

Introduction

In this report, the dataset containing daily stock price fluctuations for Tesla will be explored. Tesla Inc. designs and manufactures electric vehicles globally. Tesla has emerged as a leading force within the industry. As a result, their stock prices have reflected this rise. Within this dataset is daily stock data from 2014 - 2020 pertinent to classifying a market trend. This report uses different methods in k-means, random forest, and SVM to analyze the market based on the movement of the previous day.

The data used for this report can be downloaded from the following link:

<https://finance.yahoo.com/quote/TSLA/history?period1=1413158400&period2=1602547200&interval=1d&filter=history&frequency=1d&includeAdjustedClose=true>

Below we describe the different features from the dataset, and additional features derived:

Features description and meaning:

Date: Date of recording

Open: Day's opening price

High: Day's peak price

Low: Day's lowest price

Close: Day's closing price

Adj. Close: Day's closing price including dividends and splits

Volume: Day's trading volume

Daily Return: Days Open - [Day - 1]'s opening price

EarnLossGapTrade: Day's Opening - [Day - 1]'s closing price

threshold <- [Day - 1]'s Opening price * 0.006 = 0.6% threshold

Class <- Up, Down, or Stable indicating the stocks performance compared to previous day

factored as 1 for Up, 2 for Down, and 3 for Stable

IClass <- The previous day's class factored as 1 for Up, 2 for Down, and 3 for Stable

IMovement <- The previous day's daily return

IClassifier <- The previous day's threshold

We compare our Daily Return with our “Threshold” to determine the trend of the market for that given day. That is if our Daily Return is greater than our threshold value for that day, it is the case is classified as “UP”. Conversely, if the daily return is less than our threshold value, it is classified as “DOWN”. Any other instance would fall between -0.6% to 0.6% and is classified as “STABLE”.

Number of Cases: 2588

Number of features: 12

Number of classes: 3

The breakdown of our classes are:

Class 1 (Up): 1094 cases

Class 2 (Down): 976 cases

Class 3 (Stable): 1035 cases

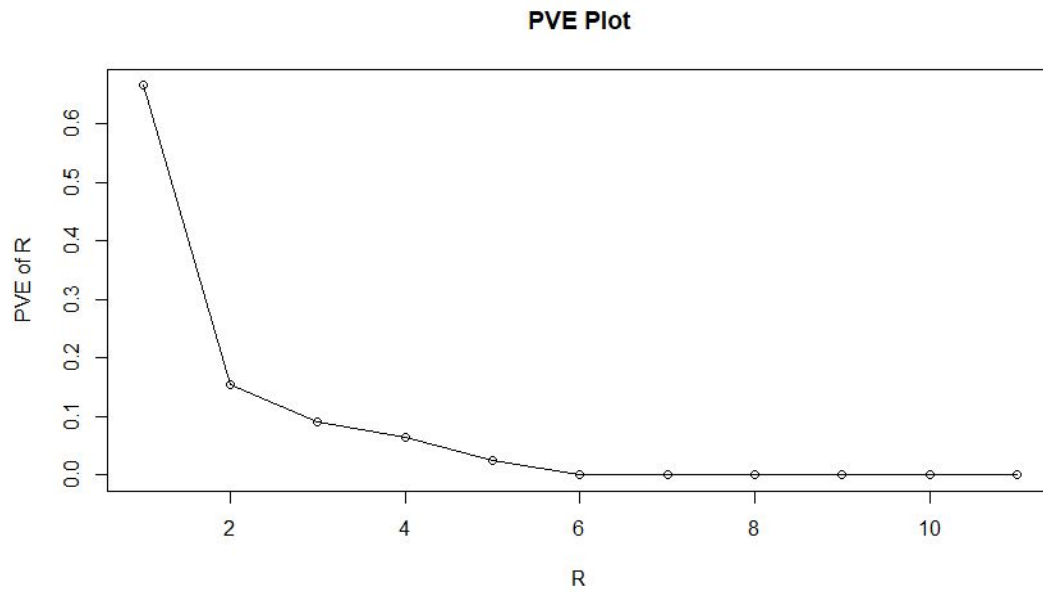
Question 1:

The data is standardized with the use of the `scale()` function and assigned to `SDATA`.

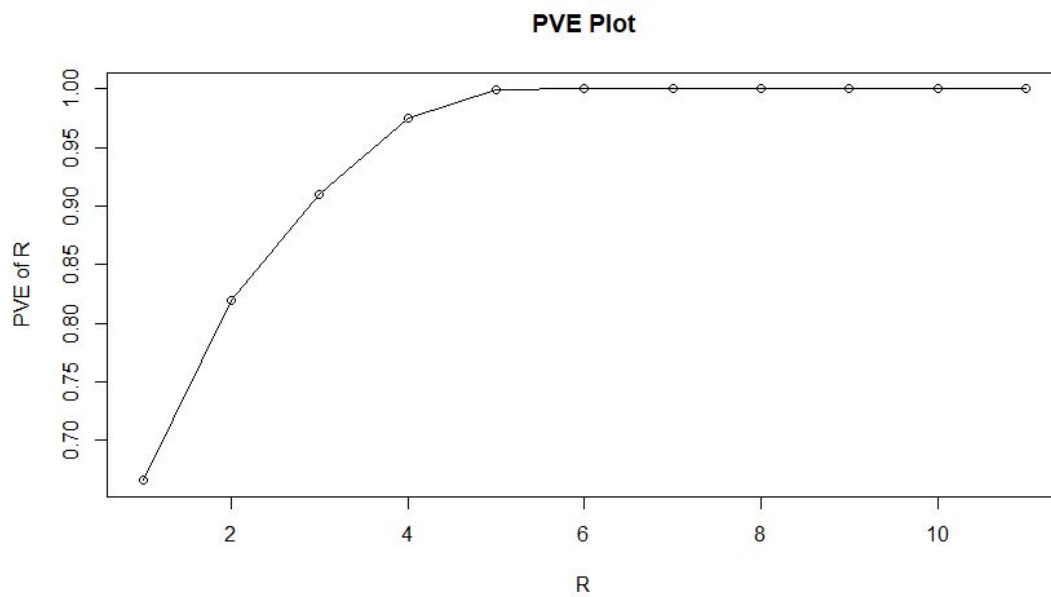
Principal Component Analysis is then used on `SDATA` to obtain our principal components.

