

Practice Final Exam

Instructions

Open book and open notes exam (textbooks (print or pdf), lecture slides, notes, practice exam, homework solutions, and TA slides, including all Rmd's*).

You are allowed to use RStudio Cloud (<https://rstudio.cloud>) , Microsoft Word, Power Point and PDF reader, and canvas on your laptop.

Proctorio is required to start this exam. If you are prompted for an access code, you must Configure Proctorio on your machine.

Please read the list of recording and restrictions provided by Proctorio carefully before taking the exam

Please pay attention any timing and technical warnings that popped up your screen

The exam will be available from Monday December 14th at 8 am EST through Tuesday December 15th at 8:00pm EST. Once you start the exam, you have to complete the exam in 3 hours or by Tuesday December 15th at 8:00pm EST, whichever comes first.

In order to receive full credit, please provide full explanations and calculations for each questions

Make sure that you are familiar with the procedures for troubleshooting exam issues Preview the document Make sure you submit both .Rmd and (knitted) pdf or html files.

You need to have a camera on your laptop.

Problem 1

Use the question1 data, fit the regression model on Y by using all the variables (X6 and X7 are categorical variables). Create development sample (70% of the data) and hold-out sample (30% of the data). Perform statistical tests, use graphs or calculate the measures (e.g. VIF, Leverage Points, Cook's Distance) for questions below. Use the development sample for part a to d. Use the hold-out sample for part e.

a-) Is the model significant? Is there a Multicollinearity in the data? Are the errors Normally distributed with constant variance?

X_8 is significantly correlated with X_9 and X_{10}

X_1 , X_2 , X_7 , X_8 , and X_9 are significant variables. R^2 is 70%. QQ plot indicates S function and suggests that Y should be transformed. Errors have constant variance. There are outliers in the data set. There is a multicollinearity issue in the data set. X_8 should be dropped from the model.

```
library(datasets)
```

```
PF.Q1.Dat <- read.csv("/cloud/project/Practice Final Question 1.csv")
round(cor(PF.Q1.Dat),2)
```

```
##          Y      X1      X2      X3      X4      X5      X6      X7      X8      X9      X10
## Y      1.00  0.19  0.53  0.33  0.38  0.41 -0.30 -0.49  0.47  0.34  0.36
```

```
## X1  0.19  1.00  0.00 -0.23 -0.02 -0.06  0.15 -0.02 -0.05 -0.08 -0.04
## X2  0.53  0.00  1.00  0.56  0.45  0.36 -0.23 -0.19  0.38  0.39  0.41
## X3  0.33 -0.23  0.56  1.00  0.42  0.14 -0.24 -0.31  0.14  0.20  0.19
## X4  0.38 -0.02  0.45  0.42  1.00  0.05 -0.09 -0.30  0.06  0.08  0.11
## X5  0.41 -0.06  0.36  0.14  0.05  1.00 -0.59 -0.11  0.98  0.92  0.79
## X6 -0.30  0.15 -0.23 -0.24 -0.09 -0.59  1.00  0.10 -0.61 -0.59 -0.52
## X7 -0.49 -0.02 -0.19 -0.31 -0.30 -0.11  0.10  1.00 -0.15 -0.11 -0.21
## X8  0.47 -0.05  0.38  0.14  0.06  0.98 -0.61 -0.15  1.00  0.91  0.78
## X9  0.34 -0.08  0.39  0.20  0.08  0.92 -0.59 -0.11  0.91  1.00  0.78
## X10 0.36 -0.04  0.41  0.19  0.11  0.79 -0.52 -0.21  0.78  0.78  1.00
```

```
set.seed(1234)
n<-dim(PF.Q1.Dat)[1]
set.seed(994)
#dummy variables
table(PF.Q1.Dat$X6)
```

```
##
##  1  2
## 17 96
```

```
table(PF.Q1.Dat$X7)
```

```
##
##  1  2  3  4
## 28 32 37 16
```

```
#X7 gas 4 levels and need to create dummy variables or use factor command
#creating the dummy variables
```

```
library('fastDummies')
Q1.Dat<-dummy_cols(PF.Q1.Dat, select_columns = 'X7')
#this will create dummy variables for all levels, but we need to drop one level and drop X7
head(Q1.Dat)
```

```
##      Y   X1  X2   X3   X4  X5 X6 X7  X8  X9 X10 X7_1 X7_2 X7_3 X7_4
## 1  7.13 55.7 4.1  9.0  39.6 279  2  4 207 241  60    0    0    0    1
## 2  8.82 58.2 1.6  3.8  51.7  80  2  2  51  52  40    0    1    0    0
## 3  8.34 56.9 2.7  8.1  74.0 107  2  3  82  54  20    0    0    1    0
## 4  8.95 53.7 5.6 18.9 122.8 147  2  4  53 148  40    0    0    0    1
## 5 11.20 56.5 5.7 34.5  88.9 180  2  1 134 151  40    1    0    0    0
## 6  9.76 50.9 5.1 21.9  97.0 150  2  2 147 106  40    0    1    0    0
```

```
Q1.Dat<-Q1.Dat[,-c(8,12)]
head(Q1.Dat)
```

```
##      Y   X1  X2   X3   X4  X5 X6  X8  X9 X10 X7_2 X7_3 X7_4
## 1  7.13 55.7 4.1  9.0  39.6 279  2 207 241  60    0    0    1
## 2  8.82 58.2 1.6  3.8  51.7  80  2  51  52  40    1    0    0
## 3  8.34 56.9 2.7  8.1  74.0 107  2  82  54  20    0    1    0
## 4  8.95 53.7 5.6 18.9 122.8 147  2  53 148  40    0    0    1
## 5 11.20 56.5 5.7 34.5  88.9 180  2 134 151  40    0    0    0
## 6  9.76 50.9 5.1 21.9  97.0 150  2 147 106  40    1    0    0
```

```
IND=sample(c(1:n),n*0.7)
Q1.Dev<-Q1.Dat[IND,]
Q1.Hold<-Q1.Dat[-IND,]
f1<-lm(Y~ X1+X2+X3+X4+X5+X6+X7_2+X7_3+X7_4+X8+X9+X10 ,data=Q1.Dev)
summary(f1)
```

