

Modern Statistical Methods

Example sheet 2

$$1. \quad x, x' \in \mathbb{R}^P \quad \gamma \in \{-1, 1\}^P$$

$$\begin{aligned} E(\gamma^T x \gamma^T x') &= E((\gamma^T x)^T \gamma^T x') = E(x^T \gamma \gamma^T x') \\ &= x^T E(\gamma \gamma^T) x' = x^T x' \end{aligned}$$

$$\therefore (E(\gamma \gamma^T))_{ij} = E(\gamma_i \gamma_j) = \delta_{ij} \quad \text{as } \gamma_i, \gamma_j \text{ independent when } i \neq j. \\ \text{and } E(\gamma_i) = 0.$$

Note that $[E(\gamma^T x \gamma^T x')]^2 = (x^T x')^2$

$$\begin{aligned} (x^T x')^2 &= E(\gamma^T x \gamma^T x') E(\cancel{\gamma^T x} \tilde{\gamma}^T x \tilde{\gamma}^T x') \quad \text{where } \gamma, \tilde{\gamma} \in \{-1, 1\}^P \\ &= E(\gamma^T x \tilde{\gamma}^T x \gamma^T x' \tilde{\gamma}^T x') \\ &= E(x^T \gamma \tilde{\gamma}^T x \ x'^T \gamma \tilde{\gamma}^T x') \end{aligned}$$

as before

Hence $\hat{\phi}: \mathbb{R}^P \rightarrow \mathbb{R}$

$$x \longmapsto x^T \gamma \tilde{\gamma}^T x$$

2. ~~True~~ $p=1$ ✓

$$\text{Suppose } P_p(M = M') = \frac{|z \cap z'|}{|z \cup z'|} \quad \text{if } z, z' \neq \emptyset, z, z' \in \mathcal{X}_p$$

Take $z, z' \in \mathcal{X}_p$. Examine $k \in \{1, \dots, p\}: \pi(k) = 1$,

if $k \in z, z'$ then $M = M' = 1$

otherwise if $k \notin z, z'$ then $z, z' \in \mathcal{X}_{p \setminus \{k\}} \sim \mathcal{X}_{p-1}$ and can completely ignore element k

otherwise $M \neq M'$.

Hence,

~~$$P_p(M = M') = P(k \in z \cap z') + P(k \notin z \cup z') P_{p-1}(M = M')$$~~

$$\begin{aligned} P_p(M = M') &= P(k \in z \cap z') + P(k \notin z \cup z') P_{p-1}(M = M') \\ &= \frac{|z \cap z'|}{p} + \frac{p - |z \cup z'|}{p} \frac{|z \cap z'|}{|z \cup z'|} = \frac{|z \cap z'|}{|z \cup z'|} \end{aligned}$$

□

$$\mathbb{E}(\gamma_n \gamma_{n'}) = \mathbb{P}(M = M') = k(z, z') \quad \because \text{Condition 2 of } M \neq M' \\ \gamma_M \text{ and } \gamma_{M'} \text{ independent}$$

\Rightarrow Jaccard similarity kernel corresponds to random feature map

$$\hat{\phi}: z \mapsto \gamma_M \text{ where } \gamma \in \{-1, 1\}^P \\ \text{and } M = \bigcup_{\pi \text{ random}} \{ \pi(k) : k \in z \} \text{ for random } \pi \\ \text{over } X \text{ over the input space.}$$

For large n , beneficial to approx using

$$\Phi_{ij} = \frac{1}{\sqrt{L}} \hat{\phi}_j(x_i) \quad \text{where } i=1, \dots, n, j=1, \dots, L \\ \text{and } L \text{ is the number of random feature maps } \phi_j \text{ used.}$$

$$\text{then } \hat{\beta} = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T Y.$$

$$3. \log \left(\frac{P(Y_i = 1)}{P(Y_i = -1)} \right) = x_i^T \beta^* \quad \log P(Y_i = 1) - \log P(Y_i = -1) = x_i^T \beta^*$$