The eigenvectors and eigenvalues of a covariance (or correlation) matrix represent the “core” of a PCA: The eigenvectors determine the directions of the new feature space, and the eigenvalues determine their magnitude. In other words, the eigenvalues explain the variance of the data along the new feature axis. Eigenvalues are simply the coefficients attached to eigenvectors, which give the amount of variance carried in each Principal Component. The goal of a PCA is to reduce the dimensionality of feature spaces by projecting it onto smaller subspace, these eigenvectors will form the axis, but they only define the directions of the new axis, to decide which eigenvector(s) can be dropped without losing too much information. In other words, we are attempting to reduce the dimensionality while keeping as much of the variation as possible.

A percent variance is the change in an account during a period of from one period to the next expressed as a ratio. In other words, it shows the increase or decrease in an account over time as a percentage of the total account value. The higher the percentage of variance a proposed model manages to explain, the more valid the model seems to be. Below, we display the percentage of variance explained by each principal component



Our PVE plot expressed the percentage of variance explained for the given R, denoted L_r , in which our highest L_r values are closer to 0. Below, we observe our increasing function of cumulative PVE plot, in which we can clearly observe that the majority of the variance within our dataset occurs within the first 5 or so R values.

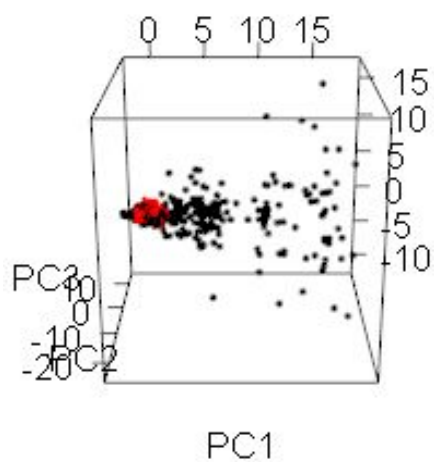


We further explore the eigenvalues associated with each principal component and examine a clear pattern in the percentage of variance explained and the eigenvalue itself.

First 3 eigenvalues: 7.33, 1.69, 1.00

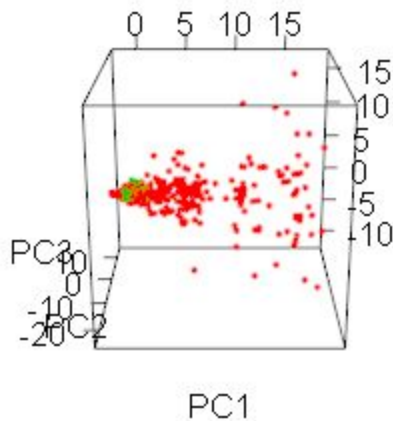
First 3 eigenvalues as a percentage of total variance: 66.58%, 15.35%, 9.08%

We then plot the 3D projection of the first 3 PCA eigenvectors by class below:

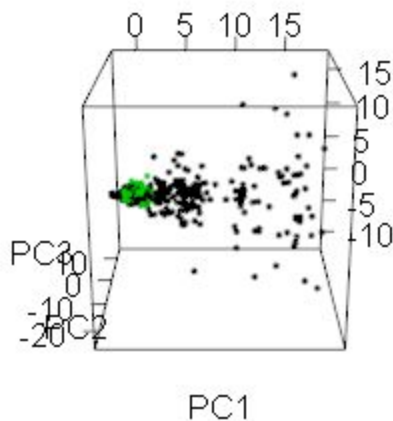


In the above 3D plot, we observe two classes, Up and Down on three PCA eigenvectors with the highest variability. In that sense, the three eigenvectors represent about 90% of the projection.

We can see that there is low separability between the classes on lower values for PC1, as the data points tend to cluster between 0-5 on the PC1 scale.



In the above 3D plot, we observe two classes, Down and Stable on three PCA eigenvectors with the highest variability. In that sense, the three eigenvectors represent about 90% of the projection. Much like our projection for Up and Down, we can see that there is low separability between the classes on lower values for PC1, as the data points tend to cluster between 0-5 on the PC1 scale.



Finally, in the above 3D plot, we observe Down and Stable on three PCA eigenvectors with the highest variability. In that sense, the three eigenvectors represent about 90% of the projection. We can see that there is low separability between the classes on lower values for PC1, as the data points tend to cluster between 0-5 on the PC1 scale.

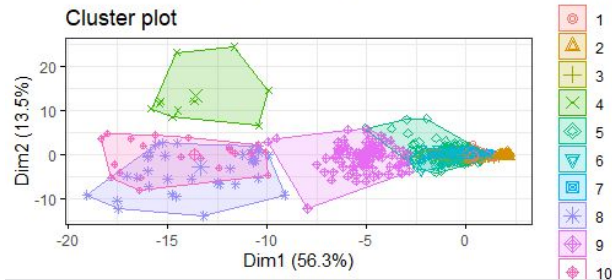
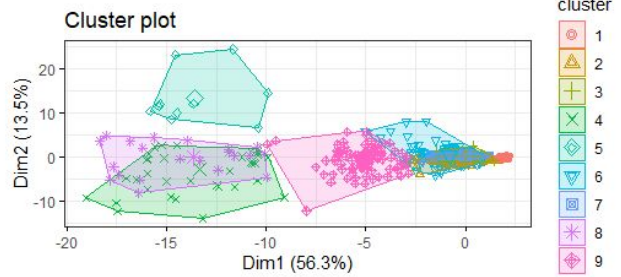
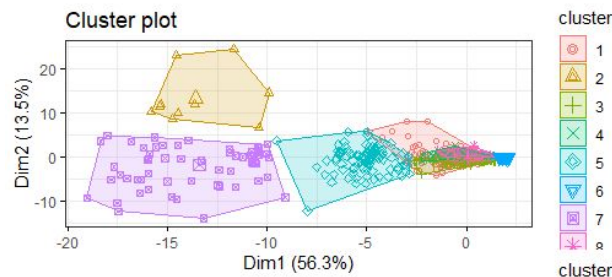
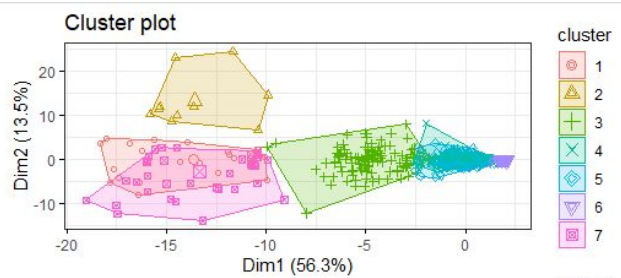
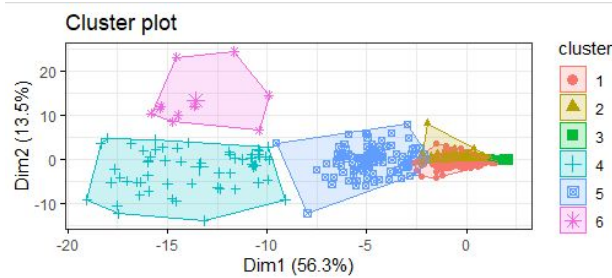
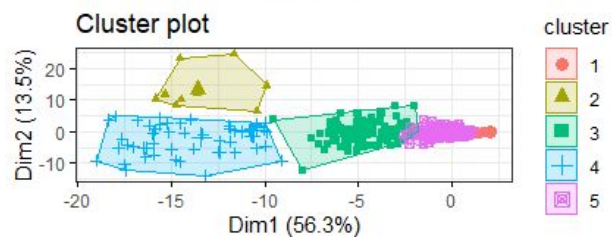
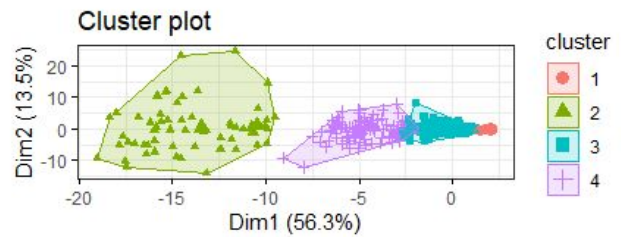
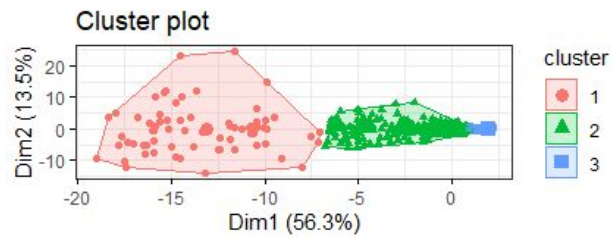
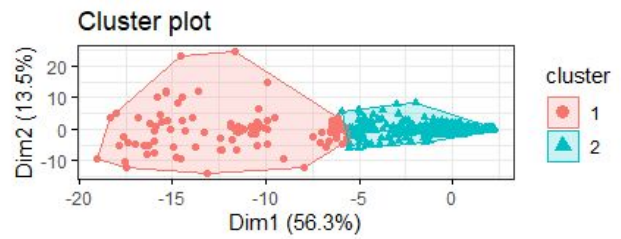
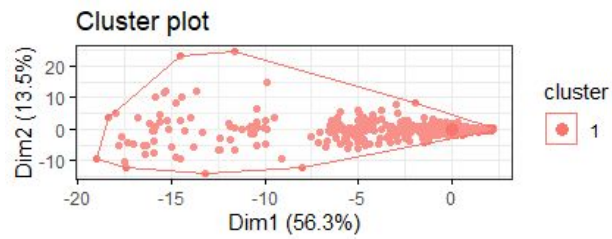
Later in the report, we will see that these PCA projections support the clustering we get from k-means.

Question 2

K-means clustering is a method used for clustering analysis, especially in data mining and statistics. It aims to partition a set of observations into a number of clusters (k), resulting in the partitioning of the data into k clusters. It can be considered a method of finding out which group a certain object really belongs to.

The K-means algorithm divides a set of samples into disjoint clusters, each described by the mean of the samples in the cluster. The means are commonly called the cluster “centroids”. The K-means algorithm aims to choose centroids that minimise the inertia, or within-cluster sum-of-squares criterion

In attempting to find our best number of clusters, we apply k-means clustering on our standardized features from $k=1$ to $k=10$ resulting the following disjointed clusters.



Having one cluster virtual makes no sense because we will not be able to differentiate the different classes.

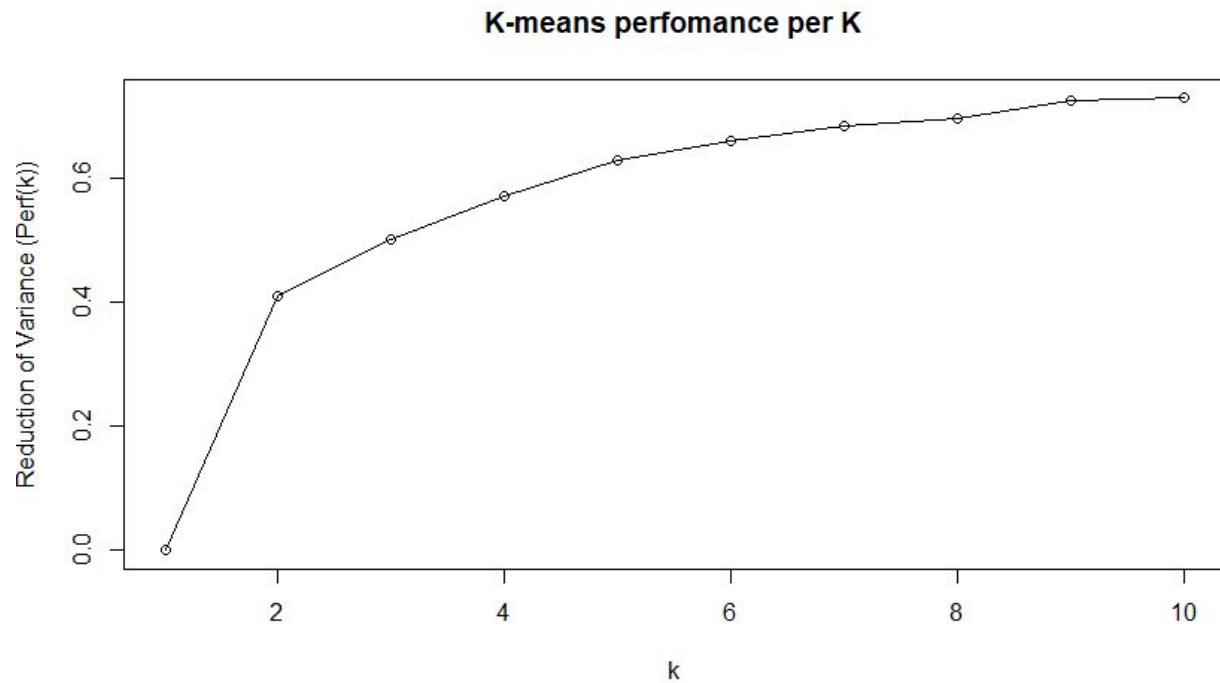
Having two clusters we can see the partitioning of the data set into non-overlapping coordinates in a 2d plot which is good but we want to see the best number of clusters so we keep partitioning.

