

STAT 239 HW 3
 THEORETICAL QUESTIONS

$$1) \text{Var}(\hat{\beta}_{RR}) = \text{Var}\left(\sum_{j=1}^p u_j \frac{d_j^{-2}}{d_j^{-2} + \lambda} u_j^\top Y\right)$$

$$\Leftrightarrow \text{Var}(x) = \sum_{j=1}^p \text{Var}\left(u_j \frac{d_j^{-2}}{d_j^{-2} + \lambda} u_j^\top Y\right)$$

$$\Leftrightarrow \text{Var}(x) = \sum_{j=1}^p \left(\frac{Y^\top u_j \frac{d_j^{-2}}{d_j^{-2} + \lambda} u_j^\top u_j \frac{d_j^{-2}}{d_j^{-2} + \lambda} u_j^\top Y}{N} \right)$$

$$\Leftrightarrow \text{Var}(x) = \sum_{j=1}^p \frac{d_j^{-4}}{(d_j^{-2} + \lambda)^2} \frac{Y^\top u_j u_j^\top Y}{N}$$

Since $X^\top X$ is diagonal, it is a scaling matrix, therefore its eigenbasis is equal to its basis $\Rightarrow V = V^\top = I \Rightarrow X^\top X = D^2 \Rightarrow x_i^{-2} = d_i^{-2}$
 Also, $X^\top X$ is full rank, therefore $n = s \Rightarrow U \in O(r)$ $\forall i \in \{1, \dots, p\}$
 $\Rightarrow u_i^\top u_j^\top = I$.

$$\Leftrightarrow \text{Var}(x) = \sum_{j=1}^p \frac{d_j^{-4}}{(d_j^{-2} + \lambda)^2} \frac{Y^\top Y}{N}$$

$$\Leftrightarrow \text{Var}(x) = \sigma^2 \sum_{j=1}^p \frac{d_j^{-4}}{(d_j^{-2} + \lambda)^2}$$

$$B(x) = E[\hat{F}(x)] - F(x)$$

$$\Leftrightarrow B(x) = E[x^T \hat{\beta}_{RR}] - x^T \beta$$

$$\Leftrightarrow B(x) = E \left[\sum_{j=1}^p u_j \frac{d_j^{-2}}{d_j^{-2} + \lambda} u_j^T y \right] - \sum_{i=1}^p x_i \beta_i$$

$$\Leftrightarrow B(x) = \sum_{j=1}^p E \left[u_j \frac{d_j^{-2}}{d_j^{-2} + \lambda} u_j^T y \right] - \sum_{i=1}^p x_i \beta_i$$

$$\Leftrightarrow B(x) = \sum_{j=1}^p x_i \beta_i \left(\frac{d_j^{-2}}{d_j^{-2} + \lambda} - 1 \right)$$

since $E[\varepsilon] = 0$, $E[y] = F(x) = x^T \beta$.

(2)

2) $M = (\beta - \hat{\beta}_{OLS})^T (\beta - \hat{\beta}_{OLS}) + \lambda \sum_{i=1}^p |\beta_i|$ is the quantity to be minimized.

Expanding the problem in Lagrangian form:

$$\hat{\beta} = \arg \min_{\beta} \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

Ignoring β_0 and assuming orthogonality, we get

$$M = \frac{1}{2} (\beta - \hat{\beta}_{OLS})^T (\beta - \hat{\beta}_{OLS}) + \lambda \sum_{i=1}^p |\beta_i|$$

For M to be minimal, β_{LASSO} can't be of a different sign than $\hat{\beta}_{OLS}$ or else $\|\beta - \hat{\beta}_{OLS}\|_2$ would not be minimal. (*)

~~Case 1~~ Fix $i \in \{1, \dots, p\}$

Case 1 : $\hat{\beta}_{OLS,i} > 0$

$$M' = \frac{1}{2} \beta_i^2 - \hat{\beta}_{OLS,i} \beta_i + \frac{1}{2} \hat{\beta}_{OLS,i}^2 + \lambda \hat{\beta}_{OLS,i} \quad \text{where } M' = \sum_{i=1}^p M'_i$$

$$\Leftrightarrow M' = \frac{1}{2} \beta_i^2 - \beta_i (\hat{\beta}_{OLS,i} - \lambda) + \frac{1}{2} \hat{\beta}_{OLS,i}^2$$

$$\frac{\partial M'}{\partial \beta_i} = \beta_i - (\hat{\beta}_{OLS,i} - \lambda)$$

Setting the derivative to 0 (which can't be a maximum by the argument given in (*)) gives us

$$\beta_i = (\hat{\beta}_{OLS,i} - \lambda)^+ = (|\hat{\beta}_{OLS,i}| - \lambda)^+ \text{ since } \hat{\beta}_{OLS,i} > 0 \text{ and } \hat{\beta}_{OLS,i} - \lambda$$

$\beta_i = \text{sign}(\hat{\beta}_{OLS,i}) (|\hat{\beta}_{OLS,i}| - \lambda)^+$ must be positive for M' to be as small as possible.

