



Semester 1 Assessment, 2021

School of Mathematics and Statistics

MAST30025 Linear Statistical Models Assignment 2

Submission deadline: **Friday April 30, 5pm**

This assignment consists of 4 pages (including this page)

Instructions to Students

Writing

- There are 5 questions with marks as shown. The total number of marks available is 40.
- This assignment is worth 7% of your total mark.
- You may choose to either typeset your assignment in L^AT_EX or handwrite and scan it to produce an electronic version.
- You may use R for this assignment, including the `lm` function unless specified. If you do, include your R commands and output.
- Write your answers on A4 paper. Page 1 should only have your student number, the subject code and the subject name. Write on one side of each sheet only. Each question should be on a new page. The question number must be written at the top of the page.

Scanning

- Put the pages in question order and all the same way up. Use a scanning app to scan all pages to PDF. Scan directly from above. Crop pages to A4. Check PDF is readable.

Submitting

- Go to the Gradescope window. Choose the Canvas assignment for this assignment. Submit your file as a single PDF document only. Get Gradescope confirmation on email.
- It is your responsibility to ensure that your assignments are submitted correctly and on time, and problems with online submissions are not a valid excuse for submitting a late or incorrect version of an assignment.

Question 1 (4 marks)

Prove Theorem 4.8: show that the maximum likelihood estimator of the error variance σ^2 is

$$\hat{\sigma}^2 = \frac{SS_{Res}}{n}.$$

Question 2 (11 marks)

We wish to predict the price of apartments in Melbourne using some of their features. Let y be the apartment price per square metre, x_1 be the apartment age (in years), x_2 be the distance (in metres) to the nearest train station, and x_3 be the number of convenience stores nearby.

The following data is collected:

<i>age</i> x_1	<i>dist</i> x_2	<i>#stores</i> x_3	y ($\times 10^3$)
32	84.9	10	37.9
19.5	306.6	9	42.2
13.3	562.0	5	47.3
13.3	562.0	5	43.1
5	390.6	5	54.8
7.1	2175.0	3	47.1
34.5	623.5	7	40.3

For this question, you may NOT use the `lm` function in R.

- (a) Fit a linear model to the data and estimate the parameters and variance.
- (b) Find a 90% confidence interval for the expected price per square metre of a 10 year old apartment that is 100 meters away from the train station and has 6 convenience stores nearby.
- (c) Find the standard error of $\beta_1 - \beta_3$.
- (d) Test the hypothesis that the price per square metre falls by \$1000 for every year that the apartment ages, at the 5% significance level.
- (e) Test for model relevance using a corrected sum of squares.

Question 3 (5 marks)

Consider two full rank linear models $\mathbf{y} = X_1\boldsymbol{\gamma}_1 + \boldsymbol{\varepsilon}_1$ and $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}_2$, where all predictors in the first model ($\boldsymbol{\gamma}_1$) are also contained in the second model ($\boldsymbol{\beta}$). Show that the SS_{Res} for the first model is at least the SS_{Res} for the second model.

Q1

$$y = x\beta + \varepsilon \sim MVN(x\beta, \sigma^2 I)$$

$$SS_{\text{Res}} = \sum_{i=1}^n \varepsilon_i^2$$

the log likelihood is

$$L(\beta, \sigma^2 | y) = -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma^2)$$

$$- \frac{1}{2\sigma^2} \|y - x\beta\|^T \|y - x\beta\|$$

$$\frac{\partial L(\beta, \sigma^2 | y)}{\partial \sigma^2} = \frac{1}{2} \|y - x\beta\|^T \|y - x\beta\| \cdot \sigma^{-4}$$

$$- \frac{N}{2} \sigma^{-2} = 0$$

$$\Rightarrow \frac{1}{2} \sigma^{-2} \left(\|y - x\beta\|^T \|y - x\beta\| \sigma^{-2} - N \right) = 0$$

$$\hat{\sigma}^2 = \frac{1}{N} \|y - x\beta\|^T \|y - x\beta\|$$

As $y - x\beta$ is the residual vector ε .

