

Supplement to “Distributed Estimation and Inference for Semi-parametric Binary Response Models”

A Extensions: Data Heterogeneity and High-dimensional Settings

In this section, we discuss two natural extensions of proposed estimators (Avg-SMSE) and (mSMSE).

First, an important question in distributed environments is the heterogeneity of datasets on different local machines. In Section A.1, we consider heterogeneous datasets with a shared parameter of interest β^* and different distributions of covariates (X, \mathbf{Z}) across the machines. We show that (mSMSE) performs better than (Avg-SMSE) with a weaker condition. While handling heterogeneous data, the performance of the divide-and-conquer algorithm relies on the smallest sample size among the local machines, while the (mSMSE) method does not rely on such conditions.

We further consider a high-dimension extension in Section A.2 where the parameter of interest $\beta^* \in \mathbb{R}^p$ is a sparse vector with s non-zero elements and $s < n < p$. We modify (mSMSE) to adapt the idea of the Dantzig Selector (Candes and Tao, 2007) to reach the convergence rate $\sqrt{s}(\log p/n)^{\frac{\alpha}{2\alpha+1}}$ in a distributed environment, which is very close to the minimax optimal rate $(s \log p/n)^{\frac{\alpha}{2\alpha+1}}$ for the linear binary response model established by Feng et al. (2022) in a non-distributed environment. Compared to the low-dimensional settings above, the algorithm in this setting reduces the per-iteration communication cost to $p \times 1$ vectors.

A.1 Data Heterogeneity

Until now, we assumed homogeneity among the data stored on different machines, which means the distribution of $(x_i, \mathbf{z}_i, \epsilon_i)$ are the same for all i . It is of practical interest to consider the heterogeneous setting, since data on different machines may not be identically distributed. Therefore, we establish the theoretical results in the presence of heterogeneity in this section.

First, we remove the restriction that the sample size on each machine is the same. Denote the number of observations on the machine \mathcal{H}_ℓ to be m_ℓ , which satisfies $\sum_{\ell=1}^L m_\ell = n$. Then we have to modify Assumptions 2–4 for different distributions on different machines.

Assumption 1. For X and \mathbf{Z} on \mathcal{H}_ℓ , define $\zeta := X + \mathbf{Z}^\top \boldsymbol{\beta}^*$, and assume that the conditional distribution density function of ζ , denoted by $\rho_\ell(\cdot | \mathbf{Z})$, is positive and bounded for almost every \mathbf{Z} . Further, for any integer $1 \leq k \leq \alpha$, assume that within a neighborhood of 0, $\rho_\ell^{(k)}(\cdot | \mathbf{Z})$ exists and is uniformly bounded for all ℓ and almost every \mathbf{Z} , i.e., $\exists M_{\rho,k}$ such that $\sup_{\zeta, \ell} \left| \rho_\ell^{(k)}(\zeta | \mathbf{Z}) \right| \leq M_{\rho,k}$.

Assumption 2. For X, \mathbf{Z}, ϵ on \mathcal{H}_ℓ , let $F_\ell(\cdot | \mathbf{Z})$ denote the conditional distribution function of the noise ϵ , and assume that ϵ and X are independent given \mathbf{Z} . For any integer $1 \leq k \leq \alpha + 1$, assume that $F_\ell^{(k)}(\cdot | \mathbf{Z})$ exists and is uniformly bounded within a neighborhood of 0 for all ℓ and almost every \mathbf{Z} , i.e., $\exists M_{F,k}$ such that $\sup_{\epsilon, \ell} \left| F_\ell^{(k)}(\epsilon | \mathbf{Z}) \right| \leq M_{F,k}$. Still, we assume $\text{median}(\epsilon | \mathbf{Z}) = 0$ on each machine.

Assumption 3. There exist constants $c_1, c_2 > 0$ such that $c_1^{-1} < \Lambda_{\min}(V_\ell) < \Lambda_{\max}(V_\ell) < c_1$, $c_2^{-1} < \Lambda_{\min}(V_{s,\ell}) < \Lambda_{\max}(V_{s,\ell}) < c_2$, $\forall \ell$, where $V_\ell := 2\mathbb{E}_{\mathbf{Z} \in \mathcal{H}_\ell}(\rho_\ell(0 | \mathbf{Z}) F_\ell'(0 | \mathbf{Z}) \mathbf{Z} \mathbf{Z}^\top)$ and $V_{s,\ell} := \pi_V \mathbb{E}_{\mathbf{Z} \in \mathcal{H}_\ell}(\rho_\ell(0 | \mathbf{Z}) \mathbf{Z} \mathbf{Z}^\top)$.

Assumptions 1–3 are parallel to Assumptions 2–4, requiring the uniform boundedness of the high-order derivatives and the eigenvalues of V_ℓ and $V_{s,\ell}$. Additionally, similar to (10), we define

$$U_\ell := \pi_U \mathbb{E}_{\mathbf{Z} \in \mathcal{H}_\ell} \left(\sum_{k=1}^{\alpha} \frac{2(-1)^{k+1}}{k! (\alpha - k)!} F_\ell^{(k)}(0 | \mathbf{Z}) \rho_\ell^{(\alpha-k)}(0 | \mathbf{Z}) \mathbf{Z} \right),$$

which is related to the bias of SMSE on each machine.

Under the modified assumptions, the data on each machine are no longer identically distributed, and therefore it is natural to allocate a different weight matrix W_ℓ to each machine, with $\sum_{\ell=1}^L W_\ell = I_{p \times p}$. Formally, the *weighted-Averaged SMSE* (**wAvg-SMSE**) is defined as follows:

$$\widehat{\boldsymbol{\beta}}_{(\text{wAvg-SMSE})} := \sum_{\ell=1}^L W_\ell \widehat{\boldsymbol{\beta}}_{\text{SMSE}, \ell}, \quad (1)$$

where $\widehat{\boldsymbol{\beta}}_{\text{SMSE}, \ell}$ is the SMSE on the ℓ -th machine that minimizes the objective function

$$F_{h,\ell}(\boldsymbol{\beta}) := \frac{1}{m_\ell} \sum_{i \in \mathcal{H}_\ell} (-y_i) H \left(\frac{x_i + \mathbf{z}_i^\top \boldsymbol{\beta}}{h} \right) \mathbf{z}_i \mathbf{z}_i^\top. \quad (2)$$

For the multi-round method, we aim to apply the iterative smoothing to minimize a weighted sum of the objective functions on each machine, i.e., $\sum_{\ell=1}^L W_\ell F_{h,\ell}(\beta)$, which leads to updating the *weighted* mSMSE (**wmSMSE**) in the t -th step by

$$\hat{\beta}_{(\text{wmSMSE})}^{(t)} = \hat{\beta}_{(\text{wmSMSE})}^{(t-1)} - \left(\sum_{\ell=1}^L W_\ell \nabla^2 F_{h,\ell}(\hat{\beta}_{(\text{wmSMSE})}^{(t-1)}) \right)^{-1} \left(\sum_{\ell=1}^L W_\ell \nabla F_{h,\ell}(\hat{\beta}_{(\text{wmSMSE})}^{(t-1)}) \right). \quad (3)$$

A natural choice of weights is proportional to the local sample size, i.e., $W_\ell = \frac{m_\ell}{n} I_{p \times p}$. Using this weight, the variances in the asymptotic distribution in (9) and (14) become $\frac{1}{n} \sum_{\ell=1}^L m_\ell V_\ell^{-1} V_{s,\ell} V_\ell^{-1}$ and $\left(\frac{1}{n} \sum_{\ell=1}^L m_\ell V_\ell \right)^{-1} \left(\frac{1}{n} \sum_{\ell=1}^L m_\ell V_{s,\ell} \right) \left(\frac{1}{n} \sum_{\ell=1}^L m_\ell V_\ell \right)^{-1}$, respectively, which can be seen as a special case of Theorems A.1 and A.2 below. Nonetheless, such a choice is by no means optimal. We could further decrease both asymptotic variances by choosing a different weight matrix W_ℓ for each machine. To illustrate the choice of weights, we first derive the theoretical results for general weight matrices W_ℓ that satisfy the following restriction:

Assumption 4. *There exist constants $c_w, C_W > 0$ such that $c_w m_\ell / n \leq \|W_\ell\|_2 \leq C_W m_\ell / n$, with $n = \sum_{\ell=1}^L m_\ell$ and $\sum_{\ell=1}^L W_\ell = I_{p \times p}$.*

Assumption 4 requires that the 2-norm of W_ℓ is not too far away from m_ℓ / n , violation of which may lead to a low convergence rate. An extreme example is $W_1 = I_{p \times p}$ and $W_2 = \dots = W_L = \mathbf{0}$, in which case only the data on a single machine will be used. Following the procedures in Sections 3.1 and 3.2, we establish the asymptotic normality for (**wAvg-SMSE**) and (**wmSMSE**) in Theorems A.1 and A.2, whose proof is given in Section C.4 in supplementary material.

Theorem A.1 (**wAvg-SMSE**). *Suppose Assumptions 1 and 5-4 hold and the sample size of the smallest local batch $\min_\ell m_\ell \gtrsim n^{3/(2\alpha+1)}$. By taking $h = n^{-1/(2\alpha+1)}$, we have*

$$n^{\frac{\alpha}{2\alpha+1}} \left(\hat{\beta}_{(\text{wAvg-SMSE})} - \beta^* \right) \xrightarrow{d} \mathcal{N} \left(\sum_{\ell=1}^L W_\ell V_\ell^{-1} U_\ell, \sum_{\ell=1}^L \frac{n}{m_\ell} W_\ell V_\ell^{-1} V_{s,\ell} V_\ell^{-1} W_\ell^\top \right). \quad (4)$$

Theorem A.2 (**wmSMSE**). *Suppose Assumptions 1 and 5-4 hold and $\|\hat{\beta}^{(0)} - \beta^*\|_2 = O_{\mathbb{P}}(\delta_0)$. By taking $h_t = \max \{n^{-1/(2\alpha+1)}, \delta_0^{2^t/\alpha}\}$ at iteration $t = 1, 2, \dots, T$, we have*

$$n^{\frac{\alpha}{2\alpha+1}} \left(\hat{\beta}_{(\text{wmSMSE})}^{(T)} - \beta^* \right) \xrightarrow{d} \mathcal{N} \left(\bar{V}_W^{-1} \bar{U}_W, \bar{V}_W^{-1} \sum_{\ell=1}^L \frac{n}{m_\ell} W_\ell V_{s,\ell} W_\ell^\top \bar{V}_W^{-1} \right), \quad (5)$$

for sufficiently large T , where $\bar{U}_W := \sum_{\ell=1}^L W_\ell U_\ell$, $\bar{V}_W := \sum_{\ell=1}^L W_\ell V_\ell$.

Assumptions 3 and 4 ensure the variance matrices in both (4) and (5) are finite. In particular, in the homogeneous setting, $W_\ell = (1/L)I_{p \times p}$, then (4) and (5) are identical to (9) and (14).

Remark 1. In Theorem A.1, the condition $\min_\ell m_\ell \gtrsim n^{3/(2\alpha+1)}$ is placed to ensure $\log m_\ell / (m_\ell h^3) = o(1)$ for all ℓ , which is necessary to guarantee the convergence of $\hat{\beta}_{\text{SMSE},\ell}$ to β^* on each machine. In the homogeneous setting, this is equivalent to the constraint $L = o(m^{\frac{2}{3}(\alpha-1)})$ in Theorem 3.1. This condition requires the sample size of the smallest batch should increase at a certain rate as $n \rightarrow \infty$. On the other hand, for (wmSMSE), there is no restriction on the smallest local sample size.

Based on the above results, we are able to artificially choose the weight matrices $\{W_\ell\}$ to minimize the covariance matrices of the two methods in (4) and (5). By choosing $W_\ell^{*,(\text{wAvg-SMSE})} = \left(\sum_{\ell=1}^L m_\ell V_\ell V_{s,\ell}^{-1} V_\ell \right)^{-1} m_\ell V_\ell V_{s,\ell}^{-1} V_\ell$, both the trace and the Frobenius norm of the variance in (4) are minimized, and the corresponding minimum variance is

$$\Sigma_{(\text{wAvg-SMSE})}^* := n \left(\sum_{\ell=1}^L m_\ell V_\ell V_{s,\ell}^{-1} V_\ell \right)^{-1}. \quad (6)$$

For (wmSMSE), if one chooses $W_\ell^{*,(\text{wmSMSE})} = \left(\sum_{\ell=1}^L m_\ell V_\ell V_{s,\ell}^{-1} \right)^{-1} m_\ell V_\ell V_{s,\ell}^{-1}$, the asymptotic variance in (5) will be the same as $\Sigma_{(\text{wAvg-SMSE})}^*$ in (6). Therefore, the multi-round method (wmSMSE) is at least as efficient as (wAvg-SMSE) by specifying certain weight matrices. Note that it is easy to verify the above optimal weights $W_\ell^{*,(\text{wAvg-SMSE})}$ and $W_\ell^{*,(\text{wmSMSE})}$ satisfy Assumption 4. The detailed derivation is given in Section C.4 in supplementary material.

A.2 High-dimensional Multi-round SMSE

In this section, we extend (mSMSE) to high-dimensional settings, where the dimension p is much larger than n . We assume that $\beta^* \in \mathbb{R}^p$ is a sparse vector with s non-zero elements. Recall (7),

$$\hat{\beta}^{(t)} = \hat{\beta}^{(t-1)} - \left[\nabla^2 F_{h_t}(\hat{\beta}^{(t-1)}) \right]^{-1} \nabla F_{h_t}(\hat{\beta}^{(t-1)}) = \hat{\beta}^{(t-1)} - \left(V_n^{(t)} \right)^{-1} U_n^{(t)}. \quad (7)$$

It is generally infeasible to compute the inverse of the Hessian matrix $V_n^{(t)} \in \mathbb{R}^{p \times p}$ in the high-dimensional case. Furthermore, it requires unacceptably high complexity to compute and communicate L high-dimensional matrices. To tackle these problems, we first note that (7) is equivalent to solving the following quadratic optimization problem:

$$\hat{\beta}^{(t)} := \underset{\beta}{\operatorname{argmin}} \quad \frac{1}{2} \beta^\top V_n^{(t)} \beta - \beta^\top \left(V_n^{(t)} \hat{\beta}^{(t-1)} - U_n^{(t)} \right). \quad (8)$$

Due to high communication complexity, we only estimate the Hessian matrix $V_n^{(t)}$ using the samples on a single machine, e.g., $V_{m,1}^{(t)}$ on the first machine. Then (8) can be written as

$$\hat{\beta}^{(t)} := \underset{\beta}{\operatorname{argmin}} \quad \frac{1}{2} \beta^\top V_{m,1}^{(t)} \beta - \beta^\top \left(V_{m,1}^{(t)} \hat{\beta}^{(t-1)} - U_n^{(t)} \right). \quad (9)$$

We adapt the idea of the Dantzig Selector proposed by Candes and Tao (2007), an ℓ_1 -regularization approach known for estimating high-dimensional sparse parameters. Formally, in the t -th iteration, given $\hat{\beta}^{(t-1)}$, the bandwidth h_t and a regularization parameter $\lambda_n^{(t)}$, we compute $\hat{\beta}^{(t)}$ by

$$\hat{\beta}^{(t)} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \|\beta\|_1 : \left\| V_{m,1}^{(t)} \beta - \left(V_{m,1}^{(t)} \hat{\beta}^{(t-1)} - U_n^{(t)} \right) \right\|_\infty \leq \lambda_n^{(t)} \right\}. \quad (10)$$

Note that a feasible solution of (10) can be obtained by linear programming. A complete algorithm is presented in Algorithm 1. It is worthwhile to mention that in Algorithm 1, we do not directly compute the $p \times p$ matrix $V_{m,1}^{(t)}$, but instead compute the vector $V_{m,1}^{(t)} \hat{\beta}^{(t-1)}$ to save computation and storage cost.

Algorithm 1 High-dimensional Multi-round Maximum Score Estimator

Input: Datasets distributed on local machines $\{x_i, z_i, y_i\}_{i \in \mathcal{H}_\ell}$ ($\ell = 1, 2, \dots, L$), an initial estimator $\hat{\beta}^{(0)}$, the total number of iterations T , bandwidth sequence $\{h_t\}$ and parameters $\{\lambda_n^{(t)}\}$.

- 1: **for** $t = 1, 2, \dots, T$ **do**
 - 2: Send $\hat{\beta}^{(t-1)}$ to each machine;
 - 3: **for** $\ell = 1, 2, \dots, L$ **do**
 - 4: Compute $U_{m,\ell}^{(t)}$ by equation (8) and send $U_{m,\ell}^{(t)}$ back to \mathcal{H}_1 ;
 - 5: **end for**
 - 6: Compute $U_n^{(t)} = \frac{1}{L} \sum_{\ell=1}^L U_{m,\ell}^{(t)}$ and $V_{m,1}^{(t)} \hat{\beta}^{(t-1)} = \frac{1}{mh_t^2} \sum_{i \in \mathcal{H}_1} (-y_i) H'' \left(\frac{x_i + z_i^\top \hat{\beta}^{(t-1)}}{h_t} \right) (z_i^\top \hat{\beta}^{(t-1)}) z_i$.
 - 7: Obtain $\hat{\beta}^{(t)}$ by solving (10);
 - 8: **end for**
 - 9: Output $\hat{\beta}^{(T)}$.
-

Now we give the convergence rate of the estimator $\hat{\beta}^{(t)}$ at iteration t .

Theorem A.3. Assume the assumptions in Theorem 3.4 hold. Further, assume that the dimension $p = O(n^\nu)$ for some $\nu > 0$, the sparsity $s = O(m^{1/4})$, and the initial value¹ $\hat{\beta}^{(0)}$ satisfies $\|\hat{\beta}^{(0)} -$

¹An initial estimator $\hat{\beta}^{(0)}$ can be obtained by existing high-dimensional MSE methods on a single machine in literature such as Mukherjee et al. (2019) and Feng et al. (2022).

$\beta^*\|_2 = O_{\mathbb{P}}(\delta_{m,0})$ and $\|\hat{\beta}^{(0)} - \beta^*\|_1 = O_{\mathbb{P}}(\sqrt{s}\delta_{m,0})$, for some $\delta_{m,0} = o(1)$. By choosing proper bandwidth h_t , parameter $\lambda_n^{(t)}$ and kernel function $H(\cdot)$, we can obtain that for $1 \leq t \leq T$,

$$\|\hat{\beta}^{(t)} - \beta^*\|_2 = O_{\mathbb{P}}\left(\sqrt{s}\left(\frac{\log p}{n}\right)^{\frac{\alpha}{2\alpha+1}} + (r_m)^t \delta_{m,0}\right), \quad (11)$$

and $\|\hat{\beta}^{(t)} - \beta^*\|_1 = O_{\mathbb{P}}(\sqrt{s}\|\hat{\beta}^{(t)} - \beta^*\|_2)$, where r_m is an infinitesimal quantity.

Theorem A.3 summarizes the ℓ_2 and ℓ_1 error bounds of $\hat{\beta}^{(t)}$ in Algorithm 1. The details, including the explicit choice for h_t , $\lambda_n^{(t)}$ and $H(\cdot)$ and the formal definition of r_m , are relegated to Section B of supplementary material. The upper bound in (11) contains two terms. The second term comes from the error of the initial estimator, and it decreases exponentially as t increases. As the algorithm operates, this quantity finally gets dominated by the first term within at most $O(\log n)$ iterations. Furthermore, the conditions on m , n , and s in Theorem A.3 are placed to guarantee that the algorithm converges in finite iterations. The first term, $\sqrt{s}(\log p/n)^{\frac{\alpha}{2\alpha+1}}$, represents the statistical convergence rate of our proposed estimator, which is very close to the optimal rate that one can obtain without a distributed environment. Feng et al. (2022) recently established the minimax optimal rate $(s \log p/n)^{\frac{\alpha}{2\alpha+1}}$ of the maximum score estimator in high-dimensional settings. Compared to that, the established rate in Theorem A.3, $\sqrt{s}(\log p/n)^{\frac{\alpha}{2\alpha+1}}$, is slightly slower due to the different designs in the algorithms, with a difference of $s^{\frac{1}{4\alpha+2}}$. Since our proposed algorithm is designed based on the Dantzig Selector that directly controls the infinity norm of the gradient, which is different from the path-following algorithm used in Feng et al. (2022), their techniques cannot be directly applied to our proposed distributed estimator. It would be a potentially interesting future work to improve the estimator in distributed settings to match the optimal rate.

B Theoretical Results of the High-dimensional (mSMSE)

In this section, we give the complete theoretical analysis of the high-dimensional multi-round SMSE in Algorithm 1. First, we restate the conditions as the following two assumptions.

Assumption 5. *The initial value $\hat{\beta}^{(0)}$ satisfies*

$$\|\hat{\beta}^{(0)} - \beta^*\|_2 = O_{\mathbb{P}}(\delta_{m,0}), \quad \|\hat{\beta}^{(0)} - \beta^*\|_1 = O_{\mathbb{P}}(\sqrt{s}\delta_{m,0}), \quad (12)$$

for some $\delta_{m,0} = o(1)$.

Assumption 6. The dimension $p = O(n^\nu)$ for some $\nu > 0$ and the sparsity $s = O(m^r)$ for some $0 < r < 1/4$.

