

Supplementary Material: Sparse Hilbert–Schmidt Independence Criterion Regression

A Asymptotics for U-statistics

U-statistics aim at estimating a parameter, say θ , based on the sample X_1, \dots, X_n . Let us denote by $\psi(\cdot)$ a symmetric kernel that will enter in the U-statistic, which is defined as follows.

Definition. Given a real-valued measurable kernel $\psi(\cdot)$, which is of degree r and symmetric in its argument, for a sample X_1, \dots, X_n of size $n \geq r$, then a U-statistic with kernel $\psi(\cdot)$ is defined as

$$U_n = \frac{1}{\binom{n}{r}} \sum_{C_{n,r}} \psi(X_{i_1}, \dots, X_{i_r}), \quad (6)$$

where the summation is over the set $C_{n,r}$ of all $\binom{n}{r}$ combinations of r integers $i_1 < \dots < i_r$ chosen from $\{1, \dots, n\}$.

Note that this definition concerns the particular case of symmetric kernel. Should the kernel be non-symmetric, then the U-statistic would be defined as

$$U_n = \frac{1}{(n)_r} \sum_{i_r^n} \psi(X_{i_1}, \dots, X_{i_r}),$$

where $(n)_r = \frac{n!}{(n-r)!}$ is the Pochhammer symbol and i_r^n is the index set, which is the set of all r -tuples drawn without replacement from $\{1, \dots, n\}$.

When deriving the asymptotic distribution of U-statistic based estimator, the variance and covariance of U-statistics are key quantities entering the asymptotic variance covariance matrix. For a given parameter θ and corresponding symmetric kernel $\psi(x_1, \dots, x_n)$, we consider the class of distribution \mathcal{F} so that $\text{Var}(\psi(X_1, \dots, X_n)) < \infty$. We then define

$$\psi_r(x_1, \dots, x_r) = \mathbb{E}_{X_{r+1}, \dots, X_n} [\psi(x_1, \dots, x_r, X_{r+1}, \dots, X_n)],$$

where X_{r+1}, \dots, X_n are i.i.d. random variables. Then $\psi_0 = \theta$, $\psi_n(x_1, \dots, x_n) = \psi(x_1, \dots, x_n)$ and

$$\forall r, \mathbb{E}[\psi_r(X_1, \dots, X_r)] = \theta.$$

Then to define the variance of the U-statistic, denote by $\sigma_r^2 = \text{Var}(\psi_r(X_1, \dots, X_r))$.

Theorem A.1. Variance of a U-statistic, (Serfling, 1980)

The variance of a U-statistic given by Eq. (6) is

$$\text{Var}(U_n) = \binom{n}{r}^{-1} \sum_{s=1}^r \binom{r}{s} \binom{n-r}{r-s} \sigma_s^2. \quad (7)$$

Theorem A.2. Asymptotic distribution of a U-statistic, (Serfling, 1980)

If $\text{Var}(\psi(X_1, \dots, X_n)) < \infty$ and $\sigma_1^2 > 0$, then

$$\sqrt{n}(U_n - \theta) \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}_{\mathbb{R}}(0, r^2 \sigma_1^2)$$

Remark. The asymptotic distribution depends on whether $\sigma_1^2 > 0$, in which case we obtain a normal distribution, or $\sigma_1^2 = 0$, in which case the asymptotic law is an infinite linear combination of $\chi^2(1)$ distributions.

Theorem A.3. Covariance of U-statistics, (Serfling, 1980)

Let $U_n^{(k)}$ and $U_n^{(l)}$ be two U-statistics based on the same sample $\mathbf{X} = (X_1, \dots, X_n)$, respectively equipped with kernel ψ_k and ψ_l of degree r_k, r_l , respectively ($r_k \leq r_l$). Then the covariance is

$$\text{cov}(U_n^{(k)}, U_n^{(l)}) = \binom{n}{r_k}^{-1} \sum_{s=1}^{r_k} \binom{r_l}{s} \binom{n-r_l}{r_k-s} \sigma_s^2, \quad (8)$$

where

$$\sigma_s^2 = \text{cov}(\mathbb{E}_{X_{s+1}, \dots, X_{r_k}}[\psi_k(x_1, \dots, x_s, X_{s+1}, \dots, X_{r_k})], \mathbb{E}_{X_{s+1}, \dots, X_{r_l}}[\psi_l(x_1, \dots, x_s, X_{s+1}, \dots, X_{r_l})]).$$

Theorem A.4. Multivariate asymptotic distribution, (Hoeffding, 1948)

Let $(U_n^{(k)})$, $k = 1, \dots, d$ be a sequence of U-statistics, with mean $U^{(k)}$ and kernel $\psi^{(k)}$ of degree r_k . Let $\mathbf{U}_n = (U_n^{(1)}, \dots, U_n^{(d)})^\top$ and $\mathbf{U} = (U^{(1)}, \dots, U^{(d)})^\top$ both d -dimensional vectors. If $\forall k \leq d$, $\text{Var}(\psi^{(k)}(X_1, \dots, X_{r_k})) < \infty$, then

$$\sqrt{n}(\mathbf{U}_n - \mathbf{U}) \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}_{\mathbb{R}^d}(\mathbf{0}, \Sigma),$$

where $\Sigma_{kk}, k \leq d$ (resp. $\Sigma_{kl}, k \leq l \leq d$) is given by Eq. (7) (resp. Eq. (8)).

Theorem A.5. Uniform Law of Large Numbers, (Yeo and Johnson, 2001)

Let $\Theta \subset \mathbb{R}^d$ a compact set and consider the kernel

$$U_n(\theta) = \frac{1}{\binom{n}{r}} \sum_{C_{n,r}} \psi(X_{i_1}, \dots, X_{i_r}; \theta),$$

where $\theta \in \Theta$. Let

$$1 \leq j \leq r, \psi_j(x_1, \dots, x_j; \theta) = \mathbb{E}_{X_{j+1}, \dots, X_r}[\psi(x_1, \dots, x_j, X_{j+1}, \dots, X_r; \theta)],$$

and assume

- (i) there is an integrable and symmetric kernel $g(\cdot)$ such that $\forall \theta \in \Theta$ and $\forall \mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_r) \in \mathbb{R}^{r \times p}$, $|\psi(\mathbf{x}; \theta)| \leq g(\mathbf{x})$,
- (ii) there is a sequence S_M^r of measurable sets such that $\mathbb{P}(\mathbb{R}^{r \times p} - \bigcup_{r=1}^{\infty} S_M^r) = 0$,
- (iii) for each M and $\forall j \leq r$, $\psi_j(x_1, \dots, x_j; \theta)$ is equicontinuous in θ for $(\mathbf{x}_1, \dots, \mathbf{x}_j) \in S_M^j$, where $S_M^r = S_M^j \times S_M^{r-j}$,

then

$$\sup_{\theta \in \Theta} |U_n(\theta) - \mathbb{E}[\psi(\mathbf{X}; \theta)]| \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0.$$

In this Supplementary material, we report some technical results we relied on to derive the asymptotic distributions

B Large sample distribution

We used the *convexity argument* to derive the asymptotic distribution in the Lasso case. (Chernozhukov and Hong, 2004), (Chernozhukov, 2005) use this convexity argument to obtain the asymptotic distribution of quantile regression type estimators. This argument relies on the convexity Lemma, which is a key result to obtain an asymptotic distribution when the objective function is not differentiable. It only requires the lower-semicontinuity and convexity of the empirical criterion. The convexity Lemma, as in (Chernozhukov, 2005), proof of Theorem 4.1, can be stated as follows.

Lemma B.1. Convexity Lemma, (Chernozhukov, 2005)

Suppose

(i) a sequence of convex lower-semicontinuous $\mathbb{F}_T : \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ marginally converges to $\mathbb{F}_\infty : \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ over a dense subset of \mathbb{R}^d ;

(ii) \mathbb{F}_∞ is finite over a nonempty open set $E \subset \mathbb{R}^d$;

(iii) \mathbb{F}_∞ is uniquely minimized at a random vector \mathbf{u}_∞ .

Then

$$\operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^d} \mathbb{F}_n(\mathbf{z}) \xrightarrow{d} \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^d} \mathbb{F}_\infty(\mathbf{z}), \text{ that is } \mathbf{u}_n \xrightarrow{d} \mathbf{u}_\infty.$$

As for the SCAD and MCP, due to the non-convexity of the penalty function, we used Lemma 3 of (Umezu *et al.*, 2018), which generalises Lemma 2 of (Hjort and Pollard, 1993) to the case of convex non-penalized loss functions with non-convex penalties. Lemma 2 of (Hjort and Pollard, 1993) allows for deriving consistency and asymptotic normality of estimators that are defined by minimisation of convex criterion functions.

Lemma B.2. (Umezu *et al.*, 2018)

Suppose that $\mathbb{G}_n(\mathbf{u})$ is a strictly convex random function that is approximated by $\tilde{\mathbb{G}}_n(\mathbf{u})$. Let $\bar{\mathbf{u}}$ be a subvector of \mathbf{u} , and let $\zeta(\mathbf{u})$ and $\eta(\bar{\mathbf{u}})$ be continuous functions such that $\zeta_n(\mathbf{u})$ and $\eta_n(\bar{\mathbf{u}})$ converge to $\zeta(\mathbf{u})$ and $\eta(\bar{\mathbf{u}})$ uniformly over \mathbf{u} and $\bar{\mathbf{u}}$ in any compact set, respectively, and assume that $\zeta(\mathbf{u})$ is convex and $\eta(\mathbf{0}) = 0$. In addition, for

$$\nu_n(\mathbf{u}) = \mathbb{G}_n(\mathbf{u}) + \zeta_n(\mathbf{u}) + \eta_n(\bar{\mathbf{u}}), \text{ and } \tilde{\nu}_n(\mathbf{u}) = \tilde{\mathbb{G}}_n(\mathbf{u}) + \zeta(\mathbf{u}) + \eta(\bar{\mathbf{u}}),$$

let \mathbf{u}_n and $\tilde{\mathbf{u}}_n$ be the argmin of $\nu_n(\mathbf{u})$ and $\tilde{\nu}_n(\mathbf{u})$, respectively, and assume that $\tilde{\mathbf{u}}_n$ is unique and $\tilde{\mathbf{u}}_n = \mathbf{0}$. Then, for any $\epsilon > 0, \delta > 0, \mu > \delta$, there exists $\gamma > 0$ such that

$$\mathbb{P}(\|\mathbf{u}_n - \tilde{\mathbf{u}}_n\| \geq \delta) \leq \mathbb{P}(2\Delta_n(\delta) + \epsilon \geq \Upsilon_n(\delta)) + \mathbb{P}(\|\mathbf{u}_n - \tilde{\mathbf{u}}_n\| \geq \mu) + \mathbb{P}(\|\tilde{\mathbf{u}}_n\| \geq \gamma),$$

where

$$\Delta_n(\delta) = \sup_{\mathbf{u}: \|\mathbf{u} - \tilde{\mathbf{u}}_n\| \leq \delta} |\nu_n(\mathbf{u}) - \tilde{\nu}_n(\mathbf{u})|, \quad \Upsilon_n(\delta) = \inf_{\mathbf{u}: \|\mathbf{u} - \tilde{\mathbf{u}}_n\| = \delta} |\tilde{\nu}_n(\mathbf{u}) - \tilde{\nu}_n(\tilde{\mathbf{u}}_n)|.$$

Finally, for the large sample distribution of the Bridge penalized estimator, we relied on Theorem 2.7 of (Kim and Pollard, 1990).