$$\|y - x\beta\|^T \|y - x\beta\| = \varepsilon^T \varepsilon = SS_{\text{Res}}$$

$$\text{Thus, } \hat{\sigma}^2 = \frac{1}{N} SS_{\text{Res}}.$$

182

assign2_lsm

Yuxin Ma

27/04/2021

#Q2 #(a)

```
y<- c(37.9, 42.2, 47.3, 43.1, 54.8, 47.1, 40.3)
x <- matrix(c(rep(1,7), 32, 19.5, 13.3, 13.3, 5, 7.1, 34.5, 84.9, 306.6, 562.0,
      562.0, 390.6, 2175.0, 623.5, 10, 9, 5, 5, 3, 7), 7, 4)
(b <- solve(t(x)%*%x, t(x)%*%y))
```

```
##           [,1]
## [1,] 58.369312708
## [2,] -0.346291960
## [3,] -0.002900359
## [4,] -0.887671692
```

```
#estimate variance
e <- y - x%*%b
SSRes <- sum(e^2)
(s2 <- SSRes/(7-4))
```

```
## [1] 13.06871
```

#{(b)}

```
s <- sqrt(s2)
xst <- c(1, 10, 100, 6)
ta <- qt(0.95, df = 7-4)
(xst%*%b - ta*s*sqrt(t(xst)%*%solve(t(x)%*%x)%*%xst))
```

```
##           [,1]
## [1,] 43.27252
```

```
(xst%*%b + ta*s*sqrt(t(xst)%*%solve(t(x)%*%x)%*%xst))
```

```
##           [,1]
## [1,] 55.30814
```

```
# the 90% CI is (43.27252, 55.30814)
```

#{(c)}

```

#  $se(\beta_1 - \beta_3) = \sqrt{var(\beta_1 - \beta_3)} = \sqrt{var(b_1 - b_3)}$ 
# matrix of variance b
varb <- solve(t(x) %*% x) * s2
varb1b3 <- varb[2, 2] + varb[4, 4] - 2 * varb[2, 4]
(seB1B3 <- sqrt(varb1b3))

## [1] 1.388968

#the standard error for  $\beta_1 - \beta_3$  is 1.388968

#(d)

#H0 :  $\beta_1 = -1$ 
c <- matrix(c(0, 1, 0, 0), 1, 4)
dst <- -1
r <- 1
num <- t(c %*% b - dst) %*% solve(c %*% solve(t(x) %*% x) %*% t(c)) %*% (c %*% b - dst) / r
Fstat <- num / (SSRes / (7 - 4))
(pf(Fstat, r, 7 - 4, lower = F))

## [,1]
## [1,] 0.04945829

#Since p-value = 0.0494 < 0.05, we can reject H0, that there's no evidence to show that price per square meter is irrelevant.

#(e)

# H0: model irrelevance with alpha = 5%
SSReg <- t(y) %*% x %*% b - 1 / 7 * sum(y)^2
MSReg <- SSReg / (4 - 1)
SSRes <- t(y) %*% y - t(y) %*% x %*% b
MSRes <- SSRes / (7 - 3 - 1)
Fstat <- MSReg / MSRes
(pf(Fstat, 3, 7 - 4, lower.tail = FALSE))

## [,1]
## [1,] 0.1500833

# We can NOT reject H0 as p-value = 0.15 > 0.05

knitr:::opts_chunk$set(echo = TRUE)

```

R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

Question 3 (5 marks)

Consider two full rank linear models $\mathbf{y} = \mathbf{X}_1\boldsymbol{\gamma}_1 + \boldsymbol{\varepsilon}_1$ and $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}_2$, where all predictors in the first model ($\boldsymbol{\gamma}_1$) are also contained in the second model ($\boldsymbol{\beta}$). Show that the SS_{Res} for the first model is at least the SS_{Res} for the second model.

As $\boldsymbol{\gamma}_1$ is a reduced model of $\boldsymbol{\beta}$, and $\boldsymbol{\beta}$ can be interpreted as $\begin{bmatrix} \boldsymbol{\gamma}_1 \\ \boldsymbol{\varepsilon}_2 \end{bmatrix}$, the order is not important.

$\boldsymbol{\gamma}_1$ is the reduced model of $\boldsymbol{\beta}$ that the 1st model can be written as $\mathbf{y} = \mathbf{X} \begin{bmatrix} \boldsymbol{\gamma}_1 \\ 0 \end{bmatrix} + \boldsymbol{\varepsilon}_1$.

As $\boldsymbol{\beta}$ is the best linear unbiased estimator, it minimize the error term, which also minimize the residual sum of squares, which is:

$$SS_{Res,2} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \boldsymbol{\varepsilon}_2^2$$

$$SS_{Res,1} = (\mathbf{y} - \mathbf{X}_1\boldsymbol{\gamma}_1)^T (\mathbf{y} - \mathbf{X}_1\boldsymbol{\gamma}_1) = \boldsymbol{\varepsilon}_1^2$$

Thus, $\boldsymbol{\varepsilon}_2^2$ has to be smaller or equal to $\boldsymbol{\varepsilon}_1^2$

Since $\boldsymbol{\gamma}_1$ is an estimator that assumed to be not as good as $\boldsymbol{\beta}$.