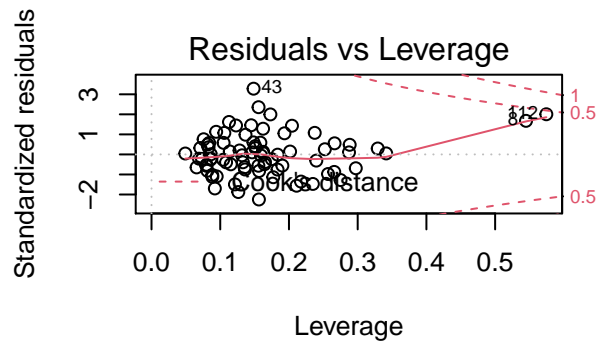
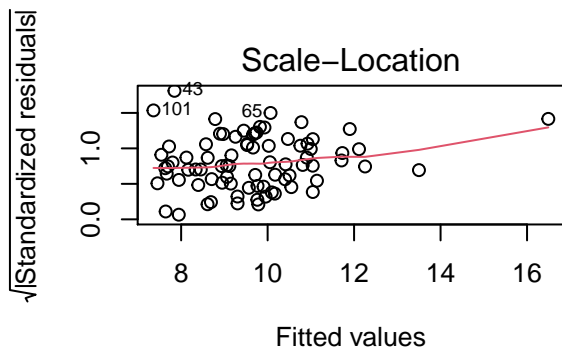
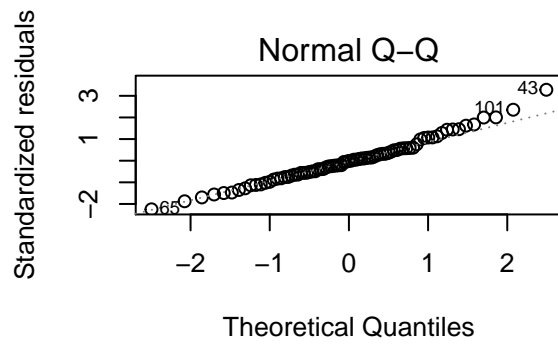
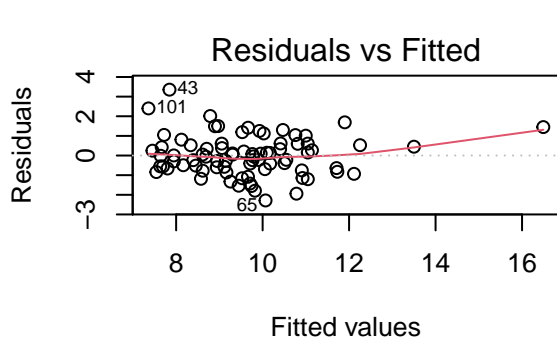
```
##
## Call:
## lm(formula = Y ~ X1 + X2 + X3 + X4 + X5 + X6 + X7_2 + X7_3 +
##      X7_4 + X8 + X9 + X10, data = Q1.Dev)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.2878 -0.6504 -0.0122  0.5620  3.3538
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.878977   2.238905   1.286 0.202976
## X1           0.085897   0.034214   2.511 0.014511 *
## X2           0.345730   0.131489   2.629 0.010633 *
## X3           0.021896   0.017803   1.230 0.223113
## X4           0.009959   0.008381   1.188 0.238974
## X5          -0.008119   0.003792  -2.141 0.035973 *
## X6          -0.015940   0.469810  -0.034 0.973037
## X7_2         -0.661010   0.384994  -1.717 0.090681 .
## X7_3         -0.742174   0.393805  -1.885 0.063885 .
## X7_4        -1.642060   0.522938  -3.140 0.002527 **
## X8           0.018730   0.004614   4.060 0.000133 ***
## X9          -0.004127   0.002443  -1.689 0.095896 .
## X10         -0.017292   0.014916  -1.159 0.250516
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.108 on 66 degrees of freedom
## Multiple R-squared:  0.6842, Adjusted R-squared:  0.6268
## F-statistic: 11.92 on 12 and 66 DF, p-value: 2.436e-12

par(mfrow=c(2,2))
plot(f1)

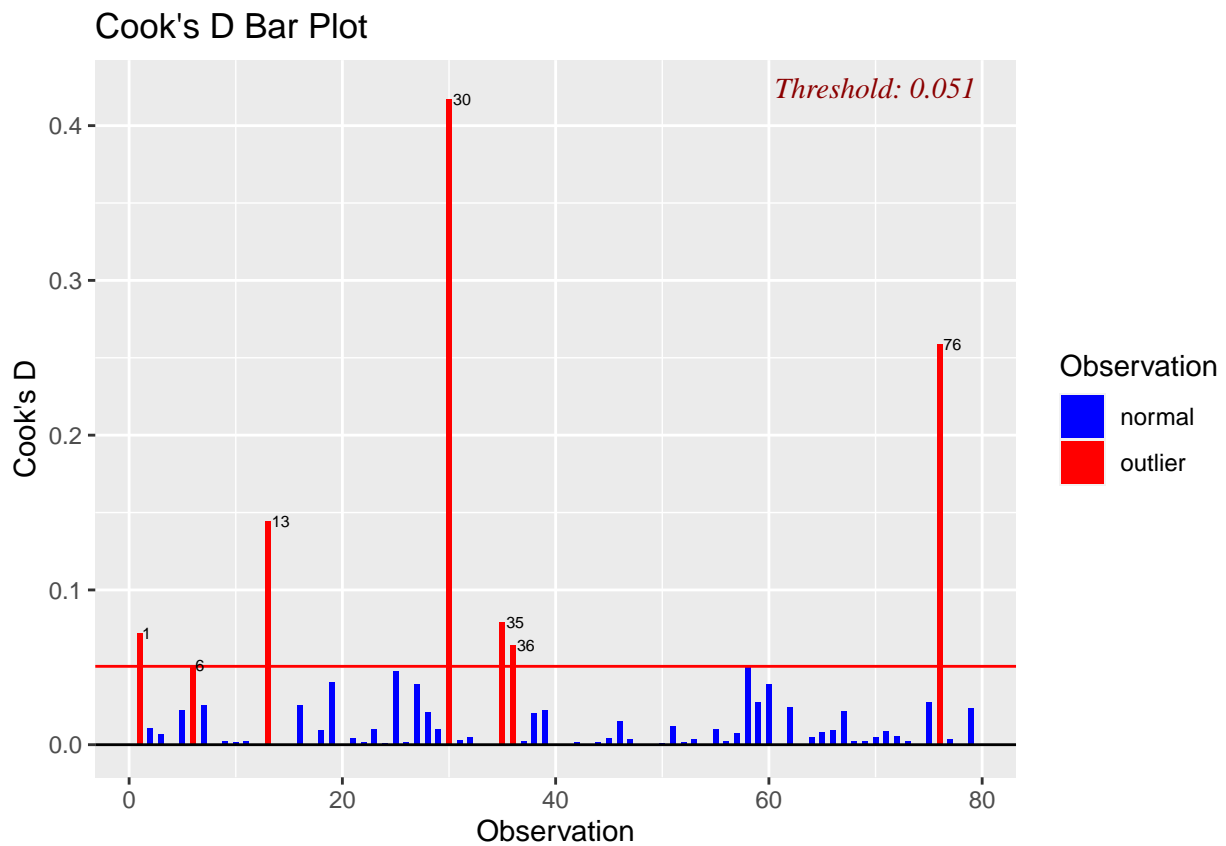
library(olsrr)

##
## Attaching package: 'olsrr'

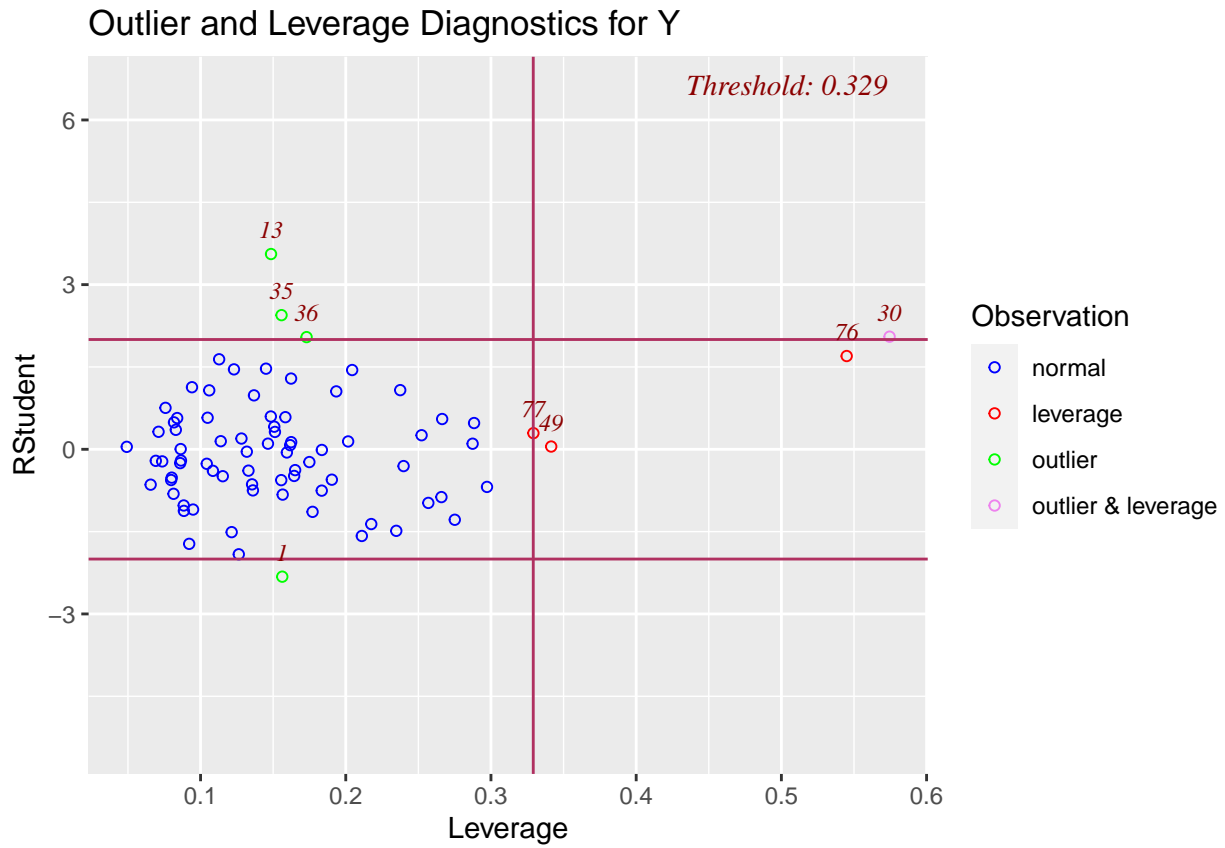
## The following object is masked from 'package:datasets':
##
##      rivers
```



