\mathbf{w}}_0\|_1$ we have $\|\widehat{\mathbf{w}}_0\|_1^2 \leq (\|\mathbf{w}^*\|_1 + \|\widehat{\mathbf{w}}_0 - \mathbf{w}^*\|_1)^2 = [\|\mathbf{w}^*\|_1 + O_{\mathbb{P}}(s_{\mathbf{w}}^* \cdot \lambda)]^2$, where the equality holds because by Lemma I.3 we have $\|\widehat{\mathbf{w}}_0 - \mathbf{w}^*\|_1 = O_{\mathbb{P}}(s_{\mathbf{w}}^* \cdot \lambda)$. Meanwhile, on the right-hand side of (I.52) we have

$$\begin{aligned}
\left\| T_n(\widehat{\beta}_0) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\|_{\infty, \infty} & \leq \left\| T_n(\widehat{\beta}_0) - T_n(\beta^*) \right\|_{\infty, \infty} + \left\| T_n(\beta^*) - \mathbb{E}_{\beta^*} [T_n(\beta^*)] \right\|_{\infty, \infty} \\
& = O_{\mathbb{P}}(\zeta^L \cdot \zeta^{\text{EM}} + \zeta^T).
\end{aligned}$$

Here the last equality is from (I.46) and (I.47). Hence, term (iii).a in (I.51) is $O_{\mathbb{P}}[(\|\mathbf{w}^*\|_1 + s_{\mathbf{w}}^* \cdot \lambda)^2 \cdot (\zeta^L \cdot \zeta^{\text{EM}} + \zeta^T)]$. Note that

$$\begin{aligned}
& (\|\mathbf{w}^*\|_1 + s_{\mathbf{w}}^* \cdot \lambda)^2 \cdot (\zeta^L \cdot \zeta^{\text{EM}} + \zeta^T) \\
& = \underbrace{\|\mathbf{w}^*\|_1^2 \cdot (\zeta^L \cdot \zeta^{\text{EM}} + \zeta^T)}_{\text{(i)}} + 2 \cdot \underbrace{s_{\mathbf{w}}^* \cdot \lambda \cdot \|\mathbf{w}^*\|_1 \cdot (\zeta^L \cdot \zeta^{\text{EM}} + \zeta^T)}_{\text{(ii)}} + \underbrace{(s_{\mathbf{w}}^* \cdot \lambda)^2 \cdot (\zeta^L \cdot \zeta^{\text{EM}} + \zeta^T)}_{\text{(iv)}}.
\end{aligned}$$

From (4.6) in Assumption 4.5 we have, for λ specified in (4.5), terms (i)-(iv) are all upper bounded by $\max\{\|\mathbf{w}^*\|_1, 1\} \cdot s_{\mathbf{w}}^* \cdot \lambda = o(1)$. Hence, we conclude term (iii).a in (I.51) is $o_{\mathbb{P}}(1)$. Also, for term (iii).b in (I.51), we have

$$\begin{aligned}
& \left| \widehat{\mathbf{w}}_0^\top \cdot [I(\beta^*)]_{\gamma, \gamma} \cdot \widehat{\mathbf{w}}_0 - (\mathbf{w}^*)^\top \cdot [I(\beta^*)]_{\gamma, \gamma} \cdot \mathbf{w}^* \right| \\
& \leq \left| (\widehat{\mathbf{w}}_0 - \mathbf{w}^*)^\top \cdot [I(\beta^*)]_{\gamma, \gamma} \cdot (\widehat{\mathbf{w}}_0 - \mathbf{w}^*) \right| + 2 \cdot \left| \mathbf{w}^* \cdot [I(\beta^*)]_{\gamma, \gamma} \cdot (\widehat{\mathbf{w}}_0 - \mathbf{w}^*) \right| \\
& \leq \|\widehat{\mathbf{w}}_0 - \mathbf{w}^*\|_1^2 \cdot \left\| [I(\beta^*)]_{\gamma, \gamma} \right\|_{\infty, \infty} + 2 \cdot \|\widehat{\mathbf{w}}_0 - \mathbf{w}^*\|_1 \cdot \|\mathbf{w}^*\|_1 \cdot \left\| [I(\beta^*)]_{\gamma, \gamma} \right\|_{\infty, \infty} \\
& = O_{\mathbb{P}}[(s_{\mathbf{w}}^* \cdot \lambda)^2 + \|\mathbf{w}^*\|_1 \cdot s_{\mathbf{w}}^* \cdot \lambda],
\end{aligned}$$

where the last equality follows from Lemma I.3 and (I.50). Moreover, by (4.6) in Assumption 4.5 we have $\max\{s_{\mathbf{w}}^* \cdot \lambda, \|\mathbf{w}^*\|_1 \cdot s_{\mathbf{w}}^* \cdot \lambda\} \leq \max\{\|\mathbf{w}^*\|_1, 1\} \cdot s_{\mathbf{w}}^* \cdot \lambda = o(1)$. Therefore, we conclude that term (iii).b in (I.51) is $o_{\mathbb{P}}(1)$. Combining the above analysis for terms (i)-(iii) in (I.42), we obtain

$$\left| [T_n(\widehat{\beta}_0)]_{\alpha|\gamma} + [I(\beta^*)]_{\alpha|\gamma} \right| = o_{\mathbb{P}}(1).$$

Thus we conclude the proof of Lemma G.4. \square

I.5 Proof of Lemma F.1

Proof. According to Algorithm 4, the final estimator $\widehat{\beta} = \beta^{(T)}$ has \widehat{s} nonzero entries. Meanwhile, we have $\|\beta^*\|_0 = s^* \leq \widehat{s}$. Hence, we have $\|\widehat{\beta} - \beta^*\|_1 \leq 2 \cdot \sqrt{\widehat{s}} \cdot \|\widehat{\beta} - \beta^*\|_2$. Invoking (E.6) in Theorem E.3, we obtain ζ^{EM} .

For Gaussian mixture model, the maximization implementation of the M-step takes the form

$$M_n(\beta) = \frac{2}{n} \sum_{i=1}^n \omega_\beta(\mathbf{y}_i) \cdot \mathbf{y}_i - \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i, \quad \text{and} \quad M(\beta) = \mathbb{E}[2 \cdot \omega_\beta(\mathbf{Y}) \cdot \mathbf{Y} - \mathbf{Y}],$$

where $\omega_\beta(\cdot)$ is defined in (A.2). Meanwhile, we have

$$\nabla_1 Q_n(\beta; \beta) = \frac{1}{n} \sum_{i=1}^n [2 \cdot \omega_\beta(\mathbf{y}_i) - 1] \cdot \mathbf{y}_i - \beta, \quad \text{and} \quad \nabla_1 Q(\beta; \beta) = \mathbb{E}[2 \cdot \omega_\beta(\mathbf{Y}) - \mathbf{Y}] - \beta.$$

Hence, we have $\|M_n(\beta) - M(\beta)\|_\infty = \|\nabla_1 Q_n(\beta; \beta) - \nabla_1 Q(\beta; \beta)\|_\infty$. By setting $\delta = 2/d$ in Lemma E.2, we obtain ζ^{G} . \square

