

רגרסיה ומודלים סטטיסטיים- פתרון תרגיל 4

שאלה 1:

נתון שהמטריצה $X \in R^{n \times (p+1)}$ מדרגה מלאה. הניחו שהנחות המודל הלינארי מתקיימות.

חשבו את $cov(e)$, $cov(\epsilon, e)$, $cov(\hat{\beta}, \epsilon)$, $Cov(Y_i, \hat{Y}_i)$ וציינו את המימד של כל אחד מהגדלים.

פתרון:

$$cov(e) = cov((I - P_X)Y) = (I - P_X)cov(Y)(I - P_X)^T = \sigma^2(I - P_X) \in R^{n \times n}$$

$$\begin{aligned} cov(\epsilon, e) &= cov(\epsilon, (I - P_X)Y) = cov(\epsilon, Y)(I - P_X)^T = cov(\epsilon, X\beta + \epsilon)(I - P_X)^T \\ &= cov(\epsilon)(I - P_X) = \sigma^2(I - P_X) \in R^{n \times n} \end{aligned}$$

$$cov(\epsilon, \hat{\beta}) = cov(\epsilon, (X^T X)^{-1} X^T Y) = cov(\epsilon, X\beta + \epsilon) X (X^T X)^{-1} = \sigma^2 X (X^T X)^{-1} \in R^{n \times (p+1)}$$

$$cov(Y_i, \hat{Y}_i) = cov(Y, P_X Y)_{ii} = [cov(Y) P_X^T]_{ii} = \sigma^2 P_{Xii} \in R^{n \times n}$$

שאלה 2:

נתון מודל ליניארי, $Y = X\beta + \varepsilon$, כאשר $\varepsilon \sim (0, \sigma^2 I_n)$.

עבור $u, v \in \mathbb{R}^n$ כלשהם נגדיר את נורמת מהלנוביס (Mahalanobis) באופן הבא:

$$\|u - v\|_{\Sigma} = \sqrt{(u - v)^T \Sigma^{-1} (u - v)}$$

כאשר $\Sigma \in \mathbb{R}^{n \times n}$ מטריצה סימטרית חיובית.

כמו כן נגדיר את $\hat{\beta}^{\Sigma}$ באופן הבא:

$$\hat{\beta}^{\Sigma} = \underset{\beta \in \mathbb{R}^{p+1}}{\operatorname{argmin}} \|Y - X\beta\|_{\Sigma}^2$$

א. הראו שמתקיים

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

והסבירו את האומד שמתקבל במקרה $\Sigma = I_n$.

הדרכה: אפשר (לא חייב ש) להעזר בעובדה הבאה: $A \in \mathbb{R}^{m \times m}$ מטריצה סימטרית חיובית, אזי קיימת מטריצה $B \in \mathbb{R}^{m \times m}$ סימטרית חיובית כך ש: $B^T B = A$. ניתן להפעיל מסקנה זו על Σ כך שיהיה ניתן לייצג את $\Sigma^{-1} = C^T C$ באופן הבא: עבור מטריצה C ריבועית כלשהי. מצאו מי היא C , הסבירו מדוע היא מוגדרת היטב, ואז הגדירו משתנים חדשים

$$\tilde{Y} = CY, \quad \tilde{X} = CX$$

מדוע).

ב. הראו כי $\hat{\beta}^{\Sigma}$ הינו אומד חסר הטיות ל β , וכן מצאו את $\operatorname{cov}(\hat{\beta}^{\Sigma})$.

ג. מה אומר משפט גאוס מרקוב לגבי המטריצה $\operatorname{cov}(\hat{\beta}) - \operatorname{cov}(\hat{\beta}^{\Sigma})$, כאשר $\hat{\beta}$ אומד הריבועים הפחותים?

ד. כעת הניחו מודל ליניארי כללי יותר: $Y = X\beta + \varepsilon$ כאשר $\varepsilon \sim (0, \sigma^2 \Sigma)$ עבור אותה מטריצה סימטרית חיובית $\Sigma \in \mathbb{R}^{n \times n}$ שביחס אליה מוגדר האומד $\hat{\beta}^{\Sigma}$ (שימו לב ש Σ מטריצה סימטרית חיובית כלשהי, זהו אכן מודל כללי יותר שכן עבור $\Sigma = I_n$ מתקבל המודל "הרגיל").