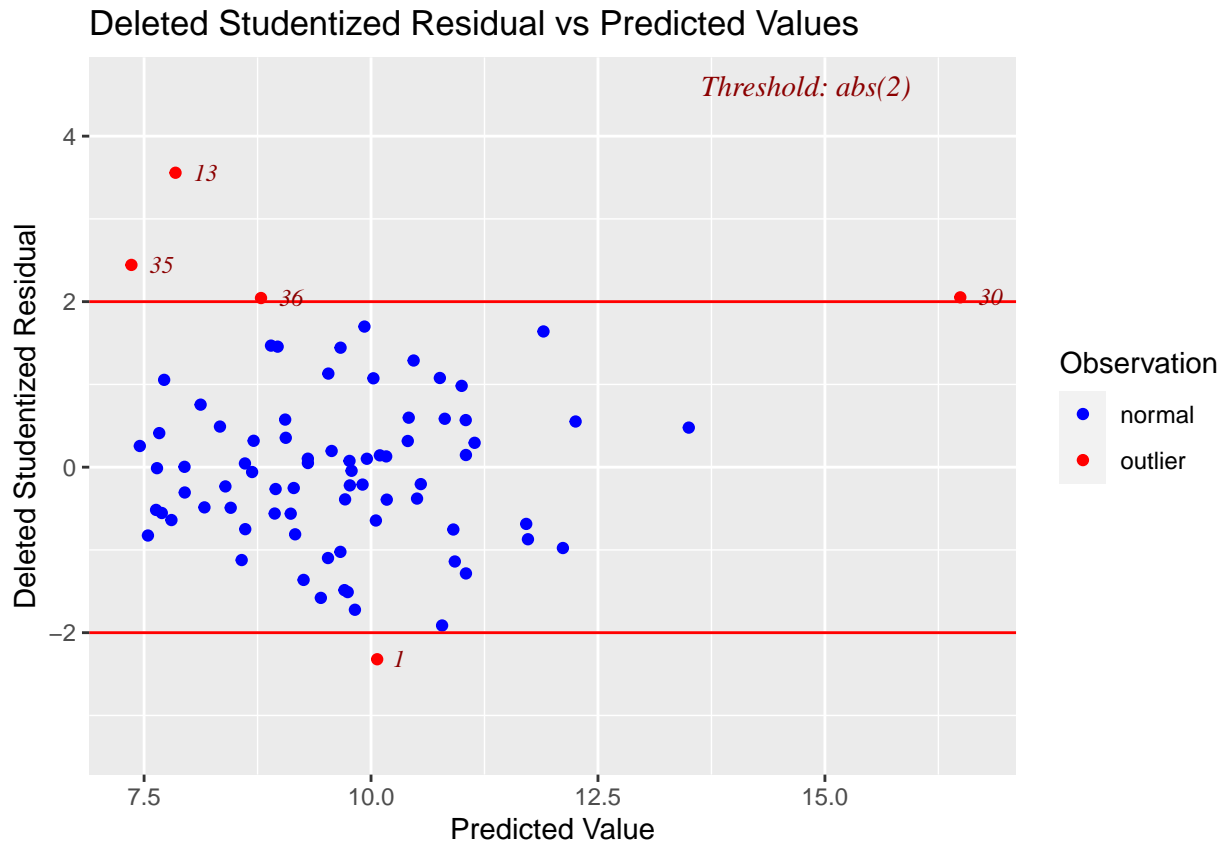
```
ols_plot_cooksd_bar(f1)
```



```
ols_plot_resid_lev(f1)
```



```
ols_plot_resid_stud_fit(f1)
```



```
library(faraway)
```

```
## Registered S3 methods overwritten by 'lme4':
##   method                      from
##   cooks.distance.influence.merMod car
##   influence.merMod             car
##   dfbeta.influence.merMod      car
##   dfbetas.influence.merMod     car
##
## Attaching package: 'faraway'
## The following object is masked from 'package:olsrr':
##
##   hsb
```

```
vif(f1)
```

```
##      X1      X2      X3      X4      X5      X6      X7_2      X7_3
## 1.293860 1.874300 2.063140 1.506713 33.389374 1.951906 1.741561 2.349673
##      X7_4      X8      X9      X10
## 1.775693 32.417458 6.983610 3.038857
```