Q4

lsm_ass2Q4

Yuxin Ma

28/04/2021

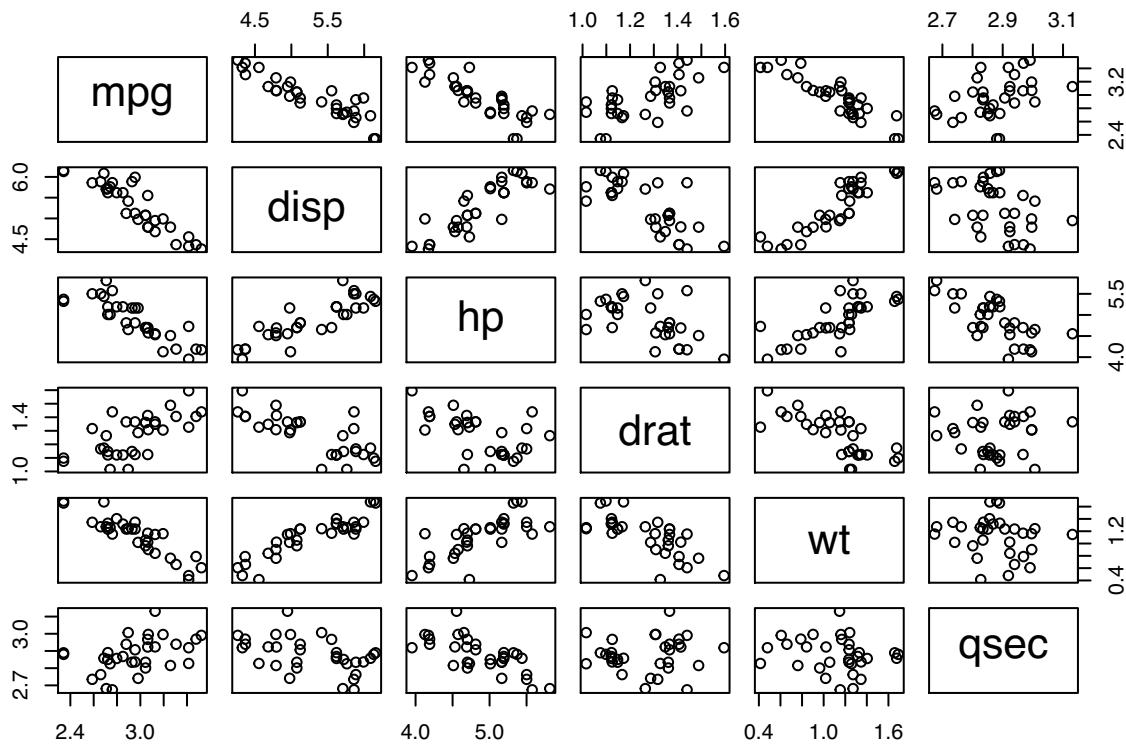
R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
#(a)
```

```
data(mtcars)
mtcars <- log(mtcars[, c(1, 3:7)])
y <- mtcars$mpg
disp <- mtcars$disp
hp <- mtcars$hp
drat <- mtcars$drat
wt <- mtcars$wt
qsec <- mtcars$qsec
pairs(mtcars)
```



#comment: the plots shows relationship between log-x-variables and response variable mpg. From the graphs, disp, #hp and wt present linear relationship with mpg whereas drat and qsec has no obvious linear relationship with mpg.

#{(b)}

```
basemodel <- lm(y~1, data = mtcars)
add1(basemodel, scope = ~.+disp+hp+drat+wt+qsec, test = "F")
```

```
## Single term additions
##
## Model:
## y ~ 1
##          Df Sum of Sq    RSS      AIC F value    Pr(>F)
## <none>        2.74874 -76.547
## disp     1   2.25596 0.49277 -129.550 137.3427 1.006e-12 ***
## hp       1   1.96733 0.78140 -114.797 75.5310 1.080e-09 ***
## drat     1   1.23131 1.51742  -93.559 24.3435 2.807e-05 ***
## wt       1   2.21452 0.53422 -126.966 124.3596 3.406e-12 ***
## qsec     1   0.47755 2.27119  -80.654  6.3079  0.01763 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