Assumption 5 requires that the error of the initial value can be upper bounded in both ℓ_1 and ℓ_2 norm. Moreover, the bound of the ℓ_1 error is of the same order as \sqrt{s} times the bound of the ℓ_2 error. This could be achieved by the path-following method proposed by Feng et al. (2022), with $\delta_{m,0} = (s \log p / m)^{\alpha / (2\alpha + 1)}$. Assumption 6 restricts the dimension p and the sparsity s . The first constraint $p = O(n^\nu)$ implies that $\log p = O(\log n)$, so we will simply use $\log p$ instead of $\log \max\{n, p\}$ in the following theorems. The sparsity condition $s = O(m^r)$ is also necessary and standard to guarantee convergence. Under these assumptions, we first state the one-step improvement in the following theorem:

Theorem B.1. Assume the assumptions in Theorem 3.4 and Assumptions 5, 6 hold. Further assume that $\sqrt{s}\delta_{m,0} = O(h_1^{3/2})$ and $\frac{s^2 \log p}{mh_1^3} = o(1)$. By taking

$$\lambda_n^{(1)} = C_0 \left(\sqrt{\frac{\log p}{nh_1}} + \sqrt{\frac{s \log p}{mh_1^3}} \delta_{m,0} + s \delta_{m,0}^2 + h_1^\alpha \right),$$

with sufficiently large constant C_0 , we have

$$\|\hat{\beta}^{(1)} - \beta^*\|_2 = O_{\mathbb{P}} \left(\sqrt{s} \lambda_n^{(1)} \right) = O_{\mathbb{P}} \left[\sqrt{\frac{s \log p}{nh_1}} + \sqrt{\frac{s^2 \log p}{mh_1^3}} \delta_{m,0} + s^{3/2} \delta_{m,0}^2 + \sqrt{s} h_1^\alpha \right], \quad (13)$$

and

$$\|\hat{\beta}^{(1)} - \beta^*\|_1 \leq 2\sqrt{s} \|\hat{\beta}^{(1)} - \beta^*\|_2, \quad (14)$$

with probability tending to one.

The proof of Theorem B.1 is given in Section C.5. The additional assumption $\frac{s^2 \log p}{mh_1^3} = o(1)$ is necessary to guarantee the so-called restricted eigenvalue condition of $V_{m,1}^{(1)}$, which is a standard assumption to obtain the convergence rate of the Dantzig Selector in theory. Another assumption $\sqrt{s}\delta_{m,0} = O(h_1^{3/2})$ is a technical condition to determine the dominant term in the convergence rate (See the proof for details). The convergence rate in (13) contains four terms. The sum of the second and the third terms can be rewritten as,

$$\left(\sqrt{\frac{s^2 \log p}{mh_1^3}} + s^{3/2} \delta_{m,0} \right) \delta_{m,0}, \quad (15)$$

which is related to the initial error $\delta_{m,0}$. If we further suppose that $s^{3/2}\delta_{m,0} = o(1)$, then (15) will be $o(\delta_{m,0})$, which can further be iteratively refined in the following iterations. The remaining terms (the first and the fourth terms) can be minimized by letting $\sqrt{\frac{s \log p}{nh_1}} \asymp \sqrt{s}h_1^\alpha$, leading to the optimal bandwidth $h^* := (\log p/n)^{1/(2\alpha+1)}$. Based on these results, we are ready to give the convergence rate of the estimator $\hat{\beta}^{(t)}$ at iteration t .

Theorem B.2. For $t = 1, 2, \dots, T$, define $h^* := (\log p/n)^{1/(2\alpha+1)}$, $r_m := \sqrt{\frac{s^2 \log p}{m(h^*)^3}} + s^{3/2}\delta_{m,0}$ and

$$\delta_{m,t} := \sqrt{s} \left(\frac{\log p}{n} \right)^{\frac{\alpha}{2\alpha+1}} + (r_m)^t \delta_{m,0}.$$

Assume the assumptions in Theorem 3.4 and Assumptions 5,6 hold. Then there exists a constant α_0 such that, by choosing a kernel $H'(\cdot)$ with order $\alpha > \alpha_0$, bandwidth $h_t \equiv h^*$ for $t = 1, 2, \dots, T$ and parameters

$$\lambda_n^{(t)} = C_\lambda \left[\left(\frac{\log p}{n} \right)^{\frac{\alpha}{2\alpha+1}} + \sqrt{\frac{s \log p}{m(h^*)^3}} \delta_{m,t-1} + s \delta_{m,t-1}^2 \right],$$

with a sufficiently large constant C_λ , we have $r_m = o(1)$,

$$\left\| \hat{\beta}^{(t)} - \beta^* \right\|_2 = O_{\mathbb{P}}(\delta_{m,t}) \quad \text{and} \quad \left\| \hat{\beta}^{(t)} - \beta^* \right\|_1 \leq 2\sqrt{s} \left\| \hat{\beta}^{(t)} - \beta^* \right\|_2, \quad (16)$$

with probability tending to one.

The proof of Theorem B.2 is given in Section C.5. The $\delta_{m,t}$ in Theorem B.2 denotes the upper bound on the ℓ_2 -error of $\hat{\beta}^{(t)}$, which contains two terms. The first term $\sqrt{s}(\log p/n)^{\alpha/(2\alpha+1)}$ is the best rate our method can achieve and the second term comes from the error of the initial estimator. As long as $r_m = o(1)$, the second term decreases exponentially as t increases, which implies that $\hat{\beta}^{(t)}$ achieves the optimal rate after at most $O(\log n)$ iterations. Concretely, the number of required iterations is

$$\frac{\alpha \left[(\log n - \log \log n) - \frac{1}{2} \log s + \log \delta_{m,0} \right]}{(2\alpha + 1) [-\log(r_m)]}, \quad (17)$$

which is larger than that in the low-dimensional case. Under the Assumption 6 and the assumption $m > n^c$ in Theorem 3.4, it's easy to see that (17) can be upper bounded by a finite number.

The condition $r_m = o(1)$ can be ensured by choosing a kernel function with order higher than a constant certain α_0 . See Remark 2 for explanation.

Remark 2. To make sure that $r_m = \frac{s^2 \log m}{m(h^*)^3} + s^{3/2}\delta_{m,0} = o(1)$ and $\sqrt{s}\delta_{m,0} = O\left((h^*)^{3/2}\right)$, we need to choose a kernel $H'(\cdot)$ with order α such that $\alpha > \alpha_0 := \max \left\{ \frac{3}{2c(1-2r)} + \frac{r}{2(1-2r)}, \frac{3r}{2(1-4r)} \right\}$, where

$c = (\log m)/(\log n) < 1$ and $r = (\log s)/(\log m) < 1/4$ are supposed in Theorem 3.4 and Assumption 6. Here we plug in the rate $\delta_{m,0} = (s \log p/m)^{\alpha/(2\alpha+1)}$ obtained by the path-following algorithm by Feng et al. (2022). If the order of the kernel is not high enough (less than α_0), Algorithm 1 still works by choosing $h_m = m^{-\frac{2\alpha(1-2r)}{3(2\alpha+1)} + \varepsilon}$ for some small $\varepsilon > 0$, and the corresponding convergence rate will be $\sqrt{s}h_m^\alpha$.

C Technical Proof of the Theoretical Results

C.1 Proof of the Results for the ℓ_2 Error Bound of (mSMSE)

Proof of Proposition 3.2

We first restate Proposition 3.2 in a detailed version. The first step of (mSMSE) can be written as

$$\hat{\beta}^{(1)} - \beta^* = \left(V_{n,h_1}(\hat{\beta}^{(0)}) \right)^{-1} U_{n,h_1}(\hat{\beta}^{(0)}), \quad (18)$$

where

$$V_{n,h}(\beta) = \nabla^2 F_h(\beta) = \frac{1}{nh^2} \sum_{i=1}^n (-y_i) H'' \left(\frac{x_i + z_i^\top \beta}{h} \right) z_i z_i^\top, \quad (19)$$

$$U_{n,h}(\beta) = \frac{1}{nh^2} \sum_{i=1}^n (-y_i) H'' \left(\frac{x_i + z_i^\top \beta}{h} \right) z_i z_i^\top (\beta - \beta^*) - \frac{1}{nh} \sum_{i=1}^n (-y_i) H' \left(\frac{x_i + z_i^\top \beta}{h} \right) z_i. \quad (20)$$

Proposition 3.2. Assume Assumptions 1–5 hold. Further assume that $\|\hat{\beta}^{(0)} - \beta^*\|_2 = O_{\mathbb{P}}(\delta_{m,0})$, $\delta_{m,0} = O(h_1)$, $h_1 = o(1)$ and $\frac{\log n}{nh_1^3} = o(1)$. We have

$$\|U_{n,h_1}(\hat{\beta}^{(0)})\|_2 = O_{\mathbb{P}} \left(\delta_{m,0}^2 + h_1^\alpha + \sqrt{\frac{1}{nh_1}} + \delta_{m,0} \sqrt{\frac{\log n}{nh_1^3}} \right), \quad (21)$$

$$\|V_{n,h_1}(\hat{\beta}^{(0)}) - V\|_2 = O_{\mathbb{P}} \left(\sqrt{\frac{\log n}{nh_1^3}} + \delta_{m,0} + h_1 \right), \quad (22)$$

and therefore

$$\|\hat{\beta}^{(1)} - \beta^*\|_2 = O_{\mathbb{P}} \left(\delta_{m,0}^2 + h_1^\alpha + \sqrt{\frac{1}{nh_1}} + \delta_{m,0} \sqrt{\frac{\log n}{nh_1^3}} \right). \quad (23)$$

Proof of Proposition 3.2. Throughout the whole proof, without loss of generality, we assume that $\|\hat{\beta}^{(0)} - \beta^*\|_2 \leq \delta_{m,0}$ with probability approaching one, i.e., we assume the constant in $O_{\mathbb{P}}(\delta_{m,0})$ to be 1. For simplicity, we replace the notation h_1 with h . Also, note that the dimension p is fixed.

Proof of (21)

First prove (21). It suffices to show that

$$\sup_{\beta: \|\beta - \beta^*\|_2 \leq \delta_{m,0}} \|U_{n,h}(\beta)\|_2 = O_{\mathbb{P}} \left(\delta_{m,0}^2 + h^\alpha + \sqrt{\frac{1}{nh}} + \delta_{m,0} \sqrt{\frac{\log n}{nh^3}} \right), \quad (24)$$

which implies (21) since $\|\hat{\beta}^{(0)} - \beta^*\|_2 \leq \delta_{m,0}$ with probability approaching one.

By the proof of Lemma 3 in Cai et al. (2010), there exists $\mathbf{v}_1, \dots, \mathbf{v}_{5^p} \in \mathbb{R}^p$, s.t. for any \mathbf{v} in the unit sphere $\mathbb{S}^{p-1} = \{\mathbf{v} \in \mathbb{R}^p : \|\mathbf{v}\|_2 = 1\}$, there exists $j_v \in [5^p]$ satisfying $\|\mathbf{v} - \mathbf{v}_{j_v}\|_2 \leq 1/2$. Then we have

$$\|U_{n,h}(\beta)\|_2 = \sup_{\mathbf{v} \in \mathbb{S}^{p-1}} \left| \mathbf{v}^\top U_{n,h}(\beta) \right| \leq \sup_{j_v \in [5^p]} \left| \mathbf{v}_{j_v}^\top U_{n,h}(\beta) \right| + \frac{1}{2} \|U_{n,h}(\beta)\|_2,$$

and thus

$$\|U_{n,h}(\beta)\|_2 \leq \sup_{j_v \in [5^p]} 2 \left| \mathbf{v}_{j_v}^\top U_{n,h}(\beta) \right|.$$

Therefore, to show (24), it suffices to show that

$$\sup_{\beta: \|\beta - \beta^*\|_2 \leq \delta_{m,0}} \sup_{j_v \in [5^p]} \left| \mathbf{v}_{j_v}^\top U_{n,h}(\beta) \right| = O_{\mathbb{P}} \left(\sqrt{\frac{1}{nh}} + \delta_{m,0} \sqrt{\frac{\log n}{nh^3}} + \delta_{m,0}^2 + h^\alpha \right). \quad (25)$$

From now on, we let \mathbf{v} be an arbitrarily fixed vector in \mathbb{S}^{p-1} . Recall the definition

$$\begin{aligned} & U_{n,h}(\beta) \\ &= \frac{1}{nh^2} \sum_{i=1}^n (-y_i) H'' \left(\frac{x_i + \mathbf{z}_i^\top \beta}{h} \right) \mathbf{z}_i \mathbf{z}_i^\top (\beta - \beta^*) - \frac{1}{nh} \sum_{i=1}^n (-y_i) H' \left(\frac{x_i + \mathbf{z}_i^\top \beta}{h} \right) \mathbf{z}_i \\ &=: \frac{1}{n} \sum_{i=1}^n U_{h,i}(\beta), \end{aligned}$$

where

$$U_{h,i}(\beta) := \frac{1}{h^2} (-y_i) H'' \left(\frac{x_i + \mathbf{z}_i^\top \beta}{h} \right) \mathbf{z}_i \mathbf{z}_i^\top (\beta - \beta^*) - \frac{1}{h} (-y_i) H' \left(\frac{x_i + \mathbf{z}_i^\top \beta}{h} \right) \mathbf{z}_i. \quad (26)$$

We have the following decomposition:

$$\begin{aligned} & \mathbf{v}^\top U_{n,h}(\beta) \\ &= (1 - \mathbb{E}) \left[\mathbf{v}^\top U_{n,h}(\beta) - \mathbf{v}^\top U_{n,h}(\beta^*) \right] + (1 - \mathbb{E}) \mathbf{v}^\top U_{n,h}(\beta^*) + \mathbb{E} \mathbf{v}^\top U_{n,h}(\beta) \\ &= \frac{1}{n} \sum_{i=1}^n \phi_i^U(\beta) + (1 - \mathbb{E}) \mathbf{v}^\top U_{n,h}(\beta^*) + \mathbb{E} \mathbf{v}^\top U_{n,h}(\beta) \end{aligned} \quad (27)$$

where $\phi_i^U(\boldsymbol{\beta}) := (1 - \mathbb{E}) [\mathbf{v}^\top U_{h,i}(\boldsymbol{\beta}) - \mathbf{v}^\top U_{h,i}(\boldsymbol{\beta}^*)]$. We will separately bound the three terms in (27). through the following three steps.

Step 1

We will show that, for some constant $\gamma > 0$, there exists $C_\phi > 0$ such that

$$\sup_{\|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2 \leq \delta_{m,0}} \left| \frac{1}{n} \sum_{i=1}^n \phi_i^U(\boldsymbol{\beta}) \right| \leq C_\phi \delta_{m,0} \sqrt{\frac{\log n}{nh^3}}, \quad (28)$$

with probability $1 - 2n^{-\gamma p}$.

For any positive γ and each $j \in \{1, \dots, p\}$, divide the interval $[\beta_j^* - \delta_{m,0}, \beta_j^* + \delta_{m,0}]$ into n^γ small intervals, each with length $\frac{2\delta_{m,0}}{n^\gamma}$. This division creates n^γ intervals on each dimension, and the direct product of those intervals divides the hypercube $\{\boldsymbol{\beta} : \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_\infty \leq \delta_{m,0}\}$ into $n^{\gamma p}$ small hypercubes. By arbitrarily picking a point on each small hypercube, we could find $\{\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_{n^{\gamma p}}\} \subset \mathbb{R}^p$, such that for all $\boldsymbol{\beta}$ in the ball $\{\boldsymbol{\beta} : \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2 \leq \delta_{m,0}\}$ (which is a subset of $\{\boldsymbol{\beta} : \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_\infty \leq \delta_{m,0}\}$), there exists $j_\beta \in [n^{\gamma p}]$ such that $\|\boldsymbol{\beta} - \boldsymbol{\beta}_{j_\beta}\|_\infty \leq \frac{2\delta_{m,0}}{n^\gamma}$.

By Assumption 5 which requires $\sup_i \|\mathbf{z}_i\|_\infty \leq \bar{B}$, we have $|\mathbf{v}^\top \mathbf{z}_i| \leq \bar{B} \|\mathbf{v}\|_1 \leq \bar{B} \sqrt{p} \|\mathbf{v}\|_2$ for any $\mathbf{v} \in \mathbb{R}^p$. By Assumption 1, $H''(x)$, $H'(x)$ are both bounded and Lipschitz continuous, and thus we have

$$\begin{aligned} & \left| \mathbf{v}^\top U_{h,i}(\boldsymbol{\beta}) - \mathbf{v}^\top U_{h,i}(\boldsymbol{\beta}_{j_\beta}) \right| \\ & \leq \frac{|\mathbf{v}^\top \mathbf{z}_i|}{h^2} \left| \mathbf{z}_i^\top (\boldsymbol{\beta} - \boldsymbol{\beta}^*) \left[H'' \left(\frac{x_i + \mathbf{z}_i^\top \boldsymbol{\beta}}{h} \right) - H'' \left(\frac{x_i + \mathbf{z}_i^\top \boldsymbol{\beta}_{j_\beta}}{h} \right) \right] + \mathbf{z}_i^\top (\boldsymbol{\beta} - \boldsymbol{\beta}_{j_\beta}) H'' \left(\frac{x_i + \mathbf{z}_i^\top \boldsymbol{\beta}_{j_\beta}}{h} \right) \right| \\ & \quad + \frac{|\mathbf{v}^\top \mathbf{z}_i|}{h} \left| H' \left(\frac{x_i + \mathbf{z}_i^\top \boldsymbol{\beta}_{j_\beta}}{h} \right) - H' \left(\frac{x_i + \mathbf{z}_i^\top \boldsymbol{\beta}}{h} \right) \right| \\ & \leq C \left(\frac{\bar{B}^3 p^{3/2} \delta_{m,0}^2}{n^\gamma h^3} + \frac{\bar{B}^2 p \delta_{m,0}}{n^\gamma h^2} \right), \end{aligned}$$

for some constant C depending on H but not on $\boldsymbol{\beta}$ or $\boldsymbol{\beta}_{j_\beta}$. Therefore,

$$|\phi_i^U(\boldsymbol{\beta}) - \phi_i^U(\boldsymbol{\beta}_{j_\beta})| \leq \left| (1 - \mathbb{E}) [\mathbf{v}^\top U_{h,i}(\boldsymbol{\beta}) - \mathbf{v}^\top U_{h,i}(\boldsymbol{\beta}_{j_\beta})] \right| \lesssim \frac{\bar{B}^3 p^{3/2} \delta_{m,0}^2}{n^\gamma h^3} + \frac{\bar{B}^2 p \delta_{m,0}}{n^\gamma h^2},$$

Choose γ to be sufficiently large such that

$$\frac{\bar{B}^3 p^{3/2} \delta_{m,0}^2}{n^\gamma h^3} + \frac{\bar{B}^2 p \delta_{m,0}}{n^\gamma h^2} \ll \delta_{m,0} \sqrt{\frac{\log n}{nh^3}},$$

and then we have

$$\begin{aligned} & \sup_{\|\beta - \beta^*\|_2 \leq \delta_{m,0}} \left| \frac{1}{n} \sum_{i=1}^n \phi_i^U(\beta) \right| - \sup_{j_\beta \in [n^{\gamma p}]} \left| \frac{1}{n} \sum_{i=1}^n \phi_i^U(\beta_{j_\beta}) \right| \\ &= o \left(\delta_{m,0} \sqrt{\frac{\log n}{nh^3}} \right). \end{aligned} \quad (29)$$

Now let β be a fixed vector in $\{\beta : \|\beta - \beta^*\|_2 \leq \delta_{m,0}\}$. Again using $|\mathbf{v}^\top \mathbf{z}_i| \leq \bar{B}\sqrt{p}$, $|\mathbf{z}_i^\top (\beta - \beta^*)| \leq \bar{B}\sqrt{p}\delta_{m,0}$ and the boundedness of $H'(x)$ and $H''(x)$, we have

$$\sup_i |\phi_i^U(\beta)| = O \left(\frac{\bar{B}^2 p \delta_{m,0}}{h^2} \right), \quad (30)$$

Recall that $\zeta = X + \mathbf{Z}^\top \beta^*$ and $\rho(\cdot | \mathbf{Z})$ denotes the density of ζ given \mathbf{Z} . Let $\mathbb{E}_{|\mathbf{Z}}$ denote the expectation conditional on \mathbf{Z} . By Assumption 2, $\rho(\cdot | \mathbf{Z})$ is bounded in a neighborhood of zero, and hence we have, when n is large enough,

$$\begin{aligned} & \mathbb{E}_{|\mathbf{Z}} \left[H'' \left(\frac{X + \mathbf{Z}^\top \beta}{h} \right) \right]^2 = \mathbb{E}_{|\mathbf{Z}} \left[H'' \left(\frac{\mathbf{Z}^\top (\beta - \beta^*) + \zeta}{h} \right) \right]^2 \\ &= h \int_{-1}^1 H''(\xi)^2 \rho(\xi h - \mathbf{Z}^\top (\beta - \beta^*) | \mathbf{Z}) d\xi \\ &= O(h), \end{aligned} \quad (31)$$

where we use that $h = o(1)$ and $\mathbf{Z}^\top (\beta - \beta^*) \leq \bar{B}\sqrt{p}\delta_{m,0} = o(1)$. Similarly,

$$\begin{aligned} & \mathbb{E}_{|\mathbf{Z}} \left[H' \left(\frac{X + \mathbf{Z}^\top \beta}{h} \right) - H' \left(\frac{X + \mathbf{Z}^\top \beta^*}{h} \right) \right]^2 = \mathbb{E}_{|\mathbf{Z}} \left[\frac{\mathbf{Z}^\top (\beta - \beta^*)}{h} H'' \left(\frac{X + \mathbf{Z}^\top \beta}{h} \right) \right]^2 \\ &= O \left(\frac{\bar{B}^2 p \delta_{m,0}^2}{h} \right). \end{aligned} \quad (32)$$

By (31) and (32),

$$\begin{aligned} \text{var}(\phi_i^U(\beta)) &\leq \mathbb{E} \left(\mathbf{v}^\top U_{h,i}(\beta) - \mathbf{v}^\top U_{h,i}(\beta^*) \right)^2 \\ &= O \left(\frac{\bar{B}^4 p^2 \delta_{m,0}^2}{h^3} \right). \end{aligned} \quad (33)$$

Neither of the constant hidden in the big O notation of (30) or (33) depends on β . By (30) and (33), we can apply Bernstein inequality to show that there exists a large enough $C_1 > 0$, which does not depend on β , such that

$$\mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n \phi_i^U(\beta) \right| > C_1 \delta_{m,0} \sqrt{\frac{\log n}{nh^3}} \right) \leq 2n^{-2\gamma p}.$$

The above inequality is true for any $\beta \in \mathbb{R}^p$ that satisfies $\|\beta - \beta^*\|_2 \leq \delta_{m,0}$. In particular, for any $j_\beta \in [n^{\gamma p}]$, with probability at least $1 - 2n^{-2\gamma p}$,

$$\left| \frac{1}{n} \sum_{i=1}^n \phi_i^U(\beta_{j_\beta}) \right| \leq C_1 \delta_{m,0} \sqrt{\frac{\log n}{nh^3}},$$

which implies that with probability at least $1 - 2n^{-\gamma p}$,

$$\sup_{j_\beta \in [n^{\gamma p}]} \left| \frac{1}{n} \sum_{i=1}^n \phi_i^U(\beta_{j_\beta}) \right| \leq C_1 \delta_{m,0} \sqrt{\frac{\log n}{nh^3}}.$$

Combining with (29), we obtain (28).