I.6 Proof of Lemma F.2

Proof. Recall that for Gaussian mixture model we have

$$Q_n(\beta'; \beta) = -\frac{1}{2n} \sum_{i=1}^n \left\{ \omega_\beta(\mathbf{y}_i) \cdot \|\mathbf{y}_i - \beta'\|_2^2 + [1 - \omega_\beta(\mathbf{y}_i)] \cdot \|\mathbf{y}_i + \beta'\|_2^2 \right\},$$

where $\omega_\beta(\cdot)$ is defined in (A.2). Hence, by calculation we have

$$\nabla_1 Q_n(\beta'; \beta) = \frac{1}{n} \sum_{i=1}^n [2 \cdot \omega_\beta(\mathbf{y}_i) - 1] \cdot \mathbf{y}_i - \beta', \quad \nabla_{1,1}^2 Q_n(\beta'; \beta) = -\mathbf{I}_d, \quad (\text{I.53})$$

$$\nabla_{1,2}^2 Q_n(\beta'; \beta) = \frac{4}{n} \sum_{i=1}^n \frac{\mathbf{y}_i \cdot \mathbf{y}_i^\top}{\sigma^2 \cdot [1 + \exp(-2 \cdot \langle \beta, \mathbf{y}_i \rangle / \sigma^2)] \cdot [1 + \exp(2 \cdot \langle \beta, \mathbf{y}_i \rangle / \sigma^2)]}. \quad (\text{I.54})$$

For notational simplicity we define

$$\nu_\beta(\mathbf{y}) = \frac{4}{\sigma^2 \cdot [1 + \exp(-2 \cdot \langle \beta, \mathbf{y} \rangle / \sigma^2)] \cdot [1 + \exp(2 \cdot \langle \beta, \mathbf{y} \rangle / \sigma^2)]}. \quad (\text{I.55})$$

Then by the definition of $T_n(\cdot)$ in (2.4), from (I.53) and (I.54) we have

$$\{T_n(\beta^*) - \mathbb{E}_{\beta^*}[T_n(\beta^*)]\}_{j,k} = \frac{1}{n} \sum_{i=1}^n \nu_{\beta^*}(\mathbf{y}_i) \cdot y_{i,j} \cdot y_{i,k} - \mathbb{E}_{\beta^*}[\nu_{\beta^*}(\mathbf{Y}) \cdot Y_j \cdot Y_k].$$

Applying the symmetrization result in Lemma J.4 to the right-hand side, we have that for $u > 0$,

$$\mathbb{E}_{\beta^*} \left\{ \exp \left[u \cdot \left| \{T_n(\beta^*) - \mathbb{E}_{\beta^*}[T_n(\beta^*)]\}_{j,k} \right| \right] \right\} \leq \mathbb{E}_{\beta^*} \left\{ \exp \left[u \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot \nu_{\beta^*}(\mathbf{y}_i) \cdot y_{i,j} \cdot y_{i,k} \right| \right] \right\}, \quad (\text{I.56})$$

where ξ_1, \dots, ξ_n are i.i.d. Rademacher random variables that are independent of $\mathbf{y}_1, \dots, \mathbf{y}_n$. Then we invoke the contraction result in Lemma J.5 by setting

$$f(y_{i,j} \cdot y_{i,k}) = y_{i,j} \cdot y_{i,k}, \quad \mathcal{F} = \{f\}, \quad \psi_i(v) = \nu_{\beta^*}(\mathbf{y}_i) \cdot v, \quad \text{and} \quad \phi(v) = \exp(u \cdot v),$$

where u is the variable of the moment generating function in (I.56). By the definition in (I.55) we have $|\nu_{\beta^*}(\mathbf{y}_i)| \leq 4/\sigma^2$, which implies

$$|\psi_i(v) - \psi_i(v')| \leq |\nu_{\beta^*}(\mathbf{y}_i) \cdot (v - v')| \leq 4/\sigma^2 \cdot |v - v'|, \quad \text{for all } v, v' \in \mathbb{R}.$$

Therefore, applying the contraction result in Lemma J.5 to the right-hand side of (I.56), we obtain

$$\mathbb{E}_{\beta^*} \left\{ \exp \left[u \cdot \left| \{T_n(\beta^*) - \mathbb{E}_{\beta^*}[T_n(\beta^*)]\}_{j,k} \right| \right] \right\} \leq \mathbb{E}_{\beta^*} \left\{ \exp \left[u \cdot 4/\sigma^2 \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_{i,j} \cdot y_{i,k} \right| \right] \right\}. \quad (\text{I.57})$$

Note that $\mathbb{E}_{\beta^*}(\xi_i \cdot y_{i,j} \cdot y_{i,k}) = 0$, since ξ_i is a Rademacher random variable independent of $y_{i,j} \cdot y_{i,k}$. Recall that in Gaussian mixture model we have $y_{i,j} = z_i \cdot \beta_j^* + v_{i,j}$, where z_i is a Rademacher

random variable and $v_{i,j} \sim N(0, \sigma^2)$. Hence, by Example 5.8 in [28], we have $\|z_i \cdot \beta_j^*\|_{\psi_2} \leq |\beta_j^*|$ and $\|v_{i,j}\|_{\psi_2} \leq C \cdot \sigma$. Therefore, by Lemma J.1 we have

$$\|y_{i,j}\|_{\psi_2} = \|z_i \cdot \beta_j^* + v_{i,j}\|_{\psi_2} \leq C' \cdot \sqrt{|\beta_j^*|^2 + C'' \cdot \sigma^2} \leq C' \cdot \sqrt{\|\beta^*\|_\infty^2 + C'' \cdot \sigma^2}. \quad (\text{I.58})$$

Since $|\xi_i \cdot y_{i,j} \cdot y_{i,k}| = |y_{i,j} \cdot y_{i,k}|$, by definition $\xi_i \cdot y_{i,j} \cdot y_{i,k}$ and $y_{i,j} \cdot y_{i,k}$ have the same ψ_1 -norm. By Lemma J.2 we have

$$\|\xi_i \cdot y_{i,j} \cdot y_{i,k}\|_{\psi_1} \leq C \cdot \max\{\|y_{i,j}\|_{\psi_2}^2, \|y_{i,k}\|_{\psi_2}^2\} \leq C' \cdot (\|\beta^*\|_\infty^2 + C'' \cdot \sigma^2).$$

Then by Lemma 5.15 in [28], we obtain

$$\mathbb{E}_{\beta^*} [\exp(u' \cdot \xi_i \cdot y_{i,j} \cdot y_{i,k})] \leq \exp\left[(u')^2 \cdot C \cdot (\|\beta^*\|_\infty^2 + C'' \cdot \sigma^2)\right] \quad (\text{I.59})$$

for all $|u'| \leq C'' / (\|\beta^*\|_\infty^2 + C'' \cdot \sigma^2)$. Note that on the right-hand side of (I.57), we have

$$\begin{aligned} & \mathbb{E}_{\beta^*} \left\{ \exp \left[u \cdot 4/\sigma^2 \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_{i,j} \cdot y_{i,k} \right| \right] \right\} \\ & \leq \mathbb{E}_{\beta^*} \left(\max \left\{ \exp \left[u \cdot 4/\sigma^2 \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_{i,j} \cdot y_{i,k} \right], \exp \left[-u \cdot 4/\sigma^2 \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_{i,j} \cdot y_{i,k} \right] \right\} \right) \\ & \leq \mathbb{E}_{\beta^*} \left\{ \exp \left[u \cdot 4/\sigma^2 \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_{i,j} \cdot y_{i,k} \right] \right\} + \mathbb{E}_{\beta^*} \left\{ \exp \left[-u \cdot 4/\sigma^2 \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \cdot y_{i,j} \cdot y_{i,k} \right] \right\}. \end{aligned} \quad (\text{I.60})$$

