

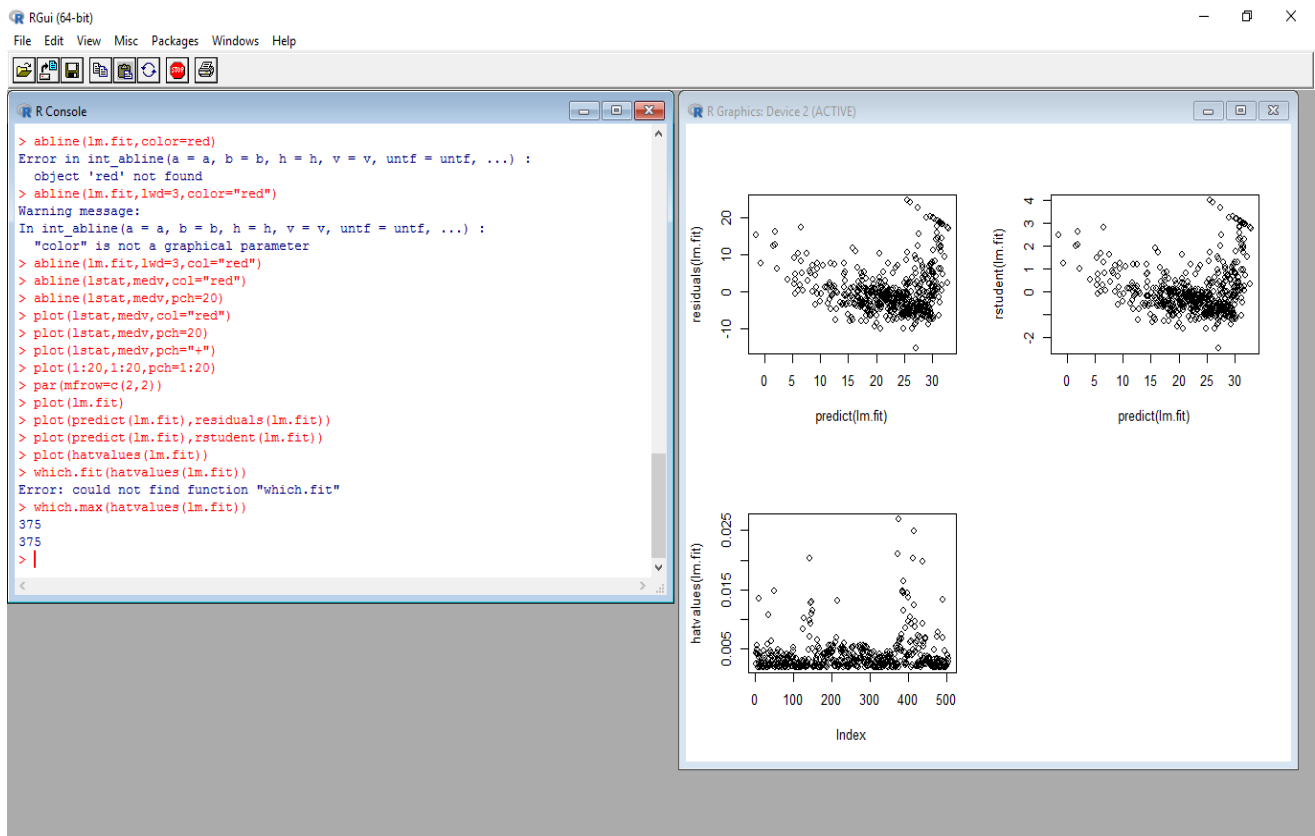
HOMEWORK-2(Big Data)

Lalitha Madhuri Putchala

My Central ID: 700657631

3.6.2 ~ 3.6.6 (5 Captures)

3.6.2 (Capture 1):



3.6.3 (Capture 2):

```
RGU (64-bit) - [R Console]
File Edit View Misc Packages Windows Help

> vif(lm.fit)
Error: could not find function "vif"
> library(car)
> vif(lm.fit)
      crim      zn      indus      chas      nox      rm      age      dis      rad      tax      ptratio      black      lstat
1.792192 2.298758 3.991596 1.073995 4.393720 1.933744 3.100826 3.955945 7.484496 9.008554 1.799084 1.348521 2.941491
> lm.fit1=lm(medv~.,data=Boston)
> summary(lm.fit1)

Call:
lm(formula = medv ~ ., data = Boston)

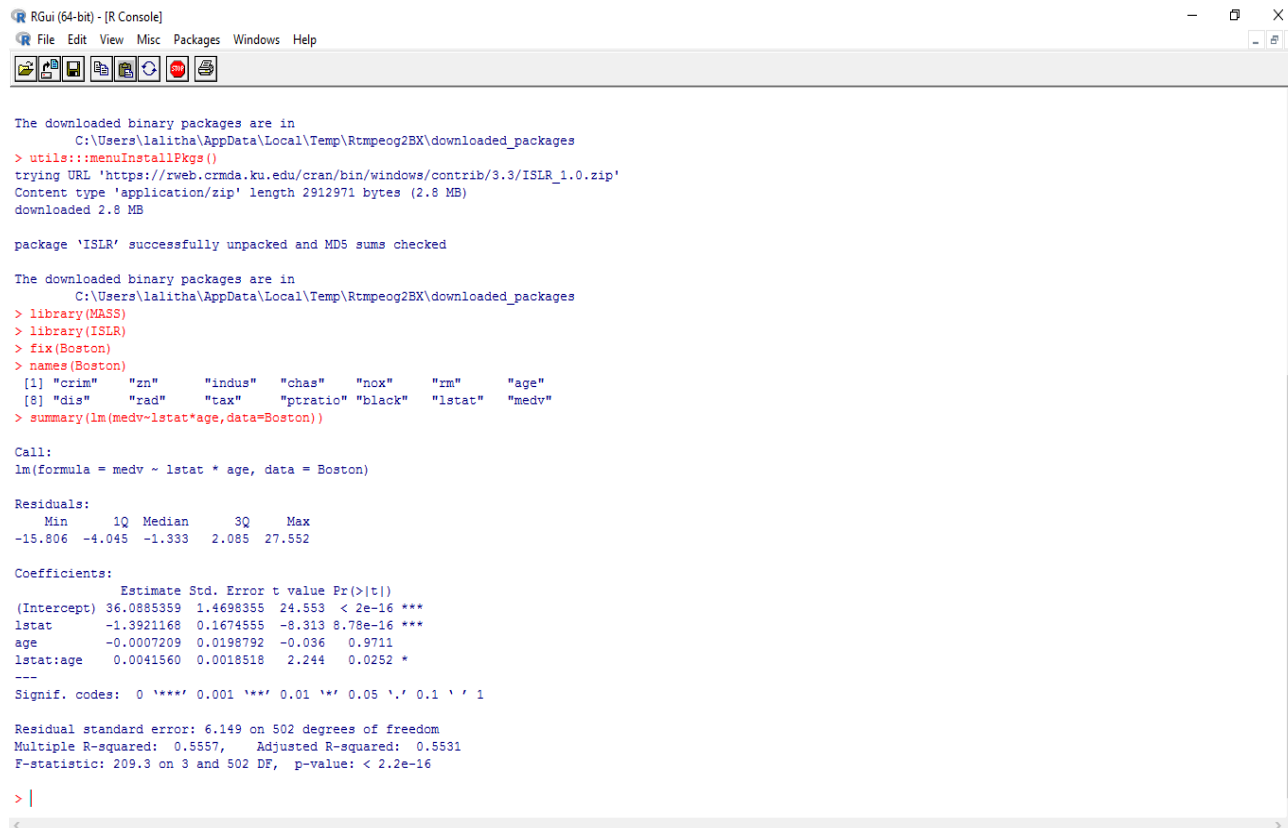
Residuals:
    Min       1Q   Median       3Q      Max
-15.595  -2.730  -0.518   1.777   26.199

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  3.646e+01  5.103e+00   7.144 3.28e-12 ***
crim        -1.080e-01  3.286e-02  -3.287 0.001087 **
zn          4.642e-02  1.373e-02   3.382 0.000778 ***
indus       2.056e-02  6.150e-02   0.334 0.738288
chas       2.687e+00  8.616e-01   3.118 0.001925 **
nox        -1.777e+01  3.820e+00  -4.651 4.25e-06 ***
rm          3.810e+00  4.179e-01   9.116 < 2e-16 ***
age         6.922e-04  1.321e-02   0.052 0.958229
dis        -1.476e+00  1.995e-01  -7.398 6.01e-13 ***
rad         3.060e-01  6.635e-02   4.613 5.07e-06 ***
tax        -1.233e-02  3.760e-03  -3.280 0.001112 **
ptratio     -9.527e-01  1.308e-01  -7.283 1.31e-12 ***
black       9.312e-03  2.686e-03   3.467 0.000573 ***
lstat      -5.248e-01  5.072e-02 -10.347 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.745 on 492 degrees of freedom
Multiple R-squared:  0.7406,    Adjusted R-squared:  0.7338
F-statistic: 108.1 on 13 and 492 DF,  p-value: < 2.2e-16

> lm.fit1=update(lm.fit,~.-age)
> |
```

3.6.4 (Capture 3):



```
RGui (64-bit) - [R Console]
File Edit View Misc Packages Windows Help

The downloaded binary packages are in
  C:\Users\lalitha\AppData\Local\Temp\Rtmpeog2BX\downloaded_packages
> utils::menuInstallPkgs()
trying URL 'https://rweb.cmcda.ku.edu/cran/bin/windows/contrib/3.3/ISLR_1.0.zip'
Content type 'application/zip' length 2912971 bytes (2.8 MB)
downloaded 2.8 MB

package 'ISLR' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
  C:\Users\lalitha\AppData\Local\Temp\Rtmpeog2BX\downloaded_packages
> library(MASS)
> library(ISLR)
> fix(Boston)
> names(Boston)
 [1] "crim"    "zn"      "indus"   "chas"    "nox"     "rm"      "age"
 [8] "dis"     "rad"     "tax"     "ptratio" "black"   "lstat"   "medv"
> summary(lm(medv~lstat*age,data=Boston))

Call:
lm(formula = medv ~ lstat * age, data = Boston)

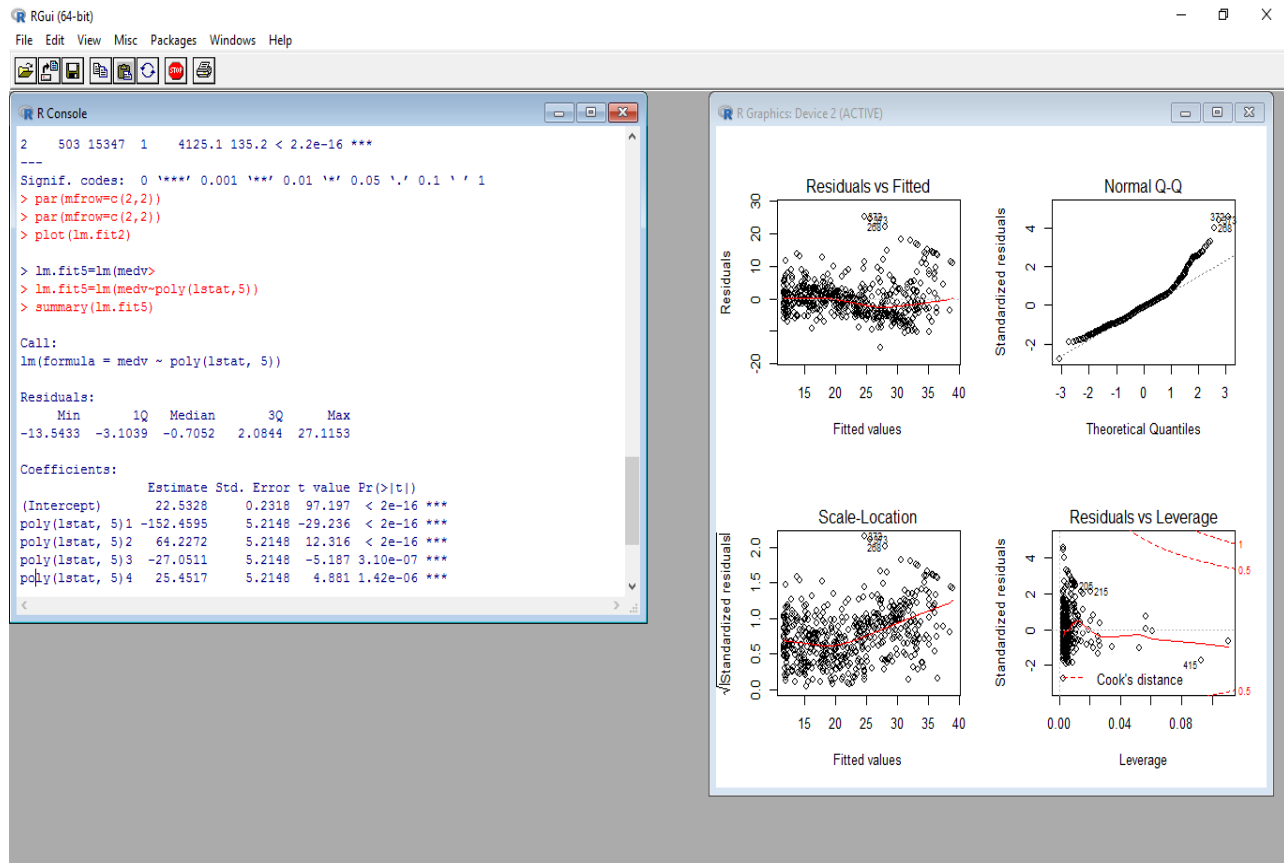
Residuals:
    Min       1Q   Median       3Q      Max
-15.806  -4.045  -1.333   2.085  27.552

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 36.0885359  1.4698355   24.553 < 2e-16 ***
lstat      -1.3921168  0.1674555   -8.313 8.78e-16 ***
age        -0.0007209  0.0198792   -0.036  0.9711
lstat:age    0.0041560  0.0018518    2.244  0.0252 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.149 on 502 degrees of freedom
Multiple R-squared:  0.5557,    Adjusted R-squared:  0.5531
F-statistic: 209.3 on 3 and 502 DF,  p-value: < 2.2e-16

> |
```

3.6.5 (Capture 4):



3.6.6 (Capture 5):

```
RGU (64-bit) - [R Console]
File Edit View Misc Packages Windows Help

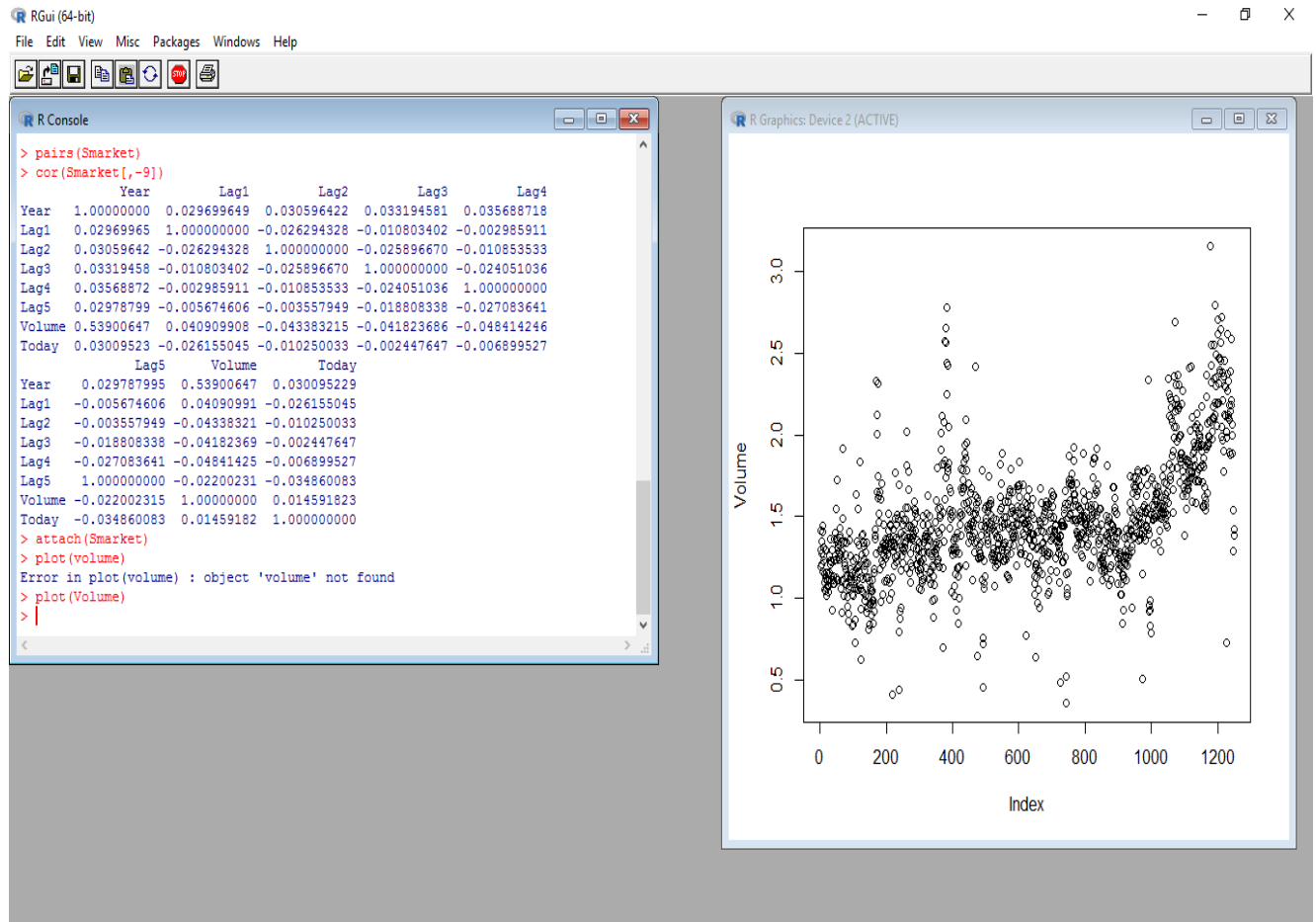
Coefficients:
      Estimate Std. Error t value Pr(>|t|)
(Intercept)    6.5755654   1.0087470    6.519 2.22e-10 ***
CompPrice      0.0929371   0.0041183   22.567 < 2e-16 ***
Income         0.0108940   0.0026044    4.183 3.57e-05 ***
Advertising    0.0702462   0.0226091    3.107 0.002030 **
Population     0.0001592   0.0003679    0.433 0.665330
Price         -0.1008064   0.0074399   -13.549 < 2e-16 ***
ShelveLocGood  4.8486762   0.1528378   31.724 < 2e-16 ***
ShelveLocMedium 1.9532620   0.1257682   15.531 < 2e-16 ***
Age           -0.0579466   0.0159506   -3.633 0.000318 ***
Education     -0.0208525   0.0196131   -1.063 0.288361
UrbanYes      0.1401597   0.1124019    1.247 0.213171
USYes        -0.1575571   0.1489234   -1.058 0.290729
Income:Advertising 0.0007510 0.0002784    2.698 0.007290 **
Price:Age      0.0001068 0.0001333    0.801 0.423812
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.011 on 386 degrees of freedom
Multiple R-squared:  0.8761,    Adjusted R-squared:  0.8719
F-statistic: 210 on 13 and 386 DF,  p-value: < 2.2e-16

> attach(Carseats)
> contrasts(ShelveLoc)
      Good Medium
Bad      0      0
Good     1      0
Medium   0      1
> LoadLibraries
Error: object 'LoadLibraries' not found
> LoadLibraries()
Error: could not find function "LoadLibraries"
> LoadLibraries=function(){
+ library(ISLR)
+ library(MASS)
+ print("The libraries have been loaded.")
+ }
> LoadLibraries()
[1] "The libraries have been loaded."
> LoadLibraries()
```

4.6.1 ~ 4.6.4 (4 Captures)

4.6.1 (Capture 1):

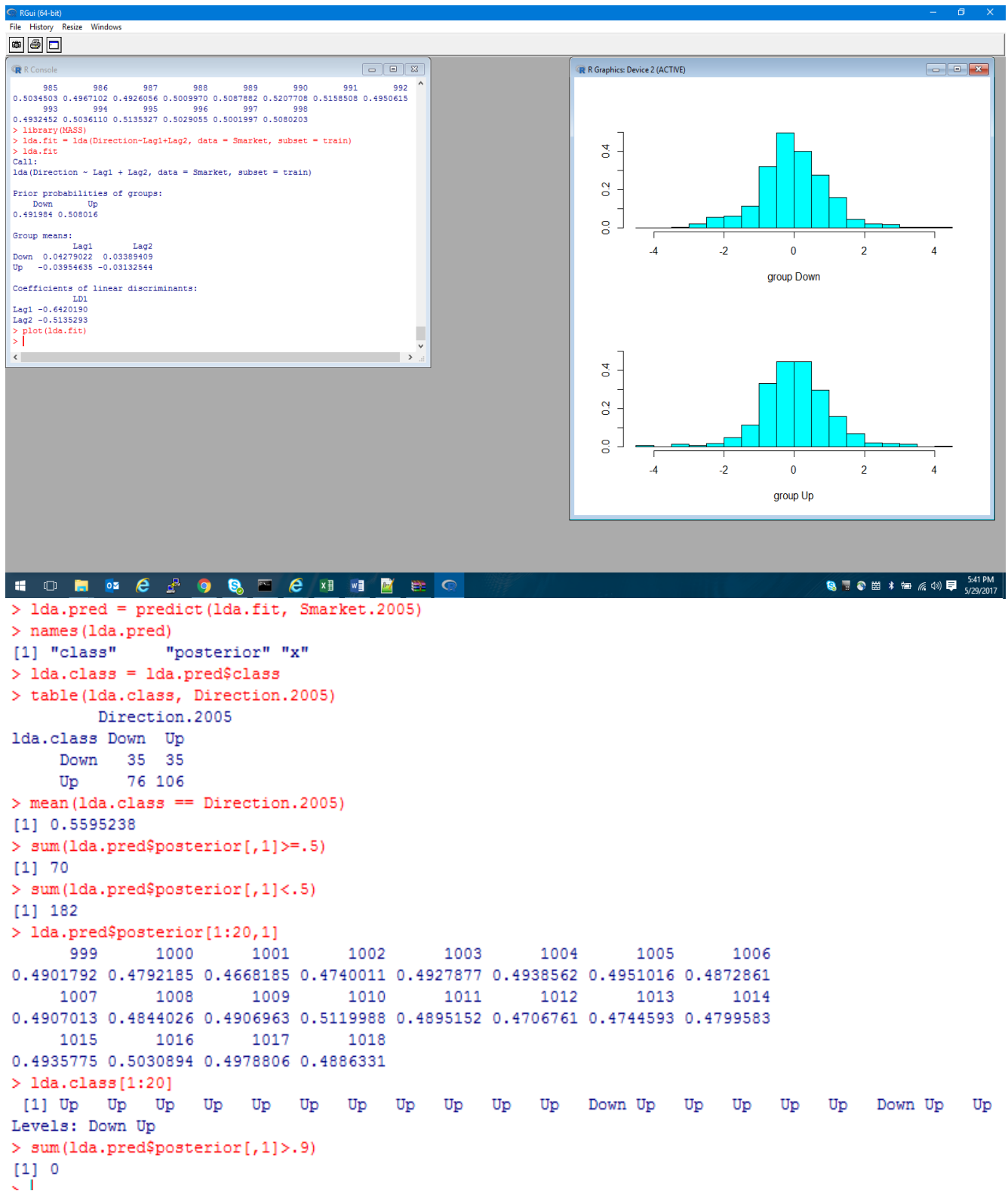


4.6.2 (Capture 2):

```
RGU (64-bit) - [R Console]
File Edit View Misc Packages Windows Help

glm.pred Down Up
Down 145 141
Up 457 507
> (507+145)/1250
[1] 0.5216
> mean(glm.pred==Direction)
[1] 0.5216
> train=(Year<2005)
> Smarket.2005=Smarket[!train,]
> dim(Smarket.2005)
[1] 252 9
> Direction.2005=Direction[!train]
> glm.fit=glm(Direction~Lag1+Lag2+Lag3+Lag4+Lag5+Volume,family=binomial,data=Smarket,subset=train)
> glm.probs=predict(glm.fit,Smarket.2005,type="response")
> glm.pred=rep("Down",252)
> glm.pred[glm.probs>.5]="Up"
> table(glm.pred,Direction.2005)
      Direction.2005
glm.pred Down Up
Down   77  97
Up     34   44
> mean(glm.pred==Direction.2005)
[1] 0.4801587
> mean(glm.pred!=Direction.2005)
Error: unexpected '=' in "mean(glm.pred!="
> mean(glm.pred!=Direction.2005)
[1] 0.5198413
> glm.fit=glm(Direction~Lag1+Lag2+Lag3+Lag4+Lag5+Volume,family=binomial,data=Smarket,subset=train)
> glm.probs=predict(glm.fit,Smarket.2005,type="response")
> glm.pred=rep("Down",252)
> glm.pred[glm.probs>.5]="Up"
> table(glm.pred,Direction.2005)
      Direction.2005
glm.pred Down Up
Down   77  97
Up     34   44
> mean(glm.pred==Direction.2005)
[1] 0.4801587
> 106/(106+76)
[1] 0.5824176
> predict(glm.fit,newdata=data.frame(Lag1=c(1.2,1.5),Lag2=c(1.1,-0.8)),type="response")
```