Step 2 We will show that there exists a constant $C^* > 0$, such that

$$(1 - \mathbb{E}) \mathbf{v}^\top U_{n,h}(\beta^*) \geq C^* \sqrt{\frac{1}{nh}}, \quad (34)$$

with probability at least $1 - 2n^{-\gamma p}$, where γ has been specified in Step 1.

Note that $(1 - \mathbb{E}) \mathbf{v}^\top U_{n,h}(\beta^*) = \frac{1}{n} \sum_{i=1}^n (1 - \mathbb{E}) [\mathbf{v}^\top U_{h,i}(\beta^*)]$ and $(1 - \mathbb{E}) [\mathbf{v}^\top U_{h,i}(\beta^*)]$ are i.i.d. among different i . By Assumption 2, the density function of $\zeta = X + \mathbf{Z}^\top \beta^*$ satisfies that $\rho^{(1)}(\cdot | \mathbf{Z})$ is bounded uniformly for all \mathbf{Z} in a neighborhood of 0, which implies that $\rho(t | \mathbf{Z}) = \rho(0 | \mathbf{Z}) + O(t)$.

The constant in $O(t)$ is the same for all t in the neighborhood. Therefore,

$$\begin{aligned} & \mathbb{E}_{\cdot | \mathbf{Z}} \left[\mathbf{v}^\top U_{h,i}(\beta^*) \right]^2 \\ &= \frac{(\mathbf{v}^\top \mathbf{Z})^2}{h^2} \mathbb{E}_{\cdot | \mathbf{Z}} \left[H' \left(\frac{X + \mathbf{Z}^\top \beta^*}{h} \right) \right]^2 \\ &= \int_{-1}^1 \frac{(\mathbf{v}^\top \mathbf{Z})^2}{h} [H'(\xi)]^2 \rho(\xi h | \mathbf{Z}) d\xi \\ &= \frac{(\mathbf{v}^\top \mathbf{Z})^2}{h} \int_{-1}^1 [H'(\xi)]^2 \rho(0 | \mathbf{Z}) d\xi + O\left((\mathbf{v}^\top \mathbf{Z})^2\right). \end{aligned} \quad (35)$$

Recall that $V_s := \pi_V \mathbb{E} \rho(0 | \mathbf{Z}) \mathbf{Z} \mathbf{Z}^\top$ and $\pi_V := \int_{-1}^1 [H'(\xi)]^2 d\xi$, which leads to

$$\mathbb{E} \left[\mathbf{v}^\top U_{h,i}(\beta^*) \right]^2 = \frac{1}{h} \left(\mathbf{v}^\top V_s \mathbf{v} + O(h) \right).$$

In Step 3, we also show that $\sup_{\beta: \|\beta - \beta^*\|_2 \leq \delta_{m,0}} \mathbb{E} [\mathbf{v}^\top U_{h,i}(\beta)] = o(1)$. Therefore,

$$\begin{aligned} & \text{var} \left[\sqrt{h} \mathbf{v}^\top U_{h,i}(\beta^*) \right] \\ &= h \left\{ \mathbb{E} \left[\mathbf{v}^\top U_{h,i}(\beta^*) \right]^2 - \left(\mathbb{E} \left[\mathbf{v}^\top U_{h,i}(\beta^*) \right] \right)^2 \right\} \\ &= \mathbf{v}^\top V_s \mathbf{v} + o(1), \end{aligned} \quad (36)$$

By CLT and Slutsky's Theorem,

$$\sqrt{nh} (1 - \mathbb{E}) \left[\mathbf{v}^\top U_{n,h}(\boldsymbol{\beta}^*) \right] \xrightarrow{d} \mathcal{N} \left(0, \mathbf{v}^\top V_s \mathbf{v} \right),$$

which implies (34).

Step 3

We will show that,

$$\sup_{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2 \leq \delta_{m,0}} \mathbb{E} \left[\mathbf{v}^\top U_{n,h}(\boldsymbol{\beta}) \right] \leq C_E (\delta_{m,0}^2 + h^\alpha), \quad (37)$$

where C_E does not depend on \mathbf{v} .

By Assumption 2 and 3, for almost every \mathbf{Z} ,

$$\begin{aligned} \rho(t | \mathbf{Z}) &= \sum_{k=0}^{\alpha-1} \frac{1}{k!} \rho^{(k)}(0 | \mathbf{Z}) t^k + \frac{1}{\alpha!} \rho^{(\alpha)}(t' | \mathbf{Z}) t^\alpha, \\ F(-t | \mathbf{Z}) &= \frac{1}{2} + \sum_{k=1}^{\alpha} \frac{1}{k!} F^{(k)}(0 | \mathbf{Z}) (-t)^k + \frac{1}{(\alpha+1)!} F^{(\alpha+1)}(t'' | \mathbf{Z}) (-t)^{\alpha+1}, \end{aligned}$$

where t', t'' are between 0 and t . Therefore,

$$(2F(-t | \mathbf{Z}) - 1) \rho(t | \mathbf{Z}) = \sum_{k=1}^{2\alpha+1} M_k(\mathbf{Z}) t^k, \quad (38)$$

where $M_k(\mathbf{Z})$'s are constants depending on ρ , F , \mathbf{Z} , and t . Since $\rho^{(k)}(\cdot | \mathbf{Z})$ and $F^{(k)}(\cdot | \mathbf{Z})$ are bounded around 0 for all k and \mathbf{Z} , we know there exists a constant M such that $\sup_{k, \mathbf{Z}, t} |M_k(\mathbf{Z})| \leq M$ for all t in a certain neighborhood of 0.

(In particular,

$$M_\alpha(\mathbf{Z}) = \sum_{k=1}^{\alpha} \frac{2(-1)^{-k}}{(\alpha-k)!k!} F^{(k)}(0 | \mathbf{Z}) \rho^{(\alpha-k)}(0 | \mathbf{Z}),$$

which will be used in the proof of Theorem 3.4.)

By Assumption 1, when $x > 1$ or $x < -1$, $H'(x) = H''(x) = 0$. The kernel $H'(x) = \int_{-1}^x H''(t) dt$ is bounded, satisfying $\int_{-1}^1 H'(x) dx = 1$ and $\int_{-1}^1 x^k H'(x) dx = 0$ for any $1 \leq k \leq \alpha - 1$. Using integration by parts, we have $\int_{-1}^1 x H''(x) dx = -1$ and $\int_{-1}^1 x^k H''(x) dx = 0$ for $k = 0$ and $2 \leq k \leq \alpha$.

Now we are ready to compute the expectation of $\mathbf{v}^\top U_{n,h}(\boldsymbol{\beta})$. Since $\mathbb{E}[\mathbf{v}^\top U_{n,h}(\boldsymbol{\beta})] = \mathbb{E}[\mathbf{v}^\top U_{h,i}(\boldsymbol{\beta})]$ for all i , we omit i in the following computation of expectation. Let $\mathbb{E}_{|\mathbf{Z}}$ denote the expectation

conditional on \mathbf{Z} . Define $\Delta(\beta) := \beta - \beta^*$ and recall that $\zeta = X + \mathbf{Z}^\top \beta^*$.

$$\begin{aligned}
& \mathbb{E}_{\cdot|\mathbf{Z}} \left[\mathbf{v}^\top U_{n,h}(\beta) \right] \\
&= \mathbf{Z}^\top \mathbf{v} \cdot \mathbb{E}_{\cdot|\mathbf{Z}} \left[\frac{\mathbf{Z}^\top \Delta(\beta)}{h^2} \left[2\mathbb{I}(X + \mathbf{Z}^\top \beta^* + \epsilon < 0) - 1 \right] H'' \left(\frac{X + \mathbf{Z}^\top \beta}{h} \right) \right. \\
&\quad \left. - \frac{1}{h} \left[2\mathbb{I}(X + \mathbf{Z}^\top \beta^* + \epsilon < 0) - 1 \right] H' \left(\frac{X + \mathbf{Z}^\top \beta}{h} \right) \right] \\
&= \mathbf{Z}^\top \mathbf{v} \cdot \mathbb{E}_{\cdot|\mathbf{Z}} \left\{ \left[2\mathbb{I}(\zeta + \epsilon < 0) - 1 \right] \left[\frac{\mathbf{Z}^\top \Delta(\beta)}{h^2} H'' \left(\frac{\zeta + \mathbf{Z}^\top \Delta(\beta)}{h} \right) - \frac{1}{h} H' \left(\frac{\zeta + \mathbf{Z}^\top \Delta(\beta)}{h} \right) \right] \right\} \\
&= (\mathbf{Z}^\top \mathbf{v}) \int_{-1}^1 \left[2F(\mathbf{Z}^\top \Delta(\beta) - \xi h | \mathbf{Z}) - 1 \right] \rho(\xi h - \mathbf{Z}^\top \Delta(\beta) | \mathbf{Z}) \\
&\quad \cdot \left(\frac{\mathbf{Z}^\top \Delta(\beta)}{h} H''(\xi) - H'(\xi) \right) d\xi \quad \left(\text{by changing variable } \xi = \frac{\zeta + \mathbf{Z}^\top \Delta(\beta)}{h} \right) \\
&= (\mathbf{Z}^\top \mathbf{v}) \sum_{k=1}^{2\alpha+1} M_k(\mathbf{Z}) \int_{-1}^1 (\xi h - \mathbf{Z}^\top \Delta(\beta))^k \left(\frac{\mathbf{Z}^\top \Delta(\beta)}{h} H''(\xi) - H'(\xi) \right) d\xi.
\end{aligned} \tag{39}$$

When $1 \leq k \leq \alpha - 1$,

$$\begin{aligned}
& \int_{-1}^1 (\xi h - \mathbf{Z}^\top \Delta(\beta))^k \left(\frac{\mathbf{Z}^\top \Delta(\beta)}{h} H''(\xi) - H'(\xi) \right) d\xi \\
&= \sum_{k'=0}^k \binom{k}{k'} h^{k'} (-\mathbf{Z}^\top \Delta(\beta))^{k-k'} \left[\left(\mathbf{Z}^\top \Delta(\beta) / h \right) \int_{-1}^1 \xi^{k'} H''(\xi) d\xi - \int_{-1}^1 \xi^{k'} H'(\xi) d\xi \right] \\
&= (k-1) (-\mathbf{Z}^\top \Delta(\beta))^k,
\end{aligned}$$

When $k = \alpha$,

$$\begin{aligned}
& \int_{-1}^1 (\xi h - \mathbf{Z}^\top \Delta(\beta))^\alpha \left(\frac{\mathbf{Z}^\top \Delta(\beta)}{h} H''(\xi) - H'(\xi) \right) d\xi \\
&= \sum_{k'=0}^{\alpha} \binom{\alpha}{k'} h^{k'} (-\mathbf{Z}^\top \Delta(\beta))^{\alpha-k'} \left[\left(\mathbf{Z}^\top \Delta(\beta) / h \right) \int_{-1}^1 \xi^{k'} H''(\xi) d\xi - \int_{-1}^1 \xi^{k'} H'(\xi) d\xi \right] \\
&= (\alpha-1) (-\mathbf{Z}^\top \Delta(\beta))^\alpha - \pi_U h^\alpha,
\end{aligned} \tag{40}$$

where $\pi_U = \int_{-1}^1 \xi^\alpha H'(\xi) d\xi$ is defined in Assumption 1.

When $\alpha+1 \leq k \leq 2\alpha+1$, using $|\mathbf{Z}^\top \Delta(\beta)| \leq \bar{B} \sqrt{p} \delta_{m,0}$, $\delta_{m,0} = O(h)$ and that H', H'' are both bounded, we have

$$\sup_{\beta: \|\beta - \beta^*\|_2 \leq \delta_{m,0}} \left| \int_{-1}^1 (\xi h - \mathbf{Z}^\top \Delta(\beta))^k \left(\frac{\mathbf{Z}^\top \Delta(\beta)}{h} H''(\xi) - H'(\xi) \right) d\xi \right| = O \left[h^{\alpha+1} + \delta_{m,0}^{\alpha+1} \right].$$

Since $M_k(\mathbf{Z})$'s are uniformly bounded, we finally obtain

$$\sup_{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2 \leq \delta_{m,0}} \left| \mathbb{E} \left[\mathbf{v}^\top U_{n,h}(\boldsymbol{\beta}) \right] \right| = O(\delta_{m,0}^2 + h^\alpha), \quad (41)$$

which leads to (37).

Combining the three steps leads to, for any fixed vector $\mathbf{v} \in \mathbb{S}^{p-1}$, there exists a constant $C_{\mathbf{v}}$, with probability at least $1 - 4(n^{-\gamma})^p$,

$$\sup_{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2 \leq \delta_{m,0}} \left| \mathbf{v}^\top U_{n,h}(\boldsymbol{\beta}) \right| \leq C_{\mathbf{v}} \left(\sqrt{\frac{1}{nh}} + \delta_{m,0} \sqrt{\frac{\log n}{nh^3}} + \delta_{m,0}^2 + h^\alpha \right),$$

The constant $C_{\mathbf{v}}$ depends on \mathbf{v} . Recall in equation (25), we only consider finite number of \mathbf{v} 's.

Take $C_{\text{sup}} = \sup_{j \in [5p]} C_{\mathbf{v}_j}$, and then we have

$$\sup_{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2 \leq \delta_{m,0}} \sup_{j_v \in [5p]} \left| \mathbf{v}_{j_v}^\top U_{n,h}(\boldsymbol{\beta}) \right| \leq C_{\text{sup}} \left(\sqrt{\frac{1}{nh}} + \delta_{m,0} \sqrt{\frac{\log n}{nh^3}} + \delta_{m,0}^2 + h^\alpha \right),$$

with probability at least $1 - 4(5n^{-\gamma})^p$. This proves (24) and completes the proof of (21).

Proof of (22)

The proof of (22) is similar. We will use the same \mathbf{v}_{j_v} and $\boldsymbol{\beta}_{j_\beta}$ as before. By the proof of Lemma 3 in Cai et al. (2010),

$$\|A\|_2 \leq 4 \sup_{j_v \in [5p]} \left| \mathbf{v}_{j_v}^\top A \mathbf{v}_{j_v} \right|,$$

for any symmetric $A \in \mathbb{R}^{p \times p}$. Therefore, we only need to bound $\sup_{j_v \in [5p]} \left| \mathbf{v}_{j_v}^\top [V_{n,h}(\boldsymbol{\beta}^{(0)}) - V] \mathbf{v}_{j_v} \right|$.

By the choice of $\{\boldsymbol{\beta}_{j_\beta}\}$, for all $\boldsymbol{\beta}$ in the ball $\{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2 \leq \delta_{m,0}\}$, there exists $j_\beta \in [n^{\gamma p}]$ such that $\|\boldsymbol{\beta} - \boldsymbol{\beta}_{j_\beta}\|_\infty \leq \frac{2\delta_{m,0}}{n^\gamma}$. Recall

$$V_{n,h}(\boldsymbol{\beta}) = \frac{1}{nh^2} \sum_{i=1}^n (-y_i) H'' \left(\frac{x_i + \mathbf{z}_i^\top \boldsymbol{\beta}}{h} \right) \mathbf{z}_i \mathbf{z}_i^\top =: \frac{1}{n} \sum_{i=1}^n V_{h,i}(\boldsymbol{\beta}),$$

where

$$V_{h,i}(\boldsymbol{\beta}) := \frac{1}{h^2} (-y_i) H'' \left(\frac{x_i + \mathbf{z}_i^\top \boldsymbol{\beta}}{h} \right) \mathbf{z}_i \mathbf{z}_i^\top. \quad (42)$$

We have

$$\left| \mathbf{v}_{j_v}^\top V_{n,h}(\boldsymbol{\beta}) \mathbf{v}_{j_v} - \mathbf{v}_{j_v}^\top V_{n,h}(\boldsymbol{\beta}_{j_\beta}) \mathbf{v}_{j_v} \right| \leq \sup_i \frac{(\mathbf{z}_i^\top \mathbf{v}_{j_v})^2}{h^3} \mathbf{z}_i^\top (\boldsymbol{\beta} - \boldsymbol{\beta}_{j_\beta}) = O \left(\frac{\bar{B}^3 p^{3/2} \delta_{m,0}}{n^\gamma h^3} \right).$$

By taking $\gamma > 0$ large enough, it suffices to show that

$$\sup_{j_v \in [5p]} \sup_{j_\beta \in [n^{\gamma p}]} \left| \mathbf{v}_{j_v}^\top [V_{n,h}(\boldsymbol{\beta}_{j_\beta}) - V] \mathbf{v}_{j_v} \right| = O_{\mathbb{P}} \left(\sqrt{\frac{\log n}{nh^3}} + \delta_{m,0} + h \right).$$

In the following proof, we will again consider any fixed $\beta \in \mathbb{R}^p, \mathbf{v} \in \mathbb{R}^p$ that satisfy $\|\beta - \beta^*\|_2 \leq \delta_{m,0}$ and $\|\mathbf{v}\|_2 = 1$, and then apply the result to the specific \mathbf{v}_{j_v} and β_{j_β} . The computation of $\mathbb{E}[\mathbf{v}^\top [V_{n,h}(\beta) - V] \mathbf{v}]$ is similar to before, but this time we only need the following expansion:

$$\rho(t | \mathbf{Z}) = \rho(0 | \mathbf{Z}) + O(t),$$

$$F(-t | \mathbf{Z}) = 1/2 - F'(0 | \mathbf{Z})t + O(t^2).$$

Recall that $\int_{-1}^1 x H''(x) dx = -1$ and $\int_{-1}^1 x^k H''(x) dx = 0$ for $k = 0$ and $2 \leq k \leq \alpha$.

$$\begin{aligned} & \mathbb{E}_{|\mathbf{Z}} \left(\mathbf{v}^\top V_{n,h}(\beta) \mathbf{v} \right) \\ &= \frac{(\mathbf{v}^\top \mathbf{z})^2}{h^2} \mathbb{E}_{|\mathbf{Z}} [2\mathbb{I}(\zeta + \epsilon < 0) - 1] H'' \left(\frac{\mathbf{Z}^\top \Delta(\beta) + \zeta}{h} \right) \\ &= \frac{(\mathbf{v}^\top \mathbf{Z})^2}{h} \int_{-1}^1 \left(2F(\mathbf{Z}^\top \Delta(\beta) - \xi h | \mathbf{Z}) - 1 \right) \rho(\xi h - \mathbf{Z}^\top \Delta(\beta) | \mathbf{Z}) H''(\xi) d\xi \\ &= -2 (\mathbf{v}^\top \mathbf{Z})^2 F'(0 | \mathbf{Z}) \rho(0 | \mathbf{Z}) \int_{-1}^1 \xi H''(\xi) d\xi + O \left(h + \left| \mathbf{Z}^\top \Delta(\beta) \right|^2 / h \right) \\ &= 2 (\mathbf{v}^\top \mathbf{Z})^2 F'(0 | \mathbf{Z}) \rho(0 | \mathbf{Z}) + O(h + \delta_{m,0}), \end{aligned} \quad (43)$$

hence

$$\mathbb{E} \left(\mathbf{v}^\top [V_{n,h}(\beta) - V] \mathbf{v} \right) = O(\delta_{m,0} + h). \quad (44)$$

Since $H''(x)$ is bounded and $|\mathbf{v}^\top \mathbf{z}_i| \leq \bar{B}\sqrt{p}$ (p is fixed), we have

$$\left| \mathbf{v}^\top V_{h,i}(\beta) \mathbf{v} \right| = O \left(\frac{1}{h^2} \right). \quad (45)$$

Adding (31), we also have

$$\mathbb{E} \left(\mathbf{v}^\top V_{h,i}(\beta) \mathbf{v} \right)^2 \leq \mathbb{E} \left[\frac{(\mathbf{v}^\top \mathbf{z}_i)^4}{h^4} O(h) \right] = O \left(\frac{1}{h^3} \right). \quad (46)$$

By (45) and (46), we could apply Bernstein inequality to show that

$$\mathbb{P} \left(\left| (1 - \mathbb{E}) \frac{1}{n} \sum_{i=1}^n \mathbf{v}^\top V_{h,i}(\beta) \mathbf{v} \right| > C_2 \sqrt{\frac{p \log n}{nh^3}} \right) \leq 2n^{-2\gamma p},$$

for some constant $C_2 > 0$. By the same procedure in the proof of (21), we obtain, with probability at least $1 - 2(5n^{-\gamma})^p$,

$$\sup_{j_\beta \in [n^{\gamma p}]} \sup_{j_v \in [5p]} \left| \mathbf{v}_{j_v}^\top (V_{n,h}(\beta_{j_\beta}) - V) \mathbf{v}_{j_v} \right| = O \left(\sqrt{\frac{\log n}{nh^3}} + \delta_{m,0} + h \right),$$

which completes the proof of (22).

