

Supplementary Materials of “Transfer Learning with High-dimensional Generalized Linear Models”

Contents

| | | |
|------------|--|-----------|
| S.1 | More details | 2 |
| S.1.1 | A schematic to illustrate \mathcal{A} -Trans-GLM | 2 |
| S.1.2 | Theory | 3 |
| S.1.2.1 | Explicit forms of convergence rates in Assumption 5 | 3 |
| S.1.2.2 | Bound of a prediction error measure | 5 |
| S.1.3 | Additional numerical results | 6 |
| S.1.3.1 | More details about the implementation of numerical experiments . . | 6 |
| S.1.3.2 | Transfer learning on \mathcal{A}_h | 7 |
| S.1.3.3 | Transferable source detection | 8 |
| S.1.3.4 | Additional results of real data analysis | 11 |
| S.2 | Proofs | 13 |
| S.2.1 | Some lemmas | 16 |
| S.2.2 | Proof of lemmas | 18 |
| S.2.2.1 | Proof of Lemma 1 | 18 |
| S.2.2.2 | Proof of Lemma 2 | 18 |
| S.2.2.3 | Proof of Lemma 4 | 19 |
| S.2.2.4 | Proof of Lemma 5 | 23 |
| S.2.3 | Proof of theorems | 24 |
| S.2.3.1 | Proof of Theorem 1 | 24 |

| | |
|--|----|
| S.2.3.2 Proof of Theorem 2 | 31 |
| S.2.3.3 Proof of Theorem 3 | 32 |
| S.2.3.4 Proof of Theorem 4 | 34 |
| S.2.3.5 Proof of Theorem 5 | 37 |
| S.2.3.6 Proof of Theorem 6 | 44 |
| S.2.3.7 Proof of Theorem 7 | 44 |
| S.2.4 Proof of propositions | 44 |
| S.2.4.1 Proof of Proposition 1 | 44 |
| S.2.4.2 Proof of Proposition 2 | 53 |

S.1 More details

S.1.1 A schematic to illustrate \mathcal{A} -Trans-GLM

To better illustrate Algorithm 1, we draw a schematic in Figure 5. The blue point represents the target coefficient $\beta = \mathbf{w}^{(0)}$ and the surrounding blue circle represents the estimation error $\mathcal{O}_p\left(\sqrt{\frac{s \log p}{n_0}}\right)$. The purple point denotes the estimator $\hat{\beta}_{\text{naive-Lasso}}$ from the classical Lasso with only target data. The orange point represents $\mathbf{w}^{\mathcal{A}_h}$, which is the population version of the rough estimator we obtain from the transferring step by pooling target and source data in \mathcal{A} (see Section 2.3), and the surrounding orange circle denotes its estimation error. It can be seen that $\mathbf{w}^{\mathcal{A}_h}$ is a pooled version of $\{\mathbf{w}^{(k)}\}_{k \in \{0\} \cup \mathcal{A}_h}$, which is close to β when h is small. Starting from an initial estimate of β , the transferring step of \mathcal{A} -Trans-GLM algorithm updates the estimate to $\hat{\mathbf{w}}^{\mathcal{A}_h}$ (an estimate of $\mathbf{w}^{\mathcal{A}_h}$ based on source data in \mathcal{A} and the target data), then the debiasing step yields the final estimator $\hat{\beta}_{\mathcal{A}\text{-Trans-GLM}}$.

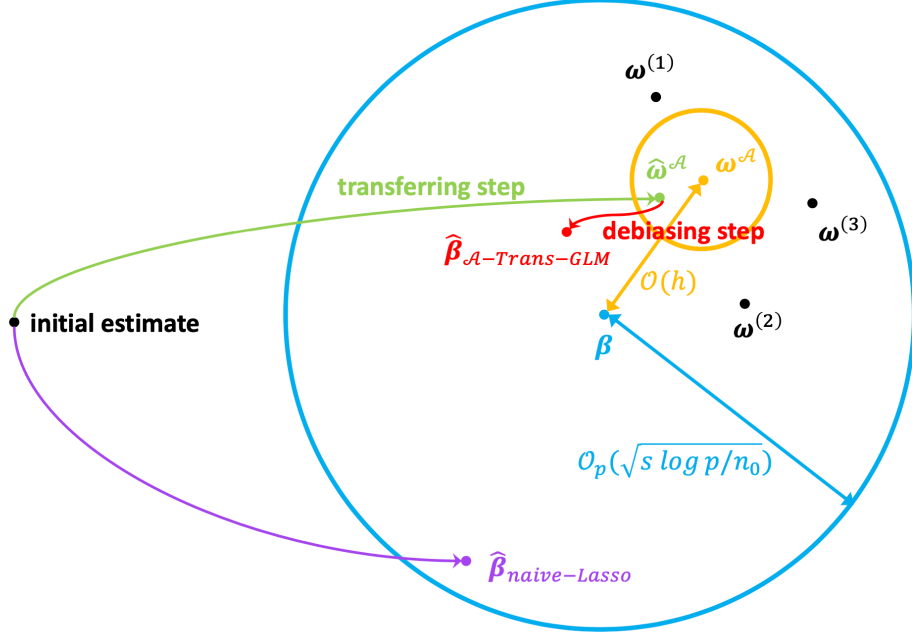


Figure 5: A schematic of \mathcal{A} -Trans-GLM (Algorithm 1). $\mathcal{A} = \{1, 2, 3\}$.

S.1.2 Theory

S.1.2.1 Explicit forms of convergence rates in Assumption 5

Proposition 1 (Explicit forms of $\Upsilon_1^{(k)}$, $\Gamma_1^{(k)}$, $\Gamma_2^{(k)}$, $g_1^{(k)}$ and $g_2^{(k)}$ for certain families). *Denote*

$$\Omega_k = \begin{cases} \sqrt{\frac{s \log p}{n_0}}, & k = 0, \\ \sqrt{\frac{s \log p}{n_k + n_0}} + \left(\frac{\log p}{n_k + n_0}\right)^{1/4} \sqrt{h} + \sqrt{s}h, & k \in \mathcal{A}_h \\ h' \sqrt{\frac{\log p}{n_k + n_0}} + \sqrt{\frac{s' \log p}{n_k + n_0}} \cdot W_k + \left(\frac{\log p}{n_k + n_0}\right)^{1/4} \sqrt{h' W_k}, & k \in \mathcal{A}_h^c, \end{cases}$$

where $W_k = 1 \vee \|\beta^{(k)} - \beta\|_2 \vee \|\beta^{(k)} - \mathbf{w}^{(k)}\|_2$. Assume Assumptions 1 and 2 hold. For Poisson model, it is further required that $h \leq U^{-1} \bar{U}$ and $U \sup_{k \in \mathcal{A}^c} \{\|\beta^{(k)} - \beta\|_1 \vee \|\beta^{(k)} - \mathbf{w}^{(k)}\|_1\} \leq \bar{U}$.

With $\lambda^{(k)[r]} = C \left(\sqrt{\frac{\log p}{n_k + n_0}} + h \right)$ when $k \in \mathcal{A}_h$, $\lambda^{(k)[r]} = C \sqrt{\frac{\log p}{n_k + n_0}} \cdot W_k$ when $k \in \mathcal{A}_h^c$ and

$\lambda^{(0)} = C\sqrt{\frac{\log p}{n_0}}$ for some sufficiently large constant $C > 0$, we have the following explicit forms of $\Gamma_1^{(k)}$, $\Gamma_2^{(k)}$, $\Upsilon_1^{(k)}$, $g_1^{(k)}$ and $g_2^{(k)}$ for logistic, linear and Poisson regression models.

(i) For the logistic regression model:

$$\begin{aligned}\Gamma_1^{(0)} &= \sqrt{\frac{s \log p}{n_0}}, & \Gamma_2^{(0)} &= \|\boldsymbol{\beta}\|_2 / \sqrt{n_0}, \\ \Upsilon_1^{(k)} &= \Omega_k, & \Gamma_1^{(k)} &= \sqrt{\frac{1}{n_0}} \Omega_k, & \Gamma_2^{(k)} &= \sqrt{\frac{1}{n_0}} \cdot [\|\mathbf{w}^{(k)}\|_2 \mathbf{1}(k \in \mathcal{A}_h) + \|\boldsymbol{\beta}^{(k)}\|_2 \mathbf{1}(k \in \mathcal{A}_h^c)],\end{aligned}$$

$$g_1^{(k)}(\zeta) = g_2^{(k)}(\zeta) = \exp(-\zeta^2).$$

(ii) For the linear model:

$$\begin{aligned}\Gamma_1^{(0)} &= \sqrt{\frac{s \log p}{n_0}} \cdot \|\boldsymbol{\beta}\|_2, & \Gamma_2^{(0)} &= (\|\boldsymbol{\beta}\|_2^2 \vee \|\boldsymbol{\beta}\|_2) / \sqrt{n_0}, \\ \Upsilon_1^{(k)} &= \Omega_k \cdot [\|\mathbf{w}^{(k)}\|_2 \mathbf{1}(k \in \mathcal{A}_h) + \|\boldsymbol{\beta}^{(k)}\|_2 \mathbf{1}(k \in \mathcal{A}_h^c)], \\ \Gamma_1^{(k)} &= \sqrt{\frac{1}{n_0}} \Omega_k \cdot [\|\mathbf{w}^{(k)}\|_2 \mathbf{1}(k \in \mathcal{A}_h) + \|\boldsymbol{\beta}^{(k)}\|_2 \mathbf{1}(k \in \mathcal{A}_h^c)], \\ \Gamma_2^{(k)} &= \sqrt{\frac{1}{n_0}} [(\|\mathbf{w}^{(k)}\|_2^2 \vee \|\mathbf{w}^{(k)}\|_2) \mathbf{1}(k \in \mathcal{A}_h) + (\|\boldsymbol{\beta}^{(k)}\|_2^2 \vee \|\boldsymbol{\beta}^{(k)}\|_2) \mathbf{1}(k \in \mathcal{A}_h^c)], \\ g_1^{(k)}(\zeta) &= g_2^{(k)}(\zeta) = \exp(-\zeta^2), k \neq 0; & g_1^{(0)}(\zeta) &= \exp(-\zeta^2), g_2^{(0)}(\zeta) = \exp\{-n_0\} + \exp(-\zeta^2).\end{aligned}$$

(iii) For the Poisson regression model with bounded predictors ($\sup_k \|\mathbf{x}^{(k)}\|_\infty \leq U < \infty$):

$$\begin{aligned}\Gamma_1^{(0)} &= \sqrt{\frac{s \log p}{n_0}} \cdot \exp(U\|\boldsymbol{\beta}\|_1), & \Gamma_2^{(0)} &= \exp(U\|\boldsymbol{\beta}\|_1) \cdot \frac{1 + \|\boldsymbol{\beta}\|_2 + U\|\boldsymbol{\beta}\|_1}{\sqrt{n_0}}, \\ \Upsilon_1^{(k)} &= \Omega_k \cdot \exp\{U\|\mathbf{w}^{(k)}\|_1 \cdot \mathbf{1}(k \in \mathcal{A}_h) + U\|\boldsymbol{\beta}^{(k)}\|_1 \cdot \mathbf{1}(k \in \mathcal{A}_h^c)\}, \\ \Gamma_1^{(k)} &= \sqrt{\frac{1}{n_0}} \Omega_k \cdot \exp\{U\|\mathbf{w}^{(k)}\|_1 \cdot \mathbf{1}(k \in \mathcal{A}_h) + U\|\boldsymbol{\beta}^{(k)}\|_1 \cdot \mathbf{1}(k \in \mathcal{A}_h^c)\},\end{aligned}$$

$$\begin{aligned}\Gamma_2^{(k)} &= \sqrt{\frac{1}{n_0}} \left[\exp(U\|\mathbf{w}^{(k)}\|_1) (1 + \|\mathbf{w}^{(k)}\|_2 + U\|\mathbf{w}^{(k)}\|_1) \cdot \mathbf{1}(k \in \mathcal{A}_h) \right. \\ &\quad \left. + \exp(U\|\boldsymbol{\beta}^{(k)}\|_1) (1 + \|\boldsymbol{\beta}^{(k)}\|_2 + U\|\boldsymbol{\beta}^{(k)}\|_1) \cdot \mathbf{1}(k \in \mathcal{A}_h^c) \right], \\ g_1^{(k)}(\zeta) &= g_2^{(k)}(\zeta) = \exp(-\zeta^2), k \neq 0; \quad g_1^{(0)}(\zeta) = \exp(-\zeta^2), g_2^{(0)}(\zeta) = \zeta^{-2}.\end{aligned}$$

S.1.2.2 Bound of a prediction error measure

Denote $L_{n_0}^{(0)}(\mathbf{w}) = -\frac{1}{n_0} \sum_{i=1}^{n_0} \log \rho(\mathbf{x}_i^{(0)}) - \frac{1}{n_0} (\mathbf{y}^{(0)})^T \mathbf{X}^{(0)} \mathbf{w} + \frac{1}{n_0} \sum_{i=1}^{n_0} \psi(\mathbf{w}^T \mathbf{x}_i^{(0)})$. Suggested by [Loh and Wainwright \(2015\)](#), we consider a special measure of the prediction error, which is defined by $\langle \nabla L_{n_0}^{(0)}(\hat{\boldsymbol{\beta}}) - \nabla L_{n_0}^{(0)}(\boldsymbol{\beta}), \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \rangle$, where $\nabla L_{n_0}^{(0)}(\mathbf{w}) = -\frac{1}{n_0} (\mathbf{X}^{(0)})^T \mathbf{y}^{(0)} + \frac{1}{n_0} \sum_{i=1}^{n_0} \mathbf{x}_i^{(0)} \psi(\mathbf{w}^T \mathbf{x}_i^{(0)}) \in \mathbb{R}^p$. Note that the loss function $L_{n_0}^{(0)}$ is convex, therefore this quantity is non-negative. As argued in their paper, $\langle \nabla L_{n_0}^{(0)}(\hat{\boldsymbol{\beta}}) - \nabla L_{n_0}^{(0)}(\boldsymbol{\beta}), \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \rangle$ can be easily interpreted in GLMs. For example, in the case of linear models where $\psi(u) = u^2/2$, this measure equals to the in-sample square loss $\|\mathbf{X}^{(0)}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\|_2^2$. For general GLMs, it is equivalent to the symmetrized Bregman divergence defined by ψ .

Next we would like to present the bounds of $\langle \nabla L_{n_0}^{(0)}(\hat{\boldsymbol{\beta}}) - \nabla L_{n_0}^{(0)}(\boldsymbol{\beta}), \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \rangle$ with and without Assumption 4 for \mathcal{A}_h -Trans-GLM.

Theorem 6 (Bound of a prediction error measure for \mathcal{A}_h -Trans-GLM with Assumption 4). *By imposing the same conditions in Theorem 1, we have*

$$\begin{aligned}\sup_{\boldsymbol{\xi} \in \Xi(s, h)} \langle \nabla L_{n_0}^{(0)}(\hat{\boldsymbol{\beta}}) - \nabla L_{n_0}^{(0)}(\boldsymbol{\beta}), \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \rangle &\lesssim s \left(\frac{\log p}{n_0} \right)^{1/2} \left(\frac{\log p}{n_{\mathcal{A}_h} + n_0} \right)^{1/2} + h \left(\frac{\log p}{n_0} \right)^{1/2} \\ &\quad + \left(\frac{s \log p}{n_0} \right)^{1/2} \left(\frac{\log p}{n_{\mathcal{A}_h} + n_0} \right)^{1/4} h^{1/2},\end{aligned}$$

with probability at least $1 - n_0^{-1}$.

Remark 11. When $h \ll s \sqrt{\frac{\log p}{n_0}}$ and $n_{\mathcal{A}_h} \gg n_0$, the upper bounds in (i) are better than the classical Lasso bound $\mathcal{O}_p\left(\frac{s \log p}{n_0}\right)$ with target data ([Loh and Wainwright, 2015](#)).

Theorem 7 (Bound of a prediction error measure for \mathcal{A}_h -Trans-GLM without Assumption 4). *By imposing the same conditions in Theorem 3, we have*

$$\begin{aligned} \sup_{\xi \in \Xi(s, h)} \langle \nabla L_{n_0}^{(0)}(\hat{\boldsymbol{\beta}}) - \nabla L_{n_0}^{(0)}(\boldsymbol{\beta}), \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \rangle &\lesssim s \left(\frac{\log p}{n_0} \right)^{1/2} \left(\frac{\log p}{n_{\mathcal{A}_h} + n_0} \right)^{1/2} + sh \left(\frac{\log p}{n_0} \right)^{1/2} \\ &\quad + \left(\frac{s \log p}{n_0} \right)^{1/2} \left(\frac{\log p}{n_{\mathcal{A}_h} + n_0} \right)^{1/4} h^{1/2} \end{aligned}$$

with probability at least $1 - n_0^{-1}$.

Remark 12. *When $h \ll \sqrt{\frac{\log p}{n_0}}$ and $n_{\mathcal{A}_h} \gg n_0$, the upper bounds in (i) are better than the classical Lasso bound $\mathcal{O}_p\left(\frac{s \log p}{n_0}\right)$ with target data (Loh and Wainwright, 2015).*

S.1.3 Additional numerical results

S.1.3.1 More details about the implementation of numerical experiments

All experiments in this paper are conducted in R. The GLM Lasso is implemented via R package `glmnet` (Friedman et al., 2010). We summarize R codes for GLM transfer learning algorithms in a new R package `glmtrans`, which is available on CRAN. We use 10-fold cross-validation to choose the penalty parameter for naïve-Lasso and our GLM transfer learning algorithms. The largest λ which achieves one standard error within the minimum cross-validation error will be chosen for the transferring step, which is sometimes called “*lambda.1se*” (Friedman et al., 2010). To effectively debias the transferring step, we choose the lambda achieving minimal cross-validation error, which is often called “*lambda.min*”. Since in transferable source detection, the first step is the same as the transferring step of $\{k\}$ -Trans-GLM, therefore we keep the same setting as the transferring step in Algorithm 1, i.e. take “*lambda.1se*”. And in Algorithm 2, we set the constant $C_0 = 2$. In the two-step

transfer learning procedure of Algorithm 3, we use “lambda.min” in both transferring and debiasing steps.

In real-data studies, SVM with RBF kernel is implemented by package `e1071`, and decision trees are implemented through package `rpart`. We fit the random forest via package `randomForest`, and implement boosting trees through package `fastAdaboost`. The number of weak classifiers in boosting trees is set to be 50. Since the sample size of each state is small and some states have very imbalanced responses, we change the cross-validation folds from default 10 to 5 for all Lasso-based methods. All the other parameters are kept the same as the default settings.

S.1.3.2 Transfer learning on \mathcal{A}_h

In this section, we supplement more numerical results about the performance of Algorithm 1 under different h and $(\{n_k\}_{k=0}^K, p, s)$ settings. In addition to the previous $(\{n_k\}_{k=0}^K, p, s)$ setting studied in Section 4.1.1, the following two settings of $(\{n_k\}_{k=0}^K, p, s)$ are considered:

- (i) $n_k = 150$ for all $k = 0, \dots, K$, $p = 1000$, $s = 15$;
- (ii) $n_k = 200$ for all $k = 0, \dots, K$, $p = 2000$, $s = 20$.

Given each $(\{n_k\}_{k=0}^K, p, s)$ setting, consider the same setting we use in Section 4.1.1. All the experiments are replicated 200 times and the average ℓ_2 -estimation errors of \mathcal{A}_h -Trans-GLM and naïve-Lasso under linear, logistic and Poisson regression models are shown in Figure 6 and 7.

The trends in Figures 6 and 7 are similar to that in Figure 1. As K increases, the estimation error of $K_{\mathcal{A}_h}$ -Trans-GLM continues declining and is lower than that of naïve-Lasso on target data only.

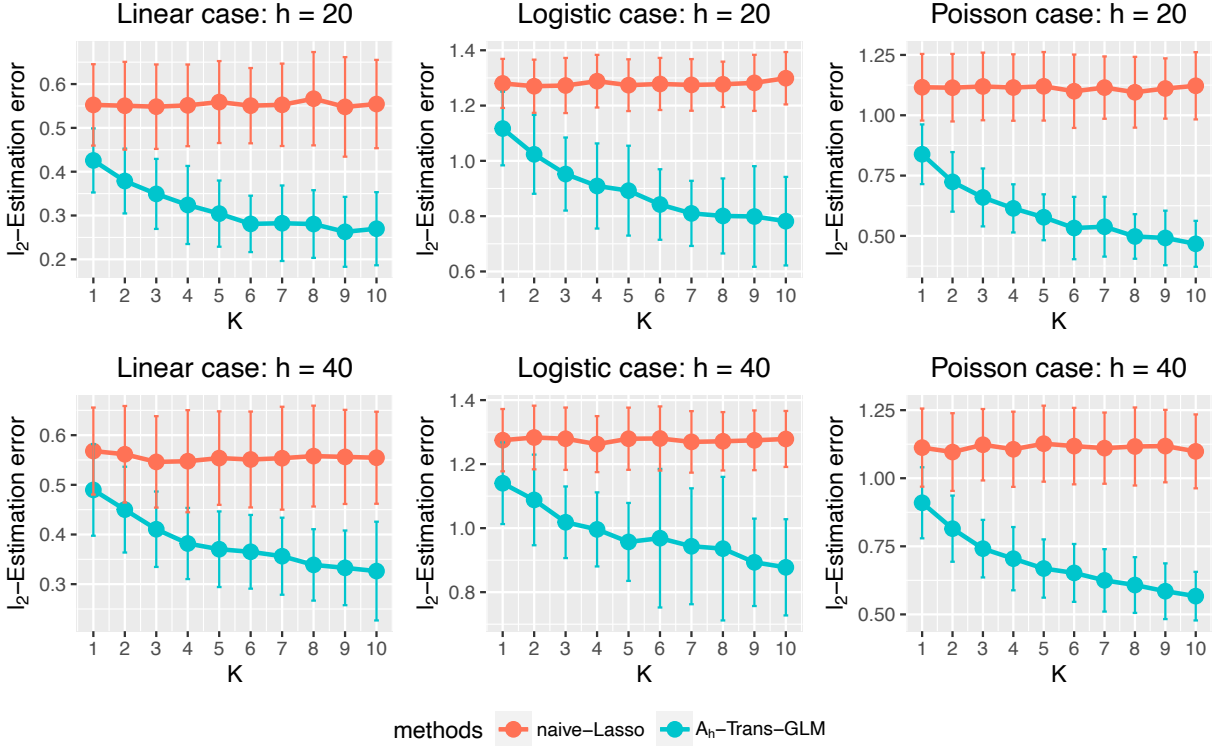


Figure 6: The average ℓ_2 -estimation of \mathcal{A}_h -Trans-GLM and naïve-Lasso under linear, logistic and Poisson regression models with different settings of h and K . $n_k = 150$ for all $k = 0, \dots, K$, $p = 1000$, $s = 15$. Error bars denote the standard deviations.

S.1.3.3 Transferable source detection

In this section we supplement more experimental results in the case that some sources are not in the level- h transferring set \mathcal{A}_h . The model settings are the same as those in Section 4.1.2. In addition to the setting used in Section 4.1.2, two more ones are considered:

- (i) $n_k = 100$ for all $k = 0, \dots, K$, $p = 500$, $s = 10$;
- (ii) $n_k = 150$ for all $k = 0, \dots, K$, $p = 1000$, $s = 15$.