הראו שבמקרה זה $\hat{\beta}^{\Sigma}$ הינו האומד הליניארי חסר ההטיות הטוב ביותר מבחינת שגיאה ריבועית מבין האומדים הליניאריים חסרי ההטיות ל β . כלומר, שלכל צירוף ליניארי $\theta = a^T \beta$ האומד $\hat{\theta}^{\Sigma} = a^T \hat{\beta}^{\Sigma}$ משיג שונות מינימאלית מבין כל האומדים הליניאריים חסרי ההטיות ל θ .

הדרכה: השתמשו בהדרכה מסעיף א', כתבו את המודל במונחי \tilde{Y} , \tilde{X} . הראו שמתקבל המודל הליניארי "הרגיל" ומשם ההוכחה שקולה להוכחת משפט גאוס מרקוב "הרגיל".

פתרון:

ללא 1.

6. Σ הוא מטריצה (ללא ריבוע) שלילית-החזקה, כלומר $\Sigma^T = -\Sigma$, $\Sigma^2 = -D$, $D = \text{diag}(d_1, \dots, d_n)$, $d_i > 0$.

$$(1) \quad \Sigma = U D U^T$$

כאשר $U \in \mathbb{R}^{n \times n}$ מטריצה יחידה, וכלומר $D = \text{diag}(d_1, \dots, d_n)$ ו- $d_i > 0$.

נציב

$$C := U D^{-1/2} U^T, \quad D^{-1/2} := \text{diag}(d_1^{-1/2}, \dots, d_n^{-1/2})$$

אז נקבל:

$$\begin{aligned} C^T C &= (U D^{-1/2} U^T)^T U D^{-1/2} U^T = U D^{-1/2} U^T U D^{-1/2} U^T = \\ &= U D^{-1} U^T = (U D U^T)^{-1} = \Sigma^{-1} \end{aligned}$$

ובנוסף נראה שהמטריצה U היא יחידה, כלומר $U^T = U^{-1}$.

נציב $\tilde{\Sigma} := C^T \Sigma C$ ו- $\tilde{X} := C X$.

$$(2) \quad \tilde{\Sigma} := C^T \Sigma C, \quad \tilde{X} := C X$$

אז נקבל:

$$\begin{aligned} \|Y - Xb\|_2^2 &= (Y - Xb)^T \Sigma^{-1} (Y - Xb) = \\ &= (Y - Xb)^T C^T C (Y - Xb) = \\ &= \|C(Y - Xb)\|_2^2 = \\ &= \|\tilde{Y} - \tilde{X}b\|_2^2 \end{aligned}$$

והכלל המכיל את המטריצה Σ^{-1} (המכילה) הוא $\Sigma^{-1} = -\Sigma$, ולכן נכתוב:

במקרה של הריבועים המכילים, אנו מקבלים:

$$\begin{aligned} \hat{\beta}^\Sigma &= \arg \min_b \|Y - Xb\|_2^2 = \arg \min_b \|\tilde{Y} - \tilde{X}b\|_2^2 = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \tilde{Y} = \\ &= [(CX)^T CX]^{-1} (CX)^T CY \end{aligned}$$

$$= (X^T C C X)^{-1} X^T C C Y$$

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

תוצאה. הוכחה: מילוי של המטרה (ההוכחה) מקובל
 שם $\hat{\beta}$ הוא הצורה הכללית של $\hat{\beta}$ — כלומר $X\beta$ —
הוכחה.

$$E[\hat{\beta}] = E[(X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y] = \quad (D)$$

$$= (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} E[Y] =$$

$$= (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} X \beta = \beta \quad (\forall \beta \in \mathbb{R}^{p+1}).$$

$$\text{cov}(\hat{\beta}) = \text{cov}[(X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y] =$$

$$= (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} [\text{cov}(Y)] \Sigma^{-1} X (X^T \Sigma^{-1} X)^{-1}$$

$$= \Sigma^2 (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} \Sigma^{-1} X (X^T \Sigma^{-1} X)^{-1}$$

$$= \Sigma^2 (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-2} X (X^T \Sigma^{-1} X)^{-1}$$