4.6.3 (Capture 3):



4.6.4 (Capture 4):

```
> qda.fit = qda(Direction~Lag1+Lag2, data = Smarket, subset = train)
> qda.fit
Call:
qda(Direction ~ Lag1 + Lag2, data = Smarket, subset = train)

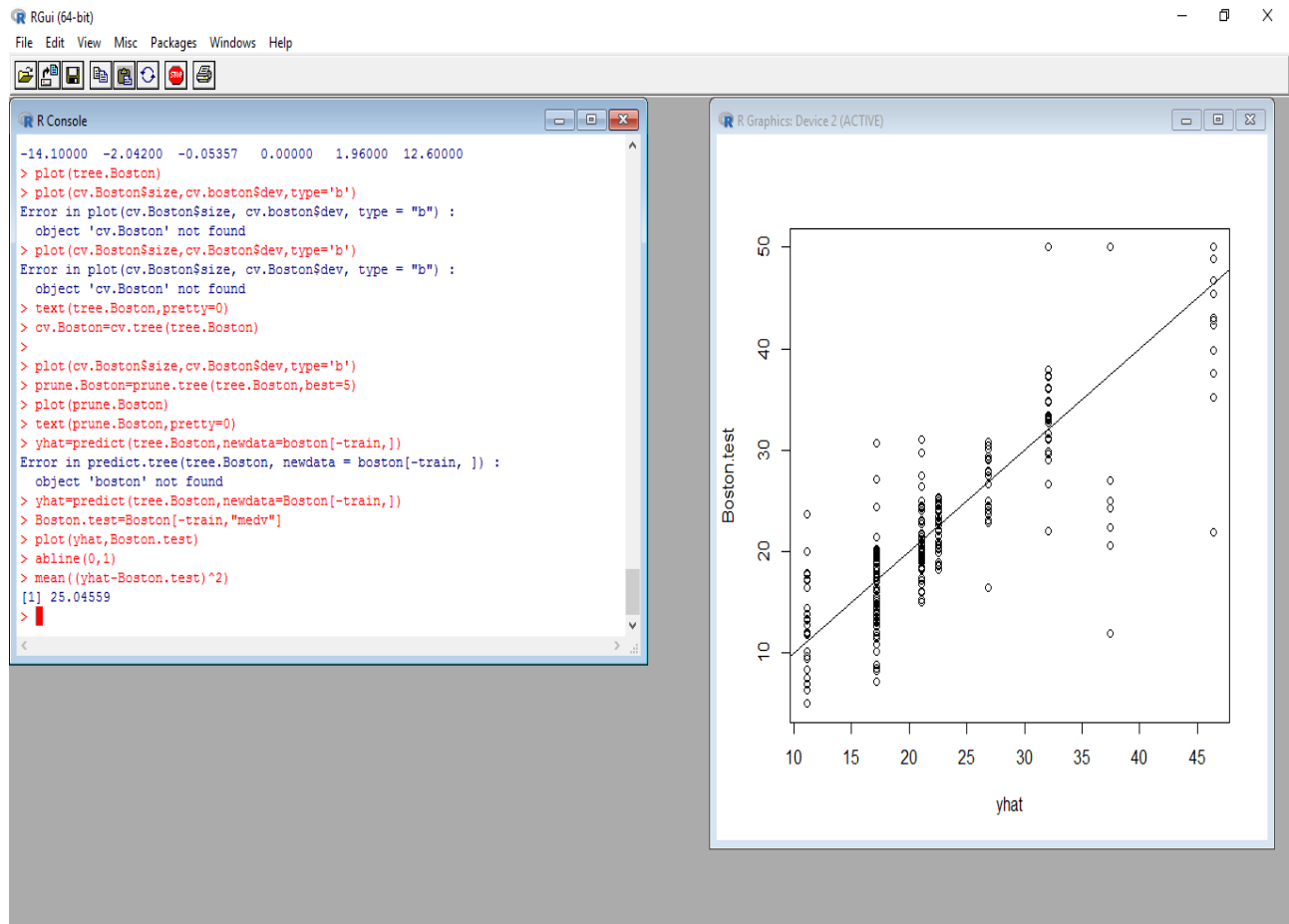
Prior probabilities of groups:
      Down      Up 
0.491984 0.508016 

Group means:
      Lag1      Lag2 
Down 0.04279022 0.03389409 
Up   -0.03954635 -0.03132544 
> qda.class = predict(qda.fit, Smarket.2005)$class
> table(qda.class, Direction.2005)
      Direction.2005
qda.class Down  Up 
      Down   30  20 
      Up    81 121 
> mean(qda.class == Direction.2005)
[1] 0.5992063
|
```

8.3.1 (Capture 1):

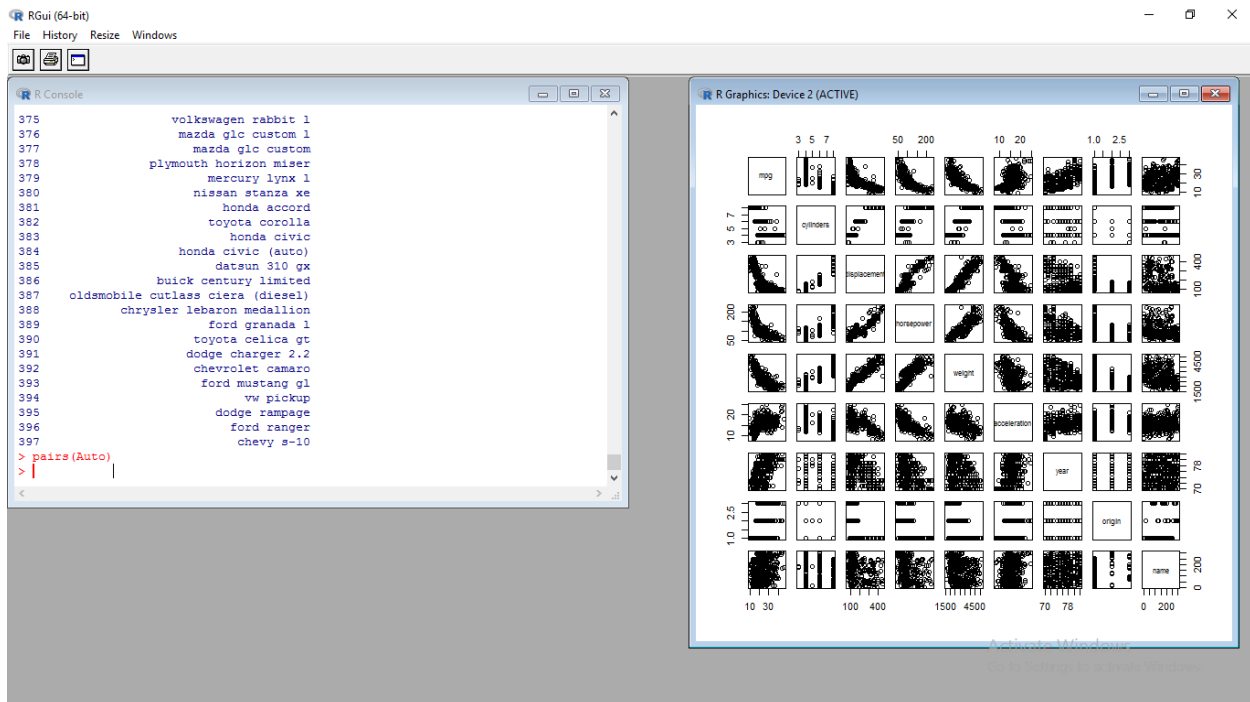


8.3.2 (Capture 2):

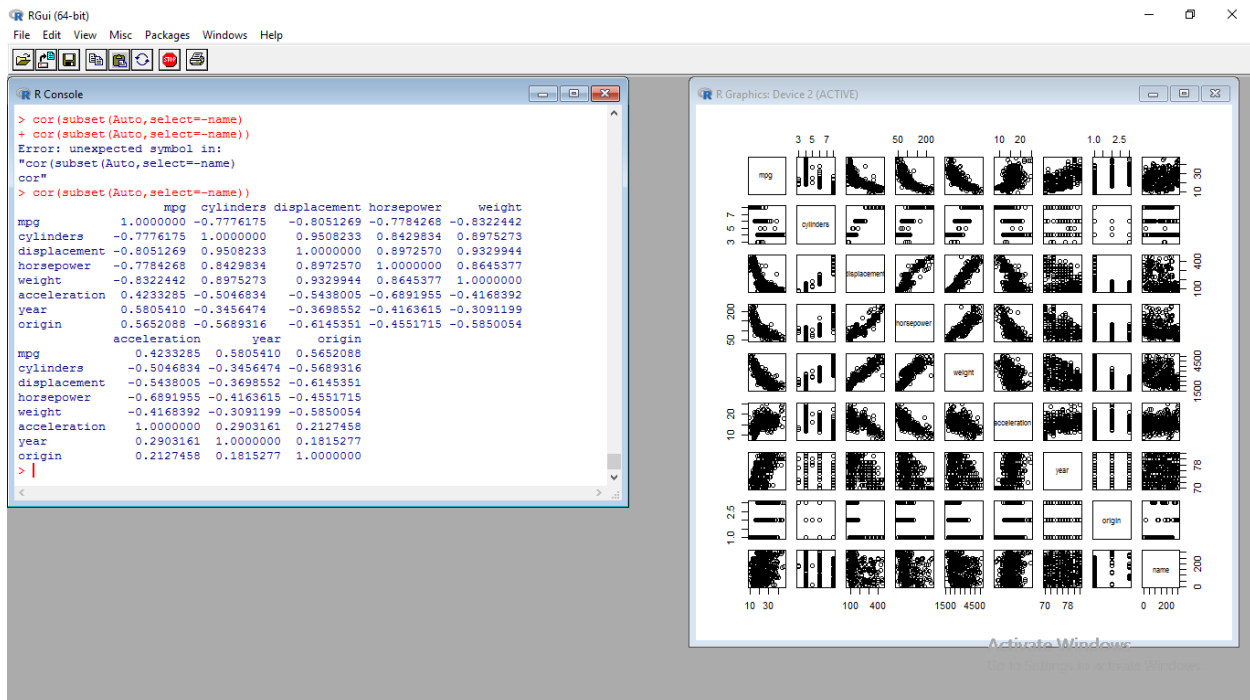


Chapter-3 (Exercise 3.7)

