

Question 1:

1. $X^T X \beta = X^T y$, $X = S \Sigma_r \Phi^T$, $X^T X = \Phi \Sigma_r^2 \Phi^T$

$$X^T y = \Phi \Sigma_r S^T y$$

$$\hat{\beta}_{LS} = (X^T X)^{-1} X^T y = \Phi \Sigma_r^{-2} \Phi^T \Phi \Sigma_r S^T y = \Phi \Sigma_r^{-1} S^T y$$

$$= \sum_{i=1}^p \frac{S_i^T y}{\sigma_i} \phi_i = \sum_{i=1}^p \hat{b}_i$$

$$\hat{b}_i = \frac{S_i^T y}{\sigma_i} \phi_i$$

2. The ridge regression estimator is given by

$$\hat{\beta}_{RR} = \arg \min_{\beta} \|y - X\beta\|_2^2 + \lambda \| \beta \|_2^2$$

$$\hat{\beta}_{RR} = \frac{\partial}{\partial \beta} (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

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

$$(X^T X + \lambda I_p) \beta = X^T y \rightarrow \hat{\beta}_{RR} = (X^T X + \lambda I_p)^{-1} X^T y$$

$$\hat{\beta}_{RR} = (\Phi \Sigma_r^2 \Phi^T + \lambda \Phi \Phi^T)^{-1} \Phi \Sigma_r S^T y$$

$$= [\Phi (\Sigma_r^2 + \lambda I_r) \Phi^T]^{-1} \Phi \Sigma_r S^T y$$

$$= \Phi (\Sigma_r^2 + \lambda I_r)^{-1} \Sigma_r S^T y = \sum_{i=1}^p \frac{\sigma_i}{\sigma_i^2 + \lambda} S_i^T y \phi_i$$

$$= \sum_{i=1}^p \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \frac{S_i^T y}{\sigma_i} \phi_i = \sum_{i=1}^p \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \hat{b}_i$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

3: $f(\hat{v}_i) = \begin{cases} 1 & \text{if } i \leq k < r \\ 0 & \text{if } i > k \end{cases}$

$$f(\hat{v}_i) = \frac{\hat{v}_i^2}{\hat{v}_i^2 + 1}$$

4:

$$\begin{aligned} E(\hat{\beta}_{RR}) &= E[(X^T X + \lambda I_p)^{-1} X^T y] = (X^T X + \lambda I_p)^{-1} X^T E(y) \\ &= (X^T X + \lambda I_p)^{-1} X^T X \beta = (X^T X + \lambda I_p)^{-1} (X^T X + \lambda I_p - \lambda I_p) \beta \\ &= \beta - \lambda (X^T X + \lambda I_p)^{-1} \beta \end{aligned}$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

5:

$$\text{Var}[\hat{\beta}_{RR}] = \text{Var}[(X^T X + \lambda I_p)^{-1} X^T y]$$

$$\begin{aligned} &= (X^T X + \lambda I_p)^{-1} X^T \text{Var}(y) X (X^T X + \lambda I_p)^{-1} \\ &= \sigma^2 (X^T X + \lambda I_p)^{-1} X^T X (X^T X + \lambda I_p)^{-1} = \sigma^2 W_\lambda (X^T X) W_\lambda^T \end{aligned}$$

$$W_\lambda = (X^T X + \lambda I_p)^{-1} X^T X, \quad W_\lambda^{-1} = (X^T X)^{-1} (X^T X + \lambda I_p) = I_p + \lambda (X^T X)^{-1}$$

$$\text{Var}(\hat{\beta}_{LS}) - \text{Var}(\hat{\beta}_{RR}) = \sigma^2 [(X^T X)^{-1} W_\lambda (X^T X) W_\lambda^T]$$

$$\sigma^2 W_\lambda \{ [I + \lambda (X^T X)^{-1}] (X^T X)^{-1} [I + \lambda (X^T X)^{-1}]^T - (X^T X)^{-1} \} W_\lambda^T$$

$$\sigma^2 W_\lambda \{ 2\lambda (X^T X)^{-2} + \lambda^2 (X^T X)^{-3} \} W_\lambda^T$$

$$= \sigma^2 (X^T X + \lambda I_p)^{-1} [2\lambda I_p + \lambda^2 (X^T X)^{-1}] (X^T X + \lambda I_p)^{-T}$$

