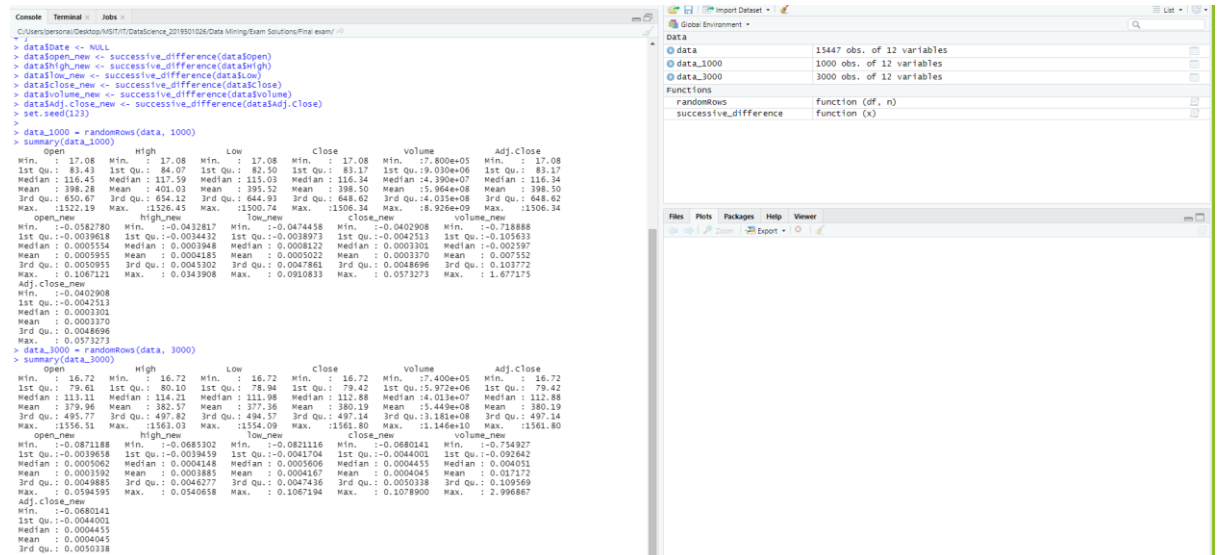# Data Mining - Lab Exam – Remedial

1. a) For the dataset BSE_Sensex_Index.csv, create an extra column of successive differences for each column of numeric values in this data file. Extract two simple random samples with replacement of 1000 and 3000 observations (rows). Show your R commands for doing this. Do the same thing by using Excel. Show your Excel commands.

   **Note:** Successive difference for date d1= (date d1 value-previous date of d1 value)/previous date of d1. For the last row fill up values with mean of its immediate three previous row values.



   b) For your samples, use the functions mean(), max(), var() and quantile(,.25) to compute the mean, maximum, variance and 1st quartile respectively for each column which has successive differences. Show your R code and the resulting values.

