

Hw1

Runcheng Liu 2018010316

1.1

For SS_{tot} ,

$$\begin{aligned} SS_{tot} &= \sum_{i=1}^n (y_i - \bar{y})^2 = \begin{pmatrix} y_1 - \bar{y} & \cdots & y_n - \bar{y} \end{pmatrix} \begin{pmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{pmatrix} \\ &= (\mathbf{y}^T - \bar{y}\mathbf{1}^T) (\mathbf{y} - \bar{y}\mathbf{1}) \\ &= \mathbf{y}^T (\mathbf{I}_n - \mathbf{1}\mathbf{1}^T/n) \mathbf{y} \end{aligned}$$

For SS_{reg} ,

$$\begin{aligned} SS_{reg} &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \begin{pmatrix} \hat{y}_1 - \bar{y} & \cdots & \hat{y}_n - \bar{y} \end{pmatrix} \begin{pmatrix} \hat{y}_1 - \bar{y} \\ \vdots \\ \hat{y}_n - \bar{y} \end{pmatrix} \\ &= \hat{\mathbf{y}}^T (\mathbf{I}_n - \mathbf{1}\mathbf{1}^T/n) \hat{\mathbf{y}} \\ &= (\mathbf{P}\mathbf{y})^T (\mathbf{I}_n - \mathbf{1}\mathbf{1}^T/n) (\mathbf{P}\mathbf{y}) \\ &= \mathbf{y}^T \mathbf{P}^T (\mathbf{I}_n - \mathbf{1}\mathbf{1}^T/n) \mathbf{P} \mathbf{y} \end{aligned}$$

For SS_{res} ,

$$\begin{aligned} SS_{res} &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \begin{pmatrix} y_1 - \hat{y}_1 & \cdots & y_n - \hat{y}_n \end{pmatrix} \begin{pmatrix} y_1 - \hat{y}_1 \\ \vdots \\ y_n - \hat{y}_n \end{pmatrix} \\ &= (\mathbf{y}^T - \hat{\mathbf{y}}^T) (\mathbf{y} - \hat{\mathbf{y}}) \\ &= \mathbf{y}^T (\mathbf{I}_n - \mathbf{P}) \mathbf{y} \end{aligned}$$

Use the above three fomulas, thus:

$$\begin{aligned} SS_{reg} + SS_{res} &= \mathbf{y}^T \mathbf{P}^T (\mathbf{I}_n - \mathbf{1}\mathbf{1}^T/n) \mathbf{P} \mathbf{y} + \mathbf{y}^T (\mathbf{I}_n - \mathbf{P}) \mathbf{y} \\ &= \mathbf{y}^T (\mathbf{P}^T \mathbf{P} - \mathbf{P}^T \mathbf{1}\mathbf{1}^T/n \mathbf{P} + \mathbf{I}_n - \mathbf{P}) \mathbf{y} \\ &= \mathbf{y}^T (\mathbf{P} - (\mathbf{P}\mathbf{1})^T (\mathbf{P}\mathbf{1})/n + \mathbf{I}_n - \mathbf{P}) \mathbf{y} \\ &= \mathbf{y}^T (\mathbf{I}_n - \mathbf{1}\mathbf{1}^T/n) \mathbf{y} \\ &= SS_{tot} \end{aligned}$$

1.2

We know,

$$R^2 = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}$$

TSS measures the total variance in the response Y , and can be squares thought of as the amount of variability inherent in the response before the regression is performed. In contrast, RSS measures the amount of variability that is left unexplained after performing the regression. Hence, $TSS - RSS$ measures the amount of variability in the response that is explained (or removed) by performing the regression, and R^2 measures the proportion of variability in Y that can be explained using X . An R^2 statistic that is close to 1 indicates that a large proportion of the variability in the response has been explained by the regression. A number near 0 indicates that the regression did not explain much of the variability in the response; this might occur because the linear model is wrong, or the inherent error σ^2 is high, or both.

