

1. a. Ridge regression parameter estimate

$$\hat{\beta}_r = \arg \min_{\beta} \sum_{i=1}^n (\gamma_i - \beta^T x_i)^2 + \lambda \|\beta\|^2 \quad (1)$$

Show that

$$\hat{\beta}_r = [X^T X + \lambda I]^{-1} X^T Y$$

$$\hat{\beta}_r = \arg \min_{\beta} (\gamma - \beta^T X)^T (\gamma - \beta^T X) + \lambda \beta^T \beta$$

minimizing above equation

$$= \frac{\partial}{\partial \beta} (\gamma - \beta^T X)^T (\gamma - \beta^T X)$$

$$= -2 X^T (\gamma - \beta^T X)$$

$$\frac{\partial}{\partial \beta} \lambda \beta^T \beta = 2 \lambda \beta$$

$$0 = -2 X^T (\gamma - \beta^T X) + 2 \lambda \beta$$

$$0 = -X^T \gamma + X^T \beta^T X + \lambda \beta$$

$$X^T \gamma = (X^T X + \lambda I) \beta$$

$$\hat{\beta}_r = (X^T X + \lambda I)^{-1} X^T \gamma$$

$$b. \quad \hat{\beta} = [X^T X]^{-1} X^T Y$$

$\hat{\beta}_r$ can be viewed as MLE?

Yes, it can be when $\lambda = 0$ $\Rightarrow p < n$

$$\hat{\beta}_\lambda = [X^T X + \lambda I]^{-1} X^T Y$$

$$= [X^T X]^{-1} X^T Y$$

$$\|\hat{\beta}_r\| < \|\hat{\beta}\|$$

$$\hat{\beta} = [X^T X]^{-1} X^T Y$$

$$X^T X$$

$$X^T Y$$

$$\hat{\beta}_r = (X^T X + \lambda I)^{-1} X^T Y$$

$$= \frac{1}{1+\lambda} X^T Y$$

$$\hat{\beta}_r = \frac{1}{1+\lambda} \hat{\beta}$$

1. c. S.T.

$$\hat{\beta}_{\text{MAP}} = \hat{\beta}_{\text{r}} \text{AM}$$

$$\hat{\beta}_{\text{MAP}} = \arg \max_{\beta} f_{\beta|y}(\beta|y)$$

$$P f_{\beta|y}(\beta|y) = \text{constant} \times f_{y|\beta}(y|\beta) f_{\beta}(\beta)$$

$$\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} \sum_{n=1}^N \frac{1}{2} (y_n - \beta^T x_n)^2 + \lambda \sum_{i=1}^p \beta_i^2$$

$$\hat{\beta}^{\text{MAP}} = \arg \max_{\beta} \{ \log P(\beta | x_{1:N}, y_{1:N}, \lambda) \}$$

$$= \arg \max_{\beta} \left\{ \log (P(y_{1:N} | x_{1:N}, \beta)) \sum_{i=1}^p P(\beta_i / \lambda) \right\}$$

$$= \arg \max_{\beta} \left\{ \log P(y_{1:N} | x_{1:N}, \beta) + \sum_{i=1}^p \log(\beta_i / \lambda) \right\}$$

$$P(y_{1:N} | x_{1:N}, \beta) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \beta^T x_n)^2 \right\}$$

$$= \frac{1}{(2\pi\sigma^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma^2} \text{RSS}(\beta) \right\}$$

$$\arg \max_{\beta} \{ \log P(y_{1:N} | x_{1:N}, \beta) \} = \arg \max_{\beta} \{ -\text{RSS}(\beta) \}$$

$$\beta_i \sim N(0, 1/2\lambda) \quad P(\beta_i / \lambda) = \frac{1}{\sqrt{2\pi/\lambda}} \exp \left(-\frac{\beta_i^2 \lambda}{2} \right)$$

variance of β is $\lambda/2$.

$\lambda \uparrow$ MAP estimate diverges from MLE

$$\hat{\beta}_{\text{LASSO}} = \arg \min_{\beta} \sum_{i=1}^n (y_i - \beta^T x_i)^2 + \lambda \|\beta\|$$

Lasso L_1 penalty is equivalent to assuming Laplace distribution of β values

$$p(\beta_i) \propto \exp\{-\lambda |\beta_i|\}$$

$$e^{-\lambda \|\beta\|} \|\hat{\beta}_r\| < \|\hat{\beta}\|$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$= \frac{X^T Y}{X^T X}$$

$$\hat{\beta} = \frac{X^T Y}{X^T X} \quad X^T X = I$$

$$\hat{\beta}_{\text{ridge}} = \frac{(X^T X + \lambda I)^{-1} X^T Y}{1 + \lambda}$$

$$\hat{\beta}_{\text{ridge}} = \frac{1}{1 + \lambda} \hat{\beta}$$

$$\|\hat{\beta}_{\text{ridge}}\| < \|\hat{\beta}\|$$

$$e \text{ ii)} \quad \text{Cov}[\hat{\beta}_r] < \text{Cov}[\hat{\beta}]$$

$$\text{Cov}[\hat{\beta}_r] = \sigma_e^2 [X^T X + \lambda I]^{-1} [X^T X] [X^T X + \lambda I]^{-1}$$

$$= \sigma_e^2 \frac{X^T X}{[X^T X + \lambda I][X^T X + \lambda I]}$$

$$\because X^T X = I \quad = \frac{\sigma_e^2}{(1+\lambda)(1+\lambda)}$$

$$\text{Cov}[\hat{\beta}_r] = (1+\lambda)^{-2} \sigma_e^2$$

$$\text{Cov}[\hat{\beta}] = \sigma^2 [X^T X]^{-1} = \sigma^2$$

$$\frac{\sigma_e^2}{(1+\lambda)^2} < \sigma^2$$

2.

