

# Data Mining

## Home work 10

### Machine Learning III Regression Analysis

Aqeel Labash  
**Lecturer:** Jaak Vilo

13 April 2016

## First Question

For this task I used python and here is the code :

```
1 # #Question 1
2 import numpy as np
3 get_ipython().magic(u'matplotlib inline')
4 import matplotlib.pyplot as plt
5 points = [(9,3,1),(2,4,1),(3,3,1),(4,1,1),(1,6,1),(3,9,0),(5,6,0),(6,4,0),(6,2,0),(3,7,0)]
6 #Draw the figure
7 plt.figure('points.jpg')
8 plt.plot([x[0] for x in points],[x[1] for x in points], 'bo')
9 plt.plot([3,4],[5,6], 'rv')
10 plt.ylabel('Y')
11 plt.xlabel('X')
12 plt.title('Points')
13 plt.savefig('points.jpg')
14
15 #Print Probabilities for classes
16 def GetClosePoints(centerpoint,k=1):
17     indx =0
18     distances={}
19     for point in points:
20
21         #Calculate Euclidean distance
22         distance = np.linalg.norm(np.array(centerpoint)-np.array((point[0],point[1])))
23
24         #Store all points with the same distance under the same
25         if distance in distances.keys():
26             distances[distance].append(indx)
27         else:
28             distances[distance] = []
29             distances[distance].append(indx)
30     indx+=1
31
32 #Sort list by distance
33 keys = distances.keys()
34 keys.sort(key = lambda x:x,reverse=False)
35 #print the list
36 #for key in keys:
37 #    print key , distances[key]
38
39 #Get Points in K distance
40 pointindx=[]
41 for i in range(0,k):
42     for index in distances[keys[i]]:
43         pointindx.append(index)
44
45 #Print Selected Points indexes
46 #print pointindx
47
48 #Calculate Probabilities
49 Totalpoints = len(pointindx)
50 Class0=0
51 Class1=0
52 for i in pointindx:
```

```

53     if points[i][2]:
54         Class1+=1
55     else:
56         Class0+=1
57     print 'For Point ({},{}) with K={},Probabilities is Class 0:{},Class 1:{}'.format(
58         centerpoint[0],centerpoint[1],k,
59         Class0/float(Totalpoints),Class1/float(Totalpoints))
60 p1 = (3,5)
61 p2 = (4,6)
62 GetClosePoints(p1,1)
63 GetClosePoints(p1,2)
64 GetClosePoints(p1,3)
65 GetClosePoints(p2,1)
66 GetClosePoints(p2,2)
67 GetClosePoints(p2,3)

```

What I did in the previous code is:

1. Store the indexes of points with same distance from our point under same hash label.
2. Select the indexes depending on K
3. Count total points,points in class0 , points in class1
4. Print the probabilities

Here is an image showing how the points look like :

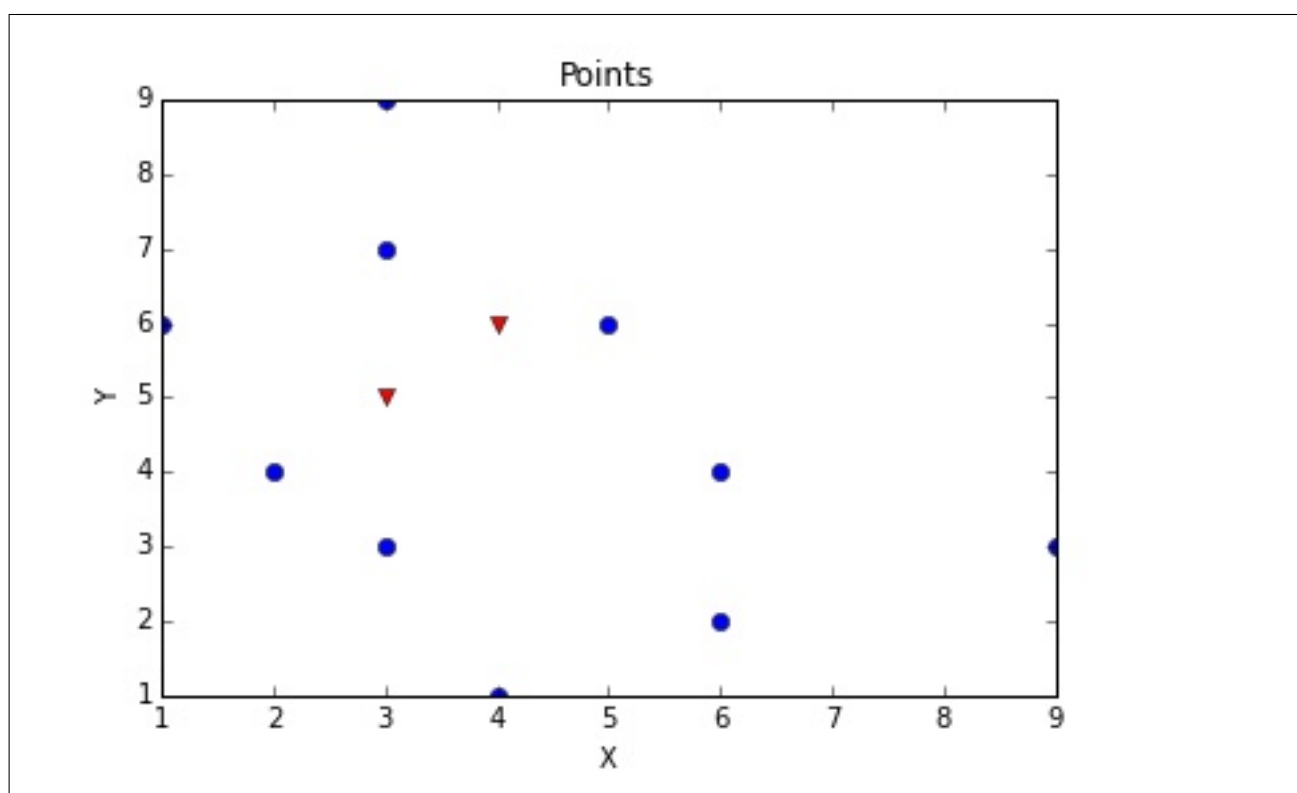


Figure 1: Points in the plane.

And here I report the output of the previous code :

```

For Point (3,5) with K=1,Probabilities is Class0=0.0,Class1=1.0
For Point (3,5) with K=2,Probabilities is Class0=0.333333333333,Class1=0.666666666667
For Point (3,5) with K=3,Probabilities is Class0=0.4,Class1=0.6
For Point (4,6) with K=1,Probabilities is Class0=1.0,Class1=0.0
For Point (4,6) with K=2,Probabilities is Class0=1.0,Class1=0.0
For Point (4,6) with K=3,Probabilities is Class0=0.75,Class1=0.25

```

## Second Question

For this question I build my own function to calculate the measurements from confusion matrix , here is the code:

```

1 measuresprint<-function(cm)
2 {
3   Accuracy<-(cm[2,2]+cm[1,1])/sum(cm)
4   Precision<-(cm[2,2])/(cm[2,2]+cm[1,2])
5   Recall<-cm[2,2]/(cm[2,2]+cm[2,1])
6   F1<-2*(Precision*Recall)/(Precision+Recall)
7   print(c(Accuracy, Precision, Recall, F1))
8 }
9 ##### Second Question #####
10 diabetes <- read.csv('pima-indians-diabetes.data',header = FALSE)
11 colnames(diabetes)<-c('PregnantTimes', 'glucos_constr', 'Blood_pressure', 'Triceps', 'insulin', '
    BMI', 'diabts_pedigree', 'Age', 'Class')
12 #Shuffle The list so we pick randomly
13 diabetes<- diabetes[sample(nrow(diabetes)),]
14 traindata<-diabetes[seq(1,floor(nrow(diabetes)*0.8),1),]
15 test<-diabetes[seq(floor(nrow(diabetes)*0.8)+1,nrow(diabetes)),]
16 module <- glm(Class~.,data = traindata)
17 summary(module)
18 png('correlationhm.png')
19 heatmap(cor(diabetes),symm = TRUE, Colv=NA, Rowv=NA,col=colorRampPalette(c("red", "yellow", "
    green"))(n = 299))
20 dev.off()
21
22 #prediction_prop<-
23 prediction_bin<-ifelse(predict(module, test)<=0.5,0,1)
24 measuresprint(table(real=test$Class, predictions=prediction_bin))