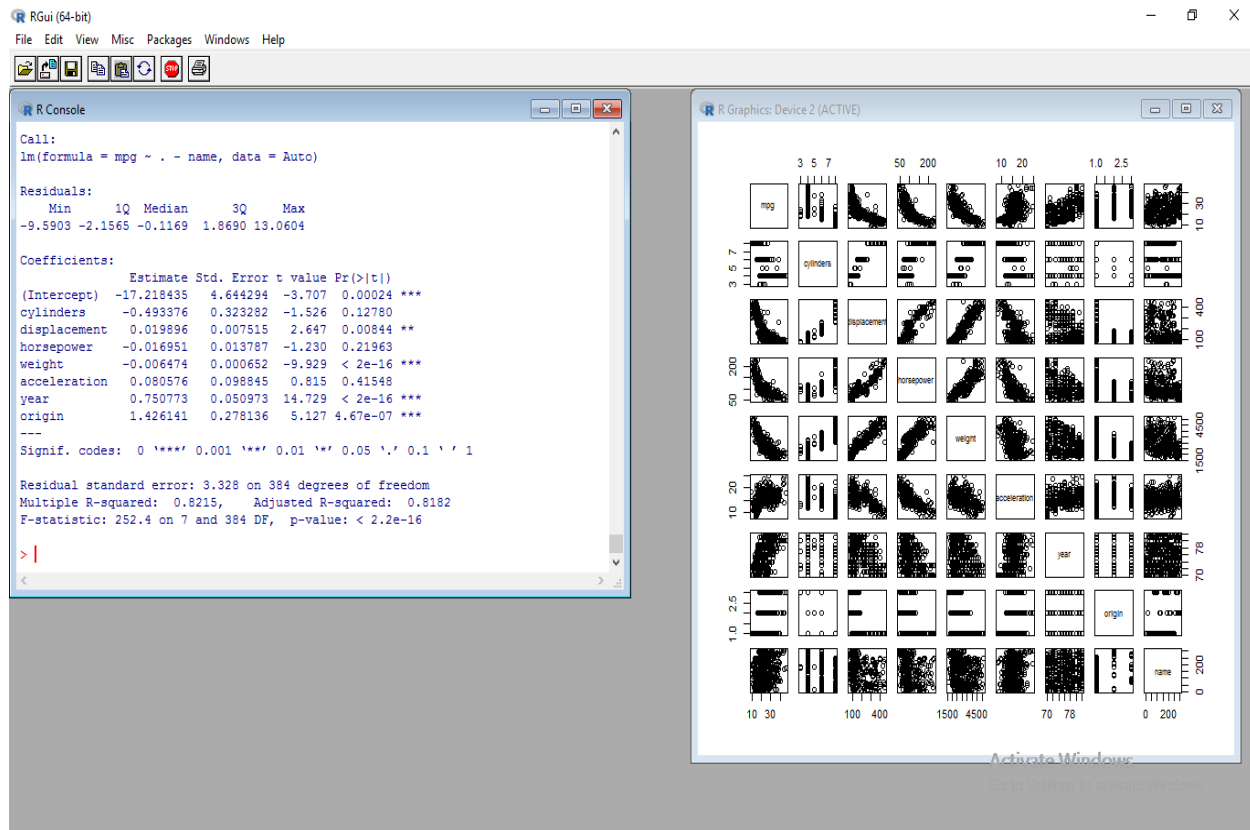
Q9. (a)



Q 9. (b)

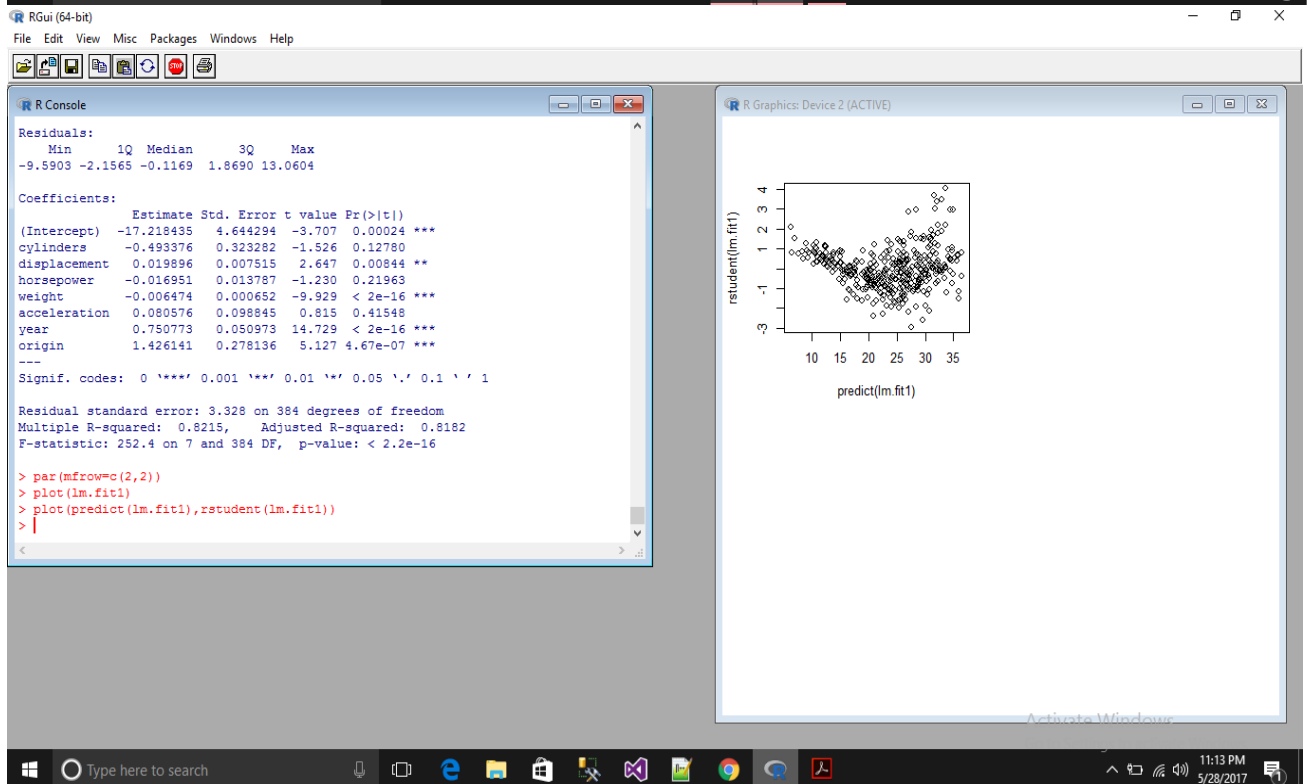
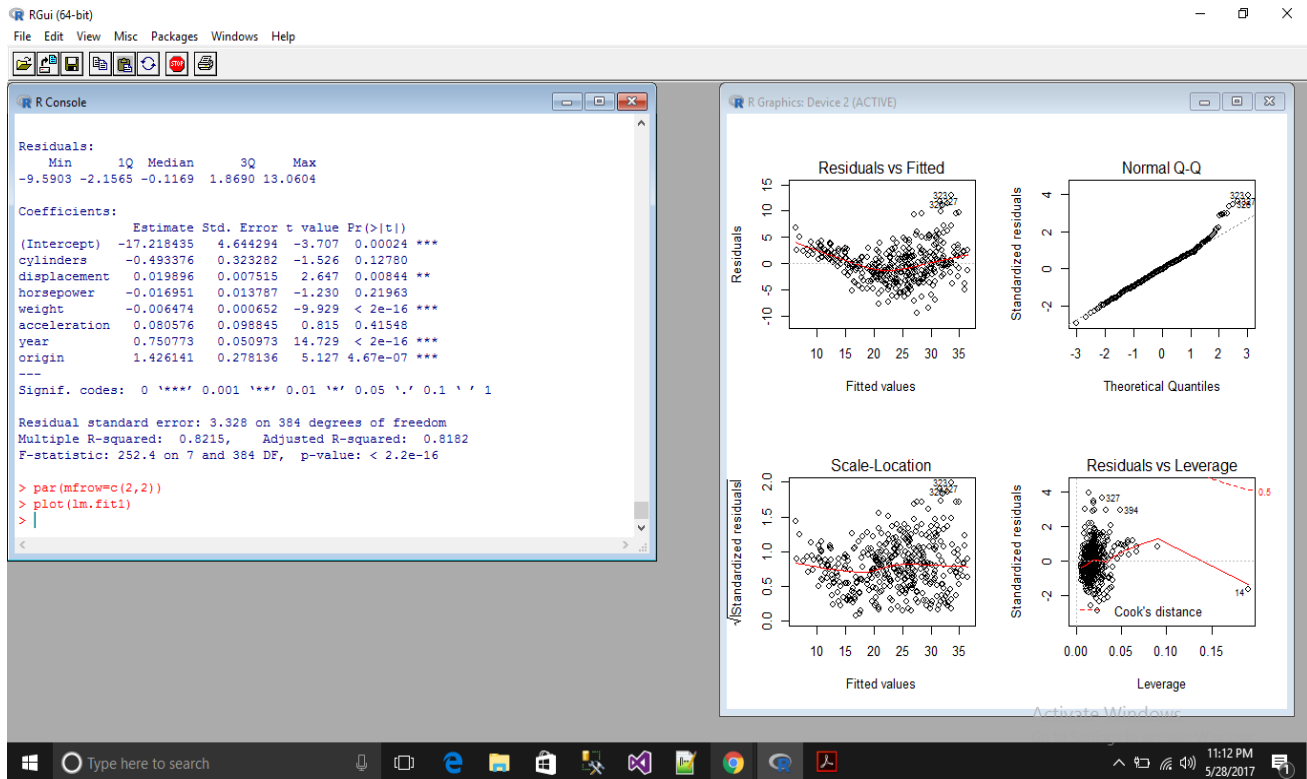


Q.9 ©



- (i) Here all the regression coefficients are zero by testing the null hypothesis. So, we can say that there is a relationship between predictors and response.
- (ii) From the p-values and t-static values, we can see that origin, weight, year and displacement have statistically significant relationship.
- (iii) Regression coefficient for the year, 0.7508, says that mpg increases by coefficient for every one year.

Q.9(d)



Q.9(e)

```
RGui (64-bit) - [R Console]
File Edit View Misc Packages Windows Help

386      buick century limited
387      oldsmobile cutlass ciera (diesel)
388      chrysler lebaron medallion
389      ford granada l
390      toyota celica gt
391      dodge charger 2.2
392      chevrolet camaro
393      ford mustang gl
394      vw pickup
395      dodge rampage
396      ford ranger
397      chevy s-10
> lm.fit2 = lm(mpg~cylinders*displacement+displacement*weight)
Error in eval(expr, envir, enclos) : object 'mpg' not found
> lm.fit2 = lm(mpg~cylinders*displacement+displacement*weight,data=Auto)
> summary(lm.fit2)

Call:
lm(formula = mpg ~ cylinders * displacement + displacement *
    weight, data = Auto)

Residuals:
    Min       1Q   Median       3Q      Max
-13.2934  -2.5184  -0.3476   1.8399  17.7723

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  5.262e+01  2.237e+00  23.519  < 2e-16 ***
cylinders     7.606e-01  7.669e-01   0.992   0.322
displacement  -7.351e-02  1.669e-02  -4.403  1.38e-05 ***
weight       -9.888e-03  1.329e-03  -7.438  6.69e-13 ***
cylinders:displacement -2.986e-03  3.426e-03  -0.872   0.384
displacement:weight  2.128e-05  5.002e-06  4.254  2.64e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.103 on 386 degrees of freedom
Multiple R-squared:  0.7272,    Adjusted R-squared:  0.7237
F-statistic: 205.8 on 5 and 386 DF,  p-value: < 2.2e-16

> |
```

Q.9(f) (6 captures)

```
RGui (64-bit) - [R Console]
File Edit View Misc Packages Windows Help

(Intercept)    5.262e+01  2.237e+00  23.519  < 2e-16 ***
cylinders       7.606e-01  7.669e-01   0.992   0.322
displacement   -7.351e-02  1.669e-02  -4.403  1.38e-05 ***
weight        -9.888e-03  1.329e-03  -7.438  6.69e-13 ***
cylinders:displacement -2.986e-03  3.426e-03  -0.872   0.384
displacement:weight  2.128e-05  5.002e-06  4.254  2.64e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.103 on 386 degrees of freedom
Multiple R-squared:  0.7272,    Adjusted R-squared:  0.7237
F-statistic: 205.8 on 5 and 386 DF,  p-value: < 2.2e-16

> lm.fit3 = lm(mpg~log(weight)+sqrt(horsepower)+acceleration+I(acceleration^2))
Error in eval(expr, envir, enclos) : object 'mpg' not found
> lm.fit3 = lm(mpg~log(weight)+sqrt(horsepower)+acceleration+I(acceleration^2),data=Auto)
> summary(lm.fit3)

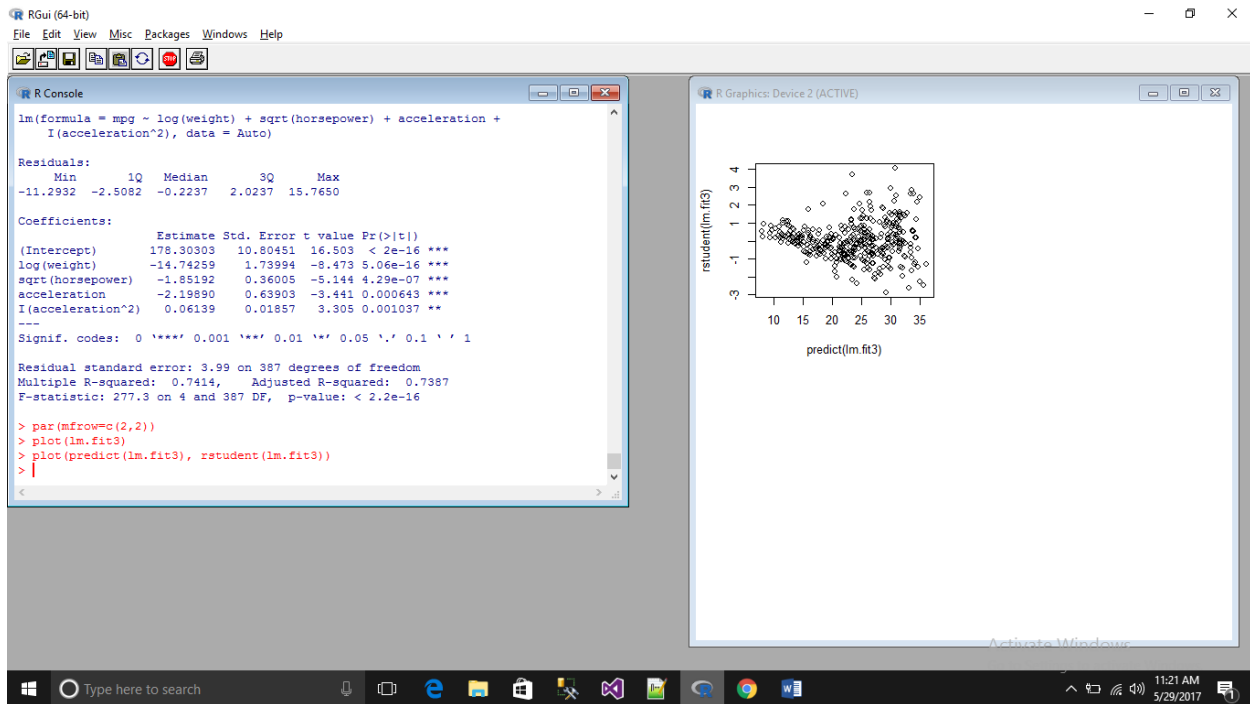
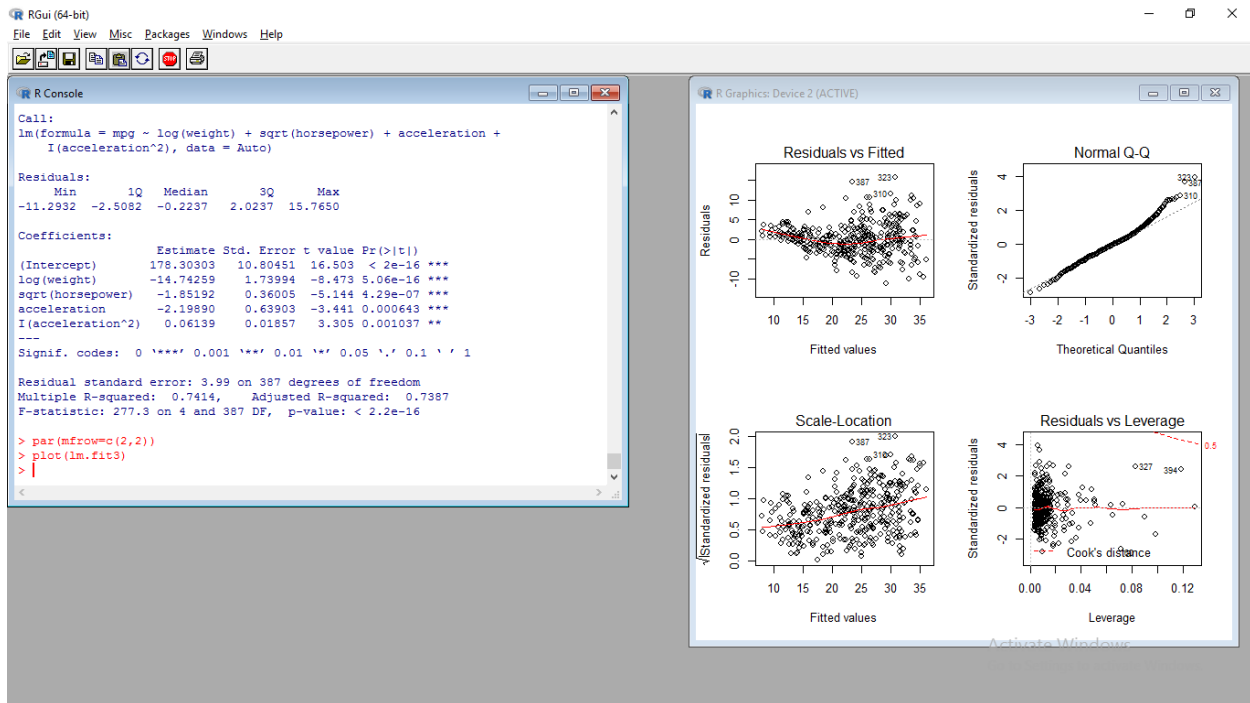
Call:
lm(formula = mpg ~ log(weight) + sqrt(horsepower) + acceleration +
    I(acceleration^2), data = Auto)

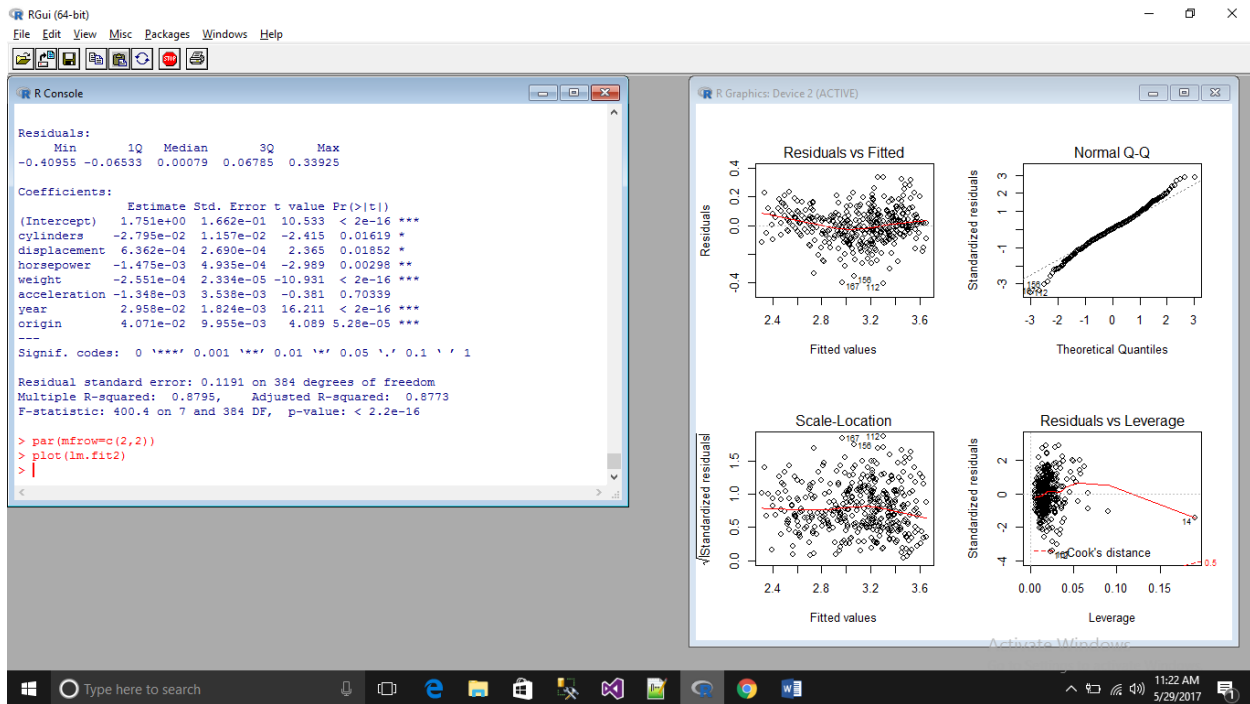
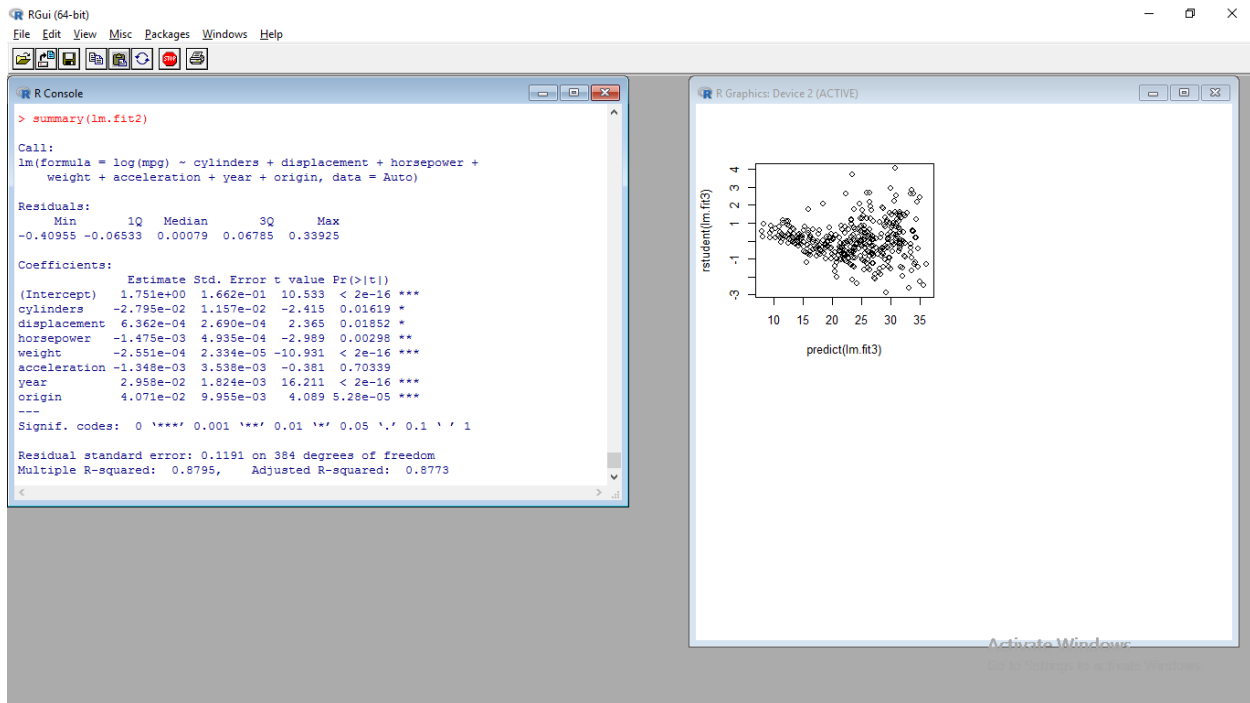
Residuals:
    Min       1Q   Median       3Q      Max
-11.2932  -2.5082  -0.2237   2.0237  15.7650

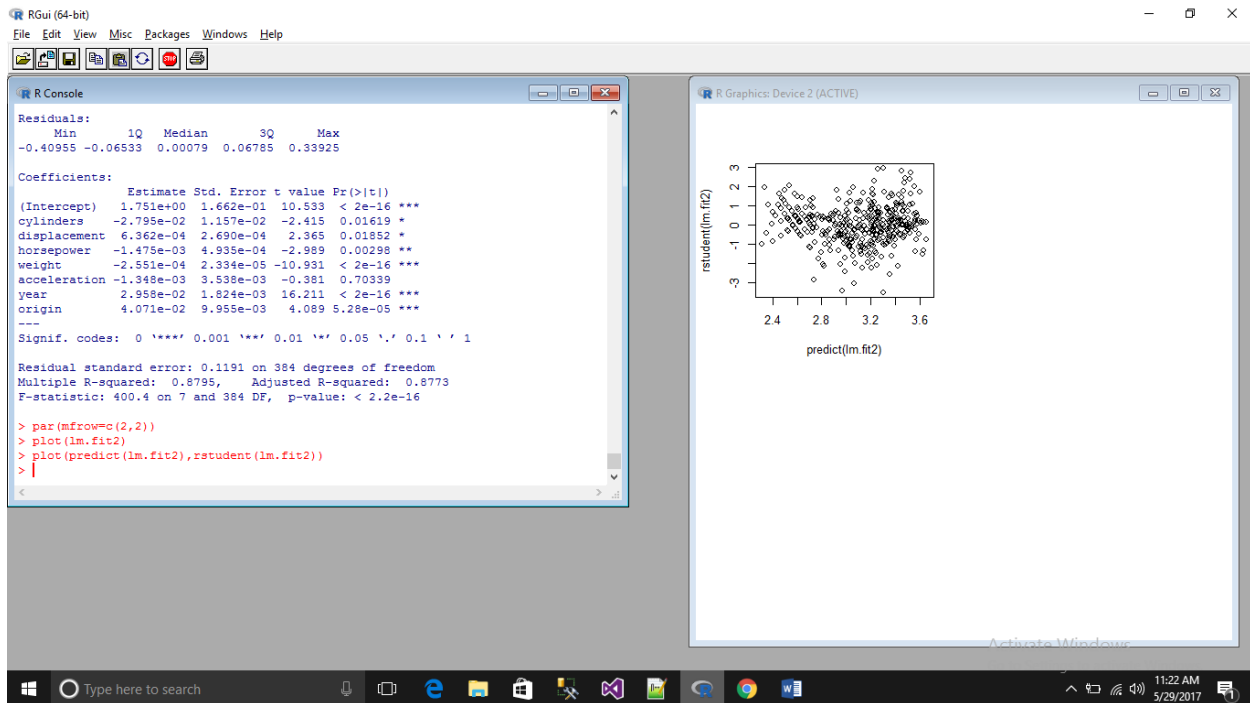
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  178.30303   10.80451   16.503  < 2e-16 ***
log(weight)  -14.74259    1.73994   -8.473  5.06e-16 ***
sqrt(horsepower) -1.85192    0.36005   -5.144  4.29e-07 ***
acceleration   -2.19890    0.63903   -3.441  0.000643 ***
I(acceleration^2)  0.06139    0.01857    3.305  0.001037 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.99 on 387 degrees of freedom
Multiple R-squared:  0.7414,    Adjusted R-squared:  0.7387
F-statistic: 277.3 on 4 and 387 DF,  p-value: < 2.2e-16

> |
```

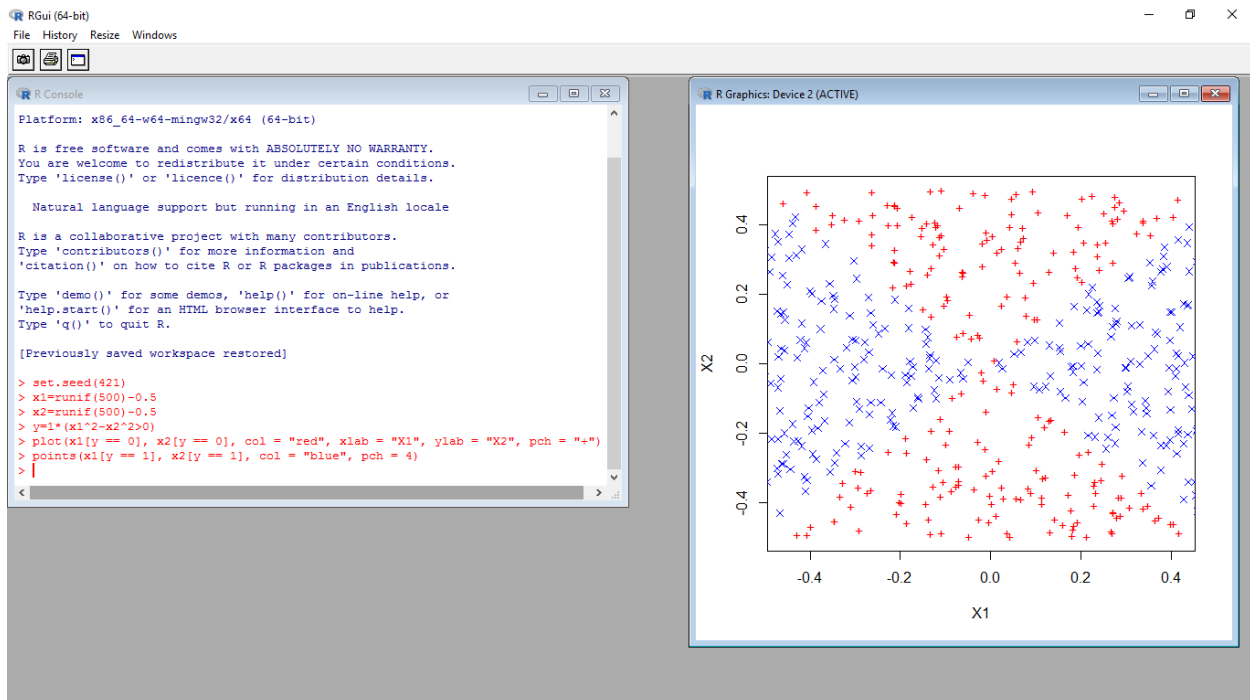




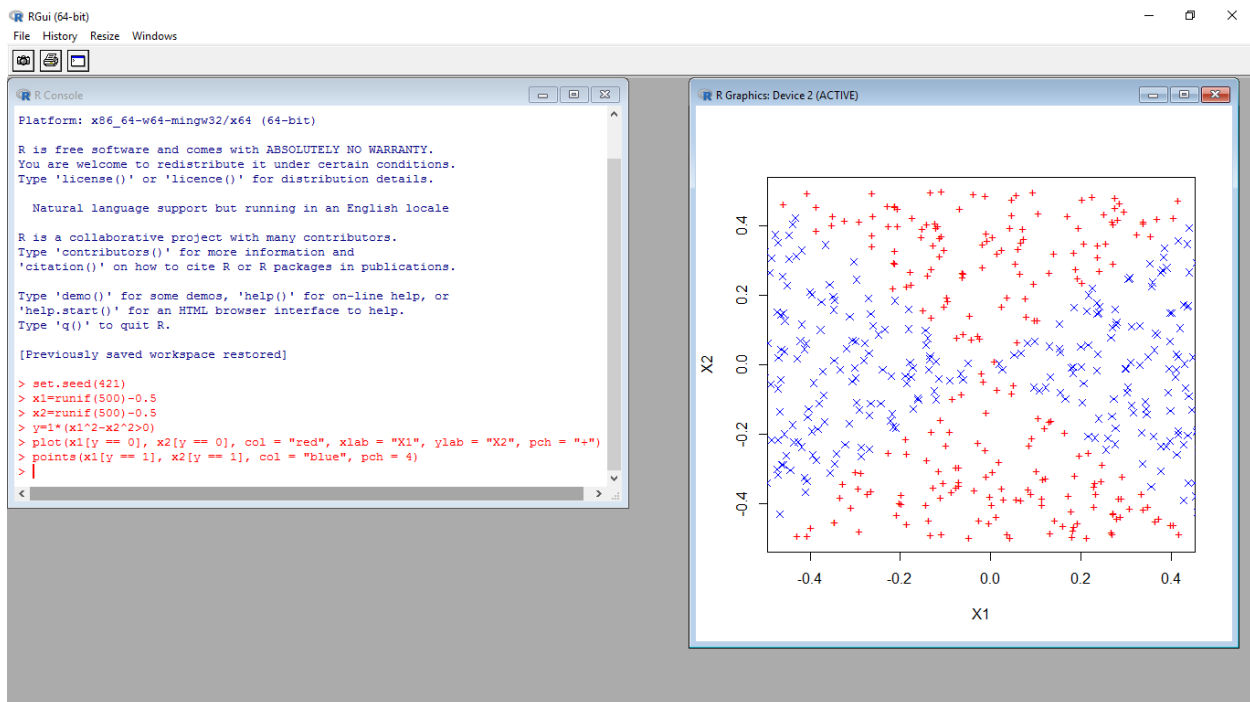


Chapter-9(Exercise 9.7)

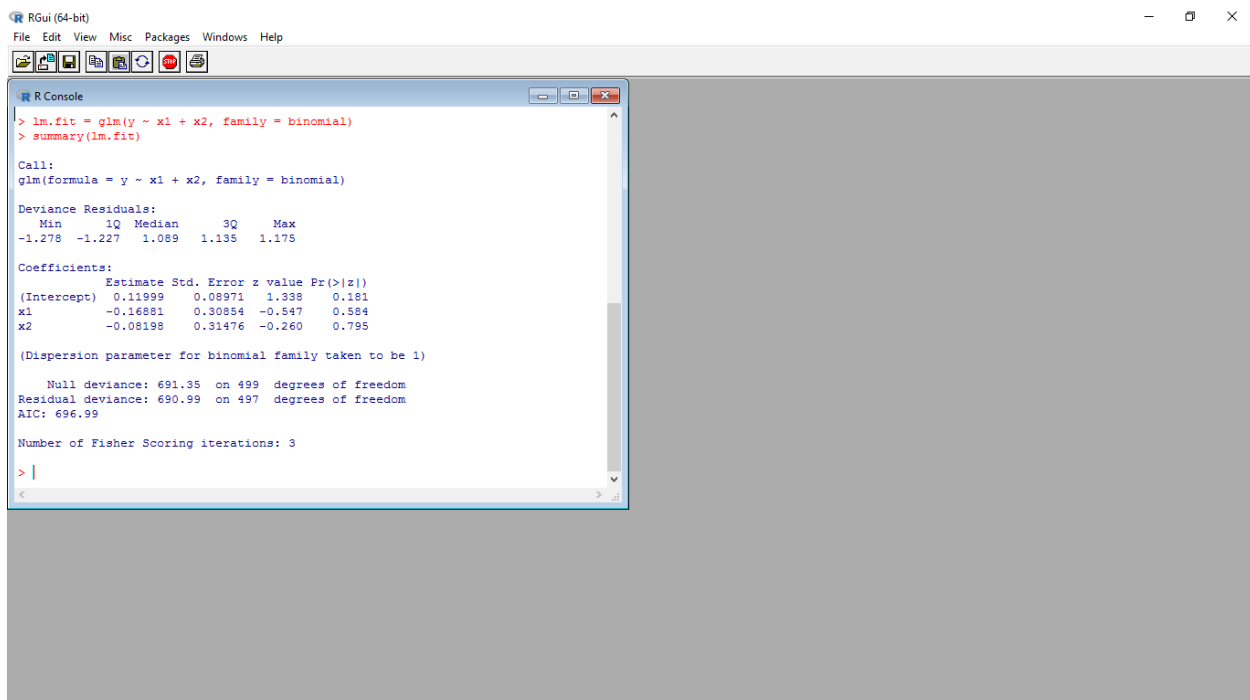
Q.5(a)