#disp is added to model2 as it has the largest F-value and has p-value < 0.05

```
model2 <- lm(y~disp, data = mtcars)
add1(model2, scope = ~.+hp+drat+wt+qsec, test = "F")
```

```

## Single term additions
##
## Model:
## y ~ disp
##      Df Sum of Sq    RSS    AIC F value Pr(>F)
## <none>          0.49277 -129.55
## hp     1  0.045531 0.44724 -130.65  2.9523 0.09641 .
## drat   1  0.001383 0.49139 -127.64  0.0816 0.77711
## wt     1  0.098796 0.39398 -134.71  7.2722 0.01154 *
## qsec   1  0.000308 0.49247 -127.57  0.0181 0.89382
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

# wt is added as it has the largest F-value among the remaining variables and has p-value < 0.05
model3 <- lm(y~disp+wt, data = mtcars)
add1(model3, scope = ~.+hp+drat+qsec, test = "F")

## Single term additions
##
## Model:
## y ~ disp + wt
##      Df Sum of Sq    RSS    AIC F value Pr(>F)
## <none>          0.39398 -134.71
## hp     1  0.078605 0.31537 -139.83  6.9789 0.01334 *
## drat   1  0.007358 0.38662 -133.31  0.5329 0.47146
## qsec   1  0.057788 0.33619 -137.79  4.8130 0.03671 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

#hp is added for the same reason as above
model4 <- lm(y~disp+wt+hp, data = mtcars)
add1(model4, scope = ~.+drat+qsec, test = "F")

## Single term additions
##
## Model:
## y ~ disp + wt + hp
##      Df Sum of Sq    RSS    AIC F value Pr(>F)
## <none>          0.31537 -139.83
## drat   1  0.0000095 0.31536 -137.83  0.0008 0.9774
## qsec   1  0.0033067 0.31206 -138.17  0.2861 0.5971

#no more variable will be added as they both have p-value>0.05 #the final model is y~disp+wt+hp
#(c)

fullmodel <- lm(y~disp+hp+drat+wt+qsec, data = mtcars)
model2 <- step(fullmodel, scope = ~., steps = 1)

## Start:  AIC=-136.21
## y ~ disp + hp + drat + wt + qsec
##
```

```

##          Df Sum of Sq      RSS      AIC
## - drat   1  0.000402 0.31207 -138.17
## - disp   1  0.002104 0.31377 -138.00
## - qsec   1  0.003699 0.31536 -137.83
## <none>           0.31166 -136.21
## - hp     1  0.023697 0.33536 -135.87
## - wt     1  0.103076 0.41474 -129.07
##
## Step:  AIC=-138.17
## y ~ disp + hp + wt + qsec

step(model2, scope=~.+drat)

## Start:  AIC=-138.17
## y ~ disp + hp + wt + qsec
##
##          Df Sum of Sq      RSS      AIC
## - qsec   1  0.003307 0.31537 -139.83
## - disp   1  0.004372 0.31644 -139.72
## <none>           0.31207 -138.17
## - hp     1  0.024123 0.33619 -137.79
## + drat   1  0.000402 0.31166 -136.21
## - wt     1  0.103779 0.41584 -130.98
##
## Step:  AIC=-139.83
## y ~ disp + hp + wt
##
##          Df Sum of Sq      RSS      AIC
## - disp   1  0.006635 0.32201 -141.16
## <none>           0.31537 -139.83
## + qsec   1  0.003307 0.31207 -138.17
## + drat   1  0.000010 0.31536 -137.83
## - hp     1  0.078605 0.39398 -134.71
## - wt     1  0.131870 0.44724 -130.65
##
## Step:  AIC=-141.17
## y ~ hp + wt
##
##          Df Sum of Sq      RSS      AIC
## <none>           0.32201 -141.16
## + disp   1  0.00664 0.31537 -139.83
## + qsec   1  0.00557 0.31644 -139.72
## + drat   1  0.00112 0.32089 -139.28
## - hp     1  0.21221 0.53422 -126.97
## - wt     1  0.45939 0.78140 -114.80

##
## Call:
## lm(formula = y ~ hp + wt, data = mtcars)
##
## Coefficients:
## (Intercept)          hp          wt
##        4.8347     -0.2553     -0.5623