(3)

The difference is non-negative definite as each component in the matrix product is non-negative definite. Then

$$\text{Var}(\hat{\beta}_{LS}) \geq \text{Var}(\hat{\beta}_{RR})$$

The variance of the ML estimator is larger than that of the ridge estimator.

$$\lim_{\lambda \rightarrow 0} \text{Var}(\hat{\beta}_{RR}) = \text{Var}(\hat{\beta}_{LS})$$

$$\lim_{\lambda \rightarrow \infty} \text{Var}(\hat{\beta}_{RR}) = \lim_{\lambda \rightarrow \infty} \sigma^2 W_\lambda (X^T X) W_\lambda^T = 0_p$$

<https://powcoder.com>

The variance of the ridge estimator decreases towards zero as the penalty parameter becomes large.

6: The ridge regression estimator is linear in y

$$\hat{\beta}_{RR} \sim N((X^T X + \lambda I_p)^{-1} X^T X \beta, \sigma^2 (X^T X + \lambda I_p)^{-1} X^T X [(X^T X + \lambda I_p)^{-1}]^T)$$

7: The MSE of the ridge regression estimator is

$$\begin{aligned} \text{MSE}(\hat{\beta}_{RR}) &= E[(W_\lambda \hat{\beta}_{LS} - \beta)(W_\lambda \hat{\beta}_{LS} - \beta)^T] \text{ we drop LS} \\ &= E[\hat{\beta}^T W_\lambda^T W_\lambda \hat{\beta}] - E[\beta^T W_\lambda^T W_\lambda \hat{\beta}] - E[\hat{\beta}^T W_\lambda^T W_\lambda \beta] + E[\beta^T \beta] \\ &= E[\hat{\beta}^T W_\lambda^T W_\lambda \hat{\beta}] - E[\beta^T W_\lambda^T W_\lambda \hat{\beta}] - E[\hat{\beta}^T W_\lambda^T W_\lambda \beta] + E[\beta^T W_\lambda^T W_\lambda \beta] \\ &\quad - E[\beta^T W_\lambda^T W_\lambda \beta] + E[\beta^T W_\lambda^T W_\lambda \hat{\beta}] + E[\hat{\beta}^T W_\lambda^T W_\lambda \beta] - E[\beta^T W_\lambda^T W_\lambda \hat{\beta}] \end{aligned}$$

$$- E[\hat{\beta}^T W_1^T B] + E[B^T B]$$

(4)

$$= E[(\hat{\beta} - \beta)^T W_1^T W_1 (\hat{\beta} - \beta)] - \beta^T W_1^T W_1 \beta + \beta^T W_1^T W_1 \beta + \beta^T W_1^T W_1 \beta - \beta^T W_1 \beta - \beta^T W_1^T B + \beta^T B$$

$$= E[(\hat{\beta} - \beta)^T W_1^T W_1 (\hat{\beta} - \beta)] + \beta^T (W_1 - I_p)^T (W_1 - I_p) \beta$$

$$= \sigma^2 \text{tr}[W_1 (X^T X)^{-1} W_1^T] + \beta^T (W_1 - I_p)^T (W_1 - I_p) \beta$$

where we used $E[\xi^T A \xi] = \text{tr}(A \Sigma_\xi) + \mu_\xi^T A \mu_\xi$ for $\xi \sim N(\mu_\xi, \Sigma_\xi)$.
 $\lim_{d \rightarrow \infty} \text{MSE}(\hat{\beta}_{RR}) = \beta^T B$ which is the minimal bound.
 8: In the case of orthonormal design matrix: $X^T X = I_p$

$$\bullet \hat{\beta}_{RR} = (X^T X + d I_p)^{-1} X^T y$$

$$= (X^T X + d I_p)^{-1} (X^T X) X^T y$$

$$= (I_p + d I_p)^{-1} \hat{\beta} = (1+d)^{-1} \hat{\beta}$$

$$\bullet E(\hat{\beta}_{RR}) = E\{(1+d)^{-1} \hat{\beta}\} = (1+d)^{-1} E(\hat{\beta}) = (1+d)^{-1} \beta \neq \beta$$