b-) Are there any influential or outlier observations?

Yes, observation 6 is an outlier.

c-) Can X5, X6, and X7 be dropped from the model? Perform the statistical test and state your final model.

Ho: Variables can be dropped Ha: Variables cannot be dropped

Reject H_0 . Variables cannot be dropped from the model.

```
f1.r<-lm(Y~ X1+X2+X3+X4+X8+X9+X10 ,data=Q1.Dev)
anova(f1.r,f1)
```

```
## Analysis of Variance Table
```

```
##
```

```
## Model 1: Y ~ X1 + X2 + X3 + X4 + X8 + X9 + X10
```

```
## Model 2: Y ~ X1 + X2 + X3 + X4 + X5 + X6 + X7_2 + X7_3 + X7_4 + X8 + X9 +
```

```
##      X10
```

```
##   Res.Df      RSS Df Sum of Sq      F   Pr(>F)
```

```
## 1      71 106.336
```

```
## 2      66  81.057  5    25.279 4.1166 0.002582 **
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

d-) Develop an alternative model by using the Regression Tree and compare the performance against the regression model built in part a.

Regression models has a lower SSE, performs better than the tree method.

```
library(rpart)
```

```
##
```

```
## Attaching package: 'rpart'
```

```
## The following object is masked from 'package:faraway':
```

```
##
```

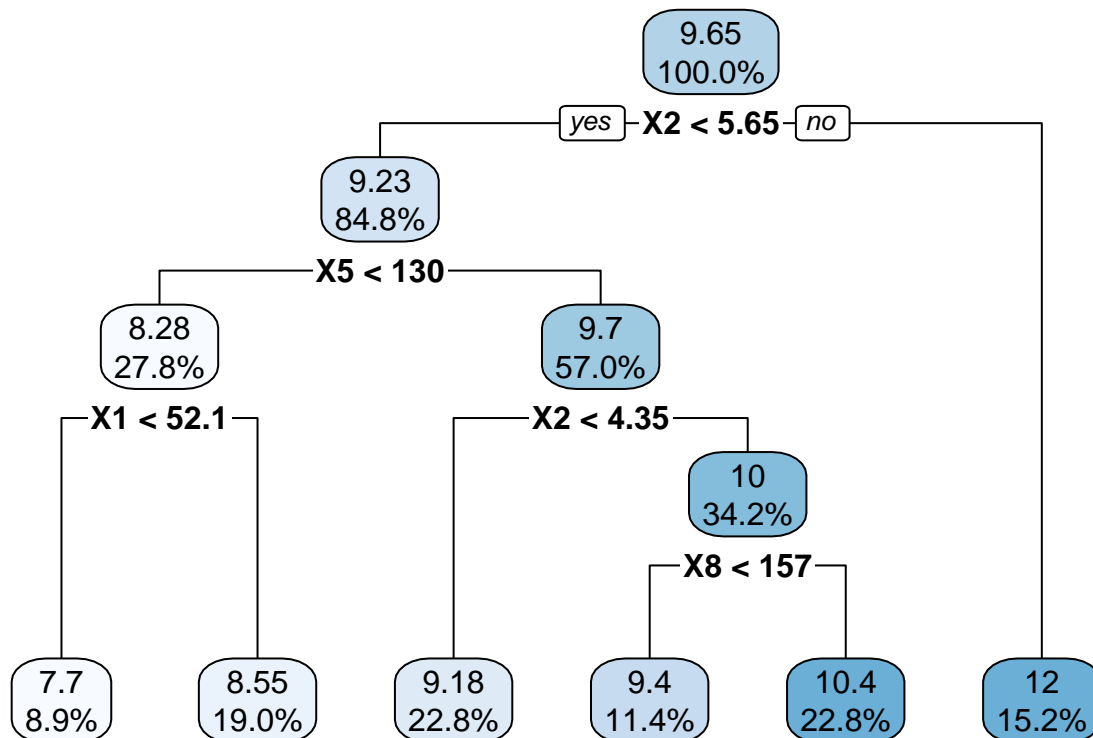
```
##      solder
```

```
q1.tr<-rpart(Y~X1+X2+X3+X4+X5+X6+X7_2+X7_3+X7_4+X8+X9+X10,data=Q1.Dev)
```

```
library(rpart.plot)
```

```
par(mfrow=c(1,1))
```

```
rpart.plot(q1.tr,digits = 3)
```



```
SSE.Tree.Dev<-sum((predict(q1.tr)-Q1.Dev$Y)^2)
SSE.Tree.Dev
```

```
## [1] 132.0195
```

```
SSE.Reg.Dev<-anova(f1)$`Sum Sq`[length(anova(f1)$`Sum Sq`)]
cbind(SSE.Tree.Dev,SSE.Reg.Dev)
```

```
##      SSE.Tree.Dev SSE.Reg.Dev
## [1,]      132.0195      81.05688
```

f-) Score the model on hold-out sample, and recalibrate the model on the holdout sample. Compare the results against the final model derived part c.

Regression model performs better than the tree model, it has also smallest SSE on the holdout sample.