```

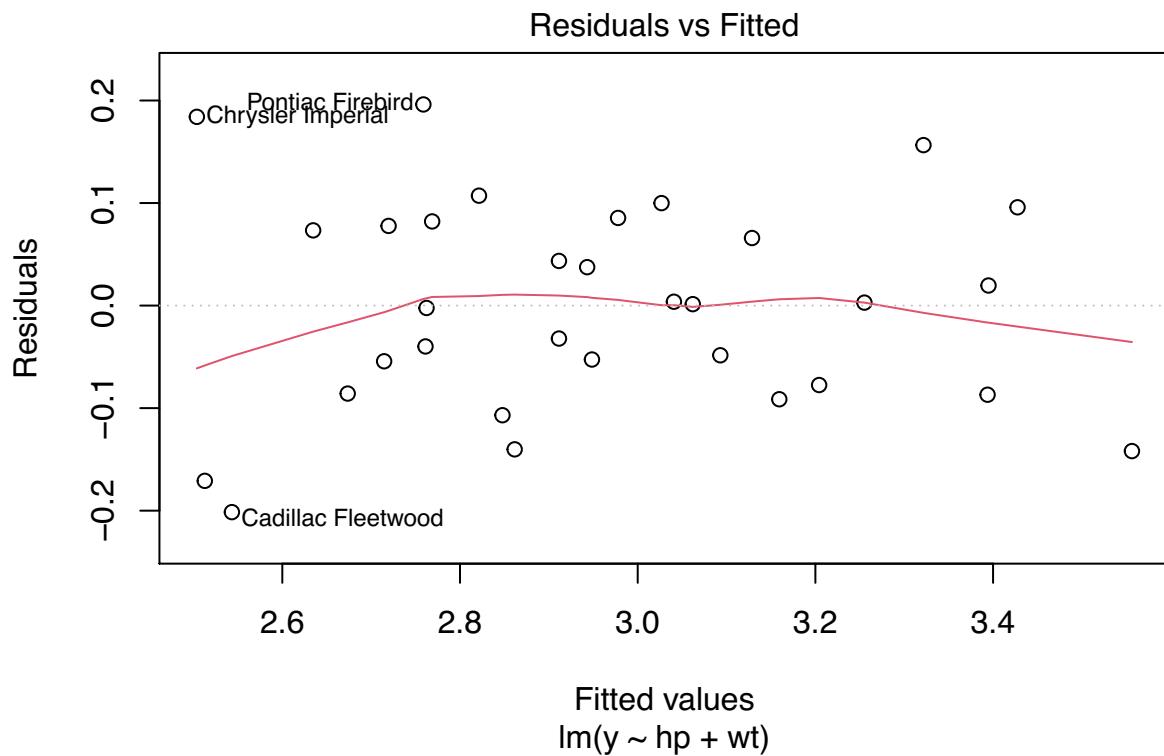
At last we have hp and wt included in the final linear model with stepwise function, starting from full model.

#(d)

```
# with the coefficient listed above, we have:  
#mpg = 4.8347 - 0.2553*hp - 0.5623*wt
```

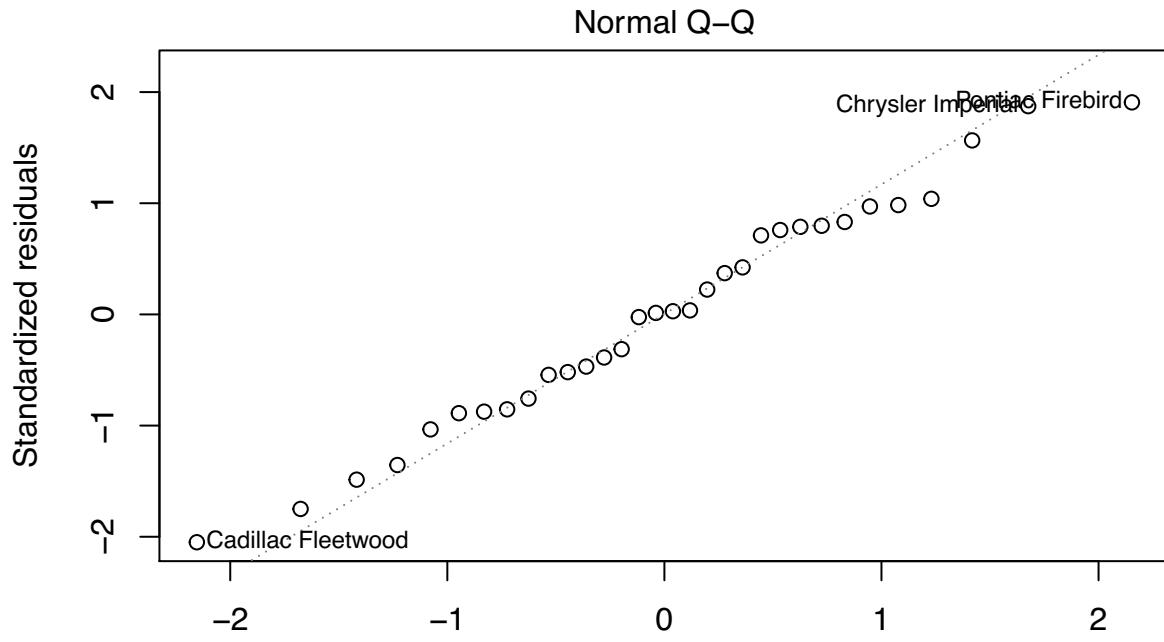
#(e)

```
finalmodel <- lm(y ~ hp + wt, data = mtcars)  
plot(finalmodel, which = 1)
```