(E) . תוצאה: הוכחה של המטרה (ההוכחה) מקובל

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

הוא המאריך הנמוך ביותר של Y במרחב U , $U \subseteq \mathbb{R}^n$,
 שם U הוא המרחב הנמוך ביותר של Y במרחב \mathbb{R}^n , $U \subseteq \mathbb{R}^n$,
 $\alpha \in \mathbb{R}^{p+1}$

$$\text{Var}(\alpha^T \hat{\beta}) \leq \text{Var}(\alpha^T \tilde{\beta})$$

$$\Leftrightarrow \text{cov}(\alpha^T \hat{\beta}) \leq \text{cov}(\alpha^T \tilde{\beta})$$

$$\Leftrightarrow \alpha^T \text{cov}(\hat{\beta}) \alpha \leq \alpha^T \text{cov}(\tilde{\beta}) \alpha$$

$$\Leftrightarrow 0 \leq \alpha^T [\text{cov}(\tilde{\beta}) - \text{cov}(\hat{\beta})] \alpha$$

ואם $\alpha^T [\text{cov}(\tilde{\beta}) - \text{cov}(\hat{\beta})] \alpha \geq 0$ הרי $\alpha^T [\text{cov}(\tilde{\beta}) - \text{cov}(\hat{\beta})] \alpha \geq 0$
 שם $\alpha^T [\text{cov}(\tilde{\beta}) - \text{cov}(\hat{\beta})] \alpha \geq 0$ הרי $\alpha^T [\text{cov}(\tilde{\beta}) - \text{cov}(\hat{\beta})] \alpha \geq 0$

$$\text{cov}(\hat{\beta}^{\tilde{Y}}) - \text{cov}(\hat{\beta})$$

מכאן נובע:

(3) . דבריו המלאים:

$$Y = X\beta + e, \quad e \sim (0, \sigma^2 \Sigma)$$

הם נגזרים מ- Σ שבו נמצא המידע.

הם נגזרים מ- Σ שבו נמצא המידע, ו- \tilde{Y} כפי שהוצג בהשאלה (2).

לכן:

$$\tilde{Y} = CY = C(X\beta + e) = CX\beta + Ce = \tilde{X}\beta + \tilde{e}$$

כאן $\tilde{e} = Ce$, $E[\tilde{e}] = CE[E[e]] = 0$, $\text{cov}(\tilde{e}) = \text{cov}(Ce) = C[\text{cov}(e)]C^T =$

$$\text{cov}(\tilde{e}) = \text{cov}(Ce) = C[\text{cov}(e)]C^T =$$

$$= \sigma^2 C \Sigma C^T =$$

$$= \sigma^2 U D^{-1/2} U^T (U D U^T) (U D^{-1/2} U^T)^T$$

$$= \sigma^2 U D^{-1/2} U^T (U D U^T) U D^{-1/2} U^T$$

$$= \sigma^2 U \underbrace{(D^{-1/2} D D^{-1/2})}_{= I_n} U^T$$

$$= \sigma^2 U U^T = \sigma^2 I_n$$

המשפט המלא, דבריו המלאים, הם:

$$(3) \quad \tilde{Y} = \tilde{X}\beta + \tilde{e}, \quad \tilde{e} \sim (0, \sigma^2 I_n)$$

אם נניח שהמידע "המלא" הוא Σ , אז $\text{cov} = \sigma^2 I_n$.

אם β הוא וקטור $n \times 1$ ו- \tilde{X} הוא $n \times p$ (כפי שהוצג).

אם \tilde{e} הוא וקטור $n \times 1$ (כפי שהוצג) אז $\text{cov}(\tilde{e}) = \sigma^2 I_n$.

הם נגזרים מ- Σ .

$$(4) \quad \tilde{\beta} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \tilde{Y}$$

Q.2.