```
p.tree.hold<-predict(q1.tr,Q1.Hold)
p.reg.hold<-predict(f1,Q1.Hold)
```

```
SSE.Tree.Hold<-sum((p.tree.hold-Q1.Hold$Y)^2)
SSE.Reg.Hold<-sum((p.reg.hold-Q1.Hold$Y)^2)
cbind(SSE.Tree.Hold,SSE.Reg.Hold)
```

```
##      SSE.Tree.Hold SSE.Reg.Hold
## [1,]      88.15951      79.86798
```

Problem 2

Use the question2 data set to answer this question. We are interested in predicting (Y) the number of customers who complained about the service.

a-) Build a model to predict the number of complaints, perform the statistical tests that shows that model is significant

It is a poisson regression model since the dependent variable is a count data. All variables and model are significant.

```
PF.Q2.Dat <- read.csv("/cloud/project/Practice Final Question 2.csv")
f2<-glm(Y~X1+X2+X3+X4+X5,data=PF.Q2.Dat,family=poisson)
summary(f2)
```

```
##
## Call:
## glm(formula = Y ~ X1 + X2 + X3 + X4 + X5, family = poisson, data = PF.Q2.Dat)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.93195  -0.58868  -0.00009   0.59269   2.23441
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  2.942e+00  2.072e-01  14.198  < 2e-16 ***
## X1           6.058e-04  1.421e-04   4.262 2.02e-05 ***
## X2          -1.169e-05  2.112e-06  -5.534 3.13e-08 ***
## X3          -3.726e-03  1.782e-03  -2.091  0.0365 *
## X4           1.684e-01  2.577e-02   6.534 6.39e-11 ***
## X5          -1.288e-01  1.620e-02  -7.948 1.89e-15 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 422.22  on 109  degrees of freedom
## Residual deviance: 114.99  on 104  degrees of freedom
## AIC: 571.02
##
## Number of Fisher Scoring iterations: 4
```

b-) Find the predicted number complaints given the independent variables below and predict 95% confidence interval

X1=606 X2=41393 X3=3 X4=3.04 X5=6.32

Please see below

```
test.dat<-data.frame(X1=606,X2=41393,X3=3,X4=3.04,X5=6.32)
pred<-predict(f2,test.dat,type="link",se.fit = TRUE)
exp(pred$fit)
```

```
##      1
## 12.33778
critval <- 1.96 ## approx 95% CI
upr <- exp(pred$fit + (critval * pred$se.fit))
lwr <- exp(pred$fit - (critval * pred$se.fit))
cbind(lwr,upr)

##      lwr      upr
## 1 11.08404 13.73332
```

Problem 3

Use question 3 data sets. Monthly data on amount of billings (Y) and on number of hours of staff time (X) for the 20 most recent months are recorded.

a-) Build a model to predict Y based on the independent variables and test if there is an autocorrelation persists in the data. If autocorrelation persists, remediate the autocorrelation.

There is autocorrelation in the data set. I will use Cochrane-Orcutt procedure to eliminate it. The suggested ρ is 0.33 and autocorrelation is remediated.

```
PF.Q3.Dat <- read.csv("/cloud/project/Practice Final Question 3.csv")
f3<-lm(Y~X,data=PF.Q3.Dat)
summary(f3)
```

```
##
## Call:
## lm(formula = Y ~ X, data = PF.Q3.Dat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.55515 -0.23700  0.05229  0.56250  0.80657
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  93.6865     0.8229   113.8  <2e-16 ***
## X            50.8801     0.2634   193.1  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.631 on 18 degrees of freedom
## Multiple R-squared:  0.9995, Adjusted R-squared:  0.9995
## F-statistic: 3.73e+04 on 1 and 18 DF,  p-value: < 2.2e-16
```

```
library(lmtest)
```

```
## Loading required package: zoo
##
## Attaching package: 'zoo'
##
## The following objects are masked from 'package:base':
##
##      as.Date, as.Date.numeric
```

```
dwtest(f3)
```

```
##
## Durbin-Watson test
##
## data: f3
## DW = 0.97374, p-value = 0.002891
## alternative hypothesis: true autocorrelation is greater than 0
```

```
#manual solution
```

```
library(Hmisc)
```

```
## Loading required package: lattice
```

```
##
## Attaching package: 'lattice'

## The following object is masked from 'package:faraway':
##
##      melanoma

## Loading required package: survival

##
## Attaching package: 'survival'

## The following objects are masked from 'package:faraway':
##
##      rats, solder

## Loading required package: Formula

## Loading required package: ggplot2

##
## Attaching package: 'Hmisc'

## The following objects are masked from 'package:base':
##
##      format.pval, units

et<-f3$residuals
et1<-Lag(et, shift = 1)

d1<-sum(na.omit(et1*et))
d2<-sum(na.omit(et1)^2)
rho<-d1/d2

Ytnew=PF.Q3.Dat$Y - rho*Lag(PF.Q3.Dat$Y , shift = 1)
Xtnew=PF.Q3.Dat$X - rho*Lag(PF.Q3.Dat$X , shift = 1)

f3.1<-lm(Ytnew~Xtnew)
summary(f3.1)

##
## Call:
## lm(formula = Ytnew ~ Xtnew)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.95813 -0.29553 -0.02312  0.34451  0.60490
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   63.3840     0.5592   113.4 <2e-16 ***
## Xtnew         50.5470     0.2622   192.8 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4546 on 17 degrees of freedom
## (1 observation deleted due to missingness)
## Multiple R-squared:  0.9995, Adjusted R-squared:  0.9995
## F-statistic: 3.715e+04 on 1 and 17 DF, p-value: < 2.2e-16
```

```

dwtest(f3.1)

##
## Durbin-Watson test
##
## data: f3.1
## DW = 1.7612, p-value = 0.2337
## alternative hypothesis: true autocorrelation is greater than 0

#Alternatively
#use the function
library(orcutt)
coch<- cochrane.orcutt(f3)
summary(coch)