$$y_i = \frac{\gamma_0 x_i^0}{\gamma_1 + x_i^0} + \epsilon_i$$

$$y_i^0 = \beta_0 + \beta_1 x_i^0$$

$$\text{If } y_i^0 = \frac{1}{y_i^0} \quad \beta_0 = \frac{1}{\gamma_0} \quad \beta_1 = \frac{\gamma_1}{\gamma_0} \quad x_i^0 = \frac{1}{x_i^0}$$

$$y_i^0 = \beta_0 + \beta_1 x_i^0$$

$$= \frac{1}{\gamma_0} + \frac{\gamma_1}{\gamma_0} \frac{1}{x_i^0}$$

$$y_i^0 = \frac{x_i^0 + \gamma_1}{\gamma_0 x_i^0}$$

$$= \frac{1 + x_i^0 \gamma_1}{\gamma_0}$$

$$y_i^0 = \frac{1 + x_i^0 \hat{\beta}_1 / \hat{\beta}_0}{1/\hat{\beta}_0}$$

$$\hat{y}_i^0 = \hat{\beta}_0 + x_i^0 \hat{\beta}_1$$

q. a

$$Errin = \frac{1}{n} \sum_{i=1}^n E \left[L(y_i, \hat{g}(x_i)) \right]_{x=x_i}$$

$$= \frac{1}{n} \sum_{i=1}^n E \left[[y_i - \hat{g}(x_i)]^2 \right]_{x=x_i}$$

$$= \frac{1}{n} \sum_{i=1}^n E \left[(y_i - \hat{g}(x_i)) (y_i - \hat{g}(x_i))^T \right]_{x=x_i}$$

$$= \frac{1}{n} \sum_{i=1}^n \left(E(y_i) - \hat{g}(x_i) \right) \left(E(y_i) - \hat{g}(x_i) \right)^T + \\ E \left[(y_i - E[\hat{g}(x_i)]) (y_i - E[\hat{g}(x_i)])^T \right]$$

Thus,

$$Errin = \frac{1}{n} E[\bar{err}] + \frac{2}{n} \sum_{i=1}^n \text{Cov}(y_i, \hat{g}(x_i)) \quad [\text{from ① and ②}]$$

$$\therefore E[\bar{err}] = \frac{1}{n} E \left[(y - \hat{y})^T (y - \hat{y}) \right] \quad \text{①}$$

$$\therefore \text{Cov}(y_i, \hat{g}(x_i)) = E \left[(y_i - E[\hat{g}(x_i)]) (y_i - E[\hat{g}(x_i)])^T \right] \quad \text{②}$$

①

q. b.

$$y_i = g_i^T \theta + \epsilon_i$$

$$y_i = g^T(x_i) \theta + \epsilon_i$$

In sample
test data

$$y^* = G^* \theta + \epsilon^* \quad (1)$$

Training data

$$y = G\theta + \epsilon \quad (2)$$

G is same for (1) and (2)

Training fits

$$\begin{aligned} \hat{y} &= G\hat{\theta} \\ &= G[G^T G]^{-1} G^T y \\ &= G[G^T G]^{-1} G^T [G\theta + \epsilon] \end{aligned}$$

$$= G + P\epsilon$$

↑
Projection matrix $n \times n$

$$\begin{aligned} y^* - \hat{y} &= G\theta + \epsilon^* - G\theta + P\epsilon \\ &= \epsilon^* - P\epsilon \end{aligned}$$

Training error

$$\begin{aligned} y - \hat{y} &= [G\theta + \epsilon] - [G\theta + P\epsilon] \\ &= [I - P] \epsilon \end{aligned}$$

$$Err_{in} = \frac{1}{n} E[(y^* - \hat{y})^T (y^* - \hat{y})]$$

$$= \frac{1}{n} E[(\epsilon^* - P\epsilon)^T (\epsilon^* - P\epsilon)]$$

$$= \frac{1}{n} E[\epsilon^{*T} \epsilon^* - 2\epsilon^{*T} P\epsilon + \epsilon^T P^T P\epsilon]$$

(2)

$$= \frac{1}{n} \{ n \sigma_e^2 - 0 \} + \frac{1}{n} E \{ e^T P e \}$$

for projection
matrix
 $P^T = P$

$$\text{trace } E \{ \sigma^2 I P \}$$

$$\sigma^2 \text{ trace}(P)$$

P
no of regression
coefficient

$$= \frac{1}{n} \{ n \sigma_e^2 + \sigma_e^2 P \}$$

$$\text{Err}_{in} = \sigma_e^2 + \frac{\sigma_e^2 P}{n} \quad (1)$$

Expected value of in sample error

$$E[\bar{err}] = \frac{1}{n} E[(Y - \hat{Y})^T (Y - \hat{Y})]$$

$$= \frac{1}{n} E[e^T [I - P] e]$$

$$= \frac{1}{n} \text{tr } E[e e^T [I - P]]$$

$$= \frac{1}{n} \text{tr} \{ \sigma_e^2 I [I - P] \}$$

$$= \frac{\sigma_e^2}{n} \text{tr} [I - P]$$

$$\text{Err}_{in} - E[\bar{err}] = \frac{2 \sigma_e^2 P}{n} \quad (2)$$

from (1) $\sigma_e^2 = \frac{SSE}{n}$

$$C_p \equiv \hat{\text{Err}}_{in} = \frac{SSE}{n} + \frac{2p \sigma_e^2}{n}$$