Log-likelihood

$$\ell(\beta) = \sum_{i=1}^n \log P(Y_i | x_i; \beta) = \sum_{i=1}^n \frac{Y_i + 1}{2} \log P(Y_i = 1) - \frac{Y_i - 1}{2} \log P(Y_i = -1) \\ = \frac{1}{2} \sum_{i=1}^n Y_i \left(\log P(Y_i = 1) - \log P(Y_i = -1) \right) + \frac{1}{2} \log P(Y_i = 1) + \frac{1}{2} \log P(Y_i = -1) \\ = \sum_{i=1}^n \frac{1}{2} Y_i x_i^T \beta - \frac{1}{2} x_i^T \beta + \frac{1}{2} \log P(Y_i = 1) \\ = \sum_{i=1}^n \frac{1}{2} Y_i x_i^T \beta - \frac{1}{2} x_i^T \beta + \frac{1}{2} \log \frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)} \\ = \sum_{i=1}^n \frac{1}{2} Y_i x_i^T \beta + \frac{1}{2} x_i^T \beta - \frac{1}{2} \log (1 + e^{x_i^T \beta})$$

Maximize w.r.t. β

$$\frac{\partial \ell(\beta)}{\partial \beta} = \sum_{i=1}^n \cancel{\frac{1}{2} Y_i x_i^T} \cancel{x_i^T} \cancel{x_i^T} \cancel{x_i^T} \cancel{\frac{1}{2} x_i^T} \cancel{\frac{1}{2} x_i^T} \cancel{\frac{1}{2} x_i^T} \cancel{\frac{1}{2} x_i^T} \cancel{\frac{1}{2} x_i^T} = 0$$

$$= \sum_{i=1}^n \log \frac{\exp(\frac{1}{2} Y_i x_i^T \beta)}{(1 + \exp(x_i^T \beta))} = \sum_{i=1}^n \log \frac{\exp(\frac{1}{2} Y_i x_i^T \beta)}{\left(\exp(-\frac{1}{2} x_i^T \beta) + \exp(\frac{1}{2} x_i^T \beta) \right)}$$

Equivalent measure

$$\sum_{i=1}^n \log \left(\exp \left(-\frac{1}{2} \gamma_i (Y_i + 1) \mathbf{x}_i^T \beta \right) + \exp \left(-\frac{1}{2} (Y_i - 1) \mathbf{x}_i^T \beta \right) \right)$$

$$= \sum_{i=1}^n \log \left(1 + \exp (-Y_i \mathbf{x}_i^T \beta) \right) \quad \text{for } Y_i \in \{1, -1\}$$

$$\therefore \hat{\beta} = \arg \min_{\beta \in \mathbb{R}^P} \sum_{i=1}^n \log \{ 1 + \exp (-Y_i \mathbf{x}_i^T \beta) \}$$

4. Suppose this algorithm at iteration m is equivalent to m steps of forward selection such that

$$R = Y - X \hat{\beta}_{S_m}^{(k)}$$

$$\hat{\beta}_{S_m}^{(k)} = \underset{\beta_{S_m}}{\arg \min} \|Y - X \beta_{S_m}^{(k)}\|_2^2 \quad \text{with } S_m = \{k \in \{1, \dots, P\} : \beta_{S_m, k} \neq 0\}$$

The algorithm involves choosing $k^* = \underset{k \in M^c}{\arg \max} \operatorname{Corr}(R, X_{k^*})$, then

$$R \rightarrow R - X_{k^*} \hat{\beta}_{k^*} \quad \text{where } \hat{\beta}_{k^*} = \underset{\beta_{k^*}}{\arg \min} \|R - X_{k^*} \beta_{k^*}\|_2^2$$

$$\rightarrow \hat{\beta}_{k^*} = (X_{k^*}^T X_{k^*})^{-1} X_{k^*}^T R$$

$$\begin{aligned} R^* &= R - X_{k^*} \hat{\beta}_{k^*} = Y - X \hat{\beta}_{S_m} - X_{k^*} (X_{k^*}^T X_{k^*})^{-1} X_{k^*}^T (Y - X \hat{\beta}_{S_m}) \\ &= Y - X \hat{\beta}_{S_m} - X_{k^*} (X_{k^*}^T X_{k^*})^{-1} X_{k^*}^T Y \end{aligned}$$