## Call:
## lm(formula = Y ~ X, data = PF.Q3.Dat)
##
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept) 95.16377    0.91343   104.18 < 2.2e-16 ***
## X           50.46593    0.28415   177.61 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4514 on 17 degrees of freedom
## Multiple R-squared:  0.9995 , Adjusted R-squared:  0.9994
## F-statistic: 31543.9 on 1 and 17 DF, p-value: < 3.137e-29
##
## Durbin-Watson statistic
## (original): 0.97374 , p-value: 2.891e-03
## (transformed): 1.96762 , p-value: 4.079e-01

b-) X (Staff time) in month 21 is expected to be 3.625 thousand hours. Predict the amount of billings in
constant dollars for month 21, using a 99 percent prediction interval. Interpret your interval.

b0 <- summary(f3.1)[[4]][1,1]/(1-rho)
b1 <- summary(f3.1)[[4]][2,1]
MSE<-summary(f3.1)$sigma^2

X.prime<-Xtnew
X.bar.prime <- mean(X.prime[-1])

X.n.plus.1 <- 3.625
X.n <- rev(PF.Q3.Dat$X)[1]
X.n.plus.1.prime <- X.n.plus.1 - rho*X.n

# Point forecast:

Y.hat.n.plus.1 <- b0 + b1*X.n.plus.1
Y.n <- rev(PF.Q3.Dat$X)[1]
e.n <- Y.n - (b0 + b1*X.n)
Y.hat.FORECAST.n.plus.1 <- Y.hat.n.plus.1 + rho*e.n

print(paste("forecasted response at time n+1 is:", round(Y.hat.FORECAST.n.plus.1,4) ))

## [1] "forecasted response at time n+1 is: 187.1193"

```

```

# Prediction interval:

alpha <- 0.01
n<-length(PF.Q3.Dat$X)
s.pred <- sqrt(MSE*(1 + (1/n) + (X.n.plus.1.prime -X.bar.prime)^2/(sum((X.prime[-1]-X.bar.prime)^2))))
s.pred

## [1] 0.4737689

pred.L <- Y.hat.FORECAST.n.plus.1 - qt(1-alpha/2,df=n-3)*s.pred
pred.U <- Y.hat.FORECAST.n.plus.1 + qt(1-alpha/2,df=n-3)*s.pred

print(paste(100*(1-alpha) ,"percent PI for response at time n+1 is:", round(pred.L,4), ",", round(pred.U,4)))

## [1] "99 percent PI for response at time n+1 is: 185.7462 , 188.4924"

```

Problem 4

Use question 4 data set, Create development sample (70% of the data) and hold-out sample (30% of the data) use set.seed(1023) before creating the samples.

a-) Use the development sample , fit a linear regression model, regression tree and Neural Network Model, and calculate the SSE for each model, which method has the lowest SSE?

Reg Model: X_1, X_2, X_3 and X_4 are significant and R^2 square is 56%. QQ plot indicates S shape suggesting that transformation is needed. Residual vs Fitted graph suggest that further testing for unequal variances are needed. However, there is no multicollinearity in the data. SSE is 10128646733.

Tree: It is based only 2 variables (X_1 and X_2). SSE is 7450047035, significantly lower than Regression model.

Neural Network:

I used 2 hidden layers with 5 nodes each, you can also try single layer or multiple layer SSE is 22389967988. It has the highest SSE.

```

PF.Q4.Dat<- read.csv("/cloud/project/Practice Final Question 4.csv")
n<-dim(PF.Q4.Dat)[1]
IND=sample(c(1:n),n*0.7)
set.seed(1023)
Q4.Dev<-PF.Q4.Dat[IND,]
Q4.Hold<-PF.Q4.Dat[-IND,]

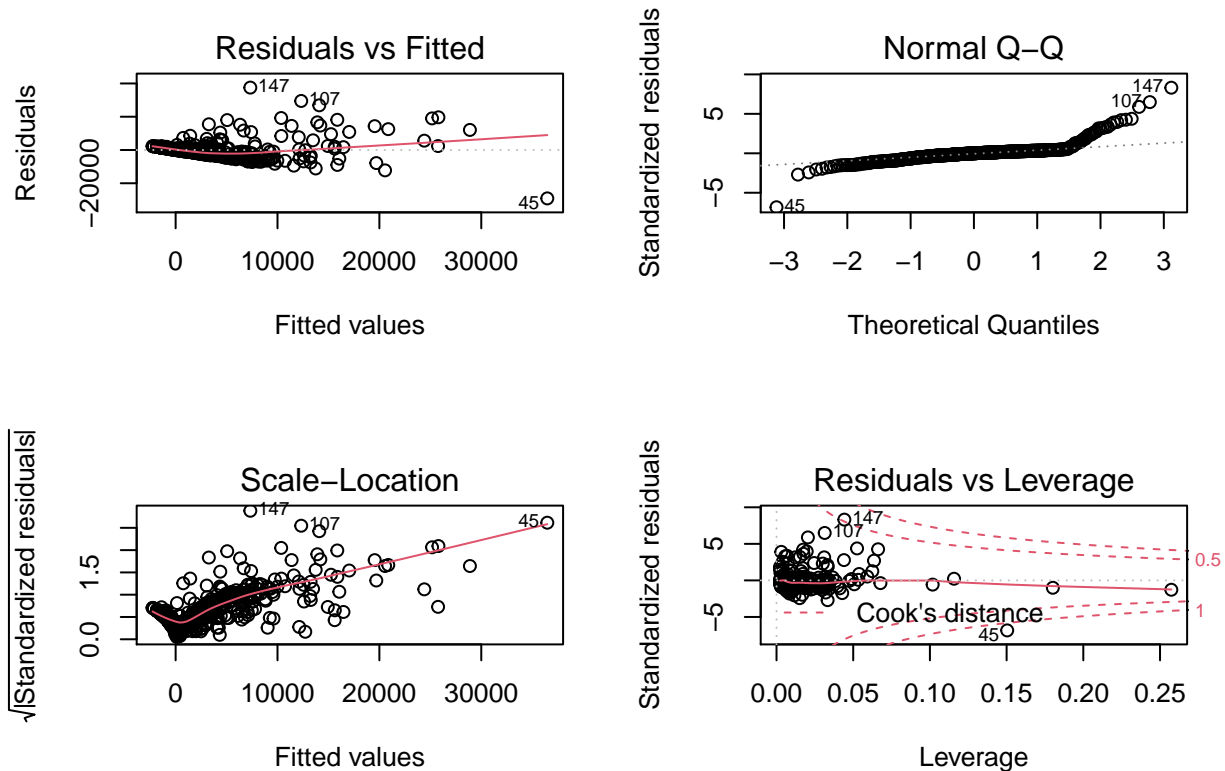
#regression model
f4<-lm(Y~X1+X2+X3+X4+X5,data=Q4.Dev)
summary(f4)