By plugging (I.59) into the right-hand side of (I.60) with $u' = u \cdot 4/(\sigma^2 \cdot n)$ and $u' = -u \cdot 4/(\sigma^2 \cdot n)$, from (I.57) we have that for any $j, k \in \{1, \dots, d\}$,

$$\begin{aligned} & \mathbb{E}_{\beta^*} \left\{ \exp \left[u \cdot \left| \{T_n(\beta^*) - \mathbb{E}_{\beta^*} [T_n(\beta^*)]\}_{j,k} \right| \right] \right\} \\ & \leq 2 \cdot \exp \left[C \cdot u^2/n \cdot (\|\beta^*\|_\infty^2 + C'' \cdot \sigma^2)^2/\sigma^4 \right]. \end{aligned} \quad (\text{I.61})$$

Therefore, by Chernoff bound we have that, for all $v > 0$ and $|u| \leq C'' / (\|\beta^*\|_\infty^2 + C'' \cdot \sigma^2)$,

$$\begin{aligned} & \mathbb{P} \left[\|T_n(\beta^*) - \mathbb{E}_{\beta^*} [T_n(\beta^*)]\|_{\infty, \infty} > v \right] \\ & \leq \mathbb{E}_{\beta^*} \left\{ \exp \left[u \cdot \|T_n(\beta^*) - \mathbb{E}_{\beta^*} [T_n(\beta^*)]\|_{\infty, \infty} \right] \right\} / \exp(u \cdot v) \\ & \leq \sum_{j=1}^d \sum_{k=1}^d \mathbb{E}_{\beta^*} \left\{ \exp \left[u \cdot \left| \{T_n(\beta^*) - \mathbb{E}_{\beta^*} [T_n(\beta^*)]\}_{j,k} \right| \right] \right\} / \exp(u \cdot v) \\ & \leq 2 \cdot \exp \left[C \cdot u^2/n \cdot (\|\beta^*\|_\infty^2 + C'' \cdot \sigma^2)^2/\sigma^4 - u \cdot v + 2 \cdot \log d \right], \end{aligned}$$

where the last inequality is obtained from (I.61). By minimizing its right-hand side over u , we conclude that for $0 < v \leq C'' \cdot (\|\beta^*\|_\infty^2 + C'' \cdot \sigma^2)$,

$$\mathbb{P} \left[\|T_n(\beta^*) - \mathbb{E}_{\beta^*} [T_n(\beta^*)]\|_{\infty, \infty} > v \right] \leq 2 \cdot \exp \left\{ -n \cdot v^2 / \left[C \cdot (\|\beta^*\|_\infty^2 + C'' \cdot \sigma^2)^2/\sigma^4 \right] + 2 \cdot \log d \right\}.$$

Setting the right-hand side to be δ , we have that

$$\|T_n(\beta^*) - \mathbb{E}_{\beta^*} [T_n(\beta^*)]\|_{\infty, \infty} \leq v = C \cdot (\|\beta^*\|_\infty^2 + C'' \cdot \sigma^2)/\sigma^2 \cdot \sqrt{\frac{\log(2/\delta) + 2 \cdot \log d}{n}}$$

holds with probability at least $1 - \delta$. By setting $\delta = 2/d$, we conclude the proof of Lemma F.2. \square

I.7 Proof of Lemma F.3

Proof. For any $j, k \in \{1, \dots, d\}$, by the mean-value theorem we have

$$\begin{aligned} \|T_n(\boldsymbol{\beta}) - T_n(\boldsymbol{\beta}^*)\|_{\infty, \infty} &= \max_{j, k \in \{1, \dots, d\}} \left| [T_n(\boldsymbol{\beta})]_{j, k} - [T_n(\boldsymbol{\beta}^*)]_{j, k} \right| \\ &= \max_{j, k \in \{1, \dots, d\}} \left| (\boldsymbol{\beta} - \boldsymbol{\beta}^*)^\top \cdot \nabla [T_n(\boldsymbol{\beta}^\#)]_{j, k} \right| \\ &\leq \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_1 \cdot \max_{j, k \in \{1, \dots, d\}} \left\| \nabla [T_n(\boldsymbol{\beta}^\#)]_{j, k} \right\|_{\infty}, \end{aligned} \quad (\text{I.62})$$

where $\boldsymbol{\beta}^\#$ is an intermediate value between $\boldsymbol{\beta}$ and $\boldsymbol{\beta}^*$. According to (I.53), (I.54) and the definition of $T_n(\cdot)$ in (2.4), by calculation we have

$$\nabla [T_n(\boldsymbol{\beta}^\#)]_{j, k} = \frac{1}{n} \sum_{i=1}^n \bar{\nu}_{\boldsymbol{\beta}^\#}(\mathbf{y}_i) \cdot y_{i, j} \cdot y_{i, k} \cdot \mathbf{y}_i,$$

where

$$\begin{aligned} \bar{\nu}_{\boldsymbol{\beta}^\#}(\mathbf{y}) &= \frac{8/\sigma^4}{\left[1 + \exp(-2 \cdot \langle \boldsymbol{\beta}, \mathbf{y} \rangle / \sigma^2)\right] \cdot \left[1 + \exp(2 \cdot \langle \boldsymbol{\beta}, \mathbf{y} \rangle / \sigma^2)\right]^2} \\ &\quad - \frac{8/\sigma^4}{\left[1 + \exp(-2 \cdot \langle \boldsymbol{\beta}, \mathbf{y} \rangle / \sigma^2)\right]^2 \cdot \left[1 + \exp(2 \cdot \langle \boldsymbol{\beta}, \mathbf{y} \rangle / \sigma^2)\right]}. \end{aligned} \quad (\text{I.63})$$

For notational simplicity, we define the following event

$$\mathcal{E} = \{\|\mathbf{y}_i\|_{\infty} \leq \tau, \text{ for all } i = 1, \dots, n\},$$

where $\tau > 0$ will be specified later. By maximal inequality we have

$$\begin{aligned} \mathbb{P}\left\{\left\|\nabla [T_n(\boldsymbol{\beta}^\#)]_{j, k}\right\|_{\infty} > v\right\} &\leq d \cdot \mathbb{P}\left\{\left|\nabla [T_n(\boldsymbol{\beta}^\#)]_{j, k}\right|_l > v\right\} \\ &= d \cdot \mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{\nu}_{\boldsymbol{\beta}^\#}(\mathbf{y}_i) \cdot y_{i, j} \cdot y_{i, k} \cdot y_{i, l}\right| > v\right]. \end{aligned} \quad (\text{I.64})$$

Let $\bar{\mathcal{E}}$ be the complement of \mathcal{E} . On the right-hand side of (I.64) we have

$$\mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{\nu}_{\boldsymbol{\beta}^\#}(\mathbf{y}_i) \cdot y_{i, j} \cdot y_{i, k} \cdot y_{i, l}\right| > v\right] = \underbrace{\mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{\nu}_{\boldsymbol{\beta}^\#}(\mathbf{y}_i) \cdot y_{i, j} \cdot y_{i, k} \cdot y_{i, l}\right| > v, \mathcal{E}\right]}_{(i)} + \underbrace{\mathbb{P}(\bar{\mathcal{E}})}_{(ii)}. \quad (\text{I.65})$$