because $X_{k^*}^T X \hat{\beta}_{S_m} = 0$:

$$X_{k^*}^T X_k = 0 \quad \forall k \in M$$

Suppose $X_k^T X_\ell = 0$ at iteration $m \forall k \in M, \ell \in M^c$

then for $X_{k^*}, k^* \in M^c$

$$\begin{aligned} X_\ell &\rightarrow X_\ell - X_{k^*} \tilde{\beta}_{k^*} \quad \text{with } \tilde{\beta}_{k^*} = (X_{k^*}^T X_{k^*})^{-1} X_{k^*}^T X_\ell \\ &= X_\ell - X_{k^*} (X_{k^*}^T X_{k^*})^{-1} X_{k^*}^T X_\ell \end{aligned}$$

$$\text{Now } X_k^T X_\ell^* = X_k^T X_\ell - X_k^T X_{k^*} (X_{k^*}^T X_{k^*})^{-1} X_{k^*}^T X_\ell = 0$$

as $X_k^T X_\ell = 0$ and $X_k^T X_{k^*} = 0$ as $k^* \in M^c$

$$\text{and } X_{k^*}^T X_{k^*} = X_{k^*}^{*T} X_k - X_{k^*}^T X_{k^*} (X_{k^*}^T X_{k^*})^{-1} X_{k^*}^T X_k \\ = X_{k^*}^{*T} X_k - X_{k^*}^T X_k = 0$$

Hence, $X_k^T X_k = 0$ at iteration $m+1 \quad \forall k \in M^* \setminus \{k^*\}, k \in M^{*c}$ \square

{ Sufficient to show that }

$$R^* \stackrel{?}{=} Y - X \hat{\beta}_{S_m^*} \quad \text{where } S_m^* = S_m \cup \{k^*\}$$

$\{k^*$ chosen by forward selection

{ Know that }

$$\hat{\beta}_{S_m} = (X_{S_m}^T X_{S_m})^{-1} X_{S_m}^T Y$$

$$\|R^*\|_2^2 = (R - X_{k^*} \hat{\beta}_{k^*})^T (R - X_{k^*} \hat{\beta}_{k^*}) = \|R\|_2^2 - 2\hat{\beta}_{k^*} X_{k^*}^T R + \|X_{k^*} \hat{\beta}_{k^*}\|_2^2 \\ = \|R\|_2^2 - \|X_{k^*} \hat{\beta}_{k^*}\|_2^2 \quad \text{because } X_{k^*}^T R = X_{k^*}^T Y = (X_{k^*}^T X_{k^*})(X_{k^*}^T X_{k^*})^{-1} X_{k^*}^T Y \\ = \|X_{k^*} \hat{\beta}_{k^*}\|_2^2$$

However,

$$k^* = \underset{k^* \in M^c}{\operatorname{argmax}} \operatorname{Cov}(R, X_{k^*}) \Rightarrow \frac{\operatorname{Cov}(R, X_{k^*})}{\|R\|_2 \|X_{k^*}\|_2} \geq \frac{\operatorname{Cov}(R, X_k)}{\|R\|_2 \|X_k\|_2} \quad \forall k \in M^c$$

$$\Rightarrow \frac{X_{k^*}^T R}{\|R\|_2 \|X_{k^*}\|_2} \geq \frac{X_k^T R}{\|R\|_2 \|X_k\|_2}$$

$$\frac{\|X_{k^*}\|_2^2 \hat{\beta}_{k^*}}{\|X_{k^*}\|_2} \geq \frac{\|X_k\|_2^2 \hat{\beta}_k}{\|X_k\|_2}$$

$$\Rightarrow \|X_{k^*} \hat{\beta}_{k^*}\|_2 \geq \|X_k \hat{\beta}_k\|_2$$

$$\therefore \|R^*\|_2^2 = \|R\|_2^2 - \|X_{k^*} \hat{\beta}_{k^*}\|_2^2 \leq \|R\|_2^2 - \|X_k \hat{\beta}_k\|_2^2 \quad \forall k \in M^c$$

Hence, k^* minimizes $\|R - X_k \beta_k\|_2^2$, $k \in M^c$.