Finally, (18), (21), and (22) directly leads to (11), given the assumption that $\Lambda_{\min}(V) > c_1^{-1}$ for some $c_1 > 0$ (Assumption 4). □

Proof of Theorem 3.3

Proof. Without loss of generality, we assume $\lambda_h = 1$. Then we have

$$h_t = \max\{n^{-1/(2\alpha+1)}, m^{-2^t/(3\alpha)}\} \geq n^{-\frac{1}{2\alpha+1}},$$

which implies the assumption $\frac{\log n}{nh_t^3} = o(1)$ holds for any t , since $\frac{\log n}{nh_t^3} \leq \log n \cdot n^{-\frac{2\alpha-2}{2\alpha+1}} \rightarrow 0$ if $\alpha > 1$. Moreover, for any t , $h_t^\alpha = \max\{m^{-2^t/3}, n^{-\alpha/(2\alpha+1)}\}$ and $1/\sqrt{nh_t} \leq n^{-\alpha/(2\alpha+1)}$.

We first show that, for any t ,

$$\|\hat{\beta}^{(t)} - \beta^*\|_2 = O_{\mathbb{P}}\left(n^{-\alpha/(2\alpha+1)}\sqrt{\log n} + m^{-2^t/3}\right). \quad (47)$$

Recall that by (11), if $\|\hat{\beta}^{(t-1)} - \beta^*\|_2 = O_{\mathbb{P}}(\delta_{m,t-1})$ and $\delta_{m,t-1} = O(h_t)$, we have that

$$\|\hat{\beta}^{(t)} - \beta^*\|_2 = O_{\mathbb{P}}\left(\delta_{m,t-1}^2 + h_t^\alpha + \sqrt{\frac{1}{nh_t}} + \delta_{m,t-1}\sqrt{\frac{\log n}{nh_t^3}}\right), \quad (48)$$

which implies that

$$\|\hat{\beta}^{(t)} - \beta^*\|_2 = O_{\mathbb{P}}\left(\delta_{m,t-1}^2 + h_t^\alpha + \sqrt{\frac{\log n}{nh_t}}\right), \quad (49)$$

since $\delta_{m,t-1}\sqrt{\frac{\log n}{nh_t^3}} = O(\sqrt{\frac{\log n}{nh_t}})$.

When $t = 1$, $\delta_{m,0} := m^{-1/3} \leq h_1 = m^{-\frac{2}{3\alpha}}$, which verifies $\delta_{m,0} = O(h_1)$. Then (47) holds since

$$\|\hat{\beta}^{(1)} - \beta^*\|_2 = O_{\mathbb{P}}\left(\sqrt{\frac{\log n}{nh_1}} + m^{-2/3} + h_1^\alpha\right) = O_{\mathbb{P}}\left(n^{-\alpha/(2\alpha+1)}\sqrt{\log n} + m^{-2/3}\right).$$

Assume (47) holds for $t-1$, i.e.,

$$\|\hat{\beta}^{(t-1)} - \beta^*\|_2 = O_{\mathbb{P}}\left(n^{-\alpha/(2\alpha+1)}\sqrt{\log n} + m^{-2^{t-1}/3}\right).$$

Then $\delta_{m,t-1} = n^{-\alpha/(2\alpha+1)}\sqrt{\log n} + m^{-2^{t-1}/3}$. Since $n^{-\alpha/(2\alpha+1)}\sqrt{\log n} \ll n^{-\frac{1}{2\alpha+1}} \leq h_t$ and $m^{-2^{t-1}/3} \leq m^{-\frac{2^t}{3\alpha}} \leq h_t$, it holds that $\delta_{m,t-1} = O(h_t)$. Then by (49),

$$\|\hat{\beta}^{(t)} - \beta^*\|_2 = O_{\mathbb{P}}\left(\sqrt{\frac{\log n}{nh_t}} + \delta_{m,t-1}^2 + h_t^\alpha\right) = O_{\mathbb{P}}\left(n^{-\alpha/(2\alpha+1)}\sqrt{\log n} + m^{-2^t/3}\right),$$

since $h_t^\alpha = \max \left\{ m^{-2^t/3}, n^{-\alpha/(2\alpha+1)} \right\}$ and $1/\sqrt{nh_t} \leq n^{-\alpha/(2\alpha+1)}$. Therefore, we have proved that (47) holds for all t by induction.

To see (12), note that plugging $\delta_{m,t-1} = n^{-\alpha/(2\alpha+1)}\sqrt{\log n} + m^{-2^{t-1}/3}$ into (48) yields that

$$\left\| \widehat{\beta}^{(t)} - \beta^* \right\|_2 = O_{\mathbb{P}} \left(\frac{\log n}{n^{2\alpha/(2\alpha+1)}} + m^{-2^t/3} + h_t^\alpha + \sqrt{\frac{1}{nh_t}} + \left(\frac{\sqrt{\log n}}{n^{\alpha/(2\alpha+1)}} + m^{-2^{t-1}/3} \right) \sqrt{\frac{\log n}{nh_t^3}} \right). \quad (50)$$

Since $\alpha > 1$,

$$\log n \cdot \sqrt{\frac{1}{nh_t^3}} \leq \log n \cdot n^{-\frac{\alpha-1}{2\alpha+1}} = o(1).$$

Therefore, the rate in (50) is upper bounded by

$$O \left(n^{-\alpha/(2\alpha+1)} + m^{-2^t/3} + m^{-2^{t-1}/3} n^{-\frac{\alpha-1}{2\alpha+1}} \sqrt{\log n} \right),$$

which completes the proof. \square

C.2 Theoretical Results for the Inference of (mSMSE)

To show the asymptotic normality, we first prove an important Lemma about $U_{n,h}(\beta)$ defined by (20).

Lemma C.1. *Under Assumptions 1–5, if $h = o(1)$, $\|\beta - \beta^*\|_2 = O_{\mathbb{P}}(\delta)$, $\delta = o(\max\{h^{\alpha/2}, h/\sqrt{\log n}\})$ then*

$$\mathbb{E}[U_{n,h}(\beta)] = Uh^\alpha + o_{\mathbb{P}}(h^\alpha), \quad (51)$$

and

$$(1 - \mathbb{E})\sqrt{nh}U_{n,h}(\beta) \xrightarrow{d} \mathcal{N}(0, V_s). \quad (52)$$

Proof of Lemma C.1. Recall that

$$M_\alpha(\mathbf{Z}) = \sum_{k=1}^{\alpha} \frac{2(-1)^k}{(\alpha-k)!k!} F^{(k)}(0|\mathbf{Z}) \rho^{(\alpha-k)}(0|\mathbf{Z}),$$

and

$$U := \pi_U \mathbb{E} \left(\sum_{k=1}^{\alpha} \frac{2(-1)^{k+1}}{k!(\alpha-k)!} F^{(k)}(0|\mathbf{Z}) \rho^{(\alpha-k)}(0|\mathbf{Z}) \mathbf{Z} \right),$$

where π_U is defined in Assumption 1. For any $\mathbf{v} \in \mathbb{S}^{p-1}$, the computation in (39) yields that

$$\begin{aligned} \mathbb{E} \left[\mathbf{v}^\top U_{n,h}(\beta) \right] &= -\mathbb{E} \left[\left(\mathbf{Z}^\top \mathbf{v} \right) \pi_U M_\alpha(\mathbf{Z}) h^\alpha \right] + O \left(\|\beta - \beta^*\|_2^2 + h^{\alpha+1} \right) \\ &= \mathbf{v}^\top U \cdot h^\alpha + o_{\mathbb{P}}(h^\alpha), \end{aligned}$$

where the big O notation hides a constant not depending on β . This proves (51).

To show (52), further recall that in **Step 2** of the proof of (21), we show that

$$(1 - \mathbb{E}) \sqrt{nh} U_{n,h}(\beta^*) \xrightarrow{d} \mathcal{N}(0, V_s).$$

Furthermore, in **Step 1** of the same proof, we prove that

$$\sup_{\beta: \|\beta - \beta^*\| \leq \delta} |(1 - \mathbb{E}) [U_{n,h}(\beta) - U_{n,h}(\beta^*)]| = O_{\mathbb{P}} \left(\delta \sqrt{\frac{\log n}{nh^3}} \right),$$

which yields (52) if $\delta = o(h/\sqrt{\log n})$.

□

Proof of Theorem 3.4

We now restate Theorem 3.4 and give the proof using Lemma C.1.

Theorem 3.4. *Assume the local size $m > n^c$ for some constant $0 < c < 1$ and the assumptions in Theorem 3.3 hold. When T satisfies (13), we have*

$$n^{\frac{\alpha}{2\alpha+1}} (\hat{\beta}^{(T)} - \beta^*) \xrightarrow{d} \mathcal{N} \left(\lambda_h^{\frac{\alpha}{2\alpha+1}} V^{-1} U, \lambda_h^{\frac{-1}{2\alpha+1}} V^{-1} V_s V^{-1} \right), \quad (53)$$

where U and V_s are defined in (10). The asymptotic mean squared error is given by

$$\mathbb{E} \left[\left(\hat{\beta}^{(T)} - \beta^* \right)^\top \left(\hat{\beta}^{(T)} - \beta^* \right) \right] = n^{-\frac{2\alpha}{2\alpha+1}} \cdot \text{trace} \left[\lambda_h^{-\frac{1}{2\alpha+1}} V^{-1} V_s V^{-1} + \lambda_h^{-\frac{2\alpha}{2\alpha+1}} U^\top V^{-1} V^{-1} U \right],$$

by minimizing which we obtain the optimal λ_h^* as follows.

$$\lambda_h^* := \frac{\text{trace}(V^{-1} V_s V^{-1})}{2\alpha U^\top V^{-1} V^{-1} U}.$$

Proof. By Theorem 3.3, when

$$T > \log_2 \left(\frac{6\alpha}{2\alpha+1} \cdot \frac{\log n - \log \lambda_h}{\log m} \right),$$

$h_T = \left(\frac{\lambda_h}{n} \right)^{\frac{1}{2\alpha+1}}$ and $m^{-2^{T-1}/3} \lesssim n^{-\alpha/(2\alpha+1)}$. Hence, $\left\| \hat{\beta}^{(T-1)} - \beta^* \right\|_2 = O_{\mathbb{P}} \left(\sqrt{\log n} \cdot n^{-\frac{\alpha}{2\alpha+1}} \right) = o_{\mathbb{P}} \left(\max\{h_T^{\alpha/2}, h_T/\sqrt{\log n}\} \right)$, and thus the assumptions of Lemma C.1 hold. Since $\sqrt{nh_T} = \sqrt{\lambda_h} h_T^{-\alpha}$, we obtain by Lemma C.1 that

$$\sqrt{nh_T} \cdot U_{n,h_T} \left(\hat{\beta}^{(T-1)} \right) \xrightarrow{d} \mathcal{N} \left(\sqrt{\lambda_h} U, V_s \right),$$

and hence

$$\sqrt{nh_T} \cdot V^{-1} U_{n,h_T} \left(\widehat{\beta}^{(T-1)} \right) \xrightarrow{d} \mathcal{N} \left(\sqrt{\lambda_h} V^{-1} U, V^{-1} V_s V \right).$$

By (22) and Assumption 4, $\left\| V_{n,h_T}^{-1} \left(\widehat{\beta}^{(T-1)} \right) - V^{-1} \right\|_2 = o_{\mathbb{P}}(1)$, which yields

$$\sqrt{nh_T} \cdot V_{n,h_T}^{-1} \left(\widehat{\beta}^{(T-1)} \right) U_{n,h_T} \left(\widehat{\beta}^{(T-1)} \right) \xrightarrow{d} \mathcal{N} \left(\sqrt{\lambda_h} V^{-1} U, V^{-1} V_s V^{-1} \right).$$

Plugging in $h_T = \left(\frac{\lambda_h}{n} \right)^{\frac{1}{2\alpha+1}}$ and $\widehat{\beta}^{(T)} - \beta^* = V_{n,h_T}^{-1} \left(\widehat{\beta}^{(T-1)} \right) U_{n,h_T} \left(\widehat{\beta}^{(T-1)} \right)$ leads to (14). Using the asymptotic mean and variance obtained by Lemma C.1, we have

$$\begin{aligned} & \mathbb{E} \left[\left(\widehat{\beta}^{(T)} - \beta^* \right)^{\top} \left(\widehat{\beta}^{(T)} - \beta^* \right) \right] \\ &= \mathbb{E} \text{trace} \left[\left(\widehat{\beta}^{(T)} - \beta^* \right)^{\top} \left(\widehat{\beta}^{(T)} - \beta^* \right) \right] \\ &= \mathbb{E} \text{trace} \left[\left(\widehat{\beta}^{(T)} - \beta^* \right) \left(\widehat{\beta}^{(T)} - \beta^* \right)^{\top} \right] \\ &\rightarrow n^{-\frac{2\alpha}{2\alpha+1}} \cdot \text{trace} \left[\lambda_h^{-\frac{1}{2\alpha+1}} V^{-1} V_s V^{-1} + \lambda_h^{-\frac{2\alpha}{2\alpha+1}} U^{\top} V^{-1} V^{-1} U \right]. \end{aligned}$$

Then it's direct to obtain the optimal λ_h^* given in (15). \square

Estimators for V , U and V_s

Now we formally define the estimators for V , U and V_s . Proposition C.1 in Section C.1 has already implies that $V_{n,h_T} \left(\widehat{\beta}^{(T)} \right) \xrightarrow{p} V$, where $V_{n,h}(\beta)$ is defined in (19), so we can use $\widehat{V} := V_{n,h_T} \left(\widehat{\beta}^{(T)} \right)$ to estimate V . It remains to provide estimators for U and V_s .

Theorem C.2. Assume assumptions in Theorem 3.3 hold. Let $h_{\kappa} = n^{-\frac{\kappa}{2\alpha+1}}$ for some $0 < \kappa < 1$.

When

$$T \geq \log_2 \left(\frac{6\alpha}{2\alpha+1} \cdot \frac{\log n - \log \lambda_h}{\log m} \right),$$

we have

$$\widehat{U} := \frac{1}{nh_{\kappa}^{\alpha+1}} \sum_{i=1}^n y_i H' \left(\frac{x_i + \mathbf{z}_i^{\top} \widehat{\beta}^{(T)}}{h_{\kappa}} \right) \mathbf{z}_i \xrightarrow{p} U, \quad (54)$$

$$\widehat{V}_s := \frac{1}{nh_T} \sum_{i=1}^n \left[H' \left(\frac{x_i + \mathbf{z}_i^{\top} \widehat{\beta}^{(T)}}{h_T} \right) \right]^2 \mathbf{z}_i \mathbf{z}_i^{\top} \xrightarrow{p} V_s. \quad (55)$$

Proof. When $t > \log_2 \left(\frac{3\alpha}{2\alpha+1} \cdot \frac{\log n - \log \lambda_h}{\log m} \right)$, we have

$$n^{-\frac{\kappa}{2\alpha+1}} = h_{\kappa} \gg h_t = \left(\frac{\lambda_h}{n} \right)^{\frac{1}{2\alpha+1}},$$

and

$$\left\| \widehat{\boldsymbol{\beta}}^{(t)} - \boldsymbol{\beta}^* \right\|_2 = O_{\mathbb{P}} \left(\sqrt{\log n} \cdot n^{-\frac{\alpha}{2\alpha+1}} \right) = o_{\mathbb{P}} (h_{\kappa}^{\alpha}).$$

Replacing h_t with h^* in the proof of Theorem 3.4 leads to

$$\mathbb{E} \mathbf{v}^{\top} U_{n, h_{\kappa}} \left(\widehat{\boldsymbol{\beta}}^{(t)} \right) = \mathbf{v}^{\top} U h_{\kappa}^{\alpha} + o_{\mathbb{P}} (h_{\kappa}^{\alpha}),$$

and

$$\sqrt{nh_{\kappa}} (1 - \mathbb{E}) \mathbf{v}^{\top} U_{n, h_{\kappa}} \left(\widehat{\boldsymbol{\beta}}^{(t)} \right) \xrightarrow{d} \mathcal{N} \left(0, \mathbf{v}^{\top} V_s \mathbf{v} \right),$$

for any $\mathbf{v} \in \mathbb{S}^p$. Also, $V_{n, h_{\kappa}} \left(\widehat{\boldsymbol{\beta}}^{(t)} \right) \xrightarrow{p} V$ by Proposition 3.2. Note that

$$\widehat{U} = (h_{\kappa})^{-\alpha} \left[U_{n, h_{\kappa}} \left(\widehat{\boldsymbol{\beta}}^{(t)} \right) - V_{n, h_{\kappa}} \left(\widehat{\boldsymbol{\beta}}^{(t)} \right) \left(\widehat{\boldsymbol{\beta}}^{(t)} - \boldsymbol{\beta}^* \right) \right],$$

and $(h_{\kappa})^{-\alpha} \ll \sqrt{nh_{\kappa}}$. Then

$$\begin{aligned} \mathbf{v}^{\top} \widehat{U} &= (h_{\kappa})^{-\alpha} \left[\mathbb{E} \mathbf{v}^{\top} U_{n, h_{\kappa}} \left(\widehat{\boldsymbol{\beta}}^{(t)} \right) + (1 - \mathbb{E}) \mathbf{v}^{\top} U_{n, h_{\kappa}} \left(\widehat{\boldsymbol{\beta}}^{(t)} \right) - \mathbf{v}^{\top} V_{n, h_{\kappa}} \left(\widehat{\boldsymbol{\beta}}^{(t)} \right) \left(\widehat{\boldsymbol{\beta}}^{(t)} - \boldsymbol{\beta}^* \right) \right] \\ &= \mathbf{v}^{\top} U + o_{\mathbb{P}} (1), \end{aligned}$$

which yields (54).

To prove (55), recall

$$\widehat{V}_s := \frac{1}{nh_t} \sum_{i=1}^n \left[H' \left(\frac{x_i + \mathbf{z}_i^{\top} \widehat{\boldsymbol{\beta}}^{(t)}}{h_t} \right) \right]^2 \mathbf{z}_i \mathbf{z}_i^{\top}.$$

For any $\mathbf{v} \in \mathbb{S}^p$,

$$\begin{aligned} \mathbb{E} \mathbf{v}^{\top} \widehat{V}_s \mathbf{v} &= \mathbb{E} \left(\mathbb{E}_{\cdot | \mathbf{Z}} \mathbf{v}^{\top} \widehat{V}_s \mathbf{v} \right) \\ &= \mathbb{E} \left[\left(\mathbf{Z}^{\top} \mathbf{v} \right)^2 \int_{-1}^1 [H'(\xi)]^2 \rho \left(\xi h_t - \mathbf{Z}^{\top} \left(\widehat{\boldsymbol{\beta}}^{(t)} - \boldsymbol{\beta}^* \right) \mid \mathbf{Z} \right) d\xi \right] \\ &= \mathbb{E} \left(\mathbf{Z}^{\top} \mathbf{v} \right)^2 \pi_V \rho(0 \mid \mathbf{Z}) + O \left(h_t + \left\| \widehat{\boldsymbol{\beta}}^{(t)} - \boldsymbol{\beta}^* \right\|_2 \right) \\ &= \mathbf{v}^{\top} V_s \mathbf{v} + o(1). \end{aligned}$$

Also, since $H'(x)$ and $\|\mathbf{z}_i\|_{\infty}$ are both bounded, $\text{var} \left(\mathbf{v}^{\top} \widehat{V}_s \mathbf{v} \right) = O \left(\frac{1}{nh_t^2} \right) = o(1)$. By Chebyshev's inequality, (55) is proved. □

C.3 Proof of the Results for (Avg-SMSE)

Proof of Theorem 3.1

Proof. Since $\widehat{\beta}_{\text{SMSE},\ell}$ is the minimizer of $F_{h,\ell}(\beta) = \sum_{i \in \mathcal{H}_\ell} (-y_i) H' \left(\frac{x_i + \mathbf{z}_i^\top \beta}{h} \right)$, we have $\nabla_{\beta} F_{h,\ell} \left(\widehat{\beta}_{\text{SMSE},\ell} \right) = 0$. By Taylor's expansion of $\nabla_{\beta} F_{h,\ell}$ at β^* , we have

$$0 = \nabla_{\beta} F_{h,\ell}(\beta^*) + \nabla_{\beta}^2 F_{h,\ell}(\check{\beta}_\ell) \left(\widehat{\beta}_{\text{SMSE},\ell} - \beta^* \right),$$

where $\check{\beta}_\ell$ is between $\widehat{\beta}_{\text{SMSE},\ell}$ and β^* .