Analysis of Term (i): For term (i) in (I.65), we have

$$\begin{aligned} &\mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{\nu}_{\boldsymbol{\beta}^\#}(\mathbf{y}_i) \cdot y_{i, j} \cdot y_{i, k} \cdot y_{i, l}\right| > v, \mathcal{E}\right] \\ &= \mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{\nu}_{\boldsymbol{\beta}^\#}(\mathbf{y}_i) \cdot y_{i, j} \cdot y_{i, k} \cdot y_{i, l} \cdot \mathbf{1}\{\|\mathbf{y}_i\|_{\infty} \leq \tau\}\right| > v, \mathcal{E}\right] \\ &\leq \mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{\nu}_{\boldsymbol{\beta}^\#}(\mathbf{y}_i) \cdot y_{i, j} \cdot y_{i, k} \cdot y_{i, l} \cdot \mathbf{1}\{\|\mathbf{y}_i\|_{\infty} \leq \tau\}\right| > v\right] \\ &\leq \sum_{i=1}^n \mathbb{P}\left[\left|\bar{\nu}_{\boldsymbol{\beta}^\#}(\mathbf{y}_i) \cdot y_{i, j} \cdot y_{i, k} \cdot y_{i, l} \cdot \mathbf{1}\{\|\mathbf{y}_i\|_{\infty} \leq \tau\}\right| > v\right], \end{aligned}$$

where the last inequality is from union bound. By (I.63) we have $|\bar{\nu}_{\boldsymbol{\beta}^\#}(\mathbf{y}_i)| \leq 16/\sigma^4$. Thus we obtain

$$\mathbb{P}\left[\left|\bar{\nu}_{\boldsymbol{\beta}^\#}(\mathbf{y}_i) \cdot y_{i, j} \cdot y_{i, k} \cdot y_{i, l} \cdot \mathbf{1}\{\|\mathbf{y}_i\|_{\infty} \leq \tau\}\right| > v\right] \leq \mathbb{P}\left[|y_{i, j} \cdot y_{i, k} \cdot y_{i, l} \cdot \mathbf{1}\{\|\mathbf{y}_i\|_{\infty} \leq \tau\}| > \sigma^4/16 \cdot v\right].$$

Taking $v = 16 \cdot \tau^3/\sigma^4$, we have that the right-hand side is zero and hence term (i) in (I.65) is zero.

Analysis of Term (ii): Meanwhile, for term (ii) in (I.65) by maximal inequality we have

$$\mathbb{P}(\bar{\mathcal{E}}) = \mathbb{P}\left(\max_{i \in \{1, \dots, n\}} \|\mathbf{y}_i\|_{\infty} > \tau\right) \leq n \cdot \mathbb{P}(\|\mathbf{y}_i\|_{\infty} > \tau) \leq n \cdot d \cdot \mathbb{P}(|y_{i, j}| > \tau).$$

Furthermore, by (I.58) in the proof of Lemma F.1, we have that $y_{i,j}$ is sub-Gaussian with $\|y_{i,j}\|_{\psi_2} = C \cdot \sqrt{\|\beta^*\|_\infty^2 + C' \cdot \sigma^2}$. Therefore, by Lemma 5.5 in [28] we have

$$\mathbb{P}(\bar{\mathcal{E}}) \leq n \cdot d \cdot \mathbb{P}(|y_{i,j}| > \tau) \leq n \cdot d \cdot 2 \cdot \exp[-C \cdot \tau^2 / (\|\beta^*\|_\infty^2 + C' \cdot \sigma^2)].$$

To ensure the right-hand side is upper bounded by δ , we set τ to be

$$\tau = C \cdot \sqrt{\|\beta^*\|_\infty^2 + C' \cdot \sigma^2} \cdot \sqrt{\log d + \log n + \log(2/\delta)}. \quad (\text{I.66})$$

Finally, by (I.64), (I.65) and maximal inequality we have

$$\mathbb{P}\left\{\max_{j,k \in \{1, \dots, d\}} \left\| \nabla [T_n(\beta^\#)]_{j,k} \right\|_\infty > v\right\} \leq d^2 \cdot d \cdot \mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{\nu}_{\beta^\#}(\mathbf{y}_i) \cdot y_{i,j} \cdot y_{i,k} \cdot y_{i,l}\right| > v\right] \leq d^3 \cdot \delta$$

for $v = 16 \cdot \tau^3 / \sigma^4$ with τ specified in (I.66). By setting $\delta = 2 \cdot d^{-4}$ and plugging (I.66) into (I.62), we conclude the proof of Lemma F.3. \square

I.8 Proof of Lemma F.5

By the same proof of Lemma F.1 in §I.5, we obtain ζ^{EM} by invoking Theorem E.6. To obtain ζ^{G} , note that for the gradient implementation of the M-step (Algorithm 3), we have

$$M_n(\beta) = \beta + \eta \cdot \nabla_1 Q_n(\beta; \beta), \quad \text{and} \quad M(\beta) = \beta + \eta \cdot \nabla_1 Q(\beta; \beta).$$

Hence, we obtain $\|\nabla_1 Q_n(\beta^*; \beta^*) - \nabla_1 Q(\beta^*; \beta^*)\|_\infty = 1/\eta \cdot \|M_n(\beta^*) - M(\beta^*)\|_\infty$. Setting $\delta = 4/d$ and $s = s^*$ in (E.9) of Lemma E.5, we then obtain ζ^{G} .

I.9 Proof of Lemma F.6

Proof. Recall that for mixture of regression model, we have

$$Q_n(\beta'; \beta) = -\frac{1}{2n} \sum_{i=1}^n \left\{ \omega_\beta(\mathbf{x}_i, y_i) \cdot (y_i - \langle \mathbf{x}_i, \beta' \rangle)^2 + [1 - \omega_\beta(\mathbf{x}_i, y_i)] \cdot (y_i + \langle \mathbf{x}_i, \beta' \rangle)^2 \right\},$$

where $\omega_\beta(\cdot)$ is defined in (A.5). Hence, by calculation we have

$$\nabla_1 Q_n(\beta', \beta) = \frac{1}{n} \sum_{i=1}^n [2 \cdot \omega_\beta(\mathbf{x}_i, y_i) \cdot y_i \cdot \mathbf{x}_i - \mathbf{x}_i \cdot \mathbf{x}_i^\top \cdot \beta'], \quad \nabla_{1,1}^2 Q_n(\beta', \beta) = -\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{x}_i^\top, \quad (\text{I.67})$$

$$\nabla_{1,2}^2 Q_n(\beta', \beta) = \frac{4}{n} \sum_{i=1}^n \frac{y_i^2 \cdot \mathbf{x}_i \cdot \mathbf{x}_i^\top}{\sigma^2 \cdot [1 + \exp(-2 \cdot y_i \cdot \langle \beta, \mathbf{x}_i \rangle / \sigma^2)] \cdot [1 + \exp(2 \cdot y_i \cdot \langle \beta, \mathbf{x}_i \rangle / \sigma^2)]}. \quad (\text{I.68})$$