On the other hand, forward selection on this step would minimize

~~$$\|Y - X \hat{\beta}_{S_m} - X_k \hat{\beta}_k\|_2^2$$~~

~~$$= \|Y - X \hat{\beta}_{S_m} - X_k \hat{\beta}_k\|_2^2$$~~

$= \|Y - X \hat{\beta}_{S_m^*}\|_2^2$ with cols of X modified according to this algorithm. Hence equivalent to forward selection on transformed cols

$$5. \mathbb{E}W=0 \quad \mathbb{E}e^{\alpha W} \leq e^{\alpha^2 \sigma^2/2} \quad \forall \alpha \in \mathbb{R}$$

$$\text{Var}(W) = \mathbb{E}((W - \mathbb{E}W)^2) = \mathbb{E}(W^2)$$

$$\frac{d}{d\alpha} \mathbb{E}e^{\alpha W} = \mathbb{E} \frac{d}{d\alpha} e^{\alpha W} = \mathbb{E}We^{\alpha W}$$

$$\frac{d}{d\alpha} \mathbb{E}e^{\alpha W} \leq \frac{d}{d\alpha} (e^{\alpha^2 \sigma^2/2}) = \alpha \sigma^2 e^{\alpha^2 \sigma^2/2}$$

$$\frac{d^2}{d\alpha^2} \mathbb{E}e^{\alpha W} = \mathbb{E}W^2 e^{\alpha W} \leq (\sigma^2 + \alpha^2 \sigma^4) e^{\alpha^2 \sigma^2/2}$$

At $\alpha = 0$,

$$\mathbb{E}W^2 \leq \sigma^2 \quad \square$$

6. Hoeffding's lemma: $\mathbb{E}W=0$, $W \in [a, b]$ \Rightarrow sub-G, $\sigma = (b-a)/2$

If $W \in \{-1, 1\}$ with equal probability,

$$\mathbb{E}W = 0$$

$$\mathbb{E}e^{\alpha W} = \frac{1}{2} e^{\alpha(-1)} + \frac{1}{2} e^{\alpha} = \frac{1}{2}(e^{-\alpha} + e^{\alpha}) \leq e^{\alpha^2/2}$$

$\Rightarrow \sigma = 1$ as required

~~7(a) $W \sim \chi_d^2 \sim \sum_i^d \chi_i^2$, $W = \sum_{i=1}^d W_i$, $W_i \sim \chi_i^2$~~

~~$P(\frac{1}{d}W/d - 1 \geq t) \leq \inf_{\alpha > 0} e^{-\alpha t} \mathbb{E}e^{\alpha(W-d)}$~~

~~$\mathbb{E}e^{\alpha(W-d)} = \mathbb{E}e^{\sum_{i=1}^d \alpha(W_i-1)} = \prod_{i=1}^d \mathbb{E}e^{\alpha(W_i-1)} = \prod_{i=1}^d e^{\alpha^2/2} = (1-2\alpha)^{-d/2} e^{\alpha^2 d/2}$ for $\alpha < \frac{1}{2}$~~

~~$\Rightarrow P(\frac{1}{d}W/d - 1 \geq t) \leq \inf_{\alpha > 0} e^{-\alpha t} (1-2\alpha)^{-d/2} e^{\alpha^2 d/2} = \inf_{\alpha > 0} e^{-\alpha t} (1-2\alpha)^{-d/2} e^{\alpha^2 d/2} = e^{-dt^2/8}$~~