1.3

First, we know that,

$$\text{Var}(\hat{\beta}_j) = \sigma^2 \left((X^T X)^{-1} \right)_{jj}$$

Now let $r = X^T X$, and without losing generality, we reorder the columns of X to set the first column to be X_j ,

$$r^{-1} = \begin{bmatrix} r_{jj} & r_{j,-j} \\ r_{-j,j} & r_{-j,-j} \end{bmatrix}^{-1}$$

$$r_{jj} = X_j^T X_j, r_{j,-j} = X_j^T X_{-j}, r_{-j,j} = X_{-j}^T X_j, r_{-j,-j} = X_{-j}^T X_{-j}$$

By using Schur complement, the element in the first row and first column in r^{-1} is,

$$r_{1,1}^{-1} = \left[r_{jj} - r_{j,-j} r_{-j,-j}^{-1} r_{-j,j} \right]^{-1}$$

Then we have,

$$\text{Var}(\hat{\beta}_j) = \sigma^2 \left[(X^T X)^{-1} \right]_{jj} = \sigma^2 r_{1,1}^{-1}$$

$$= \sigma^2 \left[X_j^T X_j - X_j^T X_{-j} (X_{-j}^T X_{-j})^{-1} X_{-j}^T X_j \right]^{-1}$$

And we have,

$$RSS_j = \sum_{i=1}^n (X_{ij} - \bar{X}_j)^2 = X_j^T (\mathbf{I}_n - P) X_j = X_j^T X_j - X_j^T P X_j = X_j^T X_j - X_j^T X_j - X_j^T X_{-j} (X_{-j}^T X_{-j})^{-1} X_{-j}^T X_j$$

Put this into the last formula, so,

$$\text{Var}(\hat{\beta}_j) = \sigma^2 \left[X_j^T X_j - X_j^T X_{-j} (X_{-j}^T X_{-j})^{-1} X_{-j}^T X_j \right]^{-1}$$

$$= \sigma^2 \frac{1}{RSS_j}$$

$$= \frac{\sigma^2}{(n-1)S_{X_j}^2} \cdot \frac{1}{1 - R_j^2}$$

2

Prove $\hat{\beta}_0$ is BLUE is equivalent as proving $\hat{\beta}$ is BLUE. First, prove $\hat{\beta}$ is unbiased, use definition and the law of total expectation, we have,

$$\begin{aligned}
 E[\hat{\beta}] &= E\left[(X^T X)^{-1} X^T y\right] \\
 &= E\left[(X^T X)^{-1} X^T (X\beta + \varepsilon)\right] \\
 &= \beta + E\left[(X^T X)^{-1} X^T \varepsilon\right] \\
 &= \beta + E\left[E\left[(X^T X)^{-1} X^T \varepsilon | X\right]\right] \\
 &= \beta + E\left[(X^T X)^{-1} X^T E[\varepsilon | X]\right] \\
 &= \beta
 \end{aligned}$$

Where $E[\varepsilon | X] = 0$ by the assumption of the model. So $\hat{\beta}$ is unbiased.

Then prove $\hat{\beta}$'s variance is minimal, let $\tilde{\beta} = Cy$ be another linear estimator of β with $C = (X^T X)^{-1} X^T + D$ where D is a $K \times n$ non-zero matrix. As we're restricting to unbiased estimators, minimum mean squared error implies minimum variance. The goal is therefore to show that such an estimator has a variance no smaller than that of $\hat{\beta}$, the OLS estimator. We calculate,

$$\begin{aligned}
 E[\tilde{\beta}] &= E[Cy] \\
 &= E\left[\left((X'X)^{-1} X' + D\right) (X\beta + \varepsilon)\right] \\
 &= \left((X'X)^{-1} X' + D\right) X\beta + \left((X'X)^{-1} X' + D\right) E[\varepsilon] \\
 &\stackrel{E[\varepsilon]=0}{=} \left((X'X)^{-1} X' + D\right) X\beta \\
 &= (X'X)^{-1} X' X\beta + DX\beta \\
 &= (I_K + DX) \beta
 \end{aligned}$$

Therefore, since β is unobservable, and $\tilde{\beta}$ is unbiased if and only if $DX = 0$. Then,

$$\begin{aligned}
 \text{Var}(\tilde{\beta}) &= \text{Var}(Cy) \\
 &= C \text{Var}(y) C^T \\
 &= \sigma^2 C C^T \\
 &= \sigma^2 \left((X^T X)^{-1} X^T + D \right) \left(X (X^T X)^{-1} + D^T \right) \\
 &= \sigma^2 \left((X^T X)^{-1} X^T X (X^T X)^{-1} + (X^T X)^{-1} X^T D^T + DX (X^T X)^{-1} + DD^T \right) \\
 &= \sigma^2 (X^T X)^{-1} + \sigma^2 (X^T X)^{-1} (DX)^T + \sigma^2 DX (X^T X)^{-1} + \sigma^2 DD^T \\
 &\stackrel{DX=0}{=} \sigma^2 (X^T X)^{-1} + \sigma^2 DD^T \\
 &\stackrel{\text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}}{=} \text{Var}(\hat{\beta}) + \sigma^2 DD^T
 \end{aligned}$$

Since DD^T is positive semidefinite matrix. $\text{Var}(\tilde{\beta})$ is always bigger than $\text{Var}(\hat{\beta})$.

3.1

First, the ordinary least squares model posits that the conditioning distribution of the response y is,

$$y | X, \beta \sim \mathcal{N}(X\beta, \sigma^2 I)$$

Also we know the prior distribution of β ,

$$\beta \sim \mathcal{N}(0, \sigma^2 I / \lambda)$$

Thus,

$$\begin{aligned}\hat{\beta}^{MAP} &= \arg \max_{\beta} \prod_{i=1}^n P(Y_i | X_i, \beta) P(\beta) \\ &= \arg \max_{\beta} \log \left(\prod_{i=1}^n P(Y_i | X_i, \beta) P(\beta) \right) \\ &= \arg \max_{\beta} \left(\sum_{i=1}^n \log P(Y_i | X_i, \beta) + \log P(\beta) \right) \\ &= \arg \max_{\beta} \left(\sum_{i=1}^n -\frac{(Y_i - X_i^\top \beta)^2}{2\sigma^2 \mathbf{I}} - \frac{\lambda \|\beta\|_2^2}{2\mathbf{I}\sigma^2} \right) \\ &= \arg \min_{\beta} (Y_i - X_i^\top \beta)^2 + \lambda \|\beta\|_2^2 \\ &= \arg \min_{\beta} (X\beta - Y)^\top (X\beta - Y) + \lambda \|\beta\|_2^2\end{aligned}$$

3.2

As $\lambda \rightarrow 0$, the prior distribution $P(\beta)$ is very wide, which means that β has huge variance, thus is not constrained. However, as $\lambda \rightarrow \infty$, the prior distribution $P(\beta)$ is very narrow, which means that β has small variance, thus is constrained.