Clusters = 3 looks good even though there seems to be some overlapping this seems like the best cluster especially knowing that there are 3 different groups but again we keep going.

Clusters = 4 is now starting to show a lot of overlap on the 2-d plot. But in general the more cluster we have the less the dispersion will be so we keep going to find the right balance.

Cluster 5-10 clearly shows a lot of overlapping, but that is expected as we are attempting to plot a 400 dimensional cluster on a 2 dimensional plane.

We further explore the performance of each k number of clusters to help determine the best number of clusters for our model.



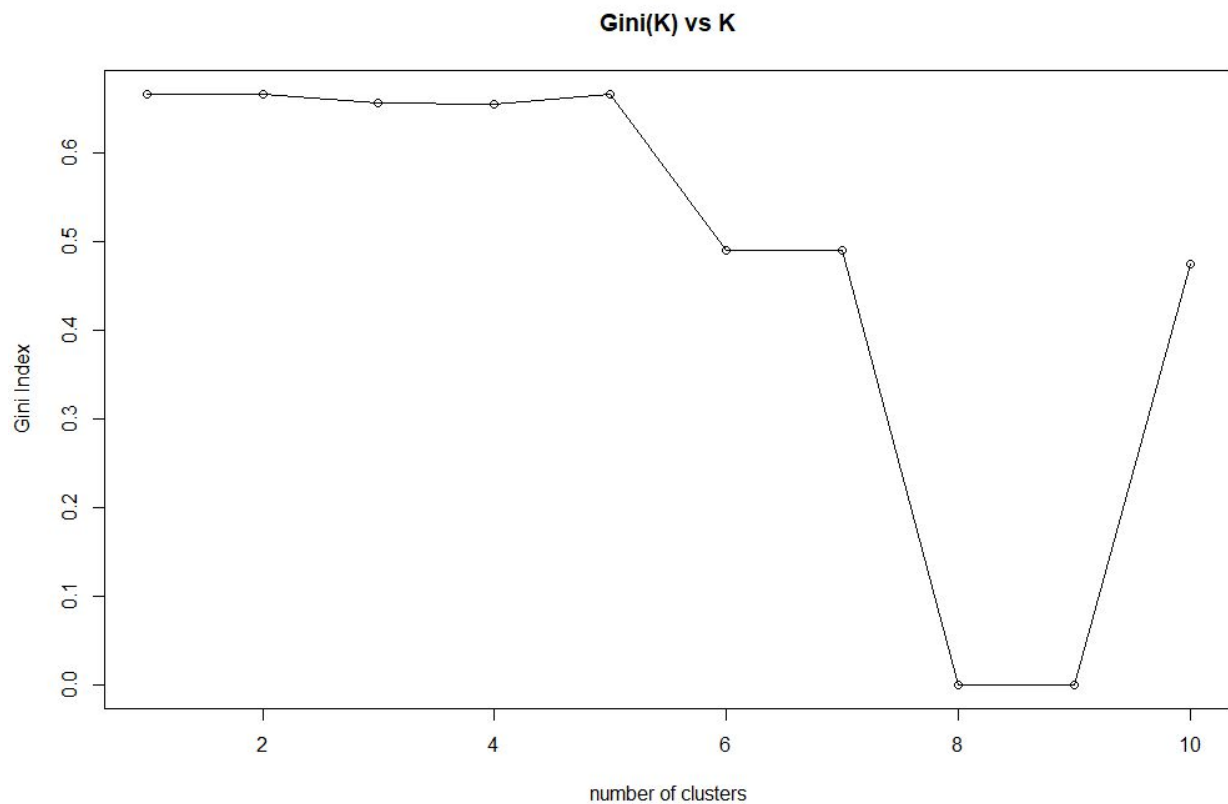
One method to validate the number of clusters is the elbow method. The idea of the elbow method is to run k-means clustering on the dataset for a range of values of k , and for each value of k calculate the sum of squared errors (SSE). We then plot a line chart of the SSE for each value of k . If the line chart looks like an arm, then the "elbow" on the arm is the value of k that is the best. The idea is that we want a small SSE, but that the SSE tends to decrease toward 0 as we increase k . Our goal is to choose a small value of k that still has a low SSE, and the elbow usually represents where we start to have diminishing returns by increasing k . In referencing our plotted k graph against performance, we identify our best k as 6.

To further select our best k value, we explore our gini index for each k number of clusters below.

Gini index or Gini impurity measures the degree or probability of a particular variable being wrongly classified when it is randomly chosen. If all the elements belong to a single class, then it can be called pure. The degree of Gini index varies between 0 and 1, where 0 denotes that all elements belong to a certain class or if there exists only one class, and 1 denotes that the

elements are randomly distributed across various classes. A Gini Index of 0.5 denotes equally distributed elements into some classes.

Gini Impurity is a measurement of the likelihood of an incorrect classification of a new instance of a random variable, if that new instance were randomly classified according to the distribution of class labels from the data set. We aim to construct clusters with smaller impurity.



Clearly, we observe our best gini index that is 0 comes at $k = 6$, so we select 6 as our best k .

Question 3

The K-means algorithm divides a set of samples into disjoint clusters, each described by the mean of the samples in the cluster. The means are commonly called the cluster “centroids”. The K-means algorithm aims to choose centroids that minimise the inertia, or within-cluster sum-of-squares criterion

