

LEVERAGE CLASSIFIER: ANOTHER LOOK AT SUPPORT VECTOR MACHINE

Yixin Han¹, Jun Yu², Nan Zhang³, Cheng Meng⁴, Ping Ma⁵, Wenxuan Zhong⁵, and Changliang Zou¹

¹*School of Statistics and Data Science, LPMC & KLMDASR, Nankai University, Tianjin, P.R. China*

²*School of Mathematics and Statistics, Beijing Institute of Technology, Beijing, P.R.China*

³*School of Data Science, Fudan University, Shanghai, P.R.China*

⁴*Institute of Statistics and Big Data, Renmin University, Beijing, P.R.China*

⁵*Department of Statistics, University of Georgia, Athens, GA, USA*

Supplementary Material

This supplementary material contains the proofs of technical results and some additional simulation results.

Appendix A: Useful Lemma

The following Lemma is a multivariate extension of the martingale central limit theorem, see Lemma 4 in Zhang et al. (2021) for details.

Lemma S.1 (Multivariate version of martingale CLT). *Let $\{\boldsymbol{\eta}_{ki}, i = 1, \dots, N_k\}$ be a martingale difference sequence in \mathbb{R}^p relative to the filtration $\{\mathcal{F}_{ki}, i = 0, 1, \dots, N_k\}$ and let $\mathbf{Z}_k \in \mathbb{R}^p$ be an \mathcal{F}_{k0} -measurable random vector for $k = 1, 2, 3, \dots$. Denote $\mathbf{R}_k = \sum_{i=1}^{N_k} \boldsymbol{\eta}_{ki}$. Assume the following conditions hold.*

$$(i) \lim_{k \rightarrow \infty} \sum_{i=1}^{N_k} \mathbb{E}(\|\boldsymbol{\eta}_{ki}\|^4) = 0.$$

(ii) $\lim_{k \rightarrow \infty} \mathbb{E} \left\{ \left\| \sum_{i=1}^{N_k} \mathbb{E}(\boldsymbol{\eta}_{ki} \boldsymbol{\eta}_{ki}^\top \mid \mathcal{F}_{k,i-1}) - \mathbf{B}_k \right\|^2 \right\} = 0$ for some sequence of positive-definite matrices $\{\mathbf{B}_k\}_{k=1}^\infty$ with $\sup_k \lambda_{\max}(\mathbf{B}_k) < \infty$, say that the largest eigenvalue is uniformly bounded.

(iii) For a probability distribution \mathbf{L}_0 , $*$ denotes convolution and $\mathbf{L}(\cdot)$ denotes the law of random variables, $\mathbf{L}(\mathbf{Z}_k) * \mathcal{N}(\mathbf{0}, \mathbf{B}_k) \rightarrow \mathbf{L}_0$, where the convergence is in distribution.

Then we have

$$\mathbf{L}(\mathbf{Z}_k + \mathbf{R}_k) \rightarrow \mathbf{L}_0.$$

Appendix B: Proof of Theorem 1

Proof. Denote

$$L_n(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{N\pi_i^*} [1 - Y_i^* f(\mathbf{X}_i^*, \boldsymbol{\beta})]_+, \quad L_N(\boldsymbol{\beta}) = \frac{1}{N} \sum_{j=1}^N [1 - Y_j f(\mathbf{X}_j, \boldsymbol{\beta})]_+,$$

$$l_{\lambda,n}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{N\pi_i^*} [1 - Y_i^* f(\mathbf{X}_i^*, \boldsymbol{\beta})]_+ + \frac{\lambda}{2} \|\boldsymbol{\beta}_1\|^2.$$

The proof can be divided into the following intermediate parts.

First, we consider the influence of a fixed λ . For a fixed $\boldsymbol{\theta} = (1, \boldsymbol{\theta}_1^\top)^\top \in \mathbb{R}^{p+1}$,

define

$$\Lambda_n(\boldsymbol{\theta}) = n \left\{ l_{\lambda,n} \left(\boldsymbol{\beta}^\dagger + \frac{\boldsymbol{\theta}}{\sqrt{n}} \right) - l_{\lambda,n}(\boldsymbol{\beta}^\dagger) \right\}, \quad T_n(\boldsymbol{\theta}) = \mathbb{E} \{ \Lambda_n(\boldsymbol{\theta}) \}.$$

Observe that

$$\begin{aligned} \Lambda_n(\boldsymbol{\theta}) &= \sum_{i=1}^n \frac{1}{N\pi_i^*} \left\{ \left[1 - Y_i^* f \left(\mathbf{X}_i^*, \boldsymbol{\beta}^\dagger + \frac{\boldsymbol{\theta}}{\sqrt{n}} \right) \right]_+ - [1 - Y_i^* f(\mathbf{X}_i^*, \boldsymbol{\beta}^\dagger)]_+ \right\} \\ &\quad + n \frac{\lambda}{2} \left(\left\| \boldsymbol{\beta}_1^\dagger + \frac{\boldsymbol{\theta}_1}{\sqrt{n}} \right\|^2 - \|\boldsymbol{\beta}_1^\dagger\|^2 \right), \end{aligned}$$

and $\mathbb{E} \{ L_n(\boldsymbol{\beta}) \} = \mathbb{E} [\mathbb{E} \{ L_n(\boldsymbol{\beta}) \mid \mathcal{D}_N \}] = L(\boldsymbol{\beta}) = \mathbb{E} [1 - Y f(\mathbf{X}, \boldsymbol{\beta})]_+$. Under Assumption 3, we assume $\boldsymbol{\beta}_1^\dagger \neq 0$ without loss of generality. By Lemma 3 in Koo et al. (2008),

we have

$$\begin{aligned} T_n(\boldsymbol{\theta}) &= n \left\{ L\left(\boldsymbol{\beta}^\dagger + \frac{\boldsymbol{\theta}}{\sqrt{n}}\right) - L(\boldsymbol{\beta}^\dagger) \right\} + \frac{\lambda}{2} \left(\|\boldsymbol{\theta}_1\|^2 + 2\sqrt{n}\boldsymbol{\theta}_1^\top \boldsymbol{\beta}_1^\dagger \right), \\ &= \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{H}(\check{\boldsymbol{\beta}}) \boldsymbol{\theta} + \frac{\lambda}{2} \left(\|\boldsymbol{\theta}_1\|^2 + 2\sqrt{n}\boldsymbol{\theta}_1^\top \boldsymbol{\beta}_1^\dagger \right), \end{aligned}$$

by applying Taylor expansion of $L(\boldsymbol{\beta})$ around $\boldsymbol{\beta}^\dagger$, where $\check{\boldsymbol{\beta}} = \boldsymbol{\beta}^\dagger + (\boldsymbol{\theta}/\sqrt{n})t$ for some $0 < t < 1$.