Theorem B.3. (Kim and Pollard, 1990)

Let $\{\mathbb{F}_n\}$ be a random function into the space of all locally bounded real functions on \mathbb{R}^d , and \mathbf{u}_n random mapps into \mathbb{R}^d such that

- (i) $\mathbb{F}_n \xrightarrow{d} Q$ for a Borel measure Q concentrated on $C_{\max}(\mathbb{R}^d)^1$;
- (ii) $\mathbf{u}_n = O_p(1)$;
- (iii) $\mathbb{F}_n(\mathbf{u}_n) \geq \sup_{\mathbf{u}} \{\mathbb{F}_n(\mathbf{u})\} - \alpha_n$ for random variables (α_n) of order $o_p(1)$.

Then $\mathbf{u}_n \xrightarrow{d} \operatorname{argmax}_{\mathbf{u}} \{\mathbb{F}(\mathbf{u})\}$ for a $\mathbb{F}(\mathbf{u})$ with distribution Q .

C Proofs

Proof of Proposition 4.1

Proof. We first expand the V-statistic based quantity

$$\begin{aligned} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) &= \frac{1}{n^2} \sum_{i,j=1}^n \left((\mathbf{H}_n \mathbf{L} \mathbf{H}_n)_{ij} - \sum_{k=1}^d \theta_k (\mathbf{H}_n \mathbf{L}_X^k \mathbf{H}_n)_{ij} \right)^2 \\ &= \widehat{\text{HSIC}}(Y, Y) - 2 \sum_{k=1}^d \theta_k \widehat{\text{HSIC}}(X^{(k)}, Y) + \sum_{k,l=1}^d \theta_k \theta_l \widehat{\text{HSIC}}(X^{(k)}, X^{(l)}), \end{aligned}$$

¹See their page 195 for a definition of this set

where $\widehat{\text{HSIC}}(Y, Y)$ is the empirical estimator of $\text{HSIC}(Y, Y)$ given in Eq. (2). Then using Theorem 3 of (Song *et al.*, 2012), which enables to write the estimator of $\text{HSIC}(Y, X)$ as a U-statistic, we have $\widehat{\text{HSIC}}(Y, Y)$ can be rewritten as

$$\widehat{\text{HSIC}}(Y, Y) = \frac{1}{(n)_4} \sum_{(i,j,q,r) \in \mathcal{I}_4^n} h_1(i, j, q, r), \quad \text{with } h_1(i, j, q, r) = \sum_{(s,t,u,v)}^{(i,j,q,r)} L_{st}(L_{st} + L_{uv} - 2L_{su}),$$

where the sum represents all ordered quadruples (s, t, u, v) selected without replacement from (i, j, q, r) . In the same manner, denoting $K_{st}^k = \psi(X_s^{(k)}, X_t^{(k)})$, we have for any $k, l \leq d$

$$\begin{aligned} \widehat{\text{HSIC}}(X^{(k)}, Y) &= \frac{1}{(n)_4} \sum_{(i,j,q,r) \in \mathcal{I}_4^n} h_2(i, j, q, r), \quad h_2(i, j, q, r) = \frac{1}{4!} \sum_{(s,t,u,v)}^{(i,j,q,r)} K_{st}^k(L_{st}^Y + L_{uv}^Y - 2L_{su}^Y), \\ \widehat{\text{HSIC}}(X^{(k)}, X^{(l)}) &= \frac{1}{(n)_4} \sum_{(i,j,q,r) \in \mathcal{I}_4^n} h_3(i, j, q, r), \quad h_3(i, j, q, r) = \frac{1}{4!} \sum_{(s,t,u,v)}^{(i,j,q,r)} K_{st}^k(K_{st}^l + K_{uv}^l - 2K_{su}^l). \end{aligned}$$

Thus, plugging these U-statistic based estimators of HSIC, we obtain

$$\begin{aligned} &\mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) \\ &= \frac{1}{(n)_4} \sum_{(i,j,q,r) \in \mathcal{I}_4^n} \frac{1}{4!} \sum_{(s,t,u,v)}^{(i,j,q,r)} \{h_1(i, j, q, r) - 2 \sum_{k=1}^d \theta_k h_2(i, j, q, r) + \sum_{k,l=1}^d \theta_k \theta_l h_3(i, j, q, r)\} \\ &= \frac{1}{(n)_4} \sum_{(i,j,q,r) \in \mathcal{I}_4^n} \frac{1}{4!} \sum_{(s,t,u,v)}^{(i,j,q,r)} \{L_{st}(L_{st} + L_{uv} - 2L_{su}) - 2 \sum_{k=1}^d \theta_k K_{st}^k (L_{st} + L_{uv} - 2L_{su}) \\ &\quad + \sum_{k,l=1}^d \theta_k \theta_l K_{st}^k (K_{st}^l + K_{uv}^l - 2K_{su}^l)\} \\ &= \frac{1}{(n)_4} \sum_{(i,j,q,r) \in \mathcal{I}_4^n} \frac{1}{4!} \sum_{(s,t,u,v)}^{(i,j,q,r)} \{(L_{st} - \sum_{k=1}^d \theta_k K_{st}^k)(L_{st} - \sum_{l=1}^d \theta_l K_{st}^k) \\ &\quad + (L_{st} - \sum_{k=1}^d \theta_k K_{st}^k)(L_{uv} - \sum_{l=1}^d \theta_l K_{uv}^k) + \sum_{k=1}^d \theta_k K_{uv}^k L_{st} - \sum_{k=1}^d \theta_k K_{st}^k L_{uv} \\ &\quad - 2[(L_{st} - \sum_{k=1}^d \theta_k K_{st}^k)(L_{su} - \sum_{l=1}^d \theta_l K_{su}^l) + \sum_{k=1}^d \theta_k K_{su}^k L_{st} - \sum_{k=1}^d \theta_k K_{st}^k L_{su}]\} \\ &:= \frac{1}{(n)_4} \sum_{(i,j,q,r) \in \mathcal{I}_4^n} \ell(\mathbf{Z}_i, \mathbf{Z}_j, \mathbf{Z}_q, \mathbf{Z}_r; \boldsymbol{\theta}). \end{aligned}$$

□

Proof of Proposition 4.2

Proof. Let us denote $\mathbf{z}_1 = (\mathbf{x}_1, \mathbf{y}_1)^\top \in \mathbb{R}^{p+q}$. Then we rewrite

$$\begin{aligned} &\nabla_{\theta_k} \ell(\mathbf{z}_1, \mathbf{Z}_j, \mathbf{Z}_q, \mathbf{Z}_r; \boldsymbol{\theta}_0) \\ &= \frac{1}{4!} \left[\sum_{s=1, (t,u,v)}^{(2,3,4)} w_{1,tuv}(\mathbf{z}_1; \boldsymbol{\theta}_0) + \sum_{t=1, (s,u,v)}^{(2,3,4)} w_{2,suv}(\mathbf{z}_1; \boldsymbol{\theta}_0) + \sum_{u=1, (s,t,v)}^{(2,3,4)} w_{3,stv}(\mathbf{z}_1; \boldsymbol{\theta}_0) + \sum_{v=1, (s,t,u)}^{(2,3,4)} w_{4,stu}(\mathbf{z}_1; \boldsymbol{\theta}_0) \right] \\ &:= \frac{1}{4!} [\Upsilon_1(\boldsymbol{\theta}_0) + \Upsilon_2(\boldsymbol{\theta}_0) + \Upsilon_3(\boldsymbol{\theta}_0) + \Upsilon_4(\boldsymbol{\theta}_0)], \end{aligned}$$

where

$$\begin{aligned}
 & w_{1,tuv}(\mathbf{z}; \boldsymbol{\theta}_0) \\
 &= -2\psi(x_1^{(k)}, X_t^{(k)})(\phi(y_1, Y_t) - \sum_{l=1}^d \theta_{0,l} \psi(x_1^{(l)}, X_t^{(l)})) - \psi(x_1^{(k)}, X_t^{(k)})(\phi(Y_u, Y_v) - \sum_{l=1}^d \theta_{0,l} \psi(X_u^{(l)}, X_v^{(l)})) \\
 &- \psi(X_u^{(k)}, X_v^{(k)})(\phi(y_1, Y_t) - \sum_{l=1}^d \theta_{0,l} \psi(x_1^{(l)}, X_t^{(l)})) + \psi(X_u^{(k)}, X_v^{(k)})\phi(y_1, Y_t) - \psi(x_1^{(k)}, X_t^{(k)})\phi(Y_u, Y_v) \\
 &+ 2\psi(x_1^{(k)}, X_t^{(k)})(\phi(y_1, Y_u) - \sum_{l=1}^d \theta_{0,l} \psi(x_1^{(l)}, X_u^{(l)})) + 2\psi(x_1^{(k)}, X_u^{(k)})(\phi(y_1, Y_t) - \sum_{l=1}^d \theta_{0,l} \psi(x_1^{(l)}, X_t^{(l)})) \\
 &+ \psi(x_1^{(k)}, X_u^{(k)})\phi(y_1, Y_t) - \psi(x_1^{(k)}, X_t^{(k)})\phi(y_1, Y_u),
 \end{aligned}$$

and

$$\begin{aligned}
 & w_{2,suv}(\mathbf{z}; \boldsymbol{\theta}_0) \\
 &= -2\psi(X_s^{(k)}, x_1^{(k)})(\phi(Y_s, y_1) - \sum_{l=1}^d \theta_{0,l} \psi(X_s^{(l)}, x_1^{(l)})) - \psi(X_s^{(k)}, x_1^{(k)})(\phi(Y_u, Y_v) - \sum_{l=1}^d \theta_{0,l} \psi(X_u^{(l)}, X_v^{(l)})) \\
 &- \psi(X_u^{(k)}, X_v^{(k)})(\phi(Y_s, y_1) - \sum_{l=1}^d \theta_{0,l} \psi(X_s^{(l)}, x_1^{(l)})) + \psi(X_u^{(k)}, X_v^{(k)})\phi(Y_s, y_1) - \psi(X_s^{(k)}, x_1^{(k)})\phi(Y_u, Y_v) \\
 &+ 2\psi(X_s^{(k)}, x_1^{(k)})(\phi(Y_s, Y_u) - \sum_{l=1}^d \theta_{0,l} \psi(X_s^{(l)}, X_u^{(l)})) + 2\psi(X_s^{(k)}, X_u^{(k)})(\phi(Y_s, y_1) - \sum_{l=1}^d \theta_{0,l} \psi(X_s^{(l)}, x_1^{(l)})) \\
 &+ \psi(X_s^{(k)}, X_u^{(k)})\phi(Y_s, y_1) - \psi(X_s^{(k)}, x_1^{(k)})\phi(Y_s, Y_u),
 \end{aligned}$$