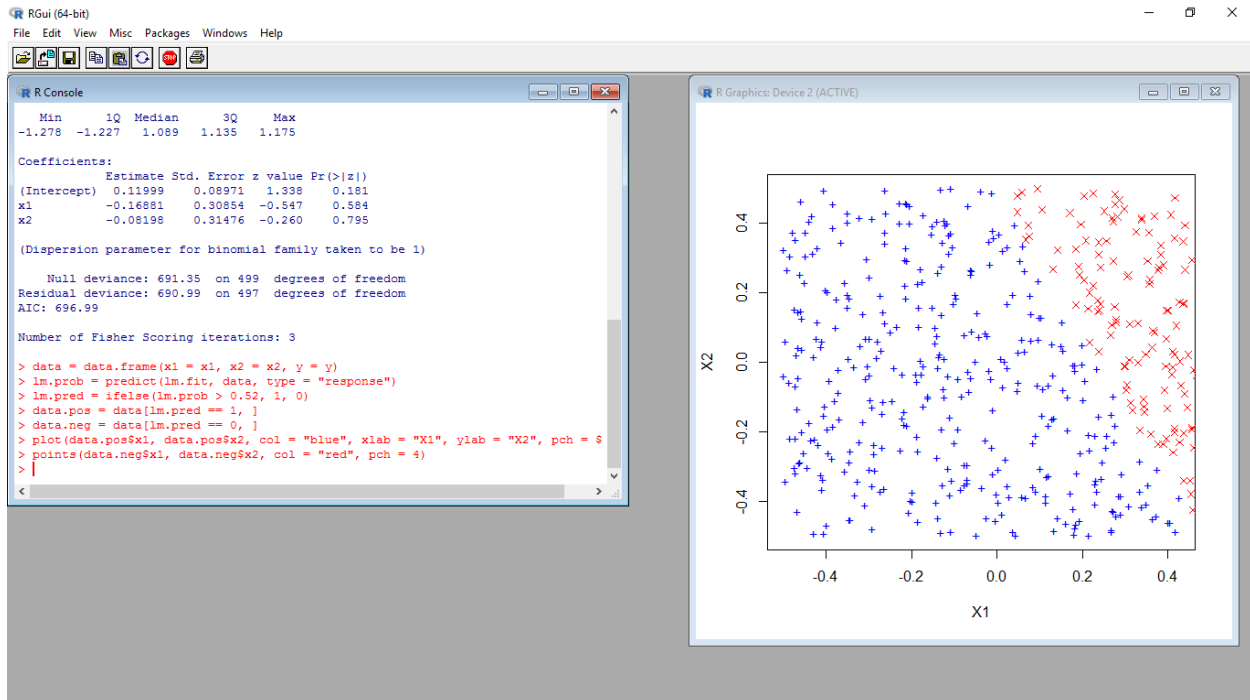
Q.5(b)



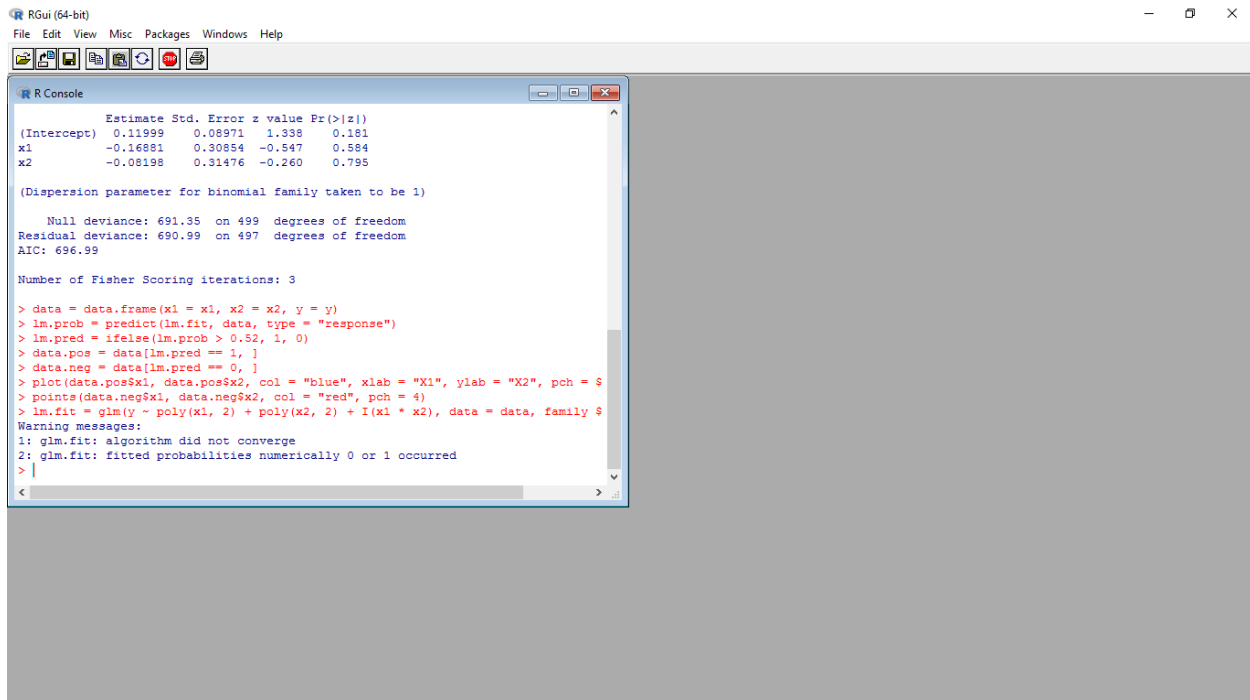
Q.5©



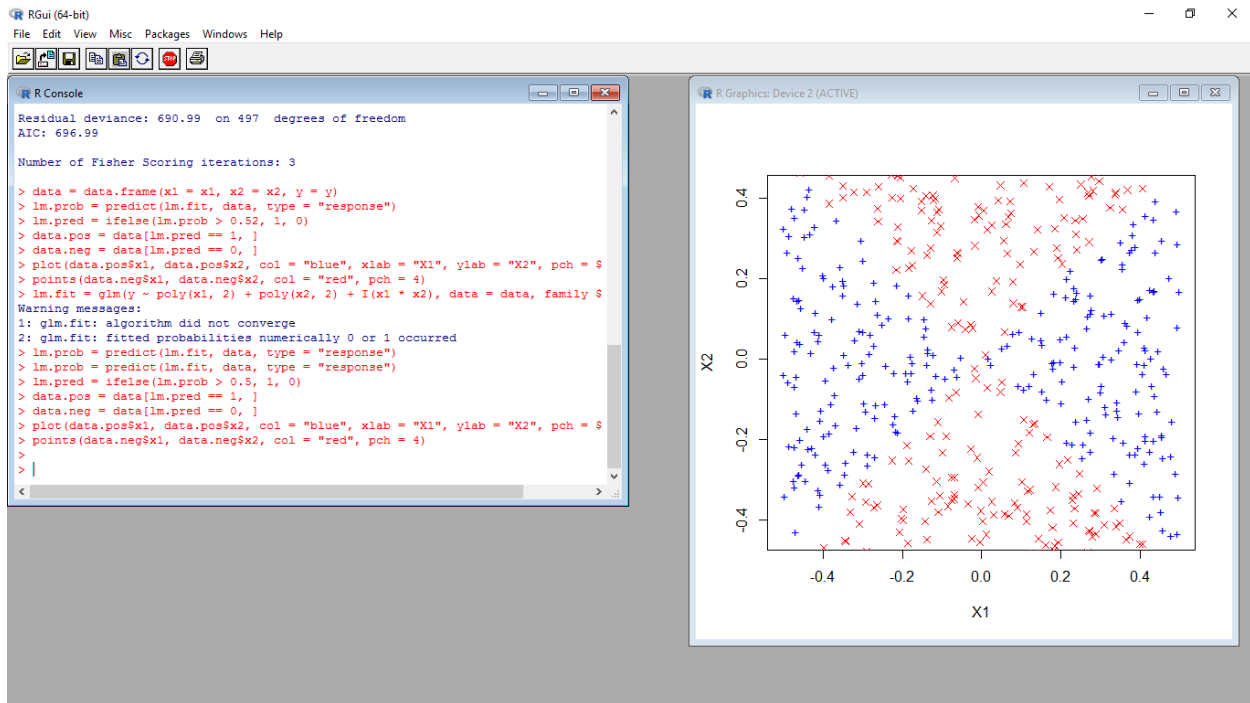
Q.5(d)



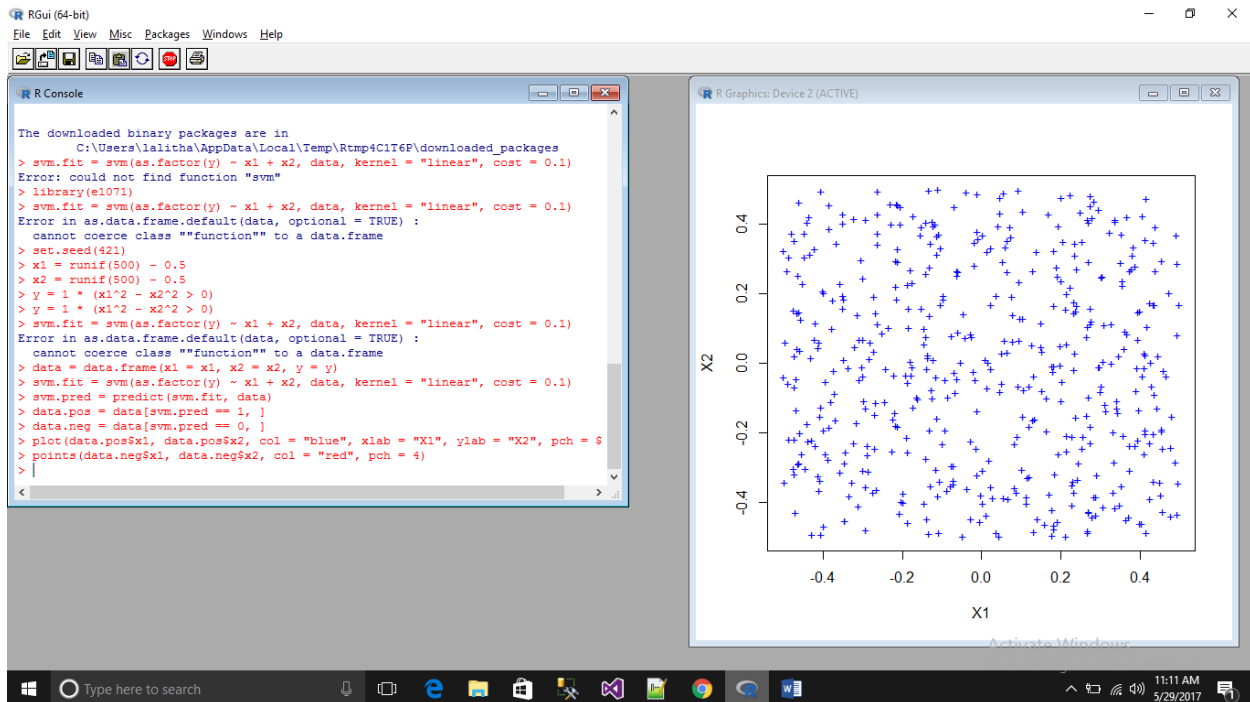
Q.5(e)



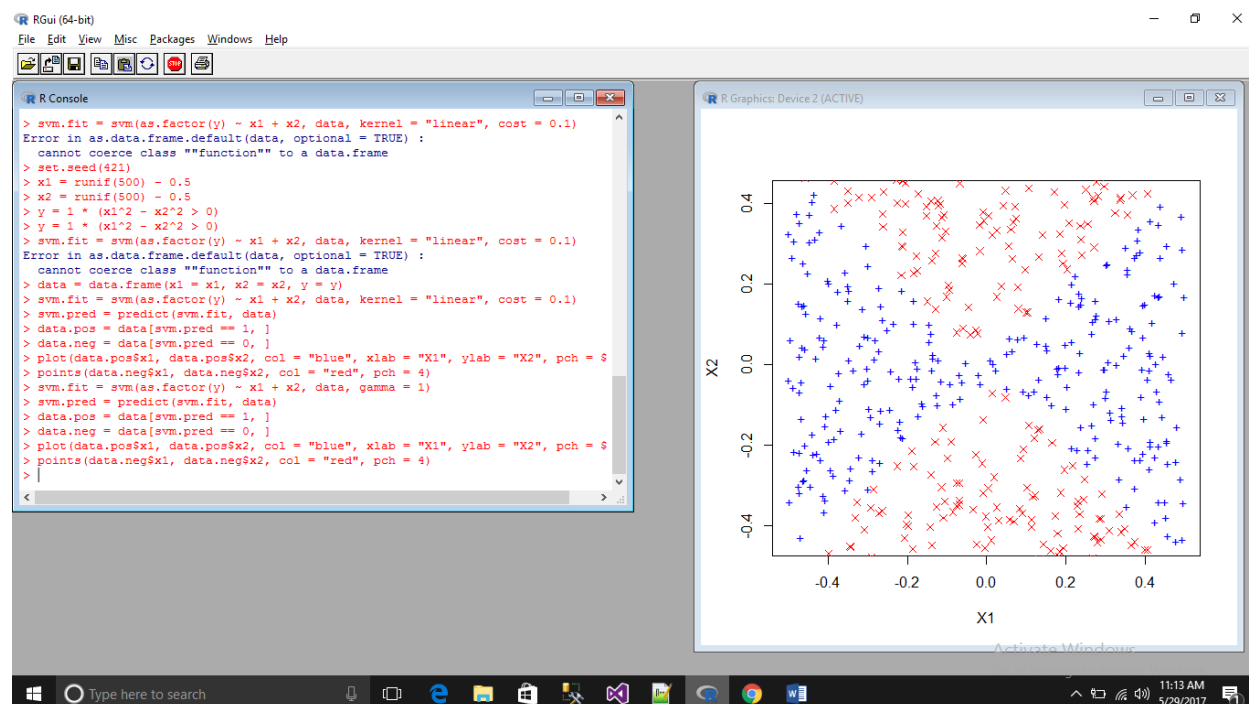
Q.5(f)