The estimator and its expectation vanishes as $d \rightarrow \infty$.

$$\bullet \text{Var}(\hat{\beta}_{RR}) = \sigma^2 W_1 (X^T X)^{-1} W_1^T = \sigma^2 (I_p + d I_p)^{-1} I_p [(I_p + d I_p)^{-1}]^T = \sigma^2 (1+d)^{-2} I_p$$

which vanishes as $d \rightarrow \infty$.

$$\bullet \text{MSE}(\hat{\beta}_{RR}) = \rho \sigma^2 (1+d)^{-2} + d (1+d)^{-2} \beta^T B$$

which is minimized for $d = \sigma^2 \beta^T B / \rho$

9 Increasing norm with decreasing λ

$$\begin{aligned}\|\hat{\beta}_{rr}\|_2^2 &= y^T S \Sigma_r (\Sigma_r^2 + \lambda I)^{-1} \Phi_r^T \Phi_r (\Sigma_r^2 + \lambda I)^{-1} \Sigma_r S^T y \\&= y^T S \Sigma_r (\Sigma_r^2 + \lambda I)^{-2} \Sigma_r S^T y \\&= (S^T y)^T \Sigma_r (\Sigma_r^2 + \lambda I)^{-2} \Sigma_r (S^T y) \\&= \sum_{i=1}^r \frac{\sigma_i^2}{(\sigma_i^2 + \lambda)^2} (S_i^T y)^2\end{aligned}$$

As $\lambda \rightarrow 0$, $\frac{\sigma_i^2}{(\sigma_i^2 + \lambda)^2}$ increases and also $\|\hat{\beta}_{rr}\|_2^2$.

Assignment Project Exam Help

Question 2:

1: The number of features of \hat{y} is affected

<https://powcoder.com>

Add WeChat powcoder

for least square: $\hat{y} = X(X^T X)^{-1} X^T y$

$$\begin{aligned}\text{cov}(\hat{y}, y) &= \text{cov}(X(X^T X)^{-1} X^T y, y) = X(X^T X)^{-1} X^T \text{cov}(y, y) \\&= X(X^T X)^{-1} X^T \sigma^2\end{aligned}$$

The values $\text{cov}(\hat{y}_i, y_i)$ are the diagonal values of the matrix given by $\text{cov}(\hat{y}, y)$

$$\sum_{i=1}^n \text{cov}(\hat{y}_i, y_i) = \text{tr}[X(X^T X)^{-1} X^T] \sigma^2 = \text{tr}[(X^T X)^{-1} X^T X] \sigma^2 = p \sigma^2$$

In the orthonormal design matrix case: $X^T X = I_p$

$$\sum_{i=1}^n \text{cov}(\hat{y}_i, y_i) = p \sigma^2$$

2. for ridge regression: $\hat{y} = X(X^T X + \lambda I)^{-1} X^T y$ (6)
 $\text{cov}(\hat{y}, y) = X(X^T X + \lambda I)^{-1} X^T \Sigma^{-1}$

$$\sum_{i=1}^n \text{cov}(\hat{y}_i, y_i) = \text{tr} \{ X(X^T X + \lambda I)^{-1} X^T \} \Sigma^{-1} = \text{tr} \{ (X^T X + \lambda I)^{-1} X^T X \} \Sigma^{-1}$$

For orthonormal design matrix

$$\sum_{i=1}^n \text{cov}(\hat{y}_i, y_i) = \frac{p \Sigma^{-1}}{1 + \lambda}$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder