

Predicting IMDB Movie Ratings Using Classification Techniques

Group ID - G5

Introduction

The IMDB database contains a list of movies released up to date. Based on votes from users, each movie is rated on a scale of 10 for its overall performance and likeability. The success of a movie is based on various factors, such as the cast of the movie, director, budget, length of the movie, year of release, genre etc.

In this project we hope to identify the impact of various factors on the rating of a movie and make a prediction on the possible rating of a new movie when released. The ideal model will be able to take any of the above factors or combination of factors as input and estimate the likely rating of a movie.

For the purpose of analysis, the rating of the movies has been classified into three categories. Movies that fall below 6.5 in the IMDB rating have been classified as 3, movies in the range of 6.5 to 7.3 have been classified as 2 and movies above 7.3 in the rating are classified as 1. This classification was performed keeping in mind the quantiles of the ratings. The three categories are, movies below the median, between median and third quantile and above third quantile respectively.

As this is a classification problem, some of the methods that we chose for analysis are Multinomial Logistic Regression, Tree based methods and Support Vector Machines. Using, each of these methods, we were able to identify important predictors that impacted the rating. The results from all the three models were then compared enabling us to draw conclusions about the dataset and the methodologies.

Methods

Dataset Preparation

The dataset used for the analysis is taken from the ggplot2 library of R. It contains data about movies with information such as title, year of release, length, budget, rating, number of votes, MPAA rating and genre. Information about actors and directors for the above movies was extracted from the IMDB database available online.

Actors and Directors were then compared with people's choice of popular actors and directors (Ranker Website). This is a user rating based on number of awards and accolades won by actors and directors and their overall presence on the screen. Based on this list, actors and directors were classified into three categories, with '1' being the most popular, '2' being moderately popular and '3' being unpopular or new to the field. In the same manner as the ratings, these categories were made based on quantiles.

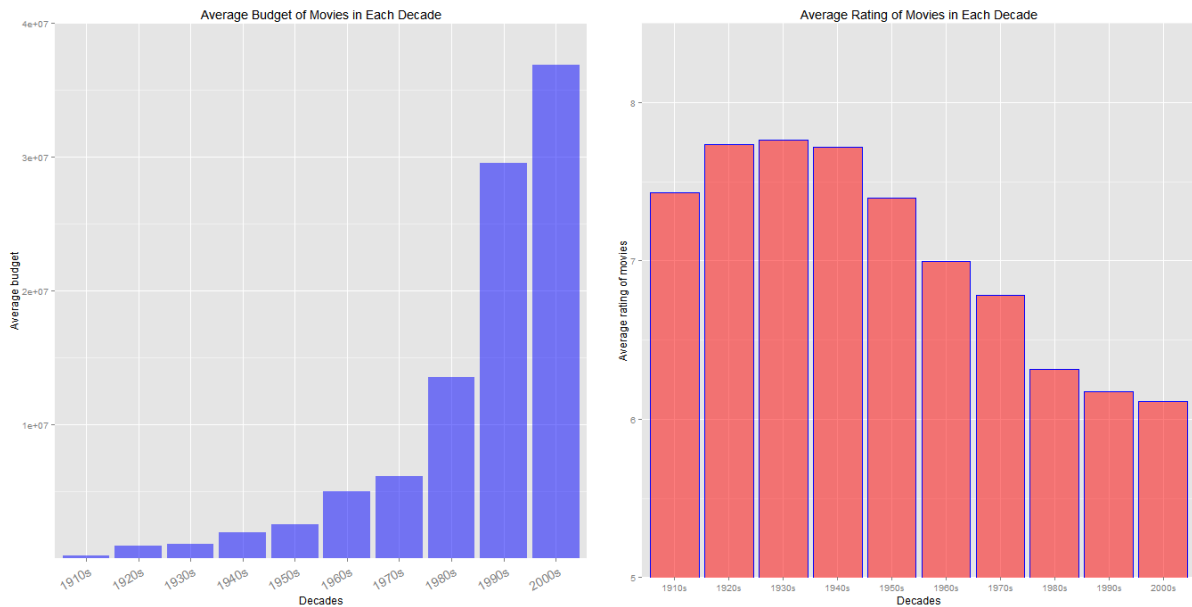
For the purpose of consistency, only full length movies were considered for the analysis. All short movies were removed. Also, all rows with missing values of budget were removed. Thus the final dataset consisted of approximately 2500 rows and 13 predictors.

Exploratory Data Analysis

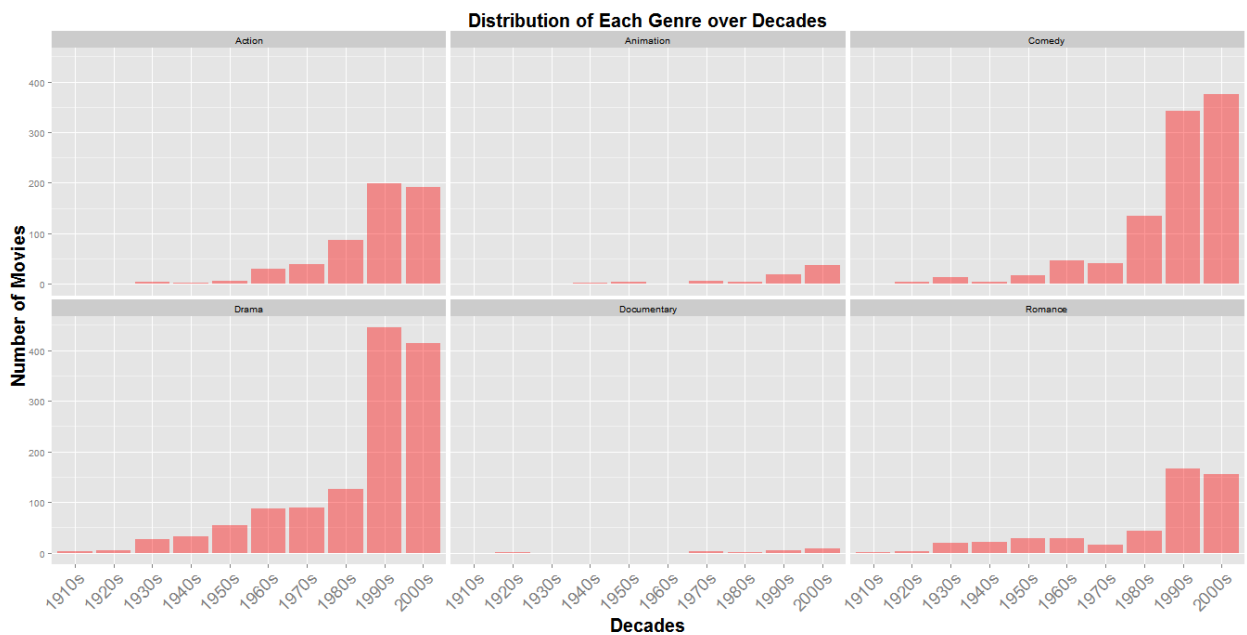
The first step is to understand the data and to see how certain predictors are related with each other and determine if there are any visual patterns between the predictors and the response variable. Below are various plots illustrating the same.

In the below figure, the left histogram shows how the average budget of movies changed over the years. It can be seen that there is a sudden increase in budget from 1980s to 1990s. The right histogram shows how the average ratings of movies changed over time. It shows that movies released in the first half of the century have better average ratings than the movies released in the latter half of the century.