```

The summary of the module output was :

```

1 Coefficients:
2
3      Estimate Std. Error t value Pr(>|t|)
4 (Intercept)   -0.8270580  0.0952931  -8.679  < 2e-16 ***
5 PregnantTimes    0.0227378  0.0055950   4.064 5.46e-05 ***
6 glucos_constr    0.0058791  0.0005769  10.192 < 2e-16 ***
7 Blood_pressure  -0.0027814  0.0009006  -3.088  0.00211 **
8 Triceps         -0.0001970  0.0012530  -0.157  0.87511
9 insulin         -0.0002141  0.0001702  -1.258  0.20893
10 BMI            0.0135468  0.0022537   6.011 3.19e-09 ***
11 diabts_pedigree 0.1493781  0.0511308   2.921  0.00361 **
12 Age            0.0024483  0.0017115   1.430  0.15310

```

From the summary we can notice that "Plasma glucose concentration" which equals to "glucos\_constr" in table , that it's high significant to predict the class. Also we can see that it's highly correlated with the class.

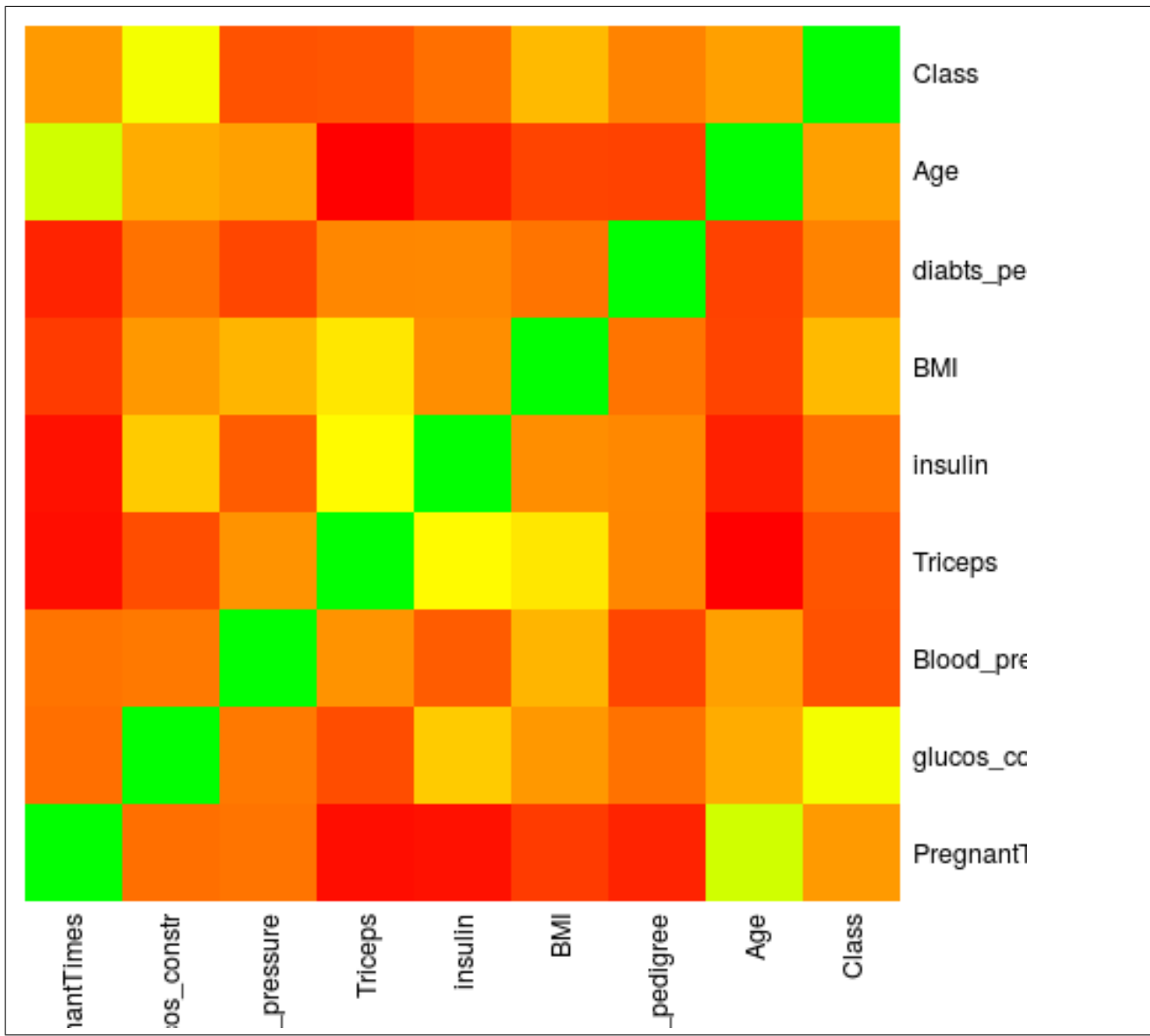


Figure 2: Heatmap for correlation

From the previous figure we can see that plasma glucos is the most correlated to the class. For "diabetes pedigree" we can see from the table that it's 2 stars significant, and from the correlation picture we can see it's less correlated than plasma glucos.

To do the calculation I printed out the confusion matrix

| real_ predic | 0  | 1  |
|--------------|----|----|
| 0            | 84 | 15 |
| 1            | 22 | 33 |

The rules for the measurements :

$$Accuracy = \frac{TP + TN}{Total}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

F1 "is the harmonic mean of precision and sensitivity" [1]

$$F1 = \frac{2TP}{2TP + FP + FN}$$

Another way to calculate it

$$F1 = 2 \cdot \frac{precision * recall}{precision + recall}$$

The measurements :

$Accuracy = 0.7597403, Precision = 0.6875000, Recall = 0.6000000, F1 = 0.6407767$

## Third Question

For this question I used the function knn in R and here is the code :

```
1 ##### Third Question#####
2 library(class)
3 knn1prediction<-knn(traindata,test = test,cl=traindata$class)
4 measuresprint(table(real=test$class, predictions=knn1prediction))
5 knn3prediction<-knn(traindata,test = test,cl=traindata$class,k = 3)
6 measuresprint(table(real=test$class, predictions=knn3prediction))
```

From the previous code we get the measurements for k1 and k3 and in the following table I put all the measurements for k1,k3 and glm to compare them

| Method | Accuracy  | Precision | Recall    | F1        |
|--------|-----------|-----------|-----------|-----------|
| glm    | 0.7597403 | 0.6875000 | 0.6000000 | 0.6407767 |
| K1     | 0.6948052 | 0.5769231 | 0.5454545 | 0.5607477 |
| K3     | 0.7272727 | 0.6181818 | 0.6181818 | 0.6181818 |

From the previous table we can see that logit regression got the best Accuracy and the highest F1 score.

## Fourth Question

For this question I used the following code :

```
1 ##### Fourth Question #####
2 rm(list=ls())
3 setwd('/home/ageel/Study/DM/HW10/')
4 library(ggplot2)
5 data("diamonds")
6 diamonds<- diamonds[sample(nrow(diamonds)),]
7 trainset<-diamonds[seq(1,floor(0.8*nrow(diamonds))),]
8 testset<-diamonds[seq(floor(0.8*nrow(diamonds))+1,nrow(diamonds)),]
9 module1<-lm(price~.,data = trainset)
10 module2<-lm(price~.+poly(carat,2)+poly(depth,2)-carat-depth,data = trainset)
11 module3<-lm(price~.+poly(carat,3)+poly(depth,3)-carat-depth,data = trainset)
12 module4<-lm(price~.+poly(carat,3)+poly(depth,3)+poly(x,2)+poly(y,2)+poly(z,2)-carat-depth-x-y-
    z,data = trainset)
13 module1predtrn<-predict(module1, trainset)
14 module2predtrn<-predict(module2, trainset)
15 module3predtrn<-predict(module3, trainset)
16 module4predtrn<-predict(module4, trainset)
17 module1predtst<-predict(module1, testset)
18 module2predtst<-predict(module2, testset)
19 module3predtst<-predict(module3, testset)
20 module4predtst<-predict(module4, testset)
21 #install.packages("qpcR")
22 library(qpcR)
23 trainRMSE<-c(sqrt(sum((module1predtrn-trainset$price)^2)/length(trainset$price)),
24 sqrt(sum((module2predtrn-trainset$price)^2)/length(trainset$price)),
25 sqrt(sum((module3predtrn-trainset$price)^2)/length(trainset$price)),
26 sqrt(sum((module4predtrn-trainset$price)^2)/length(trainset$price)))
27
28 testRMSE<-c(sqrt(sum((module1predtst-testset$price)^2)/length(testset$price)),
29 sqrt(sum((module2predtst-testset$price)^2)/length(testset$price)),
30 sqrt(sum((module3predtst-testset$price)^2)/length(testset$price)),
31 sqrt(sum((module4predtst-testset$price)^2)/length(testset$price)))
32 numbers<-seq(1,4,1)
33 png('train-test.png')
34 ggplot(data.frame(cbind(numbers,trainRMSE)),aes(numbers, trainRMSE))+geom_line(col="red")+
35   geom_line(aes(numbers,testRMSE),data =data.frame(cbind(numbers,testRMSE)),col="blue")+ xlab
    ("Module")+
36   ylab("RSME")
```

The previous code will output the following picture which represent the Train vs Test RMSE value.