and

$$\begin{aligned}
 & w_{3,stv}(\mathbf{z}; \boldsymbol{\theta}_0) \\
 &= -2\psi(X_s^{(k)}, X_t^{(k)})(\phi(Y_s, Y_t) - \sum_{l=1}^d \theta_{0,l} \psi(X_s^{(l)}, X_t^{(l)})) - \psi(X_s^{(k)}, X_t^{(k)})(\phi(y_1, Y_v) - \sum_{l=1}^d \theta_{0,l} \psi(x_1^{(l)}, X_v^{(l)})) \\
 &- \psi(x_1^{(k)}, X_v^{(k)})(\phi(Y_s, Y_t) - \sum_{l=1}^d \theta_{0,l} \psi(X_s^{(l)}, X_t^{(l)})) + \psi(x_1^{(k)}, X_v^{(k)})\phi(Y_s, Y_t) - \psi(X_s^{(k)}, X_t^{(k)})\phi(y_1, Y_v) \\
 &+ 2\psi(X_s^{(k)}, X_t^{(k)})(\phi(Y_s, y_1) - \sum_{l=1}^d \theta_{0,l} \psi(X_s^{(l)}, x_1^{(l)})) + 2\psi(X_s^{(k)}, x_1^{(k)})(\phi(Y_s, Y_t) - \sum_{l=1}^d \theta_{0,l} \psi(X_s^{(l)}, X_t^{(l)})) \\
 &+ \psi(X_s^{(k)}, x_1^{(k)})\phi(Y_s, Y_t) - \psi(X_s^{(k)}, X_t^{(k)})\phi(Y_s, y_1),
 \end{aligned}$$

and finally

$$\begin{aligned}
 & w_{4,stu}(\mathbf{z}; \boldsymbol{\theta}_0) \\
 &= -2\psi(X_s^{(k)}, X_t^{(k)})(\phi(Y_s, Y_t) - \sum_{l=1}^d \theta_{0,l} \psi(X_s^{(l)}, X_t^{(l)})) - \psi(X_s^{(k)}, X_t^{(k)})(\phi(Y_u, y_1) - \sum_{l=1}^d \theta_{0,l} \psi(X_u^{(l)}, x_1^{(l)})) \\
 &- \psi(X_u^{(k)}, x_1^{(k)})(\phi(Y_s, Y_t) - \sum_{l=1}^d \theta_{0,l} \psi(x_1^{(l)}, X_t^{(l)})) + \psi(X_u^{(k)}, x_1^{(k)})\phi(Y_s, Y_t) - \psi(X_s^{(k)}, X_t^{(k)})\phi(Y_u, y_1) \\
 &+ 2\psi(X_s^{(k)}, X_t^{(k)})(\phi(Y_s, Y_u) - \sum_{l=1}^d \theta_{0,l} \psi(X_s^{(l)}, X_u^{(l)})) + 2\psi(X_s^{(k)}, X_u^{(k)})(\phi(Y_s, Y_t) - \sum_{l=1}^d \theta_{0,l} \psi(X_s^{(l)}, X_t^{(l)})) \\
 &+ \psi(X_s^{(k)}, X_u^{(k)})\phi(Y_s, Y_t) - \psi(X_s^{(k)}, X_t^{(k)})\phi(Y_s, Y_u).
 \end{aligned}$$

Then taking the expectation with respect to $\mathbf{Z}_j, \mathbf{Z}_q, \mathbf{Z}_r$ and using the independence between \mathbf{Y} and \mathbf{X} , we have

$$\begin{aligned}
 & \mathbb{E}[\Upsilon_1(\boldsymbol{\theta}_0)] \\
 &= -2\mathbb{E}[\psi(x_1^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(y_1, Y_t)] + 2\sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(x_1^{(k)}, X_t^{(k)})\psi(x_1^{(l)}, X_t^{(l)})] - \mathbb{E}[\psi(x_1^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(Y_u, Y_v)] \\
 &+ \sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(x_1^{(k)}, X_t^{(k)})\psi(X_u^{(l)}, X_v^{(l)})] - \mathbb{E}[\psi(X_u^{(k)}, X_v^{(k)})]\mathbb{E}[\phi(y_1, Y_t)] \\
 &+ \sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_u^{(k)}, X_v^{(k)})\psi(x_1^{(l)}, X_t^{(l)})] + \mathbb{E}[\psi(X_u^{(k)}, X_v^{(k)})]\mathbb{E}[\phi(y_1, Y_t)] - \mathbb{E}[\psi(x_1^{(k)}, X_t^{(k)})\phi(Y_u, Y_v)] \\
 &+ 2\mathbb{E}[\psi(x_1^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(y_1, Y_u)] - 2\sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(x_1^{(k)}, X_t^{(k)})\psi(x_1^{(l)}, X_u^{(l)})] \\
 &+ 2\mathbb{E}[\psi(x_1^{(k)}, X_u^{(k)})]\mathbb{E}[\phi(y_1, Y_t)] - 2\sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(x_1^{(k)}, X_u^{(k)})\psi(x_1^{(l)}, X_t^{(l)})] \\
 &+ \mathbb{E}[\psi(x_1^{(k)}, X_u^{(k)})]\mathbb{E}[\phi(y_1, Y_t)] - \mathbb{E}[\psi(x_1^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(y_1, Y_u)],
 \end{aligned}$$

Proceeding the same way for the other expectations, we obtain

$$\begin{aligned}
 & \mathbb{E}[\Upsilon_2(\boldsymbol{\theta}_0)] \\
 &= -2\mathbb{E}[\psi(X_s^{(k)}, x_1^{(k)})]\mathbb{E}[\phi(Y_s, y_1)] + 2\sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_s^{(k)}, x_1^{(k)})\psi(X_s^{(l)}, x_1^{(l)})] - \mathbb{E}[\psi(X_s^{(k)}, x_1^{(k)})]\mathbb{E}[\phi(Y_u, Y_v)] \\
 &+ \sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_s^{(k)}, x_1^{(k)})\psi(X_u^{(l)}, X_v^{(l)})] - \mathbb{E}[\psi(X_u^{(k)}, X_v^{(k)})]\mathbb{E}[\phi(Y_s, y_1)] \\
 &+ \sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_u^{(k)}, X_v^{(k)})\psi(X_s^{(l)}, x_1^{(l)})] + \mathbb{E}[\psi(X_u^{(k)}, X_v^{(k)})]\mathbb{E}[\phi(Y_s, y_1)] - \mathbb{E}[\psi(X_s^{(k)}, x_1^{(k)})]\mathbb{E}[\phi(Y_u, Y_v)] \\
 &+ 2\mathbb{E}[\psi(X_s^{(k)}, x_1^{(k)})]\mathbb{E}[\phi(Y_s, Y_u)] - 2\sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_s^{(k)}, x_1^{(k)})\psi(X_s^{(l)}, X_u^{(l)})] \\
 &+ 2\mathbb{E}[\psi(X_s^{(k)}, X_u^{(k)})]\mathbb{E}[\phi(Y_s, y_1)] - 2\sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_s^{(k)}, X_u^{(k)})\psi(X_s^{(l)}, x_1^{(l)})] \\
 &+ \mathbb{E}[\psi(X_s^{(k)}, X_u^{(k)})]\mathbb{E}[\phi(Y_s, y_1)] - \mathbb{E}[\psi(X_s^{(k)}, x_1^{(k)})\phi(Y_s, Y_u)],
 \end{aligned}$$

and

$$\begin{aligned}
 & \mathbb{E}[\Upsilon_3(\boldsymbol{\theta}_0)] \\
 &= -2\mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(Y_s, Y_t)] + 2\sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})\psi(X_s^{(l)}, X_t^{(l)})] \\
 &- \mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(y_1, Y_v)] + \sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})\psi(x_1^{(l)}, X_v^{(l)})] - \mathbb{E}[\psi(x_1^{(k)}, X_v^{(k)})]\mathbb{E}[\phi(Y_s, Y_t)] \\
 &+ \sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(x_1^{(k)}, X_v^{(k)})\psi(X_s^{(l)}, X_t^{(l)})] + \mathbb{E}[\psi(x_1^{(k)}, X_v^{(k)})]\mathbb{E}[\phi(Y_s, Y_t)] - \mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(y_1, Y_v)] \\
 &+ 2\mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(Y_s, y_1)] - 2\sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})\psi(X_s^{(l)}, x_1^{(l)})] + 2\mathbb{E}[\psi(X_s^{(k)}, x_1^{(k)})]\mathbb{E}[\phi(Y_s, Y_t)] \\
 &- 2\sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_s^{(k)}, x_1^{(k)})\psi(X_s^{(l)}, X_t^{(l)})] + \mathbb{E}[\psi(X_s^{(k)}, x_1^{(k)})]\mathbb{E}[\phi(Y_s, Y_t)] - \mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(Y_s, y_1)],
 \end{aligned}$$

and finally

$$\begin{aligned}
 & \mathbb{E}[\Upsilon_4(\boldsymbol{\theta}_0)] \\
 &= -2\mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(Y_s, Y_t)] + 2\sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})\psi(X_s^{(l)}, X_t^{(l)})] \\
 &- \mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(Y_u, y_1)] + \sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})\psi(X_u^{(l)}, x_1^{(l)})] - \mathbb{E}[\psi(X_u^{(k)}, x_1^{(k)})]\mathbb{E}[\phi(Y_s, Y_t)] \\
 &+ \sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_u^{(k)}, x_1^{(k)})\psi(X_s^{(l)}, X_t^{(l)})] + \mathbb{E}[\psi(X_u^{(k)}, x_1^{(k)})]\mathbb{E}[\phi(Y_s, Y_t)] - \mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(Y_u, y_1)] \\
 &+ 2\mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(Y_s, Y_u)] - 2\sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})\psi(X_s^{(l)}, X_u^{(l)})] + 2\mathbb{E}[\psi(X_s^{(k)}, X_u^{(k)})]\mathbb{E}[\phi(Y_s, Y_t)] \\
 &- 2\sum_{l=1}^d \theta_{0,l} \mathbb{E}[\psi(X_s^{(k)}, X_u^{(k)})\psi(X_s^{(l)}, X_t^{(l)})] + \mathbb{E}[\psi(X_s^{(k)}, X_u^{(k)})]\mathbb{E}[\phi(Y_s, Y_t)] - \mathbb{E}[\psi(X_s^{(k)}, X_t^{(k)})]\mathbb{E}[\phi(Y_s, Y_u)].
 \end{aligned}$$

Thus by symmetry of the kernel, i.i.d. random variables and under $\mathbb{P}_{\mathbf{Y}\mathbf{X}} = \mathbb{P}_{\mathbf{X}}\mathbb{P}_{\mathbf{Y}}$, then we obtain $\forall \mathbf{z}_1 \in \mathbb{R}^{p+q}$ by summing up all these expectations

$$\forall k \leq d, \mathbb{E}_{\mathbf{Z}_j \mathbf{Z}_q \mathbf{Z}_r}[\nabla_{\theta_k} \ell(\mathbf{z}_1, \mathbf{Z}_j, \mathbf{Z}_q, \mathbf{Z}_r; \boldsymbol{\theta}_0)] = 0.$$

□

Proof of Theorem 4.3

Proof. In a first step, we prove the uniform convergence of $\mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \cdot)$ to the limit quantity $\mathbb{L}_\infty^{pen}(\cdot)$ on any compact set $\mathcal{B} \subset \Theta$, idest

$$\sup_{\boldsymbol{\theta} \in \mathcal{B}} |\mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) - \mathbb{L}_\infty^{pen}(\boldsymbol{\theta})| \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0, \quad (9)$$

where $\mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) = \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) + \sum_{k=1}^d \varphi(\frac{\lambda_n}{n}, |\theta_k|)$. We define $\mathcal{C} \subset \Theta$ an open convex set and pick $\boldsymbol{\theta} \in \mathcal{C}$. Then we need to apply the uniform law of large numbers on the non-penalized term, where we rely on Theorem 1 of (Yeo and Johnson, 2001), provided in the technical results (see the Supplementary material). To do so, let $k \leq 4$ and define

$$\ell_k(\mathbf{z}_1, \dots, \mathbf{z}_k; \boldsymbol{\theta}) = \mathbb{E}_{\mathbf{Z}_{k+1}, \dots, \mathbf{Z}_4}[\ell(\mathbf{z}_1, \dots, \mathbf{z}_k, \mathbf{Z}_{k+1}, \dots, \mathbf{Z}_4; \boldsymbol{\theta})].$$