Define

$$U_{m,\ell,h}(\beta) := \frac{1}{m} \sum_{i \in \mathcal{H}_\ell} U_{h,i}(\beta),$$

$$V_{m,\ell,h}(\beta) := \frac{1}{m} \sum_{i \in \mathcal{H}_\ell} V_{h,i}(\beta),$$

where $U_{h,i}(\beta)$ and $V_{h,i}(\beta)$ are defined in (26) and (42). By definition, $\nabla_{\beta}^2 F_{h,\ell}(\check{\beta}_\ell) = V_{m,\ell,h}(\check{\beta}_\ell)$ and $\nabla_{\beta} F_{h,\ell}(\beta^*) = -U_{m,\ell,h}(\beta^*)$. Horowitz (1992) showed that

$$\left\| \widehat{\beta}_{\text{SMSE},\ell} - \beta^* \right\|_2 = O_{\mathbb{P}} \left(h^\alpha + \frac{1}{\sqrt{mh}} \right),$$

and thus by the proof of (22), if $\log m/mh^3 \rightarrow 0$, we have $\left\| V^{-1} - V_{m,\ell,h}^{-1}(\check{\beta}_\ell) \right\|_2 = o(1)$ with probability at least $1 - 2(5m^{-\gamma})^p$, where γ is a large positive constant. Hence, we have

$$\sup_{\ell} \left\| V^{-1} - V_{m,\ell,h}^{-1}(\check{\beta}_\ell) \right\|_2 = o_{\mathbb{P}}(1).$$

Therefore,

$$\widehat{\beta}_{(\text{Avg-SMSE})} - \beta^* = \frac{1}{L} \sum_{\ell=1}^L \left(\widehat{\beta}_{\text{SMSE},\ell} - \beta^* \right) = V^{-1} U_{n,h}(\beta^*) + U_{n,h}(\beta^*) o_{\mathbb{P}}(1),$$

using that $U_{n,h_m}(\beta) = \frac{1}{n} \sum_{i=1}^n U_{h,i}(\beta)$ and $n = mL$.

By Lemma C.1,

$$\mathbb{E} U_{n,h}(\beta^*) = U h^\alpha + o_{\mathbb{P}}(h^\alpha),$$

and

$$\sqrt{mLh} V^{-1} [U_{n,h}(\beta^*) - \mathbb{E} U_{n,h}(\beta^*)] \xrightarrow{d} \mathcal{N}(0, V^{-1} V_s V^{-1}).$$

If $L = o\left(m^{\frac{2}{3}(\alpha-1)}/(\log m)^{\frac{2\alpha+1}{3}}\right)$ and $h = \left(\frac{\lambda_h}{n}\right)^{\frac{1}{2\alpha+1}}$, we have $\frac{\log m}{mh^3} = o(1)$ and $\sqrt{\lambda_h}h^{-\alpha} = \sqrt{mLh}$, and thus

$$\begin{aligned} & \sqrt{mLh} \left(\hat{\beta}_{(\text{Avg-SMSE})} - \beta^* \right) \\ &= \sqrt{mLh} V^{-1} (U_{n,h}(\beta^*) - \mathbb{E}U_{n,h}(\beta^*) + \mathbb{E}U_{n,h}(\beta^*)) \\ & \quad + o_{\mathbb{P}}(1) \sqrt{mLh} (U_{n,h}(\beta^*) - \mathbb{E}U_{n,h}(\beta^*) + \mathbb{E}U_{n,h}(\beta^*)) \\ & \xrightarrow{d} \mathcal{N} \left(\sqrt{\lambda_h} V^{-1} U, V^{-1} V_s V^{-1} \right), \end{aligned}$$

which already proves (9). Furthermore, if $h \gtrsim \left(\frac{1}{mL}\right)^{\frac{1}{2\alpha+1}}$, we have $h^{-\alpha} \lesssim \sqrt{mLh}$, and thus

$$h^{-\alpha} \left(\hat{\beta}_{(\text{Avg-SMSE})} - \beta^* - V^{-1} \mathbb{E}U_{n,h}(\beta^*) \right) = o_{\mathbb{P}}(1).$$

If $h \lesssim \left(\frac{1}{mL}\right)^{\frac{1}{2\alpha+1}}$ and still satisfies $\frac{\log m}{mh^3} = o(1)$, we have $h^{-\alpha} \gtrsim \sqrt{mLh}$, and thus

$$\sqrt{mLh} \mathbb{E}U_{n,h}(\beta^*) = o_{\mathbb{P}}(1),$$

which yields

$$\sqrt{mLh} \left(\hat{\beta}_{(\text{Avg-SMSE})} - \beta^* \right) \xrightarrow{d} \mathcal{N} \left(0, V^{-1} V_s V^{-1} \right).$$

□

C.4 Proof of the Results for (wAvg-SMSE) and (wmSMSE)

Proof of Theorem A.1

Proof. Analogous to the proof of Theorem 3.1,

$$\hat{\beta}_{(\text{wAvg-SMSE})} - \beta^* = \sum_{\ell=1}^L W_{\ell} V_{\ell}^{-1} U_{m_{\ell}, \ell, h_{\ell}}(\beta^*) + o_{\mathbb{P}}(1) \sum_{\ell=1}^L W_{\ell} U_{m_{\ell}, \ell, h_{\ell}}. \quad (56)$$

By the uniformness, the bias satisfies

$$\mathbb{E} \left[\sum_{\ell=1}^L W_{\ell} V_{\ell}^{-1} U_{m_{\ell}, \ell, h_{\ell}}(\beta^*) \right] = \sum_{\ell=1}^L W_{\ell} V_{\ell}^{-1} U_{\ell} h^{\alpha} + o(h^{\alpha}) = B_{(\text{wAvg-SMSE})} + o(h^{\alpha}), \quad (57)$$

where $B_{(\text{wAvg-SMSE})} := \sum_{\ell=1}^L W_{\ell} V_{\ell}^{-1} U_{\ell} h^{\alpha}$. For any \mathbf{v} satisfying $\|\mathbf{v}\|_2 = 1$, define $\Phi_i := (1 - \mathbb{E}) m_{\ell}^{-1} \mathbf{v}^{\top} W_{\ell} V_{\ell}^{-1} U_{h_{\ell}, i}(\beta^*)$ where $i \in \mathcal{H}_{\ell}$. (Note that the definitions of Φ_i are different in the proof of different theorems, but they are parallel.) Then

$$(1 - \mathbb{E}) \left[\sum_{\ell=1}^L \mathbf{v}^{\top} W_{\ell} V_{\ell}^{-1} U_{m_{\ell}, \ell, h_{\ell}}(\beta^*) \right] = (1 - \mathbb{E}) \left[\sum_{\ell=1}^L \sum_{i \in \mathcal{H}_{\ell}} m_{\ell}^{-1} \mathbf{v}^{\top} W_{\ell} V_{\ell}^{-1} U_{h_{\ell}, i}(\beta^*) \right] = \sum_{i=1}^n \Phi_i. \quad (58)$$

Recall that by (30) and (36),

$$\sup_i |\Phi_i| = O\left(\max_\ell \frac{\|W_\ell\|_2}{m_\ell h}\right),$$

$$s_n^2 := \sum_{i=1}^n \text{var} [\Phi_i] = h^{-1} \sum_{\ell=1}^L m_\ell^{-1} \mathbf{v}^\top W_\ell V_\ell^{-1} [V_{s,\ell} + o(1)] V_\ell^{-1} W_\ell^\top \mathbf{v}.$$

The Lindeberg's condition is satisfied, since

$$\frac{1}{s_n^2} \sum_{i=1}^n \mathbb{E} \left[(\Phi_i)^2 I(|\Phi_i| > \varepsilon s_n) \right] \lesssim \frac{\sup_i |\Phi_i|^2}{s_n^2} \sum_{i=1}^n \frac{E[\Phi_i^2]}{\varepsilon^2 s_n^2} \lesssim \frac{\max_\ell [\|W_\ell\|_2^2 / m_\ell^2]}{h \sum_{\ell=1}^L [\|W_\ell\|_2^2 / m_\ell]} \asymp \frac{1}{nh} \rightarrow 0.$$

Therefore, using the fact that $s_n^2 \rightarrow \mathbf{v}^\top \Sigma_{(\mathbf{w}\text{Avg-SMSE})} \mathbf{v}$, we have

$$\left(\mathbf{v}^\top \Sigma_{(\mathbf{w}\text{Avg-SMSE})} \mathbf{v} \right)^{-1/2} \sum_{i=1}^n \Phi_i \xrightarrow{d} \mathcal{N}(0, 1),$$

where $\Sigma_{(\mathbf{w}\text{Avg-SMSE})} := h^{-1} \sum_{\ell=1}^L m_\ell^{-1} W_\ell V_\ell^{-1} V_{s,\ell} V_\ell^{-1} W_\ell^\top$. Now we do the following decomposition.

$$\begin{aligned} & \left(\mathbf{v}^\top \Sigma_{(\mathbf{w}\text{Avg-SMSE})} \mathbf{v} \right)^{-1/2} \mathbf{v}^\top \left(\hat{\beta}_{(\mathbf{w}\text{Avg-SMSE})} - \beta^* - B_{(\mathbf{w}\text{Avg-SMSE})} \right) \\ &= \left(\mathbf{v}^\top \Sigma_{(\mathbf{w}\text{Avg-SMSE})} \mathbf{v} \right)^{-1/2} \sum_{i=1}^n \Phi_i + \left(\mathbf{v}^\top \Sigma_{(\mathbf{w}\text{Avg-SMSE})} \mathbf{v} \right)^{-1/2} o(h^\alpha) \\ & \quad + o_{\mathbb{P}}(1) \left(\mathbf{v}^\top \Sigma_{(\mathbf{w}\text{Avg-SMSE})} \mathbf{v} \right)^{-1/2} \sum_{\ell=1}^L W_\ell U_{m_\ell, \ell, h_\ell}(\beta^*). \end{aligned} \tag{59}$$

When $\left(\mathbf{v}^\top \Sigma_{(\mathbf{w}\text{Avg-SMSE})} \mathbf{v} \right)^{1/2} \asymp h^\alpha$, the second term on the RHS of (59) is $o(1)$. Meanwhile, it is straightforward to show that $\left(\mathbf{v}^\top \Sigma_{(\mathbf{w}\text{Avg-SMSE})} \mathbf{v} \right)^{-1/2} \sum_{\ell=1}^L W_\ell U_{m_\ell, \ell, h_\ell}(\beta^*)$ is also asymptotically normal by repeating the previous procedure. Hence,

$$\left(\mathbf{v}^\top \Sigma_{(\mathbf{w}\text{Avg-SMSE})} \mathbf{v} \right)^{-1/2} \mathbf{v}^\top \left(\hat{\beta}_{(\mathbf{w}\text{Avg-SMSE})} - \beta^* - B_{(\mathbf{w}\text{Avg-SMSE})} \right) = \left(\mathbf{v}^\top \Sigma_{(\mathbf{w}\text{Avg-SMSE})} \mathbf{v} \right)^{-1/2} \sum_{i=1}^n \Phi_i + o_{\mathbb{P}}(1),$$

which converges to $\mathcal{N}(0, 1)$ in distribution. In particular, when Assumption 4 is satisfied, the matrix $\sum_{\ell=1}^L \frac{n}{m_\ell} W_\ell V_\ell^{-1} V_{s,\ell} V_\ell^{-1} W_\ell^\top$ is finite and $\left(\mathbf{v}^\top \Sigma_{(\mathbf{w}\text{Avg-SMSE})} \mathbf{v} \right)^{1/2} \asymp 1/\sqrt{nh}$. Then the above results directly lead to Theorem A.1.

Additionally, we show that the minimum of the asymptotic variance in (4),

$$\sum_{\ell=1}^L \frac{n}{m_\ell} W_\ell V_\ell^{-1} V_{s,\ell} V_\ell^{-1} W_\ell^\top,$$

is minimized at

$$W_\ell^{*,(\text{wAvg-SMSE})} = \left(\sum_{\ell=1}^L m_\ell V_\ell V_{s,\ell}^{-1} V_\ell \right)^{-1} m_\ell V_\ell V_{s,\ell}^{-1} V_\ell,$$

in the sense of both trace and Frobenius norm. For trace, we want to solve the following optimization problem:

$$\min_{W_\ell} \text{trace} \left(\sum_{\ell=1}^L \frac{n}{m_\ell} W_\ell V_\ell^{-1} V_{s,\ell} V_\ell^{-1} W_\ell^\top \right), \text{ s.t. } \sum_{\ell=1}^L W_\ell = I_{p \times p}.$$

The Lagrangian is

$$\text{trace} \left(\sum_{\ell=1}^L \frac{n}{m_\ell} W_\ell V_\ell^{-1} V_{s,\ell} V_\ell^{-1} W_\ell^\top \right) + \left\langle \Lambda, \sum_{\ell=1}^L W_\ell - I_{p \times p} \right\rangle.$$

By taking derivative w.r.t W_ℓ and letting the derivative be zero, we obtain $2 \frac{n}{m_\ell} W_\ell V_\ell^{-1} V_{s,\ell} V_\ell^{-1} + \Lambda = \mathbf{0}$ or $W_\ell \propto (m_\ell^{-1} V_\ell^{-1} V_{s,\ell} V_\ell^{-1})^{-1} = m_\ell V_\ell V_{s,\ell}^{-1} V_\ell$. By the constraint that $\sum_{\ell=1}^L W_\ell = I_{p \times p}$, we find the minimizer $W_\ell^{*,(\text{wAvg-SMSE})}$ defined above. The proof for Frobenius norm is similar. \square

Proof of Theorem A.2

Proof. For weighted (mSMSE), we have

$$\hat{\beta}^{(t)} - \beta^* = \left(\sum_{\ell=1}^L \sum_{i \in \mathcal{H}_\ell} m_\ell^{-1} W_\ell V_{h,i} \left(\hat{\beta}^{(t-1)} \right) \right)^{-1} \left(\sum_{\ell=1}^L \sum_{i \in \mathcal{H}_\ell} m_\ell^{-1} W_\ell U_{h,i} \left(\hat{\beta}^{(t-1)} \right) \right).$$

Define the weighted sample size $\bar{n}_W := 1 / \left(\sum_{\ell=1}^L \frac{\|W_\ell\|_2^2}{m_\ell} \right)$. We will show the one-step error under the assumptions in Proposition 3.2 (replacing n by \bar{n}_W),

$$\hat{\beta}^{(1)} - \beta^* = O_{\mathbb{P}} \left(\sqrt{\frac{\log \bar{n}_W}{\bar{n}_W h}} + \delta_{m,0}^2 + h^\alpha \right). \quad (60)$$

The proof is parallel to the proof of Proposition 3.2, and we will only show the difference in the remaining part.

By the uniformness in the assumptions, similar to (39), we can obtain

$$\mathbb{E} \left[\mathbf{v}^\top \left(\sum_{\ell=1}^L \sum_{i \in \mathcal{H}_\ell} m_\ell^{-1} W_\ell U_{h,i} \left(\hat{\beta}^{(0)} \right) \right) \right] = O(\delta_{m,0}^2 + h^\alpha).$$

Let $\Phi_i = m_\ell^{-1} \mathbf{v}^\top W_\ell U_{h,i} \left(\hat{\beta}^{(0)} \right)$ with $i \in \mathcal{H}_\ell$. By (30) and (36),

$$\sup_i |\Phi_i| = O \left(\max_\ell \frac{\|W_\ell\|_2}{m_\ell h} \right),$$

$$\sum_{i=1}^n \text{var} [\Phi_i] = h^{-1} \sum_{\ell=1}^L m_{\ell}^{-1} \mathbf{v}^{\top} W_{\ell} [V_{s,\ell} + o(1)] W_{\ell}^{\top} \mathbf{v} = O\left(\frac{1}{\bar{n}_W h}\right).$$

By Bernstein's inequality, for any γ , there exists a constant C_W such that

$$\mathbb{P}\left(\sum_i \Phi_i \geq C_W \sqrt{\frac{p \log \bar{n}_W}{\bar{n}_W h}}\right) \leq 2(\bar{n}_W)^{-2\gamma p}.$$

Then following the same procedure as the proof of Proposition 3.2, we can obtain that

$$\left\| \sum_{\ell=1}^L \sum_{i \in \mathcal{H}_{\ell}} m_{\ell}^{-1} W_{\ell} U_{h,i} \left(\hat{\beta}^{(0)} \right) \right\|_2 = O_{\mathbb{P}} \left(\sqrt{\frac{\log \bar{n}_W}{\bar{n}_W h}} + \delta_{m,0}^2 + h^{\alpha} \right).$$

Similarly, we can also obtain the convergence of the Hessian matrix in the same way, given by

$$\left\| \sum_{\ell=1}^L \sum_{i \in \mathcal{H}_{\ell}} m_{\ell}^{-1} W_{\ell} V_{h,i} \left(\hat{\beta}^{(0)} \right) - \bar{V}_W \right\|_2 = O_{\mathbb{P}} \left(\sqrt{\frac{\log \bar{n}_W}{\bar{n}_W h^3}} + \delta_{m,0} + h \right).$$

Hence, (60) holds, and we use the same strategy to choose the bandwidth h_t as in Theorem 3.3. When t is large enough, $h \asymp \bar{n}_W^{-1/(2\alpha+1)}$ and the error $\delta_{t-1} = \left\| \hat{\beta}_{(\text{wmSMSE})}^{(t-1)} - \beta^* \right\|_2 = O\left(\bar{n}_W^{-\alpha/(2\alpha+1)} \log \bar{n}_W\right)$, which satisfies $\delta_{t-1}^2 = o(h^{\alpha})$. Therefore, similar to the proof of (51), we have

$$\mathbb{E} \left(\sum_{\ell=1}^L \sum_{i \in \mathcal{H}_{\ell}} m_{\ell}^{-1} W_{\ell} U_{h,i} \left(\hat{\beta}^{(t-1)} \right) \right) = \bar{U}_W h^{\alpha} + o(h^{\alpha}).$$

Also, similar to the proof of Theorem A.1, when $h^{\alpha} \asymp \left(h^{-1} \sum_{\ell=1}^L m_{\ell}^{-1} W_{\ell} V_{s,\ell} W_{\ell}^{\top} \right)^{1/2}$, which is equivalent to $h \asymp \bar{n}_W^{-1/(2\alpha+1)}$, we have

$$\left(h^{-1} \sum_{\ell=1}^L m_{\ell}^{-1} W_{\ell} V_{s,\ell} W_{\ell}^{\top} \right)^{-1/2} \left(\sum_{\ell=1}^L \sum_{i \in \mathcal{H}_{\ell}} m_{\ell}^{-1} W_{\ell} U_{h,i} \left(\hat{\beta}^{(t-1)} \right) - \bar{U}_W h^{\alpha} \right) \xrightarrow{d} \mathcal{N}(0, I_{p \times p}).$$

In particular, when Assumption 4 is satisfied, $\bar{n}_W \asymp n$, and $n \sum_{\ell=1}^L m_{\ell}^{-1} W_{\ell} V_{s,\ell} W_{\ell}^{\top}$ is a finite matrix. Then we obtain

$$n^{\alpha/(2\alpha+1)} \sum_{\ell=1}^L \sum_{i \in \mathcal{H}_{\ell}} m_{\ell}^{-1} W_{\ell} U_{h,i} \left(\hat{\beta}^{(t-1)} \right) \xrightarrow{d} \mathcal{N} \left(\bar{U}_W, n \sum_{\ell=1}^L m_{\ell}^{-1} W_{\ell} V_{s,\ell} W_{\ell}^{\top} \right).$$

Combining with

$$\hat{\beta}^{(t)} - \beta^* = \left(\sum_{\ell=1}^L \sum_{i \in \mathcal{H}_{\ell}} m_{\ell}^{-1} W_{\ell} V_{h,i} \left(\hat{\beta}^{(t-1)} \right) \right)^{-1} \left(\sum_{\ell=1}^L \sum_{i \in \mathcal{H}_{\ell}} m_{\ell}^{-1} W_{\ell} U_{h,i} \left(\hat{\beta}^{(t-1)} \right) \right),$$

and

$$\left\| \sum_{\ell=1}^L \sum_{i \in \mathcal{H}_\ell} m_\ell^{-1} W_\ell V_{h,i} \left(\widehat{\beta}^{(t-1)} \right) - \bar{V}_W \right\|_2 = O_{\mathbb{P}} \left(\sqrt{\frac{\log \bar{n}_W}{\bar{n}_W h^3}} + \delta_{m,0} + h_t \right),$$

by Slutsky's Theorem, (5) is true. □

C.5 Proof of the Results for the High-dimensional (mSMSE)

Before starting the proof, we first formalize our notation. Define

$$V_n(\beta) = \frac{1}{nh^2} \sum_{i=1}^n (-y_i) H'' \left(\frac{x_i + \mathbf{z}_i^\top \beta}{h} \right) \mathbf{z}_i \mathbf{z}_i^\top. \quad (61)$$

$$V_{m,\ell}(\beta) = \frac{1}{nh^2} \sum_{i \in \mathcal{H}_\ell} (-y_i) H'' \left(\frac{x_i + \mathbf{z}_i^\top \beta}{h} \right) \mathbf{z}_i \mathbf{z}_i^\top. \quad (62)$$

$$U_n(\beta) = \frac{1}{nh} \sum_{i=1}^n (-y_i) H' \left(\frac{x_i + \mathbf{z}_i^\top \beta}{h} \right) \mathbf{z}_i. \quad (63)$$

As claimed before, we will not compute $V_n(\beta)$ in the algorithm, but it is an important intermediate quantity in the theoretical analysis. Without loss of generality, we only consider $V_{m,1}(\beta)$ in the sequel. We first give the convergence property of $V_n(\widehat{\beta}^{(0)})$, $V_{m,1}(\widehat{\beta}^{(0)})$ and $U_n(\widehat{\beta}^{(0)})$, which is crucial for deriving the convergence rate of $\widehat{\beta}^{(1)}$. Note that we omit the dependence of these quantities on the bandwidth h in the notation.