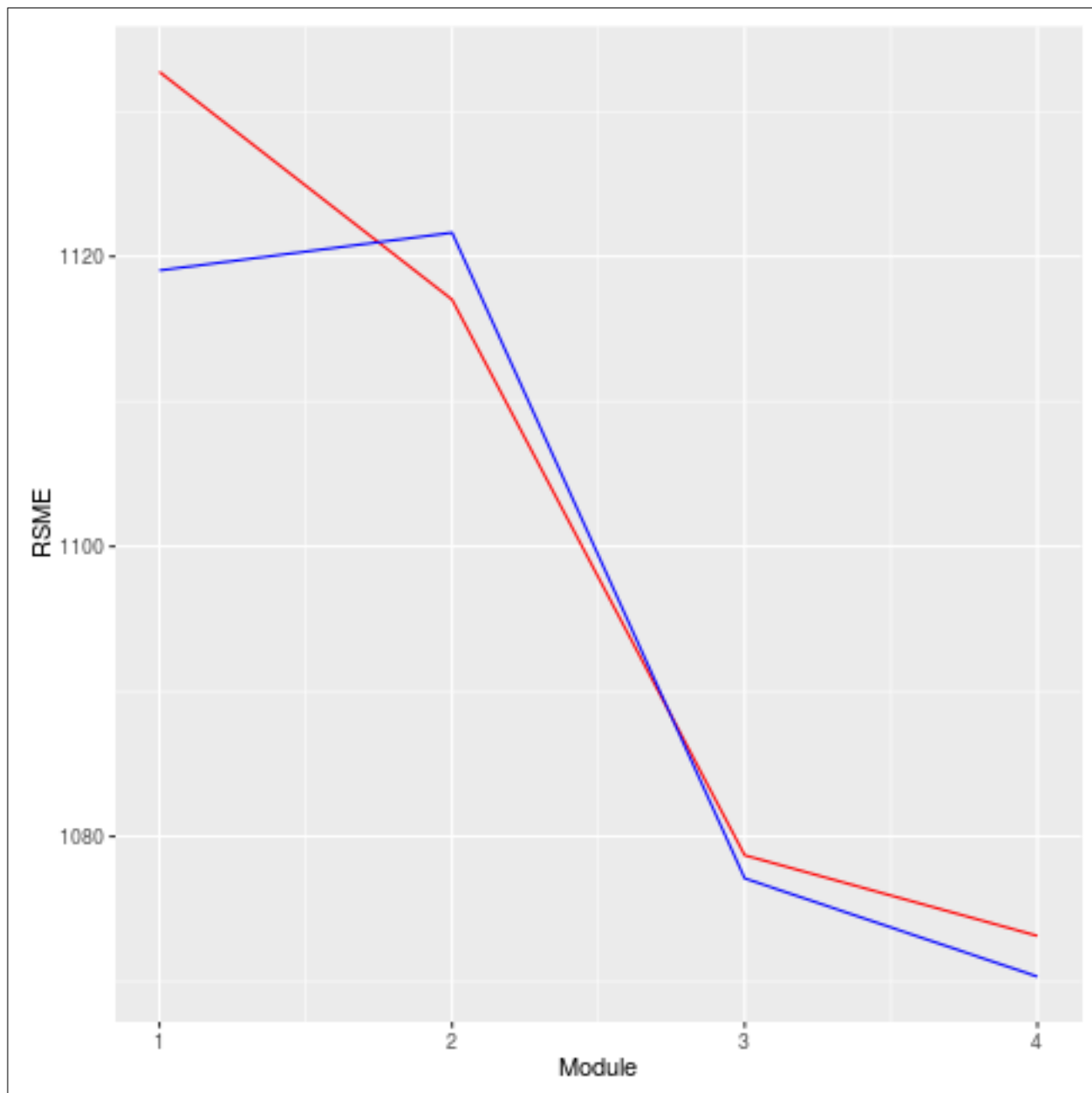


Figure 3: Train RMSE in Red , Test in blue

I would detect overfitting when the training RMSE keep decreasing while the testing RMSE start increasing. I believe we can see that from the plot after doing many iteration for the same module. With only one iteration it's hard to decide if we have overfitting mmmm maybe we can if the RMSE was pretty low for training and pretty high for testing.