Define $\mathbf{D}_{ij}(\boldsymbol{\alpha}) = \mathbf{H}(\boldsymbol{\beta}^\dagger + \boldsymbol{\alpha})_{ij} - \mathbf{H}(\boldsymbol{\beta}^\dagger)_{ij}$ for $0 \leq i, j \leq p+1$. By Assumption 1, $\mathbf{H}(\boldsymbol{\beta})$ is continuous in $\boldsymbol{\beta}$. Then, for any $\varepsilon_1 > 0$, there exist $\delta_1 > 0$ such that $\mathbf{D}_{ij}(\boldsymbol{\alpha}) < \varepsilon_1$ if $\|\boldsymbol{\alpha}\| < \delta_1$ for all $0 \leq i, j \leq p+1$. Thus, for sufficiently large n such that $\|(\boldsymbol{\theta}/\sqrt{n})t\| < \delta_1$

$$\left| \boldsymbol{\theta}^\top \left(\mathbf{H}(\check{\boldsymbol{\beta}}) - \mathbf{H}(\boldsymbol{\beta}^\dagger) \right) \boldsymbol{\theta} \right| \leq \sum_{i,j} |\boldsymbol{\theta}_i| |\boldsymbol{\theta}_j| \left| \mathbf{D}_{ij} \left(\frac{\boldsymbol{\theta}}{\sqrt{n}} t \right) \right| \leq 2\varepsilon_1 \|\boldsymbol{\theta}\|^2,$$

then $\boldsymbol{\theta}^\top \mathbf{H}(\check{\boldsymbol{\beta}}) \boldsymbol{\theta} / 2 = \boldsymbol{\theta}^\top \mathbf{H}(\boldsymbol{\beta}^\dagger) \boldsymbol{\theta} / 2 + o(1)$ as $n \rightarrow \infty$. Combining the assumption that $\lambda = o(n^{-1/2})$, we have

$$T_n(\boldsymbol{\theta}) = \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{H}(\boldsymbol{\beta}^\dagger) \boldsymbol{\theta} + o(1).$$

Next, we would like to provide an expansion of $\Lambda_n(\boldsymbol{\theta})$ under Assumptions 1–3. Let $\mathbf{W}_n = -n^{-1} \sum_{i=1}^n (N\pi_i^*)^{-1} \xi_i^* Y_i^* \widetilde{\mathbf{X}}_i^*$, where $\xi_i^* = \mathbb{I}(Y_i^* f(\mathbf{X}_i^*, \boldsymbol{\beta}^\dagger) \leq 1)$. If we define

$$\begin{aligned} R_{i,n}(\boldsymbol{\theta}) &= \frac{1}{N\pi_i^*} \left\{ \left[1 - Y_i^* f\left(\mathbf{X}_i^*, \boldsymbol{\beta}^\dagger + \frac{\boldsymbol{\theta}}{\sqrt{n}} t\right) \right]_+ - [1 - Y_i^* f(\mathbf{X}_i^*, \boldsymbol{\beta}^\dagger)]_+ + \xi_i^* Y_i^* f\left(\mathbf{X}_i^*, \frac{\boldsymbol{\theta}}{\sqrt{n}}\right) \right\}, \\ R_{j,N}(\boldsymbol{\theta}) &= \left[1 - Y_j f\left(\mathbf{X}_j, \boldsymbol{\beta}^\dagger + \frac{\boldsymbol{\theta}}{\sqrt{n}} t\right) \right]_+ - [1 - Y_j f(\mathbf{X}_j, \boldsymbol{\beta}^\dagger)]_+ + \xi_j Y_j f\left(\mathbf{X}_j, \frac{\boldsymbol{\theta}}{\sqrt{n}}\right), \end{aligned}$$

where $i = 1, \dots, n$ and $j = 1, \dots, N$. Recall that $\mathbb{E}\{(N\pi_i^*)^{-1} \xi_i^* Y_i^* \widetilde{\mathbf{X}}_i^*\} = \mathbf{S}(\boldsymbol{\beta}^\dagger) = \mathbf{0}$.

Recall the definitions of $T_n(\boldsymbol{\theta})$ and \mathbf{W}_n , we have

$$\begin{aligned}
\Lambda_n(\boldsymbol{\theta}) &= \sum_{i=1}^n \frac{1}{N\pi_i^*} \left[1 - Y_i^* f \left(\mathbf{X}_i^*, \boldsymbol{\beta}^\dagger + \frac{\boldsymbol{\theta}}{\sqrt{n}} \right) \right]_+ - nL \left(\boldsymbol{\beta}^\dagger + \frac{\boldsymbol{\theta}}{\sqrt{n}} \right) \\
&\quad - \sum_{i=1}^n \frac{1}{N\pi_i^*} [1 - Y_i^* f(\mathbf{X}_i^*, \boldsymbol{\beta}^\dagger)]_+ + nL(\boldsymbol{\beta}^\dagger) + \frac{\lambda}{2} \left(\|\boldsymbol{\theta}_1\|^2 + 2\sqrt{n}\boldsymbol{\theta}_1^\top \boldsymbol{\beta}_1^\dagger \right) \\
&\quad + \sum_{i=1}^n \frac{1}{N\pi_i^*} \xi_i^* Y_i^* (\widetilde{\mathbf{X}}_i^*)^\top \frac{\boldsymbol{\theta}}{\sqrt{n}} - \sum_{i=1}^n \frac{1}{N\pi_i^*} \xi_i^* Y_i^* (\widetilde{\mathbf{X}}_i^*)^\top \frac{\boldsymbol{\theta}}{\sqrt{n}} \\
&= T_n(\boldsymbol{\theta}) + \sqrt{n} \mathbf{W}_n^\top \boldsymbol{\theta} + \sum_{i=1}^n [R_{i,n}(\boldsymbol{\theta}) - \mathbb{E}\{R_{i,n}(\boldsymbol{\theta})\}]. \tag{S.1}
\end{aligned}$$

Recall that $[\cdot]_+$ denotes the hinge loss. We define $\varphi = \mathbb{I}(a \leq 1)$ and $D = [1 - z]_+ - [1 - a]_+ + \varphi(z - a)$. Then we have

$$\begin{aligned}
D &= (1 - z)\mathbb{I}(a > 1, z \leq 1) + (z - 1)\mathbb{I}(a < 1, z > 1) \\
&\leq |z - a| \mathbb{I}(a > 1, z \leq 1) + |z - a| \mathbb{I}(a < 1, z > 1) \\
&= |z - a| \{ \mathbb{I}(a > 1, z \leq 1) + \mathbb{I}(a < 1, z > 1) \} \\
&\leq |z - a| \mathbb{I}(|1 - a| \leq |z - a|). \tag{S.2}
\end{aligned}$$