4

```
#Load data
prostate <- read.csv("prostate.csv")
train = subset(prostate, train==TRUE)
test = subset(prostate, train==FALSE)

#Split data into training data and test data
train_x=model.matrix(lpsa~. ,train[,c("lcavol","lweight","age","lbph","svi","lcp","gleason","pgg45","lpsa")])[, -1]
train_y=train$lpsa
test_x=model.matrix(lpsa~. ,test[,c("lcavol","lweight","age","lbph","svi","lcp","gleason","pgg45","lpsa")])[, -1]
test_y=test$lpsa
library(glmnet)
```

```
## Loading required package: Matrix
```

```
## Loaded glmnet 4.0-2
```

```
#Ridge regression
grid=seq(0,1, length =100)
ridge.mod=glmnet (train_x,train_y,alpha=0, lambda=grid, thresh =1e-12)

#5-fold Cross-validation
set.seed(1)
cv.out=cv.glmnet(train_x,train_y,alpha=0, nfolds = 5, lambda=grid, type.measure = "mse")
cvm_5 = rev(cv.out$cvm)

#Minimum MSE of lambda for 5-fold
bestlam_5 =cv.out$lambda.min
bestlam_5
```

```
## [1] 0.05050505
```

```
#10-fold Cross-validation
set.seed(1)
cv.out=cv.glmnet(train_x,train_y,alpha=0, nfolds = 10, lambda=grid, type.measure = "mse"
)
cvm_10 = rev(cv.out$cvm)

#Minimum MSE of lambda for 10-fold
bestlam_10 =cv.out$lambda.min
bestlam_10
```

```
## [1] 0.06060606
```

```

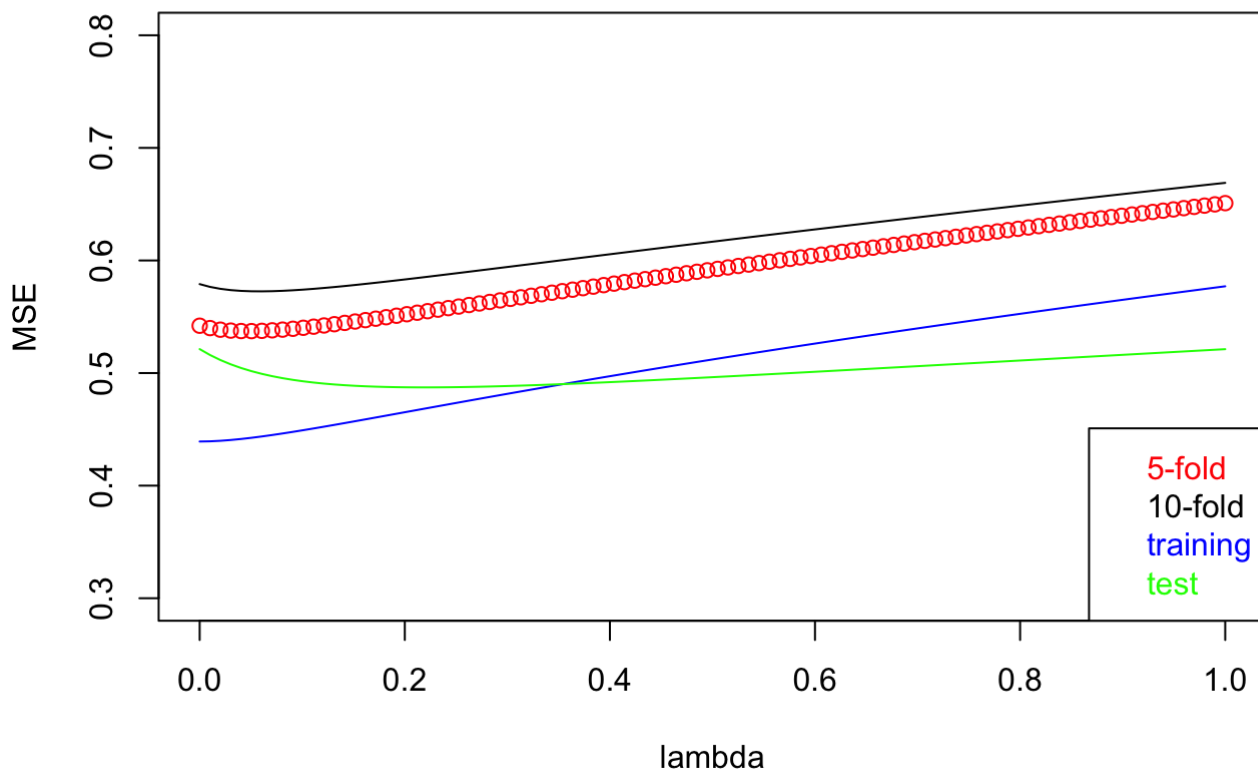
#Training error
train_error = c()
for (p in grid){
  ridge.pred=predict (ridge.mod ,s=p ,newx=train_x)
  train_error = c(mean((ridge.pred -train_y)^2), train_error)
}
train_error = rev(train_error)

#Test error
test_error = c()
for (p in grid){
  ridge.pred=predict (ridge.mod ,s=p ,newx=test_x)
  test_error = c(mean((ridge.pred -test_y)^2), test_error)
}
test_error = rev(test_error)

# PLOT
plot(grid, cvm_5, col="red", ylim = c(0.3, 0.8), main = "Prostate", xlab = "lambda", ylab="MSE")
lines(grid, cvm_10, col="black")
lines(grid,train_error,col="blue")
lines(grid,test_error,col="green")
legend("bottomright",legend=c("5-fold","10-fold","training","test"), text.col = c("red",
"black","blue","green"))

```

Prostate



From the above results, we can see that the training error, test error, 5-fold cross-validation error, 10-fold cross-validation for training, test, 5-fold cross-validation, 10-fold cross-validation respectively with respect to λ . I find that 5-fold cross-validation error, 10-fold cross-validation error are always larger than training error, test error. I think the reason for this is that the size of test and training data are larger than size of cross-validation data, thus the error is lower. I also find out that sometimes 5-fold cross-validation error is bigger than 10-fold cross-validation error, sometimes is reverse, unless you fix the random seed. For this random seed(which controls the splits of k-fold), the optimal λ for 5-fold and 10-fold cross-validation are 0.05 and 0.06 respectively.