For notational simplicity we define

$$\nu_\beta(\mathbf{x}, y) = \frac{4}{\sigma^2 \cdot [1 + \exp(-2 \cdot y \cdot \langle \beta, \mathbf{x} \rangle / \sigma^2)] \cdot [1 + \exp(2 \cdot y \cdot \langle \beta, \mathbf{x} \rangle / \sigma^2)]}. \quad (\text{I.69})$$

Then by the definition of $T_n(\cdot)$ in (2.4), from (I.67) and (I.68) we have

$$\{T_n(\beta^*) - \mathbb{E}_{\beta^*} [T_n(\beta^*)]\}_{j,k} = \frac{1}{n} \sum_{i=1}^n \nu_{\beta^*}(\mathbf{x}_i, y_i) \cdot x_{i,j} \cdot x_{i,k} \cdot y_i^2 - \mathbb{E}_{\beta^*} [\nu_{\beta^*}(Y, \mathbf{X}) \cdot X_j \cdot X_k \cdot Y_i^2].$$

Applying the symmetrization result in Lemma J.4 to the right-hand side, we have that for $u > 0$,

$$\begin{aligned} & \mathbb{E}_{\beta^*} \left\{ \exp \left[u \cdot \left| \{T_n(\beta^*) - \mathbb{E}_{\beta^*} [T_n(\beta^*)]\}_{j,k} \right| \right] \right\} \\ & \leq \mathbb{E}_{\beta^*} \left\{ \exp \left[u \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot \nu_{\beta^*}(\mathbf{x}_i, y_i) \cdot x_{i,j} \cdot x_{i,k} \cdot y_i^2 \right| \right] \right\}, \end{aligned} \quad (\text{I.70})$$

where ξ_1, \dots, ξ_n are i.i.d. Rademacher random variables, which are independent of $\mathbf{x}_1, \dots, \mathbf{x}_n$ and y_1, \dots, y_n . Then we invoke the contraction result in Lemma J.5 by setting

$$f(x_{i,j} \cdot x_{i,k} \cdot y_i^2) = x_{i,j} \cdot x_{i,k} \cdot y_i^2, \quad \mathcal{F} = \{f\}, \quad \psi_i(v) = \nu_{\beta^*}(\mathbf{x}_i, y_i) \cdot v, \quad \text{and} \quad \phi(v) = \exp(u \cdot v),$$

where u is the variable of the moment generating function in (I.70). By the definition in (I.69) we have $|\nu_{\beta^*}(\mathbf{x}_i, y_i)| \leq 4/\sigma^2$, which implies

$$|\psi_i(v) - \psi_i(v')| \leq |\nu_{\beta^*}(\mathbf{x}_i, y_i) \cdot (v - v')| \leq 4/\sigma^2 \cdot |v - v'|, \quad \text{for all } v, v' \in \mathbb{R}.$$

Therefore, applying Lemma J.5 to the right-hand side of (I.70), we obtain

$$\begin{aligned} & \mathbb{E}_{\beta^*} \left\{ \exp \left[u \cdot \left| \{T_n(\beta^*) - \mathbb{E}_{\beta^*} [T_n(\beta^*)]\}_{j,k} \right| \right] \right\} \\ & \leq \mathbb{E}_{\beta^*} \left\{ \exp \left[u \cdot 4/\sigma^2 \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot x_{i,j} \cdot x_{i,k} \cdot y_i^2 \right| \right] \right\}. \end{aligned} \quad (\text{I.71})$$

For notational simplicity, we define the following event

$$\mathcal{E} = \{ \|\mathbf{x}_i\|_\infty \leq \tau, \text{ for all } i = 1, \dots, n \}.$$

Let $\bar{\mathcal{E}}$ be the complement of \mathcal{E} . We consider the following tail probability

$$\mathbb{P} \left[4/\sigma^2 \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot x_{i,j} \cdot x_{i,k} \cdot y_i^2 \right| > v \right] \leq \underbrace{\mathbb{P} \left[4/\sigma^2 \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot x_{i,j} \cdot x_{i,k} \cdot y_i^2 \right| > v, \mathcal{E} \right]}_{(i)} + \underbrace{\mathbb{P}(\bar{\mathcal{E}})}_{(ii)}. \quad (\text{I.72})$$

Analysis of Term (i): For term (i) in (I.72), we have

$$\begin{aligned} \mathbb{P} \left[4/\sigma^2 \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot x_{i,j} \cdot x_{i,k} \cdot y_i^2 \right| > v, \mathcal{E} \right] &= \mathbb{P} \left[4/\sigma^2 \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot x_{i,j} \cdot x_{i,k} \cdot y_i^2 \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\} \right| > v, \mathcal{E} \right] \\ &\leq \mathbb{P} \left[4/\sigma^2 \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot x_{i,j} \cdot x_{i,k} \cdot y_i^2 \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\} \right| > v \right]. \end{aligned}$$

Here note that $\mathbb{E}_{\beta^*} (\xi_i \cdot x_{i,j} \cdot x_{i,k} \cdot y_i^2 \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\}) = 0$, because ξ_i is a Rademacher random variable independent of \mathbf{x}_i and y_i . Recall that for mixture of regression model we have $y_i = z_i \cdot \langle \beta^*, \mathbf{x}_i \rangle + v_i$, where z_i is a Rademacher random variable, $\mathbf{x}_i \sim N(0, \mathbf{I}_d)$ and $v_i \sim N(0, \sigma^2)$. According to Example 5.8 in [28], we have $\|z_i \cdot \langle \beta^*, \mathbf{x}_i \rangle \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\}\|_{\psi_2} = \|\langle \beta^*, \mathbf{x}_i \rangle \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\}\|_{\psi_2} \leq \tau \cdot \|\beta^*\|_1$ and $\|v_i \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\}\|_{\psi_2} \leq \|v_i\|_{\psi_2} \leq C \cdot \sigma$. Hence, by Lemma J.1 we have

$$\begin{aligned} \|y_i \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\}\|_{\psi_2} &= \|z_i \cdot \langle \beta^*, \mathbf{x}_i \rangle \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\} + v_i \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\}\|_{\psi_2} \\ &\leq C \cdot \sqrt{\tau^2 \cdot \|\beta^*\|_1^2 + C' \cdot \sigma^2}. \end{aligned} \quad (\text{I.73})$$

By the definition of ψ_1 -norm, we have $\|\xi_i \cdot x_{i,j} \cdot x_{i,k} \cdot y_i^2 \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\}\|_{\psi_1} \leq \tau^2 \cdot \|y_i^2 \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\}\|_{\psi_1}$. Further by applying Lemma J.2 to its right-hand side with $Z_1 = Z_2 = y_i \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\}$, we obtain

$$\begin{aligned} \|\xi_i \cdot x_{i,j} \cdot x_{i,k} \cdot y_i^2 \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\}\|_{\psi_1} &\leq C \cdot \tau^2 \cdot \|y_i \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\}\|_{\psi_2}^2 \\ &\leq C' \cdot \tau^2 \cdot (\tau^2 \cdot \|\beta^*\|_1^2 + C'' \cdot \sigma^2), \end{aligned}$$

where the last inequality follows from (I.73). Therefore, by Bernstein's inequality (Proposition 5.16 in [28]), we obtain

$$\begin{aligned} & \mathbb{P} \left[4/\sigma^2 \cdot \left| \frac{1}{n} \sum_{i=1}^n \xi_i \cdot x_{i,j} \cdot x_{i,k} \cdot y_i^2 \cdot \mathbb{1}\{\|\mathbf{x}_i\|_\infty \leq \tau\} \right| > v \right] \\ & \leq 2 \cdot \exp \left[- \frac{C \cdot n \cdot v^2 \cdot \sigma^4}{\tau^4 \cdot (\tau^2 \cdot \|\beta^*\|_1^2 + C' \cdot \sigma^2)^2} \right], \end{aligned} \quad (\text{I.74})$$

for all $0 \leq v \leq C \cdot \tau^2 \cdot (\|\beta^*\|_1^2 + C' \cdot \sigma^2)$ and a sufficiently large sample size n .