Let $z_i = Y_i^* f(\mathbf{X}_i^*, \boldsymbol{\beta}^\dagger + \boldsymbol{\theta}/\sqrt{n})$ and $a_i = Y_i^* f(\mathbf{X}_i^*, \boldsymbol{\beta}^\dagger)$ in (S.2), we have

$$|R_{i,n}(\boldsymbol{\theta})| \leq \frac{1}{N\pi_i^*} \left| \frac{f(\mathbf{X}_i^*, \boldsymbol{\theta})}{\sqrt{n}} \right| U_i \left(\left| \frac{f(\mathbf{X}_i^*, \boldsymbol{\theta})}{\sqrt{n}} \right| \right), \tag{S.3}$$

where $U_i(t) = \mathbb{I}(|1 - Y_i^* f(\mathbf{X}_i^*, \boldsymbol{\beta}^\dagger)| \leq t)$ with respect to the i -th subsample point for

$t \in \mathbb{R}$. By (S.3), for each fixed $\boldsymbol{\theta}$ we obtain

$$\begin{aligned}
\mathbb{E} \left[\sum_{i=1}^n \{R_{i,n}(\boldsymbol{\theta}) - \mathbb{E}(R_{i,n}(\boldsymbol{\theta}))\} \right]^2 &= \mathbb{E} \left\{ \mathbb{E} \left[\sum_{i=1}^n \{R_{i,n}(\boldsymbol{\theta}) - \mathbb{E}(R_{i,n}(\boldsymbol{\theta}))\} \right]^2 \middle| \mathcal{D}_N \right\} \\
&= \frac{n}{N^2} \sum_{j=1}^N \mathbb{E} \left[\frac{1}{\pi_j} \{R_{j,N}(\boldsymbol{\theta}) - \mathbb{E}(R_{j,N}(\boldsymbol{\theta}))\}^2 \right] \\
&\leq \frac{n}{N^2} \sum_{j=1}^N \mathbb{E} \left\{ \frac{1}{\pi_j} R_{j,N}^2(\boldsymbol{\theta}) \right\} \\
&\leq \frac{n}{N^2} \sum_{j=1}^N \mathbb{E} \left\{ \frac{1}{\pi_j} (1 + \|\mathbf{X}_j\|^2) \frac{\|\boldsymbol{\theta}\|^2}{n} U_j \left(\sqrt{1 + \|\mathbf{X}_j\|^2} \frac{\|\boldsymbol{\theta}\|}{\sqrt{n}} \right) \right\} \\
&\leq \frac{\|\boldsymbol{\theta}\|^2}{N^2} \sum_{j=1}^N \mathbb{E} \left\{ \frac{1}{\pi_j} (1 + \|\mathbf{X}_j\|^2) U_j \left(\sqrt{1 + \|\mathbf{X}_j\|^2} \frac{\|\boldsymbol{\theta}\|}{\sqrt{n}} \right) \right\}.
\end{aligned}$$

By Assumption 1 implies that $\mathbb{E}(\|\mathbf{X}\|^4) < \infty$, there exists c_1 such that

$$\mathbb{E} \{ (1 + \|\mathbf{X}\|^4) \mathbb{I}(\|\mathbf{X}\| > c_1) \} < \varepsilon_2/2,$$

for any $\varepsilon_2 > 0$. Let $U(t) = \mathbb{I}(|1 - Yf(\mathbf{X}, \boldsymbol{\beta}^\dagger)| \leq t)$ for $t \in \mathbb{R}$. By Assumption 4 and holder inequality, we have

$$\begin{aligned}
&\frac{1}{N^2} \sum_{j=1}^N \mathbb{E} \left\{ \frac{1}{\pi_j} (1 + \|\mathbf{X}_j\|^2) U_j \left(\sqrt{1 + \|\mathbf{X}_j\|^2} \frac{\|\boldsymbol{\theta}\|}{\sqrt{n}} \right) \right\} \\
&\leq \frac{1}{N^2} \sum_{j=1}^N \mathbb{E} \left\{ \frac{1}{\pi_j} (1 + \|\mathbf{X}_j\|^2) \mathbb{I}(\|\mathbf{X}_j\| > c_1) \right\} + \frac{1}{N^2} \sum_{j=1}^N \mathbb{E} \left\{ \frac{1 + c_1^2}{\pi_j} U \left(\sqrt{1 + c_1^2} \frac{\|\boldsymbol{\theta}\|}{\sqrt{n}} \right) \right\} \\
&\leq \sqrt{\mathbb{E} \left(\frac{1}{N^3} \sum_{j=1}^N \frac{1}{\pi_j^2} \right)} \sqrt{\mathbb{E} \left\{ \frac{1}{N} \sum_{j=1}^N (1 + \|\mathbf{X}_j\|^2)^2 \mathbb{I}(\|\mathbf{X}_j\| > c_1) \right\}} \\
&\quad + (1 + c_1^2) \sqrt{\mathbb{E} \left(\frac{1}{N^3} \sum_{j=1}^N \frac{1}{\pi_j^2} \right)} \sqrt{\frac{1}{N} \sum_{j=1}^N \mathbb{P} \left\{ U \left(\sqrt{1 + c_1^2} \frac{\|\boldsymbol{\theta}\|}{\sqrt{n}} \right) = 1 \right\}},
\end{aligned}$$

By Assumption 1, the conditional distribution of \mathbf{X} given Y is not degenerate,

which implies $\lim_{t \rightarrow 0} \mathbb{P}(U(t) = 1) = 0$. We can take a large c_2 such that

$$\mathbb{P} \left\{ U \left(\sqrt{1 + c_1^2} \frac{\|\boldsymbol{\theta}\|}{\sqrt{n}} \right) = 1 \right\} < \varepsilon_2 / \{2(1 + c_1^2)\},$$

for $n > c_2$. By Assumption 4, it proves that $\mathbb{E} [\sum_{i=1}^n \{R_{i,n}(\boldsymbol{\theta}) - \mathbb{E}(R_{i,n}(\boldsymbol{\theta}))\}]^2 \rightarrow 0$.