As seen in class, the multivariate MSE is given by:

$$E\|\hat{\theta} - \theta\|^2 = \text{tr}(\text{Var}(\hat{\theta})) + \|E(\hat{\theta}) - \theta\|^2$$

We already know that $\hat{\beta}_{OLS}$ is an unbiased estimator, and its variance matrix is given by $\sigma^2(X^T X)^{-1}$. From the previous calculations, we can compute the same terms for $\hat{\beta}_{Ridge}$:

$$\text{Var}(\hat{\beta}_{Ridge}) = \sigma^2(X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}$$

Where its bias^2 is:

$$\|(X^T X + \lambda I)^{-1} X^T X \beta - \beta\|^2$$

Now we'll read and scale the data for estimations:

```
df <- read.csv("C:/Users/nivbr/Downloads/quiz2_df.csv")

#scaling
n = length(df$Y)
fac = sqrt(n / (n - 1))
for (i in c(3:9)){
  df[,i] = as.numeric(scale(df[,i])) * fac
}
x_matrix = cbind(df$X0,df$X1,df$X2,df$X3,df$X4,df$X5,df$X6,df$X7)
```

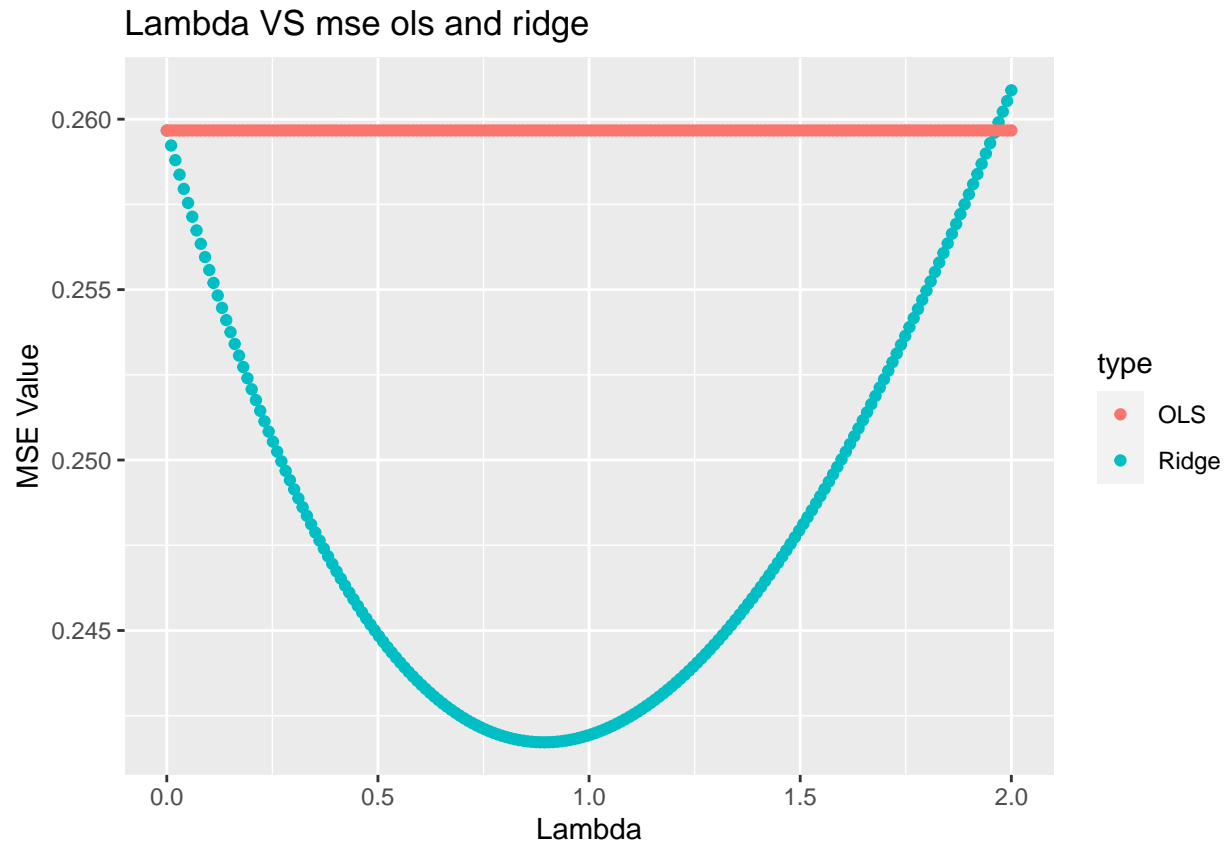
Estimations of for each value of λ with its corresponding MSE. Then, create the plot:

```
mse_ridge <- c()
for (i in lambda_seq){
  mse_ridge <- c(mse_ridge,ridge_aux_functions(x_matrix,df$Y,i,sigma_sq,beta)[2])
}

mse_ols <- c()
for (i in c(1:length(lambda_seq))){
  mse_ols <- c(mse_ols,ridge_aux_functions(x_matrix,df$Y,lambda_seq[i],sigma_sq,beta)[3])
}

type = c(rep("Ridge",200),rep("OLS",200))
value = c(unlist(mse_ridge),unlist(mse_ols))
data_for_ggplot <- data.frame(lambda = c(lambda_seq,lambda_seq),type = type,value=value )

ggplot(data_for_ggplot, aes(x=lambda, y = value, color=type)) +
  geom_point() + ylab("MSE Value") + xlab("Lambda") +
  labs(title="Lambda VS mse ols and ridge ")
```