A centroid is a data point at the center of a cluster. The resulting classifier is used to classify the data and thereby produce an initial randomized set of clusters. Each centroid is thereafter set to the arithmetic mean of the cluster it defines. In the section below, we explore the prominent features of k-means clustering: centers, size, dispersion, gini impurity, and frequencies.

```
> km6$centers
  class iclass      open      high      low      close  Adj.Close      volume  DailyReturn  EarnLossGapTrade
1 1.152461 1.822329 -0.01012226 -0.01324337 -0.008577549 -0.01257705 -0.01257705  0.49846775  0.16766657  0.04026537
2 2.575859 1.763621  0.01609792  0.01139315  0.021672508  0.01621069  0.01621069 -0.07563972 -0.15773976 -0.08470241
3 1.995418 1.790378 -0.72142148 -0.71428057 -0.728900891 -0.71938440 -0.71938440 -0.82091729 -0.03300873 -0.03995241
4 1.916667 1.805556  4.99327655  5.00283831  4.973159010  5.00893134  5.00893134  1.61613718  1.86456645  2.39520704
5 1.758065 1.637097  1.66189254  1.66695432  1.661833087  1.65824171  1.65824171  2.06316843  0.21266194  0.20678663
6 2.000000 1.500000  5.34373826  5.40999430  5.348552288  5.28948822  5.28948822  1.94805844 -8.32849721 -9.57090718
  Classifier  IMOVEMENT  IClassifier
1 -0.02197853 -0.055662219 -0.01832313
2  0.02745355  0.007464967  0.02682416
3 -0.72390288 -0.035042297 -0.72479229
4  4.89581979  0.737665828  4.86964350
5  1.65739711  0.156604441  1.65509658
6  5.95069815 -0.447804179  6.01593102
```

Size

Cluster	1	2	3	4	5	6
Size	833	1193	873	72	124	10

We observe the size of each cluster, and identify our biggest cluster as 4. We then observe the biggest cluster to distinguish if there is a distinct class that clearly holds a higher proportion in that cluster.

Dispersion

Cluster	1	2	3	4	5	6
Dispersion	2250.42	1888.75	1158.83	4599.76	2397.36	675.57

Dispersion reflects the compactness of a cluster when employed at the intra-cluster level and reveals the separation when measured at the inter-cluster level. So it results in a vector with a number for each cluster. One expects this ratio to be as low as possible for each cluster, since we would like to have homogeneity within the clusters. We can observe that cluster 3, the most balanced class has a low dispersion score.

Gini

Cluster	1	2	3	4	5	6
Gini	.26	.49	.67	.63	.64	0

We obtained the above gini values for each cluster and observed that most of the clusters have decent gini values with the exception of cluster 1 and cluster 6. Upon further investigation, we recognize the classes with lower gini scores are more unbalanced. This notion is supported by the frequency table in the next section.

Impurity is the sum of the gini index for every cluster. We obtain an impurity score of 2.68.

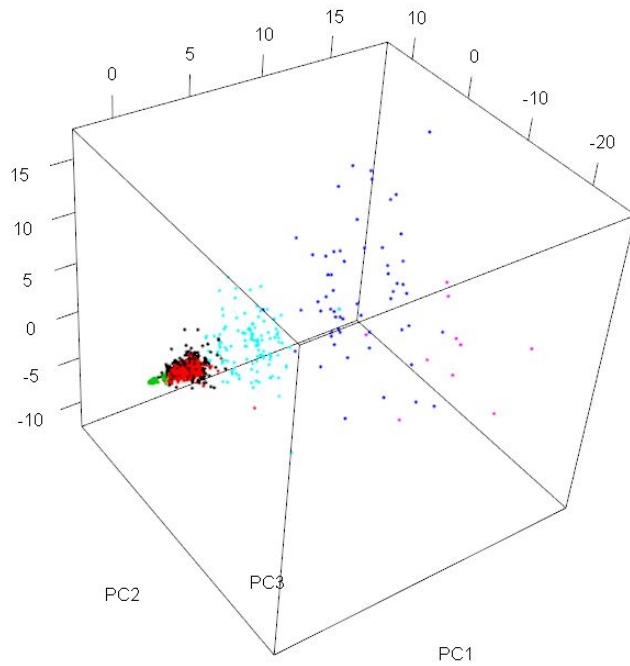
Frequency

Class/Cluster	Up	Down	Stable
1	708	123	2
2	0	506	687
3	300	277	296
4	32	14	26
5	54	46	24
6	0	10	0

On our frequency table, we observe what we expect in relation to our gini indexes for each cluster. That is, Cluster 1 is extremely unbalanced so we expect a low gini score, followed by

cluster 2, with cluster 3,4,5 being relatively more balanced so we observe high gini scores.

Finally cluster 6 is not balanced whatsoever, so we observe a gini score of 0.



Lastly, we display the 3D projection of our kmeans for 6 clusters, on the three most prominent PCA eigenvectors. That is, we can observe a large number of green, red, and black data points; which is reflected by our size chart.

Question 4

In preparation for further classification methods, we obtain training and test sets for each class within our dataset.

The sample split used for deriving the training and test set is 80/20% split with 80% of the random sampling being for training purposes and 20% of the random sampling being for testing purposes. After obtaining the training and test sets for each class we observe the dimensions of each:

Class Up Train size: 876 rows

Class Down Train size: 781 rows

Class Stable Train size: 828 rows

Class Up Test Size: 218 rows

Class Down Test Size: 195 rows

Class Stable Test Size: 207 rows

Furthermore, each class is relatively balanced, no cloning needs to be done on the training and test sets.