By (S.1), for each fixed $\boldsymbol{\theta}$

$$\Lambda_n(\boldsymbol{\theta}) = \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{H}(\boldsymbol{\beta}^\dagger) \boldsymbol{\theta} + \sqrt{n} \mathbf{W}_n^\top \boldsymbol{\theta} + o_P(1).$$

Last, we devote to giving the Bahadur representation of $\tilde{\boldsymbol{\beta}}$. Let $\boldsymbol{\kappa}_n = -\sqrt{n} \mathbf{H}(\boldsymbol{\beta}^\dagger)^{-1} \mathbf{W}_n$ and Θ be a convex open subset in \mathbb{R}^{p+1} . By Convexity Lemma in Pollard (1991), we have

$$\Lambda_n(\boldsymbol{\theta}) = \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\kappa}_n)^\top \mathbf{H}(\boldsymbol{\beta}^\dagger) (\boldsymbol{\theta} - \boldsymbol{\kappa}_n) - \frac{1}{2} \boldsymbol{\kappa}_n^\top \mathbf{H}(\boldsymbol{\beta}^\dagger) \boldsymbol{\kappa}_n + r_n(\boldsymbol{\theta}),$$

where for each compact set K of Θ , the aforementioned part is shown for every $\boldsymbol{\theta} \in \Theta$, and then we have $\sup_{\boldsymbol{\theta} \in K} |r_n(\boldsymbol{\theta})| \rightarrow 0$ in probability. Lemma S.4 shows that $\boldsymbol{\kappa}_n$ is asymptotically normal which will be proved in the next section, then there exists a compact set $K \in \mathcal{B}_\rho$ with probability close to one, where \mathcal{B}_ρ is a closed ball with center $\boldsymbol{\kappa}_n$ and radius ρ . Let $\Delta_n = \sup_{\boldsymbol{\theta} \in \mathcal{B}_\rho} |r_n(\boldsymbol{\theta})|$. Then we have

$$\Delta_n \rightarrow 0 \text{ in probability.} \tag{S.4}$$

Next, we discuss the behavior of $\Lambda_n(\boldsymbol{\theta})$ outside the closed ball \mathcal{B}_ρ . Consider $\boldsymbol{\theta} = \boldsymbol{\kappa}_n + \gamma \mathbf{e}$, with $\gamma > \rho$ and the unit vector \mathbf{e} . A boundary point $\boldsymbol{\theta}^\dagger = \boldsymbol{\kappa}_n + \rho \mathbf{e}$. Under Assumptions 1–3 and a similar discussion in Lemma 5 of Koo et al. (2008), there exists a constant c_3 such that $\boldsymbol{\beta}^\top \mathbf{H}(\boldsymbol{\beta}^\dagger) \boldsymbol{\beta} \geq c_3 \|\boldsymbol{\beta}\|^2$. Then, by the convexity of $\Lambda_n(\boldsymbol{\theta})$ and

the definition of Δ_n , we have

$$\begin{aligned}
\frac{\rho}{\gamma}\Lambda_n(\boldsymbol{\theta}) + \left(1 - \frac{\rho}{\gamma}\right)\Lambda_n(\boldsymbol{\kappa}_n) &\geq \Lambda_n\left(\frac{\rho}{\gamma}\boldsymbol{\theta} + \left(1 - \frac{\rho}{\gamma}\right)\boldsymbol{\kappa}_n\right) \\
&= \Lambda_n(\boldsymbol{\theta}^\dagger) \\
&\geq \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\kappa}_n)^\top \mathbf{H}(\boldsymbol{\beta}^\dagger)(\boldsymbol{\theta} - \boldsymbol{\kappa}_n) - \frac{1}{2}\boldsymbol{\kappa}_n^\top \mathbf{H}(\boldsymbol{\beta}^\dagger)\boldsymbol{\kappa}_n - \Delta_n \\
&\geq \frac{c_3}{2}\rho^2 + \Lambda_n(\boldsymbol{\kappa}_n) - 2\Delta_n,
\end{aligned}$$

which implies that

$$\inf_{\|\boldsymbol{\theta} - \boldsymbol{\kappa}_n\| > \rho} \Lambda_n(\boldsymbol{\theta}) \geq \Lambda_n(\boldsymbol{\kappa}_n) + \left(\frac{c_3}{2}\rho^2 - 2\Delta_n\right).$$

By (S.4), we can take Δ_n such that $2\Delta_n < c_3\rho^2/2$ with probability tending to one. Thus $\inf_{\|\boldsymbol{\theta} - \boldsymbol{\kappa}_n\| > \rho} \Lambda_n(\boldsymbol{\theta}) \geq \Lambda_n(\boldsymbol{\kappa}_n)$. This implies the minimum of $\Lambda_n(\boldsymbol{\theta})$ cannot occur at any $\boldsymbol{\theta}$ with $\|\boldsymbol{\theta} - \boldsymbol{\kappa}_n\| > \rho$. Hence for each $\rho > 0$ and let $\tilde{\boldsymbol{\theta}}_n = \sqrt{n}(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^\dagger)$, we have $\mathbb{P}(\|\tilde{\boldsymbol{\theta}}_n - \boldsymbol{\kappa}_n\| > \rho) \rightarrow 0$. Thus

$$\sqrt{n}(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^\dagger) = -\sqrt{n}\mathbf{H}(\boldsymbol{\beta}^\dagger)^{-1}\mathbf{W}_n + o_P(1).$$

The theorem follows the above arguments. \square

Appendix C: Proof of asymptotic normality

Recall that

$$\begin{aligned}
\mathbf{M} &= \sum_{i=1}^n M_i = \sum_{i=1}^n \frac{1}{nN\pi_i^*} \xi_i^* Y_i^* \tilde{\mathbf{X}}_i^* - \sum_{i=1}^n \left(\frac{1}{nN} \sum_{j=1}^N \xi_j Y_j \tilde{\mathbf{X}}_j \right), \\
\mathbf{Q} &= \frac{1}{N} \sum_{j=1}^N \xi_j Y_j \tilde{\mathbf{X}}_j, \quad \mathbf{T} = \frac{1}{n} \sum_{i=1}^n \frac{1}{N\pi_i^*} \xi_i^* Y_i^* \tilde{\mathbf{X}}_i^*, \quad \mathbf{B}_N = \mathbf{V}_T^{-1/2} \mathbf{V}_M \mathbf{V}_T^{-1/2},
\end{aligned} \tag{S.5}$$

where \mathbf{V}_T and \mathbf{V}_M are the variances of \mathbf{T} and \mathbf{M} .

Lemma S.2. $\{\mathbf{M}_i, i = 1, \dots, n\}$ in (S.5) is a martingale difference sequence relative to the filtration $\{\mathcal{F}_{N,i}, i = 1, \dots, n\}$.

Proof. The $\mathcal{F}_{N,i}$ -measurability follows from the definition of \mathbf{M}_i and the definition of the filtration $\{\mathcal{F}_{N,i}, i = 1, \dots, n\}$. Moreover, we have

$$\begin{aligned} \mathbb{E}\{\mathbf{M}_i \mid \mathcal{F}_{N,i-1}\} &= \mathbb{E}_{Y|\mathbf{X}} \left\{ \frac{1}{nN\pi_i^*} \xi_i^* Y_i^* \widetilde{\mathbf{X}}_i^* \right\} - \frac{1}{nN} \sum_{j=1}^N \xi_j Y_j \widetilde{\mathbf{X}}_j \\ &= \frac{1}{nN} \sum_{i=1}^N \xi_i Y_i \widetilde{\mathbf{X}}_i - \frac{1}{nN} \sum_{j=1}^N \xi_j Y_j \widetilde{\mathbf{X}}_j \\ &= 0, \end{aligned}$$

where $\mathbb{E}_{Y|\mathbf{X}}$ is the expectation with respect to sampling randomness or the conditional expectation of Y given \mathbf{X}_1^N with $\mathbf{X}_1^N = (\mathbf{X}_1, \dots, \mathbf{X}_N)$. Then $\{\mathbf{M}_i, i = 1, \dots, n\}$ is a martingale difference sequence. \square

Lemma S.3. Suppose Assumptions 1 and 4 hold. Let \mathbf{V}_T and \mathbf{V}_Q denote the variances of \mathbf{T} and \mathbf{Q} . For any $\mathbf{t} \in \mathbb{R}^{p+1}$, we have

$$\left| \mathbb{E} \left\{ \exp \left(i \mathbf{t}^\top \mathbf{V}_T^{-1/2} \mathbf{Q} \right) \right\} - \mathbb{E} \left\{ \exp \left(i \mathbf{t}^\top \mathbf{V}_T^{-1/2} \mathbf{V}_Q^{1/2} \mathbf{A}_0 \right) \right\} \right| \rightarrow 0,$$

as $N \rightarrow \infty$, where $\mathbf{A}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{p+1})$.

Proof. Note \mathbf{Q} is a sum of i.i.d mean zero random vectors, $\xi_j Y_j \widetilde{\mathbf{X}}_j$. The Linderberg-Feller conditions are satisfied by Assumption 1 and Assumption 4, then we have

$$\mathbf{V}_Q^{-1/2} \mathbf{Q} \rightarrow \mathcal{N}(\mathbf{0}, \mathbf{I}_{p+1}). \quad (\text{S.6})$$

Furthermore, for any $\boldsymbol{\varsigma} \in \mathbb{R}^{p+1}$ and as $N \rightarrow \infty$

$$\left| \mathbb{E} \left\{ \exp \left(i \boldsymbol{\varsigma}^\top \mathbf{V}_Q^{-1/2} \mathbf{Q} \right) \right\} - \mathbb{E} \left\{ \exp \left(i \boldsymbol{\varsigma}^\top \mathbf{A}_0 \right) \right\} \right| \rightarrow 0.$$

Let $\boldsymbol{\varsigma} = \mathbf{V}_Q^{1/2} \mathbf{V}_T^{-1/2} \mathbf{t}^\top$. For any fixed \mathbf{t} , we need to verify the following condition to prove this lemma

$$\sup_N \|\boldsymbol{\varsigma}\| < \infty.$$

We note that $\|\boldsymbol{\varsigma}\| \leq \sigma_{\max}(\mathbf{V}_Q^{1/2} \mathbf{V}_T^{-1/2}) \cdot \|\mathbf{t}\|$, where $\sigma_{\max}(\cdot)$ denotes the maximum eigenvalue of the corresponding matrix. Hence it is enough to show $\sigma_{\max}(\mathbf{V}_Q^{1/2} \mathbf{V}_T^{-1/2}) \leq 1$. Since the covariance matrix \mathbf{V}_Q and \mathbf{V}_T are positive-defined, the following equation holds

$$\mathbf{V}_Q^{1/2} \mathbf{V}_T^{-1/2} = \mathbf{V}_T^{1/4} \left(\mathbf{V}_T^{-1/4} \mathbf{V}_Q^{1/2} \mathbf{V}_T^{-1/4} \right) \mathbf{V}_T^{-1/4},$$

thus $\mathbf{V}_Q^{1/2} \mathbf{V}_T^{-1/2}$ is similar to $\mathbf{V}_T^{-1/4} \mathbf{V}_Q^{1/2} \mathbf{V}_T^{-1/4}$. It only needs to show $\sigma_{\max}(\mathbf{V}_T^{-1/4} \mathbf{V}_Q^{1/2} \mathbf{V}_T^{-1/4}) \leq 1$, which is equal to show

$$\mathbf{I}_{p+1} - \mathbf{V}_T^{-1/4} \mathbf{V}_Q^{1/2} \mathbf{V}_T^{-1/4} = \mathbf{V}_T^{-1/4} \left(\mathbf{V}_T^{1/2} - \mathbf{V}_Q^{1/2} \right) \mathbf{V}_T^{-1/4} > 0,$$

that is equivalent to show $\mathbf{V}_T^{1/2} - \mathbf{V}_Q^{1/2}$ is positive-defined.

Recall that $\mathbf{M} = \mathbf{T} - \mathbf{Q}$ and by Lemma S.1, we have $\mathbf{V}_T - \mathbf{V}_Q = \mathbf{V}_M > 0$. Then by the Löwner-Heinz theorem in Zhan (2004), we get $\mathbf{V}_T^{1/2} - \mathbf{V}_Q^{1/2} > 0$ which completes the proof of this lemma. \square

Lemma S.4. *Suppose Assumptions 1 and 4 hold. Then we have*

$$\mathbf{V}_T^{-1/2} \mathbf{T} \rightarrow \mathcal{N}(\mathbf{0}, \mathbf{I}_{p+1}).$$

Proof. Recall the conditions in Lemma S.1 with

$$\boldsymbol{\eta}_{ki} = \boldsymbol{\eta}_{Ni}, \mathbf{Z}_k = \mathbf{V}_T^{-1/2} \mathbf{Q}, \mathbf{B}_k = \mathbf{B}_N, \mathbf{L}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{p+1}).$$

By Lemma S.2, $\{M_i, i = 1, \dots, n\}$ is a martingale difference sequence, then the first two conditions in Lemma S.2 are easily satisfied by Assumption 1. It suffices to show the third condition in Lemma S.1 holds.