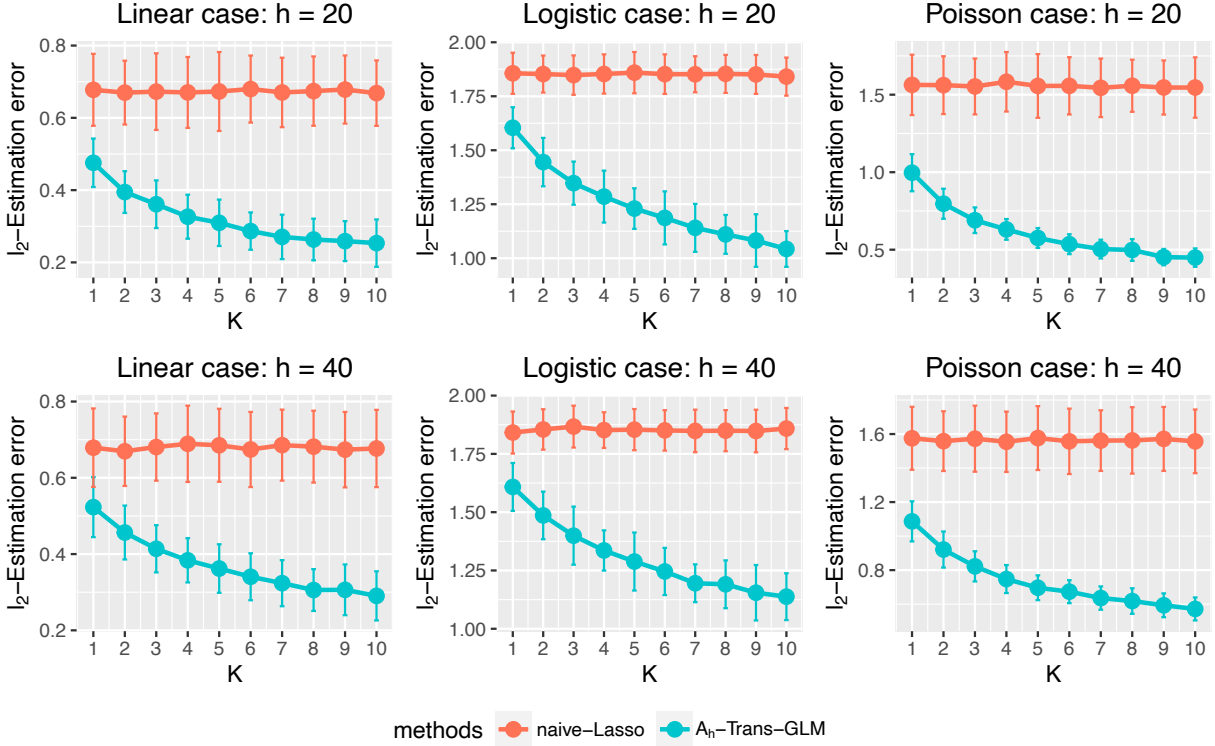


Figure 7: The average ℓ_2 -estimation of \mathcal{A}_h -Trans-GLM and naïve-Lasso under linear, logistic, and Poisson regression models with different settings of h and K . $n_k = 200$ for all $k = 0, \dots, K$, $p = 2000$, $s = 20$. Error bars denote the standard deviations.

We vary the values of $|K_{\mathcal{A}_h}|$ and h , and repeat each setting for 200 times. The average ℓ_2 -estimation errors are summarized in Figures 8 and 9.

Similar to Figure 3, it can be seen that \mathcal{A}_h -Trans-GLM always achieves the best performance. Trans-GLM mimics the behavior of \mathcal{A}_h -Trans-GLM very well in most cases, implying that the detection algorithm can accurately identify \mathcal{A} . We also observe that for linear models and logistic regression models, when $h = 40$, there is a gap between the estimation error of \mathcal{A}_h -Trans-GLM and Trans-GLM, meaning that when h increases,

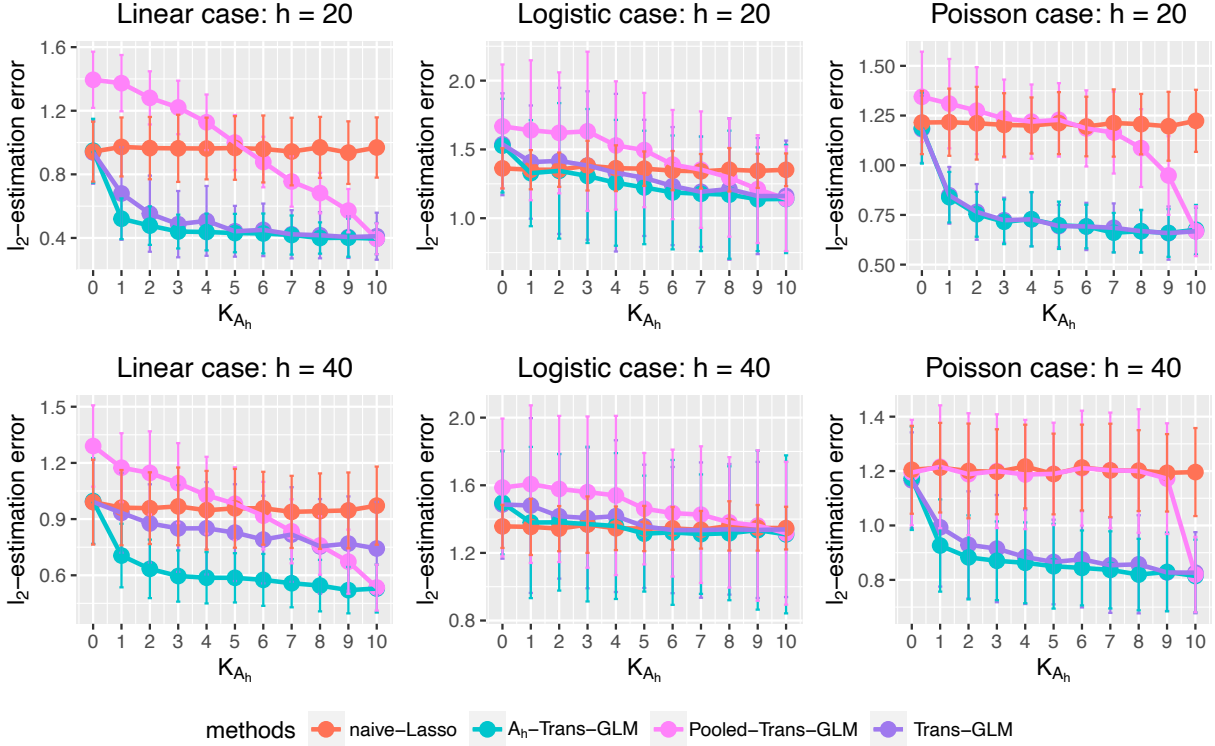


Figure 8: The average ℓ_2 -estimation error of various models with different settings of h and $K_{\mathcal{A}_h}$ when $K = 10$. $n_k = 100$ for all $k = 0, \dots, K$, $p = 500$, $s = 10$. Error bars denote the standard deviations.

Trans-GLM might begin missing sources in \mathcal{A}_h or wrongly including sources in \mathcal{A}_h .

Furthermore, we pick the setting when $p = 2000$ and try different numbers of folds in the cross-validation procedure of Algorithm 2 (steps 2-5). The results are displayed in Figure 10. The findings suggest that more cross-validation folds may lead to better performance of Trans-GLM. When the cross-validation folds are large, the detection is likely to be more accurate. In the meantime, this may cause more computational burdens. Therefore, we may choose a moderate fold number like 3 or 5 in practice to achieve a good trade-off

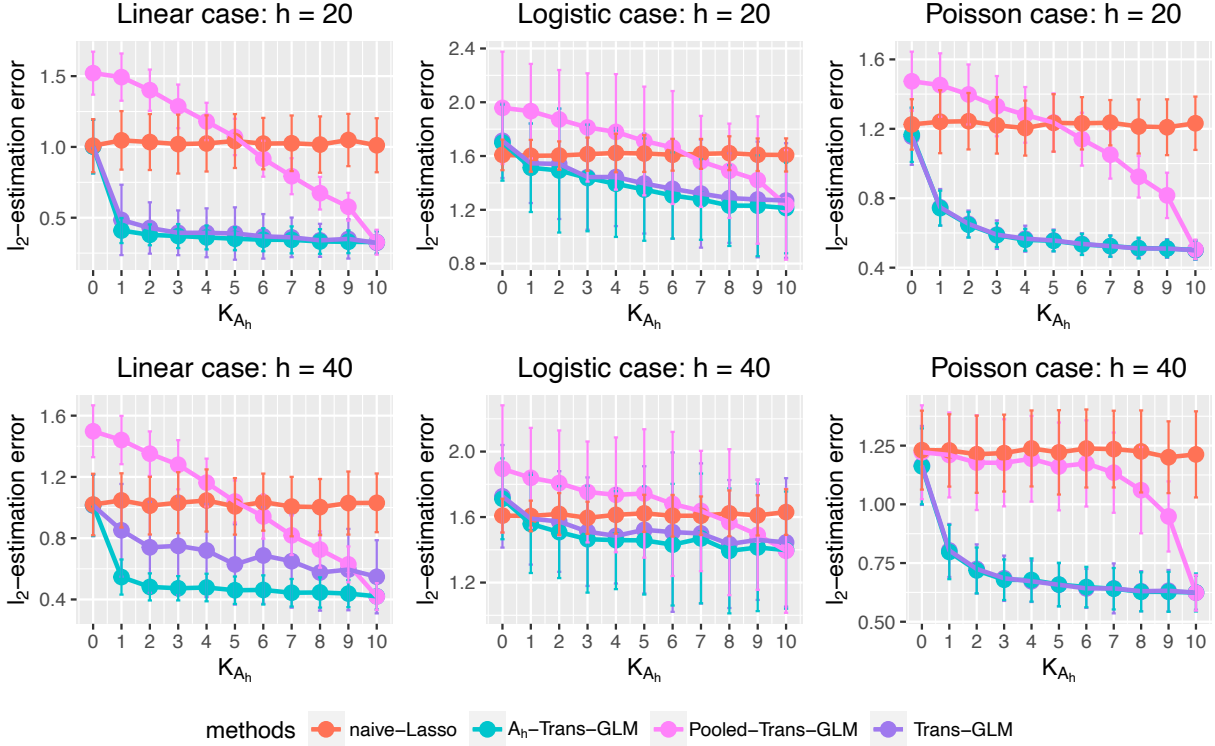


Figure 9: The average ℓ_2 -estimation error of various models with different settings of h and K_{A_h} when $K = 10$. $n_k = 150$ for all $k = 0, \dots, K$, $p = 1000$, $s = 15$. Error bars denote the standard deviations.

between the accuracy and computational costs.

S.1.3.4 Additional results of real data analysis

In this section, we aim to identify variables with significant effects for different targets in the real-data study (Section 4.2), by applying Algorithm 3. Taking the randomness caused by the cross-validation procedure in algorithms, we repeat the experiment 500 times. In

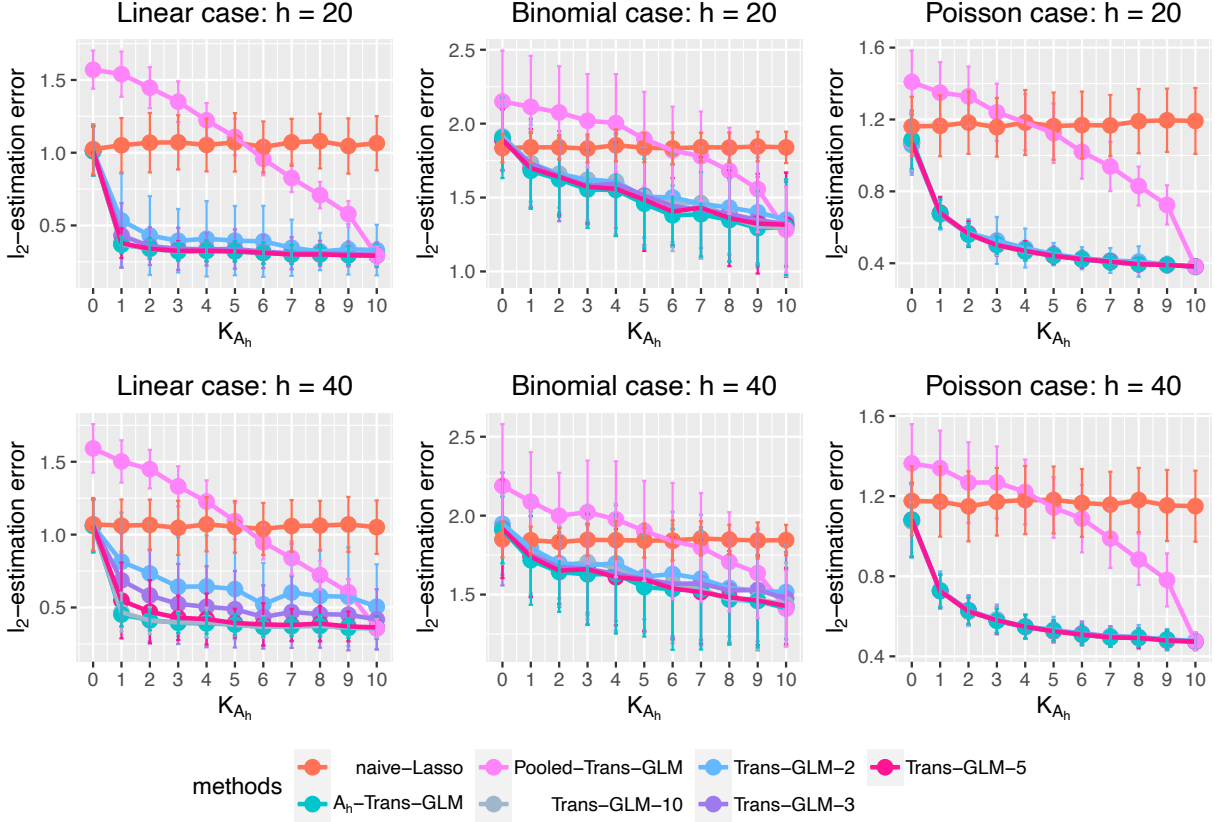


Figure 10: The average ℓ_2 -estimation error of various models with different settings of h and K_{A_h} when $K = 10$. $n_k = 200$ for all $k = 0, \dots, K$, $p = 2000$, $s = 20$. Error bars denote the standard deviations. The numbers after Trans-GLM indicate the number of folds used in cross-validation procedure of Algorithm 2 (steps 2-5)

each replication, for each target state, we first run Algorithm 2 to get a point estimate $\hat{\beta}$ of the coefficient and the estimated informative source index set $\hat{\mathcal{A}}$. Then, we run Algorithm 3 with $\hat{\beta}$ and $\hat{\mathcal{A}}$ to get the significant variables under 90% confidence level. Equivalently speaking, we identify coefficient components whose 90% confidence interval (CI) does not

cover zero and divided them into two parts based on the sign of CI center. Recall our recoding rule of the response: 0 denotes Republicans and 1 denotes Democrats. Therefore, given all other variables fixed, increasing the variable with a positive CI center gives rise to the chance of a county to vote Democrats. In contrast, increasing the variable with a negative CI center gives rise to the chance of a county to vote Republicans. The results are summarized in Figure 11, where we list all variables which are significant under 90% level in at least 5% of 500 replications. We use different colors and shapes to distinguish the effect of these variables (the sign of their CI centers). The x -axis and y -axis display the target state ¹ and the variable names. The description of some main effects is presented in Table 2. It reveals that RHI825214 is positively significant in 6 of 8 target states, which means that when other variables are fixed, a larger White percentage in a county leads to a higher chance to vote Republicans. In opposition to this, EDU685213 is negatively significant in 4 of 8 target states, showing that when other variables are fixed, a larger Bachelor or higher degree holder percentage benefits Democrats under county-level. More interesting findings can be obtained from Figure 11, which are expected to provide some insights to better understand the election results.

S.2 Proofs

Define $\hat{\mathbf{u}}^{\mathcal{A}_h} = \hat{\mathbf{w}}^{\mathcal{A}_h} - \mathbf{w}^{\mathcal{A}_h}$ and $\mathcal{D} = \{(\mathbf{X}^{(k)}, \mathbf{y}^{(k)})\}_{k \in \{0\} \cup \mathcal{A}_h}$. In the following, we will use bolded $\boldsymbol{\psi}'$ to represent the vector whose each component comes from the scalar function

¹Target state AR is removed because none of the variables are significant in 5% of 500 replications

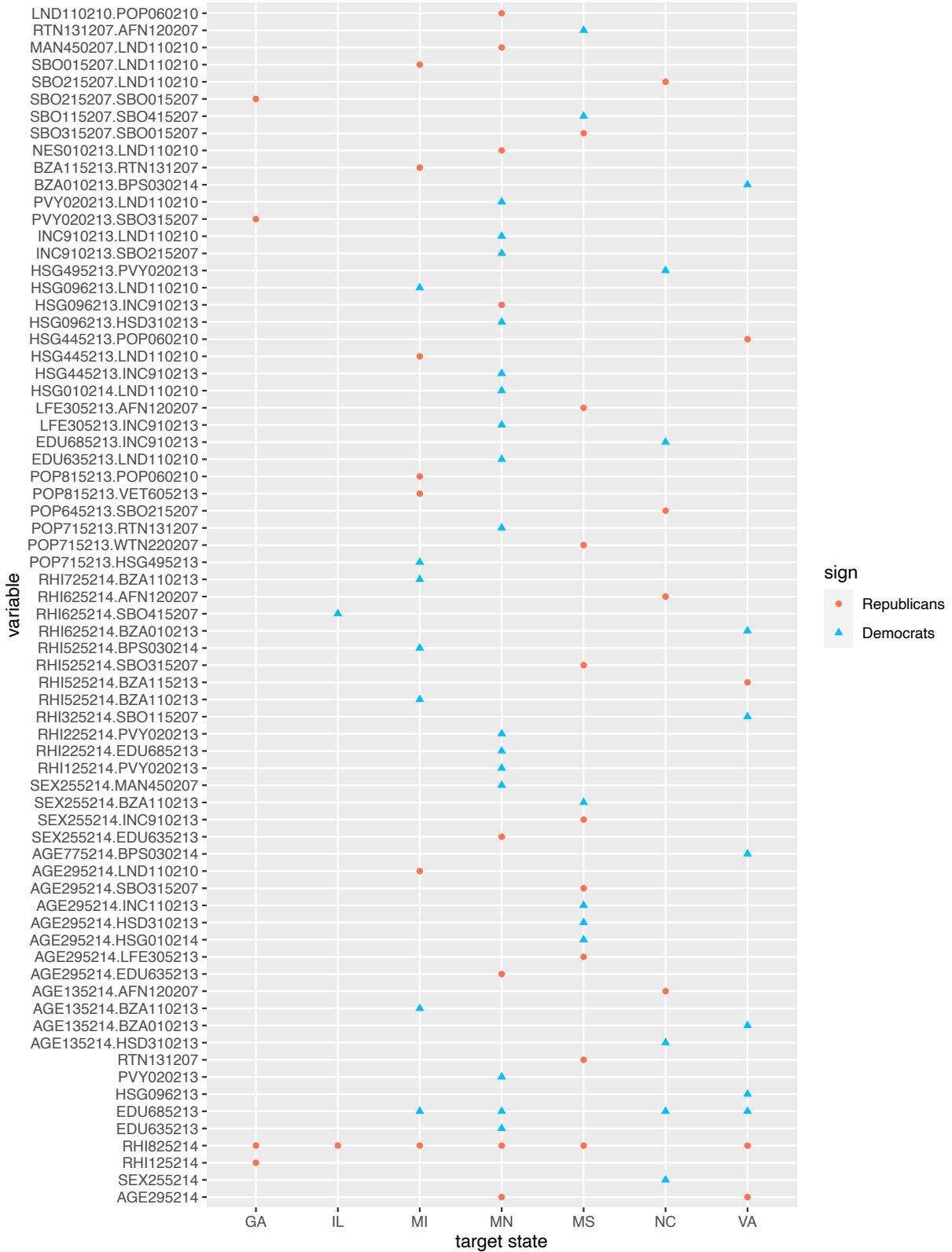


Figure 11: Variables that are significant under 90% level confidence for different targets in at least 5% of 500 replications. Provided by Algorithms 2 and 3.

| Variable name | Description |
|---------------|---|
| AGE135214 | Persons under 5 years, percent, 2014 |
| AGE295214 | Persons under 18 years, percent, 2014 |
| AGE775214 | Persons 65 years and over, percent, 2014 |
| SEX255214 | Female persons, percent, 2014 |
| RHI125214 | White alone, percent, 2014 |
| RHI225214 | Black or African American alone, percent, 2014 |
| RHI325214 | American Indian and Alaska Native alone, percent, 2014 |
| RHI525214 | Native Hawaiian and Other Pacific Islander alone, percent, 2014 |
| RHI625214 | Two or More Races, percent, 2014 |
| RHI725214 | Hispanic or Latino, percent, 2014 |
| RHI825214 | White alone, not Hispanic or Latino, percent, 2014 |
| POP715213 | Living in same house 1 year & over, percent, 2009-2013 |
| POP645213 | Foreign born persons, percent, 2009-2013 |
| POP815213 | Language other than English spoken at home, pct age 5+, 2009-2013 |
| EDU635213 | High school graduate or higher, percent of persons age 25+, 2009-2013 |
| EDU685213 | Bachelor's degree or higher, percent of persons age 25+, 2009-2013 |
| VET605213 | Veterans, 2009-2013 |
| LFE305213 | Mean travel time to work (minutes), workers age 16+, 2009-2013 |
| HSG010214 | Housing units, 2014 |
| HSG445213 | Homeownership rate, 2009-2013 |
| HSG096213 | Housing units in multi-unit structures, percent, 2009-2013 |
| HSG495213 | Median value of owner-occupied housing units, 2009-2013 |
| HSD310213 | Persons per household, 2009-2013 |
| INC910213 | Per capita money income in past 12 months (2013 dollars), 2009-2013 |
| INC110213 | Median household income, 2009-2013 |
| PVY020213 | Persons below poverty level, percent, 2009-2013 |
| BZA010213 | Private nonfarm establishments, 2013 |
| BZA110213 | Private nonfarm employment, 2013 |
| BZA115213 | Private nonfarm employment, percent change, 2012-2013 |
| NES010213 | Nonemployer establishments, 2013 |
| SBO315207 | Black-owned firms, percent, 2007 |
| SBO115207 | American Indian- and Alaska Native-owned firms, percent, 2007 |
| SBO215207 | Asian-owned firms, percent, 2007 |
| SBO415207 | Hispanic-owned firms, percent, 2007 |
| SBO015207 | Women-owned firms, percent, 2007 |
| MAN450207 | Manufacturers shipments, 2007 (\$1,000) |
| WTN220207 | Merchant wholesaler sales, 2007 (\$1,000) |
| RTN131207 | Retail sales per capita, 2007 |
| AFN120207 | Accommodation and food services sales, 2007 (\$1,000) |
| BPS030214 | Building permits, 2014 |
| LND110210 | Land area in square miles, 2010 |
| POP060210 | Population per square mile, 2010 |

Table 2: Description of some variables in the original dataset.

ψ' with corresponding predictors. Denote

$$\begin{aligned}\hat{L}(\mathbf{w}, \mathcal{D}) &= -\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{y}^{(k)})^T \mathbf{X}^{(k)} \mathbf{w} + \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} \psi(\mathbf{w}^T \mathbf{x}_i^{(k)}), \\ \nabla \hat{L}(\mathbf{w}, \mathcal{D}) &= -\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}^{(k)})^T \mathbf{y}^{(k)} + \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}^{(k)})^T \psi'(\mathbf{w}^T \mathbf{x}^{(k)}), \\ \delta \hat{L}(\mathbf{u}, \mathcal{D}) &= \hat{L}(\mathbf{w}^{\mathcal{A}_h} + \mathbf{u}, \mathcal{D}) - \hat{L}(\mathbf{w}^{\mathcal{A}_h}) - \nabla \hat{L}(\mathbf{w}^{\mathcal{A}_h})^T \mathbf{u}.\end{aligned}$$

Denote $\partial \|\mathbf{w}\|_1$ as the subgradient of $\|\mathbf{w}\|_1$ w.r.t. $\mathbf{w} \in \mathbb{R}^p$, which falls between -1 and 1 .