We then need to prove

- (i) there is an integrable and symmetric kernel $g(\cdot)$ such that $\forall \boldsymbol{\theta} \in \Theta$ and $\forall \mathbf{z} \in \mathbb{R}^{(p+q) \times 4}$, $|\ell(\mathbf{z}; \boldsymbol{\theta})| \leq g(\mathbf{z})$.
- (ii) There is a sequence S_M^4 of measurable sets such that $\mathbb{P}(\mathbb{R}^{(p+q) \times 4} - \bigcup_{M=1}^{\infty} S_M^4) = 0$.
- (iii) For each M and $\forall k \leq 4$, $\ell_k(\mathbf{z}_1, \dots, \mathbf{z}_k; \boldsymbol{\theta})$ is equicontinuous in $\boldsymbol{\theta}$ for $\mathbf{z}_1, \dots, \mathbf{z}_k \in S_M^k$, where $S_M^4 = S_M^k \times S_M^{4-k}$

First, $\ell(\mathbf{z}_1, \dots, \mathbf{z}_4; \boldsymbol{\theta})$ is a symmetric kernel that is bounded by assumption 4 and continuous for all $\boldsymbol{\theta}, (\mathbf{z}_1, \dots, \mathbf{z}_4) \in \Theta \times \mathbb{R}^{4(p+q)}$. Let $S_M^4 = \times_{k=1}^4 S_{k,M}$, where $S_{k,M} = \{\mathbf{z}_k : |\mathbf{z}_{k,i}| \leq M\}$, with M a positive integer. Then we have $\mathbb{P}(\mathbb{R}^{4(p+q)} - \bigcup_{M=1}^{\infty} S_M^4) = 0$. For each fixed M , if $(\mathbf{z}_1, \dots, \mathbf{z}_k) \in S_M^k$, then the coefficients of $\ell_k(\mathbf{z}_1, \dots, \mathbf{z}_k; \boldsymbol{\theta})$ are bounded, hence $\ell_k(\mathbf{z}_1, \dots, \mathbf{z}_k; \boldsymbol{\theta})$ is equicontinuous in $\boldsymbol{\theta}$. As a consequence, we obtain by Theorem 1 of (Yeo and Johnson, 2001) the uniform convergence

$$\sup_{\boldsymbol{\theta} \in \mathcal{B}} |\mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) - \mathbb{L}_\infty(\boldsymbol{\theta})| \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0.$$

Thus, using $\lambda_n/n \rightarrow \lambda_0$ and since the parameter is taken over a compact set, we obtain

$$\sup_{\boldsymbol{\theta} \in \mathcal{B}} |\mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) - \mathbb{L}_\infty^{pen}(\boldsymbol{\theta})| \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0.$$

Now we would like

$$\operatorname{argmin} \left\{ \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \cdot) + \sum_{k=1}^d \varphi\left(\frac{\lambda_n}{n}, \cdot\right) \right\} \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \operatorname{argmin} \left\{ \mathbb{L}_\infty(\cdot) + \sum_{k=1}^d \varphi(\lambda_0, \cdot) \right\}.$$

First, $\mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) + \sum_{k=1}^d \varphi\left(\frac{\lambda_n}{n}, |\theta_k|\right) \geq \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta})$, and $\operatorname{argmin}_{\boldsymbol{\theta} \in \Omega} \{\mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta})\} = O_p(1)$ by convexity of the criterion, it thus follows that $\operatorname{argmin}_{\boldsymbol{\theta} \in \Omega} \{\mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) + \sum_{k=1}^d \varphi\left(\frac{\lambda_n}{n}, |\theta_k|\right)\} = O_p(1)$ with probability one. \square

Proof of Theorem 4.4

Proof. Let $\nu_n = n^{-1/2} + \sqrt{\operatorname{card}(\mathcal{A})}A_{1,n}$. We would like to prove that for any $\epsilon > 0$, there exists $C_\epsilon > 0$ such that

$$\mathbb{P}\left(\frac{1}{\nu_n} \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| > C_\epsilon\right) < \epsilon.$$

Following the reasoning of (Fan and Li, 2001) Theorem 1, and denoting $\mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) = \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) + \sum_{k=1}^d \varphi\left(\frac{\lambda_n}{n}, |\theta_k|\right)$, we have

$$\mathbb{P}\left(\frac{1}{\nu_n} \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| > C_\epsilon\right) \leq \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0 + \nu_n \mathbf{u}) \leq \mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)),$$

which implies that there is a local minimum in the ball $\{\boldsymbol{\theta}_0 + \nu_n \mathbf{u}, \|\mathbf{u}\|_2 \leq C_\epsilon\}$ so that the minimum $\hat{\boldsymbol{\theta}}$ satisfies $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| = O_p(\nu_n)$. Now by a Taylor expansion of the penalized loss function, we obtain

$$\begin{aligned} \mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0 + \nu_n \mathbf{u}) - \mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) &= \nu_n \mathbf{u}^\top \nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) + \frac{\nu_n^2}{2} \mathbf{u}^\top \nabla_{\boldsymbol{\theta}\boldsymbol{\theta}^\top}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) \mathbf{u} \\ &+ \sum_{k=1}^d \left\{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \right\}, \end{aligned}$$

since the third derivative vanishes. We want to prove

$$\begin{aligned} \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : & \mathbf{u}^\top \nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) + \frac{\nu_n}{2} \mathbf{u}^\top \mathbb{H} \mathbf{u} + \frac{\nu_n}{2} \mathcal{R}_n(\boldsymbol{\theta}_0) \\ &+ \nu_n^{-1} \left(\sum_{k=1}^d \left\{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \right\} \right) \leq 0) < \epsilon, \end{aligned} \quad (10)$$

where $\mathcal{R}_n(\boldsymbol{\theta}_0) = \mathbf{u}^\top \{\nabla_{\boldsymbol{\theta}\boldsymbol{\theta}^\top}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) - \mathbb{H}\} \mathbf{u}$. First, element-by-element $\forall k \leq d$, the score is

$$\begin{aligned} \nabla_{\theta_k} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) &= (n)_4^{-1} \sum_{(i,j,q,r) \in \mathcal{I}_4^n} \nabla_{\theta} \ell(\mathbf{Z}_i, \mathbf{Z}_j, \mathbf{Z}_q, \mathbf{Z}_r; \boldsymbol{\theta}_0) \\ &:= (n)_4^{-1} \sum_{(i,j,q,r) \in \mathcal{I}_4^n} \frac{1}{24} \sum_{(s,t,u,v)}^{(i,j,q,r)} \left(-2K_{st}^k(L_{st} - \sum_{l=1}^d \theta_{0,l} K_{st}^l) - K_{st}^k(L_{uv} - \sum_{l=1}^d \theta_{0,l} K_{uv}^l) \right. \\ &- K_{uv}^k(L_{st} - \sum_{l=1}^d \theta_{0,l} K_{st}^l) + K_{uv}^k L_{st} - K_{st}^k L_{uv} + 2K_{st}^k(L_{su} - \sum_{l=1}^d \theta_{0,l} K_{su}^l) + 2K_{su}^k(L_{st} - \sum_{l=1}^d \theta_{0,l} K_{st}^l) \\ &\left. + K_{su}^k L_{st} - K_{st}^k L_{su} \right). \end{aligned}$$

By the Central Limit Theorem 8, since the score based symmetric kernel $\nabla_{\theta} \ell(\mathbf{Z}_i, \mathbf{Z}_j, \mathbf{Z}_q, \mathbf{Z}_r; \boldsymbol{\theta}_0)$ is of degree 4 and is non-degenerate by assumption, we have

$$\mathbf{u}^\top \nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) = O_p(n^{-1/2} \mathbf{u}^\top \mathbb{M} \mathbf{u}),$$

with $\mathbb{M} = \mathbb{E}[\nabla_{\theta} \ell(\mathbf{Z}; \boldsymbol{\theta}_0) \nabla_{\theta^\top} \ell(\mathbf{Z}; \boldsymbol{\theta}_0)]$ assumed well defined by assumption 5. Moreover, the Hessian is formed with the non-redundant $d(d+1)/2$ elements $1 \leq k, l \leq d$

$$\nabla_{\theta_k \theta_l}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) = (n)_4^{-1} \sum_{(i,j,q,r) \in \mathcal{I}_4^n} \frac{1}{24} \sum_{(s,t,u,v)}^{(i,j,q,r)} (2K_{st}^k K_{st}^l + K_{st}^k K_{uv}^l + K_{uv}^k K_{st}^l - 2K_{st}^k K_{su}^l - 2K_{su}^k K_{st}^l).$$

Thus by the law of large numbers for U-statistics (see e.g. Theorem 3 of (Lee, 1990), subsection 3.4.2)

$$\nabla_{\boldsymbol{\theta} \boldsymbol{\theta}^\top}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \mathbb{H},$$

with $\mathbb{H} = \mathbb{E}[\nabla_{\boldsymbol{\theta} \boldsymbol{\theta}^\top}^2 \ell(\mathbf{Z}; \boldsymbol{\theta}_0)]$ so that $\mathcal{R}_n(\boldsymbol{\theta}_0) = o_p(1)$.

We now focus on the penalty terms. First, we show that the penalty functions satisfy assumption 5. The Lasso satisfies

$$\forall k \in \mathcal{A}, \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) = \frac{\lambda_n}{n} \operatorname{sgn}(\theta_{0,k}), \quad \nabla_{\theta_k \theta_k}^2 \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) = 0.$$

For $0 < q < 1$, the Bridge satisfies

$$\begin{aligned} \forall k \in \mathcal{A}, \quad \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) &= \frac{\lambda_n}{n} q |\theta_{0,k}|^{q-1} \operatorname{sgn}(\theta_{0,k}), \\ \nabla_{\theta_k \theta_k}^2 \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) &= \frac{\lambda_n}{n} q(q-1) |\theta_{0,k}|^{q-2} \operatorname{sgn}(\theta_{0,k}), \end{aligned}$$

thus the second order derivative of the Bridge converges to 0 when $\lambda_n = o(n)$. As for the SCAD, we have

$$\forall k \in \mathcal{A}, \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) = \frac{\lambda_n}{n} \left(\mathbf{1}_{\{|\theta_{0,k}| \leq \frac{\lambda_n}{n}\}} + \frac{(b_{\text{scad}} \frac{\lambda_n}{n} - |\theta_{0,k}|)_+}{(b_{\text{scad}} - 1) \frac{\lambda_n}{n}} \mathbf{1}_{\{|\theta_{0,k}| > \frac{\lambda_n}{n}\}} \right).$$

As a consequence, the SCAD penalty is twice continuously differentiable for $\frac{\lambda_n}{n} < |\theta_{0,k}|$, which implies that $\forall k \in \mathcal{A}, \nabla_{\theta_k \theta_k}^2 \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) = o(1)$. In the MCP case, we have

$$\forall k \in \mathcal{A}, \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) = \left(\frac{\lambda_n}{n} - \frac{|\theta_{0,k}|}{b_{\text{mcp}}} \right) \operatorname{sgn}(\theta_{0,k}) \mathbf{1}_{\{|\theta_{0,k}| \leq b_{\text{mcp}} \frac{\lambda_n}{n}\}}.$$