Also I would notice from the plot that module 4 performed the best :).

## Fifth Question

For this question I used straight forward solution and here is the code :

```

1 rm(list=ls())
2 setwd('/home/ageel/Study/DM/HW10/')
3 train <- read.csv('train.csv', header = TRUE)
4 module<-lm(data = train, target~.)
5 summary(train)
6 test<-read.csv('test.csv', header = TRUE)
7 head(test[,c(1,2)])
8 result<-predict(module, test)
9 output<-cbind(result)
10 colnames(output)<-c('ID', 'target')
```

```
11
12 write.csv(output, file='submit 001')
```

I got score 1.95523 for this code :)

## Sixth Question

## Seventh Question

The code for this question :

```
1 ##### Seventh Question #####
2 library(MASS)
3 module4ridge<- lm.ridge(price~.+poly(carat,3)+poly(depth,3)+poly(x,2)+poly(y,2)+poly(z,2)-
  carat-depth-x-y-z,data = trainset)
4 module4ridge.trn.prd = as.matrix(model.matrix(price~.+poly(carat,3)+poly(depth,3)+poly(x,2)+
  poly(y,2)+poly(z,2)-carat-depth-x-y-z,trainset))%%coef(module4ridge)
5 module4ridge.tst.prd = as.matrix(model.matrix(price~.+poly(carat,3)+poly(depth,3)+poly(x,2)+
  poly(y,2)+poly(z,2)-carat-depth-x-y-z,testset))%%coef(module4ridge)
6
7
8 #install.packages("lars")
9 library(lars)
10 module4lasso <- lars(
11   model.matrix(price~.+poly(carat,3)+poly(depth,3)+poly(x,2)+poly(y,2)+poly(z,2)-carat-depth-x-
    -y-z,trainset),
12   trainset$price, type="lasso", trace = TRUE, max.steps=20)
13 module4lasso.trn.prd <- predict(module4lasso,
14   model.matrix(price~.+poly(carat,3)+poly(depth,3)+poly(x,2)+poly(y,2)+
    poly(z,2)-carat-depth-x-y-z,trainset),
15   s=module4lasso$df[which.min(module4lasso$RSS)], type="fit")$fit
16
17
18 module4lasso.tst.prd <- predict(module4lasso,
19   model.matrix(price~.+poly(carat,3)+poly(depth,3)+poly(x,2)+
    poly(y,2)+poly(z,2)-carat-depth-x-y-z,testset),
20   s=module4lasso$df[which.min(module4lasso$RSS)], type="fit")$
  fit
21
22 trainRMSE<-c(sqrt(sum((module1predtrn-trainset$price)^2)/length(trainset$price)),
23   sqrt(sum((module2predtrn-trainset$price)^2)/length(trainset$price)),
24   sqrt(sum((module3predtrn-trainset$price)^2)/length(trainset$price)),
25   sqrt(sum((module4predtrn-trainset$price)^2)/length(trainset$price)),
26   sqrt(sum((module4lasso.trn.prd-trainset$price)^2)/length(trainset$price)),
27   sqrt(sum((module4ridge.trn.prd-trainset$price)^2)/length(trainset$price)))
28
29 testRMSE<-c(sqrt(sum((module1predtst-testset$price)^2)/length(testset$price)),
30   sqrt(sum((module2predtst-testset$price)^2)/length(testset$price)),
31   sqrt(sum((module3predtst-testset$price)^2)/length(testset$price)),
32   sqrt(sum((module4predtst-testset$price)^2)/length(testset$price)),
33   sqrt(sum((module4lasso.tst.prd-testset$price)^2)/length(testset$price)),
34   sqrt(sum((module4ridge.tst.prd-testset$price)^2)/length(testset$price)))
35 numbers<-seq(1,6,1)
36 png('all_modules.png')
37 ggplot(data.frame(cbind(numbers,trainRMSE)),aes(numbers, trainRMSE))+geom_line(col="red") +
38   geom_line(aes(numbers,testRMSE),data =data.frame(cbind(numbers,testRMSE)),col="blue") + xlab
  ("Module")+
39   ylab("RSME")
40 dev.off()
```