For any $\mathbf{w} \in \mathbb{R}^n$, denote $\mathbf{W}_{\mathbf{w}}^{(k)} = \text{diag} \left(\sqrt{\psi''((\mathbf{x}_1^{(k)})^T \mathbf{w})}, \dots, \sqrt{\psi''((\mathbf{x}_{n_k}^{(k)})^T \mathbf{w})} \right)$ and $\mathbf{X}_{\mathbf{w}}^{(k)} = \mathbf{W}_{\mathbf{w}}^{(k)} \mathbf{X}^{(k)}$. $\mathbf{X}_{\mathbf{w},j}^{(k)}$ represents the j -th column of $\mathbf{X}_{\mathbf{w}}^{(k)}$ and $\mathbf{X}_{\mathbf{w},-j}^{(k)}$ represents the matrix $\mathbf{X}_{\mathbf{w}}^{(k)}$ without the j -th column. Denote $\mathbf{X}_{\mathbf{w},-j}^{(k)}$ as the submatrix without j -th column. $\mathbf{X}_{\mathbf{w},j}^{(k)}$ represents the j -th column of $\mathbf{X}_{\mathbf{w}}^{(k)}$ without the diagonal (j, j) elements. Denote $\mathbf{x}_{\mathbf{w}}^{(k)} = \sqrt{\psi''((\mathbf{x}^{(k)})^T \mathbf{w})} \cdot \mathbf{x}^{(k)}$. $\mathbf{x}_{\mathbf{w},j}^{(k)}$ and $\mathbf{x}_{\mathbf{w},-j}^{(k)}$ represent the j -th component of $\mathbf{x}_{\mathbf{w}}^{(k)}$ and the vector without j -th component, respectively. Define

$$\gamma_j^{(k)} = \arg \min_{\gamma \in \mathbb{R}^{p-1}} \mathbb{E} \left[\mathbf{x}_{\mathbf{w}^{(k)},j}^{(k)} - (\mathbf{x}_{\mathbf{w}^{(k)},-j}^{(k)})^T \gamma \right]^2 = \arg \min_{\gamma \in \mathbb{R}^{p-1}} \mathbb{E} \left\{ \psi''((\mathbf{w}^{(k)})^T \mathbf{x}^{(k)}) \cdot [\mathbf{x}_j^{(k)} - (\mathbf{x}_{-j}^{(k)})^T \gamma]^2 \right\}.$$

Also define $\boldsymbol{\eta}_{\mathbf{w}^{(k)},j}^{(k)} = \mathbf{X}_{\mathbf{w},j}^{(k)} - \mathbf{X}_{\mathbf{w},-j}^{(k)} \gamma_j^{(k)}$ and $(\tau_j^{(k)})^2 = \mathbb{E}(\boldsymbol{\eta}_{\mathbf{w}^{(k)},j}^{(k)})^2$. And define

$$\gamma_j^{\mathcal{A}} = \arg \min_{\gamma \in \mathbb{R}^{p-1}} \left\{ \sum_{k \in \{0\} \cup \mathcal{A}_h} \mathbb{E} \left[\psi''((\mathbf{w}^{(k)})^T \mathbf{x}^{(k)}) \cdot [\mathbf{x}_j^{(k)} - (\mathbf{x}_{-j}^{(k)})^T \gamma]^2 \right] \right\}.$$

S.2.1 Some lemmas

Lemma 1. *Under Assumptions 1 and 4,*

$$\|\boldsymbol{\delta}^{\mathcal{A}_h}\|_1 = \|\mathbf{w}^{\mathcal{A}_h} - \boldsymbol{\beta}\|_1 \leq C_1 h,$$

where $\mathbf{w}^{\mathcal{A}_h}$ is defined by equation (2) and $C_1 := \sup_{k \in \{0\} \cup \mathcal{A}_h} \|\tilde{\boldsymbol{\Sigma}}_h^{-1} \tilde{\boldsymbol{\Sigma}}_h^{(k)}\|_1 < \infty$.

Lemma 2. Under Assumptions 1 and 2, there exists some positive constants κ_1 , κ_2 , C_3 and C_4 such that,

$$\delta\hat{L}(\hat{\mathbf{u}}^{\mathcal{A}_h}, \mathcal{D}) \geq \kappa_1 \|\mathbf{u}\|_2^2 - \kappa_1 \kappa_2 \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} \|\mathbf{u}\|_1 \|\mathbf{u}\|_2, \quad \forall \mathbf{u} : \|\mathbf{u}\|_2 \leq 1$$

with probability at least $1 - C_3 \exp\{-C_4(n_{\mathcal{A}_h} + n_0)\}$.

Lemma 3. Under Assumptions 1-3 and Assumption 7.(i)-(iii),

$$\sup_{k,j} \|\gamma_j^{\mathcal{A}_h} - \gamma_j^{(0)}\|_1 \lesssim h_1.$$

Lemma 4. Assume Assumptions 1-4 and 7 hold. Let $\lambda_j \asymp \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + (\frac{1}{\sqrt{s}}(\frac{\log p}{n_0})^{1/4} h^{1/2}) \wedge \frac{h}{\sqrt{s}}$ for any $j = 1, \dots, p$. Suppose $h_1 \lesssim s^{-1/2}$. With probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$,

$$\begin{aligned} \|\hat{\gamma}_j^{\mathcal{A}} - \gamma_j^{\mathcal{A}}\|_2^2 &\lesssim h_1 \left(\frac{\log p}{n_{\mathcal{A}_h} + n_0} \right)^{1/2} + \mathfrak{R}_1^2 + \mathfrak{R}_1 \frac{h_1}{\sqrt{s}}, \\ \|\hat{\gamma}_j^{\mathcal{A}} - \gamma_j^{\mathcal{A}}\|_1 &\lesssim \sqrt{s} \mathfrak{R}_1 + s^{1/4} h_1^{1/2} \mathfrak{R}_1^{1/2} + h_1. \end{aligned}$$

Lemma 5. Let $\tilde{\lambda}_j \asymp \sqrt{\frac{\log p}{n_0}} + \mathfrak{R}_1$. Impose the same conditions assumed by Lemma 4. Suppose $h_1 \lesssim s_*^{-1/2}$. Then with probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$,

$$\begin{aligned} \|\hat{\boldsymbol{\rho}}_j - \boldsymbol{\rho}_j\|_2^2 &\lesssim h_1 \sqrt{\frac{\log p}{n_0}} + h_1 \mathfrak{R}_1 + \mathfrak{R}_1^2, \\ \|\hat{\boldsymbol{\rho}}_j - \boldsymbol{\rho}_j\|_1 &\lesssim h_1 + \mathfrak{R}_1. \end{aligned}$$

The proof of Lemma 1 will be presented later in Section S.2.2.1. Lemma 2 can be derived in the same spirit as the proof of Proposition 2 in the full-length version of Negahban et al. (2009), so we omit the full proof and only highlight the sketch in Section S.2.2.2. Lemma 3 can be proved by following the same idea in the proof of Lemma 1, therefore we omit its proof as well.

Also, it is important to point out that all the constants involved in the proofs of Theorem 1-7 are independent with $\boldsymbol{\xi} = \{\boldsymbol{\beta}, \{\mathbf{w}^{(k)}\}_{k \in \mathcal{A}}\} \in \Xi(s, h)$, therefore we can take the supremum over $\boldsymbol{\xi} \in \Xi(s, h)$ in the final conclusion without changing the rate.

S.2.2 Proof of lemmas

S.2.2.1 Proof of Lemma 1

By definition,

$$\sum_{k \in \{0\} \cup \mathcal{A}} \alpha_k \mathbb{E} \{ [\psi'((\mathbf{w}^{\mathcal{A}_h})^T \mathbf{x}^{(k)}) - \psi'((\mathbf{w}^{(k)})^T \mathbf{x}^{(k)})] \mathbf{x}^{(k)} \} = 0,$$

which implies

$$\begin{aligned} & \sum_{k \in \{0\} \cup \mathcal{A}} \alpha_k \mathbb{E} \{ [\psi'((\mathbf{w}^{\mathcal{A}_h})^T \mathbf{x}^{(k)}) - \psi'(\boldsymbol{\beta}^T \mathbf{x}^{(k)})] \mathbf{x}^{(k)} \} \\ &= \sum_{k \in \mathcal{A}} \alpha_k \mathbb{E} \{ [\psi'((\mathbf{w}^{(k)})^T \mathbf{x}^{(k)}) - \psi'(\boldsymbol{\beta}^T \mathbf{x}^{(k)})] \mathbf{x}^{(k)} \}. \end{aligned}$$

By Taylor expansion,

$$\begin{aligned} & \sum_{k \in \{0\} \cup \mathcal{A}} \alpha_k \mathbb{E} \left[\int_0^1 \psi''((\mathbf{w}^{\mathcal{A}_h})^T \mathbf{x}^{(k)} + t(\mathbf{w}^{\mathcal{A}_h} - \boldsymbol{\beta})^T \mathbf{x}^{(k)}) \mathbf{x}^{(k)} (\mathbf{x}^{(k)})^T \right] (\mathbf{w}^{\mathcal{A}_h} - \boldsymbol{\beta}) \\ &= \sum_{k \in \mathcal{A}} \alpha_k \mathbb{E} \left[\int_0^1 \psi''((\mathbf{w}^{(k)})^T \mathbf{x}^{(k)} + t(\mathbf{w}^{(k)} - \boldsymbol{\beta})^T \mathbf{x}^{(k)}) \mathbf{x}^{(k)} (\mathbf{x}^{(k)})^T \right] (\mathbf{w}^{(k)} - \boldsymbol{\beta}). \end{aligned}$$

Therefore, by Assumption 4, $\|\mathbf{w}^{\mathcal{A}_h} - \boldsymbol{\beta}\|_1 \leq \sum_{k \in \mathcal{A}} \alpha_k \|\tilde{\boldsymbol{\Sigma}}_h^{-1} \tilde{\boldsymbol{\Sigma}}_h^{(k)}\|_1 \cdot \|\mathbf{w}^{(k)} - \boldsymbol{\beta}\|_1 \leq C_1 h$.

S.2.2.2 Proof of Lemma 2

By the second-order Taylor expansion, for some $t_i^{(k)} \in [0, 1]$,

$$\delta \hat{L}(\mathbf{u}, \mathcal{D}) = \hat{L}(\mathbf{w}^{\mathcal{A}_h} + \mathbf{u}, \mathcal{D}) - \hat{L}(\mathbf{w}^{\mathcal{A}_h}) - \nabla \hat{L}^{(k)}(\mathbf{w}^{\mathcal{A}_h})^T \mathbf{u}$$

$$\begin{aligned}
&= \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} \left[\psi((\mathbf{w}^{\mathcal{A}_h} + \mathbf{u})^T \mathbf{x}_i^{(k)}) - \psi((\mathbf{w}^{\mathcal{A}_h})^T \mathbf{x}_i^{(k)}) - \psi'((\mathbf{w}^{\mathcal{A}_h})^T \mathbf{x}_i^{(k)}) \mathbf{u}^T \mathbf{x}_i^{(k)} \right] \\
&= \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} \psi''((\mathbf{w}^{\mathcal{A}_h})^T \mathbf{x}_i^{(k)} + t_i^{(k)} \mathbf{u}^T \mathbf{x}_i^{(k)}) (\mathbf{u}^T \mathbf{x}_i^{(k)})^2,
\end{aligned}$$

which is the counterpart of equation (63) in the full-length version of [Negahban et al. \(2009\)](#). Due to the independence of between $\mathbf{x}_i^{(k)}$ for any i and k , the arguments in [Negahban et al. \(2009\)](#) directly follow.

S.2.2.3 Proof of Lemma 4

Recall Theorem 1, under the assumptions, we have

$$\begin{aligned}
\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_2 &\lesssim \sqrt{\frac{s \log p}{n_{\mathcal{A}} + n_0}} + \left[\left(\frac{\log p}{n_0} \right)^{1/4} h^{1/2} \right] \wedge h, \\
\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_1 &\lesssim s \sqrt{\frac{\log p}{n_{\mathcal{A}} + n_0}} + h,
\end{aligned}$$

with probability at least $1 - n_0^{-1}$. By basic inequality,

$$\begin{aligned}
&\frac{1}{2(n_{\mathcal{A}_h} + n_0)} \sum_{k \in \{0\} \cup \mathcal{A}_h} \|\mathbf{X}_{\hat{\boldsymbol{\beta}}, j}^{(k)} - \mathbf{X}_{\boldsymbol{\beta}, -j}^{(k)} \hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h}\|_2^2 + \lambda_j \|\hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h}\|_1 \\
&\leq \frac{1}{2(n_{\mathcal{A}_h} + n_0)} \sum_{k \in \{0\} \cup \mathcal{A}_h} \|\mathbf{X}_{\hat{\boldsymbol{\beta}}, j}^{(k)} - \mathbf{X}_{\boldsymbol{\beta}, -j}^{(k)} \boldsymbol{\gamma}_j^{\mathcal{A}_h}\|_2^2 + \lambda_j \|\boldsymbol{\gamma}_j^{\mathcal{A}_h}\|_1,
\end{aligned}$$

which implies

$$\begin{aligned}
&\frac{1}{2(n_{\mathcal{A}_h} + n_0)} \sum_{k \in \{0\} \cup \mathcal{A}_h} \|\mathbf{X}_{\hat{\boldsymbol{\beta}}, -j}^{(k)} (\hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h} - \boldsymbol{\gamma}_j^{\mathcal{A}_h})\|_2^2 \\
&\leq \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \langle \mathbf{W}_{\hat{\boldsymbol{\beta}}}^{(k)} (\mathbf{W}_{\boldsymbol{\beta}}^{(k)})^{-1} \boldsymbol{\eta}_{\boldsymbol{\beta}, j}^{(k)}, \mathbf{X}_{\hat{\boldsymbol{\beta}}, -j}^{(k)} (\boldsymbol{\gamma}_j^{\mathcal{A}_h} - \hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h}) \rangle
\end{aligned}$$

$$\begin{aligned}
& + \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \langle \mathbf{X}_{\hat{\beta}, -j}^{(k)} (\boldsymbol{\gamma}_j^{(k)} - \boldsymbol{\gamma}_j^{\mathcal{A}_h}), \mathbf{X}_{\hat{\beta}, -j}^{(k)} (\boldsymbol{\gamma}_j^{\mathcal{A}_h} - \hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h}) \rangle \\
& + \lambda_j (\|\boldsymbol{\gamma}_j^{\mathcal{A}_h}\|_1 - \|\hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h}\|_1). \tag{S.2.12}
\end{aligned}$$

Denote $\boldsymbol{\Lambda}^{(k)} = \text{diag}(\{\psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta})\}_{i=1}^{n_k})$ and $\hat{\boldsymbol{\Lambda}}^{(k)} = \text{diag}(\{\psi''((\mathbf{x}_i^{(k)})^T \hat{\boldsymbol{\beta}})\}_{i=1}^{n_k})$.

$$\begin{aligned}
& \left| \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \langle \mathbf{X}_{\hat{\beta}, -j}^{(k)} (\boldsymbol{\gamma}_j^{(k)} - \boldsymbol{\gamma}_j^{\mathcal{A}_h}), \mathbf{X}_{\hat{\beta}, -j}^{(k)} (\boldsymbol{\gamma}_j^{\mathcal{A}_h} - \hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h}) \rangle \right| \\
& \leq \frac{1}{n_{\mathcal{A}_h} + n_0} \left[\left\| \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}_{-j}^{(k)})^T \boldsymbol{\Lambda}^{(k)} \mathbf{X}_{-j}^{(k)} (\boldsymbol{\gamma}_j^{(k)} - \boldsymbol{\gamma}_j^{\mathcal{A}_h}) \right\|_{\infty} \right. \\
& \quad \left. + \left\| \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}_{-j}^{(k)})^T (\hat{\boldsymbol{\Lambda}}^{(k)} - \boldsymbol{\Lambda}^{(k)}) \mathbf{X}_{-j}^{(k)} (\boldsymbol{\gamma}_j^{(k)} - \boldsymbol{\gamma}_j^{\mathcal{A}_h}) \right\|_{\infty} \right] \cdot \|\boldsymbol{\gamma}_j^{\mathcal{A}_h} - \hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h}\|_1. \tag{S.2.13}
\end{aligned}$$

It's easy to see that each component of $\sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}_{-j}^{(k)})^T \boldsymbol{\Lambda}^{(k)} \mathbf{X}_{-j}^{(k)} (\boldsymbol{\gamma}_j^{(k)} - \boldsymbol{\gamma}_j^{\mathcal{A}_h})$ is a zero-mean, sub-exponential variable with variance $C(n_{\mathcal{A}} + n_0)h^2$. By union bounds,

$$\frac{1}{n_{\mathcal{A}_h} + n_0} \left\| \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}_{-j}^{(k)})^T \boldsymbol{\Lambda}^{(k)} \mathbf{X}_{-j}^{(k)} (\boldsymbol{\gamma}_j^{(k)} - \boldsymbol{\gamma}_j^{\mathcal{A}_h}) \right\|_{\infty} \lesssim \sqrt{\frac{\log p}{n_{\mathcal{A}} + n_0}} h_1,$$

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

$$\begin{aligned}
& \frac{1}{n_{\mathcal{A}_h} + n_0} \left\| \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}_{-j}^{(k)})^T (\hat{\boldsymbol{\Lambda}}^{(k)} - \boldsymbol{\Lambda}^{(k)}) \mathbf{X}_{-j}^{(k)} (\boldsymbol{\gamma}_j^{(k)} - \boldsymbol{\gamma}_j^{\mathcal{A}_h}) \right\|_{\infty} \\
& \leq \frac{1}{n_{\mathcal{A}_h} + n_0} \left\| \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}_{-j}^{(k)})^T (\hat{\boldsymbol{\Lambda}}^{(k)} - \boldsymbol{\Lambda}^{(k)}) \mathbf{X}_{-j}^{(k)} \right\|_{\max} \cdot \|\boldsymbol{\gamma}_j^{(k)} - \boldsymbol{\gamma}_j^{\mathcal{A}_h}\|_1.
\end{aligned}$$

For $j_1, j_2 \neq j$,

$$\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} x_{i,j_1}^{(k)} x_{i,j_2}^{(k)} [\psi''((\mathbf{x}_i^{(k)})^T \hat{\boldsymbol{\beta}}) - \psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta})]$$

$$\leq \sqrt{\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} (x_{i,j_1}^{(k)} x_{i,j_2}^{(k)})^2} \cdot \sqrt{\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} \psi'''(\mathbf{a}_i^{(k)}) [(\mathbf{x}_i^{(k)})^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})]^2}$$

$$\leq C \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_2,$$

where the constant C is the same for different j_1 and j_2 , and $\mathbf{a}_i^{(k)}$ falls on the line segment between $(\mathbf{x}_i^{(k)})^T \hat{\boldsymbol{\beta}}$ and $(\mathbf{x}_i^{(k)})^T \boldsymbol{\beta}$. Therefore, the right-hand side of (S.2.13) can be upper bounded by

$$C \left(\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_2 + \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} \right) h_1 \cdot \|\gamma_j^{\mathcal{A}_h} - \hat{\gamma}_j^{\mathcal{A}_h}\|_1,$$

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

On the other hand,

$$\begin{aligned} & \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \langle \mathbf{W}_{\hat{\boldsymbol{\beta}}}^{(k)} (\mathbf{W}_{\boldsymbol{\beta}}^{(k)})^{-1} \boldsymbol{\eta}_{\boldsymbol{\beta},j}^{(k)}, \mathbf{X}_{\hat{\boldsymbol{\beta}},-j}^{(k)} (\gamma_j^{\mathcal{A}_h} - \hat{\gamma}_j^{\mathcal{A}_h}) \rangle \\ & \leq \frac{1}{n_{\mathcal{A}_h} + n_0} \left\| \sum_{k \in \{0\} \cup \mathcal{A}_h} \mathbf{X}_{\hat{\boldsymbol{\beta}},-j}^{(k)} \boldsymbol{\eta}_{\boldsymbol{\beta},j}^{(k)} \right\|_{\infty} \|\gamma_j^{\mathcal{A}_h} - \hat{\gamma}_j^{\mathcal{A}_h}\|_1 \\ & \quad + \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \left[\frac{1}{4c_0} \left\| (\mathbf{W}_{\hat{\boldsymbol{\beta}}}^{(k)} (\mathbf{W}_{\boldsymbol{\beta}}^{(k)})^{-1} - \mathbf{I}) \boldsymbol{\eta}_{\boldsymbol{\beta},j}^{(k)} \right\|_2^2 \right. \\ & \quad \left. + c_0 \left\| \mathbf{X}_{\hat{\boldsymbol{\beta}},-j}^{(k)} (\gamma_j^{\mathcal{A}_h} - \hat{\gamma}_j^{\mathcal{A}_h}) \right\|_2^2 \right], \end{aligned} \tag{S.2.14}$$

where c_0 is a positive constant smaller than $1/4$. Note that

$$\begin{aligned} & \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \left\| (\mathbf{W}_{\hat{\boldsymbol{\beta}}}^{(k)} (\mathbf{W}_{\boldsymbol{\beta}}^{(k)})^{-1} - \mathbf{I}) \boldsymbol{\eta}_{\boldsymbol{\beta},j}^{(k)} \right\|_2^2 \\ & \leq \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} \frac{\left(\sqrt{\psi''((\mathbf{x}_i^{(k)})^T \hat{\boldsymbol{\beta}})} - \sqrt{\psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta})} \right)^2}{\psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta})} \cdot |\eta_i^{(k)}|^2 \end{aligned}$$

$$\begin{aligned}
&\lesssim \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} [(\mathbf{x}_i^{(k)})^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})]^2, \\
&\lesssim \sup_{k \in \{0\} \cup \mathcal{A}_h} \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_2^2 \\
&\lesssim \mathfrak{R}_1^2,
\end{aligned}$$

with probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$, where we use Assumption 7.(i) to bound $\psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta})$. The last two inequalities follows because of the same reason as that of (S.2.27). Plugging (S.2.13) and (S.2.14) and into (S.2.12),

$$\begin{aligned}
&\frac{1}{4(n_{\mathcal{A}_h} + n_0)} \sum_{k \in \{0\} \cup \mathcal{A}_h} \|\mathbf{X}_{\hat{\boldsymbol{\beta}}, -j}^{(k)} (\hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h} - \boldsymbol{\gamma}_j^{\mathcal{A}_h})\|_2^2 \\
&\leq C \left(\sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + \mathfrak{R}_1 \right) h_1 \cdot \|\hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h} - \boldsymbol{\gamma}_j^{\mathcal{A}_h}\|_1 + \mathfrak{R}_1^2 \\
&\quad + \frac{C}{n_{\mathcal{A}_h} + n_0} \left\| \sum_{k \in \{0\} \cup \mathcal{A}_h} \mathbf{X}_{\hat{\boldsymbol{\beta}}, -j}^{(k)} \boldsymbol{\eta}_{\boldsymbol{\beta}, j}^{(k)} \right\|_{\infty} \|\boldsymbol{\gamma}_j^{\mathcal{A}_h} - \hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h}\|_1 + \lambda_j (\|\boldsymbol{\gamma}_j^{\mathcal{A}_h}\|_1 - \|\hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h}\|_1)
\end{aligned}$$

with probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$. Note that $\frac{1}{n_{\mathcal{A}_h} + n_0} \left\| \sum_{k \in \{0\} \cup \mathcal{A}_h} \mathbf{X}_{\hat{\boldsymbol{\beta}}, -j}^{(k)} \right\|_{\infty} \lesssim \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}}$ with probability at least $1 - p^{-1}$. Therefore, if $\lambda_j \geq C' \mathfrak{R}_1 / \sqrt{s}$ with some large $C' > 0$, then since $h_1 \lesssim s^{-1/2}$, we have $\lambda_j > \frac{2C}{n_{\mathcal{A}_h} + n_0} \left\| \sum_{k \in \{0\} \cup \mathcal{A}_h} \mathbf{X}_{\hat{\boldsymbol{\beta}}, -j}^{(k)} \right\|_{\infty} + C \left(\sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + \mathfrak{R}_1 \right) h_1$, which implies

$$\begin{aligned}
\frac{1}{4(n_{\mathcal{A}_h} + n_0)} \sum_{k \in \{0\} \cup \mathcal{A}_h} \|\mathbf{X}_{\hat{\boldsymbol{\beta}}, -j}^{(k)} (\hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h} - \boldsymbol{\gamma}_j^{\mathcal{A}_h})\|_2^2 &\leq \frac{3}{2} \lambda_j \|(\hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h} - \boldsymbol{\gamma}_j^{\mathcal{A}_h})_{S_j}\|_1 - \frac{1}{2} \lambda_j \|(\hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h} - \boldsymbol{\gamma}_j^{\mathcal{A}_h})_{S_j^c}\|_1 \\
&\quad + \lambda_j \|(\boldsymbol{\gamma}_j^{\mathcal{A}_h})_{S_j^c}\|_1 + C \mathfrak{R}_1^2, \tag{S.2.15}
\end{aligned}$$

with probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$. Therefore,

$$\|\hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h} - \boldsymbol{\gamma}_j^{\mathcal{A}_h}\|_1 \leq 4\sqrt{s} \|\hat{\boldsymbol{\gamma}}_j^{\mathcal{A}_h} - \boldsymbol{\gamma}_j^{\mathcal{A}_h}\|_2 + C\sqrt{s} \mathfrak{R}_1 + h_1, \tag{S.2.16}$$

with probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$. Thus, we have either $\|\hat{\gamma}_j^{\mathcal{A}_h} - \gamma_j^{\mathcal{A}_h}\|_1 \leq 2C\sqrt{s}\mathfrak{R}_1 + 2h_1 \lesssim \sqrt{s}$ or $\|\hat{\gamma}_j^{\mathcal{A}_h} - \gamma_j^{\mathcal{A}_h}\|_1 \leq 8\sqrt{s}\|\hat{\gamma}_j^{\mathcal{A}_h} - \gamma_j^{\mathcal{A}_h}\|_2$. By similar arguments to get (S.2.27), we obtain

$$\begin{aligned} \frac{1}{4(n_{\mathcal{A}_h} + n_0)} \sum_{k \in \{0\} \cup \mathcal{A}_h} \|\mathbf{X}_{\hat{\beta}, -j}^{(k)}(\hat{\gamma}_j^{\mathcal{A}_h} - \gamma_j^{\mathcal{A}_h})\|_2^2 &\gtrsim \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \|\mathbf{X}_{-j}^{(k)}(\hat{\gamma}_j^{\mathcal{A}_h} - \gamma_j^{\mathcal{A}_h})\|_2^2 \\ &\gtrsim \|\hat{\gamma}_j^{\mathcal{A}_h} - \gamma_j^{\mathcal{A}_h}\|_2^2, \end{aligned} \quad (\text{S.2.17})$$

with probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$. Combine (S.2.15), (S.2.16) and (S.2.17) to get the desired conclusions.