Analysis of Term (ii): For term (ii) in (I.72), by union bound we have

$$\mathbb{P}(\bar{\mathcal{E}}) = \mathbb{P} \left(\max_{1 \leq i \leq n} \|\mathbf{x}_i\|_\infty > \tau \right) \leq n \cdot \mathbb{P}(\|\mathbf{x}_i\|_\infty > \tau) \leq n \cdot d \cdot \mathbb{P}(|x_{i,j}| > \tau).$$

Moreover, $x_{i,j}$ is sub-Gaussian with $\|x_{i,j}\|_{\psi_2} = C$. Thus, by Lemma 5.5 in [28] we have

$$\mathbb{P}(\bar{\mathcal{E}}) \leq n \cdot d \cdot 2 \cdot \exp(-C' \cdot \tau^2) = 2 \cdot \exp(-C' \cdot \tau^2 + \log n + \log d). \quad (\text{I.75})$$

Plugging (I.74) and (I.75) into (I.72), we obtain

$$\begin{aligned} & \mathbb{P}\left[4/\sigma^2 \cdot \left|\frac{1}{n} \sum_{i=1}^n \xi_i \cdot x_{i,j} \cdot x_{i,k} \cdot y_i^2\right| > v\right] \\ & \leq 2 \cdot \exp\left[-\frac{C \cdot n \cdot v^2 \cdot \sigma^4}{\tau^4 \cdot (\tau^2 \cdot \|\beta^*\|_1^2 + C' \cdot \sigma^2)^2}\right] + 2 \cdot \exp(-C'' \cdot \tau^2 + \log n + \log d). \end{aligned} \quad (\text{I.76})$$

Note that (I.71) is obtained by applying Lemmas J.4 and J.5 with $\phi(v) = \exp(u \cdot v)$. Since Lemmas J.4 and J.5 allow any increasing convex function $\phi(\cdot)$, similar results hold correspondingly. Hence, applying Panchenko's theorem in Lemma J.6 to (I.71), from (I.76) we have

$$\begin{aligned} \mathbb{P}\left[\left|\{T_n(\beta^*) - \mathbb{E}_{\beta^*}[T_n(\beta^*)]\}_{j,k}\right| > v\right] & \leq 2 \cdot e \cdot \exp\left[-\frac{C \cdot n \cdot v^2 \cdot \sigma^4}{\tau^4 \cdot (\tau^2 \cdot \|\beta^*\|_1^2 + C' \cdot \sigma^2)^2}\right] \\ & \quad + 2 \cdot e \cdot \exp(-C'' \cdot \tau^2 + \log n + \log d). \end{aligned}$$

Furthermore, by union bound we have

$$\begin{aligned} \mathbb{P}\left[\|T_n(\beta^*) - \mathbb{E}_{\beta^*}[T_n(\beta^*)]\|_{\infty, \infty} > v\right] & \leq \sum_{j=1}^d \sum_{k=1}^d \mathbb{P}\left[\left|\{T_n(\beta^*) - \mathbb{E}_{\beta^*}[T_n(\beta^*)]\}_{j,k}\right| > v\right] \\ & \leq 2 \cdot e \cdot \exp\left[-\frac{C \cdot n \cdot v^2 \cdot \sigma^4}{\tau^4 \cdot (\tau^2 \cdot \|\beta^*\|_1^2 + C' \cdot \sigma^2)^2} + 2 \cdot \log d\right] \\ & \quad + 2 \cdot e \cdot \exp(-C'' \cdot \tau^2 + \log n + 3 \cdot \log d). \end{aligned} \quad (\text{I.77})$$

To ensure the right-hand side is upper bounded by δ , we set the second term on the right-hand side of (I.77) to be $\delta/2$. Then we obtain

$$\tau = C \cdot \sqrt{\log n + 3 \cdot \log d + \log(4 \cdot e/\delta)}.$$

Let the first term on the right-hand side of (I.77) be upper bounded by $\delta/2$ and plugging in τ , we then obtain

$$\begin{aligned} v & = C \cdot [\log n + 3 \cdot \log d + \log(4 \cdot e/\delta)] \\ & \quad \cdot \left\{[\log n + 3 \cdot \log d + \log(4 \cdot e/\delta)] \cdot \|\beta^*\|_1^2 + C' \cdot \sigma^2\right\} / \sigma^2 \cdot \sqrt{\frac{\log(4 \cdot e/\delta) + 2 \cdot \log d}{n}}. \end{aligned}$$

Therefore, by setting $\delta = 4 \cdot e/d$ we have that

$$\begin{aligned} & \|T_n(\beta^*) - \mathbb{E}_{\beta^*}[T_n(\beta^*)]\|_{\infty, \infty} \leq v \\ & = C \cdot (\log n + 4 \cdot \log d) \cdot [(\log n + 4 \cdot \log d) \cdot \|\beta^*\|_1^2 + C' \cdot \sigma^2] / \sigma^2 \cdot \sqrt{\frac{\log d}{n}} \end{aligned}$$

holds with probability at least $1 - 4 \cdot e/d$, which completes the proof of Lemma F.6. \square

I.10 Proof of Lemma F.7

Proof. For any $j, k \in \{1, \dots, d\}$, by the mean-value theorem we have

$$\begin{aligned} \|T_n(\beta) - T_n(\beta^*)\|_{\infty, \infty} & = \max_{j,k \in \{1, \dots, d\}} \left| [T_n(\beta)]_{j,k} - [T_n(\beta^*)]_{j,k} \right| \\ & = \max_{j,k \in \{1, \dots, d\}} \left| (\beta - \beta^*)^\top \cdot \nabla [T_n(\beta^\#)]_{j,k} \right| \\ & \leq \|\beta - \beta^*\|_1 \cdot \max_{j,k \in \{1, \dots, d\}} \left\| \nabla [T_n(\beta^\#)]_{j,k} \right\|_{\infty}, \end{aligned} \quad (\text{I.78})$$

where $\beta^\#$ is an intermediate value between β and β^* . According to (I.67), (I.68) and the definition of $T_n(\cdot)$ in (2.4), by calculation we have

$$\nabla [T_n(\beta^\#)]_{j,k} = \frac{1}{n} \sum_{i=1}^n \bar{v}_{\beta^\#}(\mathbf{x}_i, y_i) \cdot y_i^3 \cdot x_{i,j} \cdot x_{i,k} \cdot \mathbf{x}_i,$$

where

$$\bar{v}_{\beta}(\mathbf{x}, y) = \frac{8/\sigma^4}{\left[1 + \exp(-2 \cdot y \cdot \langle \beta, \mathbf{x} \rangle / \sigma^2)\right] \cdot \left[1 + \exp(2 \cdot y \cdot \langle \beta, \mathbf{x} \rangle / \sigma^2)\right]^2} \quad (\text{I.79})$$

$$= \frac{8/\sigma^4}{\left[1 + \exp(-2 \cdot y \cdot \langle \beta, \mathbf{x} \rangle / \sigma^2)\right]^2 \cdot \left[1 + \exp(2 \cdot y \cdot \langle \beta, \mathbf{x} \rangle / \sigma^2)\right]}.$$