##
## Call:
## lm(formula = Y ~ X1 + X2 + X3 + X4 + X5, data = Q4.Dev)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -29093  -1670      38    1110   37508
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2153.10      352.43  -6.109 1.91e-09 ***

```

```
## X1          811.24      37.93  21.388 < 2e-16 ***
## X2        -438.64     204.48  -2.145 0.032385 *
## X3         319.31      89.23   3.579 0.000376 ***
## X4         882.90     783.08   1.127 0.260040
## X5          64.89      34.43   1.885 0.059981 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4616 on 545 degrees of freedom
## Multiple R-squared:  0.5326, Adjusted R-squared:  0.5283
## F-statistic: 124.2 on 5 and 545 DF,  p-value: < 2.2e-16
```

```
par(mfrow=c(2,2))
plot(f4)
```

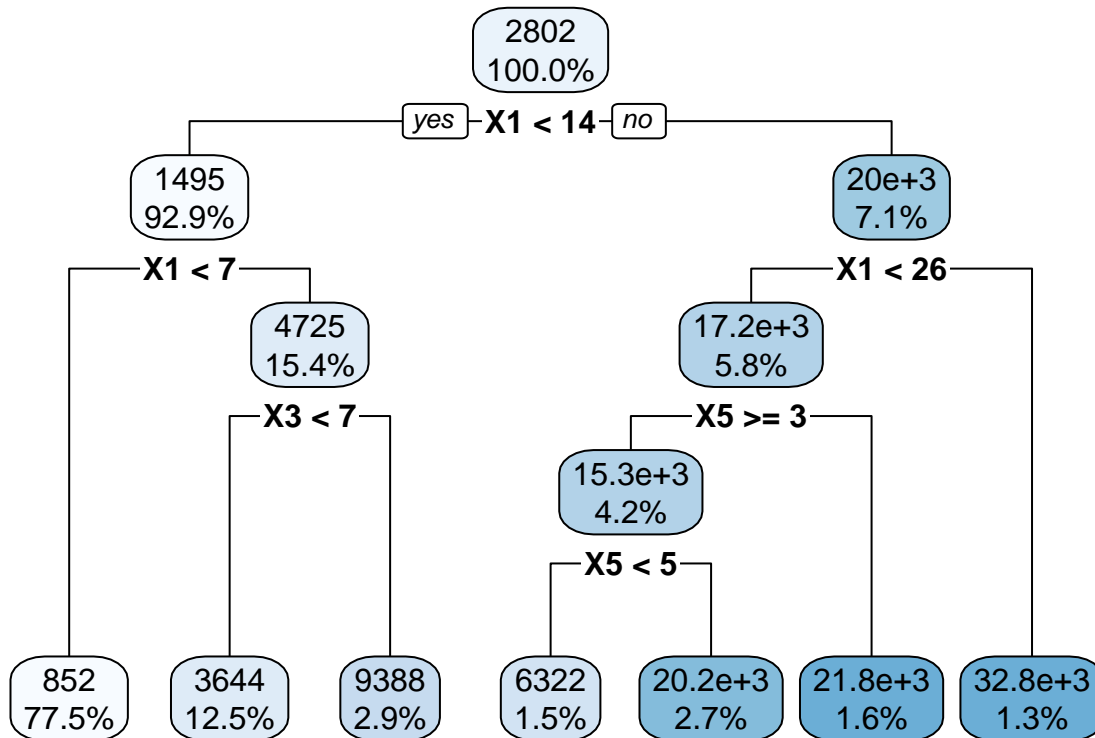


```
vif(f4)
```

```
##          X1          X2          X3          X4          X5
## 1.179818 1.479310 1.528698 1.065644 1.035085
SSE.Reg.Dev<-anova(f4)$`Sum Sq`[length(anova(f4)$`Sum Sq`)]
SSE.Reg.Dev
```

```
## [1] 11611220820
```

```
#Tree
library(rpart)
q4.tr<-rpart(Y~X1+X2+X3+X4+X5,data=Q4.Dev)
library(rpart.plot)
par(mfrow=c(1,1))
rpart.plot(q4.tr,digits = 3)
```



```
SSE.Tree.Dev<-sum((predict(q4.tr)-Q4.Dev$Y)^2)
SSE.Tree.Dev
```

```
## [1] 8323992319
```

```

#Neural Network
#install.packages("neuralnet")
library(neuralnet)
normalize <- function(x) {return((x - min(x)) / (max(x) - min(x)))}
scaled.Q4.Dat <- as.data.frame(lapply(PF.Q4.Dat, normalize))
scaled.Q4.Dev<- scaled.Q4.Dat[IND,]
scaled.Q4.Hold<- scaled.Q4.Dat[-IND,]

NN = neuralnet(Y~X1+X2+X3+X4+X5,hidden=c(5,5),scaled.Q4.Dev,linear.output= T )
plot(NN)
predict_testNN= compute(NN, scaled.Q4.Dev[, -c(1)])
#we need to transform it back to original scale
predict_testNN1 = (predict_testNN$net.result* (max(PF.Q4.Dat$Y) -min(PF.Q4.Dat$Y))) + min(PF.Q4.Dat$Y)
plot(scaled.Q4.Dev$Y, predict_testNN1, col='blue', pch=16, ylab= "Predicted Y", xlab= "Actual Y")

SSE.NN.Dev<-sum((predict_testNN1-scaled.Q4.Dev$Y)^2)
SSE.NN.Dev

```

```
## [1] 24981840718
```

b-) test the models performances on the hold out sample, which model would you choose?

Tree has the lowest SSE, I would choose the Tree approach. It outperforms other models both in and out samples.