Lemma C.3. *Assume Assumptions 1–5, 5 and 6 hold. Further assume that $\frac{\log m}{mh^3} = o(1)$, $\sqrt{s}\delta_{m,0} = O(h^{3/2})$ ($\delta_{m,0}$ is defined in Assumption 5) and $h = o(1)$, we have the following results:*

$$\left\| (1 - \mathbb{E}) \left[V_n(\widehat{\beta}^{(0)}) - V \right] \right\|_{\max} = O_{\mathbb{P}} \left(\sqrt{\frac{\log p}{nh^3}} \right), \quad (64)$$

$$\sup_{\|\mathbf{v}_1\|_2 = \|\mathbf{v}_2\|_2 = 1} \mathbf{v}_1^\top \left(\mathbb{E} \left[V_n(\widehat{\beta}^{(0)}) \right] - V \right) \mathbf{v}_2 = O_{\mathbb{P}} \left(\sqrt{s}\delta_{m,0} + h^\alpha \right), \quad (65)$$

$$\left\| (1 - \mathbb{E}) \left[V_{m,1}(\widehat{\beta}^{(0)}) - V \right] \right\|_{\max} = O_{\mathbb{P}} \left(\sqrt{\frac{\log p}{mh^3}} \right), \quad (66)$$

$$\sup_{\|\mathbf{v}_1\|_2 = \|\mathbf{v}_2\|_2 = 1} \mathbf{v}_1^\top \left(\mathbb{E} \left[V_{m,1}(\widehat{\beta}^{(0)}) \right] - V \right) \mathbf{v}_2 = O_{\mathbb{P}} \left(\sqrt{s}\delta_{m,0} + h^\alpha \right). \quad (67)$$

Additionally, define

$$\Psi_n(\widehat{\beta}^{(0)}) := U_n(\widehat{\beta}^{(0)}) - V_n(\widehat{\beta}^{(0)})(\widehat{\beta}^{(0)} - \beta^*),$$

and then we have

$$\left\| (1 - \mathbb{E}) \Psi_n(\widehat{\beta}^{(0)}) \right\|_{\infty} = O_{\mathbb{P}} \left(\sqrt{\frac{\log p}{nh}} \right), \quad (68)$$

and

$$\sup_{\|\mathbf{v}\|_2=1} \left| \mathbf{v}^{\top} \mathbb{E} \left[\Psi_n(\widehat{\beta}^{(0)}) \right] \right| = O_{\mathbb{P}} (s\delta_{m,0}^2 + h^{\alpha}). \quad (69)$$

Proof. Proof of (64):

Recall our definitions

$$V_n(\widehat{\beta}^{(0)}) = \frac{1}{nh^2} \sum_{i=1}^n (-y_i) H'' \left(\frac{x_i + \mathbf{z}_i^{\top} \widehat{\beta}^{(0)}}{h} \right) \mathbf{z}_i \mathbf{z}_i^{\top},$$

and

$$V = -2\mathbb{E} \left[\rho(0 | \mathbf{Z}) F'(0 | \mathbf{Z}) \mathbf{Z} \mathbf{Z}^{\top} \right].$$

By Assumption 5, there exists constant C_{ℓ_1}, C_{ℓ_2} such that $\mathbb{P}(\beta^{(0)} \in \Theta) \rightarrow 1$, where

$$\Theta := \{ \beta : \|\beta - \beta^*\|_2 \leq C_{\ell_1} \delta_{m,0}, \|\beta - \beta^*\|_1 \leq C_{\ell_2} \sqrt{s} \delta_{m,0} \}.$$

Without loss of generality, we assume $C_{\ell_1} = C_{\ell_2} = 1$ and $\beta^{(0)} \in \Theta$ in the following proof.

For each $(j_1, j_2) \in \{1, \dots, p\} \times \{1, \dots, p\}$, define

$$V_{n,j_1,j_2}(\beta) := \frac{1}{nh^2} \sum_{i=1}^n (-y_i) H'' \left(\frac{x_i + \mathbf{z}_i^{\top} \beta}{h} \right) z_{i,j_1} z_{i,j_2},$$

$$V_{j_1,j_2} := -2\mathbb{E} \left[\rho(0 | \mathbf{Z}) F'(0 | \mathbf{Z}) Z_{j_1} Z_{j_2} \right],$$

$$\phi_{ij_1j_2}^V(\beta) := \frac{-y_i}{h^2} H'' \left(\frac{x_i + \mathbf{z}_i^{\top} \beta}{h} \right) z_{i,j_1} z_{i,j_2} + \frac{y_i}{h^2} H'' \left(\frac{x_i + \mathbf{z}_i^{\top} \beta^*}{h} \right) z_{i,j_1} z_{i,j_2},$$

and

$$\Phi_{n,j_1,j_2}^V := \sup_{\beta \in \Theta} |(1 - \mathbb{E}) [V_{n,j_1,j_2}(\beta) - V_{n,j_1,j_2}(\beta^*)]| = \sup_{\beta \in \Theta} \left| (1 - \mathbb{E}) \frac{1}{n} \sum_{i=1}^n \phi_{ij_1j_2}^V(\beta) \right|.$$

Since

$$\begin{aligned} & \left\| (1 - \mathbb{E}) \left[V_n(\widehat{\beta}^{(0)}) - V \right] \right\|_{\max} \\ &= \sup_{j_1, j_2} \left| (1 - \mathbb{E}) \left[V_{n,j_1,j_2}(\widehat{\beta}^{(0)}) - V_{j_1,j_2} \right] \right| \\ &\leq \sup_{j_1, j_2} \Phi_{n,j_1,j_2}^V + \sup_{j_1, j_2} |(1 - \mathbb{E}) V_{n,j_1,j_2}(\beta^*)|, \end{aligned} \quad (70)$$

we break the proof of (64) into two steps, separately controlling the two terms in the last line of (70).

Step 1:

$$\sup_{j_1, j_2} \Phi_{n, j_1, j_2}^V = O_{\mathbb{P}} \left(\delta_{m,0} \sqrt{\frac{s \log p}{nh^6}} \right) \stackrel{\sqrt{s} \delta_{m,0} = O(h^{3/2})}{=} O_{\mathbb{P}} \left(\sqrt{\frac{\log p}{nh^3}} \right). \quad (71)$$

The proof in this step is analogous to the proof of Lemma B.1 in Luo et al. (2022). Since $|z_{i,j}|$ is upper bounded by \bar{B} , for any i and $\beta \in \Theta$, we have $\mathbf{z}_i^\top (\beta - \beta^*) \leq \bar{B} \|\beta - \beta^*\|_1 \leq \bar{B} \sqrt{s} \delta_{m,0}$. Since $H''(x)$ is Lipschitz, we have

$$\bar{\phi} := \sup_{i, j_1, j_2} \sup_{\beta \in \Theta} |\phi_{ij_1 j_2}^V(\beta)| = O \left(\frac{\sqrt{s} \delta_{m,0}}{h^3} \right). \quad (72)$$

Since $\rho(\cdot | \mathbf{Z})$ is bounded, we also have

$$\begin{aligned} & \sup_i \sup_{\beta \in \Theta} \mathbb{E}_{|\mathbf{Z}} \left[H'' \left(\frac{X + \mathbf{Z}^\top \beta}{h} \right) - H'' \left(\frac{X + \mathbf{Z}^\top \beta^*}{h} \right) \right]^2 \\ &= \sup_i \sup_{\beta \in \Theta} h \int_{-1}^1 \left[H'' \left(\xi + \frac{\mathbf{Z}^\top (\beta - \beta^*)}{h} \right) - H''(\xi) \right]^2 \rho(h\xi | \mathbf{Z}) \, d\xi \\ &= O(s \delta_{m,0}^2 / h), \end{aligned} \quad (73)$$

which implies that

$$\sup_{i, j_1, j_2} \sup_{\beta \in \Theta} \mathbb{E} \left[|\phi_{ij_1 j_2}^V(\beta)|^2 \right] = O(s \delta_{m,0}^2 / h^5). \quad (74)$$

Define $\sigma_1, \dots, \sigma_n$ to be independent Rademacher variables, i.e., binary variables that are uniformly distributed on $\{-1, +1\}$. By Rademacher symmetrization,

$$\mathbb{E} \Phi_{n, j_1, j_2}^V \leq 2 \mathbb{E} \sup_{\beta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \phi_{ij_1 j_2}^V(\beta) \right|.$$

Further, as

$$\phi_{ij_1 j_2}^V(\beta) := \frac{-y_i}{h^2} H'' \left(\frac{x_i + \mathbf{z}_i^\top \beta^* + \mathbf{z}_i^\top (\beta - \beta^*)}{h} \right) z_{i, j_1} z_{i, j_2} + \frac{y_i}{h^2} H'' \left(\frac{x_i + \mathbf{z}_i^\top \beta^*}{h} \right) z_{i, j_1} z_{i, j_2},$$

we can view $\phi_{ij_1 j_2}^V(\beta)$ as a function of $\mathbf{z}_i^\top (\beta - \beta^*)$ with Lipschitz constant $\asymp \bar{B}^2 / h^3$. By Talagrand's Lemma,

$$\mathbb{E}_\sigma \sup_{\beta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \phi_{ij_1 j_2}^V(\beta) \right| \lesssim \frac{1}{h^3} \mathbb{E}_\sigma \sup_{\beta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbf{z}_i^\top (\beta - \beta^*) \right| \lesssim \left(\frac{\sqrt{s} \delta_{m,0}}{nh^3} \right) \mathbb{E}_\sigma \left\| \sum_{i=1}^n \sigma_i \mathbf{z}_i \right\|_\infty.$$

Since $\sigma_i z_{i,j} \in [-\bar{B}, +\bar{B}]$ for all $j \in \{1, \dots, p\}$, by Hoeffding's inequality,

$$\mathbb{P} \left(\left| \sum_{i=1}^n \sigma_i z_{i,j} \right| \geq \sqrt{4\bar{B}^2 n \log \max \{n, p\}} \right) \leq 2 \exp \left(-\frac{4n\bar{B}^2 \log \max \{n, p\}}{2n\bar{B}^2} \right) = \frac{2}{\max \{n, p\}^2},$$

which implies that with probability at least $1 - \frac{2}{\max \{n, p\}}$,

$$\left\| \sum_{i=1}^n \sigma_i \mathbf{z}_i \right\|_{\infty} \leq 2\bar{B} \sqrt{n \log \max \{n, p\}}.$$

Assumption 6 supposes that $\log p = O(\log n)$, and thus we obtain

$$\mathbb{E}_{\sigma} \left\| \sum_{i=1}^n \sigma_i \mathbf{z}_i \right\|_{\infty} \leq 2\bar{B} \sqrt{n \log \max \{n, p\}} \left(1 - \frac{2}{\max \{n, p\}} \right) + \frac{2\bar{B}}{\max \{n, p\}} = O \left(\sqrt{n \log p} \right),$$

and hence

$$\mathbb{E} \Phi_{n,j_1,j_2}^V = O \left(\delta_{m,0} \sqrt{\frac{s \log p}{nh^6}} \right). \quad (75)$$

To show (71), we use Theorem 7.3 in Bousquet (2003), which is restated as the following Lemma C.4.

Lemma C.4 (Bousquet (2003)). *Assume $\{\mathbf{z}_i\}_{i=1}^n$ are identically distributed random variables. Let \mathcal{F} be a set of countable real-value functions such that all functions $f \in \mathcal{F}$ are measurable, square-integrable and satisfy $\mathbb{E} f(\mathbf{z}_i) = 0$. Assume $\sup_{f, \mathbf{z}} f(\mathbf{z}) \leq 1$. Define*

$$\Upsilon := \sup_{f \in \mathcal{F}} \sum_{i=1}^n f(\mathbf{z}_i).$$

If $\sum_{i=1}^n \sup_{f \in \mathcal{F}} \mathbb{E} f^2(\mathbf{z}_i) \leq n\sigma^2$, then for all $x > 0$, we have

$$\mathbb{P} \left(\Upsilon > \mathbb{E} \Upsilon + \sqrt{2x(n\sigma^2 + 2\mathbb{E} \Upsilon)} + \frac{x}{3} \right) < e^{-x}.$$

Note that Lemma C.4 requires \mathcal{F} to be countable. We first apply Lemma C.4 to prove (71) on rational β , i.e.,

$$\sup_{j_1, j_2} \sup_{\beta \in \Theta \cap \mathbb{Q}^p} \left| (1 - \mathbb{E}) \frac{1}{n} \sum_{i=1}^n \phi_{ij_1 j_2}^V(\beta) \right| = O_{\mathbb{P}} \left(\delta_{m,0} \sqrt{\frac{s \log n}{nh^6}} \right).$$

Fix j_1, j_2 and take

$$\mathcal{F} := \left\{ f_{\beta}(\mathbf{z}_i) = \frac{(1 - \mathbb{E}) \phi_{ij_1 j_2}^V(\beta)}{2\bar{\phi}} : \beta \in \Theta \cap \mathbb{Q}^p \right\}.$$

By (72), (74) and (75), we have $f(\mathbf{z}_i) \leq 1$,

$$\sum_{i=1}^n \sup_{f_{\beta} \in \mathcal{F}} \mathbb{E} f_{\beta}^2(\mathbf{z}_i) = O\left(\frac{ns\delta_{m,0}^2}{\bar{\phi}^2 h^5}\right),$$

and

$$\mathbb{E}\Upsilon = O\left(\frac{\delta_{m,0}}{\bar{\phi}h^3} \sqrt{ns \log p}\right).$$

By Lemma C.4, for all $x > 0$, with probability $1 - e^{-x}$,

$$\frac{1}{n} \sup_{\beta \in \Theta \cap \mathbb{Q}^p} (1 - \mathbb{E}) \sum_{i=1}^n \phi_{ij_1 j_2}^V(\beta) = O\left(\frac{\delta_{m,0}}{h^3} \sqrt{\frac{s \log p}{n}} + \sqrt{2x \frac{s\delta_{m,0}^2}{nh^5} + 4x \frac{\bar{\phi}\delta_{m,0} \sqrt{s \log p}}{n^{3/2}h^3}} + \frac{\bar{\phi}x}{3n}\right).$$

By taking $x = 3 \log \max\{n, p\}$, plugging (72) in and using that $\log p = O(\log n)$ again, the above bound can be written as

$$O\left[\delta_{m,0} \left(\sqrt{\frac{s \log p}{nh^6}} + \sqrt{\frac{s \log p}{nh^5} + \frac{s \log p}{n^{3/2}h^6}} + \frac{\sqrt{s \log p}}{nh^3}\right)\right] = O\left(\delta_{m,0} \sqrt{\frac{s \log p}{nh^6}}\right).$$

For the same reason, with probability $1 - 1/\max\{n, p\}^3$,

$$\frac{1}{n} \sup_{\beta \in \Theta \cap \mathbb{Q}^p} (1 - \mathbb{E}) \sum_{i=1}^n [-\phi_{ij_1 j_2}^V(\beta)] = O\left(\delta_{m,0} \sqrt{\frac{s \log p}{nh^6}}\right).$$

Therefore, with probability $1 - 2/\max\{n, p\}^3$,

$$\sup_{\beta \in \Theta \cap \mathbb{Q}^p} \left| (1 - \mathbb{E}) \frac{1}{n} \sum_{i=1}^n \phi_{ij_1 j_2}^V(\beta) \right| = O\left(\delta_{m,0} \sqrt{\frac{s \log p}{nh^6}}\right).$$

By the continuity of $\phi_{ij_1 j_2}^V(\beta)$,

$$\sup_{\beta \in \Theta} \left| (1 - \mathbb{E}) \frac{1}{n} \sum_{i=1}^n \phi_{ij_1 j_2}^V(\beta) \right| = O\left(\delta_{m,0} \sqrt{\frac{s \log p}{nh^6}}\right),$$

with the same probability. This is true for any j_1, j_2 , so

$$\sup_{j_1, j_2} \Phi_{n, j_1, j_2}^V = O\left(\delta_{m,0} \sqrt{\frac{s \log p}{nh^6}}\right) \stackrel{\sqrt{s}\delta_{m,0}=O(h^{3/2})}{=} O\left(\sqrt{\frac{\log p}{nh^3}}\right),$$

with probability at least $1 - 2/\max\{n, p\}$, which completes the proof of (71).

Step 2:

$$\sup_{j_1, j_2} |(1 - \mathbb{E}) V_{n, j_1, j_2}(\beta^*)| = O_{\mathbb{P}}\left(\sqrt{\frac{\log p}{nh^3}}\right). \quad (76)$$

Recall that

$$V_{n,j_1,j_2}(\beta^*) := \frac{1}{nh^2} \sum_{i=1}^n (-y_i) H'' \left(\frac{x_i + \mathbf{z}_i^\top \beta^*}{h} \right) z_{i,j_1} z_{i,j_2}.$$

We have

$$\sup_i \left| \frac{-y_i}{h^2} H'' \left(\frac{x_i + \mathbf{z}_i^\top \beta^*}{h} \right) z_{i,j_1} z_{i,j_2} \right| = O(1/h^2),$$

and

$$\sup_i \mathbb{E} \left| \frac{-y_i}{h^2} H'' \left(\frac{x_i + \mathbf{z}_i^\top \beta^*}{h} \right) z_{i,j_1} z_{i,j_2} \right|^2 = O(1/h^3),$$

as $H''(x)$, $|z_{i,j}|$ is bounded and

$$\mathbb{E}_{\cdot|\mathbf{Z}} \left[H'' \left(\frac{X + \mathbf{Z}^\top \beta^*}{h} \right) \right]^2 = h \int_{-1}^1 [H''(\xi)]^2 \rho(\xi h | \mathbf{Z}) d\xi = O(h). \quad (77)$$

By Bernstein's inequality, there exists a constant $C > 0$,

$$\mathbb{P} \left(|(1 - \mathbb{E}) V_{n,j_1,j_2}(\beta^*)| \geq \sqrt{\frac{C_2 \log \max\{n, p\}}{nh^3}} \right) \leq 2 \exp \left(\frac{-\frac{C}{2} \log \max\{n, p\}}{1 + \frac{1}{3} \sqrt{C \log \max\{n, p\} / (nh)}} \right).$$

Our assumptions $\log m/mh^3 = o(1)$, $m > n^c$ and $p = O(n^\gamma)$ ensure that $\log \max\{n, p\} / (nh) = o(1)$, and hence we can take large enough C to make

$$2 \exp \left(\frac{-\frac{C}{2} \log \max\{n, p\}}{1 + \frac{1}{3} \sqrt{C \log \max\{n, p\} / (nh)}} \right) \leq \frac{2}{\max\{n, p\}^3}.$$

This implies that

$$\mathbb{P} \left(\sup_{j_1, j_2} |(1 - \mathbb{E}) V_{n,j_1,j_2}(\beta^*)| \leq \sqrt{\frac{C \log \max\{n, p\}}{nh^3}} \right) \geq 1 - \frac{2}{\max\{n, p\}^3},$$

which proves (76). Together with (70) and (71), we conclude the proof of (64).

Proof of (65):

Recall that, by Equation (38), for almost every \mathbf{Z} ,

$$\begin{aligned} & (2F(-t | \mathbf{Z}) - 1) \rho(t | \mathbf{Z}) \\ &= 2F^{(1)}(0 | \mathbf{Z}) \rho(0 | \mathbf{Z}) t + \sum_{k=2}^{2\alpha+1} M_k(\mathbf{Z}) t^k, \end{aligned} \quad (78)$$

where $M_k(\mathbf{Z})$ is a constant depending on \mathbf{Z} , t' and t'' . Since $\rho^{(k)}(\cdot | \mathbf{Z})$ and $F^{(k)}(\cdot | \mathbf{Z})$ are bounded around 0 for all k , we know there exists a constant M such that $\sup_k |M_k(\mathbf{Z})| \leq M$ for all \mathbf{Z}, t', t'' .

In the following computation, we let $t = \xi h - \mathbf{Z}^\top \Delta(\beta)$, where $\Delta(\beta) := \beta - \beta^*$.

Recall that when $x > 1$ or $x < -1$, $H'(x) = H''(x) = 0$. The kernel $H'(x)$ is bounded, $\int_{-1}^1 H'(x) dx = 1$, and $\int_{-1}^1 x^k H'(x) dx = 0$ for all $1 \leq k \leq \alpha - 1$. Moreover, $\int_{-1}^1 x H''(x) dx = -1$ and $\int_{-1}^1 x^k H''(x) dx = 0$ for $k = 0$ and $2 \leq k \leq \alpha$. Also, recall that $\zeta = X + \mathbf{Z}^\top \boldsymbol{\beta}^*$ and $-y = -\text{sign}(y^*) = -\text{sign}(\mathbf{Z} + \epsilon) = 2\mathbb{I}(\mathbf{Z} + \epsilon < 0) - 1$.