   Do the same thing by using Excel. Show your Excel commands.

```
> mean(data_1000$open_new)
[1] 0.0005955025
> mean(data_1000$high_new)
[1] 0.0004184797
> mean(data_1000$low_new)
[1] 0.0005022487
> mean(data_1000$close_new)
[1] 0.0003369592
> mean(data_1000$volume_new)
[1] 0.007551912
> mean(data_1000$Adj.close_new)
[1] 0.0003369592
> var(data_1000$open_new)
[1] 8.714339e-05
> var(data_1000$high_new)
[1] 6.119132e-05
> var(data_1000$low_new)
[1] 8.313995e-05
> var(data_1000$close_new)
[1] 7.637739e-05
> var(data_1000$volume_new)
[1] 0.0327711
> var(data_1000$Adj.close_new)
[1] 7.637739e-05
> max(data_1000$open_new)
[1] 0.1067121
> max(data_1000$high_new)
[1] 0.03439077
> max(data_1000$low_new)
[1] 0.09108332
> max(data_1000$close_new)
[1] 0.05732732
> max(data_1000$volume_new)
[1] 1.677175
> max(data_1000$Adj.close_new)
[1] 0.05732732
> quantile(data_1000$open_new,0.25)
        25%
-0.003961827
> quantile(data_1000$high_new,0.25)
        25%
-0.003443228
> quantile(data_1000$low_new,0.25)
        25%
-0.003897353
> quantile(data_1000$close_new,0.25)
        25%
-0.004251294
> quantile(data_1000$volume_new,0.25)
       25%
-0.1056329
> quantile(data_1000$Adj.close_new,0.25)
        25%
-0.004251294
>
```

```
>
> mean(data_3000$open_new)
[1] 0.0003591911
> mean(data_3000$high_new)
[1] 0.0003884621
> mean(data_3000$low_new)
[1] 0.0004167
> mean(data_3000$close_new)
[1] 0.0004044752
> mean(data_3000$volume_new)
[1] 0.0171718
> mean(data_3000$Adj.close_new)
[1] 0.0004044752
> var(data_3000$open_new)
[1] 8.509529e-05
> var(data_3000$high_new)
[1] 6.81047e-05
> var(data_3000$low_new)
[1] 8.768766e-05
> var(data_3000$close_new)
[1] 8.588174e-05
> var(data_3000$volume_new)
[1] 0.03939109
> var(data_3000$Adj.close_new)
[1] 8.588174e-05
> max(data_3000$open_new)
[1] 0.05945946
> max(data_3000$high_new)
[1] 0.05406578
> max(data_3000$low_new)
[1] 0.1067194
> max(data_3000$close_new)
[1] 0.10789
> max(data_3000$volume_new)
[1] 2.996867
> max(data_3000$Adj.close_new)
[1] 0.10789
> quantile(data_3000$open_new,0.25)
       25%
-0.003965834
> quantile(data_3000$high_new,0.25)
       25%
-0.003945885
> quantile(data_3000$low_new,0.25)
       25%
-0.004170403
> quantile(data_3000$close_new,0.25)
       25%
-0.00440009
> quantile(data_3000$volume_new,0.25)
       25%
-0.09264194
> quantile(data_3000$Adj.close_new,0.25)
       25%
-0.00440009
> |
```

c) Compute the same quantities in part b on the entire data set and show your answers.
How much do they differ from your answers in part b? Do you find any significant difference

between two sample values like mean in comparison with entire data? If so what explanation you can give for that?

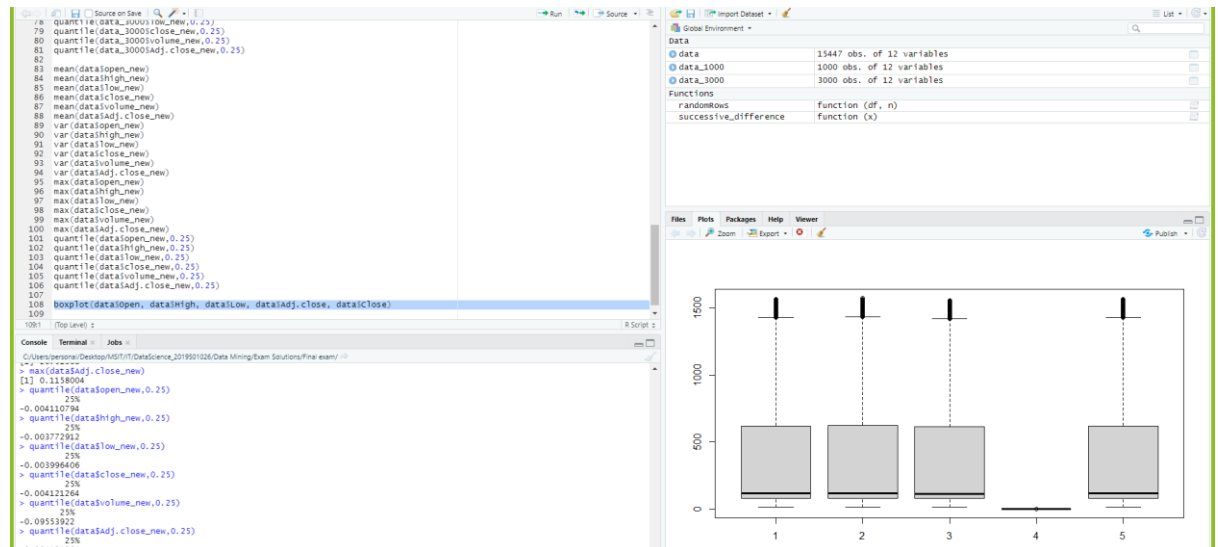Do the same thing by using Excel. Show your Excel commands.

```
 25%
-0.00440009
> mean(data$open_new)
[1] 0.000329528
> mean(data$high_new)
[1] 0.0003188991
> mean(data$low_new)
[1] 0.0003266191
> mean(data$close_new)
[1] 0.0003303709
> mean(data$volume_new)
[1] 0.02062874
> mean(data$Adj.close_new)
[1] 0.0003303709
> var(data$open_new)
[1] 9.027493e-05
> var(data$high_new)
[1] 6.939914e-05
> var(data$low_new)
[1] 8.646474e-05
> var(data$close_new)
[1] 9.350347e-05
> var(data$volume_new)
[1] 0.09080738
> var(data$Adj.close_new)
[1] 9.350347e-05
> max(data$open_new)
[1] 0.1067121
> max(data$high_new)
[1] 0.08037943
> max(data$low_new)
[1] 0.1067194
> max(data$close_new)
[1] 0.1158004
> max(data$volume_new)
[1] 26.51968
> max(data$Adj.close_new)
[1] 0.1158004
> quantile(data$open_new,0.25)
         25%
-0.004110794
> quantile(data$high_new,0.25)
         25%
-0.003772912
> quantile(data$low_new,0.25)
         25%
-0.003996406
> quantile(data$close_new,0.25)
         25%
-0.004121264
> quantile(data$volume_new,0.25)
        25%
-0.09553922
> quantile(data$Adj.close_new,0.25)
         25%
-0.004121264
>
```
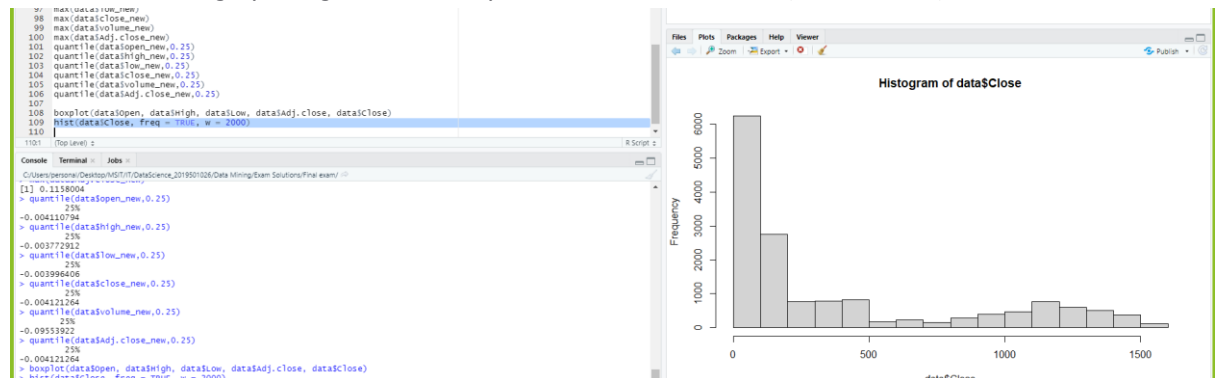
d) Use R to produce a single graph displaying a boxplot for open, close, high and low. Include the R commands and the plot.

Do the same thing by using Excel. Show your Excel commands



e) Use R to produce a frequency histogram for Close values. Use intervals of width 2000 beginning at 0. Include the R commands and the plot.

Do the same thing by using Excel. Show your Excel commands.    (10+10=20M)



2. Implement Apriori Algorithm or use built in packages to find out the frequent itemsets and generate rules for frequent itemsets. Trace program output for the following given dataset of transactions with a minimum support of 3 and submit.    (10M)

| TID, Items |
|---|
| 101, A,B,C,D,E |
| 102, A,C,D |
| 103, D,E |
| 104, B,C,E |
| 105, A,B,D,E |
| 106, A,B |
| 107, B,D,E |
| 108, A,B,D |
| 109, A,D |
| 110, D,E |

3. Build Decision Trees by using i) information gain and ii) misclassification error rate for Lenses Data Set provided at http://archive.ics.uci.edu/ml/datasets/Lenses.  In terms of tree size what do you conclude comparing these two?
   (10M)

4. Fit 1, 2 and 3-nearest-neighbor classifiers to the Liver Disorders Data Set at http://archive.ics.uci.edu/ml/datasets/Liver+Disorders for measures Euclidean and cosine. Last but one column is a decision attribute. Replace decision values in to 4 classes ($0<=c1<5$, $5<=c2<10$, $10<=c3<15$, $15<=c4<=20$). Last column is a data split column in to training and test sets. 1 means the object is used for training. 2 means the object is used for testing. Explain the input parameters you provided for the classifier.  Compute the misclassification error on the training data and also on the test data. Annotate your program. (10M)