For notational simplicity, we define the following events

$$\mathcal{E} = \{\|\mathbf{x}_i\|_{\infty} \leq \tau, \text{ for all } i = 1, \dots, n\}, \quad \text{and } \mathcal{E}' = \{|v_i| \leq \tau', \text{ for all } i = 1, \dots, n\},$$

where $\tau > 0$ and $\tau' > 0$ will be specified later. By union bound we have

$$\begin{aligned} \mathbb{P}\left\{\left\|\nabla[T_n(\beta^{\#})]_{j,k}\right\|_{\infty} > v\right\} &\leq d \cdot \mathbb{P}\left(\left\{\left|\nabla[T_n(\beta^{\#})]_{j,k}\right|_l > v\right\}\right) \\ &= d \cdot \mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{v}_{\beta^{\#}}(\mathbf{x}_i, y_i) \cdot y_i^3 \cdot x_{i,j} \cdot x_{i,k} \cdot x_{i,l}\right| > v\right]. \end{aligned} \quad (\text{I.80})$$

Let $\bar{\mathcal{E}}$ and $\bar{\mathcal{E}}'$ be the complement of \mathcal{E} and \mathcal{E}' respectively. On the right-hand side we have

$$\begin{aligned} \mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{v}_{\beta^{\#}}(\mathbf{x}_i, y_i) \cdot y_i^3 \cdot x_{i,j} \cdot x_{i,k} \cdot x_{i,l}\right| > v\right] &= \underbrace{\mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{v}_{\beta^{\#}}(\mathbf{x}_i, y_i) \cdot y_i^3 \cdot x_{i,j} \cdot x_{i,k} \cdot x_{i,l}\right| > v, \mathcal{E}, \mathcal{E}'\right]}_{(i)} \\ &\quad + \underbrace{\mathbb{P}(\bar{\mathcal{E}})}_{(ii)} + \underbrace{\mathbb{P}(\bar{\mathcal{E}}')}_{(iii)}. \end{aligned} \quad (\text{I.81})$$

Analysis of Term (i): For term (i) in (I.81), we have

$$\begin{aligned} &\mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{v}_{\beta^{\#}}(\mathbf{x}_i, y_i) \cdot y_i^3 \cdot x_{i,j} \cdot x_{i,k} \cdot x_{i,l}\right| > v, \mathcal{E}, \mathcal{E}'\right] \\ &= \mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{v}_{\beta^{\#}}(\mathbf{x}_i, y_i) \cdot y_i^3 \cdot x_{i,j} \cdot x_{i,k} \cdot x_{i,l} \cdot \mathbf{1}\{\|\mathbf{x}_i\|_{\infty} \leq \tau\} \cdot \mathbf{1}\{|v_i| \leq \tau'\}\right| > v, \mathcal{E}, \mathcal{E}'\right] \\ &\leq \mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{v}_{\beta^{\#}}(\mathbf{x}_i, y_i) \cdot y_i^3 \cdot x_{i,j} \cdot x_{i,k} \cdot x_{i,l} \cdot \mathbf{1}\{\|\mathbf{x}_i\|_{\infty} \leq \tau\} \cdot \mathbf{1}\{|v_i| \leq \tau'\}\right| > v\right]. \end{aligned}$$

To avoid confusion, note that v_i is the noise in mixture of regression model, while v appears in the tail bound. By applying union bound to the right-hand side of the above inequality, we have

$$\begin{aligned} &\mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{v}_{\beta^{\#}}(\mathbf{x}_i, y_i) \cdot y_i^3 \cdot x_{i,j} \cdot x_{i,k} \cdot x_{i,l}\right| > v, \mathcal{E}, \mathcal{E}'\right] \\ &\leq \sum_{i=1}^n \mathbb{P}\left[\left|\bar{v}_{\beta^{\#}}(\mathbf{x}_i, y_i) \cdot y_i^3 \cdot x_{i,j} \cdot x_{i,k} \cdot x_{i,l} \cdot \mathbf{1}\{\|\mathbf{x}_i\|_{\infty} \leq \tau\} \cdot \mathbf{1}\{|v_i| \leq \tau'\}\right| > v\right]. \end{aligned}$$

By (I.79) we have $|\bar{v}_{\beta^{\#}}(\mathbf{x}_i, y_i)| \leq 16/\sigma^4$. Hence, we obtain

$$\begin{aligned} &\mathbb{P}\left[\left|\bar{v}_{\beta^{\#}}(\mathbf{x}_i, y_i) \cdot y_i^3 \cdot x_{i,j} \cdot x_{i,k} \cdot x_{i,l} \cdot \mathbf{1}\{\|\mathbf{x}_i\|_{\infty} \leq \tau\} \cdot \mathbf{1}\{|v_i| \leq \tau'\}\right| > v\right] \\ &\leq \mathbb{P}\left[\left|y_i^3 \cdot x_{i,j} \cdot x_{i,k} \cdot x_{i,l} \cdot \mathbf{1}\{\|\mathbf{x}_i\|_{\infty} \leq \tau\} \cdot \mathbf{1}\{|v_i| \leq \tau'\}\right| > \sigma^4/16 \cdot v\right]. \end{aligned} \quad (\text{I.82})$$

Recall that in mixture of regression model we have $y_i = z_i \cdot \langle \beta^*, \mathbf{x}_i \rangle + v_i$, where z_i is a Rademacher random variable, $\mathbf{x}_i \sim N(0, \mathbf{I}_d)$ and $v_i \sim N(0, \sigma^2)$. Hence, we have

$$\begin{aligned} \left|y_i^3 \cdot \mathbf{1}\{\|\mathbf{x}_i\|_{\infty} \leq \tau\} \cdot \mathbf{1}\{|v_i| \leq \tau'\}\right| &\leq \left(|z_i \cdot \langle \beta^*, \mathbf{x}_i \rangle \cdot \mathbf{1}\{\|\mathbf{x}_i\|_{\infty} \leq \tau\} + |v_i \cdot \mathbf{1}\{|v_i| \leq \tau'\}|\right)^3 \\ &\leq (\tau \cdot \|\beta^*\|_1 + \tau')^3, \end{aligned}$$

$$\left|x_{i,j} \cdot x_{i,k} \cdot x_{i,l} \cdot \mathbf{1}\{\|\mathbf{x}_i\|_{\infty} \leq \tau\}\right| \leq |x_{i,j}| \cdot \mathbf{1}\{\|\mathbf{x}_i\|_{\infty} \leq \tau\} \leq \tau^3.$$

Taking $v = 16 \cdot (\tau \cdot \|\beta^*\|_1 + \tau')^3 \cdot \tau^3 / \sigma^4$, we have that the right-hand side of (I.82) is zero. Hence term (i) in (I.81) is zero.