Under $\lambda_n = o(n)$, we straightforwardly obtain $\nabla_{\theta_k \theta_k}^2 \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) = o(1)$ for non-zero components. Now for any $k \in \mathcal{A} \subset \{1, \dots, d\}$, and since the penalties are coordinate-separable, we have

$$\varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) = \nu_n u_k \operatorname{sgn}(\theta_{0,k}) \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) + \frac{\nu_n^2}{2} u_k^2 \nabla_{\theta_k \theta_k}^2 \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) (1 + o(1)).$$

Hence, using $\varphi\left(\frac{\lambda_n}{n}, 0\right) = 0$, we have

$$\left| \sum_{k \in \mathcal{A}} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \right| \leq \nu_n \|\mathbf{u}\|_1 A_{1,n} + \frac{\nu_n^2}{2} \|\mathbf{u}\|_2^2 A_{2,n} (1 + o(1)).$$

Using $\|\mathbf{u}\|_1 \leq \sqrt{\operatorname{card}(\mathcal{A})} \|\mathbf{u}\|_2$, and under assumption 5, the third derivative being dominated, we obtain

$$\left| \sum_{k \in \mathcal{A}} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \right| \leq \nu_n \sqrt{\operatorname{card}(\mathcal{A})} \|\mathbf{u}\|_2 A_{1,n} + \frac{\nu_n^2}{2} \|\mathbf{u}\|_2^2 A_{2,n}.$$

Then, denoting $\delta_n = \lambda_{\min}(\mathbb{H}) C_\epsilon^2 \nu_n$, and using $\frac{\nu_n}{2} \mathbf{u}^\top \mathbb{E}[\nabla_{\boldsymbol{\theta} \boldsymbol{\theta}^\top}^2 \ell(\mathbf{Z}; \boldsymbol{\theta}_0)] \mathbf{u} \geq \delta_n$, we deduce that Eq. (10) can be bounded as

$$\begin{aligned} \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \mathbf{u}^\top \nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) + \frac{\nu_n}{2} \mathbf{u}^\top \mathbb{E}[\nabla_{\boldsymbol{\theta} \boldsymbol{\theta}^\top}^2 \ell(\mathbf{Z}; \boldsymbol{\theta}_0)] \mathbf{u} + \frac{\nu_n}{2} \mathcal{R}_n(\boldsymbol{\theta}_0) \\ + \nu_n^{-1} \left(\sum_{k=1}^d \{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \} \right) \leq 0) \\ \leq \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \mathbf{u}^\top |\nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)| > \delta_n/6) + \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon | \frac{\nu_n}{2} \mathcal{R}_n(\boldsymbol{\theta}_0)| > \delta_n/6) \\ \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \left| \sum_{k=1}^d \{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \} \right| > \nu_n \delta_n/6). \end{aligned}$$

We have for n and C_ϵ sufficiently large enough

$$\begin{aligned} & \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \left| \sum_{k=1}^d \left\{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \right\} \right| > \nu_n \delta_n / 6) \\ & \leq \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \nu_n \sqrt{\text{card}(\mathcal{A})} \|\mathbf{u}\|_2 A_{1,n} + \frac{\nu_n^2}{2} \|\mathbf{u}\|_2^2 A_{2,n} > \nu_n \delta_n / 6) < \epsilon / 3. \end{aligned}$$

Moreover, if $\nu_n = n^{-1/2} + \sqrt{\text{card}(\mathcal{A})} A_{1,n}$, for C_ϵ large enough

$$\mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \mathbf{u}^\top |\nabla_{\theta} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)| > \delta_n / 6) \leq \frac{C_\epsilon^2 C_{st}}{n \delta_n^2} \leq \frac{C_{st}}{C_\epsilon^4} < \epsilon / 3,$$

where $C_{st} > 0$ is a generic constant. Finally, using $\mathcal{R}_n(\boldsymbol{\theta}_0) = o_p(1)$, we obtain

$$\begin{aligned} & \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \mathbf{u}^\top |\nabla_{\theta} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)| > \delta_n / 6) + \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \frac{\nu_n}{2} \mathcal{R}_n(\boldsymbol{\theta}_0) > \delta_n / 6) \\ & \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \left| \sum_{k=1}^d \left\{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \right\} \right| > \nu_n \delta_n / 6) \\ & \leq \frac{C_{st}}{C_\epsilon^4} + 2\epsilon / 3 \leq \epsilon, \end{aligned}$$

for n and C_ϵ large enough. We deduce $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| = O_p(\nu_n)$. \square

Proof of Theorem 4.5

Proof. Let $\mathbf{u} \in \mathbb{R}^d$ such that $\boldsymbol{\theta} = \boldsymbol{\theta}_0 + \mathbf{u} / \sqrt{n}$, and the empirical criterion $\mathbb{F}_n(\mathbf{u}) = n \{ \mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) - \mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) \}$. Note that $\mathbb{F}_n(\mathbf{u})$ is minimized at $\hat{u}_n = n^{1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)$ because $\hat{\boldsymbol{\theta}}$ minimizes $\mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta})$. Thus $\hat{u}_n = \underset{\mathbf{u} \in \mathbb{R}^d}{\text{argmin}} \{ \mathbb{F}_n(\mathbf{u}) \}$.

We first establish the finite distributional convergence of $\mathbb{F}_n(\cdot)$ to $\mathbb{F}_\infty(\cdot)$. Then we separate the proof depending on the penalty function under consideration. We have the expansion

$$\mathbb{F}_n(\mathbf{u}) = \sqrt{n} \nabla_{\theta} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) \mathbf{u} + \frac{1}{2} \mathbf{u}^\top \nabla_{\theta\theta}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) \mathbf{u} + n \left(\sum_{k=1}^d \left\{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + u_k / \sqrt{n}|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \right\} \right).$$

By theorem A, section 5.5.1. of (Serfling, 1980) for the Central Limit Theorem and by Theorem 3, section 3.4.2 of (Lee, 1990) for the Law of Large Numbers for U-statistics

$$\sqrt{n} \nabla_{\theta} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}_{\mathbb{R}^d}(\mathbf{0}, \mathbb{M}), \quad \nabla_{\theta\theta}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \mathbb{H}.$$

As for the penalty terms, in the MCP and SCAD cases, we proceed as follows. We have

$$n \left(\sum_{k=1}^d \left\{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + u_k / \sqrt{n}|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \right\} \right) = \zeta_n(\mathbf{u}) + \eta_n(\mathbf{u}),$$

where using the coordinate-separability property of the penalties

$$\zeta_n(\mathbf{u}) = n \sum_{k \in \mathcal{A}} \left\{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + u_k / \sqrt{n}|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \right\}, \quad \eta_n(\mathbf{u}) = n \sum_{k \in \mathcal{A}^c} \left\{ \varphi\left(\frac{\lambda_n}{n}, u_k / \sqrt{n}\right) \right\}.$$

Then we have by a Taylor expansion for the indices $k \in \mathcal{A}$

$$\zeta_n(\mathbf{u}) = \sum_{k \in \mathcal{A}} \sqrt{n} \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) u_k \text{sgn}(\theta_{0,k}) + \nabla_{\theta_k \theta_k}^2 \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) u_k^2 / 2 (1 + o(1)).$$

Under assumption 5, we have $A_{2,n} = o(1)$. We then need to treat the first order term. For both SCAD and MCP, since their derivatives respectively vanish outside $[-b_{\text{scad}} \frac{\lambda_n}{n}, b_{\text{scad}} \frac{\lambda_n}{n}]$, $[-b_{\text{mcp}} \frac{\lambda_n}{n}, b_{\text{mcp}} \frac{\lambda_n}{n}]$ we have

$$\begin{aligned} \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) &= \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \mathbf{1}_{\{|\theta_{0,k}| \leq b_{\text{scad}} \frac{\lambda_n}{n}\}}, \\ \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) &= \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \mathbf{1}_{\{|\theta_{0,k}| \leq b_{\text{mcp}} \frac{\lambda_n}{n}\}}, \end{aligned}$$

which implies for any $\epsilon > 0$ and $i \in \mathcal{A}$ that

$$\begin{aligned} \mathbb{P}(\sqrt{n}\nabla_{\theta_k}\varphi(\frac{\lambda_n}{n}, |\theta_{0,k}|)\mathbf{1}_{\{|\theta_{0,k}| \leq b_{\text{scad}} \frac{\lambda_n}{n}\}} > \epsilon) &\leq \mathbb{P}(|\theta_{0,k}| \leq b_{\text{scad}} \frac{\lambda_n}{n}) \rightarrow 0, \\ \mathbb{P}(\sqrt{n}\nabla_{\theta_k}\varphi(\frac{\lambda_n}{n}, |\theta_{0,k}|)\mathbf{1}_{\{|\theta_{0,k}| \leq b_{\text{mcp}} \frac{\lambda_n}{n}\}} > \epsilon) &\leq \mathbb{P}(|\theta_{0,k}| \leq b_{\text{mcp}} \frac{\lambda_n}{n}) \rightarrow 0. \end{aligned}$$

As a consequence, $\sqrt{n}\nabla_{\theta_k}\varphi(\frac{\lambda_n}{n}, |\theta_{0,k}|) = o_p(1)$. This is a direct consequence of the unbiasedness property when regularizing large coefficients. Thus $\zeta_n(\mathbf{u}) \rightarrow 0$ as $n \rightarrow \infty$. As for $k \in \mathcal{A}^c$, we have

$$\eta_n(\mathbf{u}) = \sum_{k \in \mathcal{A}^c} \sqrt{n}(\nabla_{\theta_k}\varphi(\frac{\lambda_n}{n}, |\theta_k|))_{\theta_k=0} |u_k| + \frac{1}{2}(\nabla_{\theta_k}^2\varphi(\frac{\lambda_n}{n}, |\theta_k|))_{\theta_k=0} u_k^2(1 + o(1)),$$

Based on the assumption that $\lim_{x \rightarrow 0^+} \nabla_x \varphi(\frac{\lambda_n}{n}, x) = \frac{\lambda_n}{n}$ and $\lambda_n = O(\sqrt{n})$, we deduce

$$\eta_n(\mathbf{u}) \rightarrow \lambda_0 \sum_{k \in \mathcal{A}^c} |u_k|.$$

As a consequence, by Lemma B.2, where in the latter we take $\mathbb{G}_n(\mathbf{u}) = \tilde{\mathbb{G}}_n(\mathbf{u}) = \sqrt{n}\nabla_{\theta}\mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)\mathbf{u} + \frac{1}{2}\mathbf{u}^\top \nabla_{\boldsymbol{\theta}\boldsymbol{\theta}^\top}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)\mathbf{u}$, we obtain $\underset{\mathbf{u}}{\operatorname{argmin}} \{\mathbb{F}_n\} \xrightarrow{d} \underset{\mathbf{u}}{\operatorname{argmin}} \{\mathbb{F}_\infty\}$.