We can see in the plot that the MSE of the OLS is constant (as it does not change with different λ values). When λ is zero, the OLS MSE and the Ridge MSE are identical. The shape of the Ridge-MSE graph resembles a parabola. This means that there is a certain λ from which (around 1) the Ridge MSE starts to increase. In other words, the trade-off between the bias and the variance flips, and the added bias becomes greater than the decreasing variance. Using Ridge regression with the given λ penalties (which are between 0 and 1) does seem to produce a smaller MSE in this particular model and might be preferable to OLS regression. (This actually happens because this particular data is highly multicollinear, which increases the estimator variances).

Next we'll examine the shrinking effect of the λ -s size on the $\|\hat{\beta}_{Ridge}\|$. As we saw previously, this norm is a monotonically decreasing function of λ :

```
beta_hat_func <- function(X, Y, lambda){
  I <- diag(ncol(X))
  beta_hat_ols <- solve(t(X)%*%X) %*% t(X)%*%Y
  beta_hat_ridge <- solve(t(X)%*%X + lambda*I) %*% t(X)%*%Y
  return(list(lambda = lambda, beta_hat_ridge = beta_hat_ridge)[2])
}

beta_hat_ridge <- c()
for (i in c(1:length(lambda_seq_mod))){
  beta_hat_ridge <- c(beta_hat_ridge, beta_hat_func(x_matrix, df$Y, lambda_seq_mod[i]))
}

#creating dataframe for the plot
```

```

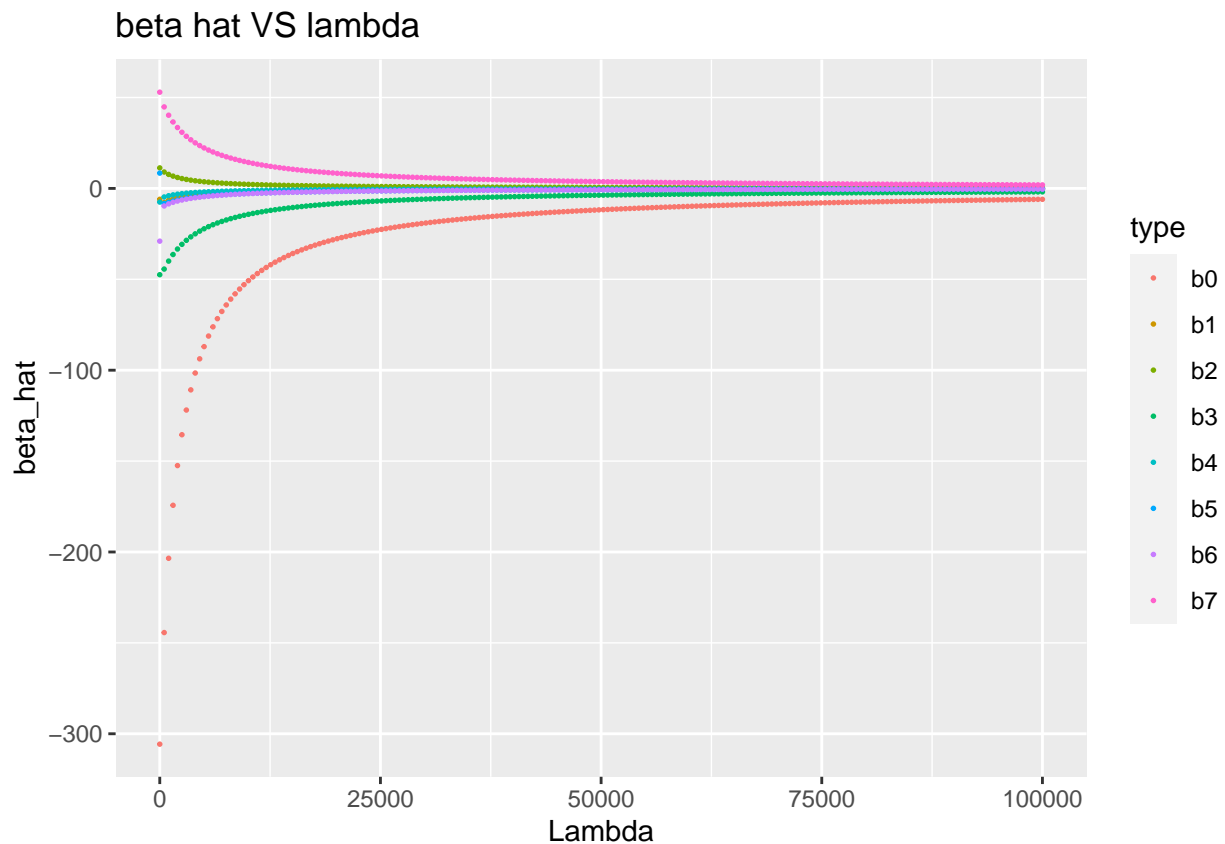
b0 <- sapply(beta_hat_ridge, "[", 1)
b1 <- sapply(beta_hat_ridge, "[", 2)
b2 <- sapply(beta_hat_ridge, "[", 3)