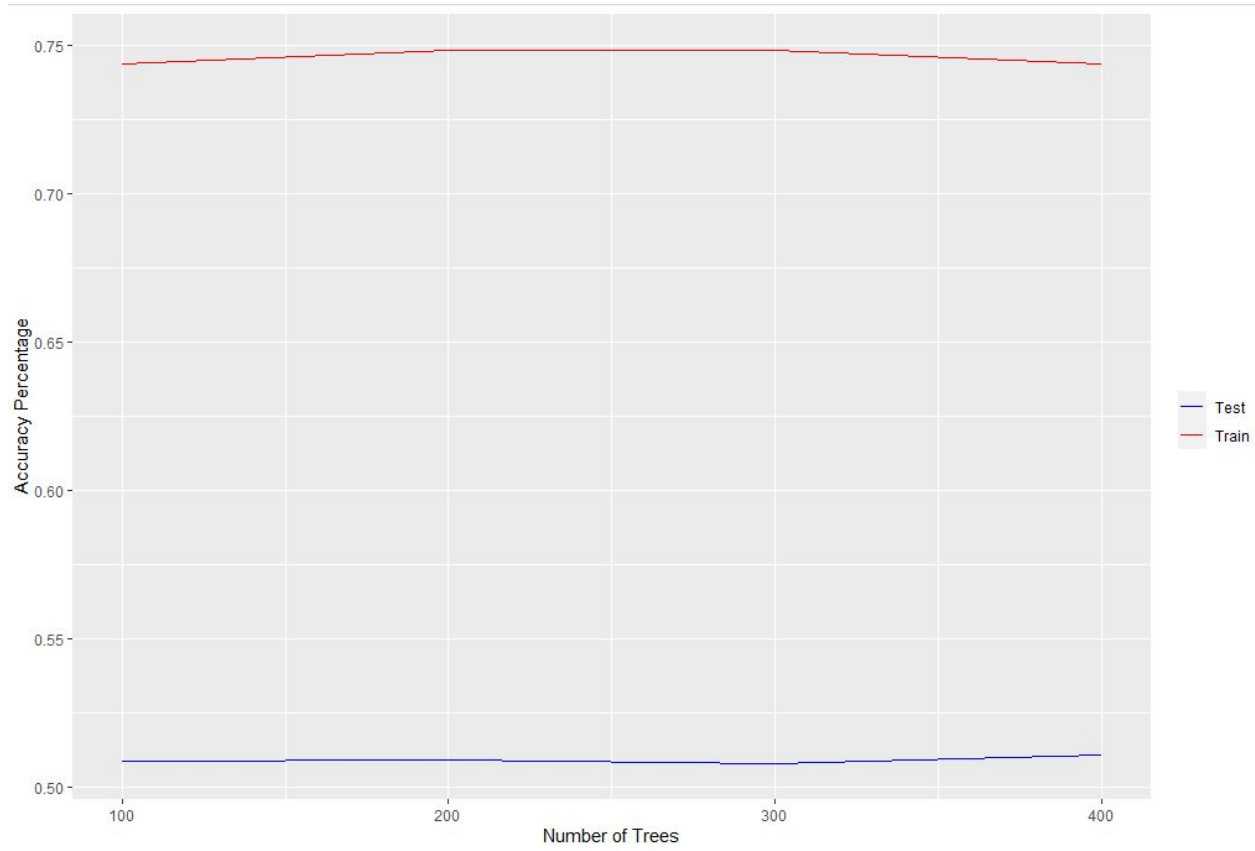
Question 5

Random Forest is a robust machine learning algorithm that can be used for a variety of tasks including regression and classification. It is an ensemble method, meaning that a random forest model is made up of a large number of small decision trees, called estimators, which each produce their own predictions. The random forest model combines the predictions of the estimators to produce a more accurate prediction. It works by aggregating the predictions made by multiple decision trees of varying depth. Every decision tree in the forest is trained on a subset of the dataset; this is called the bootstrapped dataset. The portion of samples that were left out of each decision tree in the forest is referred to as the Out-Of-Bag (OOB) dataset. As we'll see later, the model will automatically evaluate its own performance by running each of the samples in the OOB dataset through the forest.

Since the practical classification task is to determine the movement of "Class" (i.e. Today's stock movement) based on today's stock price, and yesterday's movement, we omit the following features of the formula when building the model: Class, DailyReturn, and Classifier(threshold). That is, because including those features in the model would result in 100% accuracy in predicting class, since the model will know what class is already. By excluding that, we are

effectively trying to predict class, with only the knowledge of today's stock prices, and yesterday's stock's movement (IClass)

Below, we plot the performance of the training and test set of the random forest algorithm for the following number of trees [100,200,300,400].



The random forest performance accuracy lied generally around 75% for the training set, and around 52-53% for the test set. We further explore in detail the performance of each ntrees model by class with confusion matrices.

NTREES=100	Predicted:Up	Predicted:Down	Predicted:Stable
True:Up	55%	19%	26%
True:Down	25%	46%	29%
True:Stable	28%	22%	50%

N TREES=200	Predicted:Up	Predicted:Down	Predicted:Stable
True:Up	55%	19%	26%
True:Down	25%	47%	29%
True:Stable	28%	22%	50%

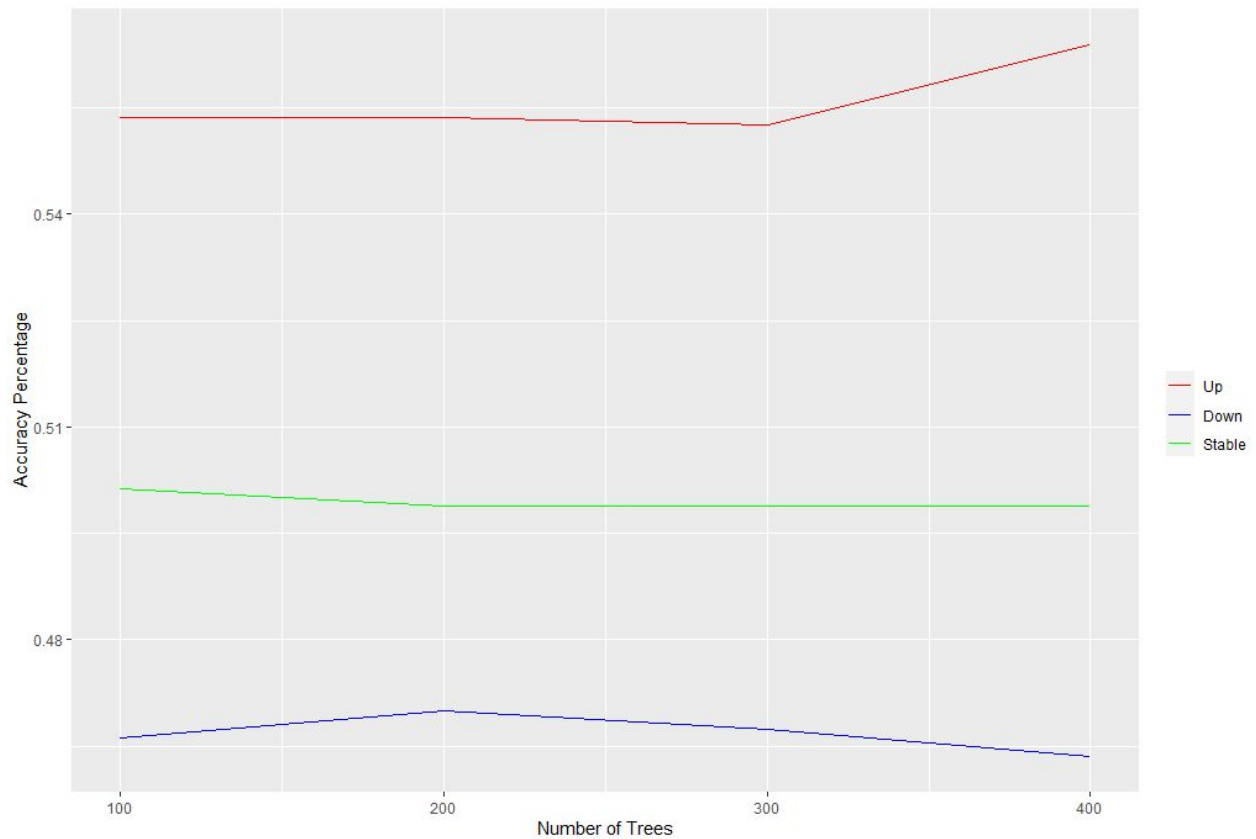
N TREES=300	Predicted:Up	Predicted:Down	Predicted:Stable
True:Up	55%	19%	26%
True:Down	25%	47%	29%
True:Stable	27%	23%	50%