For the Bridge estimator, we have for any index $i \in \mathcal{A} \cup \mathcal{A}^c$ and under the rate $\lambda_n/n^{q/2} \rightarrow \lambda_0$, then

$$\begin{aligned} n\left(\sum_{k=1}^d \left\{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + u_k/\sqrt{n}|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \right\}\right) \\ = \lambda_n \sum_{k=1}^d \left\{ |\theta_{0,k} + u_k/\sqrt{n}|^q - |\theta_{0,k}|^q \right\} \rightarrow \lambda_0 \sum_{k=1}^d |u_k|^q \mathbf{1}_{\theta_{0,k}=0}. \end{aligned}$$

Now we need to prove that $\underset{\mathbf{u}}{\operatorname{argmin}} \{\mathbb{F}_n\} \xrightarrow{d} \underset{\mathbf{u}}{\operatorname{argmin}} \{\mathbb{F}_\infty\}$ for any \mathbf{u} and n sufficiently large. To do so, we use Theorem B.3 of (Kim and Pollard, 1990) (see Supplementary material B) and show $\underset{\mathbf{u}}{\operatorname{argmin}} \{\mathbb{F}_n\} = O_p(1)$. We first have the expansion

$$\begin{aligned} \mathbb{F}_n(\mathbf{u}) &= n\{\mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0 + \mathbf{u}/\sqrt{n}) - \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)\} + \lambda_n \sum_{k=1}^d \{|\theta_{0,k} + u_k/\sqrt{n}|^q - |\theta_{0,k}|^q\} \\ &\geq n\{\mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0 + \mathbf{u}/\sqrt{n}) - \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)\} - \lambda_n \sum_{k=1}^d |u_k/\sqrt{n}|^q \\ &\geq n\{\mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0 + \mathbf{u}/\sqrt{n}) - \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)\} - (\lambda_0 + \epsilon) \sum_{k=1}^d |u_k/\sqrt{n}|^q := \tilde{\mathbb{F}}_n(\mathbf{u}), \end{aligned}$$

where, following the proof of Theorem 3 of (Knight and Fu, 2000), ϵ is such that $\lambda_n/n^{q/2} \leq \lambda_0 + \epsilon$. Then, expanding $n\{\mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0 + \mathbf{u}/\sqrt{n}) - \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)\}$ in $\tilde{\mathbb{F}}_n(\mathbf{u})$, we have

$$\tilde{\mathbb{F}}_n(\mathbf{u}) = \sqrt{n}\nabla_{\theta}\mathbb{L}_n(\boldsymbol{\theta})\mathbf{u} + \frac{1}{2}\mathbf{u}^\top \nabla_{\boldsymbol{\theta}\boldsymbol{\theta}^\top}^2 \mathbb{L}_n(\boldsymbol{\theta})\mathbf{u} - (\lambda_0 + \epsilon) \sum_{k=1}^d |u_k/\sqrt{n}|^q.$$

The second term, which is quadratic in \mathbf{u} , dominates the term with $|u_i|^q$, which is that the quadratic term grows faster $|u_i|^q$. Hence $\underset{\mathbf{u}}{\operatorname{argmin}} \{\tilde{\mathbb{F}}_n\} = O_p(1)$, which in turns implies $\underset{\mathbf{u}}{\operatorname{argmin}} \{\mathbb{F}_n\} = O_p(1)$. We then obtain for the Bridge

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) = \underset{\mathbf{u}}{\operatorname{argmin}} \{\mathbb{F}_n\} \xrightarrow{d} \underset{\mathbf{u}}{\operatorname{argmin}} \{\mathbb{F}_\infty\}.$$

Finally, for the Lasso estimator, we have

$$\begin{aligned} n\left(\sum_{k=1}^d \left\{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + u_k/\sqrt{n}|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \right\}\right) \\ = \lambda_n \sum_{k=1}^d \left\{ |\theta_{0,k} + u_k/\sqrt{n}| - |\theta_{0,k}| \right\} \rightarrow \lambda_0 \sum_{k=1}^d (u_k \operatorname{sgn}(\theta_{0,k}) \mathbf{1}_{\theta_{0,k} \neq 0} + |u_k| \mathbf{1}_{\theta_{0,k} = 0}). \end{aligned}$$

We thus proved that $\mathbb{F}_n(\mathbf{u}) \xrightarrow{d} \mathbb{F}_\infty(\mathbf{u})$, for a fixed \mathbf{u} . Let us observe that

$$\mathbf{u}_n^* = \underset{\mathbf{u}}{\operatorname{argmin}} \{\mathbb{F}_n(\mathbf{u})\},$$

and $\mathbb{F}_n(\cdot)$ admits as a minimizer $\mathbf{u}_n^* = \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)$. As \mathbb{F}_n is convex and \mathbb{F}_∞ is continuous, convex and has a unique minimum by assumption 5, then by the convexity Lemma B.1, we obtain

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) = \underset{\mathbf{u}}{\operatorname{argmin}} \{\mathbb{F}_n\} \xrightarrow{d} \underset{\mathbf{u}}{\operatorname{argmin}} \{\mathbb{F}_\infty\}.$$

□

Proof of Theorem 4.6

Proof. Let us define $\boldsymbol{\theta} = (\boldsymbol{\theta}_{\mathcal{A}}^\top, \boldsymbol{\theta}_{\mathcal{A}^c}^\top)^\top$. To prove the support recovery consistency, we show with probability tending to one when $n \rightarrow \infty$, under $\|\boldsymbol{\theta}_{\mathcal{A}} - \boldsymbol{\theta}_{0,\mathcal{A}}\| = O_p(n^{-1/2})$ and suitable regularization rates depending on the penalty, that

$$\mathbb{L}_n^{\text{pen}}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_{\mathcal{A}}, \mathbf{0}_{\mathcal{A}^c}) = \min_{\|\boldsymbol{\theta}_{\mathcal{A}^c}\| \leq Cn^{-1/2}} \{\mathbb{L}_n^{\text{pen}}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_{\mathcal{A}}, \boldsymbol{\theta}_{\mathcal{A}^c})\}. \quad (11)$$

To prove Eq. (11), for any \sqrt{n} -consistent $\boldsymbol{\theta}_{\mathcal{A}}$, we show that over the set $\{k \in \mathcal{A}^c, \theta_k : |\theta_k| \leq n^{-1/2}C\}$ for $C > 0$

$$\begin{aligned} \nabla_{\theta_k} \mathbb{L}_n^{\text{pen}}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) &> 0 && \text{when } 0 < \theta_k < n^{-1/2}C, \\ \nabla_{\theta_k} \mathbb{L}_n^{\text{pen}}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) &< 0 && \text{when } -n^{-1/2}C < \theta_k < 0, \end{aligned} \quad (12)$$

with probability converging to 1. For any index $k \in \mathcal{A}^c$, by a Taylor expansion around the true parameter, we have

$$\begin{aligned} \nabla_{\theta_k} \mathbb{L}_n^{\text{pen}}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) &= \nabla_{\theta_k} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) + \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_k|\right) \operatorname{sgn}(\theta_k) \\ &= \nabla_{\theta_k} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) + \nabla_{\theta_k \theta_k}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) (\theta_k - \theta_{0,k}) + \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_k|\right) \operatorname{sgn}(\theta_k). \end{aligned}$$

By the central limit theorem and the law of large numbers, we have

$$\sqrt{n} \nabla_{\theta_k} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) = O_p(1), \quad \nabla_{\boldsymbol{\theta} \boldsymbol{\theta}^\top}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \mathbb{E}[\nabla_{\boldsymbol{\theta} \boldsymbol{\theta}^\top}^2 \ell(\mathbf{Z}; \boldsymbol{\theta}_0)].$$

We thus obtain for the SCAD and MCP penalties

$$\begin{aligned} \nabla_{\theta_k} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \hat{\boldsymbol{\theta}}) &= O_p(n^{-1/2}) + \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_k|\right) \operatorname{sgn}(\theta_k) \\ &= \frac{\lambda_n}{n} \left\{ \frac{n}{\lambda_n} \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_k|\right) \operatorname{sgn}(\theta_k) + O_p\left(\frac{\sqrt{n}}{\lambda_n}\right) \right\}. \end{aligned}$$

As a consequence, under the condition $\lim_{n \rightarrow \infty} \liminf_{x \rightarrow 0^+} \frac{n}{\lambda_n} \nabla_x \varphi\left(\frac{\lambda_n}{n}, x\right) > 0$ and if the regularization parameter satisfies $\frac{\lambda_n}{n^{1/2}} \rightarrow \infty$, we deduce that the sign of the gradient entirely depends on the sign of θ_k . This this proves Eq. (12).

For the Bridge penalty, following the same reasoning as in the SCAD and MCP, the non-penalized terms are of order $O_p(n^{-1/2})$. As for the penalty, we have

$$\nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_k|\right) = \frac{\lambda_n}{n} q |\theta_k|^{q-1} \operatorname{sgn}(\theta_k) = \frac{\lambda_n}{n^{(q+1)/2}} q |n^{1/2} \theta_k|^{q-1} \operatorname{sgn}(\theta_k).$$

As a consequence, we obtain

$$\nabla_{\theta_k} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) = \frac{\lambda_n}{n^{(q+1)/2}} \{q |n^{1/2} \theta_k|^{q-1} \operatorname{sgn}(\theta_k) + O_p\left(\frac{n^{q/2}}{\lambda_n}\right)\}.$$

Thus, under the assumption that $\lambda_n/n^{q/2} \rightarrow \infty$, the sign of θ_k also entirely determines the sign of the gradient.

We now turn to the asymptotic distribution. We proved that $\hat{\boldsymbol{\theta}}_{\mathcal{A}^c}$ degenerates at $\mathbf{0}_{\mathcal{A}^c}$ with probability approaching one. Now by a Taylor expansion around $\theta_{0,k}$, for $k \in \mathcal{A}$, we have

$$\begin{aligned} & \nabla_{\theta_k} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \hat{\boldsymbol{\theta}}) + \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\hat{\boldsymbol{\theta}}_k| \right) \text{sgn}(\hat{\boldsymbol{\theta}}_k) \\ &= \nabla_{\theta_k} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) + \sum_{j \in \mathcal{A}} \nabla_{\theta_k \theta_j}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) (\hat{\theta}_j - \theta_{0,j}) + \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}| \right) \text{sgn}(\theta_{0,k}) \\ &+ \nabla_{\theta_k \theta_k}^2 \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}| \right) (\hat{\theta}_k - \theta_{0,k}) (1 + o(1)). \end{aligned}$$

Then inverting this relationship and multiplying by \sqrt{n} , we obtain in vector form with respect to the elements in \mathcal{A}

$$\begin{aligned} & \sqrt{n} \nabla_{\mathcal{A}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) \\ &= \left(-\nabla_{\mathcal{A}\mathcal{A}}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) - \mathbf{S}_{n,\mathcal{A}\mathcal{A}} \right) \sqrt{n} \{(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)_{\mathcal{A}}\} + \left(\nabla_{\mathcal{A}\mathcal{A}}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) + \mathbf{S}_{n,\mathcal{A}\mathcal{A}} \right)^{-1} \mathbf{b}_{n,\mathcal{A}}, \end{aligned}$$

where

$$\begin{aligned} \mathbf{b}_{n,\mathcal{A}} &= \left(\nabla_{\theta_1} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,1}| \right) \text{sgn}(\theta_{0,1}), \dots, \nabla_{\theta_{k_0}} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k_0}| \right) \text{sgn}(\theta_{0,k_0}) \right)^\top, \\ \mathbf{S}_{n,\mathcal{A}\mathcal{A}} &= \text{diag}\left(\nabla_{\theta_k \theta_k}^2 \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}| \right), k = 1, \dots, k_0\right). \end{aligned}$$

We thus deduce that by the central limit theorem for U-statistics and the Slutsky theorem

$$\left(-\nabla_{\mathcal{A}\mathcal{A}}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) - \mathbf{S}_{n,\mathcal{A}\mathcal{A}} \right) \sqrt{n} \{(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)_{\mathcal{A}}\} + \left(\nabla_{\mathcal{A}\mathcal{A}}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) + \mathbf{S}_{n,\mathcal{A}\mathcal{A}} \right)^{-1} \mathbf{b}_{n,\mathcal{A}} \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}_{\mathbb{R}^{k_0}}(\mathbf{0}, \mathbb{M}),$$

with

$$\mathbb{M} = \mathbb{E}[\nabla_{\boldsymbol{\theta}} \ell(\mathbf{Z}; \boldsymbol{\theta}_0) \nabla_{\boldsymbol{\theta}^\top} \ell(\mathbf{Z}; \boldsymbol{\theta}_0)].$$

□

Proof of Theorem 4.7

The true sparse support now depends on the sample size so that the sparsity assumptions is defined as follows:

Assumption 7. *Sparsity assumption:* $|\mathcal{A}_n| = k_{0,n} < d_n$ with $\mathcal{A}_n = \{1 \leq k \leq d_n : \theta_{0,k} \neq 0\}$.