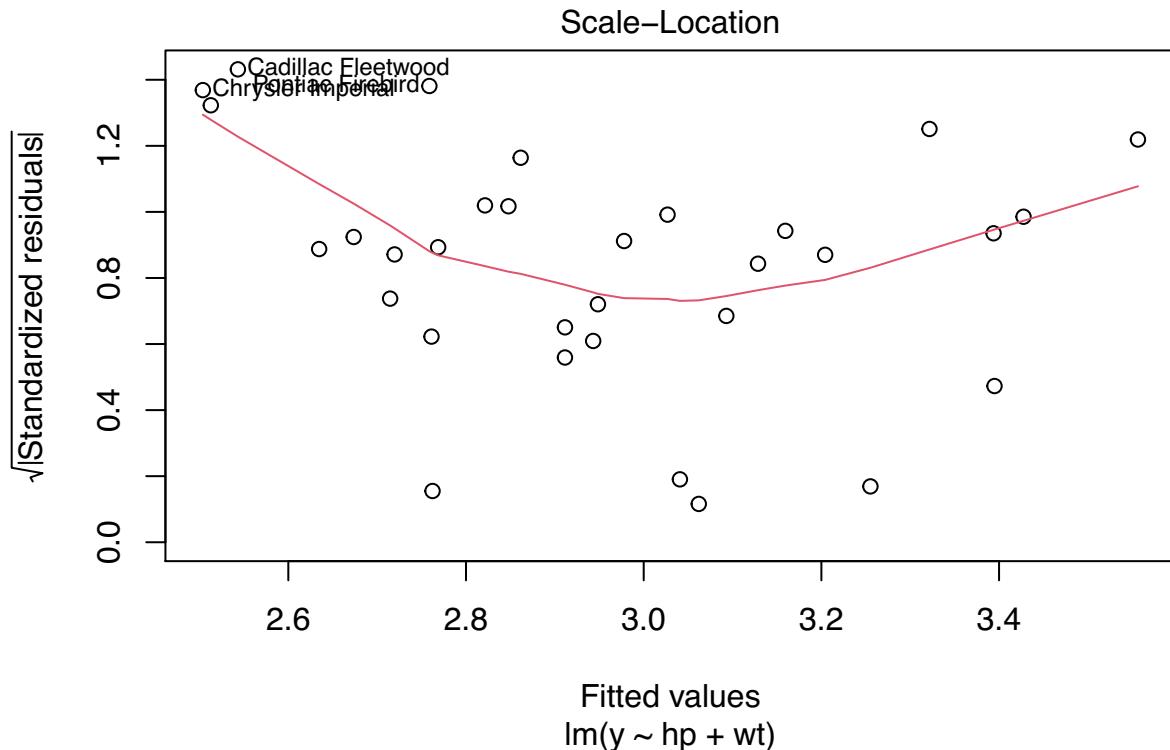
```
# We can observe equally spread residuals around a horizontal line without distinct patterns and even  
# outliers # have small residuals. The model is fitted well.
```

```
plot(finalmodel, which=2)
```