```
SSE <- function(actual, predicted) {sum((actual - predicted)^2)}
```

```

#Regression
reg.predict<-predict(f4,Q4.Hold)
#Tree
tree.predict<-predict(q4.tr,Q4.Hold)
#NN
nn.predict<-compute(NN, scaled.Q4.Hold)
nn.predict1 = (nn.predict$net.result*(max(PF.Q4.Dat$Y) -min(PF.Q4.Dat$Y))) + min(PF.Q4.Dat$Y)
#SSEs

cbind(REG=SSE(Q4.Hold$Y,reg.predict),
Tree=SSE(Q4.Hold$Y,tree.predict),
NN=SSE(Q4.Hold$Y,nn.predict1))

##           REG           Tree           NN
## [1,] 4408302064 5467759713 5725076042

```

Problem 5

Use Question 5 dataset, Y is a dichotomous response variable and X2, X3, and X4 are categorical variables.

a-) Fit a regression model containing the predictor variables in first-order terms and interaction terms for all pairs of predictor variables on development sample.

Y is a dichotomous response variable. Therefore, we will use logistic regression model. X_2 has 3 levels and other two categorical variables has two levels. X_3 is coded 1-2. X_4 is coded 0-1. We need to create two dummy variables for X_2 and change the levels of X_3 to 0-1.

Only X_2 is significant. Please see below

```

PF.Q5.Dat<- read.csv("/cloud/project/Practice Final Question 5.csv")
table(PF.Q5.Dat$X2)

```

```

##
##  1  2  3
## 77 49 70

```

```

table(PF.Q5.Dat$X3)

```

```

##
##   1   2
## 117  79

```

```

table(PF.Q5.Dat$X4)

```

```

##
##   0   1
## 139  57

```

```

library('fastDummies')
Q5.Dat<-dummy_cols(PF.Q5.Dat, select_columns = 'X2')
Q5.Dat<-dummy_cols(Q5.Dat, select_columns = 'X3')
head(Q5.Dat)

```

```

##   X1 X2 X3 X4 Y X2_1 X2_2 X2_3 X3_1 X3_2
## 1 33  1  1  0 1    1    0    0    1    0
## 2 35  1  1  0 1    1    0    0    1    0
## 3  6  1  1  0 0    1    0    0    1    0
## 4 60  1  1  0 1    1    0    0    1    0

```



```
## 5 18 3 1 1 0 0 0 1 1 0
## 6 26 3 1 0 0 0 0 1 1 0
```

```
#drop X2, X3,X2_1, and X3_1
Q5.Dat1<-Q5.Dat[,-c(2,3,6,9)]
head(Q5.Dat1)
```

```
## X1 X4 Y X2_2 X2_3 X3_2
## 1 33 0 1 0 0 0
## 2 35 0 1 0 0 0
## 3 6 0 0 0 0 0
## 4 60 0 1 0 0 0
## 5 18 1 0 0 1 0
## 6 26 0 0 0 1 0
```

```
#i used the short cut, rather typing
f5<-glm(Y~X1+X4+X2_2+X2_3+X3_2+X1:X4+X1:X2_2+X1:X2_3+X1:X3_2+X4:X2_2+X4:X2_3+X4:X3_2+X2_2:X3_2+X2_3:X3_2)
```

b-)Use the likelihood ratio test to determine whether all interaction terms can be dropped from the regression model; State the alternatives, full and reduced models, decision rule, and conclusion.

Ho: Variables can be dropped Ha: Variables cannot be dropped

Accept Ho, all interaction terms can be dropped.

```
f5r<-glm(Y~X1+X4+X2_2+X2_3+X3_2,data=Q5.Dat1,family=binomial)
anova(f5r,f5, test="Chi")
```

```
## Analysis of Deviance Table
##
## Model 1: Y ~ X1 + X4 + X2_2 + X2_3 + X3_2
## Model 2: Y ~ X1 + X4 + X2_2 + X2_3 + X3_2 + X1:X4 + X1:X2_2 + X1:X2_3 +
## X1:X3_2 + X4:X2_2 + X4:X2_3 + X4:X3_2 + X2_2:X3_2 + X2_3:X3_2
## Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1 190 215.36
## 2 181 212.84 9 2.5213 0.9803
```

c-)For logistic regression model in part (a), use backward elimination to decide which predictor variables can be dropped from the regression model. Which variables are retained in the regression model?

Please see below, all interaction terms, X_4 are dropped from the model. All variables are significant.

```
t0<-step(f5,direction="backward",trace=0)
summary(t0)
```

```
##
## Call:
## glm(formula = Y ~ X1 + X2_2 + X2_3 + X3_2, family = binomial,
## data = Q5.Dat1)
##
## Deviance Residuals:
## Min 1Q Median 3Q Max
## -2.2898 -0.8648 0.3887 0.8149 1.9887
##
## Coefficients:
## Estimate Std. Error z value Pr(>|z|)
## (Intercept) 0.054054 0.381415 0.142 0.887302
## X1 0.035471 0.009796 3.621 0.000294 ***
## X2_2 -1.174332 0.417764 -2.811 0.004939 **
```

```
## X2_3      -1.953575    0.402550   -4.853 1.22e-06 ***
## X3_2      0.789524    0.348572    2.265 0.023511 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 270.06  on 195  degrees of freedom
## Residual deviance: 215.36  on 191  degrees of freedom
## AIC: 225.36
##
## Number of Fisher Scoring iterations: 4
```

d-) Conduct the Hosmer-Lemeshow goodness of fit test for the appropriateness of the logistic regression function by forming five groups. State the alternatives, decision rule, and conclusion.

Using the model in part C, Ho: Fit is good Ha: Fit is not good

Accept null, the fit is good.

```
library(ResourceSelection)
```

```
## ResourceSelection 0.3-5    2019-07-22
```

```
hoslem.test(t0$y,fitted(t0),g=5)
```

```
##
## Hosmer and Lemeshow goodness of fit (GOF) test
##
## data:  t0$y, fitted(t0)
## X-squared = 3.8928, df = 3, p-value = 0.2733
```

e-) Make the prediction for the following two cases and calculate 95% confidence interval

X1=(33,6) X2=(1,1) X3=(1,1) X4=(0,0)

Please see below

```
test.dat<-data.frame(X1=c(33,6),X2=c(1,1),X3=c(1,1),X4=c(0,0))
```

#however we need to create dummy variables

```
test.dat<-data.frame(X1=c(33,6),X2_2=c(0,0),X2_3=c(0,0),X3_2=c(0,0),X4=c(0,0))
```

#to install inv.logit function, we need to boot library

```
library(boot)
```

```
##
## Attaching package: 'boot'
##
## The following object is masked from 'package:survival':
##
##      aml
##
## The following object is masked from 'package:lattice':
##
##      melanoma
##
## The following objects are masked from 'package:faraway':
##
##      logit, melanoma
```

```
pred<-predict(t0,test.dat,type="link",se.fit = TRUE)
```

```
inv.logit(pred$fit)
```

```
##           1           2  
## 0.7728740 0.5663273
```

```
critval <- 1.96 ## approx 95% CI
```

```
upr <- inv.logit(pred$fit + (critval * pred$se.fit))
```

```
lwr <- inv.logit(pred$fit - (critval * pred$se.fit))
```

```
cbind(lwr,upr)
```

```
##           lwr           upr  
## 1 0.6383973 0.8677038  
## 2 0.3958197 0.7224562
```