b3 <- sapply(beta_hat_ridge, "[", 4)
b4 <- sapply(beta_hat_ridge, "[", 5)
b5 <- sapply(beta_hat_ridge, "[", 6)
b6 <- sapply(beta_hat_ridge, "[", 7)
b7 <- sapply(beta_hat_ridge, "[", 8)

values <- c(b0,b1,b2,b3,b4,b5,b6,b7)
type = c(rep("b0",200),rep("b1",200),rep("b2",200),rep("b3",200),rep("b4",200),rep("b5",200),rep("b6",200),rep("b7",200))
lambda2 <- rep(lambda_seq_mod,8)

data_for_q12 <- data.frame(lambda = lambda2,type=type,values=values)

ggplot(data_for_q12, aes(x=lambda, y = values, color=type)) + geom_point(size=0.3) + ylab("beta_hat") +

```



as expected, we can see that as lambda increases, the beta's estimators converge to zero. the ridge regression puts a penalty lambda on beta. when the penalty increases, the absolute value of beta decreases, as the beta adds little contribution to the model.

Question 3. For Q3, we need to use the following results:

- Let C_1, C_2 be positive definite matrices with spectral decompositions $C_1 = UD_1U^T$ and $C_2 = UD_2U^T$ respectively. Then

$$C_1 + C_2 = U(D_1 + D_2)U^T,$$

and

$$C_1C_2 = U(D_1D_2)U^T.$$

- Moreover, write $D_1 = \text{diag}(d_1, \dots, d_p)$ and $D_2 = \text{diag}(\tilde{d}_1, \dots, \tilde{d}_p)$. Then

$$D_1 + D_2 = \text{diag}(d_1 + \tilde{d}_1, \dots, d_p + \tilde{d}_p),$$

and

$$D_1D_2 = \text{diag}(d_1\tilde{d}_1, \dots, d_p\tilde{d}_p).$$

- For any unitary U ($UU^T = I$), the following holds:

$$I = UIU^T$$

- As a consequence, $C_1^{-1} = UD_1^{-1}U^T$, and $D_1^{-1} = \text{diag}(d_1^{-1}, \dots, d_p^{-1})$.