By (S.6) in Lemma S.3, we have $\mathbf{V}_Q^{-1/2}\mathbf{Q} \rightarrow \mathcal{N}(\mathbf{0}, \mathbf{I}_{p+1})$. Next, we devote ourselves to verifying the third condition in Lemma S.1. Let \mathbf{V}_M be the variance of \mathbf{M} . For any $\mathbf{t} \in \mathbb{R}^{p+1}$, we have the following characteristic function

$$\begin{aligned} & \mathbb{E} \left\{ \exp \left(i\mathbf{t}^\top \mathbf{V}_T^{-1/2} \mathbf{Q} \right) \right\} \cdot \exp \left(-\frac{1}{2} \mathbf{t}^\top \mathbf{V}_T^{-1/2} \mathbf{V}_M \mathbf{V}_T^{-1/2} \mathbf{t} \right) \\ &= \left\{ \exp \left(i\mathbf{t}^\top \mathbf{V}_T^{-1/2} \mathbf{V}_Q \mathbf{V}_T^{-1/2} \mathbf{t} \right) + o(1) \right\} \cdot \exp \left(-\frac{1}{2} \mathbf{t}^\top \mathbf{V}_T^{-1/2} \mathbf{V}_M \mathbf{V}_T^{-1/2} \mathbf{t} \right) \\ &= \left\{ \exp \left(i\mathbf{t}^\top \mathbf{V}_T^{-1/2} \mathbf{V}_Q \mathbf{V}_T^{-1/2} \mathbf{t} \right) \right\} \cdot \exp \left(-\frac{1}{2} \mathbf{t}^\top \mathbf{V}_T^{-1/2} \mathbf{V}_M \mathbf{V}_T^{-1/2} \mathbf{t} \right) + o(1) \\ &= \exp \left(-\frac{1}{2} \mathbf{t}^\top \mathbf{t} \right) + o(1), \end{aligned}$$

where the first equality holds by Lemma S.3. And the third condition in Lemma S.1 is satisfied. Then by Lemma S.1 and (S.6) we have

$$\mathbf{V}_T^{-1/2} \mathbf{Q} + \mathbf{V}_T^{-1/2} \mathbf{M} = \mathbf{V}_T^{-1/2} \mathbf{T} \rightarrow \mathcal{N}(\mathbf{0}, \mathbf{I}_{p+1}).$$

□

Proof of Theorem 2. By Theorem 1 and Lemma S.4, we have

$$\sqrt{n}(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^\dagger) = -\sqrt{n}\mathbf{H}(\boldsymbol{\beta}^\dagger)^{-1}\mathbf{T} + o_p(1).$$

It follows that

$$\mathbf{V}_T^{-1/2} \mathbf{H}(\boldsymbol{\beta}^\dagger)(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^\dagger) + o_p(1) = -\mathbf{V}_T^{-1/2} \mathbf{T}.$$

By Lemma S.4, we have

$$\mathbf{V}^{-1/2}(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^\dagger) \rightarrow \mathcal{N}(\mathbf{0}, \mathbf{I}_{p+1}),$$

where $\mathbf{V} = \mathbf{H}(\boldsymbol{\beta}^\dagger)^{-1} \mathbf{V}_T \mathbf{H}(\boldsymbol{\beta}^\dagger)^{-1}$. □

Appendix D: Proof of Theorem 3

Proof of Theorem 3. Recall that $\mathbf{X}_1^N = (\mathbf{X}_1, \dots, \mathbf{X}_N)$ and $Y_1^N = (Y_1, \dots, Y_N)$, then $\mathcal{D}_N = \{\mathbf{X}_1^N, Y_1^N\}$. Let $\text{var}(Y | \mathbf{X})$ be the conditional variance of Y given \mathbf{X} . First we calculate $\text{var}(\mathbf{T} | \mathbf{X}_1^N)$. We have

$$\text{var}(\mathbf{T} | \mathbf{X}_1^N) = \mathbb{E}_{Y|\mathbf{X}} \{ \text{var}(\mathbf{T} | \mathcal{D}_N) \} + \text{var}_{Y|\mathbf{X}} \{ \mathbb{E}(\mathbf{T} | \mathcal{D}_N) \}.$$

Some algebra yields

$$\begin{aligned} \text{var}_{Y|\mathbf{X}} \{ \mathbb{E}(\mathbf{T} | \mathcal{D}_N) \} &= \text{var}_{Y|\mathbf{X}} \left(\frac{1}{N} \sum_{j=1}^N \xi_j Y_j \widetilde{\mathbf{X}}_j \right) \\ &= \frac{1}{N^2} \sum_{j=1}^N \mathbb{E}_{Y|\mathbf{X}} \left(\xi_j^2 Y_j^2 \widetilde{\mathbf{X}}_j \widetilde{\mathbf{X}}_j^\top \right) - \frac{1}{N^2} \sum_{j=1}^N \left\{ \mathbb{E}_{Y|\mathbf{X}} (\xi_j Y_j \widetilde{\mathbf{X}}_j) \right\}^2 \\ &= \frac{1}{N^2} \sum_{j=1}^N \mathbb{E}_{Y|\mathbf{X}} \left(\xi_j \widetilde{\mathbf{X}}_j \widetilde{\mathbf{X}}_j^\top \right) - \frac{1}{N^2} \sum_{j=1}^N \left\{ \mathbb{E}_{Y|\mathbf{X}} (\xi_j Y_j \widetilde{\mathbf{X}}_j) \right\}^2, \end{aligned} \tag{S.7}$$

where the third equality holds by the fact that $\xi_j^2 = \xi_j$ and $Y_j^2 = 1$. Next

$$\begin{aligned} \mathbb{E}_{Y|\mathbf{X}} \{ \text{var}(\mathbf{T} | \mathcal{D}_N) \} &= \frac{1}{nN^2} \sum_{j=1}^N \mathbb{E}_{Y|\mathbf{X}} \left\{ \pi_j \left(\frac{1}{\pi_j^2} \xi_j^2 Y_j^2 \widetilde{\mathbf{X}}_j \widetilde{\mathbf{X}}_j^\top \right) \right\} - \frac{1}{nN} \sum_{j=1}^N \left\{ \mathbb{E}_{Y|\mathbf{X}} (\xi_j Y_j \widetilde{\mathbf{X}}_j) \right\}^2 \\ &= \frac{1}{nN^2} \sum_{j=1}^N \mathbb{E}_{Y|\mathbf{X}} \left\{ \frac{1}{\pi_j} \xi_j \widetilde{\mathbf{X}}_j \widetilde{\mathbf{X}}_j^\top \right\} - \frac{1}{nN} \sum_{j=1}^N \left\{ \mathbb{E}_{Y|\mathbf{X}} (\xi_j Y_j \widetilde{\mathbf{X}}_j) \right\}^2. \end{aligned} \tag{S.8}$$

In view of (S.7) and (S.8), we get

$$\begin{aligned} \text{var}(\mathbf{T} \mid \mathbf{X}_1^N) &= \frac{1}{nN^2} \sum_{j=1}^N \mathbb{E}_{Y|\mathbf{X}} \left(\frac{1}{\pi_j} \xi_j \widetilde{\mathbf{X}}_j \widetilde{\mathbf{X}}_j^\top \right) + \frac{1}{N^2} \sum_{j=1}^N \mathbb{E}_{Y|\mathbf{X}} \left(\xi_j \widetilde{\mathbf{X}}_j \widetilde{\mathbf{X}}_j^\top \right) \\ &\quad - \frac{1}{N} \sum_{j=1}^N \left\{ \mathbb{E}_{Y|\mathbf{X}} \left(\xi_j Y_j \widetilde{\mathbf{X}}_j \right) \right\}^2 \left(\frac{1}{N} + \frac{1}{n} \right). \end{aligned}$$

Next we calculate \mathbf{V}_T through

$$\mathbf{V}_T = \mathbb{E} \left\{ \text{var}(\mathbf{T} \mid \mathbf{X}_1^N) \right\} + \text{var} \left\{ \mathbb{E}(\mathbf{T} \mid \mathbf{X}_1^N) \right\}.$$

A simple calculation shows that

$$\begin{aligned} \mathbb{E}(T \mid \mathbf{X}_1^N) &= \mathbb{E} \left\{ \mathbb{E}(\mathbf{T} \mid \mathbf{X}_1^N, Y_1^N) \right\} = \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{Y|\mathbf{X}} \left(\xi_j Y_j \widetilde{\mathbf{X}}_j \right), \\ \text{var} \left\{ \mathbb{E}(T \mid \mathbf{X}_1^N) \right\} &= \frac{1}{N^2} \sum_{j=1}^N \mathbb{E}_{Y|\mathbf{X}} \left(\xi_j \widetilde{\mathbf{X}}_j \widetilde{\mathbf{X}}_j^\top \right) - \frac{1}{N^2} \sum_{j=1}^N \left\{ \mathbb{E}_{Y|\mathbf{X}} \left(\xi_j Y_j \widetilde{\mathbf{X}}_j \right) \right\}^2. \end{aligned}$$

Therefore, we have

$$\mathbf{V}_T = \frac{1}{nN^2} \sum_{j=1}^N \mathbb{E}_{Y|\mathbf{X}} \left(\frac{1}{\pi_j} \xi_j \widetilde{\mathbf{X}}_j \widetilde{\mathbf{X}}_j^\top \right) + \mathbf{C},$$

where $\mathbf{C} = 2N^{-2} \sum_{j=1}^N \mathbb{E}_{Y|\mathbf{X}} \left(\xi_j \widetilde{\mathbf{X}}_j \widetilde{\mathbf{X}}_j^\top \right) - N^{-1} \sum_{j=1}^N \left\{ \mathbb{E}_{Y|\mathbf{X}} \left(\xi_j Y_j \widetilde{\mathbf{X}}_j \right) \right\}^2 (2N^{-1} + n^{-1})$

is a constant matrix that does not depend on $\boldsymbol{\pi}$.

Let $\text{tr}(\mathbf{A})$ denotes the trace of matrix \mathbf{A} . We minimize $\text{tr}(\mathbf{V}_T)$ to obtain the A-optimality subsampling probability

$$\begin{aligned} \text{tr}(\mathbf{V}_T) &= \frac{1}{nN^2} \sum_{j=1}^N \text{tr} \left\{ \mathbb{E}_{Y|\mathbf{X}} \left(\frac{1}{\pi_j} \xi_j \mathbf{H}(\boldsymbol{\beta}^\dagger)^{-1} \widetilde{\mathbf{X}}_j \widetilde{\mathbf{X}}_j^\top \mathbf{H}(\boldsymbol{\beta}^\dagger)^{-1} \right) \right\} + \text{tr}(\mathbf{C}) \\ &= \frac{1}{nN^2} \mathbb{E}_{Y|\mathbf{X}} \left\{ \sum_{j=1}^N \pi_j \sum_{j=1}^N \left(\frac{1}{\pi_j} \xi_j \|\mathbf{H}(\boldsymbol{\beta}^\dagger)^{-1} \widetilde{\mathbf{X}}_j\|^2 \right) \right\} + \text{tr}(\mathbf{C}) \\ &\geq \frac{1}{nN^2} \left\{ \sum_{j=1}^N \text{P} \left(Y_j f(\mathbf{X}_j, \boldsymbol{\beta}^\dagger) \leq 1 \right) \|\mathbf{H}(\boldsymbol{\beta}^\dagger)^{-1} \widetilde{\mathbf{X}}_j\|^2 \right\} + \text{tr}(\mathbf{C}), \end{aligned}$$

where the last inequality follows from the Cauchy-Schwarz inequality, and the equality holds if and only if

$$\pi_j^A \propto \mathbb{I}(Y_j f(\mathbf{X}_j, \boldsymbol{\beta}^\dagger) \leq 1) \|\mathbf{H}(\boldsymbol{\beta}^\dagger)^{-1} \widetilde{\mathbf{X}}_j\|.$$

Note that $\mathbf{H}(\boldsymbol{\beta}^\dagger)^{-1} \text{var}(\mathbf{T} \mid \mathbf{X}_1^N) \mathbf{H}(\boldsymbol{\beta}^\dagger)^{-1}$ depends on subsampling probability π only through $\text{var}(\mathbf{T} \mid \mathbf{X}_1^N)$. Hence, by the similar argument for minimizing $\text{tr}\{\text{var}(\mathbf{T} \mid \mathbf{X}_1^N)\}$, we get the L-optimality subsampling probability

$$\pi_j^L \propto \mathbb{I}(Y_j f(\mathbf{X}_j, \boldsymbol{\beta}^\dagger) \leq 1) \|\widetilde{\mathbf{X}}_j\|.$$