(b) $A \in \mathbb{R}^{d \times p} \sim N(0, I)$, $u \in \mathbb{R}^p$

~~$$\frac{\|Au\|_2^2}{\|u\|_2^2} = \frac{\sum_{i=1}^d \left(\sum_{j=1}^p A_{ij} u_j \right)^2}{\sum_{j=1}^p u_j^2} = \frac{\sum_{i=1}^d \|A_i u\|_2^2}{\sum_{j=1}^p u_j^2} = \frac{\sum_{i=1}^d \|A_i u\|_2^2}{\|u\|_2^2} \leq \frac{\sum_{i=1}^d \|A_i u\|_2^2}{\|u\|_2^2}$$~~

~~$$\frac{\sum_{i=1}^d \|A_i u\|_2^2}{\|u\|_2^2} = \left\| \frac{A u}{\|u\|_2} \right\|_2^2 \text{ but } \frac{A u}{\|u\|_2} \sim N(0, I)$$~~

thus each term $\|A_i u\|_2^2 / \|u\|_2^2$ is the sum of d $N(0, 1)$ vars $\sim \chi_d^2$ and the result follows from previous part.

$$(c) u_i \mapsto A u_i / \sqrt{d} = w_i$$

$$P\left(1-t \leq \frac{\|w_i - w_j\|_2^2}{\|u_i - u_j\|_2^2} \leq 1+t \quad \forall i, j \in \{1, \dots, n\}, i \neq j\right)$$

$$= P\left(1-t \leq \frac{\|A(u_i - u_j)\|_2^2}{d \|u_i - u_j\|_2^2} \leq 1+t \quad \forall i, j, i \neq j\right)$$

$$= P\left(\left|\frac{\|A(u_i - u_j)\|_2^2}{d \|u_i - u_j\|_2^2} - 1\right| \leq t \quad \forall i, j \in \{1, \dots, n\}, i \neq j\right)$$

with constraint
or equal under $u_i = u_j$

$$= 1 - P\left(\left|\frac{\|A(u_i - u_j)\|_2^2}{d \|u_i - u_j\|_2^2} - 1\right| > t \quad \exists i, j \in \{1, \dots, n\}, i \neq j\right)$$

$$\geq 1 - \sum_{\substack{i, j \in \{1, \dots, n\} \\ i \neq j}} P\left(\left|\frac{\|A(u_i - u_j)\|_2^2}{d \|u_i - u_j\|_2^2} - 1\right| > t\right) \quad (\text{as } P(\cup_i \Omega_i) \leq \sum_i P(\Omega_i))$$

$$\geq 1 - \sum_{\substack{i, j \in \{1, \dots, n\} \\ i \neq j}} P\left(\left|\frac{\|A(u_i - u_j)\|_2^2}{d \|u_i - u_j\|_2^2} - 1\right| \geq t\right) \quad \text{treating } u_i - u_j \text{ as fixed}$$

$$\geq 1 - \sum_{\substack{i, j \in \{1, \dots, n\} \\ i \neq j}} 2e^{-dt^2/8} = 1 - \frac{n(n-1)}{2} 2e^{-dt^2/8}$$

$$= 1 - n(n-1)e^{-dt^2/8} \quad \text{use } d > 16 \log(n/\epsilon) / t^2 \\ \Rightarrow dt^2/8 > 2 \log(n/\epsilon)$$

$$> 1 - n(n-1)e^{-2 \log(n/\epsilon)} = 1 - \frac{n(n-1)}{n^2} \epsilon = 1 - \epsilon + \frac{\epsilon}{n} > 1 - \epsilon$$

$$\therefore P\left(1-t \leq \frac{\|w_i - w_j\|_2^2}{\|u_i - u_j\|_2^2} \leq 1+t \quad \forall i, j \in \{1, \dots, n\}, i \neq j\right) > 1 - \epsilon \quad \square$$

$$8. Q_\lambda(\beta) = \frac{1}{2n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

Any minimiser of Q_λ , $\hat{\beta}_\lambda^L$ must also minimise

$$\|Y - X\beta\|_2^2 \text{ subject to } \|\beta\|_1 \leq \|\hat{\beta}_\lambda^L\|_1.$$

\therefore If $\|\beta\|_1 \leq \|\hat{\beta}_\lambda^L\|_1$, and $\|Y - X\beta\|_2^2 \leq \|Y - X\hat{\beta}_\lambda^L\|_2^2$, then $Q_\lambda(\beta) < Q_\lambda(\hat{\beta}_\lambda^L)$.