S.2.2.4 Proof of Lemma 5

Similar to (S.2.12), by basic inequality,

$$\begin{aligned} \frac{1}{2n_0} \left\| \mathbf{X}_{\hat{\beta}, -j}^{(0)}(\hat{\boldsymbol{\rho}}_j - \boldsymbol{\rho}_j) \right\|_2^2 &\leq \frac{1}{n_0} \langle \mathbf{W}_{\hat{\beta}}^{(0)}(\mathbf{W}_{\beta}^{(0)})^{-1} \boldsymbol{\eta}_{\beta, j}^{(0)}, \mathbf{X}_{\hat{\beta}, -j}^{(0)}(\hat{\boldsymbol{\rho}}_j - \boldsymbol{\rho}_j) \rangle \\ &\quad + \frac{1}{n_0} \langle \mathbf{X}_{\hat{\beta}, -j}^{(0)}(\hat{\boldsymbol{\rho}}_j - \boldsymbol{\rho}_j), \mathbf{X}_{\hat{\beta}, -j}^{(0)}(\gamma_j^{\mathcal{A}_h} - \hat{\gamma}_j^{\mathcal{A}_h}) \rangle \\ &\quad + \tilde{\lambda}_j (\|\boldsymbol{\rho}_j\|_1 - \|\hat{\boldsymbol{\rho}}_j\|_1). \end{aligned}$$

Similar to the analysis in the proof of Lemma 4,

$$\begin{aligned} \frac{1}{4n_0} \left\| \mathbf{X}_{\hat{\beta}, -j}^{(0)}(\hat{\boldsymbol{\rho}}_j - \boldsymbol{\rho}_j) \right\|_2^2 &\leq C \left(\sqrt{\frac{\log p}{n_0}} + \mathfrak{R}_1 \right) h_1 \cdot \|\hat{\boldsymbol{\rho}}_j - \boldsymbol{\rho}_j\|_1 + C \|\hat{\gamma}_j^{\mathcal{A}_h} - \gamma_j^{\mathcal{A}_h}\|_2^2 \\ &\quad + C \|\hat{\beta} - \beta\|_2^2 + \frac{C}{n_0} \left\| \mathbf{X}_{\beta, -j}^{(0)} \boldsymbol{\eta}_j^{(0)} \right\|_{\infty} \cdot \|\hat{\boldsymbol{\rho}}_j - \boldsymbol{\rho}_j\|_1 \\ &\quad + \tilde{\lambda}_j (\|\boldsymbol{\rho}_j\|_1 - \|\hat{\boldsymbol{\rho}}_j\|_1) \end{aligned}$$

with probability at least $1 - n_0^{-1}$. Note that $\frac{1}{n_0} \|\mathbf{X}_{\beta, -j}^{(0)} \boldsymbol{\eta}_j^{(0)}\|_{\infty} \lesssim \sqrt{\frac{\log p}{n_0}}$ with probability at least $1 - p^{-1}$. Therefore for $\tilde{\lambda}_j \geq C' \left(\sqrt{\frac{\log p}{n_0}} + \mathfrak{R}_1 \right)$ with large enough $C' > 0$, we have

$\tilde{\lambda}_j \geq C \left(\sqrt{\frac{\log p}{n_0}} + \mathfrak{R}_1 \right) h_1 + \frac{C}{n_0} \|\mathbf{X}_{\beta, -j}^{(0)} \boldsymbol{\eta}_j^{(0)}\|_\infty$. Then with probability at least $1 - n_0^{-1}$,

$$\frac{1}{4n_0} \left\| \mathbf{X}_{\hat{\beta}, -j}^{(0)} (\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j) \right\|_2^2 \leq 2\tilde{\lambda}_j \|\boldsymbol{\theta}_j\|_1 - \frac{1}{2}\tilde{\lambda}_j \|\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j\|_1 + C \|\hat{\boldsymbol{\gamma}}_j^{A_h} - \boldsymbol{\gamma}_j^{A_h}\|_2^2 + C \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_2^2,$$

which leads to

$$\begin{aligned} \|\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j\|_1 &\leq 4\|\boldsymbol{\theta}_j\|_1 + C\tilde{\lambda}_j^{-1} \|\hat{\boldsymbol{\gamma}}_j^{A_h} - \boldsymbol{\gamma}_j^{A_h}\|_2^2 + C\tilde{\lambda}_j^{-1} \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_2^2 \\ &\lesssim h_1 + \mathfrak{R}_1 + \left(\sqrt{\frac{n_0}{\log p}} \mathfrak{R}_1^2 \right) \wedge \mathfrak{R}_1, \end{aligned}$$

with probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$. Similar to the trick we used before, we can get

$$\|\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j\|_2^2 \lesssim \tilde{\lambda}_j h_1 + \|\hat{\boldsymbol{\gamma}}_j^{A_h} - \boldsymbol{\gamma}_j^{A_h}\|_2^2 + \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_2^2 \lesssim h_1 \sqrt{\frac{\log p}{n_0}} + h_1 \mathfrak{R}_1 + \mathfrak{R}_1^2$$

with probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$, which completes the proof.

S.2.3 Proof of theorems

S.2.3.1 Proof of Theorem 1

Transferring step: Define $\hat{\mathbf{u}}^{A_h} = \hat{\mathbf{w}}^{A_h} - \mathbf{w}^{A_h}$ and $\mathcal{D} = \{(\mathbf{X}^{(k)}, \mathbf{y}^{(k)})\}_{k \in \{0\} \cup \mathcal{A}_h}$. We first claim that when $\lambda_{\mathbf{w}} \geq 2\|\nabla L(\mathbf{w}^{A_h}, \mathcal{D})\|_\infty$, with probability at least $1 - C_3 \exp\{-C_4(n_{\mathcal{A}_h} + n_0)\}$, it holds that

$$\|\hat{\mathbf{u}}^{A_h}\|_2 \leq 8\kappa_2 C_1 h \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + 3 \frac{\sqrt{s}}{\kappa_1} \lambda_{\mathbf{w}} + 2 \sqrt{\frac{C_1}{\kappa_1}} h \lambda_{\mathbf{w}}. \quad (\text{S.2.18})$$

To see this, first by the definition of $\hat{\mathbf{w}}^{A_h}$, Hölder inequality and Lemma 1, we have

$$\begin{aligned} \delta \hat{L}(\hat{\mathbf{u}}^{A_h}, \mathcal{D}) &\leq \lambda_{\mathbf{w}} (\|\mathbf{w}_S^{A_h}\|_1 + \|\mathbf{w}_{S^c}^{A_h}\|_1) - \lambda_{\mathbf{w}} (\|\hat{\mathbf{w}}_S^{A_h}\|_1 + \|\hat{\mathbf{w}}_{S^c}^{A_h}\|_1) + \nabla \hat{L}(\mathbf{w}, \mathcal{D})^T \hat{\mathbf{u}}^{A_h} \\ &\leq \lambda_{\mathbf{w}} (\|\mathbf{w}_S^{A_h}\|_1 + \|\mathbf{w}_{S^c}^{A_h}\|_1) - \lambda_{\mathbf{w}} (\|\hat{\mathbf{w}}_S^{A_h}\|_1 + \|\hat{\mathbf{w}}_{S^c}^{A_h}\|_1) + \frac{1}{2} \lambda_{\mathbf{w}} \|\hat{\mathbf{u}}^{A_h}\|_1 \end{aligned}$$

$$\begin{aligned}
&\leq \frac{3}{2}\lambda_w\|\hat{\mathbf{u}}_S^{\mathcal{A}_h}\|_1 - \frac{1}{2}\lambda_w\|\hat{\mathbf{u}}_{S^c}^{\mathcal{A}_h}\|_1 + 2\lambda_w\|\mathbf{w}_{S^c}^{\mathcal{A}_h}\|_1 \\
&\leq \frac{3}{2}\lambda_w\|\hat{\mathbf{u}}_S^{\mathcal{A}_h}\|_1 - \frac{1}{2}\lambda_w\|\hat{\mathbf{u}}_{S^c}^{\mathcal{A}_h}\|_1 + 2\lambda_w C_1 h.
\end{aligned} \tag{S.2.19}$$

If the claim does not hold, consider $\mathbf{C} = \{\mathbf{u} : \frac{3}{2}\|\mathbf{u}_S\|_1 - \frac{1}{2}\|\mathbf{u}_{S^c}\|_1 + 2C_1 h \geq 0\}$. Due to (S.2.19) and the convexity of \hat{L} , $\hat{\mathbf{u}}^{\mathcal{A}_h} \in \mathbf{C}$. Then for any $t \in (0, 1)$, it's easy to see that

$$\frac{1}{2}\|t\hat{\mathbf{u}}_{S^c}^{\mathcal{A}_h}\|_1 = t \cdot \frac{1}{2}\|\hat{\mathbf{u}}_{S^c}^{\mathcal{A}_h}\|_1 \leq t \cdot \left(\frac{3}{2}\|\hat{\mathbf{u}}_S^{\mathcal{A}_h}\|_1 + 2C_w h \right) \leq \frac{3}{2}\|t\hat{\mathbf{u}}_S^{\mathcal{A}_h}\|_1 + 2C_1 h,$$

implying that $t\hat{\mathbf{u}}^{\mathcal{A}_h} \in \mathbf{C}$. We could find some t satisfying that $\|t\hat{\mathbf{u}}^{\mathcal{A}_h}\|_2 > 8\kappa_2 C_1 h \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + 3\frac{\sqrt{s}}{\kappa_1}\lambda_w + 2\sqrt{\frac{C_1}{\kappa_1}h\lambda_w}$ and $\|t\hat{\mathbf{u}}^{\mathcal{A}_h}\|_2 \leq 1$. Denote $\tilde{\mathbf{u}}^{\mathcal{A}_h} = t\hat{\mathbf{u}}^{\mathcal{A}_h}$ and $F(\mathbf{u}) = \hat{L}(\mathbf{w}^{\mathcal{A}_h} + \mathbf{u}, \mathcal{D}) - \hat{L}(\mathbf{w}^{\mathcal{A}_h}) + \lambda_w(\|\mathbf{w}^{\mathcal{A}_h} + \mathbf{u}\|_1 - \|\mathbf{w}^{\mathcal{A}_h}\|_1)$. Since $F(\mathbf{0}) = 0$ and $F(\hat{\mathbf{u}}^{\mathcal{A}_h}) \leq 0$, by convexity,

$$F(\tilde{\mathbf{u}}^{\mathcal{A}_h}) = F(t\hat{\mathbf{u}}^{\mathcal{A}_h} + (1-t)\mathbf{0}) \leq tF(\hat{\mathbf{u}}^{\mathcal{A}_h}) \leq 0. \tag{S.2.20}$$

However, by Lemma 2 and the same trick of (S.2.19),

$$\begin{aligned}
F(\tilde{\mathbf{u}}^{\mathcal{A}_h}) &\geq \delta\hat{L}(\hat{\mathbf{u}}^{\mathcal{A}_h}, \mathcal{D}) + \nabla\hat{L}(\mathbf{w}^{\mathcal{A}_h})^T\tilde{\mathbf{u}}^{\mathcal{A}_h} - \lambda_w\|\mathbf{w}^{\mathcal{A}_h}\|_1 + \lambda_w\|\mathbf{w}^{\mathcal{A}_h} + \tilde{\mathbf{u}}^{\mathcal{A}_h}\|_1 \\
&\geq \kappa_1\|\tilde{\mathbf{u}}^{\mathcal{A}_h}\|_2^2 - \kappa_1\kappa_2\sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}}\|\tilde{\mathbf{u}}^{\mathcal{A}_h}\|_1\|\tilde{\mathbf{u}}^{\mathcal{A}_h}\|_2 - \frac{3}{2}\lambda_w\|\tilde{\mathbf{u}}_S^{\mathcal{A}_h}\|_1 + \frac{1}{2}\lambda_w\|\tilde{\mathbf{u}}_{S^c}^{\mathcal{A}_h}\|_1 - 2\lambda_w C_1 h \\
&\geq \kappa_1\|\tilde{\mathbf{u}}^{\mathcal{A}_h}\|_2^2 - \kappa_1\kappa_2\sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}}\|\tilde{\mathbf{u}}^{\mathcal{A}_h}\|_1\|\tilde{\mathbf{u}}^{\mathcal{A}_h}\|_2 - \frac{3}{2}\lambda_w\|\tilde{\mathbf{u}}_S^{\mathcal{A}_h}\|_1 - 2\lambda_w C_1 h.
\end{aligned}$$

Note that since $\tilde{\mathbf{u}}^{\mathcal{A}_h} \in \mathbf{C}$, it holds that

$$\frac{1}{2}\|\tilde{\mathbf{u}}^{\mathcal{A}_h}\|_1 \leq 2\|\tilde{\mathbf{u}}_S^{\mathcal{A}_h}\|_1 + 2C_w h \leq 2\sqrt{s}\|\tilde{\mathbf{u}}^{\mathcal{A}_h}\|_2 + 2C_1 h.$$

When $n_{\mathcal{A}_h} + n_0 > 16\kappa_2^2 s \log p$, we have $2\kappa_2\sqrt{\frac{s \log p}{n_{\mathcal{A}_h} + n_0}} \leq \frac{1}{2}$. Then it follows

$$F(\tilde{\mathbf{u}}^{\mathcal{A}_h}) \geq \frac{1}{2}\kappa_1\|\tilde{\mathbf{u}}^{\mathcal{A}_h}\|_2^2 - \left[2\kappa_1\kappa_2\sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}}C_w h + \frac{3}{2}\lambda_w\sqrt{s} \right] \|\tilde{\mathbf{u}}^{\mathcal{A}_h}\|_2 - 2\lambda_w C_1 h > 0,$$

which conflicts with (S.2.20). Therefore our claim at the beginning holds.

Next, let's prove $\|\nabla \hat{L}(\mathbf{w}^{\mathcal{A}_h})\|_\infty \lesssim \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}}$ with probability at least $1 - (n_{\mathcal{A}_h} + n_0)^{-1}$.

To see this, notice that

$$\begin{aligned} \nabla \hat{L}(\mathbf{w}^{\mathcal{A}_h}) &= \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}^{(k)})^T [-\mathbf{y}^{(k)} + \boldsymbol{\psi}'(\mathbf{X}^{(k)} \mathbf{w}^{\mathcal{A}_h})] \\ &= \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}^{(k)})^T [-\mathbf{y}^{(k)} + \boldsymbol{\psi}'(\mathbf{X}^{(k)} \boldsymbol{\beta})] \\ &\quad + \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}^{(k)})^T [-\boldsymbol{\psi}'(\mathbf{X}^{(k)} \boldsymbol{\beta}) + \boldsymbol{\psi}'(\mathbf{X}^{(k)} \mathbf{w}^{\mathcal{A}_h})]. \end{aligned} \quad (\text{S.2.21})$$

Following a similar idea in the proof of Lemma 6 in [Negahban et al. \(2009\)](#), under Assumptions 1-3 and the fact $n_{\mathcal{A}_h} \geq Cs \log p$, we can show that

$$\frac{1}{n_{\mathcal{A}_h} + n_0} \left\| \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}^{(k)})^T [-\mathbf{y}^{(k)} + \boldsymbol{\psi}'(\mathbf{X}^{(k)} \boldsymbol{\beta})] \right\|_\infty \lesssim \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}},$$

with probability at least $1 - (n_{\mathcal{A}_h} + n_0)^{-1}$.

The remaining step is to bound the infinity norm of the second term in (S.2.21). Denote $V_{ij}^{(k)} = x_{ij}^{(k)} [-\boldsymbol{\psi}'((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta}) + \boldsymbol{\psi}'((\mathbf{x}_i^{(k)})^T \mathbf{w}^{\mathcal{A}_h})]$. Under Assumption 3, by mean value theorem and Lemma 1, there exists $v_i^{(k)} \in (0, 1)$ such that

$$\begin{aligned} &\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} V_{ij}^{(k)} \\ &= \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} \boldsymbol{\psi}''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta} + v_i^{(k)} (\mathbf{x}_i^{(k)})^T (\mathbf{w}^{\mathcal{A}_h} - \boldsymbol{\beta})) x_{ij}^{(k)} (\mathbf{x}_i^{(k)})^T (\mathbf{w}^{\mathcal{A}_h} - \boldsymbol{\beta}). \end{aligned}$$

$\boldsymbol{\psi}''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta} + v_i^{(k)} (\mathbf{x}_i^{(k)})^T (\mathbf{w}^{\mathcal{A}_h} - \boldsymbol{\beta})) x_{ij}^{(k)}$ is $M_\psi^2 \kappa_u^2$ -subGaussian due to the almost sure boundedness of $\boldsymbol{\psi}''$ in Assumption 3. And $(\mathbf{x}_i^{(k)})^T (\mathbf{w}^{\mathcal{A}_h} - \boldsymbol{\beta})$ is a $4C_1^2 h^2$ -subGaussian due to Lemma 1. Then the multiplication is a $4C_1^2 M_\psi^2 \kappa_u^2 h^2$ -subexponential variable. By definition

of $\mathbf{w}^{\mathcal{A}_h}$, $\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} V_{ij}^{(k)}$ has zero mean. Notice that the infinity norm of the second term in (S.2.21) equals to $\frac{1}{n_{\mathcal{A}_h} + n_0} \sup_{j=1, \dots, p} |\sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} V_{ij}^{(k)}|$, by tail bounds of subexponential variables and union bounds, we have

$$\frac{1}{n_{\mathcal{A}_h} + n_0} \sup_{j=1, \dots, p} \left| \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} V_{ij}^{(k)} \right| \lesssim C_1 M_{\psi, \kappa_u} h \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}},$$

with probability at least $1 - (n_{\mathcal{A}_h} + n_0)^{-1}$. Therefore $\|\nabla \hat{L}(\mathbf{w}^{\mathcal{A}_h})\|_{\infty} \lesssim \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}}$ holds with probability at least $1 - (n_{\mathcal{A}_h} + n_0)^{-1}$. Plugging the rate into (S.2.18), we have

$$\|\hat{\mathbf{u}}^{\mathcal{A}_h}\|_2 \lesssim h \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + \sqrt{\frac{s \log p}{n_{\mathcal{A}_h} + n_0}} + \left(\frac{\log p}{n_{\mathcal{A}_h} + n_0}\right)^{1/4} \sqrt{h}, \quad (\text{S.2.22})$$

with probability at least $1 - (n_{\mathcal{A}_h} + n_0)^{-1}$, when $\lambda_{\mathbf{w}} \asymp C_{\mathbf{w}} \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}}$ with $C_{\mathbf{w}} > 0$ sufficiently large. Since $\hat{\mathbf{u}}^{\mathcal{A}_h} \in \mathbb{C}$, (S.2.22) encloses

$$\|\hat{\mathbf{u}}^{\mathcal{A}_h}\|_1 \lesssim s \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + \left(\frac{\log p}{n_{\mathcal{A}_h} + n_0}\right)^{1/4} \sqrt{sh} + h \left(1 + \sqrt{\frac{s \log p}{n_{\mathcal{A}_h} + n_0}}\right), \quad (\text{S.2.23})$$

with probability at least $1 - (n_{\mathcal{A}_h} + n_0)^{-1}$.

Debiasing step: Denote $\mathcal{D}^{(0)} = (\mathbf{X}^{(0)}, \mathbf{y}^{(0)})$, $\hat{L}^{(0)}(\mathbf{w}, \mathcal{D}^{(0)}) = -\frac{1}{n_0} (\mathbf{y}^{(0)})^T \mathbf{X}^{(0)} \mathbf{w} + \frac{1}{n_0} \sum_{i=1}^{n_0} \psi((\mathbf{x}_i^{(0)})^T \mathbf{w})$, $\nabla \hat{L}^{(0)}(\mathbf{w}, \mathcal{D}^{(0)}) = -\frac{1}{n_0} (\mathbf{X}^{(0)})^T \mathbf{y}^{(0)} + \frac{1}{n_0} (\mathbf{X}^{(0)})^T \psi'(\mathbf{X}^{(0)} \mathbf{w})$, $\boldsymbol{\delta}^{\mathcal{A}_h} = \boldsymbol{\beta} - \mathbf{w}^{\mathcal{A}_h}$, $\hat{\boldsymbol{\beta}} = \hat{\mathbf{w}}^{\mathcal{A}_h} + \hat{\boldsymbol{\delta}}^{\mathcal{A}_h}$, $\hat{\mathbf{v}}^{\mathcal{A}_h} = \hat{\boldsymbol{\delta}}^{\mathcal{A}_h} - \boldsymbol{\delta}^{\mathcal{A}_h}$, and $\delta \hat{L}^{(0)}(\boldsymbol{\delta}, \mathcal{D}) = \hat{L}^{(0)}(\hat{\mathbf{w}}^{\mathcal{A}_h} + \boldsymbol{\delta}, \mathcal{D}^{(0)}) - \hat{L}^{(0)}(\hat{\mathbf{w}}^{\mathcal{A}_h} + \boldsymbol{\delta}^{\mathcal{A}_h}, \mathcal{D}^{(0)}) - \nabla \hat{L}^{(0)}(\hat{\mathbf{w}}^{\mathcal{A}_h} + \boldsymbol{\delta}^{\mathcal{A}_h}, \mathcal{D}^{(0)})^T \hat{\mathbf{v}}^{\mathcal{A}_h}$.