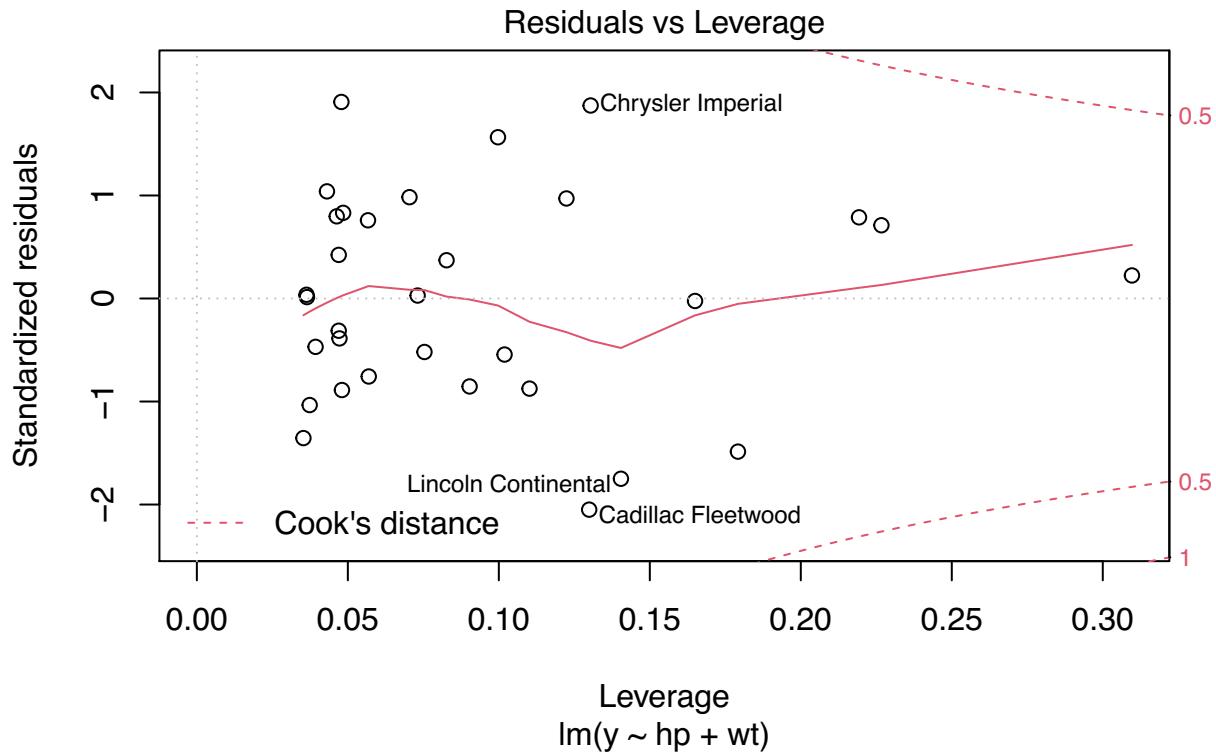
The points follow the dashed straight line with a small number of outliers. Thus it's a good fit.

```
plot(finalmodel, which=3)
```



This plot shows if residuals are spread equally along the ranges of predictors. # We can observe a generally horizontal line with equally randomly spread points with low residuals. # Thus, it is generally a good fit.

```
plot(finalmodel, which = 5)
```



the residuals, leverage and cook's distance are generally small. # Accorrding to the diagnostic plots, the final model computed from stepwise selection is acceptable.

Question 5 (10 marks)

For ridge regression, we choose parameter estimators \mathbf{b} which minimise

$$\sum_{i=1}^n e_i^2 + \lambda \sum_{j=0}^k b_j^2,$$

where λ is a constant penalty parameter.

- (a) Show that these estimators are given by

$$\mathbf{b} = (X^T X + \lambda I)^{-1} X^T \mathbf{y}.$$

- (b) Show that \mathbf{b} is biased if $\lambda \neq 0$.

- (c) One way to calculate the optimal value for the penalty parameter is to minimise the AIC.
Since the number of parameters p does not change, we use a slightly modified version:

$$AIC = n \ln \frac{SS_{Res}}{n} + 2 df,$$

where df is the “effective degrees of freedom” defined by

$$df = \text{tr}(H) = \text{tr}(X(X^T X + \lambda I)^{-1} X^T).$$

We will use the data from Q2. In order to avoid penalising some parameters unfairly, we must first standardise the variables; this also means an intercept parameter is not used. You can do this with `scale`:

```
> X <- scale(X[,-1], center=T, scale=T)
> y <- scale(y, center=T, scale=T)
> p <- 3
```

Construct a plot of λ against AIC. Thereby find the optimal value for λ .

End of Assignment — Total Available Marks = 40

$$\begin{aligned}
 a) \quad \text{Let } w &= \sum_{i=1}^n e_i^2 + \lambda \sum_{j=0}^k b_j^2 \\
 &= (y - Xb)^T (y - Xb) + \lambda b^T b \\
 &= (y^T - b^T X^T)(y - Xb) + \lambda b^T b \\
 &= \underline{y^T y - 2y^T X b + b^T X^T X b + \lambda I b^T b}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial w}{\partial b} &= 0 - 2y^T X + 2X^T X b + 2\lambda I b \\
 &= -2X^T y + 2X^T X b + 2\lambda I b \\
 &= 0
 \end{aligned}$$

According to handout & pg 0.

$$(2X^T X + 2\lambda I) b = 2X^T y$$

$$b = (X^T X + \lambda I)^{-1} X^T y$$

$$b = (X^T X + \lambda I)^{-1} X^T y$$

$$E(b) = E[(X^T X + \lambda I)^{-1} X^T y]$$

$$= E[(X^T X + \lambda I)^{-1} X^T X \beta]$$

If b is unbiased, $E(b) = \beta$.