Suppose two Lasso solutions $\hat{\beta}_\lambda^1, \hat{\beta}_\lambda^2$ have ℓ_1 -norms

$$\|\hat{\beta}_\lambda^1\|_1 \leq \|\hat{\beta}_\lambda^2\|_1.$$

However, from above, $\hat{\beta}_\lambda^2$ must also minimise

$$\|Y - X\beta\|_2^2 \text{ subject to } \|\beta\|_1 \leq \|\hat{\beta}_\lambda^1\|_1.$$

$$\therefore \|\hat{\beta}_\lambda^1\|_1 = \|\hat{\beta}_\lambda^2\|_1. \quad \square$$

$$9. \text{ Consider } S = \{-\text{sgn}(\beta_1)X_1, \dots, -\text{sgn}(\beta_p)X_p, Y\} \subseteq \mathbb{R}^{n+1}$$

and consider the convex combination

$$\alpha_0 Y - \alpha_1 \text{sgn}(\beta_1)X_1 + \dots + \alpha_p \text{sgn}(\beta_p)X_p = 0, \quad \sum_{i=0}^p \alpha_i = 1$$

then by Carathéodory's Lemma

$$v = \alpha_0 Y - \alpha_1 \text{sgn}(\beta_{j_1})X_{j_1} - \dots - \alpha_n \text{sgn}(\beta_{j_n})X_{j_n} \quad \text{where } \sum_{i=1}^n \alpha_i = 1$$

\Rightarrow Any $Y - X\beta$ can be expressed with no more than n non-zero coeffs.

$$\text{Above } \tilde{\beta}_i = \alpha_i \text{sgn}(\beta_{j_i}), \quad \|\tilde{\beta}\|_1 = \sum_i \alpha_i = 1$$

~~$$Q_\lambda(\beta) = \frac{1}{2n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1$$~~

$$= \frac{1}{2n} \|Y - X\tilde{\beta}\|_2^2 + \|\tilde{\beta}\|_1 = \frac{1}{2n\lambda} \cancel{\lambda \|Y - X\tilde{\beta}\|_2^2} + \|\tilde{\beta}\|_1$$

Identify $\lambda = \alpha_0$ to show equivalence.

$$10. \lambda \geq \lambda_{\max} := \|X^T Y\|_\infty / n$$

Then

$$\begin{aligned} Q_\lambda(\beta) &= \frac{1}{2n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 \geq \frac{1}{2n} \|Y - X\beta\|_2^2 + \frac{1}{n} \|X^T Y\|_\infty \|\beta\|_1 \\ &\geq \frac{1}{2n} \|Y - X\beta\|_2^2 + \frac{1}{n} \|(X\beta)^T Y\|_1, \quad \text{by Hölder's inequality} \\ &= \frac{1}{2n} (\|Y\|_2^2 + \|X\beta\|_2^2) - \frac{1}{n} (X\beta)^T Y + \frac{1}{n} |(X\beta)^T Y| \\ &\geq \frac{1}{2n} (\|Y\|_2^2 + \|X\beta\|_2^2) \Leftrightarrow \geq \frac{1}{2n} \|Y\|_2^2 = Q_\lambda(0) \end{aligned}$$

Hence, $\hat{\beta}_\lambda^L = 0$ always minimizes the objective and is a solution (unique by Qn 2).