Similar to (S.2.19), when $\lambda_{\boldsymbol{\delta}} \geq 2 \|\nabla \hat{L}^{(0)}(\boldsymbol{\beta}, \mathcal{D}^{(0)})\|_{\infty}$, we have

$$\begin{aligned} \delta \hat{L}^{(0)}(\hat{\boldsymbol{\delta}}^{\mathcal{A}_h}, \mathcal{D}) &\leq \lambda_{\boldsymbol{\delta}} (\|\boldsymbol{\delta}^{\mathcal{A}_h}\|_1 - \|\hat{\boldsymbol{\delta}}^{\mathcal{A}_h}\|_1) - \nabla \hat{L}^{(0)}(\hat{\mathbf{w}}^{\mathcal{A}_h} + \boldsymbol{\delta}^{\mathcal{A}_h}, \mathcal{D}^{(0)})^T \hat{\mathbf{v}}^{\mathcal{A}_h} \\ &\leq \lambda_{\boldsymbol{\delta}} (2 \|\boldsymbol{\delta}^{\mathcal{A}_h}\|_1 - \|\hat{\mathbf{v}}^{\mathcal{A}_h}\|_1) + \|\nabla \hat{L}^{(0)}(\boldsymbol{\beta}, \mathcal{D}^{(0)})\|_{\infty} \|\hat{\mathbf{v}}^{\mathcal{A}_h}\|_1 \\ &\quad - \left[\nabla \hat{L}^{(0)}(\hat{\mathbf{w}}^{\mathcal{A}_h} + \boldsymbol{\delta}^{\mathcal{A}_h}, \mathcal{D}^{(0)}) - \nabla \hat{L}^{(0)}(\boldsymbol{\beta}, \mathcal{D}^{(0)}) \right]^T \hat{\mathbf{v}}^{\mathcal{A}_h} \end{aligned}$$

$$\begin{aligned}
&\leq 2\lambda_\delta \|\boldsymbol{\delta}^{\mathcal{A}_h}\|_1 - \frac{1}{2}\lambda_\delta \|\hat{\boldsymbol{v}}^{\mathcal{A}_h}\|_1 \\
&\quad - \frac{1}{n_0} \left[\boldsymbol{\psi}'((\mathbf{X}^{(0)})^T(\hat{\boldsymbol{w}}^{\mathcal{A}_h} + \boldsymbol{\delta}^{\mathcal{A}_h})) - \boldsymbol{\psi}'((\mathbf{X}^{(0)})^T\boldsymbol{\beta}) \right]^T \hat{\boldsymbol{v}}^{\mathcal{A}_h} \\
&\leq 2\lambda_\delta \|\boldsymbol{\delta}^{\mathcal{A}_h}\|_1 - \frac{1}{2}\lambda_\delta \|\hat{\boldsymbol{v}}^{\mathcal{A}_h}\|_1 + \frac{1}{4c_0}M_\psi^2 \cdot \frac{1}{n_0} \|\mathbf{X}^{(0)}\hat{\boldsymbol{u}}^{\mathcal{A}_h}\|_2^2 + c_0 \cdot \frac{1}{n_0} \|\mathbf{X}^{(0)}\hat{\boldsymbol{v}}^{\mathcal{A}_h}\|_2^2,
\end{aligned} \tag{S.2.24}$$

where $c_0 > 0$ is an enough small constant. The last inequality holds because

$$\begin{aligned}
& - \frac{1}{n_0} \left[\boldsymbol{\psi}'((\mathbf{X}^{(0)})^T(\hat{\boldsymbol{w}}^{\mathcal{A}_h} + \boldsymbol{\delta}^{\mathcal{A}_h})) - \boldsymbol{\psi}'((\mathbf{X}^{(0)})^T\boldsymbol{\beta}) \right]^T \hat{\boldsymbol{v}}^{\mathcal{A}_h} \\
&= \frac{1}{n_0} (\hat{\boldsymbol{u}}^{\mathcal{A}_h})^T (\mathbf{X}^{(0)})^T \Lambda^{(0)} \mathbf{X}^{(0)} \hat{\boldsymbol{v}}^{\mathcal{A}_h} \\
&\leq \frac{1}{4c_0}M_\psi^2 \cdot \frac{1}{n_0} \|\mathbf{X}^{(0)}\hat{\boldsymbol{u}}^{\mathcal{A}_h}\|_2^2 + c_0 \cdot \frac{1}{n_0} \|\mathbf{X}^{(0)}\hat{\boldsymbol{v}}^{\mathcal{A}_h}\|_2^2,
\end{aligned}$$

where $\Lambda^{(0)} = \text{diag}(\{\boldsymbol{\psi}''((\boldsymbol{x}_i^{(0)})^T\boldsymbol{\beta} + t_i(\boldsymbol{x}_i^{(0)})^T\hat{\boldsymbol{u}}^{\mathcal{A}_h})\}_{i=1}^{n_0})$ is a $n_0 \times n_0$ diagonal matrix and by Assumption 3, $\|\Lambda^{(0)}\|_{\max} \leq M_\psi$. Denote $\tilde{\boldsymbol{v}}^{\mathcal{A}_h} = t\hat{\boldsymbol{v}}^{\mathcal{A}_h}$ and $F^{(0)}(\boldsymbol{v}) = \hat{L}^{(0)}(\hat{\boldsymbol{w}}^{\mathcal{A}_h} + \boldsymbol{\delta}^{\mathcal{A}_h} + \boldsymbol{v}, \mathcal{D}^{(0)}) - \hat{L}^{(0)}(\hat{\boldsymbol{w}}^{\mathcal{A}_h} + \boldsymbol{\delta}^{\mathcal{A}_h}, \mathcal{D}^{(0)}) + \lambda_\delta(\|\boldsymbol{\delta}^{\mathcal{A}_h} + \boldsymbol{v}\|_1 - \|\boldsymbol{\delta}^{\mathcal{A}_h}\|_1)$. Since $F(\mathbf{0}) = 0$ and $F^{(0)}(\hat{\boldsymbol{v}}^{\mathcal{A}_h}) \leq 0$, by convexity, for any $t \in (0, 1]$,

$$F^{(0)}(\tilde{\boldsymbol{v}}^{\mathcal{A}_h}) = F^{(0)}(t\hat{\boldsymbol{v}}^{\mathcal{A}_h} + (1-t)\mathbf{0}) \leq tF^{(0)}(\hat{\boldsymbol{v}}^{\mathcal{A}_h}) \leq 0. \tag{S.2.25}$$

We set $t \in (0, 1]$ such that $\|\tilde{\boldsymbol{v}}^{\mathcal{A}_h}\|_2 \leq 1$, which allows us to apply Proposition 1 in [Loh and Wainwright \(2015\)](#) on $\tilde{\boldsymbol{v}}^{\mathcal{A}_h}$. By basic inequality (S.2.25) and the same arguments in (S.2.24),

$$\begin{aligned}
c_1 \|\tilde{\boldsymbol{v}}^{\mathcal{A}_h}\|_2^2 - c_2 \cdot \frac{\log p}{n_0} \|\tilde{\boldsymbol{v}}^{\mathcal{A}_h}\|_1^2 &\leq F^{(0)}(\tilde{\boldsymbol{v}}^{\mathcal{A}_h}) - \nabla \hat{L}^{(0)}(\hat{\boldsymbol{w}}^{\mathcal{A}_h} + \boldsymbol{\delta}^{\mathcal{A}_h}, \mathcal{D}^{(0)})^T \tilde{\boldsymbol{v}}^{\mathcal{A}_h} \\
&\leq 2\lambda_\delta \|\boldsymbol{\delta}^{\mathcal{A}_h}\|_1 - \frac{1}{2}\lambda_\delta \|\tilde{\boldsymbol{v}}^{\mathcal{A}_h}\|_1 + \frac{1}{4c_0}M_\psi^2 \cdot \frac{1}{n_0} \|\mathbf{X}^{(0)}\hat{\boldsymbol{u}}^{\mathcal{A}_h}\|_2^2 \\
&\quad + c_0 \cdot \frac{1}{n_0} \|\mathbf{X}^{(0)}\tilde{\boldsymbol{v}}^{\mathcal{A}_h}\|_2^2.
\end{aligned} \tag{S.2.26}$$

In the discussion of transferring step, we showed that $\|\hat{\mathbf{u}}_{S^c}^{\mathcal{A}_h}\|_1 \leq 3\|\hat{\mathbf{u}}_S^{\mathcal{A}_h}\|_1 + 4C_1h$. We leverage on this fact to bound $\frac{1}{n_0}\|\mathbf{X}^{(0)}\hat{\mathbf{u}}^{\mathcal{A}_h}\|_2^2$ by $\|\mathbf{X}^{(0)}\hat{\mathbf{u}}^{\mathcal{A}_h}\|_2^2$.

If $3\|\hat{\mathbf{u}}_S^{\mathcal{A}_h}\|_1 \geq 4C_1h$, we will have $\|\hat{\mathbf{u}}_{S^c}^{\mathcal{A}_h}\|_1 \leq 6\|\hat{\mathbf{u}}_S^{\mathcal{A}_h}\|_1$. Then by Theorem 1.6 of [Zhou \(2009\)](#),

$$\frac{1}{n_0}\|\mathbf{X}^{(0)}\hat{\mathbf{u}}^{\mathcal{A}_h}\|_2^2 \lesssim \|\hat{\mathbf{u}}^{\mathcal{A}_h}\|_2^2 \lesssim \frac{s \log p}{n_{\mathcal{A}} + n_0} + h \cdot \sqrt{\frac{\log p}{n_{\mathcal{A}} + n_0}}, \quad (\text{S.2.27})$$

with probability at least $1 - C(n_{\mathcal{A}} + n_0)^{-1} - C \exp\{-n_0\}$.

If $3\|\hat{\mathbf{u}}_S^{\mathcal{A}_h}\|_1 < 4C_1h$, we will have $\|\hat{\mathbf{u}}_{S^c}^{\mathcal{A}_h}\|_1 \leq 8C_1h \leq \sqrt{s}$. Also $\|\hat{\mathbf{u}}^{\mathcal{A}_h}\|_2 \leq 1$ with probability $1 - C(n_{\mathcal{A}} + n_0)^{-1}$. Denote

$$\begin{aligned} \Pi_0(s) &= \{\mathbf{u} \in \mathbb{R}^p : \|\mathbf{u}\|_2 \leq 1, \|\mathbf{u}\|_0 \leq s\}, \\ \Pi_1(s) &= \{\mathbf{u} \in \mathbb{R}^p : \|\mathbf{u}\|_2 \leq 1, \|\mathbf{u}\|_1 \leq \sqrt{s}\}. \end{aligned}$$

By Lemma 3.1 of [Plan and Vershynin \(2013\)](#), $\Pi_1(s) \subseteq 2\overline{\text{conv}}(\Pi_0(s))$, where $\overline{\text{conv}}(\Pi_0(s))$ is the closure of convex hull of $\Pi_0(s)$. Then following a similar argument in the proof of Theorem 2.4 in [Mendelson et al. \(2008\)](#) leads to (S.2.27) as well.

Next we want to bound $\frac{1}{n_0}\|\mathbf{X}^{(0)}\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2^2$ by $\|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2^2$. By another basic inequality,

$$\begin{aligned} 0 &\leq \hat{L}^{(0)}(\hat{\boldsymbol{\beta}}, \mathcal{D}^{(0)}) - \hat{L}^{(0)}(\boldsymbol{\beta}, \mathcal{D}^{(0)}) - \nabla \hat{L}^{(0)}(\boldsymbol{\beta}, \mathcal{D}^{(0)})^T(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \\ &\leq \lambda_{\delta}(\|\boldsymbol{\beta} - \hat{\mathbf{w}}^{\mathcal{A}_h}\|_1 - \|\hat{\boldsymbol{\delta}}^{\mathcal{A}_h}\|_1) + \|\nabla \hat{L}^{(0)}(\boldsymbol{\beta}, \mathcal{D}^{(0)})\|_{\infty} \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_1 \\ &\leq \lambda_{\delta}(\|\boldsymbol{\beta} - \hat{\mathbf{w}}^{\mathcal{A}_h}\|_1 - \|\hat{\boldsymbol{\delta}}^{\mathcal{A}_h}\|_1) + \frac{1}{2}\lambda_{\delta}\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_1 \\ &\leq \frac{3}{2}\lambda_{\delta}\|\boldsymbol{\beta} - \hat{\mathbf{w}}^{\mathcal{A}_h}\|_1 - \frac{1}{2}\lambda_{\delta}\|\hat{\boldsymbol{\delta}}^{\mathcal{A}_h}\|_1 \\ &\leq \frac{3}{2}\lambda_{\delta}C_1h + \frac{3}{2}\lambda_{\delta}\|\hat{\mathbf{u}}^{\mathcal{A}}\|_1 - \frac{1}{2}\lambda_{\delta}\|\hat{\boldsymbol{\delta}}^{\mathcal{A}_h}\|_1, \end{aligned}$$

which implies

$$\|\hat{\mathbf{v}}^{\mathcal{A}_h}\|_1 \leq \|\hat{\boldsymbol{\delta}}^{\mathcal{A}_h}\|_1 + C_1h \leq 3\|\hat{\mathbf{u}}^{\mathcal{A}}\|_1 + 4C_1h.$$

Combined with (S.2.23), this leads to $\|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_1 \leq \|\hat{\mathbf{v}}^{\mathcal{A}_h}\|_1 \leq \sqrt{s}$ when $s \log p / (n_{\mathcal{A}_h} + n_0)$ and h are small enough. Due to the strict convexity, $\delta \hat{L}^{(0)}(\hat{\boldsymbol{\delta}}^{\mathcal{A}_h}, \mathcal{D}) > 0$, leading to $\|\hat{\boldsymbol{\delta}}^{\mathcal{A}_h}\|_1 \leq 3\|\hat{\mathbf{u}}^{\mathcal{A}_h}\|_1 + 3h$. Then,

$$\|\hat{\mathbf{v}}^{\mathcal{A}_h}\|_1 \leq \|\boldsymbol{\beta} - \hat{\mathbf{w}}^{\mathcal{A}_h}\|_1 + \|\hat{\boldsymbol{\delta}}^{\mathcal{A}_h}\|_1 \leq 4\|\hat{\mathbf{u}}^{\mathcal{A}_h}\|_1 + 4h \leq \sqrt{s}. \quad (\text{S.2.28})$$

Similar to the analysis of the case $3\|\hat{\mathbf{u}}_S^{\mathcal{A}_h}\|_1 < 4C_1h$ above, we can get

$$c_0 \cdot \frac{1}{n_0} \|\mathbf{X}^{(0)} \tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2^2 \leq c_0 \cdot C \|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2^2,$$

with probability at least $1 - C \exp\{-n_0\}$. As long as $c_0C < c_1/2$, by (S.2.26), we have

$$\begin{aligned} c_1 \|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2^2 - c_2 \cdot \frac{\log p}{n_0} \|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_1^2 &\leq 2\lambda_\delta \|\boldsymbol{\delta}^{\mathcal{A}_h}\|_1 - \frac{1}{2}\lambda_\delta \|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_1 + C \frac{s \log p}{n_{\mathcal{A}_h} + n_0} + Ch \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} \\ &\quad + c_1/2 \|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2^2 \end{aligned} \quad (\text{S.2.29})$$

with probability at least $1 - C'n_0^{-1}$.

If $\lambda_\delta \|\boldsymbol{\delta}^{\mathcal{A}_h}\|_1 \leq C \frac{s \log p}{n_{\mathcal{A}_h} + n_0} + Ch \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}}$:

$$\|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_1 \lesssim \left[\frac{s \log p}{n_{\mathcal{A}_h} + n_0} + h \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} \right] \cdot \sqrt{\frac{n_0}{\log p}} + \|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2^2.$$

Since $\|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2 \leq 1$, by (S.2.29), it holds

$$\|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2^2 \lesssim \frac{s \log p}{n_{\mathcal{A}_h} + n_0} + h \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} \lesssim \frac{s \log p}{n_{\mathcal{A}_h} + n_0} + \left[h \sqrt{\frac{\log p}{n_0}} \right] \wedge h^2,$$

with probability at least $1 - C'n_0^{-1}$.

If $\lambda_\delta \|\boldsymbol{\delta}^{\mathcal{A}_h}\|_1 > C \frac{s \log p}{n_{\mathcal{A}_h} + n_0} + Ch \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}}$: $\|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_1 \lesssim h + \|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2^2$, which leads to

$$\|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2^2 \lesssim 2\lambda_\delta \|\boldsymbol{\delta}^{\mathcal{A}_h}\|_1 - \frac{1}{2}\lambda_\delta \|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_1,$$

implying $\|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_1 \leq 4\|\boldsymbol{\delta}^{\mathcal{A}_h}\|_1 \leq 4C_1h$. Besides, by plugging this result into (S.2.29), we have

$$\|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2^2 \lesssim \frac{s \log p}{n_{\mathcal{A}_h} + n_0} + \left[h \sqrt{\frac{\log p}{n_0}} \right] \wedge h^2, \quad (\text{S.2.30})$$

with probability at least $1 - C'n_0^{-1}$.

When $s \log p / (n_{\mathcal{A}_h} + n_0)$ and h is small enough, and because $h \sqrt{\frac{\log p}{n_0}} = o(1)$, the right hand side of (S.2.30) can be very small. This means that we showed $\|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2 \leq c < 1$ with probability at least $1 - C'n_0^{-1}$. Note that this result holds for any $t \in (0, 1]$ such that $\|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2 \leq 1$. Now let's consider our interested vector $\hat{\mathbf{v}}^{\mathcal{A}_h}$. Suppose $\|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2^2 \geq C \frac{s \log p}{n_{\mathcal{A}_h} + n_0} + C \left[h \sqrt{\frac{\log p}{n_0}} \right] \wedge h^2$ for some constant $C > 0$ with probability at least $C'n_0^{-1}$, then we can choose $t \in (0, 1]$ such that $c < \|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2 \leq 1$, which is contradicted with the fact $\|\tilde{\mathbf{v}}^{\mathcal{A}_h}\|_2 \leq c$ with probability at least $1 - C'n_0^{-1}$. Therefore

$$\|\hat{\mathbf{v}}^{\mathcal{A}_h}\|_2^2 \lesssim \frac{s \log p}{n_{\mathcal{A}_h} + n_0} + \left[h \sqrt{\frac{\log p}{n_0}} \right] \wedge h^2 \lesssim 1,$$

with probability at least $1 - C'n_0^{-1}$. Then we can go over the analysis procedure of $\tilde{\mathbf{v}}^{\mathcal{A}_h}$ again with $\hat{\mathbf{v}}^{\mathcal{A}_h}$ (on the high-probability event $\|\hat{\mathbf{v}}^{\mathcal{A}_h}\|_2 \leq 1$) to get the ℓ_1 -bound

$$\|\hat{\mathbf{v}}^{\mathcal{A}_h}\|_1 \lesssim s \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + h,$$

with probability at least $1 - C'n_0^{-1}$.

Finally, we connect the conclusions above with the upper bounds on $\|\hat{\mathbf{u}}^{\mathcal{A}_h}\|_2$ and $\|\hat{\mathbf{u}}^{\mathcal{A}_h}\|_1$, which completes the proof.

S.2.3.2 Proof of Theorem 2

The analysis in the proof of Theorem B and Theorem 2 in Li et al. (2021), which proved the ℓ_2 minimax rate under linear models, can be directly extended to the GLMs. We only

point out some key observations. The readers can refer to the supplements of their paper for the full proof.

The proofs of Theorem B and Theorem 2 in Li et al. (2021) leverage on Fano's method, which involves Kullback-Leibler divergence. For \mathbf{w}_1 and \mathbf{w}_2 in \mathbb{R}^p , consider two GLMs where $y|\mathbf{x} \sim \mathbb{P}(y|\mathbf{x}; \mathbf{w}_1) = \rho(y) \exp\{y\mathbf{x}^T \mathbf{w}_1 - \psi(\mathbf{x}^T \mathbf{w}_1)\}$ and $y|\mathbf{x} \sim \mathbb{P}(y|\mathbf{x}; \mathbf{w}_2) = \rho(y) \exp\{y\mathbf{x}^T \mathbf{w}_2 - \psi(\mathbf{x}^T \mathbf{w}_2)\}$. Denote the corresponding joint distribution of (\mathbf{x}, y) as $f_{\mathbf{w}_1}$ and $f_{\mathbf{w}_2}$ and suppose marginal distributions of \mathbf{x} are the same. Then by definition and Assumption 3,

$$\begin{aligned} \text{KL}(f_{\mathbf{w}_1}||f_{\mathbf{w}_2}) &= \mathbb{E}_{\mathbf{x} \sim f_{\mathbf{w}_1}} [\psi''(\mathbf{x}^T(t\mathbf{w}_1 + (1-t)\mathbf{w}_2)) \cdot [\mathbf{x}^T(\mathbf{w}_1 - \mathbf{w}_2)]^2] \\ &\leq C \mathbb{E}_{\mathbf{x} \sim f_{\mathbf{w}_1}} [\mathbf{x}^T(\mathbf{w}_1 - \mathbf{w}_2)]^2, \end{aligned}$$

for some constant $C > 0$, if $\|\psi''\|_\infty < \infty$ or $\|\mathbf{x}\|_\infty \leq U$ a.s. and $\|\mathbf{w}_1 - \mathbf{w}_2\|_1 \leq C'$ with some $C' > 0$. By using this fact, all the analysis of ℓ_2 -estimation error in their proof works out for GLM.

About the ℓ_1 -estimation error, we can make slight changes to make the argument work again. In case (i) of the proof of Theorem 2 (or case (i) of Theorem B) in Li et al. (2021), replace the δ -packing under ℓ_2 norm with the δ -packing under ℓ_1 norm, and set $\delta_0 = c_0 s \sqrt{\frac{\log p}{n_{\mathcal{A}} + n_0}}$. In case (ii-1), we can do the same and set $\delta_0 = c_0 s \sqrt{\frac{\log p}{n_0}}$. In (ii-2), we suppose $s \sqrt{\log p / n_0} > h$. Set $\delta_0 = \bar{m} \sqrt{\frac{\log p}{n_0}} \asymp h$ where \bar{m} is the integer part of $h \sqrt{n_0 / \log p}$. In (ii-3), the same argument works and we set $\delta_0 = h$. So finally we can get the minimax ℓ_1 rate $s \sqrt{\log p / (n_{\mathcal{A}} + n_0)} + [s \sqrt{\log p / n_0}] \wedge h$.

S.2.3.3 Proof of Theorem 3

Throughout this proof, we denote $\mathbf{w}^{\mathcal{A}_h}$ as any vector \mathbf{w} satisfying $\|\mathbf{w} - \boldsymbol{\beta}\|_1 \leq h$. Such a $\mathbf{w}^{\mathcal{A}_h}$ indeed exists, e.g. $\mathbf{w}^{\mathcal{A}_h} = \sum_{k \in \{0\} \cup \mathcal{A}_h} \alpha_k \mathbf{w}^{(k)}$. Note that $\mathbf{w}^{\mathcal{A}_h}$ here does not necessarily

enjoy the moment condition (2), which will bring more bias. This is the price we have to pay for relaxing Assumption 4. Other notations are defined the same as in the proof of Theorem 1.

The main idea of the proof is similar to that in proof of Theorem 1. We only highlight the different parts here and do not dig into all the details.

First, the claim in (S.2.18) still holds here, i.e. when $\lambda_{\mathbf{w}} \geq 2\|\nabla\hat{L}(\mathbf{w}^{\mathcal{A}_h}, \mathcal{D})\|_{\infty}$, with probability at least $1 - C_3 \exp\{-C_4(n_{\mathcal{A}_h} + n_0)\}$, it holds that

$$\|\hat{\mathbf{w}}^{\mathcal{A}_h}\|_2 \leq 8\kappa_2 C_{\mathbf{w}} h \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + 3 \frac{\sqrt{s}}{\kappa_1} \lambda_{\mathbf{w}} + 2 \sqrt{\frac{1}{\kappa_1} h \lambda_{\mathbf{w}}}. \quad (\text{S.2.31})$$

Via the decomposition in (S.2.21), $\|\nabla L(\mathbf{w}^{\mathcal{A}_h}, \mathcal{D})\|_{\infty}$ can be bounded by two parts where the first part has rate $\mathcal{O}_p\left(\sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}}\right)$. Denote $V_{ij}^{(k)} = x_{ij}^{(k)} [-\psi'((\mathbf{x}_i^{(k)})^T \mathbf{w}^{(k)}) + \psi'((\mathbf{x}_i^{(k)})^T \mathbf{w}^{\mathcal{A}_h})]$.

$$\begin{aligned} & \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} V_{ij}^{(k)} \\ &= \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} \psi''((\mathbf{x}_i^{(k)})^T \mathbf{w}^{(k)} + v_i^{(k)} (\mathbf{x}_i^{(k)})^T (\mathbf{w}^{\mathcal{A}_h} - \mathbf{w}^{(k)})) x_{ij}^{(k)} (\mathbf{x}_i^{(k)})^T (\mathbf{w}^{\mathcal{A}_h} - \mathbf{w}^{(k)}). \end{aligned}$$

Similar as before, the multiplication $\psi''((\mathbf{x}_i^{(k)})^T \mathbf{w}^{(k)} + v_i^{(k)} (\mathbf{x}_i^{(k)})^T (\mathbf{w}^{\mathcal{A}_h} - \mathbf{w}^{(k)})) x_{ij}^{(k)} (\mathbf{x}_i^{(k)})^T (\mathbf{w}^{\mathcal{A}_h} - \mathbf{w}^{(k)})$ is a $C_1^2 M_{\psi}^2 \kappa_u^2 h^2$ -subexponential variable. And by Cauchy-Schwarz inequality and sub-Gaussian properties (Vershynin, 2018),

$$\mathbb{E}|V_{ij}^{(k)}| \leq (\mathbb{E}|x_{ij}^{(k)}|^2)^{1/2} (\mathbb{E}[(\mathbf{x}_i^{(k)})^T (\mathbf{w}^{\mathcal{A}_h} - \mathbf{w}^{(k)})]^2)^{1/2} \lesssim \kappa_u h.$$

Therefore, by tail bounds of sub-exponential variables and union bounds, we have

$$\frac{1}{n_{\mathcal{A}_h} + n_0} \sup_{j=1, \dots, p} \left| \sum_{k \in \{0\} \cup \mathcal{A}_h} \sum_{i=1}^{n_k} V_{ij}^{(k)} \right| \lesssim \kappa_u h + C_1 M_{\psi} \kappa_u h \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}},$$

with probability at least $1 - (n_{\mathcal{A}_h} + n_0)^{-1} \vee p^{-1}$, which implies that $\|\nabla \hat{L}(\mathbf{w}^{\mathcal{A}_h})\|_\infty \lesssim \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + h$ with probability at least $1 - (n_{\mathcal{A}_h} + n_0)^{-1} \vee p^{-1}$. Let $\lambda_{\mathbf{w}} = C_{\mathbf{w}} \left(\sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + h \right)$ and $\lambda_{\mathbf{w}} \geq 2\|\nabla L(\mathbf{w}^{\mathcal{A}_h}, \mathcal{D})\|_\infty$ in high probability. Plugging it into (S.2.31), we get

$$\|\hat{\mathbf{u}}^{\mathcal{A}_h}\|_2 \lesssim \sqrt{\frac{s \log p}{n_{\mathcal{A}_h} + n_0}} + \sqrt{sh} + \sqrt{h} \left(\frac{\log p}{n_{\mathcal{A}_h} + n_0} \right)^{1/4}, \quad (\text{S.2.32})$$

with probability at least $1 - (n_{\mathcal{A}_h} + n_0)^{-1} \vee p^{-1}$. Similarly, we can obtain the ℓ_1 -error bound

$$\|\hat{\mathbf{u}}^{\mathcal{A}_h}\|_1 \lesssim s \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + sh + \sqrt{sh} \left(\frac{\log p}{n_{\mathcal{A}_h} + n_0} \right)^{1/4}, \quad (\text{S.2.33})$$

with probability at least $1 - (n_{\mathcal{A}_h} + n_0)^{-1} \vee p^{-1}$. Next, consider the debiasing step. Similar as the proof of Theorem 1, let $\lambda_{\boldsymbol{\beta}} = C_{\boldsymbol{\beta}} \sqrt{\frac{\log p}{n_0}}$ which satisfies $\lambda_{\boldsymbol{\beta}} \geq 2\|\nabla L^{(0)}(\boldsymbol{\beta}, \mathcal{D}^{(0)})\|_\infty$ in high probability to get

$$\begin{aligned} \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_2^2 &\lesssim \frac{s \log p}{n_{\mathcal{A}_h} + n_0} + sh^2 + \left[h \sqrt{\frac{\log p}{n_0}} \right] \wedge h^2, \\ \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_1 &\lesssim s \sqrt{\frac{\log p}{n_{\mathcal{A}} + n_0}} + sh + \sqrt{sh} \cdot \left(\frac{\log p}{n_{\mathcal{A}_h} + n_0} \right)^{1/4}, \end{aligned}$$

with probability at least $1 - Cn_0^{-1}$.