For all $\mathbf{v}_1, \mathbf{v}_2$ that satisfies $\|\mathbf{v}_1\|_2 = \|\mathbf{v}_2\|_2 = 1$ and $\boldsymbol{\beta} \in \Theta$,

$$\begin{aligned}
& \mathbb{E}_{\cdot|\mathbf{Z}} \left[\frac{(\mathbf{v}_1^\top \mathbf{Z})(\mathbf{v}_2^\top \mathbf{Z})}{h^2} (-y) H'' \left(\frac{X + \mathbf{Z}^\top \boldsymbol{\beta}}{h} \right) \right] \\
&= \frac{(\mathbf{v}_1^\top \mathbf{Z})(\mathbf{v}_2^\top \mathbf{Z})}{h^2} \mathbb{E}_{\cdot|\mathbf{Z}} [2\mathbb{I}(\mathbf{Z} + \epsilon < 0) - 1] H'' \left(\frac{\mathbf{Z}^\top \Delta(\boldsymbol{\beta}) + \zeta}{h} \right) \\
&= \frac{(\mathbf{v}_1^\top \mathbf{Z})(\mathbf{v}_2^\top \mathbf{Z})}{h} \int_{-1}^1 \left[2F(\mathbf{Z}^\top \Delta(\boldsymbol{\beta}) - \xi h | \mathbf{Z}) - 1 \right] \rho(\xi h - \mathbf{Z}^\top \Delta(\boldsymbol{\beta}) | \mathbf{Z}) H''(\xi) d\xi \quad (79) \\
&= \frac{(\mathbf{v}_1^\top \mathbf{Z})(\mathbf{v}_2^\top \mathbf{Z})}{h} {}_2F^{(1)}(0 | \mathbf{Z}) \rho(0 | \mathbf{Z}) \int_{-1}^1 (\xi h - \mathbf{Z}^\top \Delta(\boldsymbol{\beta})) H''(\xi) d\xi \\
&\quad + \frac{(\mathbf{v}_1^\top \mathbf{Z})(\mathbf{v}_2^\top \mathbf{Z})}{h} \cdot \sum_{k=2}^{2\alpha+1} M_k(\mathbf{Z}) \int_{-1}^1 (\xi h - \mathbf{Z}^\top \Delta(\boldsymbol{\beta}))^k H''(\xi) d\xi
\end{aligned}$$

For $1 \leq k \leq \alpha$,

$$\int_{-1}^1 (\xi h - \mathbf{Z}^\top \Delta(\boldsymbol{\beta}))^k H''(\xi) d\xi = \sum_{k'=0}^k h^{k'} (\mathbf{Z}^\top \Delta(\boldsymbol{\beta}))^{k-k'} \int_{-1}^1 \xi^{k'} H''(\xi) d\xi = h (\mathbf{Z}^\top \Delta(\boldsymbol{\beta}))^{k-1}.$$

For $\alpha + 1 \leq k \leq 2\alpha + 1$, since $H''(x)$ is bounded,

$$\begin{aligned}
\left| \int_{-1}^1 (\xi h - \mathbf{Z}^\top \Delta(\boldsymbol{\beta}))^k H''(\xi) d\xi \right| &\leq \int_{-1}^1 2^{k-1} \left(|\xi h|^k + |\mathbf{Z}^\top \Delta(\boldsymbol{\beta})|^k \right) |H''(\xi)| d\xi \\
&\leq 2^{2\alpha} \sup_x |H''(x)| \left[h^{\alpha+1} + |\mathbf{Z}^\top \Delta(\boldsymbol{\beta})|^k \right] \\
&\leq 2^{2\alpha} (1 + \overline{B}^k) \sup_x |H''(x)| h^{\alpha+1}.
\end{aligned}$$

The last inequality holds because $|Z_j| \leq \overline{B}$, $\mathbf{Z}^\top \Delta(\boldsymbol{\beta}) \leq \overline{B} \sqrt{s} \delta_{m,0}$, $\sqrt{s} \delta_{m,0} = o(h^{3/2})$ and $h = o(1)$.

Hence,

$$\begin{aligned}
& \mathbb{E}_{\cdot|\mathbf{z}} \left[\frac{(\mathbf{v}_1^\top \mathbf{Z})(\mathbf{v}_2^\top \mathbf{Z})}{h} (-y) H'' \left(\frac{X + \mathbf{Z}^\top \boldsymbol{\beta}}{h} \right) \right] - \frac{(\mathbf{v}_1^\top \mathbf{Z})(\mathbf{v}_2^\top \mathbf{Z})}{h} {}_2F^{(1)}(0 | \mathbf{Z}) \rho(0 | \mathbf{Z}) \\
&\leq \left| \frac{(\mathbf{v}_1^\top \mathbf{Z})(\mathbf{v}_2^\top \mathbf{Z})}{h} \sum_{k=2}^{\alpha} M_k(\mathbf{Z}) h (\mathbf{Z}^\top \Delta(\boldsymbol{\beta}))^{k-1} \right| \\
&\quad + \left| \frac{(\mathbf{v}_1^\top \mathbf{Z})(\mathbf{v}_2^\top \mathbf{Z})}{h} 2^{2\alpha} \sup_x |H''(x)| \sum_{k=\alpha+1}^{2\alpha+1} (1 + \overline{B}^k) M_k(\mathbf{Z}) h^{\alpha+1} \right| \\
&\leq C_{EV} (\mathbf{v}_1^\top \mathbf{Z}) (\mathbf{v}_2^\top \mathbf{Z}) (\sqrt{s} \delta_{m,0} + h^\alpha),
\end{aligned}$$

where

$$C_{EV} = M \left((\alpha - 1) [\max \{1, \overline{B}\}]^{\alpha-1} + (\alpha + 1) 2^{2\alpha} \left(1 + [\max \{1, \overline{B}\}]^{2\alpha+1} \right) \sup_x |H''(x)| \right)$$

is a constant not depending on β and \mathbf{Z} . Therefore,

$$\begin{aligned} & \sup_{\beta \in \Theta} \sup_{\|\mathbf{v}_1\|_2 = \|\mathbf{v}_2\|_2 = 1} \mathbf{v}_1^\top (\mathbb{E}[V_n(\beta)] - V) \mathbf{v}_2 \\ &= \sup_{\beta \in \Theta} \sup_{\|\mathbf{v}_1\|_2 = \|\mathbf{v}_2\|_2 = 1} \mathbb{E} \left[\frac{(\mathbf{v}_1^\top \mathbf{Z})(\mathbf{v}_2^\top \mathbf{Z})}{h} (-y) H'' \left(\frac{\mathbf{Z}^\top \beta}{h} \right) - \frac{(\mathbf{v}_1^\top \mathbf{Z})(\mathbf{v}_2^\top \mathbf{Z})}{h} {}_2F^{(1)}(0 | \mathbf{Z}) \rho(0 | \mathbf{Z}) \right] \\ &\lesssim \sqrt{s} \delta_{m,0} + h^\alpha, \end{aligned} \tag{80}$$

which completes the proof of (65).

Proof of (66) and (67):

Equation (66) and (67) can be shown in the same way as above by replacing all the n with m .

Proof of (68):

The proof of (68) is analogous to that of (64). We will omit some details since they are the same. For each $j \in \{1, \dots, p\}$, define

$$\begin{aligned} U_{n,h,j}(\beta) &:= \frac{1}{nh^2} \sum_{i=1}^n (-y_i) H'' \left(\frac{x_i + \mathbf{z}_i^\top \beta}{h} \right) \mathbf{z}_i^\top (\beta - \beta^*) \mathbf{z}_{i,j} \\ &\quad - \frac{1}{nh} \sum_{i=1}^n (-y_i) H' \left(\frac{x_i + \mathbf{z}_i^\top \beta}{h} \right) \mathbf{z}_{i,j}. \end{aligned}$$

Then we have

$$\begin{aligned} & \left\| (1 - \mathbb{E}) \Psi_n(\widehat{\beta}^{(0)}) \right\|_\infty \\ &= \left\| (1 - \mathbb{E}) \left[\frac{1}{nh^2} \sum_{i=1}^n (-y_i) H'' \left(\frac{x_i + \mathbf{z}_i^\top \widehat{\beta}^{(0)}}{h} \right) \mathbf{z}_i^\top (\widehat{\beta}^{(0)} - \beta^*) \mathbf{z}_i - \frac{1}{nh} \sum_{i=1}^n (-y_i) H' \left(\frac{x_i + \mathbf{z}_i^\top \beta^{(0)}}{h} \right) \mathbf{z}_i \right] \right\|_\infty \\ &\leq \sup_j \sup_{\beta \in \Theta} |(1 - \mathbb{E})(U_{n,h,j}(\beta) - U_{n,h,j}(\beta^*))| + \sup_j |(1 - \mathbb{E}) U_{n,h,j}(\beta^*)|. \end{aligned}$$

Define

$$\phi_{i,j}^U(\beta) := \left| \frac{-y_i}{h^2} H'' \left(\frac{x_i + \mathbf{z}_i^\top \beta}{h} \right) \mathbf{z}_i^\top (\beta - \beta^*) \mathbf{z}_{i,j} - \frac{-y_i}{h} \left[H' \left(\frac{x_i + \mathbf{z}_i^\top \beta}{h} \right) - H' \left(\frac{x_i + \mathbf{z}_i^\top \beta^*}{h} \right) \right] \mathbf{z}_{i,j} \right|.$$

$$\Phi_{n,h,j}^U := \sup_{\beta \in \Theta} |(1 - \mathbb{E})(U_{n,h,j}(\beta) - U_{n,h,j}(\beta^*))| = \sup_{\beta \in \Theta} \left| (1 - \mathbb{E}) \frac{1}{n} \sum_{i=1}^n \phi_{i,j}^U(\beta) \right|.$$

Similar to the analysis of $\phi_{i,j}^V$, we have

$$\sup_{i,j} \sup_{\beta \in \Theta} |\phi_{i,j}^U(\beta)| = O\left(\frac{\sqrt{s}\delta_{m,0}}{h^2}\right),$$

and

$$\sup_{i,j} \sup_{\beta \in \Theta} \mathbb{E} [\phi_{i,j}^U(\beta)]^2 = O\left(\frac{s\delta_{m,0}^2}{h^3}\right).$$

By Rademacher symmetrization, Talagrand's concentration principle and Hoeffding's inequality,

$$\mathbb{E}\Phi_{n,h,j}^U \leq 2\mathbb{E} \sup_{\beta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \phi_{i,j}^U(\beta) \right| \lesssim \left(\frac{\sqrt{s}\delta_{m,0}}{nh^2} \right) \cdot \mathbb{E} \left(\mathbb{E}_\sigma \left\| \sum_{i=1}^n \sigma_i \mathbf{z}_i^\top \right\|_\infty \right) \lesssim \delta_{m,0} \sqrt{\frac{s \log p}{nh^4}}.$$

(Details are the same as the proof of (75), while the only difference is that $\phi_{ij}^U(\beta)$ is a Lipschitz function of $\mathbf{z}_i^\top(\beta - \beta^*)$ with Lipschitz constant $\asymp 1/h^2$, instead of $1/h^3$ for $\phi_{ij}^V(\beta)$.)

Using Lemma C.4 again, we can show that

$$\sup_j \Phi_{n,h,j}^U = O_{\mathbb{P}} \left(\delta_{m,0} \sqrt{\frac{s \log p}{nh^4}} \right) \stackrel{\sqrt{s}\delta_{m,0}=O(h^{3/2})}{=} O \left(\sqrt{\frac{\log p}{nh}} \right).$$

Similar to the proof of (76), we have

$$\sup_{\beta \in \Theta} \sup_{i,j} \left| \frac{-y_i}{h} H' \left(\frac{x_i + \mathbf{z}_i^\top \beta^*}{h} \right) x_{i,j} \right| = O(1/h),$$

and

$$\sup_{\beta \in \Theta} \sup_{i,j} \mathbb{E} \left| \frac{-y_i}{h} H' \left(\frac{x_i + \mathbf{z}_i^\top \beta^*}{h} \right) x_{i,j} \right|^2 = O(1/h).$$

By Bernstein's inequality,

$$\sup_j |(1 - \mathbb{E}) U_{n,h,j}(\beta^*)| = O_{\mathbb{P}} \left(\sqrt{\frac{\log p}{nh}} \right).$$

Proof of (69):

Recall equation (38) and $\pi_U = \int_{-1}^1 x^\alpha H'(x) dx \neq 0$. For any $\mathbf{v} \in \mathbb{R}^p$ that satisfies $\|\mathbf{v}\|_2 = 1$

and $\beta \in \Theta$, we have

$$\begin{aligned}
& \mathbb{E}_{\cdot|\mathbf{Z}} \mathbf{v}^\top \Psi_n(\beta) \\
&= (\mathbf{v}^\top \mathbf{Z}) \cdot \mathbb{E}_{\cdot|\mathbf{Z}} \left[\frac{\mathbf{Z}^\top \Delta(\beta)}{h^2} [2\mathbb{I}(\zeta + \epsilon < 0) - 1] H'' \left(\frac{X + \mathbf{Z}^\top \beta}{h} \right) - \frac{1}{h} [2\mathbb{I}(\zeta + \epsilon < 0) - 1] H' \left(\frac{X + \mathbf{Z}^\top \beta}{h} \right) \right] \\
&= (\mathbf{v}^\top \mathbf{Z}) \int_{-1}^1 \left[2F \left(\mathbf{Z}^\top \Delta(\beta) - \xi h \mid \mathbf{Z} \right) - 1 \right] \rho \left(\xi h - \mathbf{Z}^\top \Delta(\beta) \mid \mathbf{Z} \right) \\
&\quad \cdot \left(\frac{\mathbf{Z}^\top \Delta(\beta)}{h} H''(\xi) - H'(\xi) \right) d\xi \\
&= (\mathbf{v}^\top \mathbf{Z}) \sum_{k=1}^{2\alpha+1} M_k(\mathbf{Z}) \int_{-1}^1 \left(\xi h - \mathbf{Z}^\top \Delta(\beta) \right)^k \left(\frac{\mathbf{Z}^\top \Delta(\beta)}{h} H''(\xi) - H'(\xi) \right) d\xi
\end{aligned} \tag{81}$$

For $1 \leq k \leq \alpha - 1$,

$$\begin{aligned}
& \sup_{\beta \in \Theta} \left| \int_{-1}^1 \left(\xi h - \mathbf{Z}^\top \Delta(\beta) \right)^k \left(\frac{\mathbf{Z}^\top \Delta(\beta)}{h} H''(\xi) - H'(\xi) \right) d\xi \right| \\
&= \sup_{\beta \in \Theta} \left| \sum_{k'=0}^k \binom{k}{k'} h^{k'} \left(-\mathbf{Z}^\top \Delta(\beta) \right)^{k-k'} \left[\left(\mathbf{Z}^\top \Delta(\beta) / h \right) \int_{-1}^1 \xi^{k'} H''(\xi) d\xi - \int_{-1}^1 \xi^{k'} H'(\xi) d\xi \right] \right| \\
&= (k-1) \left| -\mathbf{Z}^\top \Delta(\beta) \right|^k = O(s\delta_{m,0}^2).
\end{aligned}$$

For $k = \alpha$,

$$\begin{aligned}
& \sup_{\beta \in \Theta} \left| \int_{-1}^1 \left(\xi h - \mathbf{Z}^\top \Delta(\beta) \right)^\alpha \left(\frac{\mathbf{Z}^\top \Delta(\beta)}{h} H''(\xi) - H'(\xi) \right) d\xi \right| \\
&= \sup_{\beta \in \Theta} \left| \sum_{k=0}^{\alpha} \binom{\alpha}{k} h^k \left(-\mathbf{Z}^\top \Delta(\beta) \right)^{\alpha-k} \left[\left(\mathbf{Z}^\top \Delta(\beta) / h \right) \int_{-1}^1 \xi^k H''(\xi) d\xi - \int_{-1}^1 \xi^k H'(\xi) d\xi \right] \right| \\
&\leq (\alpha-1) \left| \mathbf{Z}^\top \Delta(\beta) \right|^\alpha + |\pi_U h^\alpha| = O[h^\alpha + (\sqrt{s}\delta_{m,0})^\alpha].
\end{aligned}$$

For $\alpha + 1 \leq k \leq 2\alpha + 1$,

$$\sup_{\beta \in \Theta} \left| \int_{-1}^1 \left(\xi h - \mathbf{Z}^\top \Delta(\beta) \right)^k \left(\frac{\mathbf{Z}^\top \Delta(\beta)}{h} H''(\xi) - H'(\xi) \right) d\xi \right| = O[h^{\alpha+1} + (\sqrt{s}\delta_{m,0})^{\alpha+1}].$$

Therefore, $\mathbb{E}[\mathbf{v}^\top \Psi_n(\beta)] \lesssim \mathbb{E}(\mathbf{v}^\top \mathbf{Z}) (s\delta_{m,0}^2 + h^\alpha) \lesssim s\delta_{m,0}^2 + h^\alpha$, which completes the proof of (69).

□

Proof of Theorem B.1

Now we are ready to prove the 1-step error for $\widehat{\beta}^{(1)}$.

Proof. For simplicity, we replace $V_{m,1}(\hat{\beta}^{(0)})$, $V_n(\hat{\beta}^{(0)})$, $U_n(\hat{\beta}^{(0)})$, and $\lambda_n^{(1)}$ by $V_{m,1}$, V_n , U_n , and λ_n , respectively. Then, by Algorithm 1,

$$\hat{\beta}^{(1)} = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \left\{ \|\beta\|_1 : \left\| V_{m,1}\beta - (V_{m,1}\hat{\beta}^{(0)} - U_n) \right\|_\infty \leq \lambda_n \right\}.$$

Let $\delta := \hat{\beta}^{(1)} - \beta^*$.

Using Lemma C.3, with probability tending to 1, we have

$$\begin{aligned} \left\| V_{m,1}\beta^* - (V_{m,1}\hat{\beta}^{(0)} - U_n) \right\|_\infty &\leq \left\| U_n - V_n(\hat{\beta}^{(0)} - \beta^*) \right\|_\infty + \left\| (V_{m,1} - V_n)(\hat{\beta}^{(0)} - \beta^*) \right\|_\infty \\ &\leq \left\| U_n - V_n(\hat{\beta}^{(0)} - \beta^*) \right\|_\infty + \|V_{m,1} - V_n\|_{\max} \left\| \hat{\beta}^{(0)} - \beta^* \right\|_1 \\ &= O \left(\sqrt{\frac{\log p}{nh}} + \sqrt{\frac{s \log p}{mh^3}} \delta_{m,0} + s\delta_{m,0}^2 + h^\alpha \right) \\ &\leq \lambda_n, \end{aligned} \quad (82)$$

when C_λ is large enough. By the definition,

$$\left\| V_{m,1}\hat{\beta}^{(1)} - (V_{m,1}\hat{\beta}^{(0)} - U_n) \right\|_\infty \leq \lambda_n. \quad (83)$$

Hence, with probability tending to 1,

$$\left\| V_{m,1}(\hat{\beta}^{(1)} - \beta^*) \right\|_\infty \leq 2\lambda_n. \quad (84)$$

From (82) and the optimality of $\hat{\beta}^{(1)}$, we have $\left\| \hat{\beta}^{(1)} \right\|_1 \leq \|\beta^*\|_1$, which implies

$$\left\| \hat{\beta}_S^{(1)} \right\|_1 + \left\| \hat{\beta}_{S^c}^{(1)} \right\|_1 = \left\| \hat{\beta}^{(1)} \right\|_1 \leq \|\beta^*\|_1 = \|\beta_S^*\|_1.$$

Therefore,

$$\left\| (\beta^* - \hat{\beta}^{(1)})_{S^c} \right\|_1 = \left\| \hat{\beta}_{S^c}^{(1)} \right\|_1 \leq \left\| \beta_S^* - \hat{\beta}_S^{(1)} \right\|_1 = \left\| (\beta^* - \hat{\beta}^{(1)})_S \right\|_1.$$

Hence,

$$\left\| \beta^* - \hat{\beta}^{(1)} \right\|_1 \leq 2 \left\| (\beta^* - \hat{\beta}^{(1)})_S \right\|_1 \leq 2\sqrt{s} \left\| (\beta^* - \hat{\beta}^{(1)})_S \right\|_2 \leq 2\sqrt{s} \left\| \beta^* - \hat{\beta}^{(1)} \right\|_2,$$

Let $\delta := \hat{\beta}^{(1)} - \beta^*$. So far we have shown that, with probability tending to one, $\|\delta\|_1 \leq 2\sqrt{s} \|\delta\|_2$ and $\|V_{m,1}\delta\|_\infty \leq 2\lambda_{m,0}$. Therefore,

$$\begin{aligned} \delta^\top V_{m,1}\delta &= \delta^\top V\delta + \delta^\top (V_{m,1} - \mathbb{E}[V_{m,1}])\delta + \delta^\top (\mathbb{E}[V_{m,1}] - V)\delta \\ &\geq \Lambda_{\min}(V) \|\delta\|_2^2 - \|(1 - \mathbb{E})V_{m,1}\|_{\max} \|\delta\|_1^2 + \delta^\top (\mathbb{E}[V_{m,1}] - V)\delta \\ &\geq \Lambda_{\min}(V) \|\delta\|_2^2 - s \|(1 - \mathbb{E})V_{m,1}\|_{\max} \|\delta\|_2^2 - \left| \delta^\top (\mathbb{E}[V_{m,1}] - V)\delta \right|. \end{aligned} \quad (85)$$

By (71) and (76),

$$s \|(1 - \mathbb{E})V_{m,1}\|_{\max} = O_{\mathbb{P}} \left(\sqrt{\frac{s^2 \log p}{mh^3}} \right) = o_{\mathbb{P}}(1).$$

By (67),

$$\left| \boldsymbol{\delta}^\top (\mathbb{E}[V_{m,1}] - V) \boldsymbol{\delta} \right| \lesssim (\sqrt{s} \delta_{m,0} + h^\alpha) \|\boldsymbol{\delta}\|_2^2 = o(\|\boldsymbol{\delta}\|_2^2).$$

Therefore, (85) leads to

$$\boldsymbol{\delta}^\top V_{m,1} \boldsymbol{\delta} \geq \Lambda_{\min}(V) \|\boldsymbol{\delta}\|_2^2 - o_{\mathbb{P}}(\|\boldsymbol{\delta}\|_2^2) \geq (\Lambda_{\min}(V)/2) \|\boldsymbol{\delta}\|_2^2,$$

with probability tending to one. On the other hand,

$$\boldsymbol{\delta}^\top V_{m,1} \boldsymbol{\delta} \leq \|\boldsymbol{\delta}\|_1 \|V_{m,1} \boldsymbol{\delta}\|_\infty = O_{\mathbb{P}}(\sqrt{s} \lambda_n \|\boldsymbol{\delta}\|_2),$$

Combining the two inequalities above, we finally get $\|\boldsymbol{\delta}\|_2 = O_{\mathbb{P}}(\sqrt{s} \lambda_n)$, which completes the proof. \square

Proof of Theorem B.2

Proof. First note that $\sqrt{\frac{\log p}{nh}} + h^\alpha \asymp \left(\frac{\log p}{n}\right)^{\frac{\alpha}{2\alpha+1}}$, when $h = h^* = \left(\frac{\log p}{n}\right)^{\frac{1}{2\alpha+1}}$. Since $s = O(m^r)$ and $n = O(m^{1/c})$ for some $0 < c < 1$ and $0 < r < \frac{1}{4}$, our assumption $\alpha > \alpha_0$ guarantees $\frac{s^2 \log p}{m(h^*)^3} = o(1)$, $s^{3/2} \delta_{m,0} = o(1)$, $s \delta_{m,0} = O((h^*)^{3/2})$, and $s(h^*)^\alpha = o(1)$. Adding the requirement of $\delta_{m,0}$, the assumptions in Theorem B.1 hold, which proves Theorem B.2 when $t = 1$.