$$11. \frac{\partial Q_\lambda(\beta)}{\partial \beta_k} = -\frac{1}{n} X_k^T (Y - X\beta) + \lambda \operatorname{sgn}(\beta_k)$$

$$= -\frac{1}{n} X_k^T Y + \frac{1}{n} X_k^T X\beta + \lambda \operatorname{sgn}(\beta_k)$$

$$= -\frac{1}{n} X_k^T Y + \frac{1}{n} \|X_k\|_2^2 \beta_k + \lambda \operatorname{sgn}(\beta_k) \quad \text{as } X_k \text{ orthogonal}$$

$$= -\frac{1}{n} X_k^T Y + \beta_k + \lambda \operatorname{sgn}(\beta_k)$$

$$\left. \frac{\partial Q_\lambda(\beta)}{\partial \beta_k} \right|_{\hat{\beta}_{\lambda,k}^L} = 0 \Rightarrow -\frac{1}{n} X_k^T Y + \hat{\beta}_{\lambda,k}^L + \lambda \operatorname{sgn}(\hat{\beta}_{\lambda,k}^L) = 0$$

$$\hat{\beta}_k^{\text{OLS}} = (X^T X)^{-1} X^T Y = \frac{1}{n} X^T Y \quad \text{when } X_k \text{ orthogonal, } l_2 \text{ norm } \sqrt{n}$$

$$\hat{\beta}_k^{\text{OLS}} = \frac{1}{n} X_k^T Y$$

$$\Rightarrow -\hat{\beta}_k^{\text{OLS}} + \hat{\beta}_{\lambda,k}^L + \lambda \operatorname{sgn}(\hat{\beta}_{\lambda,k}^L) = 0$$

$$\therefore \hat{\beta}_{\lambda,k}^L = \begin{cases} \hat{\beta}_k^{\text{OLS}} - \lambda & : \hat{\beta}_{\lambda,k}^L > 0 \\ \hat{\beta}_k^{\text{OLS}} + \lambda & : \hat{\beta}_{\lambda,k}^L < 0 \end{cases} \quad \text{nonzero otherwise}$$

$$\hat{\beta}_{\lambda,k}^L = (|\hat{\beta}_k^{\text{OLS}}| - \lambda)_+ \operatorname{sgn}(\hat{\beta}_k^{\text{OLS}}) \quad \square$$

For ℓ_0 penalty

$$-\hat{\beta}_k^{OLS} + \hat{\beta}_{\lambda,k} + \cancel{\hat{\beta}_{\lambda,k}^2} = 0 \quad \text{if } \hat{\beta}_{\lambda,k} \neq 0$$

$$\lim_{\hat{\beta}_{\lambda,k} \rightarrow 0_+} Q_\lambda(\beta \neq \hat{\beta}, \hat{\beta}_{\lambda,k} \neq 0) - Q_\lambda(\beta, \hat{\beta}_{\lambda,k} = 0) = \lambda$$

Examine when the objective changes as $\beta_k = 0 \rightarrow \neq 0$,

$$\Delta Q = \frac{1}{2n} \|Y - X_p^* - X_k \beta_k\|_2^2 + \lambda - \frac{1}{2n} \|Y - X_p^*\|_2^2$$

$$= \cancel{\frac{1}{n} \beta_k X_k^T (Y - X_p^*)} + \frac{1}{2n} \|X_k \beta_k\|_2^2 + \lambda \quad \text{where } X^* \text{ is } X \text{ with } X_k \text{ removed}$$

$$= \cancel{\frac{1}{n} \beta_k X_k^T Y} + \frac{1}{2n} \beta_k^2 n + \lambda \quad \text{by orthogonality}$$

$$= \cancel{\frac{1}{2n} \beta_k^2 n} + \frac{1}{2} \beta_k^2 + \lambda$$

$$\hat{\beta}_{\lambda,k} \neq 0 \quad \text{when } -\hat{\beta}_k^{OLS} + \frac{1}{2} \hat{\beta}_k^{OLS} + \lambda \leq 0$$

$$\Rightarrow \hat{\beta}_{k,\lambda} = \begin{cases} 0 & : \hat{\beta}_k^{OLS} < \sqrt{2\lambda} \\ \hat{\beta}_k^{OLS} & : \hat{\beta}_k^{OLS} \geq \sqrt{2\lambda} \end{cases}$$