S.2.3.4 Proof of Theorem 4

Lemma 6. *Under Assumptions 1-3, $\sup_{k \in \mathcal{A}} L_0(\boldsymbol{\beta}^{(k)}) - L_0(\boldsymbol{\beta}) \lesssim \|\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}\|_2^2 \lesssim h^2$.*

Proof of Lemma 6. Note that the term involving $\mathbb{E}[\rho(y^{(0)})]$ is canceled when taking the difference, therefore we drop that term and consider $L_0(\mathbf{w}) = -\mathbb{E}[\psi'(\mathbf{w}^T \mathbf{x}^{(0)}) \mathbf{w}^T \mathbf{x}^{(0)}] + \mathbb{E}[\psi(\mathbf{w}^T \mathbf{x}^{(0)})]$. Since $\nabla L_0(\boldsymbol{\beta}) = \mathbf{0}$ and $\nabla^2 L_0(\mathbf{w}) = \mathbb{E}[\psi''(\mathbf{w}^T \mathbf{x}^{(0)}) \mathbf{x}^{(0)} (\mathbf{x}^{(0)})^T]$. By mean-theorem, $\exists t^{(k)} \in (0, 1)$, such that

$$L_0(\boldsymbol{\beta}^{(k)}) - L_0(\boldsymbol{\beta}) = (\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta})^T \mathbb{E}[\psi''(\boldsymbol{\beta}^T \mathbf{x}^{(0)} + t^{(k)}(\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta})^T \mathbf{x}^{(0)}) \mathbf{x}^{(0)} (\mathbf{x}^{(0)})^T] (\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}).$$

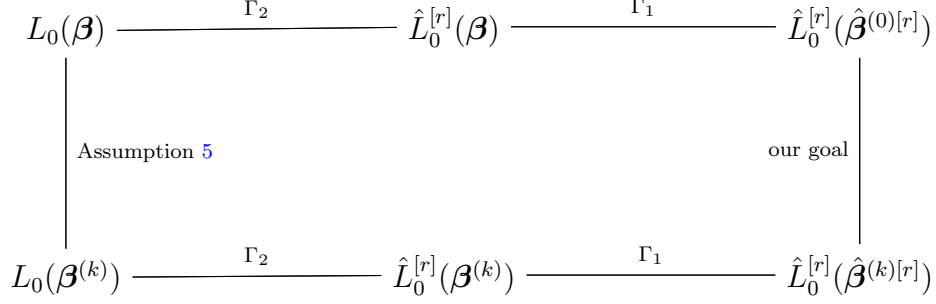


Figure 12: The idea behind Assumptions 4, 5 and Theorem 4.

Under Assumption 3.(i):

$$L_0(\boldsymbol{\beta}^{(k)}) - L_0(\boldsymbol{\beta}) \leq M\mathbb{E}[(\boldsymbol{x}^{(0)})^T(\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta})]^2] \lesssim \|\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}\|_2^2.$$

Under Assumption 3.(ii), by Cauchy-Schwarz inequality and the subGaussian moment bound (Vershynin, 2018):

$$L_0(\boldsymbol{\beta}^{(k)}) - L_0(\boldsymbol{\beta}) \leq \left[\mathbb{E} \left(\max_{z:|z|\leq 1} \psi''((\boldsymbol{x}^{(0)})^T \boldsymbol{\beta} + z) \right)^2 \right]^{1/2} [\mathbb{E}((\boldsymbol{x}^{(0)})^T(\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}))^4]^{1/4} \lesssim \|\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}\|_2^2.$$

The second half inequality automatically holds since $\boldsymbol{\beta}^{(k)}$ is a linear combination of $\boldsymbol{\beta}$ and $\boldsymbol{w}^{(k)}$. And it's easy to see that all the constants appearing in the inequalities are uniform for all $k \in \mathcal{A}$, which completes the proof. \square

Next, we prove Theorem 4. We have

$$\begin{aligned} \hat{\sigma} &= \sqrt{\sum_{r=1}^3 (\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(0)[r]}) - L_0(\boldsymbol{\beta}))^2 / 3} \leq \sqrt{\frac{2}{3}} \cdot \sum_{r=1}^3 |\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(0)[r]}) - L_0(\boldsymbol{\beta})| \\ &\leq \sqrt{\frac{2}{3}} \cdot \sum_{r=1}^3 \left[|\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(0)[r]}) - \hat{L}_0^{[r]}(\boldsymbol{\beta})| + |\hat{L}_0^{[r]}(\boldsymbol{\beta}) - L_0(\boldsymbol{\beta})| \right] \end{aligned}$$

$$\lesssim \zeta \left(\Gamma_1^{(0)} + \Gamma_2^{(0)} \right),$$

with probability at least $1 - g_1^{(0)}(\zeta) - g_2^{(0)}(\zeta)$. As Figure 12 shows, by Lemma 6, for $k \in \mathcal{A}$, there holds

$$\begin{aligned} \sup_r |\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]}) - \hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(0)[r]})| &\leq 2 \sup_r |\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]}) - \hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)[r]})| \\ &\quad + \sup_r |\hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)}) - \hat{L}_0^{[r]}(\boldsymbol{\beta}) - L_0(\boldsymbol{\beta}^{(k)}) + L_0(\boldsymbol{\beta})| \\ &\quad + |L_0(\boldsymbol{\beta}^{(k)}) - L_0(\boldsymbol{\beta})| \\ &\lesssim \zeta \left(\Gamma_1^{(k)} + \Gamma_2^{(k)} + h^2 \right) \\ &\leq C_0(\hat{\sigma} \vee 0.01), \end{aligned} \tag{S.2.34}$$

simultaneously with probability at least $1 - |\mathcal{A}| \max_{k \in \mathcal{A}} [g_1^{(k)}(\zeta) + g_2^{(k)}(\zeta)]$ for sufficiently small $\zeta > 0$ when $\min_{k \in \mathcal{A}} n_k$ and n_0 go to infinity since $\Gamma_1^{(k)} + \Gamma_2^{(k)} + h^2 = o(1)$. On the other hand, by Assumption 5 and the fact $\nabla L_0(\boldsymbol{\beta}) = \mathbf{0}$, for $k \in \mathcal{A}^c$,

$$\begin{aligned} &\inf_r \hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]}) - \hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(0)[r]}) \\ &\geq |L_0(\boldsymbol{\beta}^{(k)}) - L_0(\boldsymbol{\beta})| - \Upsilon_1^{(k)} - \zeta \Gamma_1^{(k)} - \zeta \Gamma_2^{(k)} \\ &= \|\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}\|_2^2 \cdot \lambda_{\min} \left(\int_0^1 \nabla^2 L_0(\boldsymbol{\beta} + t(\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta})) dt \right) - \Upsilon_1^{(k)} - \zeta \Gamma_1^{(k)} - \zeta \Gamma_2^{(k)} \\ &\geq \|\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}\|_2^2 \cdot \underline{\lambda} - \Upsilon_1^{(k)} - \zeta \Gamma_1^{(k)} - \zeta \Gamma_2^{(k)} \\ &\geq C_1^2 \left[(\Gamma_1^{(0)} + \Gamma_2^{(0)}) \vee 1 \right] - \zeta \Gamma_1^{(k)} - \zeta \Gamma_2^{(k)} \\ &> C_0(\hat{\sigma} \vee 0.01), \end{aligned}$$

simultaneously with probability at least $1 - |\mathcal{A}^c| \max_{k \in \mathcal{A}^c} [g_1^{(k)}(C_0^{-1}) + g_2^{(k)}(C_0^{-1})] - |\mathcal{A}^c| \max_{k \in \mathcal{A}^c}$

$[g_1^{(k)}(\zeta) + g_2^{(k)}(\zeta)]$. It entails

$$\begin{aligned}
\mathbb{P}(\widehat{\mathcal{A}} \neq \mathcal{A}_h) &\leq \mathbb{P}\left(\bigcup_{k \in \mathcal{A}} \{\hat{L}_0^{(k)} - \hat{L}_0^{(0)} > C_0(\hat{\sigma} \vee 0.01)\} \cup \bigcup_{k \in \mathcal{A}^c} \{\hat{L}_0^{(k)} - \hat{L}_0^{(0)} \leq C_0(\hat{\sigma} \vee 0.01)\}\right) \\
&\leq \sum_{k \in \mathcal{A}} \mathbb{P}\left(\hat{L}_0^{(k)} - \hat{L}_0^{(0)} > C_0(\hat{\sigma} \vee 0.01)\right) + \sum_{k \in \mathcal{A}^c} \mathbb{P}\left(\hat{L}_0^{(k)} - \hat{L}_0^{(0)} \leq C_0(\hat{\sigma} \vee 0.01)\right) \\
&\leq \sum_{k \in \mathcal{A}} \mathbb{P}\left(\inf_r |\hat{L}_0^{[r]}(\hat{\beta}^{(k)[r]}) - \hat{L}_0^{[r]}(\hat{\beta}^{(0)[r]})| > C_0(\hat{\sigma} \vee 0.01)\right) \\
&\quad + \sum_{k \in \mathcal{A}^c} \mathbb{P}\left(\sup_r |\hat{L}_0^{[r]}(\hat{\beta}^{(k)[r]}) - \hat{L}_0^{[r]}(\hat{\beta}^{(0)[r]})| \leq C_0(\hat{\sigma} \vee 0.01)\right) \\
&\leq |\mathcal{A}| \max_{k \in \mathcal{A}} [g_1^{(k)}(\zeta) + g_2^{(k)}(\zeta)] + |\mathcal{A}^c| \max_{k \in \mathcal{A}^c} [g_1^{(k)}(C_0^{-1}) + g_2^{(k)}(C_0^{-1})].
\end{aligned}$$

For any $\delta > 0$, there exist constants $C'(\delta)$ and $\zeta' > 0$ such that when $C_0 = C'(\delta)$, $K \max_{k \in \mathcal{A}^c} [g_1^{(k)}(C_0^{-1}) + g_2^{(k)}(C_0^{-1})] \leq \delta/2$, $K \max_{k \in \mathcal{A}} [g_1^{(k)}(\zeta') + g_2^{(k)}(\zeta')] < \delta/2$ and $C_1^2[(\Gamma_1^{(0)} + \Gamma_2^{(0)}) \vee 1] > \zeta' \Gamma_1^{(k)} + \zeta' \Gamma_2^{(k)}$. On the other hand, there exists $N = N(\delta) > 0$, such that when $\min_{k \in \{0\} \cup \mathcal{A}_h} n_k > N(\delta)$, $\zeta' (\Gamma_1^{(k)} + \Gamma_2^{(k)} + h^2) / [C'(\delta) \cdot (\hat{\sigma} \vee 0.01)]$ is sufficiently small to make (S.2.34) hold.

In summary, for any $\delta > 0$, there exist constants $C'(\delta)$ and $N = N(\delta) > 0$ such that when $C_0 = C'(\delta)$ and $\min_{k \in \{0\} \cup \mathcal{A}_h} n_k > N(\delta)$, $\mathbb{P}(\widehat{\mathcal{A}} \neq \mathcal{A}_h) \leq \delta$.

S.2.3.5 Proof of Theorem 5

Denote $\mathfrak{R}_1 = \sqrt{\frac{s \log p}{n_{\mathcal{A}} + n_0}} + \left[h^{1/2} \left(\frac{\log p}{n_0} \right)^{1/4} \right] \wedge h$. First we present the following proposition.

Proposition 2. *Assume Assumptions 1-4 and 7 hold (except (v)). Let $\lambda_j \asymp \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + \left[\frac{1}{\sqrt{s}} \cdot \left(\frac{\log p}{n_0} \right)^{1/4} h^{1/2} \right] \wedge \frac{h}{\sqrt{s}}$ for any $j = 1, \dots, p$, and $\tilde{\lambda}_j \asymp \sqrt{\frac{\log p}{n_0}} + \mathfrak{R}_1$. Then with probability*

at least $1 - K_{\mathcal{A}_h} n_0^{-1}$,

$$\|\hat{\gamma}_j^{(0)} - \gamma_j^{(0)}\|_2^2 \lesssim h_1 \sqrt{\frac{\log p}{n_0}} + h_1 \mathfrak{R}_1 + \mathfrak{R}_1^2, \quad (\text{S.2.35})$$

$$\|\hat{\gamma}_j^{(0)} - \gamma_j^{(0)}\|_1 \lesssim \sqrt{s} \mathfrak{R}_1 + s^{1/4} \mathfrak{R}_1^{1/2} h_1^{1/2} + h_1, \quad (\text{S.2.36})$$

$$|\hat{\tau}_j^2 - \tau_j^2| \vee |\hat{\tau}_j^{-2} - \tau_j^{-2}| \lesssim \mathfrak{R}_1 + h_1^{1/2} \left(\frac{\log p}{n_0} \right)^{1/4} + h_1^{1/2} \mathfrak{R}_1^{1/2} + h_{\max}, \quad (\text{S.2.37})$$

(S.2.35) and (S.2.36) can be proved by combining conclusions in Lemmas 4 and 5. The proof of (S.2.37) can be found in Section S.2.4.2.

Next we will leverage Proposition 2 to show (10) and (11). First, we notice that

$$\begin{aligned} \hat{\mathbf{b}} &= \hat{\boldsymbol{\beta}} + \frac{1}{n_0} \hat{\boldsymbol{\Theta}}(\mathbf{X}^{(0)})^T [\mathbf{Y}^{(0)} - \boldsymbol{\psi}'(\mathbf{X}^{(0)} \hat{\boldsymbol{\beta}})] \\ &= \hat{\boldsymbol{\beta}} + \frac{1}{n_0} \hat{\boldsymbol{\Theta}}(\mathbf{X}^{(0)})^T [\mathbf{Y}^{(0)} - \boldsymbol{\psi}'(\mathbf{X}^{(0)} \boldsymbol{\beta})] + \frac{1}{n_0} \hat{\boldsymbol{\Theta}}(\mathbf{X}^{(0)})^T [\boldsymbol{\psi}'(\mathbf{X}^{(0)} \hat{\boldsymbol{\beta}}) - \boldsymbol{\psi}'(\mathbf{X}^{(0)} \boldsymbol{\beta})] \\ &= \hat{\boldsymbol{\beta}} + \frac{1}{n_0} \hat{\boldsymbol{\Theta}}(\mathbf{X}^{(0)})^T [\mathbf{Y}^{(0)} - \boldsymbol{\psi}'(\mathbf{X}^{(0)} \boldsymbol{\beta})] + \frac{1}{n_0} \hat{\boldsymbol{\Theta}} \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\beta}}^{(0)} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \\ &\quad + \frac{1}{n_0} \hat{\boldsymbol{\Theta}}(\mathbf{X}^{(0)})^T \text{diag}(\{\psi''(\tilde{u}_i^{(0)})\}) - \psi''((\mathbf{x}_i^{(0)})^T \boldsymbol{\beta})_{i=1}^{n_0}) \mathbf{X}^{(0)} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}). \end{aligned}$$

where $\tilde{u}_i^{(0)}$ falls on the line between $(\mathbf{x}_i^{(0)})^T \boldsymbol{\beta}$ and $(\mathbf{x}_i^{(0)})^T \hat{\boldsymbol{\beta}}$. Therefore for any $j = 1, \dots, p$,

$$\begin{aligned} \hat{b}_j - \beta_j &= \underbrace{\left[\mathbf{e}_j - (\hat{\boldsymbol{\Theta}}_j)^T \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\beta}}^{(0)} \right]^T}_{(1)} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \underbrace{\frac{1}{n_0} (\hat{\boldsymbol{\Theta}}_j)^T (\mathbf{X}^{(0)})^T [\mathbf{Y}^{(0)} - \boldsymbol{\psi}'(\mathbf{X}^{(0)} \boldsymbol{\beta})]}_{(2)} \\ &\quad + \underbrace{\frac{1}{n_0} (\hat{\boldsymbol{\Theta}}_j)^T (\mathbf{X}^{(0)})^T \text{diag}(\{\psi''(\tilde{u}_i^{(0)}) - \psi''((\mathbf{x}_i^{(0)})^T \boldsymbol{\beta})_{i=1}^{n_0})\}}_{(3)} \mathbf{X}^{(0)} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}), \end{aligned}$$

For convenience, we write the j -th row of the matrix

$$\begin{pmatrix} 1 & -\hat{\gamma}_{1,2} & \cdots & -\hat{\gamma}_{1,p} \\ -\hat{\gamma}_{2,1} & 1 & \cdots & -\hat{\gamma}_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ -\hat{\gamma}_{p,1} & -\hat{\gamma}_{p,2} & \cdots & -\hat{\gamma}_{p,p} \end{pmatrix}$$

as $(\hat{\gamma}_j^{(0)})^\dagger$ and the j -th row of matrix

$$\begin{pmatrix} 1 & -\gamma_{1,2} & \cdots & -\gamma_{1,p} \\ -\gamma_{2,1} & 1 & \cdots & -\gamma_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ -\gamma_{p,1} & -\gamma_{p,2} & \cdots & -\gamma_{p,p} \end{pmatrix}$$

as $(\gamma_j^{(0)})^\dagger$.

For (i), notice that

$$\begin{aligned} & \left\| \widehat{\Sigma}_\beta^{(0)} \widehat{\Theta}_j - e_j \right\|_\infty \\ &= \left\| \widehat{\Sigma}_\beta^{(0)} \widehat{\Theta}_j - \Sigma_\beta^{(0)} \Theta_j \right\|_\infty \\ &= \left\| \frac{1}{n_0} \sum_{i=1}^{n_0} \mathbf{x}_i^{(0)} \psi''((\mathbf{x}_i^{(0)})^T \boldsymbol{\beta}) (\mathbf{x}_i^{(0)})^T \frac{(\hat{\gamma}_j^{(0)})^\dagger}{\hat{\tau}_j^2} - \mathbb{E} \left[\mathbf{x}_i^{(0)} \psi''((\mathbf{x}_i^{(0)})^T \boldsymbol{\beta}) (\mathbf{x}_i^{(0)})^T \frac{(\gamma_j^{(0)})^\dagger}{\tau_j^2} \right] \right\|_\infty \\ &= \left| \frac{1}{n_0} \sum_{i=1}^{n_0} \mathbf{x}_i^{(0)} \psi''((\mathbf{x}_i^{(0)})^T \boldsymbol{\beta}) (\mathbf{x}_i^{(0)})^T \frac{(\hat{\gamma}_j^{(0)})^\dagger}{\hat{\tau}_j^2} - \frac{1}{n_0} \sum_{i=1}^{n_0} \mathbf{x}_i^{(0)} \psi''((\mathbf{x}_i^{(0)})^T \boldsymbol{\beta}) (\mathbf{x}_i^{(0)})^T \frac{(\gamma_j^{(0)})^\dagger}{\tau_j^2} \right| \\ & \quad + \left| \frac{1}{n_0} \sum_{i=1}^{n_0} \mathbf{x}_i^{(0)} \psi''((\mathbf{x}_i^{(0)})^T \boldsymbol{\beta}) (\mathbf{x}_i^{(0)})^T \frac{(\gamma_j^{(0)})^\dagger}{\tau_j^2} - \mathbb{E} \left[\mathbf{x}_i^{(0)} \psi''((\mathbf{x}_i^{(0)})^T \boldsymbol{\beta}) (\mathbf{x}_i^{(0)})^T \frac{(\gamma_j^{(0)})^\dagger}{\tau_j^2} \right] \right| \\ &\lesssim \|\hat{\gamma}_j^{(0)} - \gamma_j^{(0)}\|_2 + |\hat{\tau}_j^{-2} - \tau_j^{-2}| + n_0^{-1/2} \\ &\lesssim \mathfrak{R}_1 + h_1^{1/2} \left(\frac{\log p}{n_0} \right)^{1/4} + h_1^{1/2} \mathfrak{R}_1^{1/2} + h_{\max} + n_0^{-1/2}, \end{aligned}$$

implying that

$$\begin{aligned}
(1) &\leq \left\| \widehat{\Sigma}_{\beta}^{(0)} \widehat{\Theta}_j - \mathbf{e}_j \right\|_{\infty} \|\widehat{\beta} - \beta\|_1 \\
&\lesssim_p \left[\mathfrak{R}_1 + h_1^{1/2} \left(\frac{\log p}{n_0} \right)^{1/4} + h_1^{1/2} \mathfrak{R}_1^{1/2} + h_{\max} + n_0^{-1/2} \right] \cdot \left(s \sqrt{\frac{\log p}{n_{\mathcal{A}_h} + n_0}} + h \right) \\
&= \mathcal{O}(n_0^{-1/2}), \tag{S.2.38}
\end{aligned}$$

under Assumption 7.

For (3),

$$(3) \leq \frac{1}{n_0} \left\| \mathbf{X}^{(0)} \widehat{\Theta}_j \right\|_{\infty} \left\| \text{diag}(\{\psi''(\tilde{u}_i^{(0)}) - \psi''((\mathbf{x}_i^{(0)})^T \beta)\}_{i=1}^{n_0}) \mathbf{X}^{(0)} (\beta - \widehat{\beta}) \right\|_1,$$

where $|\psi''(\tilde{u}_i^{(0)}) - \psi''((\mathbf{x}_i^{(0)})^T \beta)| \cdot |(\mathbf{x}_i^{(0)})^T (\beta - \widehat{\beta})| \leq \psi'''(t\tilde{u}_i^{(0)} + (1-t)(\mathbf{x}_i^{(0)})^T \beta) \cdot |\tilde{u}_i^{(0)} - (\mathbf{x}_i^{(0)})^T \beta| |(\mathbf{x}_i^{(0)})^T (\beta - \widehat{\beta})| \lesssim |(\mathbf{x}_i^{(0)})^T (\beta - \widehat{\beta})|^2$. Therefore by similar arguments to get (S.2.27), we have

$$(3) \lesssim_p (1 + \|\widehat{\Theta}_j - \Theta_j\|_1) \cdot \|\widehat{\beta} - \beta\|_2^2 \lesssim_p \mathfrak{R}_1^2 = \mathcal{O}(n_0^{-1/2}) \tag{S.2.39}$$

under Assumption 7.