(1) Note that,

$$X^TX + cI = UDU^T + cI = UDU^T + cUU^T = UD^*U^T$$

where

$$D^* = \begin{bmatrix} d_1 + c & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & d_p + c \end{bmatrix}$$

hence

$$\mathbb{E}(\hat{\beta}_{ridge}) = \mathbb{E}((X^TX + cI)^{-1}X^TY) = (X^TX + cI)^{-1}X^TX\beta + (X^TX + cI)^{-1}X^T\mathbb{E}[\varepsilon]$$

and taking $E[\varepsilon] = \mathbf{0}$ the second argument vanishes.

We get

$$(X^TX + cI)^{-1}X^TX\beta = (UD^*U^T)^{-1}UDU^T\beta = UD^{*-1}DU^T\beta,$$

where

$$D^{*-1}D = \begin{bmatrix} \frac{d_1}{d_1+c} & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \frac{d_p}{d_p+c} \end{bmatrix}.$$

(2)

$$\begin{aligned} \text{cov}(\hat{\beta}_{ridge}) &= \text{cov}((X^TX + cI)^{-1}X^T(X\beta + \varepsilon)) = \\ &= (X^TX + cI)^{-1}X^T \text{cov}(\varepsilon)((X^TX + cI)^{-1}X^T)^T = \sigma^2(X^TX + cI)^{-1}X^TX(X^TX + cI)^{-1} = \\ &= \sigma^2UD^{*-1}U^TUDU^TUD^{*-1}U^T = \sigma^2UD^{**}U^T, \end{aligned}$$

where

$$D^{**} = \begin{bmatrix} \frac{d_1}{(d_1+c)^2} & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \frac{d_p}{(d_p+c)^2} \end{bmatrix}$$

(3)

$$\text{var}(a^T \hat{\beta}_{ridge}) = \sigma^2 a^T U D^{**} U^T a$$

and

$$\text{var}(a^T \hat{\beta}_{OLS}) = \sigma^2 a^T (X^T X)^{-1} a.$$

So

$$\text{var}(a^T \hat{\beta}_{OLS}) - \text{var}(a^T \hat{\beta}_{ridge}) = \sigma^2 a^T U (D^{-1} - D^{**}) U^T a.$$

It is sufficient to prove that $U(D^{-1} - D^{**})U^T$ is positive definite. This is true if and only if the eigen-values are positive, meaning the elements on the diagonal of $(D^{-1} - D^{**})$ are positive.

The i 'th element in $(D^{-1} - D^{**})$ is of the form:

$$(D^{-1} - D^{**})_{ii} = \frac{1}{d_i} - \frac{d_i}{(d_i + c)^2}.$$

Note that for any $d_i, c > 0$:

$$\frac{d_i}{(d_i + c)^2} \leq \frac{d_i}{(d_i)^2} = \frac{1}{d_i}.$$

and the proof is complete.

According to Gauss Markov, $\hat{\beta}_{OLS}$ should be the best unbiased linear estimator for β , and in particular $a^T \hat{\beta}_{OLS}$ should be the minimal variance estimator for $a^T \beta$. However, $\hat{\beta}_{ridge}$ is biased (for any $\beta \neq \mathbf{0}$), and so is not covered by the conditions of Gauss-Markov.

To see the bias in $\hat{\beta}_{ridge}$, check the result from part 1:

$$\mathbb{E}[\hat{\beta}_{ridge}] = U D^{*-1} D U^T \beta.$$

We can therefore compute the bias:

$$\begin{aligned} \mathbb{E}[\hat{\beta}_{ridge} - \beta] &= U D^{*-1} D U^T \beta - U I U^T \beta = \\ &= U (D^{*-1} D - I) U^T \beta. \end{aligned}$$

Therefore, assuming $\beta \neq \mathbf{0}$, the bias is $\mathbf{0}$ if and only if $(D^{*-1} D - I) \equiv \mathbf{0}$. It's easy to see that the diagonal elements in $D^{*-1} D$ are $d_i/(d_i + c) \neq 1$ and the estimator is biased.