N TREES=400	Predicted:Up	Predicted:Down	Predicted:Stable
True:Up	56%	18%	25%
True:Down	25%	46%	28%
True:Stable	27%	23%	50%

Based on our confusion matrices, we can see that regardless of ntrees, the accuracy of each class prediction remains around the same value, as indicated by the number of trees vs accuracy performance plot above. That is, Up is generally predicted correctly ~ 55% of the time, Down is predicted correctly ~ 46% of the time and Stable is predicted correctly around 50% of the time.

Question 6

In this section, we observe the performance by class for each number of ntrees in 100, 200, 300, and 400.



The random forest performance accuracy is generally between 45% to 56% for ntrees = 100, 200, 300, and 400. When we explore three curves, with each curve representing the diagonal coefficients for each class for n trees = 100, 200, 300, 400, we observe that our best random forest model occurs when ntrees = 400 with a very miniscule advantage.

Question 7

There are two measures of importance given for each variable in the random forest. The first measure is based on how much the accuracy decreases when the variable is excluded. This is further broken down by outcome class. The second measure is based on the decrease of Gini impurity when a variable is chosen to split a node. Each tree has its own out-of-bag sample of data that was not used during construction. This sample is used to calculate the importance of a specific variable. First, the prediction accuracy on the out-of-bag sample is measured. Then, the values of the variable in the out-of-bag-sample are randomly shuffled, keeping all other variables the same. Finally, the decrease in prediction accuracy on the shuffled data is measured.

The mean decrease in accuracy across all trees is reported. This importance measure is also broken down by outcome class. When a tree is built, the decision about which variable to split at each node uses a calculation of the Gini impurity. For each variable, the sum of the Gini decreases across every tree of the forest and is accumulated every time that variable is chosen to split a node. The sum is divided by the number of trees in the forest to give an average. The scale is irrelevant: only the relative values matter. The importance of each feature is shown below:

Feature	MeanDecreaseGini (Importance)
IClass	35.8
Open	148.2
High	145.6
Low	145.6
Close	131.5
Adj.Close	133.8
Volume	201.8

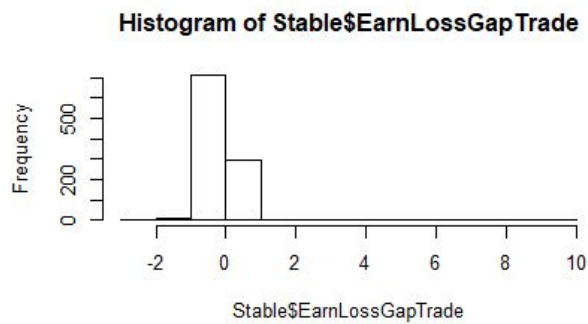
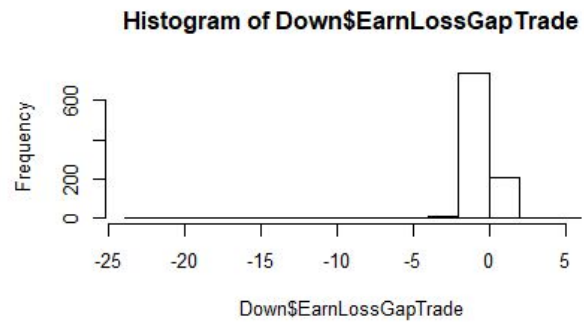
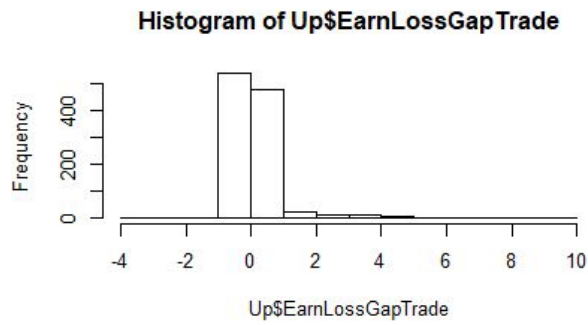
EarnLossGapTrade	305.0
IMovement	182.3
IClassifier	224.5

For better visualization, the features are sorted by importance and displayed below:

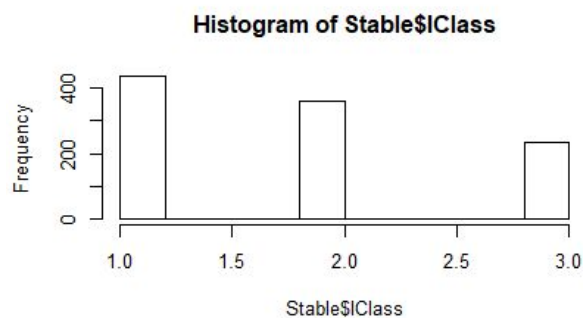
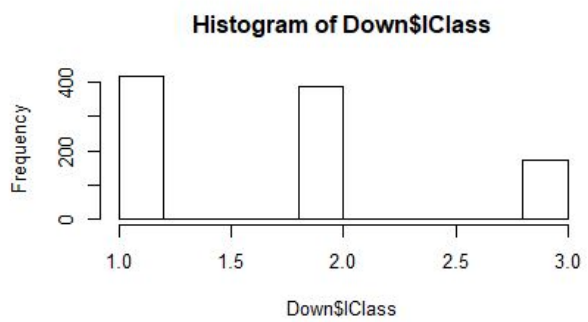
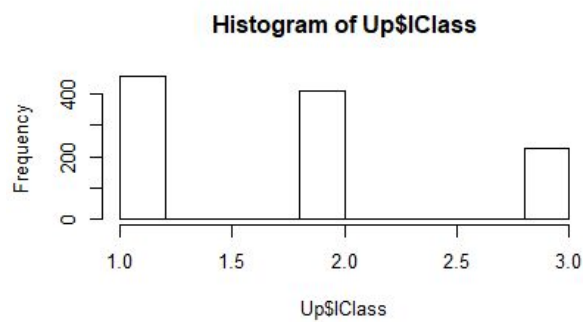
Sorted Feature	MeanDecreaseGini (Importance)
EarnLossGapTrade	305.0
IClassifier	224.5
Volume	201.8
IMovement	182.3
Open	148.2
High	145.6
Low	145.6
Adj.Close	133.8
Close	131.5
IClass	35.8

Question 8

After identifying the highest and lowest features from the previous question, we can assume that Earn Loss Gap trade has the highest importance and IClass is the lowest. To understand how the highest and lowest features differ, a Histogram and the K-S test will be performed. The histogram on each feature will display how distinct each class is in that feature and the K- S test is a numerical confirmation. The D- value of the K-S test represents the distance of each class pair. The larger the D- value, the more distinct the class pair.