□

Appendix E: Additional simulation results

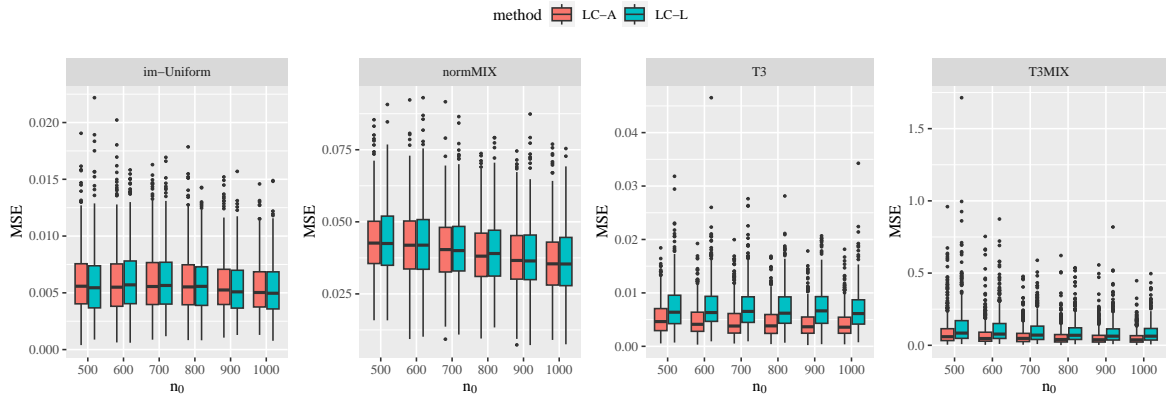


Figure S1: Comparison of MSE for approximating the full sample SVM estimator $\widehat{\boldsymbol{\beta}}$ with different pilot subsample sizes given $n = 1000$ under Scenarios I–IV.

To assess the impact of the pilot study in our proposed algorithm, we conduct the following boxplot by 500 replications on the four scenarios presented in Section 4. Figure S1 reveals that the MSE is not sensitive to the pilot subsample size n_0 .

As n_0 increases, the boxplot shows a slight decrease in MSE, suggesting that a smaller pilot subsample size can reduce computational costs without significantly compromising accuracy.

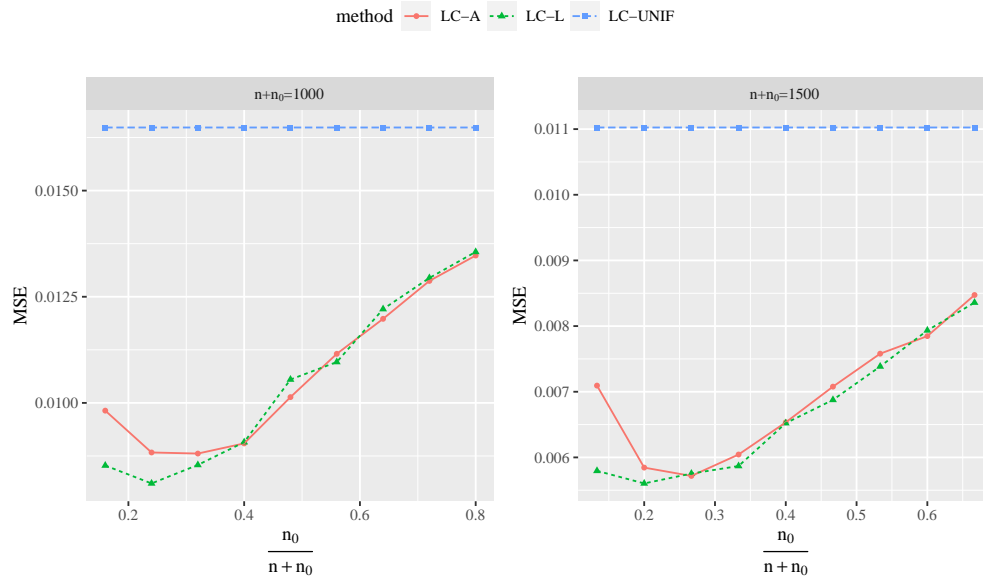


Figure S2: Comparison of mean squared errors (MSEs) for approximating the full sample SVM estimator $\hat{\beta}$ with different subsample size allocations under Scenario I.

Moreover, we fix the total subsample size of $n + n_0$ and vary the proportions of n and n_0 . It provides practical guidelines on allocating subsamples in two steps. We evaluate both $\hat{\pi}^A$ and $\hat{\pi}^L$ and the results are presented in Figure S2 under Scenario I. It illustrates that the MSEs increase when n_0 is either too small or too large. This is because that if n_0 is too small, the pilot estimate is not accurate, and thus the optimal subsampling probabilities may not be well approximated; on the other hand, if n_0 is too large, there is not enough sampling budget to select informative subsample in subsequent steps. Figure S2 shows that our methods perform well when the ratio $n_0/(n + n_0)$ is around (0.2, 0.4). Therefore, we use $n_0 = 500$ in our simulation studies

with $N = 10^5$.

Bandwidth selection is a critical issue in nonparametric estimation. In Table S1, we compare the MSE and accuracy of LC-A with three bandwidth selectors: Silverman’s rule of thumb (ROT, Silverman, 1986), Sheather and Jones method, (SJ, Sheather and Jones, 1991), and biased cross-validation, (BCV, Scott and Terrell, 1987). Clearly, The results demonstrate that the choice of bandwidth selector has a negligible impact on the empirical MSE and accuracy. To this end, we employ the commonly-used bandwidth selector, Silverman’s rule of thumb (Silverman, 1986), in our numerical analysis.

Table S1: Comparison of MSE (10^{-2}) and prediction accuracy (%) for LC-A against different bandwidth selectors under Scenarios I–II when $n = 1000$.

Scenario	n_0	ROT		SJ		BCV	
		MSE	Accuracy	MSE	Accuracy	MSE	Accuracy
	300	0.68	95.54	0.92	94.52	0.65	94.56
im-Uniform	400	0.64	94.53	0.85	94.52	0.61	94.54
	500	0.60	94.53	0.75	94.52	0.60	94.53
	300	4.84	97.52	4.89	97.52	4.87	97.52
normMIX	400	4.49	97.53	4.63	97.53	4.56	97.53
	500	4.33	97.54	4.43	97.54	4.35	97.54

References

- Koo, J.-Y., Lee, Y., Kim, Y., and Park, C. (2008). A Bahadur representation of the linear support vector machine. *Journal of Machine Learning Research*, **9**(7):1343–1368.
- Pollard, D. (1991). Asymptotics for least absolute deviation regression estimators. *Econometric Theory*, **7**(2):186–199.
- Scott, D. W. and Terrell, G. R. (1987). Biased and unbiased cross-validation in density estimation. *Journal of the American Statistical Association*, **82**(400):1131–1146.
- Sheather, S. J. and Jones, M. C. (1991). A reliable data-based bandwidth selection method for kernel density estimation. *Journal of the Royal Statistical Society: Series B (Methodological)*, **53**(3):683–690.
- Silverman, B. W. (1986). *Density estimation for statistics and data analysis*. Chapman & Hall.
- Zhan, X. (2004). *Matrix Inequalities*. Springer.
- Zhang, T., Ning, Y., and Ruppert, D. (2021). Optimal sampling for generalized linear models under measurement constraints. *Journal of Computational and Graphical Statistics*, **30**(1):106–114.