For (2), we can see that

$$(2) = \frac{1}{n_0} \widehat{\Theta}_j^T (\mathbf{X}^{(0)})^T [\mathbf{Y}^{(0)} - \psi'(\mathbf{X}^{(0)} \beta)] = \frac{1}{n_0} \sum_{i=1}^{n_0} \widehat{\Theta}_j^T \mathbf{x}_i^{(0)} [y_i^{(0)} - \psi'((\mathbf{x}_i^{(0)})^T \beta)], \tag{S.2.40}$$

where $\mathbb{E}\{(\widehat{\Theta}_j^T \mathbf{x}_i^{(0)}) [y_i^{(0)} - \psi'((\mathbf{x}_i^{(0)})^T \beta)]\} = 0$ and

$$\begin{aligned}
\mathbb{E}\{\widehat{\Theta}_j^T \mathbf{x}_i^{(0)} [y_i^{(0)} - \psi'((\mathbf{x}_i^{(0)})^T \beta)]\}^2 &= \mathbb{E}\{(\widehat{\Theta}_j^T \mathbf{x}_i^{(0)})^2 \psi''((\mathbf{x}_i^{(0)})^T \beta)\} \\
&= \Theta_j^T \mathbb{E}[\mathbf{x}^{(0)} (\mathbf{x}^{(0)})^T \psi''((\mathbf{x}^{(0)})^T \beta)] \Theta_j \\
&= \Theta_j^T \Sigma_{\beta}^{(0)} \Theta_j
\end{aligned}$$

$$\begin{aligned} &\lesssim \mathbb{E} (\boldsymbol{\Theta}_j^T \mathbf{x}^{(0)})^2 \\ &< \infty, \end{aligned}$$

by Assumption 7. Thus by Lindeberg's conditions,

$$\sqrt{n_0} \cdot (2) \xrightarrow{d} \mathcal{N}(0, \boldsymbol{\Theta}_{j,j}). \quad (\text{S.2.41})$$

On the other hand,

$$\left| \widehat{\boldsymbol{\Theta}}_j^T \widehat{\boldsymbol{\Sigma}}_{\widehat{\boldsymbol{\beta}}} \widehat{\boldsymbol{\Theta}}_j - \boldsymbol{\Theta}_{j,j} \right| \leq \left| \widehat{\boldsymbol{\Theta}}_j^T (\widehat{\boldsymbol{\Sigma}}_{\widehat{\boldsymbol{\beta}}} \widehat{\boldsymbol{\Theta}}_j - \boldsymbol{\Sigma}_{\boldsymbol{\beta}}^{(0)} \boldsymbol{\Theta}_j) \right| + \left| (\widehat{\boldsymbol{\Theta}}_j - \boldsymbol{\Theta}_j)^T \boldsymbol{\Sigma}_{\boldsymbol{\beta}}^{(0)} \boldsymbol{\Theta}_j \right|.$$

Note that

$$\begin{aligned} &\left\| \widehat{\boldsymbol{\Sigma}}_{\widehat{\boldsymbol{\beta}}} \widehat{\boldsymbol{\Theta}}_j - \boldsymbol{\Sigma}_{\boldsymbol{\beta}}^{(0)} \boldsymbol{\Theta}_j \right\|_{\infty} \\ &= \left\| \widehat{\boldsymbol{\Sigma}}_{\widehat{\boldsymbol{\beta}}} \frac{(\widehat{\boldsymbol{\gamma}}_j^{(0)})^{\dagger}}{\widehat{\tau}_j^2} - \boldsymbol{\Sigma}_{\boldsymbol{\beta}}^{(0)} \frac{(\boldsymbol{\gamma}_j^{(0)})^{\dagger}}{\tau_j^2} \right\|_{\infty} \\ &\leq \left\| \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i,k} \left[\mathbf{x}_i^{(k)} \psi''(\widehat{\boldsymbol{\beta}}^T \mathbf{x}_i^{(k)}) (\mathbf{x}_i^{(k)})^T \frac{(\widehat{\boldsymbol{\gamma}}_j^{(0)})^{\dagger}}{\widehat{\tau}_j^2} \right] - \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i,k} \left[\mathbf{x}_i^{(k)} \psi''(\boldsymbol{\beta}^T \mathbf{x}_i^{(k)}) (\mathbf{x}_i^{(k)})^T \frac{(\boldsymbol{\gamma}_j^{(0)})^{\dagger}}{\tau_j^2} \right] \right\|_{\infty} \\ &\quad + \frac{1}{n_{\mathcal{A}_h} + n_0} \left\| \sum_{i,k} \left[\mathbf{x}_i^{(k)} \psi''(\boldsymbol{\beta}^T \mathbf{x}_i^{(k)}) (\mathbf{x}_i^{(k)})^T \frac{(\boldsymbol{\gamma}_j^{(0)})^{\dagger}}{\tau_j^2} \right] - \sum_{k \in \{0\} \cup \mathcal{A}_h} n_k \mathbb{E} \left[\mathbf{x}_i^{(k)} \psi''(\boldsymbol{\beta}^T \mathbf{x}_i^{(k)}) (\mathbf{x}_i^{(k)})^T \frac{(\boldsymbol{\gamma}_j^{(0)})^{\dagger}}{\tau_j^2} \right] \right\|_{\infty} \\ &\quad + \left\| \sum_{k \in \{0\} \cup \mathcal{A}_h} \frac{n_k}{n_{\mathcal{A}_h} + n_0} \mathbb{E} \left[\mathbf{x}_i^{(k)} \psi''(\boldsymbol{\beta}^T \mathbf{x}_i^{(k)}) (\mathbf{x}_i^{(k)})^T \frac{(\boldsymbol{\gamma}_j^{(0)})^{\dagger}}{\tau_j^2} - \mathbf{x}_i^{(0)} \psi''(\boldsymbol{\beta}^T \mathbf{x}_i^{(0)}) (\mathbf{x}_i^{(0)})^T \frac{(\boldsymbol{\gamma}_j^{(0)})^{\dagger}}{\tau_j^2} \right] \right\|_{\infty} \\ &\lesssim \underbrace{\left\| \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i,k} \mathbf{x}_i^{(k)} \left[\psi''(\widehat{\boldsymbol{\beta}}^T \mathbf{x}_i^{(k)}) - \psi''(\boldsymbol{\beta}^T \mathbf{x}_i^{(k)}) \right] (\mathbf{x}_i^{(k)})^T \frac{(\widehat{\boldsymbol{\gamma}}_j^{(0)})^{\dagger}}{\widehat{\tau}_j^2} \right\|_{\infty}}_{(4)} \end{aligned}$$

$$\begin{aligned}
& + \underbrace{\left\| \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i,k} \mathbf{x}_i^{(k)} \psi''(\boldsymbol{\beta}^T \mathbf{x}_i^{(k)}) (\mathbf{x}_i^{(k)})^T [(\hat{\gamma}_j^{(0)})^\dagger - (\gamma_j^{(0)})^\dagger] \hat{\tau}_j^{-2} \right\|_\infty}_{(5)} \\
& + \underbrace{\left\| \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i,k} \mathbf{x}_i^{(k)} \psi''(\boldsymbol{\beta}^T \mathbf{x}_i^{(k)}) \cdot (\mathbf{x}_i^{(k)})^T (\hat{\gamma}_j^{(0)})^\dagger (\hat{\tau}_j^2 - \tau_j^2) \right\|_\infty}_{(6)} \\
& + \frac{1}{n_{\mathcal{A}_h} + n_0} \underbrace{\left\| \sum_{i,k} \left[\mathbf{x}_i^{(k)} \psi''(\boldsymbol{\beta}^T \mathbf{x}_i^{(k)}) (\mathbf{x}_i^{(k)})^T \frac{(\gamma_j^{(0)})^\dagger}{\tau_j^2} \right] - \sum_{k \in \{0\} \cup \mathcal{A}_h} n_k \mathbb{E} \left[\mathbf{x}_i^{(k)} \psi''(\boldsymbol{\beta}^T \mathbf{x}_i^{(k)}) (\mathbf{x}_i^{(k)})^T \frac{(\gamma_j^{(0)})^\dagger}{\tau_j^2} \right] \right\|_\infty}_{(7)} \\
& + \underbrace{\left\| \sum_{k \in \{0\} \cup \mathcal{A}_h} \frac{n_k}{n_{\mathcal{A}_h} + n_0} \mathbb{E} \left[\mathbf{x}_i^{(k)} \psi''(\boldsymbol{\beta}^T \mathbf{x}_i^{(k)}) (\mathbf{x}_i^{(k)})^T \frac{(\gamma_j^{(0)})^\dagger}{\tau_j^2} - \mathbf{x}_i^{(0)} \psi''(\boldsymbol{\beta}^T \mathbf{x}_i^{(0)}) (\mathbf{x}_i^{(0)})^T \frac{(\gamma_j^{(0)})^\dagger}{\tau_j^2} \right] \right\|_\infty}_{(8)}.
\end{aligned}$$

It's easy to see that

$$\begin{aligned}
(4) & \lesssim \left\| \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i,k} \mathbf{x}_i^{(k)} \cdot \psi'''(t_i^{(k)}) (\mathbf{x}_i^{(k)})^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) + (\mathbf{x}_i^{(k)})^T \boldsymbol{\beta} \cdot (\mathbf{x}_i^{(k)})^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \cdot (\mathbf{x}_i^{(k)})^T \frac{(\hat{\gamma}_j^{(0)})^\dagger}{\hat{\tau}_j^2} \right\|_\infty \\
& \lesssim \sqrt{\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k,i} |(\mathbf{x}_i^{(k)})^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})|^2} \\
& \lesssim \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_2 \\
& \lesssim \mathfrak{R}_1,
\end{aligned}$$

with probability at least $1 - n_0^{-1}$.

Similarly,

$$\begin{aligned}
(5) &\lesssim \|\hat{\gamma}_j^{(0)} - \gamma_j^{(0)}\|_2 + |\hat{\tau}_j^{-2} - \tau_j^{-2}| \lesssim \mathfrak{R}_1 + h_*^{1/2} \left(\frac{\log p}{n_0} \right)^{1/4} + h_1^{1/2} \mathfrak{R}_1^{1/2} + h_{\max} \\
(6) &\lesssim |\hat{\tau}_j^{-2} - \tau_j^{-2}| \lesssim \mathfrak{R}_1 + h_*^{1/2} \left(\frac{\log p}{n_0} \right)^{1/4} + h_1^{1/2} \mathfrak{R}_1^{1/2} + h_{\max} \\
(7) &\lesssim (n_{\mathcal{A}_h} + n_0)^{-1/2}, \\
(8) &\leq \sup_{k \in \mathcal{A}_h} \|\Sigma_{\beta}^{(k)} - \Sigma_{\beta}^{(0)}\|_{\max} \cdot \left\| \frac{(\gamma_j^{(0)})^\dagger}{\tau_j^2} \right\|_1 \lesssim h_{\max},
\end{aligned}$$

with probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$. Besides,

$$\|\widehat{\Theta}_j\|_1 \leq \|\Theta_j\|_1 + \|\widehat{\Theta}_j - \Theta_j\|_1 \lesssim \sqrt{s}.$$

with probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$. Therefore,

$$\begin{aligned}
\left| \widehat{\Theta}_j^T (\widehat{\Sigma}_{\beta} \widehat{\Theta}_j - \Sigma_{\beta}^{(0)} \Theta_j) \right| &\leq \|\widehat{\Theta}_j\|_1 \left\| \widehat{\Sigma}_{\beta} \widehat{\Theta}_j - \Sigma_{\beta}^{(0)} \Theta_j \right\|_{\infty} \\
&\lesssim \sqrt{s} \left[\mathfrak{R}_1 + h_*^{1/2} \left(\frac{\log p}{n_0} \right)^{1/4} + h_1^{1/2} \mathfrak{R}_1^{1/2} + h_{\max} \right],
\end{aligned}$$

with probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$. Similarly, we can obtain the same upper bound for the second term $\left| (\widehat{\Theta}_j - \Theta_j) \Sigma_{\beta}^{(0)} \Theta_j \right|$ in (S.2.40). Then finally,

$$\left| \widehat{\Theta}_j^T \widehat{\Sigma}_{\beta} \widehat{\Theta}_j - \Theta_{j,j} \right| \lesssim \sqrt{s} \left[\mathfrak{R}_1 + h_*^{1/2} \left(\frac{\log p}{n_0} \right)^{1/4} + h_1^{1/2} \mathfrak{R}_1^{1/2} + h_{\max} \right] = o(1),$$

with probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$. Then by Slutsky's lemma, equations (S.2.38), (S.2.41) and (S.2.39),

$$\frac{\sqrt{n_0}(\hat{b}_j - \beta_j)}{\sqrt{\widehat{\Theta}_j^T \widehat{\Sigma}_{\beta} \widehat{\Theta}_j}} \xrightarrow{d} \mathcal{N}(0, 1),$$

which completes the proof.

S.2.3.6 Proof of Theorem 6

Assume Assumption 3.(i) holds or Assumption 3.(ii) with $h \leq C'U^{-1}\bar{U}$ for some $C' > 0$ holds. When $\lambda_\delta = C_\delta \sqrt{\frac{\log p}{n_0}}$, where $C_\delta > 0$ is a sufficiently large constant: By the definition of $\hat{\beta}$, we have $\nabla L_{n_0}^{(0)}(\hat{\beta}) + \lambda_\delta \cdot \partial \|\hat{\beta} - \hat{\mathbf{w}}^{A_h}\|_1 = \mathbf{0}$. Then by Hölder inequality,

$$\begin{aligned} \langle \nabla L_{n_0}^{(0)}(\hat{\beta}) - \nabla L_{n_0}^{(0)}(\beta), \hat{\beta} - \beta \rangle &= \langle -\lambda_\delta \cdot \partial \|\hat{\beta} - \hat{\mathbf{w}}^{A_h}\|_1 - \nabla L_{n_0}^{(0)}(\beta), \hat{\beta} - \beta \rangle \\ &\leq \lambda_\delta \|\hat{\beta} - \beta\|_1 + \|\nabla L_{n_0}^{(0)}(\beta)\|_\infty \|\hat{\beta} - \beta\|_1. \end{aligned} \quad (\text{S.2.42})$$

Considering the fact that $\|\nabla L_{n_0}^{(0)}(\beta)\|_\infty \lesssim \sqrt{\frac{\log p}{n_0}}$ with probability at least $1 - n_0^{-1}$ (Lemma 6 in Negahban et al. (2009)) and the upper bound of $\|\hat{\beta} - \beta\|_1$ we prove in Theorem 1, the desired upper bound of $\langle \nabla L_{n_0}^{(0)}(\hat{\beta}) - \nabla L_{n_0}^{(0)}(\beta), \hat{\beta} - \beta \rangle$ follows.

S.2.3.7 Proof of Theorem 7

Assume Assumption 3.(i) holds or Assumption 3.(ii) with $h \leq C'U^{-1}\bar{U}$ for some $C' > 0$ holds. If we take $\lambda_\delta = C_\delta \sqrt{\frac{\log p}{n_0}}$, where $C_\delta > 0$ is a sufficiently large constant: Similar to (S.2.42), we can obtain

$$\langle \nabla L_{n_0}^{(0)}(\hat{\beta}) - \nabla L_{n_0}^{(0)}(\beta), \hat{\beta} - \beta \rangle \leq \lambda_\delta \|\hat{\beta} - \beta\|_1 + \|\nabla L_{n_0}^{(0)}(\beta)\|_\infty \|\hat{\beta} - \beta\|_1.$$

To bound $\|\hat{\beta} - \beta\|_1$, it suffices to combine (S.2.28) and the upper bound of $\|\hat{\mathbf{w}}^{A_h}\|_1$ in (S.2.33). Then the final upper bound follows.

S.2.4 Proof of propositions

S.2.4.1 Proof of Proposition 1

The rate of Γ_1 can be derived from the following Lemma 7 and the union bound, together with the tail inequality (S.2.32). The rate of Γ_2 comes from the following Lemma 8.

Lemma 7. *Under the same assumptions as Theorem 4, we have the following conclusions:*

(i) *For logistic regression model:*

$$\begin{aligned} \sup_r |\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]}) - \hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)})| &\leq C \left(1 + \sqrt{\frac{1}{n_0}} \cdot \zeta\right) \cdot \sup_r \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2, k \in \mathcal{A}, \\ \sup_r |\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]}) - \hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)})| &\leq C \left(1 + \sqrt{\frac{1}{n_0}} \cdot \zeta\right) \cdot \sup_r \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2, k \in \mathcal{A}^c, \\ \sup_r |\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(0)[r]}) - \hat{L}_0^{[r]}(\boldsymbol{\beta})| &\leq C \sup_r \|\hat{\boldsymbol{\beta}}^{(0)[r]} - \boldsymbol{\beta}\|_2 \cdot (1 + \zeta), \end{aligned}$$

with probability at least $1 - \exp\{-\zeta^2\}$.

(ii) *For linear model:*

$$\begin{aligned} \sup_r |\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]}) - \hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)})| &\leq C \left(1 + \sqrt{\frac{1}{n_0}} \cdot \zeta\right) \cdot (\|\mathbf{w}^{(k)}\|_2 \vee \|\boldsymbol{\beta}\|_2) \cdot \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2, k \in \mathcal{A}, \\ \sup_r |\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]}) - \hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)})| &\leq C \left(1 + \sqrt{\frac{1}{n_0}} \cdot \zeta\right) \cdot \|\boldsymbol{\beta}^{(k)}\|_2 \cdot \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2, k \in \mathcal{A}^c, \\ \sup_r |\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(0)[r]}) - \hat{L}_0^{[r]}(\boldsymbol{\beta})| &\leq C \|\boldsymbol{\beta}\|_2 \cdot \|\hat{\boldsymbol{\beta}}^{(0)[r]} - \boldsymbol{\beta}\|_2 \cdot (1 + \zeta), \end{aligned}$$

with probability at least $1 - \exp\{-\zeta^2\}$.

(iii) *For Poisson regression model with $\sup_k \|\mathbf{x}^{(k)}\|_\infty \leq U$ a.s.:*

$$\begin{aligned} \sup_r |\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]}) - \hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)})| &\leq C \left(1 + \sqrt{\frac{1}{n_0}} \cdot \zeta\right) \cdot \exp(U(\|\mathbf{w}^{(k)}\|_1 \vee \|\boldsymbol{\beta}\|_1)) \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2, k \in \mathcal{A}, \\ \sup_r |\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]}) - \hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)})| &\leq C \left(1 + \sqrt{\frac{1}{n_0}} \cdot \zeta\right) \cdot \exp(U\|\boldsymbol{\beta}^{(k)}\|_1) \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2, k \in \mathcal{A}^c, \end{aligned}$$

$$|\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(0)[r]}) - \hat{L}_0^{[r]}(\boldsymbol{\beta})| \lesssim \exp(U\|\boldsymbol{\beta}\|_1) \cdot \|\hat{\boldsymbol{\beta}}^{(0)[r]} - \boldsymbol{\beta}\|_2 \cdot (1 + \zeta),$$

with probability at least $1 - \exp\{-\zeta^2\}$.

Remark 13. It's important to point out that based on Algorithm 2, the randomness of $\hat{L}_0^{[r]}$, $\hat{\boldsymbol{\beta}}^{(k)[r]}$ ($k \neq 0$), and $\hat{\boldsymbol{\beta}}^{(0)[r]}$ is independent. Here $\hat{\boldsymbol{\beta}}^{(k)[r]}$ and $\hat{\boldsymbol{\beta}}^{(0)[r]}$ are regarded as fixed and we only consider the randomness from $\hat{L}_0^{[r]}$.

Proof of Lemma 7. For convenience, we assume n_0 is divisible by 3. Note that the term involving $\sum_{i=1}^{n_0/3} \rho(y_i^{(0)[r]})$ is canceled when we take the difference between $\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]})$ and $\hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)})$. So in the following, we drop that term from the definition of $\hat{L}_0^{[r]}$ in equation (3). We only prove the bound for $|\hat{L}_0(\hat{\boldsymbol{\beta}}^{(k)[r]}) - \hat{L}_0(\boldsymbol{\beta}^{(k)})|$ when $k \in \mathcal{A}$. The cases that $k = 0$ or $k \in \mathcal{A}^c$ can be similarly discussed. Besides, according to the proof of Theorem 3, when $k \in \mathcal{A}$, we define $\boldsymbol{\beta}^{(k)} = \frac{2n_0/3}{2n_0/3+n_k}\boldsymbol{\beta} + \frac{n_k}{2n_0/3+n_k}\boldsymbol{w}^{(k)}$, which gives us the final results shown in Proposition 1.

(i) For logistic regression model, notice that

$$\begin{aligned} |\hat{L}_0(\hat{\boldsymbol{\beta}}^{(k)[r]}) - \hat{L}_0(\boldsymbol{\beta}^{(k)})| &\leq \frac{1}{n_0/3} |(\boldsymbol{y}^{(0)[r]})^T \boldsymbol{X}^{(0)[r]} (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)})| \\ &\quad + \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} [\psi((\boldsymbol{x}_i^{(0)[r]})^T \hat{\boldsymbol{\beta}}^{(k)[r]}) - \psi((\boldsymbol{x}_i^{(0)[r]})^T \boldsymbol{\beta}^{(k)})] \right|. \end{aligned}$$

For the first term on the right-hand side, it holds that

$$\frac{1}{n_0/3} |(\boldsymbol{y}^{(0)[r]})^T \boldsymbol{X}^{(0)[r]} (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)})| \leq \frac{1}{n_0/3} \sum_{i=1}^{n_0/3} |(\boldsymbol{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)})|,$$

where $\frac{1}{n_0/3} \sum_{i=1}^{n_0/3} |(\boldsymbol{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)} - \boldsymbol{\beta}^{(k)})|$ is a $\frac{1}{n_0/3} \|\hat{\boldsymbol{\beta}}^{(k)} - \boldsymbol{\beta}^{(k)}\|_2^2$ sub-Gaussian with mean less than $C \|\hat{\boldsymbol{\beta}}^{(k)} - \boldsymbol{\beta}^{(k)}\|_2$, where $C > 0$ is a uniform constant, implying that

$$\frac{1}{n_0/3} |(\boldsymbol{y}^{(0)[r]})^T \boldsymbol{X}^{(0)[r]} (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)})| \lesssim \left(1 + \sqrt{\frac{1}{n_0}} \cdot \zeta\right) \cdot \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2,$$

with probability at least $1 - \exp\{-\zeta^2\}$. On the other hand, the second term can be bounded by $\frac{C}{n_0/3} \sum_{i=1}^{n_0/3} |(\mathbf{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)})|$, which can be similarly bounded as the first term, leading to the desired conclusion.

(ii) For linear model, note that $y_i^{(0)[r]} = \boldsymbol{\beta}^T \mathbf{x}_i^{(0)[r]} + \epsilon_i^{(0)[r]}$ and $\psi(u) = u^2/2$, leading to

$$\begin{aligned} |\hat{L}_0(\hat{\boldsymbol{\beta}}^{(k)[r]}) - \hat{L}_0(\boldsymbol{\beta}^{(k)})| &\leq \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} \epsilon_i^{(0)[r]} (\mathbf{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}) \right| \\ &\quad + \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} (\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta} \cdot (\mathbf{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}) \right| \\ &\quad + \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} (\mathbf{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} + \boldsymbol{\beta}^{(k)}) (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]} \right|. \end{aligned}$$

It is easy to see that $\frac{1}{n_0/3} \sum_{i=1}^{n_0/3} \epsilon_i^{(0)} (\mathbf{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)})$ is $\frac{1}{n_0} \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2$ -subGaussian with zero mean, while $\frac{1}{n_0/3} \sum_{i=1}^{n_0/3} (\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta} \cdot (\mathbf{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)})$ and $\frac{1}{n_0/3} \sum_{i=1}^{n_0/3} (\mathbf{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} + \boldsymbol{\beta}^{(k)}) (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]}$ are $\frac{1}{n_0/3} \|\boldsymbol{\beta}^{(k)}\|_2 \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2$ -subexponential with mean at most $C \|\boldsymbol{\beta}^{(k)}\|_2 \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2$ (Vershynin, 2018), where $C > 0$ is a uniform constant, then by tail bounds and union bounds, the conclusion follows.

(iii) For Poisson regression model with a.s. bounded covariates, it holds that

$$\begin{aligned} |\hat{L}_0^{[r]}(\hat{\boldsymbol{\beta}}^{(k)[r]}) - \hat{L}_0(\boldsymbol{\beta}^{(k)})| &\leq \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} y_i^{(0)[r]} (\mathbf{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}) \right| \\ &\quad + \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} [e^{(\mathbf{x}_i^{(0)[r]})^T \hat{\boldsymbol{\beta}}^{(k)[r]}} - e^{(\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}^{(k)}}] \right|. \end{aligned}$$

Conditioning on $\mathbf{X}^{(0)[r]}$, we know that $y_i^{(0)[r]}(\mathbf{x}_i^{(0)[r]})^T(\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}) \sim \text{Poisson}(e^{(\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}})$. $(\mathbf{x}_i^{(0)[r]})^T(\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)})$. By the fact that $\text{Poisson}(e^{(\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}})$ is a $e^{(\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}}$ -subexponential given $\mathbf{x}_i^{(0)}$, we have

$$\begin{aligned} \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} y_i^{(0)} (\mathbf{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}) \right| &\lesssim \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} e^{(\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}} (\mathbf{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}) \right| \\ &+ \frac{1}{\sqrt{n_0/3}} \max_i e^{(\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}} \sqrt{\sum_{i=1}^{n_0/3} |(\mathbf{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)})|^2} \cdot \zeta, \end{aligned}$$

with probability at least $1 - \exp\{-\zeta^2\}$. By Hölder inequality, the first term on the right hand side can be bounded by a $e^{2U\|\boldsymbol{\beta}\|_1} \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2^2$ -subexponential with mean at most $Ce^{2U\|\boldsymbol{\beta}\|_1} \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2^2$, leading to

$$\frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} e^{(\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}} (\mathbf{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}) \right| \lesssim \left(1 + \sqrt{\frac{1}{n_0}} \cdot \zeta\right) \cdot e^{U\|\boldsymbol{\beta}\|_1} \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2,$$

with probability at least $1 - \exp\{-\zeta^2\}$. On the other hand, by applying Bernstein inequality (Theorem 2.8.2 in [Vershynin \(2018\)](#)) as well as union bounds, we have

$$\begin{aligned} \frac{1}{\sqrt{n_0/3}} \max_i e^{(\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}} \sqrt{\sum_{i=1}^{n_0/3} |(\mathbf{x}_i^{(0)[r]})^T (\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)})|^2} &\lesssim_p \left(1 + \sqrt{\frac{1}{n_0}} \cdot \zeta\right) \\ &\cdot e^{U\|\boldsymbol{\beta}\|_1} \sup_k \|\hat{\boldsymbol{\beta}}^{(k)[r]} - \boldsymbol{\beta}^{(k)}\|_2, \end{aligned}$$

with probability at least $1 - \exp\{-\zeta^2\}$. Summarizing the conclusions before, we obtain the desired conclusion. \square

Lemma 8. *Under the same assumptions as Theorem 4, we have the following conclusions:*

(i) *For logistic regression model:*

$$|\hat{L}_0^{[r]}(\boldsymbol{\beta}) - L_0(\boldsymbol{\beta})| \vee |\hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)}) - \hat{L}_0^{[r]}(\boldsymbol{\beta}) - L_0(\boldsymbol{\beta}^{(k)}) + L_0(\boldsymbol{\beta})| \leq C \sqrt{\frac{1}{n_0}} \cdot \|\mathbf{w}^{(k)}\|_2 \cdot \zeta,$$

with probability at least $1 - \exp\{-\zeta^2\}$.