We assume the following conditions on the penalty functions.

Assumption 8. $\varphi\left(\frac{\lambda_n}{n}, |\cdot|\right)$ is twice continuously differentiable except at the origin. We define

$$A_{1,n} = \max_{1 \leq k \leq d_n} |\nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}| \right)|, \quad A_{2,n} = \max_{1 \leq k \leq d_n} |\nabla_{\theta_k \theta_k}^2 \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}| \right)|,$$

so that $A_{1,n} = O(n^{-1/2})$ and $A_{2,n} \rightarrow 0$.

All quantities now depend on d_n , hence on n and should be indexed by n . We denote $\mathbb{H}_n := \mathbb{H}_n(\boldsymbol{\theta}_0) = \mathbb{E}[\nabla^2 \boldsymbol{\theta} \boldsymbol{\theta}^\top \ell(\mathbf{Z}; \boldsymbol{\theta}_0)]$ and $\mathbb{M}_n := \mathbb{M}_n(\boldsymbol{\theta}_0) = \mathbb{E}[\nabla_{\boldsymbol{\theta}} \ell(\mathbf{Z}; \boldsymbol{\theta}_0) \nabla_{\boldsymbol{\theta}^\top} \ell(\mathbf{Z}; \boldsymbol{\theta}_0)]$. To make the reading easier, we do not index other quantities by n , which will be implicit. We now assume the following conditions on the loss function.

Assumption 9. \mathbb{H}_n and \mathbb{M}_n exist. \mathbb{H}_n is positive-definite and there exist b_1, b_2 with $0 < b_1 < b_2 < \infty$ and c_1, c_2 with $0 < c_1 < c_2 < \infty$ such that, for all n ,

$$b_1 < \lambda_{\min}(\mathbb{M}_n) < \lambda_{\max}(\mathbb{M}_n) < b_2, \quad c_1 < \lambda_{\min}(\mathbb{H}_n) < \lambda_{\max}(\mathbb{H}_n) < c_2,$$

where $\lambda_{\min}(K)$ (resp. $\lambda_{\max}(K)$) is the minimum (resp. maximum) eigenvalue of any positive-definite square matrix K .

Assumption 10. $\mathbb{E}[\{\nabla_{\boldsymbol{\theta}} \ell(\mathbf{Z}; \boldsymbol{\theta}_0) \nabla_{\boldsymbol{\theta}^\top} \ell(\mathbf{Z}; \boldsymbol{\theta}_0)\}^2] < \infty$ and $\mathbb{E}[\{\nabla_{\boldsymbol{\theta}_k \boldsymbol{\theta}_k}^2 \ell(\mathbf{Z}; \boldsymbol{\theta}_0)\}^2] < \infty$ for every d_n (and then of n).

Proof. We proceed as in the proof of Theorem 4.4. We denote $\nu_n = \sqrt{d_n}\{n^{-1/2} + A_{1,n}\}$ and we would like to prove that, for any $\epsilon > 0$, there exists $C_\epsilon > 0$ such that

$$\mathbb{P}(\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| > C_\epsilon) < \epsilon.$$

Now, following the same reasoning as in the proof of Theorem 4.4, denoting $\mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) = \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}) + \sum_{k=1}^{d_n} \varphi(\frac{\lambda_n}{n}, |\theta_k|)$, we have

$$\mathbb{P}(\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| > C_\epsilon) \leq \mathbb{P}(\exists \mathbf{u} \in \mathbb{R}^{d_n}, \|\mathbf{u}\|_2 = C_\epsilon : \mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0 + \nu_n \mathbf{u}) \leq \mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)),$$

which implies that there is a local minimum in the ball $\{\boldsymbol{\theta}_0 + \nu_n \mathbf{u}, \|\mathbf{u}\|_2 \leq C_\epsilon\}$ so that the minimum satisfies $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\|_2 = O_p(\nu_n)$. Now by a Taylor expansion of the penalized loss function

$$\begin{aligned} \mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0 + \nu_n \mathbf{u}) - \mathbb{L}_n^{pen}(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) &= \nu_n \mathbf{u}^\top \nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) + \frac{\nu_n^2}{2} \mathbf{u}^\top \nabla_{\boldsymbol{\theta}\boldsymbol{\theta}^\top}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) \mathbf{u} \\ &+ \sum_{k=1}^{d_n} \{\varphi(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|) - \varphi(\frac{\lambda_n}{n}, |\theta_{0,k}|\}\}, \end{aligned}$$

since the third derivative vanishes. We want to prove

$$\begin{aligned} \mathbb{P}(\exists \mathbf{u} \in \mathbb{R}^{d_n}, \|\mathbf{u}\|_2 = C_\epsilon : \mathbf{u}^\top \nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) + \frac{\nu_n}{2} \mathbf{u}^\top \mathbb{H} \mathbf{u} + \frac{\nu_n}{2} \mathcal{R}_n(\boldsymbol{\theta}_0) \\ + \nu_n^{-1} \left(\sum_{k=1}^{d_n} \{\varphi(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|) - \varphi(\frac{\lambda_n}{n}, |\theta_{0,k}|\}\right) \leq 0) < \epsilon, \end{aligned} \quad (13)$$

where $\mathcal{R}_n(\boldsymbol{\theta}_0) = \mathbf{u}^\top \{\nabla_{\boldsymbol{\theta}\boldsymbol{\theta}^\top}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) - \mathbb{H}\} \mathbf{u}$. First, for $a > 0$ and the Markov inequality, we have for the score term

$$\begin{aligned} &\mathbb{P}\left(\sup_{\mathbf{u}: \|\mathbf{u}\|_2 = C_\epsilon} |\mathbf{u}^\top \nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)| > a\right) \\ &\leq \mathbb{P}\left(\sup_{\mathbf{u}: \|\mathbf{u}\|_2 = C_\epsilon} \|\mathbf{u}\|_2 \|\nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)\|_2 > a\right) \\ &\leq \mathbb{P}\left(\|\nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)\|_2 > \frac{a}{C_\epsilon}\right) \\ &\leq \left(\frac{C_\epsilon}{a}\right)^2 \mathbb{E}[\|\nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)\|_2^2] \\ &\leq \left(\frac{C_\epsilon}{a}\right)^2 \sum_{k=1}^{d_n} \mathbb{E}[(\nabla_{\theta_k} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0))^2] \\ &= \left(\frac{C_\epsilon}{a}\right)^2 (n)_4^{-2} \sum_{(i,j,q,r) \in I_4^n} \sum_{(i',j',q',r') \in I_4^n} \sum_{k=1}^{d_n} \mathbb{E}[\{\nabla_{\boldsymbol{\theta}_k} \ell(\mathbf{Z}_i, \mathbf{Z}_j, \mathbf{Z}_q, \mathbf{Z}_r; \boldsymbol{\theta}_0) \nabla_{\boldsymbol{\theta}_k} \ell(\mathbf{Z}_{i'}, \mathbf{Z}_{j'}, \mathbf{Z}_{q'}, \mathbf{Z}_{r'}; \boldsymbol{\theta}_0)\}^2]. \end{aligned}$$

Thus, by assumption 10, we deduce

$$\mathbb{P}\left(\sup_{\mathbf{u}: \|\mathbf{u}\|_2 = C_\epsilon} |\mathbf{u}^\top \nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)| > a\right) \leq \frac{C_\epsilon^2 d_n}{a^2 n} K_1,$$

where $0 < K_1 < \infty$ is some constant.

We now focus on the Hessian. First, $\mathcal{R}_n(\boldsymbol{\theta}_0)$ can be written as

$$\mathcal{R}_n(\boldsymbol{\theta}_0) = \sum_{k,l=1}^{d_n} \mathbf{u}_k \mathbf{u}_l \{\nabla_{\boldsymbol{\theta}_k \boldsymbol{\theta}_l}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) - \mathbb{H}_n\}.$$

We have $\mathbb{E}[\mathcal{R}_n(\boldsymbol{\theta}_0)] = 0$ and the variance is defined as

$$\text{Var}(\mathcal{R}_n(\boldsymbol{\theta}_0)) = (n)_4^{-2} \sum_{(i,j,q,r) \in I_4^n} \sum_{(i',j',q',r') \in I_4^n} \sum_{k,k',l,l'=1}^{d_n} \mathbf{u}_k \mathbf{u}_l \mathbf{u}_{k'} \mathbf{u}_{l'} \mathbb{E}[\zeta_{k,l,(i,j,q,r)} \zeta_{k',l',(i',j',q',r')}],$$

where $\zeta_{k,l,(i,j,q,r)} = \nabla_{\theta_k \theta_l}^2 \ell(\mathbf{Z}_i, \mathbf{Z}_j, \mathbf{Z}_q, \mathbf{Z}_r; \boldsymbol{\theta}_0) - \mathbb{E}[\nabla_{\theta_k \theta_l}^2 \ell(\mathbf{Z}_i, \mathbf{Z}_j, \mathbf{Z}_q, \mathbf{Z}_r; \boldsymbol{\theta}_0)]$. Let $b > 0$, we have by the Markov inequality and assumption 10

$$\mathbb{P}(|\mathcal{R}_n(\boldsymbol{\theta}_0)| > b) \leq \frac{1}{b^2} \mathbb{E}[\mathcal{R}_n^2(\boldsymbol{\theta}_0)] \leq \frac{K_2 \|\mathbf{u}\|_2^4 d_n^2}{n} \leq \frac{C_\epsilon^4 K_2 d_n^2}{n},$$

for some constant $K_2 > 0$. By assumption 9,

$$\mathbf{u}^\top \mathbb{E}[\nabla_{\boldsymbol{\theta}^\top}^2 \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)] \mathbf{u} \geq \lambda_{\min}(\mathbb{H}_n) \mathbf{u}^\top \mathbf{u}.$$

We now turn to the penalty terms. For any $k \in \mathcal{A}_n \subset \{1, \dots, d_n\}$, and since the penalties are coordinate-separable, we have

$$\varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) = \nu_n u_k \operatorname{sgn}(\theta_{0,k}) \nabla_{\theta_k} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) + \frac{\nu_n^2}{2} u_k^2 \nabla_{\theta_k \theta_k}^2 \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) (1 + o(1)).$$