We will show Theorem B.2 by induction. Assume Theorem B.2 is true for t . Our assumption $\alpha > \alpha_0$ also ensures that $\sqrt{s} \left(\frac{\log p}{n}\right)^{\frac{\alpha}{2\alpha+1}} < \delta_{m,0}$, which implies $\delta_{m,t} < \delta_{m,0}$, and thus $s^{3/2} \delta_{m,t} = o(1)$, $s \delta_{m,t} = O(h^{3/2})$. Then by Theorem B.1, by taking

$$\lambda_n^{(t+1)} = C_\lambda \left[\left(\frac{\log p}{n}\right)^{\frac{\alpha}{2\alpha+1}} + \left(\sqrt{\frac{s \log p}{mh^3}} + s \delta_{m,t} \right) \delta_{m,t} \right],$$

we have

$$\begin{aligned} \left\| \widehat{\boldsymbol{\beta}}^{(t+1)} - \boldsymbol{\beta}^* \right\|_2 &= O_{\mathbb{P}} \left[\sqrt{s} \left(\frac{\log p}{n}\right)^{\frac{\alpha}{2\alpha+1}} + \left(\sqrt{\frac{s^2 \log p}{mh^3}} + s^{3/2} \delta_{m,t} \right) \delta_{m,t} \right] \\ &= O_{\mathbb{P}} \left[\sqrt{s} \left(\frac{\log p}{n}\right)^{\frac{\alpha}{2\alpha+1}} + \left(\sqrt{\frac{s^2 \log p}{mh^3}} + s^{3/2} \delta_{m,0} \right)^{t+1} \delta_{m,0} \right], \end{aligned}$$

and

$$\left\| \widehat{\boldsymbol{\beta}}^{(t+1)} - \boldsymbol{\beta}^* \right\|_1 \leq 2\sqrt{s} \left\| \widehat{\boldsymbol{\beta}}^{(t+1)} - \boldsymbol{\beta}^* \right\|_2,$$

with probability tending to 1. This completes the proof. \square

D Discussions on the Super-Efficiency Phenomenon

In this section, we show that our estimator $\widehat{\beta}^{(T)}$ achieves the same asymptotic performance over a class of underlying distributions under certain uniform assumptions. In model (1), for any β^* , denote the density function of $\zeta := X + \mathbf{Z}^\top \beta^*$ conditional on \mathbf{Z} by $\rho(\cdot | \mathbf{Z})$ and the distribution function of ϵ conditional on \mathbf{Z} by $F(\cdot | \mathbf{Z})$. We define the distribution class Θ to be the set of tuples (β^*, ρ, F) that satisfy the following assumptions:

Assumption 7. Assume that there exists a neighborhood of 0 such that, for all $(\beta^*, \rho, F) \in \Theta$ and all integers $1 \leq k \leq \alpha$, the k -th order derivative of $\rho(\cdot | \mathbf{Z})$ exists in this neighborhood for almost every \mathbf{Z} . Furthermore, there exists a constant $C_{\Theta,1} > 0$ such that $\sup_{\zeta, \mathbf{Z}, k} |\rho^{(k)}(\zeta | \mathbf{Z})| < C_{\Theta,1}$.

Assumption 8. Assume that ϵ and X are independent given \mathbf{Z} , and there exists a neighborhood of 0 such that, for all $(\beta^*, \rho, F) \in \Theta$ and all integers $1 \leq k \leq \alpha + 1$, the k -th order derivative of $F(\cdot | \mathbf{Z})$ exists in this neighborhood for almost every \mathbf{Z} . Furthermore, there exists a constant $C_{\Theta,2} > 0$ such that $\sup_{\epsilon, \mathbf{Z}, k} |F^{(k)}(\epsilon | \mathbf{Z})| < C_{\Theta,2}$.

Assumption 9. Assume that there exists a constant $c_\Theta > 0$ such that, for all $(\beta^*, \rho, F) \in \Theta$, the matrix $V = 2\mathbb{E}[\rho(0 | \mathbf{Z}) F'(0 | \mathbf{Z}) \mathbf{Z} \mathbf{Z}^\top]$ satisfies $c_\Theta^{-1} < \Lambda_{\min}(V) < \Lambda_{\max}(V) < c_\Theta$, where $\Lambda_{\min}(\Lambda_{\max})$ denotes the minimum (maximum) eigenvalue of V .

Assumptions 7–9 for the distribution class Θ are parallel to Assumptions 2–4 for a fixed distribution. Assumptions 7–8 require α -order smoothness on a fixed neighborhood for all ρ and F , which ensures that the Taylor’s expansion in the technical proof always hold when n is sufficiently large. Furthermore, the constants $C_{\Theta,1}$ and $C_{\Theta,2}$ provide uniform upper bounds for the derivatives of ρ and F over Θ . Similarly, Assumption 9 ensures that the population Hessian matrix is always positive semi-definite with eigenvalues uniformly bounded away from 0 and ∞ . Under these assumptions, replicating the analysis of (mSMSE) in Section 3 leads to the following result:

Theorem D.1. Assume Assumptions 1, 5, 7, 8 and 9 hold, and $\sup_{(\beta^*, \rho, F) \in \Theta} \|\widehat{\beta}^{(0)} - \beta^*\|_2 = O_{\mathbb{P}}(m^{-1/3})$. By choosing $h_t = \max \left\{ (\lambda_h/n)^{\frac{1}{2\alpha+1}}, m^{-\frac{2t}{3\alpha}} \right\}$ at iteration $t = 1, 2, \dots, T$, when T satis-

Table 1: The coverage rates (nominal 95%) of (mSMSE) in the first four iterations with different values of λ_h and $\log_m(n)$. The dimension $p = 1$ and the noise is homoscedastic normal.

λ_h	$\log_m(n)$	$t = 1$	$t = 2$	$t = 3$	$t = 4$	λ_h	$\log_m(n)$	$t = 1$	$t = 2$	$t = 3$	$t = 4$
1	1.35	0.90	0.95	0.96	0.96	10	1.35	0.92	0.93	0.93	0.93
	1.55	0.88	0.96	0.96	0.96		1.55	0.91	0.95	0.95	0.95
	1.75	0.83	0.95	0.95	0.95		1.75	0.89	0.93	0.93	0.93
20	1.35	0.92	0.93	0.93	0.93	$\widehat{\lambda}_h^*$	1.35	0.91	0.94	0.94	0.94
	1.55	0.90	0.96	0.96	0.96		1.55	0.93	0.95	0.95	0.95
	1.75	0.89	0.94	0.94	0.94		1.75	0.84	0.92	0.94	0.94

ties (13), we have: $\forall \varepsilon > 0, \exists M_\varepsilon, N_\varepsilon$, such that $\forall n \geq N_\varepsilon$, it holds that

$$\sup_{(\beta^*, \rho, F) \in \Theta} \mathbb{P} \left(\|\widehat{\beta}^{(T)} - \beta^*\|_2 > M_\varepsilon n^{-\frac{\alpha}{2\alpha+1}} \right) < \varepsilon. \quad (86)$$

Note that (86) is equivalent to $\|\widehat{\beta}^{(T)} - \beta^*\|_2 = O_{\mathbb{P}}(n^{-\frac{\alpha}{2\alpha+1}})$ if Θ only contains a single distribution. The proof of Theorem D.1 is almost the same as the proof of Proposition C.1 and Theorem 3.3, by noting that, for all distributions in Θ , the Taylor's expansion in (38) and the computation related to the bias in (39) are always correct. Moreover, the constant in the big O notation in (41) is uniform for all distributions in Θ , which is guaranteed by Assumptions 7 and 8.

E Additional Results in Simulations

E.1 Sensitivity Analysis

In this section, we use numerical experiments to show the sensitivity of the constant λ_h in the bandwidth h_t in Theorem 3.3. An expression of the optimal value of λ_h^* is obtained in (15) to minimize the asymptotic mean squared error. We estimate λ_h^* using \widehat{U} , \widehat{V} , and \widehat{V}_s . Under our experiment settings, the estimated constant $\widehat{\lambda}_h^*$ ranges from 1 to 17 in practice. To study the effect of λ_h on the validity of inference, we choose a wider range for λ_h , from 1 to 20, and report the coverage rates of the first four iterations of (mSMSE) in Table 1, with different λ_h and $\log_m(n)$. In the first iteration, it seems that larger λ_h 's lead to higher coverage rates, which may indicate that a larger bandwidth improves the initial estimator more aggressively at the beginning of the

algorithm. However, after the (mSMSE) converges in two or three iterations, all estimators achieve near-nominal coverage rates, no matter what value λ_h is. Therefore, (mSMSE) generally allows arbitrary choices of λ_h in a wide range, which suggests that our proposed (mSMSE) algorithm is robust with respect to λ_h .

Table 2: The cpu times (in seconds) that different methods take to compute the estimator, with $p = 10$, $\log_m(n)$ from 1.35 to 1.75 and two types of noise.

Noise Type	$\log_m(n)$	(mSMSE) $t = 2$	(mSMSE) $t = 3$	(Avg-SMSE)	pooled-SMSE
Homoscedastic Normal	1.35	0.364	0.485	0.307	0.559
	1.45	0.365	0.492	0.328	1.356
	1.55	0.426	0.605	0.370	3.260
	1.65	0.443	0.627	0.399	9.267
	1.75	0.474	0.665	0.459	23.208
Homoscedastic Uniform	1.35	0.343	0.474	0.310	0.728
	1.45	0.355	0.492	0.330	1.749
	1.55	0.418	0.598	0.370	4.202
	1.65	0.431	0.616	0.400	11.893
	1.75	0.466	0.659	0.464	29.909
Heteroscedastic Normal	1.35	0.317	0.438	0.316	0.555
	1.45	0.379	0.503	0.334	1.322
	1.55	0.433	0.608	0.357	3.228
	1.65	0.452	0.634	0.396	9.132
	1.75	0.483	0.672	0.458	22.369

E.2 Time Complexity

In this section, we compare the computational complexity of each method. The average CPU times that each method takes when $p = 10$ are reported in Table 2. The computation time is recorded in a simulated distributed environment on a RedHat Enterprise Linux cluster containing 524 Lenovo SD650 nodes interconnected by high-speed networks. On each computer node, two Intel Xeon Platinum 8268 24C 205W 2.9GHz Processors are equipped with 48 processing cores.

In Table 2, we first notice that the speed of (mSMSE) is much faster than the pooled estimator and the discrepancy greatly increases when n gets larger. Second, the computation time of (mSMSE) is comparable to (Avg-SMSE). This result may seem counterintuitive since (mSMSE) still requires

running an SMSE on the first machine for the initial estimator. However, since the computation time of (Avg-SMSE) is mainly determined by the maximum computation time of the L local machines, (Avg-SMSE) greatly suffers from the computational performance of the “worst” machine, especially when the number of machines is large. On the other hand, (mSMSE) only runs SMSE on one machine and therefore achieves comparable computation time in the experiments.

E.3 Results for Other Types of Noise

Table 3: The bias, variance and coverage rates of (mSMSE) ($t = 1, 2, 3$), (Avg-MSE), (Avg-SMSE) and pooled-SMSE, with $p = 1$, $\log_m(n)$ from 1.35 to 1.75 and homoscedastic uniform noise.

$\log_m(n)$	Bias ($\times 10^{-2}$)	Variance ($\times 10^{-4}$)	Coverage Rate	Bias ($\times 10^{-2}$)	Variance ($\times 10^{-4}$)	Coverage Rate	Bias ($\times 10^{-2}$)	Variance ($\times 10^{-4}$)	Coverage Rate
	(mSMSE) $t = 1$			(mSMSE) $t = 3$			(Avg-SMSE)		
1.35	0.02	0.70	0.88	0.31	0.60	0.91	0.04	0.44	0.92
1.45	-0.14	0.48	0.88	0.17	0.32	0.94	-0.30	0.21	0.92
1.55	-0.16	0.26	0.84	0.11	0.12	0.92	-0.76	0.09	0.62
1.65	-0.18	0.16	0.80	0.09	0.04	0.94	-1.29	0.04	0.03
1.75	-0.12	0.08	0.78	0.08	0.03	0.94	-1.85	0.01	0.00
	(mSMSE) $t = 2$			(Avg-MSE)			pooled-SMSE		
1.35	0.32	0.61	0.9	-0.33	0.84	0.90	0.30	0.60	0.90
1.45	0.17	0.34	0.93	-0.38	0.22	0.92	0.15	0.32	0.94
1.55	0.12	0.14	0.93	-0.42	0.11	0.86	0.10	0.12	0.93
1.65	0.09	0.05	0.94	-0.36	0.07	0.80	0.07	0.04	0.96
1.75	0.08	0.03	0.94	-0.37	0.02	0.65	0.06	0.03	0.94

In this section, we report the bias, variance and coverage rates in Tables 3–6 for the other two noise types, i.e., the homoscedastic uniform and heteroscedastic normal noise, with $p = 1$ and 10. From these tables, we can still see the failure of inference of (Avg-MSE) and (Avg-SMSE) when $\log_m(n)$ is large, while the (mSMSE) method with $t \geq 3$ achieves near-nominal coverage rates no matter how large $\log_m(n)$ is. These findings are all consistent with the results for the homoscedastic normal noise.

Table 4: The bias, variance and coverage rates of (mSMSE) ($t = 1, 2, 3, 4$), (Avg-SMSE) and pooled-SMSE, with $p = 10$, $\log_m(n)$ from 1.35 to 1.75 and homoscedastic uniform noise.

$\log_m(n)$	Bias ($\times 10^{-2}$)	Variance ($\times 10^{-4}$)	Coverage Rate	Bias ($\times 10^{-2}$)	Variance ($\times 10^{-4}$)	Coverage Rate	Bias ($\times 10^{-2}$)	Variance ($\times 10^{-4}$)	Coverage Rate
	(mSMSE) $t = 1$			(mSMSE) $t = 3$			(Avg-SMSE)		
1.35	-1.84	20.99	0.73	0.64	8.05	0.85	-2.30	2.63	0.91
1.45	-2.14	12.39	0.62	0.29	3.40	0.88	-4.44	1.57	0.55
1.55	-1.92	6.83	0.59	0.20	1.40	0.93	-6.13	0.74	0.02
1.65	-1.81	3.93	0.49	0.10	0.59	0.92	-7.18	0.33	0.00
1.75	-1.87	2.27	0.37	0.11	0.26	0.93	-7.97	0.10	0.00
	(mSMSE) $t = 2$			(mSMSE) $t = 4$			pooled-SMSE		
1.35	-0.12	11.71	0.78	0.55	6.27	0.84	0.62	4.76	0.93
1.45	0.52	5.06	0.86	0.62	2.65	0.89	0.29	2.63	0.93
1.55	0.12	2.03	0.85	0.27	1.26	0.94	0.08	1.24	0.94
1.65	0.06	0.95	0.86	0.17	0.58	0.92	0.07	0.56	0.92
1.75	0.11	0.31	0.89	0.11	0.26	0.94	0.05	0.26	0.93

Table 5: The bias, variance and coverage rates of (mSMSE) ($t = 1, 2, 3$), (Avg-MSE), (Avg-SMSE) and pooled-SMSE, with $p = 1$, $\log_m(n)$ from 1.35 to 1.75 and heteroscedastic normal noise.

$\log_m(n)$	Bias ($\times 10^{-2}$)	Variance ($\times 10^{-4}$)	Coverage Rate	Bias ($\times 10^{-2}$)	Variance ($\times 10^{-4}$)	Coverage Rate	Bias ($\times 10^{-2}$)	Variance ($\times 10^{-4}$)	Coverage Rate
	(mSMSE) $t = 1$			(mSMSE) $t = 3$			(Avg-SMSE)		
1.35	0.06	0.30	0.92	0.15	0.27	0.96	0.10	0.24	0.94
1.45	0.00	0.17	0.90	0.11	0.16	0.94	-0.03	0.11	0.94
1.55	-0.06	0.09	0.88	0.05	0.06	0.94	-0.29	0.07	0.86
1.65	-0.04	0.04	0.90	0.04	0.03	0.98	-0.60	0.03	0.30
1.75	-0.05	0.02	0.88	0.03	0.01	0.96	-0.96	0.01	0.00
	(mSMSE) $t = 2$			(Avg-MSE)			pooled-SMSE		
1.35	0.14	0.26	0.96	-0.40	0.51	0.89	0.14	0.26	0.96
1.45	0.11	0.16	0.94	-0.36	0.22	0.88	0.10	0.16	0.94
1.55	0.05	0.06	0.94	-0.40	0.12	0.82	0.04	0.06	0.94
1.65	0.04	0.03	0.98	-0.42	0.03	0.69	0.03	0.03	0.97
1.75	0.04	0.01	0.96	-0.41	0.01	0.40	0.02	0.01	0.96

Table 6: The bias, variance and coverage rates of (mSMSE) ($t = 1, 2, 3, 4$), (Avg-SMSE) and pooled-SMSE, with $p = 10$, $\log_m(n)$ from 1.35 to 1.75 and heteroscedastic normal noise.

$\log_m(n)$	Bias ($\times 10^{-2}$)	Variance ($\times 10^{-4}$)	Coverage Rate	Bias ($\times 10^{-2}$)	Variance ($\times 10^{-4}$)	Coverage Rate	Bias ($\times 10^{-2}$)	Variance ($\times 10^{-4}$)	Coverage Rate
	(mSMSE) $t = 1$			(mSMSE) $t = 3$			(Avg-SMSE)		
1.35	-0.67	19.36	0.82	0.74	4.86	0.97	-0.24	4.08	0.96
1.45	-1.25	12.14	0.72	0.36	2.20	0.96	-1.46	1.88	0.84
1.55	-1.39	11.72	0.63	0.18	1.22	0.94	-2.47	0.93	0.34
1.65	-1.29	6.05	0.56	0.17	0.64	0.95	-3.28	0.38	0.00
1.75	-1.18	2.91	0.49	0.14	0.35	0.92	-3.84	0.16	0.00
	(mSMSE) $t = 2$			(mSMSE) $t = 4$			pooled-SMSE		
1.35	0.78	6.04	0.96	0.79	4.61	0.97	0.69	4.59	0.97
1.45	0.31	2.75	0.93	0.36	2.17	0.96	0.27	2.15	0.95
1.55	0.18	1.63	0.91	0.19	1.20	0.94	0.10	1.20	0.94
1.65	0.13	0.74	0.92	0.17	0.63	0.95	0.09	0.63	0.95
1.75	0.12	0.39	0.90	0.14	0.35	0.93	0.05	0.35	0.92

References

- Bousquet, O. (2003). Concentration inequalities for sub-additive functions using the entropy method. In *Stochastic Inequalities and Applications*, pp. 213–247. Springer.
- Cai, T. T., C.-H. Zhang, and H. H. Zhou (2010). Optimal rates of convergence for covariance matrix estimation. *Annals of Statistics* 38(4), 2118–2144.
- Candes, E. and T. Tao (2007). The Dantzig selector: Statistical estimation when p is much larger than n . *Annals of Statistics* 35(6), 2313–2351.
- Feng, H., Y. Ning, and J. Zhao (2022). Nonregular and minimax estimation of individualized thresholds in high dimension with binary responses. *Annals of Statistics*, To appear.
- Horowitz, J. L. (1992). A smoothed maximum score estimator for the binary response model. *Econometrica* 60(3), 505–531.
- Luo, J., Q. Sun, and W.-X. Zhou (2022). Distributed adaptive Huber regression. *Computational Statistics & Data Analysis* 169, 107419.
- Mukherjee, D., M. Banerjee, and Y. Ritov (2019). Non-standard asymptotics in high dimensions: Manski’s maximum score estimator revisited. *arXiv preprint arXiv:1903.10063*.