It is shown that $(X^T X + \lambda I)^{-1} (X^T X \beta) = \beta$;

which means if $\lambda I = 0$, the estimator

is unbiased.

Thus, if $\lambda \neq 0$, b is a biased estimator.

lsm_ass2_Q5

Yuxin Ma

28/04/2021

R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

(c)

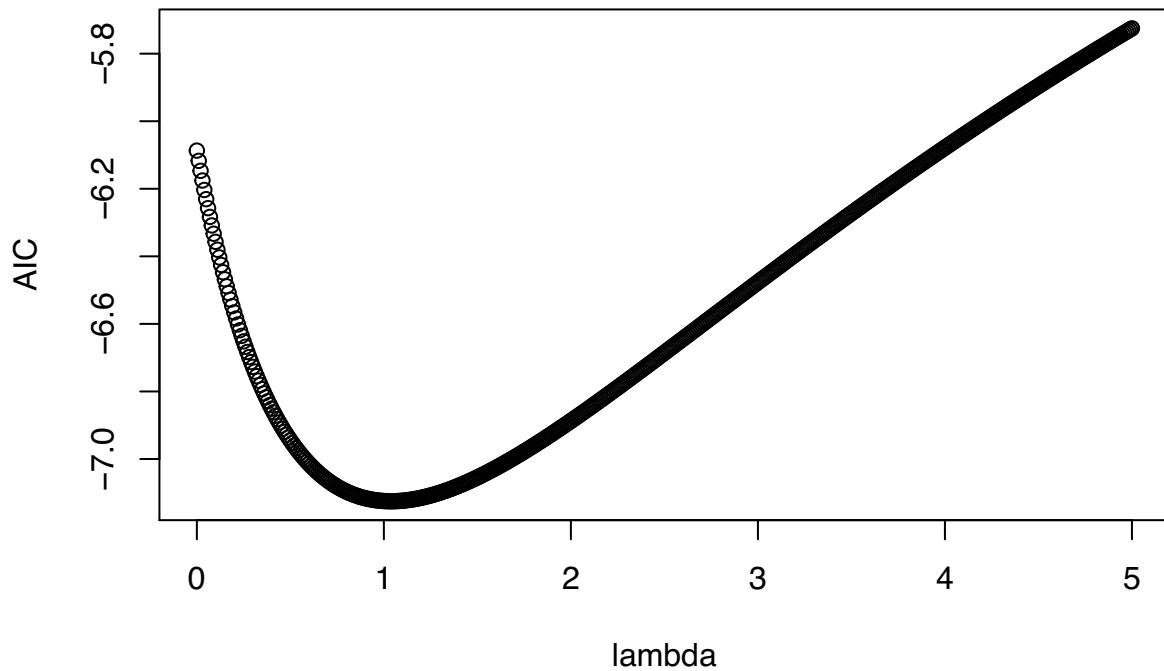
```
y<- c(37.9, 42.2, 47.3, 43.1, 54.8, 47.1, 40.3)
X <- matrix(c(rep(1,7), 32, 19.5, 13.3, 13.3, 5, 7.1, 34.5, 84.9, 306.6, 562.0,
               562.0, 390.6, 2175.0, 623.5, 10, 9, 5, 5, 3, 7), 7, 4)
X <- scale(X[,-1],center=T,scale=T)
y <- scale(y,center=T,scale=T)
p <- 3

n <- length(y)

# assume lambda are all positive.
lambda <- seq(0, 5, 0.01)
AIC <- replicate(length(lambda), 0)

for (i in 1:length(AIC)){
  b <- solve(t(X) %*% X + diag(lambda[i], p), t(X) %*% y)
  e <- y - X %*% b
  SSRes <- sum(e^2)
  H <- X %*% solve(t(X) %*% X + diag(lambda[i], p)) %*% t(X)
  df <- sum(diag(H))
  AIC[i] <- n*log(SSRes/n) + 2*df
}

plot(lambda, AIC)
```



```
#lambda with the smallest AIC
lambda[which(AIC == min(AIC))]
```

```
## [1] 1.04
```