Analysis of Term (ii): For term (ii) in (I.81), by union bound we have

$$\mathbb{P}(\bar{\mathcal{E}}) = \mathbb{P}\left(\max_{i \in \{1, \dots, n\}} \|\mathbf{x}_i\|_\infty > \tau\right) \leq n \cdot \mathbb{P}(\|\mathbf{x}_i\|_\infty > \tau) \leq n \cdot d \cdot \mathbb{P}(|x_{i,j}| > \tau).$$

Moreover, we have that $x_{i,j}$ is sub-Gaussian with $\|x_{i,j}\|_{\psi_2} = C$. Therefore, by Lemma 5.5 in [28] we have

$$\mathbb{P}(\bar{\mathcal{E}}) \leq n \cdot d \cdot \mathbb{P}(|x_{i,j}| > \tau) \leq n \cdot d \cdot 2 \cdot \exp(-C' \cdot \tau^2).$$

Analysis of Term (iii): Since v_i is sub-Gaussian with $\|v_i\|_{\psi_2} = C \cdot \sigma$, by Lemma 5.5 in [28] and union bound, for term (iii) in (I.81) we have

$$\mathbb{P}(\bar{\mathcal{E}}') \leq n \cdot \mathbb{P}(|v_i| > \tau') \leq n \cdot 2 \cdot \exp(-C' \cdot \tau'^2 / \sigma^2).$$

To ensure the right-hand side of (I.81) is upper bounded by δ , we set τ and τ' to be

$$\tau = C \cdot \sqrt{\log d + \log n + \log(4/\delta)}, \quad \text{and} \quad \tau' = C' \cdot \sigma \cdot \sqrt{\log n + \log(4/\delta)} \quad (\text{I.83})$$

to ensure terms (ii) and (iii) are upper bounded by $\delta/2$ correspondingly. Finally, by (I.80), (I.81) and union bound we have

$$\begin{aligned} & \mathbb{P}\left\{\max_{j,k \in \{1, \dots, d\}} \left\| \nabla [T_n(\boldsymbol{\beta}^\#)]_{j,k} \right\|_\infty > v\right\} \\ & \leq d^2 \cdot d \cdot \mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \bar{v}_{\boldsymbol{\beta}^\#}(\mathbf{x}_i, y_i) \cdot y_i^3 \cdot x_{i,j} \cdot x_{i,k} \cdot x_{i,l}\right| > v\right] \leq d^3 \cdot \delta \end{aligned}$$

for $v = 16 \cdot (\tau \cdot \|\boldsymbol{\beta}^*\|_1 + \tau')^3 \cdot \tau^3 / \sigma^4$ with τ and τ' specified in (I.83). Then by setting $\delta = 4 \cdot d^{-4}$ and plugging it into (I.83), we have

$$\begin{aligned} v &= 16 \cdot \left[C \cdot \sqrt{5 \cdot \log d + \log n} \cdot \|\boldsymbol{\beta}^*\|_1 + C' \cdot \sigma \cdot \sqrt{4 \cdot \log d + \log n} \right]^3 \cdot \left[C \cdot \sqrt{5 \cdot \log d + \log n} \right]^3 / \sigma^4 \\ &\leq C'' \cdot (\|\boldsymbol{\beta}^*\|_1 + C''' \cdot \sigma)^3 \cdot (5 \cdot \log d + \log n)^3, \end{aligned}$$

which together with (I.78) concludes the proof of Lemma F.7. \square

J Auxiliary Results

In this section, we lay out several auxiliary lemmas. Lemmas J.1-J.3 provide useful properties of sub-Gaussian random variables. Lemmas J.4 and J.5 establish the symmetrization and contraction results. Lemma J.6 is Panchenko's theorem. For more details of these results, see [6, 28].

Lemma J.1. Let Z_1, \dots, Z_k be the k independent zero-mean sub-Gaussian random variables, for $Z = \sum_{j=1}^k Z_j$ we have $\|Z\|_{\psi_2}^2 \leq C \cdot \sum_{j=1}^k \|Z_j\|_{\psi_2}^2$, where $C > 0$ is a constant.

Lemma J.2. For Z_1 and Z_2 being two sub-Gaussian random variables, $Z_1 \cdot Z_2$ is a sub-exponential random variable with

$$\|Z_1 \cdot Z_2\|_{\psi_1} \leq C \cdot \max\{\|Z_1\|_{\psi_2}^2, \|Z_2\|_{\psi_2}^2\},$$

where $C > 0$ is a constant.

Lemma J.3. For Z being sub-Gaussian or sub-exponential, it holds that $\|Z - \mathbb{E}Z\|_{\psi_2} \leq 2 \cdot \|Z\|_{\psi_2}$ or $\|Z - \mathbb{E}Z\|_{\psi_1} \leq 2 \cdot \|Z\|_{\psi_1}$ correspondingly.

Lemma J.4. Let $\mathbf{z}_1, \dots, \mathbf{z}_n$ be the n independent realizations of the random vector $\mathbf{Z} \in \mathcal{Z}$ and \mathcal{F} be a function class defined on \mathcal{Z} . For any increasing convex function $\phi(\cdot)$ we have

$$\mathbb{E}\left\{\phi\left[\sup_{f \in \mathcal{F}} \left|\sum_{i=1}^n f(\mathbf{z}_i) - \mathbb{E}Z\right|\right]\right\} \leq \mathbb{E}\left\{\phi\left[\sup_{f \in \mathcal{F}} \left|\sum_{i=1}^n \xi_i \cdot f(\mathbf{z}_i)\right|\right]\right\},$$

where ξ_1, \dots, ξ_n are i.i.d. Rademacher random variables that are independent of $\mathbf{z}_1, \dots, \mathbf{z}_n$.

Lemma J.5. Let $\mathbf{z}_1, \dots, \mathbf{z}_n$ be the n independent realizations of the random vector $\mathbf{Z} \in \mathcal{Z}$ and \mathcal{F} be a function class defined on \mathcal{Z} . We consider the Lipschitz functions $\psi_i(\cdot)$ ($i = 1, \dots, n$) that satisfy

$$|\psi_i(v) - \psi_i(v')| \leq L \cdot |v - v'|, \quad \text{for all } v, v' \in \mathbb{R},$$

and $\psi_i(0) = 0$. For any increasing convex function $\phi(\cdot)$ we have

$$\mathbb{E}\left\{\phi\left[\left|\sup_{f \in \mathcal{F}} \sum_{i=1}^n \xi_i \cdot \psi_i[f(\mathbf{z}_i)]\right|\right]\right\} \leq \mathbb{E}\left\{\phi\left[2 \cdot \left|L \cdot \sup_{f \in \mathcal{F}} \sum_{i=1}^n \xi_i \cdot f(\mathbf{z}_i)\right|\right]\right\},$$

where ξ_1, \dots, ξ_n are i.i.d. Rademacher random variables that are independent of $\mathbf{z}_1, \dots, \mathbf{z}_n$.

Lemma J.6. Suppose that Z_1 and Z_2 are two random variables that satisfy $\mathbb{E}[\phi(Z_2)] \leq \mathbb{E}[\phi(Z_1)]$ for any increasing convex function $\phi(\cdot)$. Assuming that $\mathbb{P}(Z_1 \geq v) \leq C \cdot \exp(-C' \cdot v^\alpha)$ ($\alpha \geq 1$) holds for all $v \geq 0$, we have $\mathbb{P}(Z_2 \geq v) \leq C \cdot \exp(1 - C' \cdot v^\alpha)$.