Hence, using $\varphi\left(\frac{\lambda_n}{n}, 0\right) = 0$, we have

$$\left| \sum_{k \in \mathcal{A}_n} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \right| \leq \nu_n \|\mathbf{u}\|_1 A_{1,n} + \frac{\nu_n^2}{2} \|\mathbf{u}\|_2^2 A_{2,n} (1 + o(1)).$$

Using $\|\mathbf{u}\|_1 \leq \sqrt{\operatorname{card}(\mathcal{A}_n)} \|\mathbf{u}\|_2$, and under assumption 8, the third derivative being dominated, we obtain

$$\left| \sum_{k \in \mathcal{A}} \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \right| \leq \nu_n \sqrt{\operatorname{card}(\mathcal{A}_n)} \|\mathbf{u}\|_2 A_{1,n} + \frac{\nu_n^2}{2} \|\mathbf{u}\|_2^2 A_{2,n}.$$

Then, denoting $\delta_n = \lambda_{\min}(\mathbb{H}_n) C_\epsilon^2 \nu_n$, using $\frac{\nu_n}{2} \mathbf{u}^\top \mathbb{E}[\nabla_{\boldsymbol{\theta}^\top}^2 \ell(\mathbf{Z}; \boldsymbol{\theta}_0)] \mathbf{u} \geq \delta_n$, we deduce that Eq. (13) can be bounded as

$$\begin{aligned} \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \mathbf{u}^\top \nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) + \frac{\nu_n}{2} \mathbf{u}^\top \mathbb{E}[\nabla_{\boldsymbol{\theta}^\top}^2 \ell(\mathbf{Z}; \boldsymbol{\theta}_0)] \mathbf{u} + \frac{\nu_n}{2} \mathcal{R}_n(\boldsymbol{\theta}_0) \\ + \nu_n^{-1} \left(\sum_{k=1}^{d_n} \{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \} \right) \leq 0) \\ \leq \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \mathbf{u}^\top |\nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)| > \delta_n/6) + \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon | \frac{\nu_n}{2} \mathcal{R}_n(\boldsymbol{\theta}_0)| > \delta_n/6) \\ \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \left| \sum_{k=1}^{d_n} \{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \} \right| > \nu_n \delta_n/6). \end{aligned}$$

We have for n and C_ϵ sufficiently large enough

$$\begin{aligned} \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \left| \sum_{k=1}^{d_n} \{ \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k} + \nu_n u_k|\right) - \varphi\left(\frac{\lambda_n}{n}, |\theta_{0,k}|\right) \} \right| > \nu_n \delta_n/6) \\ \leq \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \nu_n \sqrt{\operatorname{card}(\mathcal{A}_n)} \|\mathbf{u}\|_2 A_{1,n} + \frac{\nu_n^2}{2} \|\mathbf{u}\|_2^2 A_{2,n} > \nu_n \delta_n/6) < \epsilon/3. \end{aligned}$$

Moreover, if $\nu_n = n^{-1/2} + \sqrt{\operatorname{card}(\mathcal{A}_n)} A_{1,n}$, for C_ϵ large enough

$$\mathbb{P}(\exists \mathbf{u} \in \mathbb{R}^{d_n}, \|\mathbf{u}\|_2 = C_\epsilon : \mathbf{u}^\top |\nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0)| > \delta_n/6) \leq \frac{C_\epsilon^2 d_n C_{st}}{n \delta_n^2},$$

where $C_{st} > 0$ is a generic constant. Using

$$\mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon | \frac{\nu_n}{2} \mathcal{R}_n(\boldsymbol{\theta}_0)| > \delta_n/6) \leq \frac{L_{st} \nu_n^2 d_n^2}{n \delta_n^2},$$

where $L_{st} > 0$ is a generic constant. Thus for C_ϵ fixed sufficiently large, under the assumption $d_n^2 = O(n)$, we have

$$\begin{aligned} \mathbb{P}(\exists \mathbf{u}, \|\mathbf{u}\|_2 = C_\epsilon : \mathbf{u}^\top \nabla_{\boldsymbol{\theta}} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) + \frac{\nu_n}{2} \mathbf{u}^\top \mathbb{E}[\nabla_{\boldsymbol{\theta}^\top}^2 \ell(\mathbf{Z}; \boldsymbol{\theta}_0)] \mathbf{u} + \frac{\nu_n}{2} \mathcal{R}_n(\boldsymbol{\theta}_0) \\ \frac{C_\epsilon^2 d_n C_{st}}{n \delta_n^2} + \frac{L_{st} \nu_n^2 d_n^2}{n \delta_n^2} + \epsilon/3 < \epsilon). \end{aligned}$$

We deduce $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| = O_p(\nu_n)$. □

D Derivations

Derivation of the HSIC Lasso

$$\begin{aligned}
 n^2 J(\boldsymbol{\theta}) &= n^2 \sum_{k=1}^d \theta_k \widehat{\text{HSIC}}(X^{(k)}, Y) - \frac{n^2}{2} \sum_{k=1}^d \sum_{k'=1}^d \theta_k \theta_{k'} \widehat{\text{HSIC}}(X^{(k)}, X^{(k')}) \\
 &= \sum_{k=1}^d \theta_k \text{trace}(\mathbf{H}_n \mathbf{L} \mathbf{H}_n \mathbf{K}^{(k)}) - \frac{1}{2} \sum_{k=1}^d \sum_{k'=1}^d \theta_k \theta_{k'} \text{trace}(\mathbf{H}_n \mathbf{K}^{(k)} \mathbf{H}_n \mathbf{K}^{(k')}) \\
 &= \sum_{k=1}^d \theta_k \text{trace}(\mathbf{H}_n \mathbf{L} \mathbf{H}_n \mathbf{H}_n \mathbf{K}^{(k)} \mathbf{H}_n) - \frac{1}{2} \sum_{k=1}^d \sum_{k'=1}^d \theta_k \theta_{k'} \text{trace}(\mathbf{H}_n \mathbf{K}^{(k)} \mathbf{H}_n \mathbf{H}_n \mathbf{K}^{(k')} \mathbf{H}_n) \\
 &\propto -\frac{1}{2} \text{trace}(\mathbf{H}_n \mathbf{L} \mathbf{H}_n \mathbf{H}_n \mathbf{L} \mathbf{H}_n) + \sum_{k=1}^d \theta_k \text{trace}(\mathbf{H}_n \mathbf{L} \mathbf{H}_n \mathbf{H}_n \mathbf{K}^{(k)} \mathbf{H}_n) - \frac{1}{2} \sum_{k=1}^d \sum_{k'=1}^d \theta_k \theta_{k'} \text{trace}(\mathbf{H}_n \mathbf{K}^{(k)} \mathbf{H}_n \mathbf{H}_n \mathbf{K}^{(k')} \mathbf{H}_n), \\
 &= -\frac{1}{2} \|\text{vec}(\mathbf{H}_n \mathbf{L} \mathbf{H}_n) - \sum_{k=1}^d \theta_k \text{vec}(\mathbf{H}_n \mathbf{K}^{(k)} \mathbf{H}_n)\|_2^2, \\
 &= -\frac{1}{2} \|\mathbf{H}_n \mathbf{L} \mathbf{H}_n - \sum_{k=1}^d \theta_k \mathbf{H}_n \mathbf{K}^{(k)} \mathbf{H}_n\|_F^2
 \end{aligned}$$

where we use $\mathbf{H}_n \mathbf{H}_n = \mathbf{H}_n$, $\text{trace}(\mathbf{A}^\top \mathbf{B}) = \text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B})$, and $\text{vec}(\mathbf{A}) \in \mathbb{R}^{n^2}$, $(\mathbf{A} \in \mathbb{R}^{n \times n})$ is the vectorization operator.

Derivative of the U-statistics

Element-by-element $\forall k \leq d$, the score evaluated at the true parameter $\boldsymbol{\theta}_0$ is given by

$$\begin{aligned}
 \nabla_{\theta_k} \mathbb{L}_n(\mathbf{Z}_1, \dots, \mathbf{Z}_n; \boldsymbol{\theta}_0) &= (n)_4^{-1} \sum_{(i,j,q,r)} \nabla_{\theta_k} \ell(\mathbf{Z}_i, \mathbf{Z}_j, \mathbf{Z}_q, \mathbf{Z}_r; \boldsymbol{\theta}_0) \\
 &= (n)_4^{-1} \sum_{(i,j,q,r) \in \mathcal{I}_4^n} \frac{1}{24} \sum_{(s,t,u,v)}^{(i,j,q,r)} \left(-2K_{st}^k(L_{st} - \sum_{l=1}^d \theta_{0,l} K_{st}^l) - K_{st}^k(L_{uv} - \sum_{l=1}^d \theta_{0,l} K_{uv}^l) \right. \\
 &\quad - K_{uv}^k(L_{st} - \sum_{l=1}^d \theta_{0,l} K_{st}^l) + K_{uv}^k L_{st} - K_{st}^k L_{uv} + 2K_{st}^k(L_{su} - \sum_{l=1}^d \theta_{0,l} K_{su}^l) + 2K_{su}^k(L_{st} - \sum_{l=1}^d \theta_{0,l} K_{st}^l) \\
 &\quad \left. + K_{su}^k L_{st} - K_{st}^k L_{su} \right).
 \end{aligned}$$

Kernel expression

$$\mathbb{L}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{Z}_i \mathbf{Z}_j \mathbf{Z}_q \mathbf{Z}_r} [\ell(\mathbf{Z}_i, \mathbf{Z}_j, \mathbf{Z}_q, \mathbf{Z}_r; \boldsymbol{\theta})],$$

which can be expressed in terms of kernel as

$$\begin{aligned}
 \mathbb{L}(\boldsymbol{\theta}) &= \mathbb{E}_{Y Y'} [\phi(Y, Y')^2] - 2\mathbb{E}_{Y Y'} [\mathbb{E}_{Y'} [\phi(Y, Y')^2]] + \mathbb{E}_{Y Y'} [\phi(Y, Y')]^2 \\
 &\quad - 2 \sum_{k=1}^d \theta_k \{ \mathbb{E}_{Y Y' X^{(k)} X^{(k)'}} [\phi(Y, Y') \psi(X^{(k)}, X^{(k)'})] - 2\mathbb{E}_{Y X^{(k)}} [\mathbb{E}_{Y'} [\phi(Y, Y')] \mathbb{E}_{X^{(k)'}} [\psi(X^{(k)}, X^{(k)'})]] \} \\
 &\quad + \mathbb{E}_{Y Y'} [\phi(Y, Y')] \mathbb{E}_{X^{(k)} X^{(k)'}} [\psi(X^{(k)}, X^{(k)'})] + \sum_{k,l=1}^d \theta_k \theta_l \{ \mathbb{E}_{X^{(k)} X^{(k)' X^{(l)} X^{(l)'}} [\psi(X^{(k)}, X^{(k)'}) \psi(X^{(l)}, X^{(l)'})] \\
 &\quad - 2\mathbb{E}_{X^{(k)} X^{(l)}} [\mathbb{E}_{X^{(k)'}} [\psi(X^{(k)}, X^{(k)'})] \mathbb{E}_{X^{(l)'}} [\psi(X^{(l)}, X^{(l)'})]] + \mathbb{E}_{X^{(k)} X^{(k)'}} [\psi(X^{(k)}, X^{(k)'})] \mathbb{E}_{X^{(l)} X^{(l)'}} [\psi(X^{(l)}, X^{(l)'})] \}.
 \end{aligned} \tag{14}$$