Q.5(g)



Q.5(h)

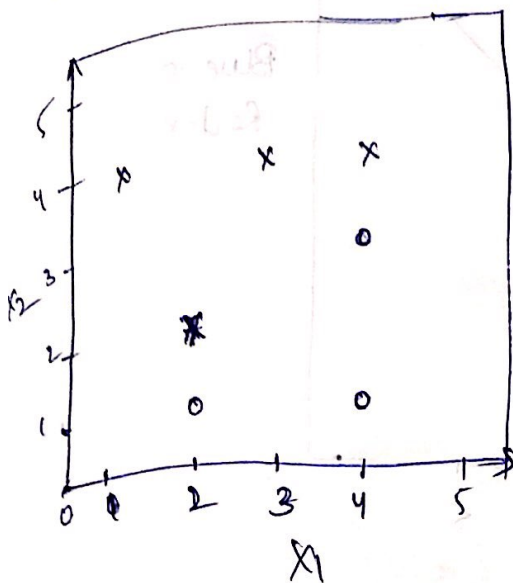


Q.5(i)

From the above results, it can be inferred that the Support Vector machines when used with non-linear kernel have greater impact in finding the non-linear boundary. However, both SVMs when used with kernels and the logistic regressions when used with non-interactions, are not successful in finding the decision boundary. If we add the interaction terms to logistic regression, it will have same impact as the radial-basis kernels. But, this requires some tuning and manual effort. This may not be fruitful if there are large number of features. On the other hand, radial basis kernels may need only tuning through cross validation of one parameter i.e gamma

chapter 9 - (Exercise 9.7)

Q.3(a) $n=7$ $p=2$

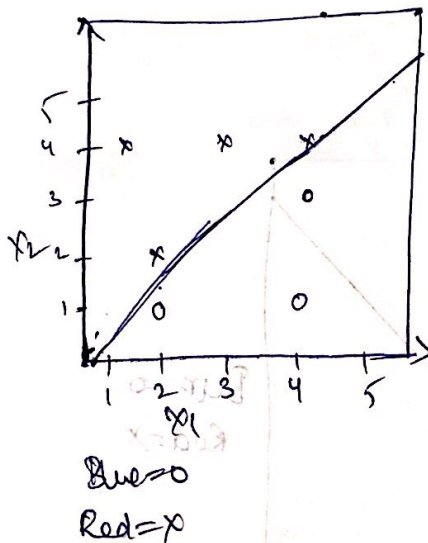


obs	X_1	X_2	Y
1	3	4	Red
2	2	2	Red
3	4	4	Red
4	1	4	Red
5	2	1	Blue
6	4	3	Blue
7	4	1	Blue

Red = x

Blue = o

b)



(Maximal Margin classifier has to be between observation 2,3 & 5,6)

So $(2, 2), (4, 4)$
 $(2, 1), (4, 3)$
 $\Rightarrow (2, 1.5), (4, 3.5)$

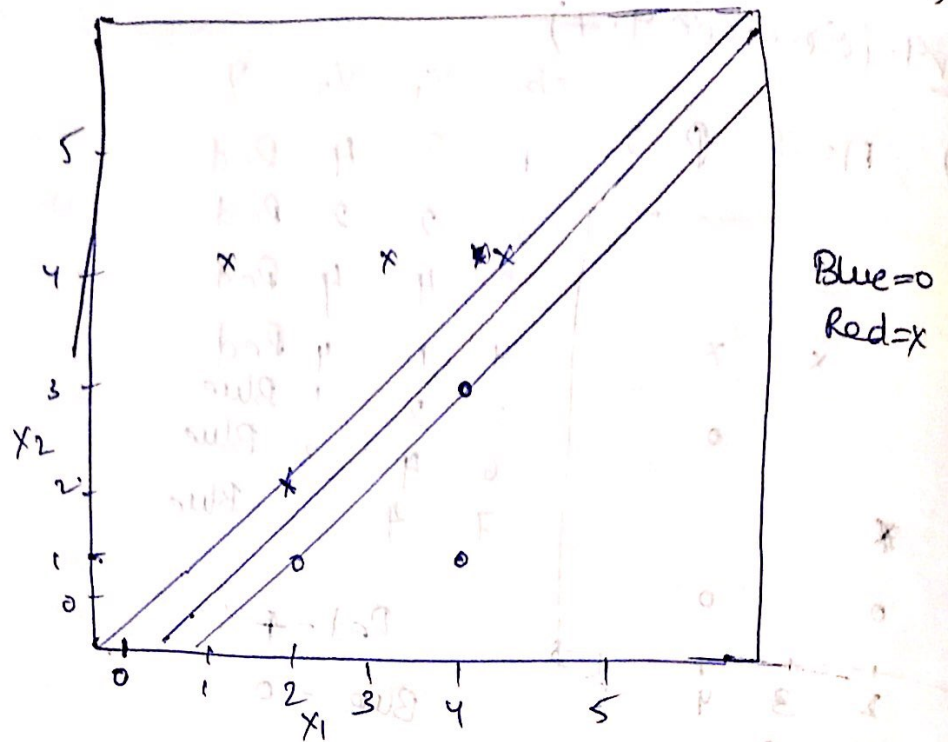
$$b = \frac{(3.5 - 1.5)}{4 - 2} = 1$$

$$a = X_2 - X_1 = 1.5 - 2 = -0.5$$

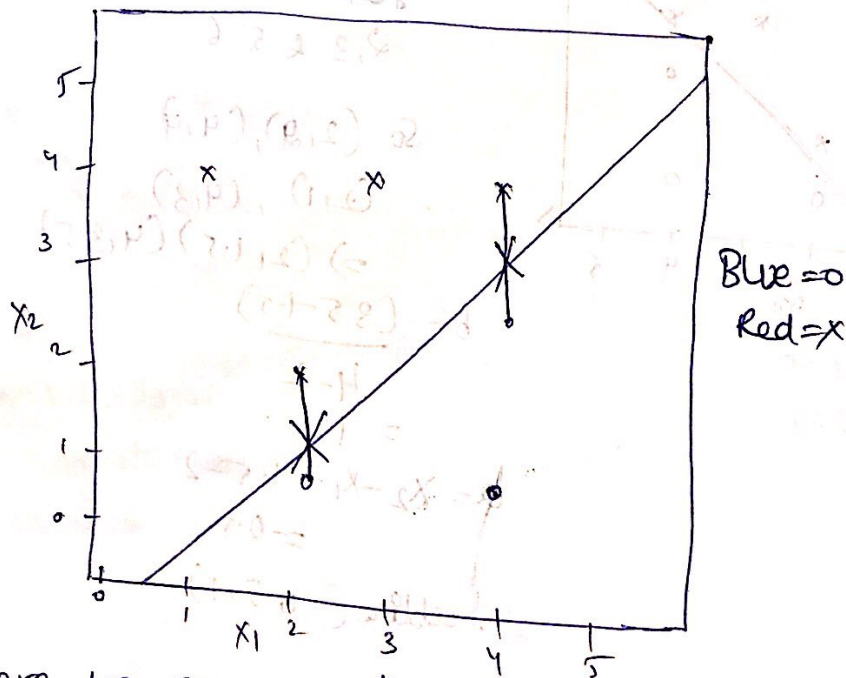
\therefore abline $(-0.5, 1)$

c) $0.5 - X_1 + X_2 \geq 0$

d)



e)

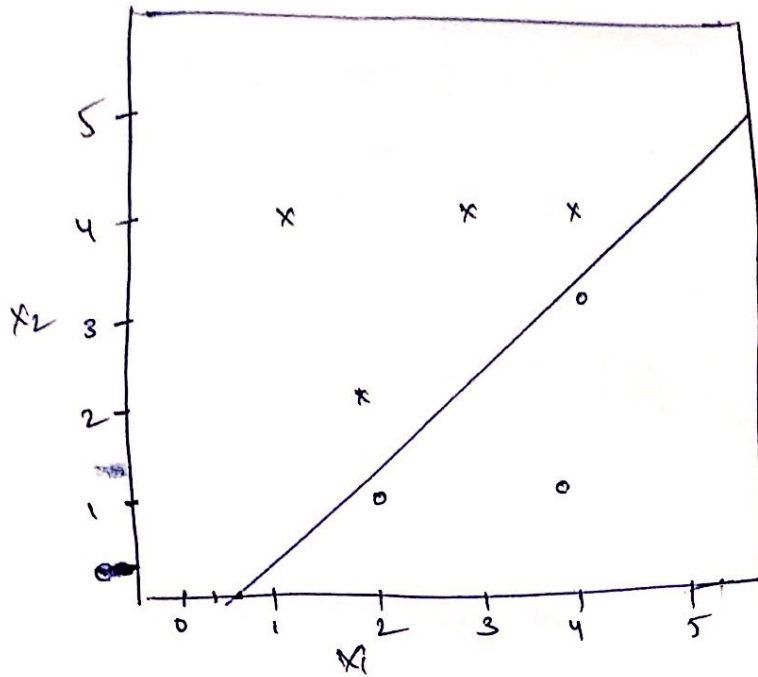


f) Here, we can see that (4,1) is outside the margin, so it would not have effect of maximal hyper plane. So, the observation 7 has no effect on maximal Margin hyperplane

abun(-0.8,1)

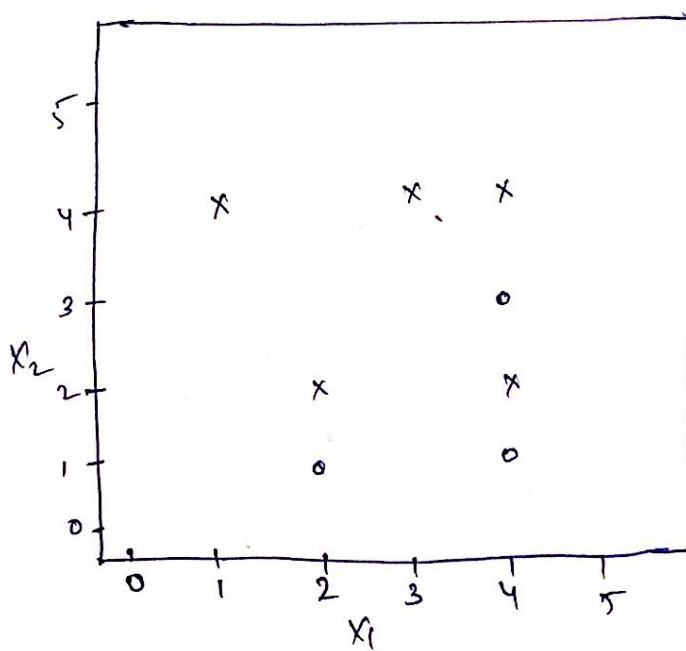
$$-0.9 - x_1 + x_2 > 0$$

g)



Blue = o
Red = x

h)



Blue = o
Red = x

Point(4,2)
col = Red