As we can see, the three histograms differ in trend. We can assume that each class is distinct in this feature. This represents that Earn Loss Gap Trade is a good feature to predict classes.



Here we can see that IClass has the same trend. We can assume that IClass features will have a difficult time in differentiating classes. It is possible that extracting this feature may improve accuracy in the model.

Comparatively, we can conclude that EarnLossGapTrade is most prominent in predicting Today's stock movement, as the formula to derive that is Today's opening price - Yesterday's closing price, which one can imagine would be a good predictor compared to Yesterday's stock movement which is denoted as T -1 Open (Yesterday's Opening price) - [T-2] Open (The day before's opening price).

```
> ks.test(Up$EarnLossGapTrade, Down$EarnLossGapTrade)

Two-sample Kolmogorov-Smirnov test

data: Up$EarnLossGapTrade and Down$EarnLossGapTrade
D = 0.33129, p-value < 2.2e-16
alternative hypothesis: two-sided
> ks.test(Up$EarnLossGapTrade, Stable$EarnLossGapTrade)

Two-sample Kolmogorov-Smirnov test

data: Up$EarnLossGapTrade and Stable$EarnLossGapTrade
D = 0.22673, p-value < 2.2e-16
alternative hypothesis: two-sided
> ks.test(Down$EarnLossGapTrade, Stable$EarnLossGapTrade)

Two-sample Kolmogorov-Smirnov test

data: Down$EarnLossGapTrade and Stable$EarnLossGapTrade
D = 0.21557, p-value < 2.2e-16
alternative hypothesis: two-sided
> ks.test(Up$Iclass, Down$Iclass)

Two-sample Kolmogorov-Smirnov test

data: Up$Iclass and Down$Iclass
D = 0.031155, p-value = 0.6986
alternative hypothesis: two-sided
```

```

> ks.test(Up$Iclass,Stable$Iclass)

      Two-sample Kolmogorov-Smirnov test

data:  Up$Iclass and Stable$Iclass
D = 0.018644, p-value = 0.9926
alternative hypothesis: two-sided

> ks.test(Down$Iclass, Stable$Iclass)

      Two-sample Kolmogorov-Smirnov test

data:  Down$Iclass and Stable$Iclass
D = 0.049799, p-value = 0.1655
alternative hypothesis: two-sided

```

KS test is used to test the null hypothesis that two datasets come from the same distribution. The p-value is <0.01 across the three pairs of classes for the feature EarnLossGapTrade, thus we conclude that the three datasets come from different distributions. The $p\text{-value} > 0.01$ across the three pairs of classes for the feature IClass, thus we conclude that the three datasets come from the same distribution. We can also see that Earn Loss Gap trade has a larger distance in each class pair compared to the IClass pairs.

Question 9

In this section we attempt to train a new random forest classifier dedicated only to cases belonging in the cluster with the lowest gini score. Initially, cluster 6 was the most prime candidate with a gini index of 0, however it seems faulty as it has only 10 cases, and all 10 of those cases exist within a single class. We then move to the next legitimate gini index cluster, cluster 1.

Since this cluster is extremely imbalanced as expected, we conduct SMOTE cloning technique to rebalance the cases for good training and test sets. This function handles unbalanced classification problems using the SMOTE (Synthetic Minority Over-sampling Technique) method. Namely, it can generate a new "SMOTEd" data set that addresses the class imbalance

problem. Now that we obtained balanced classes, we ran the best random forest function on our dataset with the balanced classes.

With the cloned balanced data, we split the training and test set in a random 80% and 20% split for our training and test set, and run same Random Forest conditions to obtain the following results:

Question 10

The RF classifier is 91% accurate in predicting class Up, with a false positive rate of 0%, and false negative rate of 9%. It is 100% accurate in predicting class Down and Stable.

New RF on Gini cluster

	Predicted:Up	Predicted:Down	Predicted:Stable
True:Up	91%	9%	0%
True:Down	0%	100%	0%
True:Stable	0%	0%	100%

Confusion matrix from section 5

N TREES=400	Predicted:Up	Predicted:Down	Predicted:Stable
True:Up	56%	18%	25%
True:Down	25%	46%	28%
True:Stable	27%	23%	50%

From the previous results in question five we can see overall the accuracy of the model in New RF is increased. We can assume that targeting the least impure cluster does influence the model by increasing accuracy. Because the Gini differ in a large range we can not assume that there is

an advantage to training each cluster separately. However, if the Gini values were similar there could be an advantage to training and testing each cluster separately.

Question 11

The support vector machine (SVM) is a supervised machine learning model that uses classification algorithms for two-group classification problems. It uses a decision boundary to classify two classes by distance. The SVM algorithm has many pro's such as less computation power. However, it is very sensitive to outliers. This could cause the accuracy to be low.

The two classes chosen for this section is Up, and Down in which we use new training sets of those classes to train a linear support vector machine model to classify whether a case is "Up" or "Down". The svm function in R is used to run this algorithm, with several cost values tested, namely {0.01,0.1,1,5}. Ultimately, cost = 0.1 is used to train the svm model.

Train performance

	Predicted:Up	Predicted:Down
True:Up	95%	5%
True:Down	53%	47%

Test performance

	Predicted:Up	Predicted:Down
True:Up	90%	10%
True:Down	60%	40%

Predicting power is lost in both classes when we move from the train set to the test set. Based on the confusion matrix, we can see that the SVM algorithm performs well in predicting Class Up with an accuracy of 90%, but not as well in predicting class Down. This result indicates that Class Up and Down are distinct. A way to possibly improve these results could be by testing a

polynomial or radial kernel. However, this is a trial- and - error process and finding the correct parameters is important.