Case 2 : $\hat{\beta}_{\text{orsi}} < 0$

$$M' = \frac{1}{2} \beta_i^2 - \beta_i (\hat{\beta}_{\text{orsi}} + \lambda) + \frac{1}{2} \hat{\beta}_{\text{orsi}}^2$$

$$\frac{\partial M'}{\partial \beta_i} = 0 \Leftrightarrow \beta_i - (\hat{\beta}_{\text{orsi}} + \lambda) = 0$$

$$\Leftrightarrow \beta_i = (\hat{\beta}_{\text{orsi}} + \lambda)^+$$

$$\Leftrightarrow \beta_i = (-|\hat{\beta}_{\text{orsi}}| + \lambda)^+ \text{ since } \hat{\beta}_{\text{orsi}} < 0$$

$$\Leftrightarrow \beta_i = -(|\hat{\beta}_{\text{orsi}}| - \lambda)^+$$

$$\Leftrightarrow \beta_i = \text{sign}(\hat{\beta}_{\text{orsi}}) (|\hat{\beta}_{\text{orsi}}| - \lambda)^+$$

$0 < \beta_i < \infty$

(3)

$$3) \text{ a) Using } L_2 \text{ loss, } L_2(y, \hat{F}(x)) = (y - \hat{F}(x))^2$$

$$F = \underset{F}{\operatorname{argmin}} \left\{ L(1, F) \times p + L(-1, F)(1-p) \right\}$$

$$= \underset{F}{\operatorname{argmin}} \left\{ (1-F)^2 \times p + (-1-F)^2 (1-p) \right\}$$

taking the derivative with respect to F :

$$-2(1-F) \times p - 2(-1-F)(1-p) = 0$$

$$\Leftrightarrow (-2+2F) \times p + (2+2F) \cancel{\times} - p \times (2+2F) = 0$$

$$\Leftrightarrow p(-2+2F) - (2+2F) + 2+2F = 0$$

$$\Leftrightarrow -4p + 2 + 2F = 0$$

$$\Leftrightarrow F = 2p - 1$$

Replacing p by its estimate, given a training set

$$\hat{F} = 2\hat{p} - 1 = 2 \left(\frac{1+y_i}{2} \right) - 1 = ay_i, y_i \in T$$

$$b) R = \frac{1}{m} \sum_{i=1}^m E[(y_i - \hat{F}(x_i))^2]$$

$$R = \frac{1}{m} \sum_{i=1}^m \left[E_{T,y^*}[(y_i - F(x_i))^2] + (E_{T,y^*}[\hat{F}(x_i)] - F(x_i))^2 + E_{T,Y}[(\hat{F}(x_i) - E_{T,Y}[\hat{F}(x_i)])^2] \right]$$

$$\begin{aligned} E_{T,y^*}[(y_i - F(x_i))^2] &= E_{T,y^*}[y_i^2 - 2y_i(2p_i - 1) + (2p_i - 1)^2] \\ &= E_{T,y^*}[y_i^2] - 2(2p_i - 1)E[y_i] + (2p_i - 1)^2 \\ &= 0 - 2(2p_i - 1)(2p_i - 1) + (2p_i - 1)^2 \\ &= -(2p_i - 1)^2 \end{aligned}$$

$$\begin{aligned} \left(E_{T,y} [\hat{F}(x_i)] - F(x_i) \right)^2 &= \left(E_{T,y} [ay_i] - (2p_{i-1}) \right)^2 \\ &= \left(a(2p_{i-1}) - (2p_{i-1}) \right)^2 \\ &= \left((2p_{i-1})(a-1) \right)^2 \end{aligned}$$

$$E_{T,y} [(\hat{F}(x_i) - E_{T,y} [\hat{F}(x_i)])^2] = -a^2 (2p_{i-1})^2$$

$$\Rightarrow R = \frac{1}{n} \sum_{i=1}^n \left(-(2p_{i-1})^2 + (a-1)^2 (2p_{i-1})^2 - a^2 (2p_{i-1})^2 \right)$$

$$R = \frac{1}{n} \sum_{i=1}^n \left(-2a(2p_{i-1})^2 \right)$$