5. Use Support Vector machine for above problem. And compare the performance of both. Explain the input parameters you provided for the classifier. (10M)
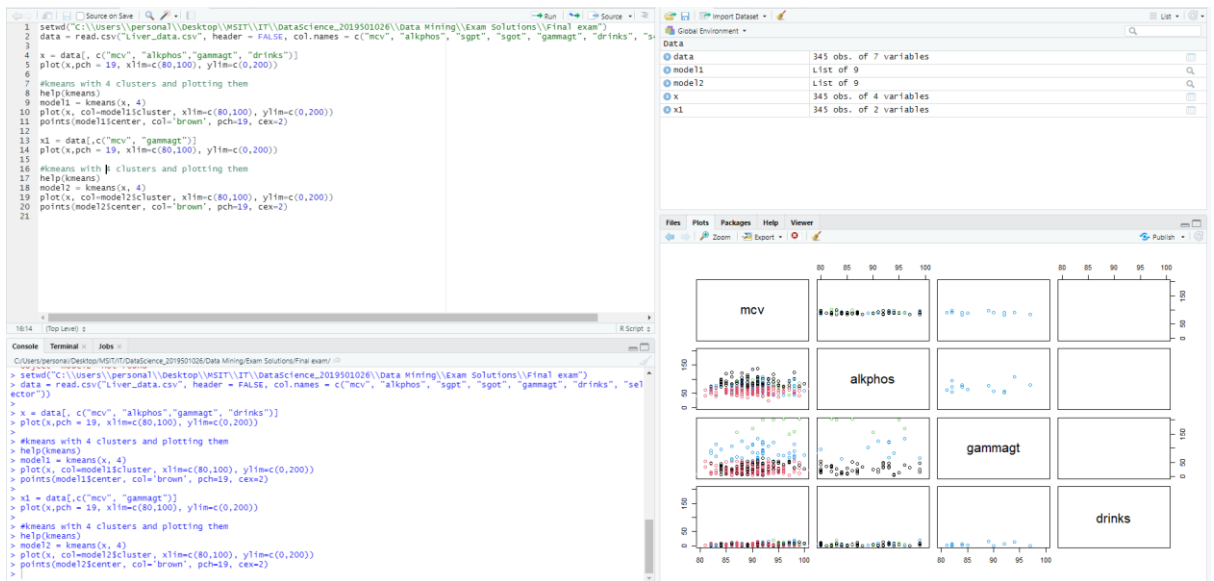


6. Create k-means clusters for k=4 for the Liver Disorders Data Set at http://archive.ics.uci.edu/ml/datasets/Liver+Disorders . Explain the input parameters you provided for the clustering algorithm. Plot the fitted cluster centers using a different color. Finally assign the cluster membership for the points to the nearest cluster center. Color the points according to their cluster membership. (10+10=20M)

```r
setwd("C:\\Users\\personal\\Desktop\\MSIT\\IT\\DataScience_2019501026\\Data Mining\\Exam Solutions\\Final exam")
data = read.csv("Liver_data.csv", header = FALSE, col.names = c("mcv", "alkphos", "sgpt", "sgot", "gammagt", "drinks", "s

x = data[, c("mcv", "alkphos","gammagt", "drinks")]
plot(x,pch = 19, xlim=c(80,100), ylim=c(0,200))

#kmeans with 4 clusters and plotting them
help(kmeans)
model1 = kmeans(x, 4)
plot(x, col=model1$cluster, xlim=c(80,100), ylim=c(0,200))
points(model1$center, col='brown', pch=19, cex=2)

x1 = data[,c("mcv", "gammagt")]
plot(x,pch = 19, xlim=c(80,100), ylim=c(0,200))

#kmeans with k clusters and plotting them
help(kmeans)
model2 = kmeans(x, 4)
plot(x, col=model2$cluster, xlim=c(80,100), ylim=c(0,200))
points(model2$center, col='brown', pch=19, cex=2)
```

7. Compute the misclassification error that would result if you used your clustering rule to classify the data by assigning the majority class of the cluster.
   (10M)

8. Consider the dataset BSE_Sensex_Index.csv. Create an extra column of successive growth rate for column close where the successive growth rate is defined as
   (value of day x- value of day x-1)/value of day x-1. Use a z score cut off of 3 to identify any outliers.  List the respective dates from the csv file on which day these outliers fall.
   (10M)