The RMSE for lasso and ridge was quite high so it won't be much clear in the graph but here is the graph :

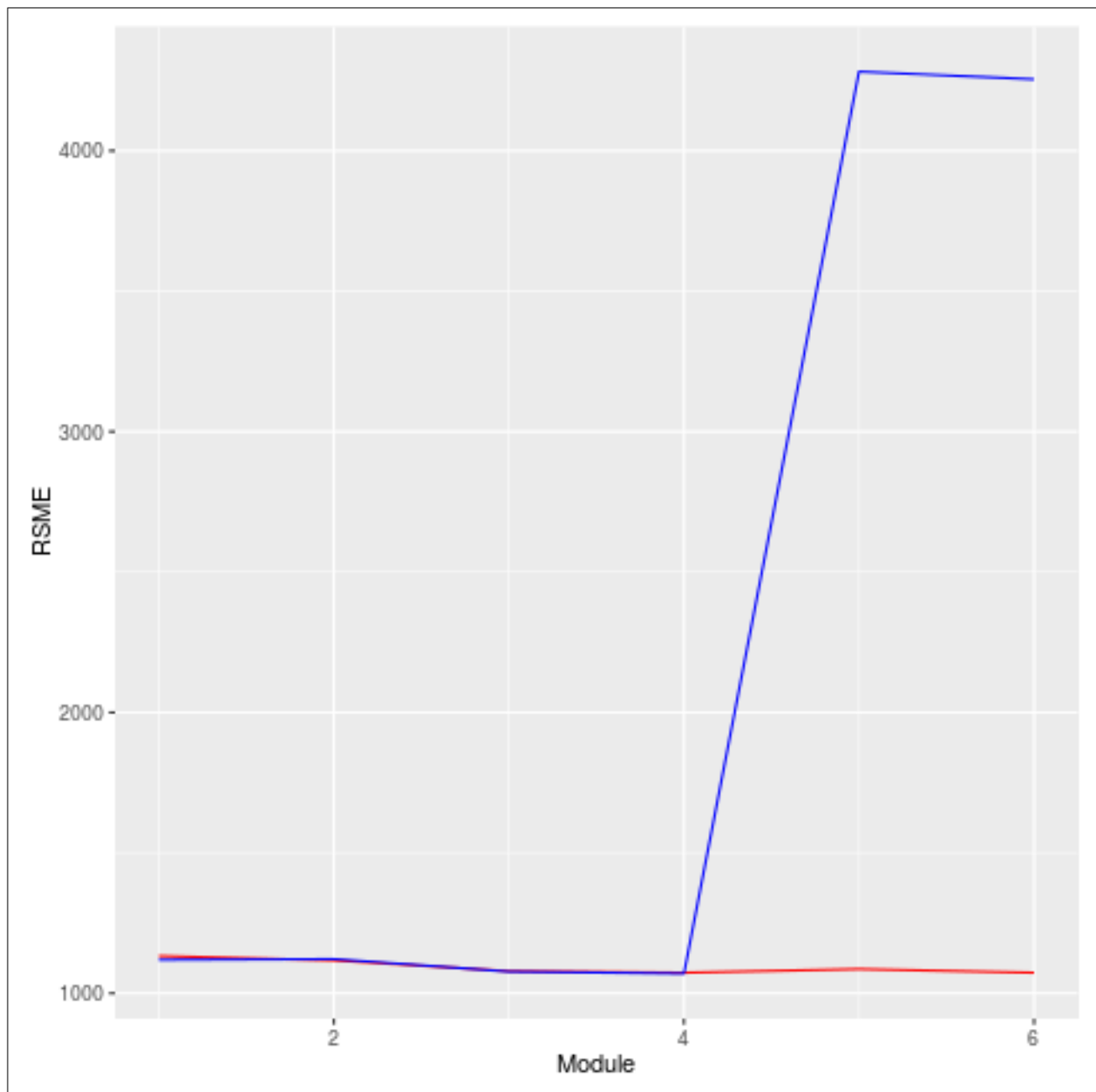


Figure 4: Last two values are lasso and ridge

**Please note:** Faiz helped a lot in this question.

## References

- [1] Precision and Recall