(ii) *For linear model:*

$$\begin{aligned} & |\hat{L}_0^{[r]}(\boldsymbol{\beta}) - L_0(\boldsymbol{\beta})| \vee |\hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)}) - \hat{L}_0^{[r]}(\boldsymbol{\beta}) - L_0(\boldsymbol{\beta}^{(k)}) + L_0(\boldsymbol{\beta})| \\ & \leq C \sqrt{\frac{1}{n_0}} \cdot (\|\mathbf{w}^{(k)}\|_2^2 \vee \|\mathbf{w}^{(k)}\|_2) \cdot \zeta, \end{aligned}$$

with probability at least $1 - \exp\{-\zeta^2\}$.

(iii) *For Poisson regression model with $\sup_k \|\mathbf{x}^{(k)}\|_\infty \leq U$ a.s.:*

$$\begin{aligned} & |\hat{L}_0^{[r]}(\boldsymbol{\beta}) - L_0(\boldsymbol{\beta})| \vee |\hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)}) - \hat{L}_0^{[r]}(\boldsymbol{\beta}) - L_0(\boldsymbol{\beta}^{(k)}) + L_0(\boldsymbol{\beta})| \\ & \leq C \sqrt{\frac{1}{n_0}} \exp(U \|\mathbf{w}^{(k)}\|_1) [1 + \|\mathbf{w}^{(k)}\|_2 + U \|\mathbf{w}^{(k)}\|_1] \cdot \zeta, \end{aligned}$$

with probability at least $1 - \zeta^{-2}$.

Proof of Lemma 8. Similar to the proofs of Lemmas 6 and 7, the terms involving $\sum_{i=1}^{n_0/3} \rho(y_i^{(0)[r]})$ and $\mathbb{E}[\rho(y^{(0)})]$ are canceled when taking the difference. Therefore without loss of generality, to prove the rate of $\sup_k |\hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)}) - \hat{L}_0^{[r]}(\boldsymbol{\beta}) - L_0(\boldsymbol{\beta}^{(k)}) + L_0(\boldsymbol{\beta})|$, throughout this proof, we discard these terms and consider

$$\begin{aligned} L_0(\mathbf{w}) &= -\mathbb{E}[\psi'(\boldsymbol{\beta}^T \mathbf{x}^{(0)}) \mathbf{w}^T \mathbf{x}^{(0)}] + \mathbb{E}[\psi(\mathbf{w}^T \mathbf{x}^{(0)})], \\ \hat{L}_0^{[r]}(\mathbf{w}) &= -\frac{1}{n_0/3} \sum_{i=1}^{n_0/3} y_i^{(0)[r]} \boldsymbol{\beta}^T \mathbf{x}_i^{(0)[r]} + \frac{1}{n_0/3} \sum_{i=1}^{n_0/3} \psi(\mathbf{w}^T \mathbf{x}_i^{(0)[r]}). \end{aligned}$$

(i) For logistic regression model:

$$\begin{aligned}
& \sup_k \left| \hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)}) - L_0(\boldsymbol{\beta}^{(k)}) \right| \\
& \leq \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} [y_i^{(0)[r]} - \psi'(\boldsymbol{\beta}^T \mathbf{x}_i^{(0)[r]})](\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]} \right| \\
& \quad + \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} \left[\psi'(\boldsymbol{\beta}^T \mathbf{x}_i^{(0)[r]})(\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]} - \mathbb{E}[\psi'(\boldsymbol{\beta}^T \mathbf{x}^{(0)})(\boldsymbol{\beta}^{(k)})^T \mathbf{x}^{(0)}] \right] \right| \\
& \quad + \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} \left[\psi((\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]}) - \mathbb{E}\psi((\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]}) \right] \right|
\end{aligned}$$

Since $\frac{1}{n_0/3} \sum_{i=1}^{n_0/3} [y_i^{(0)[r]} - \psi'(\boldsymbol{\beta}^T \mathbf{x}_i^{(0)[r]})](\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]}$ is a zero-mean $\|\boldsymbol{\beta}^{(k)}\|_2^2$ -subexponential variable, we have

$$\frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} [y_i^{(0)[r]} - \psi'(\boldsymbol{\beta}^T \mathbf{x}_i^{(0)[r]})](\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]} \right| \lesssim \sqrt{\frac{1}{n_0}} \|\mathbf{w}^{(k)}\|_2 \cdot \zeta, \quad (\text{S.2.43})$$

with probability at least $1 - \exp\{-\zeta^2\}$. For the second term, since ψ' is bounded, $\psi'(\boldsymbol{\beta}^T \mathbf{x}_i^{(0)[r]})(\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]}$ are i.i.d. $\|\boldsymbol{\beta}^{(k)}\|_2^2$ -subexponential, leading to

$$\frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} \left[\psi'(\boldsymbol{\beta}^T \mathbf{x}_i^{(0)[r]})(\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]} - \mathbb{E}[\psi'(\boldsymbol{\beta}^T \mathbf{x}^{(0)})(\boldsymbol{\beta}^{(k)})^T \mathbf{x}^{(0)}] \right] \right| \lesssim \sqrt{\frac{1}{n_0}} \|\mathbf{w}^{(k)}\|_2 \cdot \zeta, \quad (\text{S.2.44})$$

with probability at least $1 - \exp\{-\zeta^2\}$. For the last term, consider $g(u_1^{[r]}, \dots, u_{n_0/3}^{[r]}) = \sum_{i=1}^{n_0/3} \psi(u_i^{[r]})$, where $u_i^{[r]} = (\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}^{(k)}$ is an i.i.d. $\|\boldsymbol{\beta}^{(k)}\|_2^2$ -subGaussian. Since ψ is 1-Lipschitz under ℓ_1 -norm, by Theorem 1 in [Kontorovich \(2014\)](#) and union bounds,

$$\frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} \left[\psi((\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]}) - \mathbb{E}\psi((\boldsymbol{\beta}^{(k)})^T \mathbf{x}^{(0)}) \right] \right| \lesssim \sqrt{\frac{1}{n_0}} \cdot \|\mathbf{w}^{(k)}\|_2 \cdot \zeta, \quad (\text{S.2.45})$$

with probability at least $1 - \exp\{-\zeta^2\}$. By (S.2.43), (S.2.44) and (S.2.45), the conclusion follows.

(ii) For linear model: recall that $y_i^{(0)} = (\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}^{(0)} + \epsilon_i^{(0)[r]}$, then

$$\begin{aligned} & \left| \hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)}) - L_0(\boldsymbol{\beta}^{(k)}) \right| \\ & \leq \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} \epsilon_i^{(0)[r]} (\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta} \right| \\ & \quad + \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} (\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta} \cdot (\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}^{(k)} - \mathbb{E}[(\mathbf{x}^{(0)})^T \boldsymbol{\beta} \cdot (\mathbf{x}^{(0)})^T \boldsymbol{\beta}^{(k)}] \right| \\ & \quad + \frac{1}{2n_0/3} \left| \sum_{i=1}^{n_0/3} [(\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}^{(k)}]^2 - \mathbb{E}[(\mathbf{x}^{(0)})^T \boldsymbol{\beta}^{(k)}]^2 \right|. \end{aligned}$$

By subexponential tail bounds, we have

$$\frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} \epsilon_i^{(0)[r]} (\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta} \right| \lesssim \sqrt{\frac{1}{n_0}} \cdot \|\boldsymbol{\beta}^{(k)}\|_2 \cdot \zeta \lesssim \sqrt{\frac{1}{n_0}} \cdot \|\mathbf{w}^{(k)}\|_2 \cdot \zeta,$$

$$\begin{aligned} & \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} (\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}^{(0)} \cdot (\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}^{(k)} - \mathbb{E}[(\mathbf{x}^{(0)})^T \boldsymbol{\beta}^{(0)} \cdot (\mathbf{x}^{(0)})^T \boldsymbol{\beta}^{(k)}] \right| \\ & \lesssim \sqrt{\frac{1}{n_0}} \cdot \|\boldsymbol{\beta}\|_2 \sup_k \|\boldsymbol{\beta}^{(k)}\|_2 \cdot \zeta \lesssim \sqrt{\frac{1}{n_0}} \cdot \|\boldsymbol{\beta}\|_2 \sup_k \|\mathbf{w}^{(k)}\|_2, \end{aligned}$$

$$\frac{1}{2n_0/3} \left| \sum_{i=1}^{n_0/3} [(\mathbf{x}_i^{(0)[r]})^T \boldsymbol{\beta}^{(k)}]^2 - \mathbb{E}[(\mathbf{x}^{(0)})^T \boldsymbol{\beta}^{(k)}]^2 \right| \lesssim \sqrt{\frac{1}{n_0}} \cdot \|\boldsymbol{\beta}^{(k)}\|_2^2 \cdot \zeta,$$

with probability at least $1 - \exp\{-\zeta^2\}$, which leads to the desired conclusion.

(iii) For Poisson regression model: similar to the logistic regression model, it holds that

$$\begin{aligned}
& \left| \hat{L}_0^{[r]}(\boldsymbol{\beta}^{(k)}) - L_0(\boldsymbol{\beta}^{(k)}) \right| \\
& \leq \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} [y_i^{(0)[r]} - e^{\boldsymbol{\beta}^T \mathbf{x}_i^{(0)[r]}}] (\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]} \right| \\
& \quad + \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} \left[e^{\boldsymbol{\beta}^T \mathbf{x}_i^{(0)[r]}} (\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]} - \mathbb{E}[e^{\boldsymbol{\beta}^T \mathbf{x}^{(0)}} (\boldsymbol{\beta}^{(k)})^T \mathbf{x}^{(0)}] \right] \right| \\
& \quad + \frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} \left[e^{(\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]}} - \mathbb{E} e^{(\boldsymbol{\beta}^{(k)})^T \mathbf{x}^{(0)}} \right] \right|.
\end{aligned}$$

For the first term on the right-hand side, because $[y_i^{(0)[r]} - e^{\boldsymbol{\beta}^T \mathbf{x}_i^{(0)[r]}}] (\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]}$ is an i.i.d. zero-mean $e^{U\|\boldsymbol{\beta}\|_1} \|\boldsymbol{\beta}^{(k)}\|_2$ -subexponential variable, it follows

$$\frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} [y_i^{(0)[r]} - e^{\boldsymbol{\beta}^T \mathbf{x}_i^{(0)[r]}}] (\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]} \right| \lesssim \sqrt{\frac{1}{n_0}} e^{U\|\boldsymbol{\beta}\|_1} \|\boldsymbol{\beta}^{(k)}\|_2 \cdot \zeta,$$

with probability at least $1 - \exp\{-\zeta^2\}$. For the last term on the right-hand side, note that $\exp\{(\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]}\}$ is bounded by $\exp\{U\|\boldsymbol{\beta}^{(k)}\|_1\}$. Therefore by tail probability,

$$\frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} \left[e^{(\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]}} - \mathbb{E} e^{(\boldsymbol{\beta}^{(k)})^T \mathbf{x}^{(0)}} \right] \right| \lesssim \exp\{U\|\boldsymbol{\beta}^{(k)}\|_1\} \cdot \sqrt{\frac{1}{n_0}} \cdot \zeta,$$

with probability at least $1 - \exp\{-\zeta^2\}$. Denote $u_i^{(k)[r]} = (\boldsymbol{\beta}^{(k)})^T \mathbf{x}_i^{(0)[r]}$. Finally, to bound the second term on the right-hand side, we follow the same idea to get

$$\frac{1}{n_0/3} \left| \sum_{i=1}^{n_0/3} u_i^{(k)[r]} e^{u_i^{(0)[r]}} - \mathbb{E} u_i^{(k)[r]} e^{u_i^{(0)[r]}} \right| \lesssim U\|\boldsymbol{\beta}^{(k)}\|_1 \cdot \exp\{U\|\boldsymbol{\beta}^{(k)}\|_1\} \cdot \sqrt{\frac{1}{n_0}} \cdot \zeta,$$

with probability at least $1 - \exp\{-\zeta^2\}$. By combining all the conclusions above, we obtain the desired bound.

The remaining task is to calculate the rate of $|\hat{L}_0^{[r]}(\boldsymbol{\beta}) - L_0(\boldsymbol{\beta})|$ under three scenarios. For logistic case, since $\rho = 0$, the calculation in (i) naturally follows and the same bound can be derived. For the linear case, we only have to show that

$$\left| \frac{1}{n_0/3} \sum_{i=1}^{n_0/3} (y_i^{(0)[r]})^2 - \mathbb{E}(y^{(0)})^2 \right| \lesssim \sqrt{\frac{1}{n_0}} \cdot (\|\boldsymbol{\beta}^{(k)}\|_2^2 \vee \|\boldsymbol{\beta}^{(k)}\|_2),$$

with probability at least $1 - \exp\{-n_0\}$. This can be easily checked by considering $y_i^{(0)[r]} = \boldsymbol{\beta}^T \mathbf{x}_i^{(0)[r]} + \epsilon_i^{(0)[r]}$ and applying subexponential tail bounds. For Poisson regression model, notice that

$$\text{Var}(\log(y_i^{(0)[r]}!)) \leq \mathbb{E}[(y_i^{(0)[r]})^2 \log^2 y_i^{(0)[r]}] \leq \sqrt{\mathbb{E}[(y_i^{(0)[r]})^4]} \sqrt{\mathbb{E}(\log^4 y_i^{(0)[r]})}.$$

Due to moment bounds of subexponential variables and Jensen's inequality:

$$\begin{aligned} \mathbb{E}[(y_i^{(0)[r]})^4] &\lesssim \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{y^{(0)}|\mathbf{x}} (y_i^{(0)[r]})^4 \right] \leq \exp(4U\|\boldsymbol{\beta}\|_1), \\ \mathbb{E}(\log^4 y_i^{(0)[r]}) &\leq \log^4 \mathbb{E} y_i^{(0)[r]} \lesssim U^4 \|\boldsymbol{\beta}\|_1^4. \end{aligned}$$

Then by Chebyshev inequality and union bounds, it's straightforward to prove that

$$\left| \frac{1}{n_0/3} \sum_{i=1}^{n_0/3} y_i^{(0)[r]} - \mathbb{E} y^{(0)} \right| \lesssim_p \sqrt{\frac{1}{n_0}} U \|\boldsymbol{\beta}\|_1 \exp(U\|\boldsymbol{\beta}\|_1) \cdot \zeta,$$

with probability at least $1 - \zeta^{-2}$, which completes our proof. □

S.2.4.2 Proof of Proposition 2

As analyzed in the proof of Theorem 5, it remains to show (S.2.37). Recall that $\hat{\tau}_j^2 = \hat{\Sigma}_{\hat{\boldsymbol{\beta}}, j, j} - \hat{\Sigma}_{\hat{\boldsymbol{\beta}}, j, -j} \hat{\gamma}_j^{(0)} = (n_{\mathcal{A}_h} + n_0)^{-1} \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}_j^{(k)})^T \text{diag}(\{\psi''((\mathbf{x}_i^{(k)})^T \hat{\boldsymbol{\beta}})\}_{i=1}^{n_k}) \mathbf{X}_j^{(k)} - (n_{\mathcal{A}_h} +$

$n_0)^{-1} \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}_j^{(k)})^T \text{diag}(\{\psi''((\mathbf{x}_i^{(k)})^T \hat{\boldsymbol{\beta}})\}_{i=1}^{n_k}) \mathbf{X}_{-j}^{(k)} \hat{\boldsymbol{\gamma}}_j^{(0)}$, $\tau_j^2 = \boldsymbol{\Sigma}_{\boldsymbol{\beta}, j, j}^{(0)} - \boldsymbol{\Sigma}_{\boldsymbol{\beta}, j, -j}^{(0)} \boldsymbol{\gamma}_j^{(0)} = \mathbb{E}[(x_j^{(0)})^2 \psi''((\mathbf{x}^{(0)})^T \boldsymbol{\beta})] - \mathbb{E}[x_j^{(0)} \psi''((\mathbf{x}^{(0)})^T \boldsymbol{\beta}) (x_{-j}^{(0)})^T] \boldsymbol{\gamma}_j^{(0)} = (\boldsymbol{\Theta}_{j, j})^{-1}$ stays away from zero because of Assumption 7.(iii). And we have

$$\begin{aligned}
& \hat{\tau}_j^2 - \tau_j^2 \\
&= \underbrace{\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i, k} (x_{ij}^{(k)})^2 [\psi''((\mathbf{x}_i^{(k)})^T \hat{\boldsymbol{\beta}}) - \psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta})]}_{(1)} \\
&+ \underbrace{\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i, k} (x_{ij}^{(k)})^2 \psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta}) - \sum_{k \in \{0\} \cup \mathcal{A}_h} \frac{n_k}{n_{\mathcal{A}_h} + n_0} \mathbb{E}[(x_{ij}^{(k)})^2 \psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta})]}_{(2)} \\
&+ \underbrace{\sum_{k \in \{0\} \cup \mathcal{A}_h} \frac{n_k}{n_{\mathcal{A}_h} + n_0} \left\{ \mathbb{E}[(x_{ij}^{(k)})^2 \psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta})] - \mathbb{E}[(x_{ij}^{(0)})^2 \psi''((\mathbf{x}_i^{(0)})^T \boldsymbol{\beta})] \right\}}_{(3)} \\
&+ \underbrace{\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}_j^{(k)})^T \text{diag}(\{\psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta}) - \psi''((\mathbf{x}_i^{(k)})^T \hat{\boldsymbol{\beta}})\}_{i=1}^{n_k}) \mathbf{X}_{-j}^{(k)} \hat{\boldsymbol{\gamma}}_j^{(0)}}_{(4)} \\
&- \underbrace{\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} (\mathbf{X}_j^{(k)})^T \text{diag}(\{\psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta})\}_{i=1}^{n_k}) \mathbf{X}_{-j}^{(k)} (\hat{\boldsymbol{\gamma}}_j^{(0)} - \boldsymbol{\gamma}_j^{(0)})}_{(5)} \\
&+ \underbrace{\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}_h} \left\{ -(\mathbf{X}_j^{(k)})^T \text{diag}(\{\psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta})\}_{i=1}^{n_k}) \mathbf{X}_{-j}^{(k)} \boldsymbol{\gamma}_j^{(0)} + \mathbb{E}[x_j^{(k)} \psi''((\mathbf{x}^{(k)})^T \boldsymbol{\beta}) (x_{-j}^{(k)})^T] \boldsymbol{\gamma}_j^{(0)} \right\}}_{(6)} \\
&+ \underbrace{\sum_{k \in \{0\} \cup \mathcal{A}_h} \frac{n_k}{n_{\mathcal{A}_h} + n_0} \left\{ -\mathbb{E}[x_j^{(k)} \psi''((\mathbf{x}^{(k)})^T \boldsymbol{\beta}) (x_{-j}^{(k)})^T] + \mathbb{E}[x_j^{(0)} \psi''((\mathbf{x}^{(0)})^T \boldsymbol{\beta}) (x_{-j}^{(0)})^T] \right\} \boldsymbol{\gamma}_j^{(0)}}_{(7)}.
\end{aligned}$$

And we have the following control for each term:

$$\begin{aligned}
|(1)| &\leq \frac{1}{n_{\mathcal{A}_h} + n_0} \left| \sum_{i,k} (x_{ij}^{(k)})^2 \psi'''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta} + t(\mathbf{x}_i^{(k)})^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})) (\mathbf{x}_i^{(k)})^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})) \right| \\
&\leq \sqrt{\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i,k} (x_{ij}^{(k)})^4 [\psi'''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta} + t(\mathbf{x}_i^{(k)})^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}))]^2} \cdot \sqrt{\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i,k} [(\mathbf{x}_i^{(k)})^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})]^2} \\
&\lesssim \mathfrak{R}_1,
\end{aligned}$$

with probability at least $1 - n_0^{-1}$.

Since $\{(x_{ij}^{(k)})^2 \psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta})\}_{i,k}$ are independent sub-Gaussian variables with finite variance, by concentration inequality, $|(2)| \lesssim (n_{\mathcal{A}_h} + n_0)^{-1/2}$ with probability $1 - \exp\{-(n_{\mathcal{A}_h} + n_0)\}$. Similarly, $|(6)| \lesssim (n_{\mathcal{A}_h} + n_0)^{-1/2}$ with probability $1 - \exp\{-(n_{\mathcal{A}_h} + n_0)\}$.

$$|(3)| \leq \sup_{k \in \mathcal{A}_h} \left| \boldsymbol{\Sigma}_{\boldsymbol{\beta},j,j}^{(k)} - \boldsymbol{\Sigma}_{\boldsymbol{\beta},j,j}^{(0)} \right| \leq h_{\max}.$$

$$\begin{aligned}
|(4)| &\leq \left| \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i,k} x_{ij}^{(k)} (\mathbf{x}_{-j}^{(k)})^T \hat{\boldsymbol{\gamma}}_j^{(0)} \cdot [\psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta}) - \psi''((\mathbf{x}_i^{(k)})^T \hat{\boldsymbol{\beta}})] \right| \\
&\leq \left| \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i,k} x_{ij}^{(k)} (\mathbf{x}_{-j}^{(k)})^T \boldsymbol{\gamma}_j^{(0)} \cdot \psi'''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta} + t(\mathbf{x}_i^{(k)})^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})) \cdot (\mathbf{x}_i^{(k)})^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \right| \\
&\quad + \left| \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i,k} x_{ij}^{(k)} (\mathbf{x}_{-j}^{(k)})^T (\hat{\boldsymbol{\gamma}}_j^{(0)} - \boldsymbol{\gamma}_j^{(0)}) \cdot \psi'''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta} + t(\mathbf{x}_i^{(k)})^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})) \cdot (\mathbf{x}_i^{(k)})^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \right| \\
&\lesssim \mathfrak{R}_1,
\end{aligned}$$

with probability at least $1 - n_0^{-1}$.

$$\begin{aligned}
|(5)| &\leq \left| \frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i,k} \mathbf{x}_{ij}^{(k)} \psi''((\mathbf{x}_i^{(k)})^T \boldsymbol{\beta}) \cdot (\mathbf{x}_{i,-j}^{(k)})^T (\hat{\boldsymbol{\gamma}}_j^{(0)} - \boldsymbol{\gamma}_j^{(0)}) \right| \\
&\lesssim \sqrt{\frac{1}{n_{\mathcal{A}_h} + n_0} \sum_{i,k} [(\mathbf{x}_{i,-j}^{(k)})^T (\hat{\boldsymbol{\gamma}}_j^{(0)} - \boldsymbol{\gamma}_j^{(0)})]^2} \\
&\lesssim h_1^{1/2} \left(\frac{\log p}{n_0} \right)^{1/4} + h_1^{1/2} \mathfrak{R}_1^{1/2} + \mathfrak{R}_1
\end{aligned}$$

with probability at least $1 - K_{\mathcal{A}_h} n_0^{-1}$ by (S.2.35).

$$|(7)| \lesssim \sup_{k \in \mathcal{A}_h} \left| (\boldsymbol{\Sigma}_{\boldsymbol{\beta},-j,j}^{(k)} - \boldsymbol{\Sigma}_{\boldsymbol{\beta},-j,j}^{(0)}) \boldsymbol{\gamma}_j^{(0)} \right| \lesssim h_{\max}.$$

Combine all the inequalities above to finish the proof of (S.2.37). Note that the bound of $|\hat{\tau}_j^{-2} - \tau_j^{-2}|$ follows because $\inf_j \tau_j^2 = \mathcal{O}(1)$.