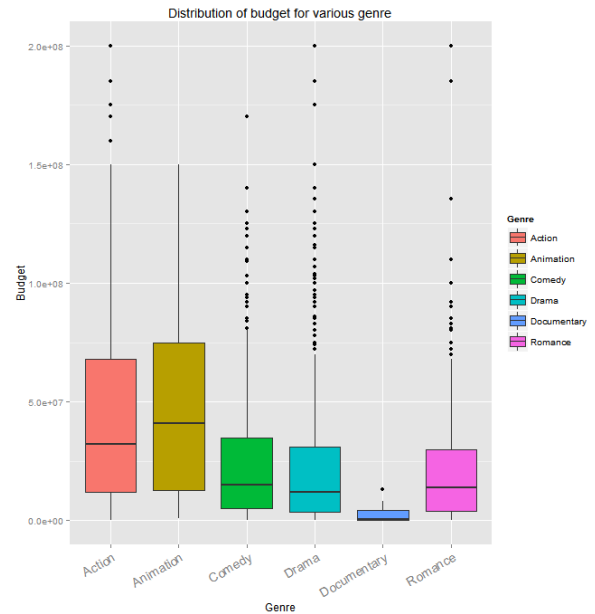
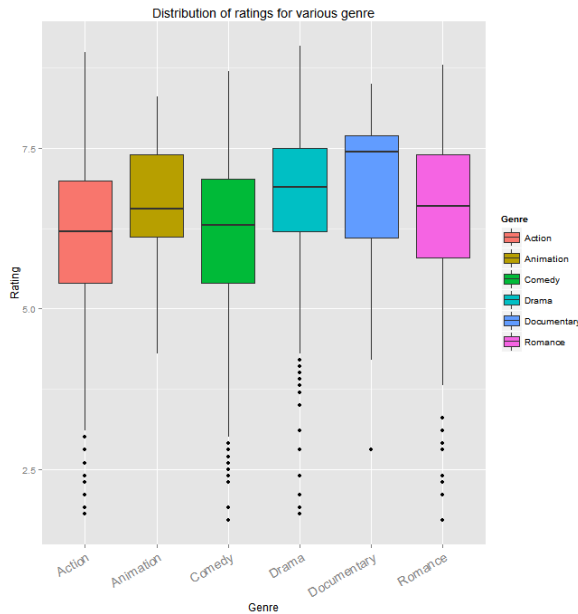
The increase in budget could be attributed to the rising use of technology in film making. Also, the ratings could be attributed to the fact that, people in general are more biased toward older movies and are likely to give it a higher rating. This clearly shows that budget and year of release of the movies are important variables for analysis.



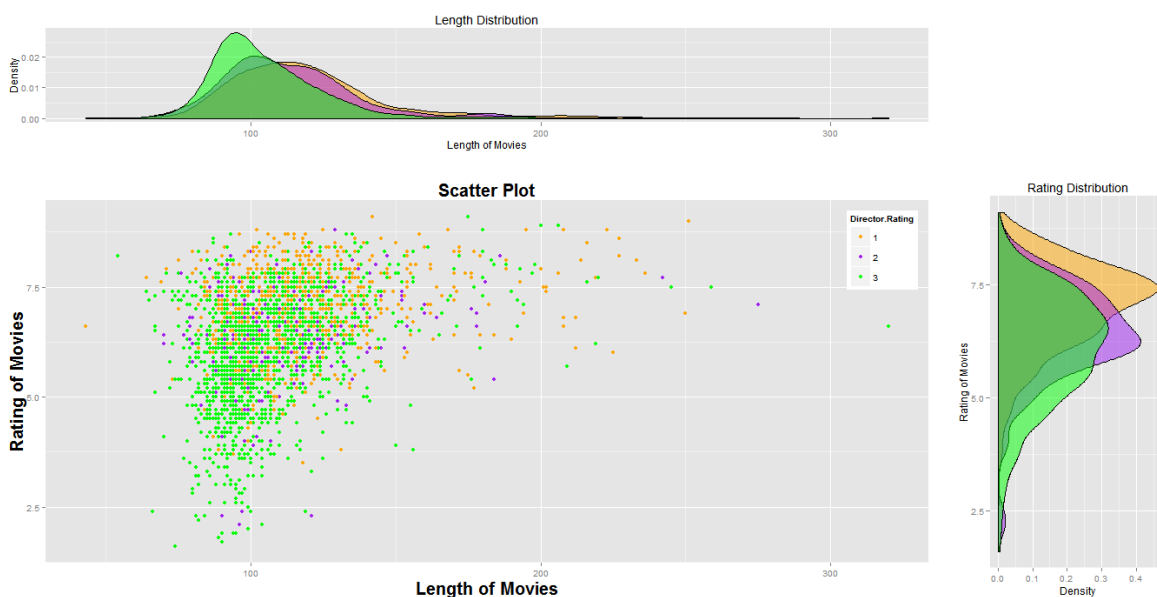
Below is a genre based classification of the total number of movies released in every decade. It can be seen that the highest number of movies are usually made in the drama and comedy genre with the least being in animation and documentary. There is also the possibility that certain movies could belong to more than one genre. For instance, most comedy movies are likely to fall within the drama genre as well. Hence, it is worthwhile to analyze the interactions between genres while building models.



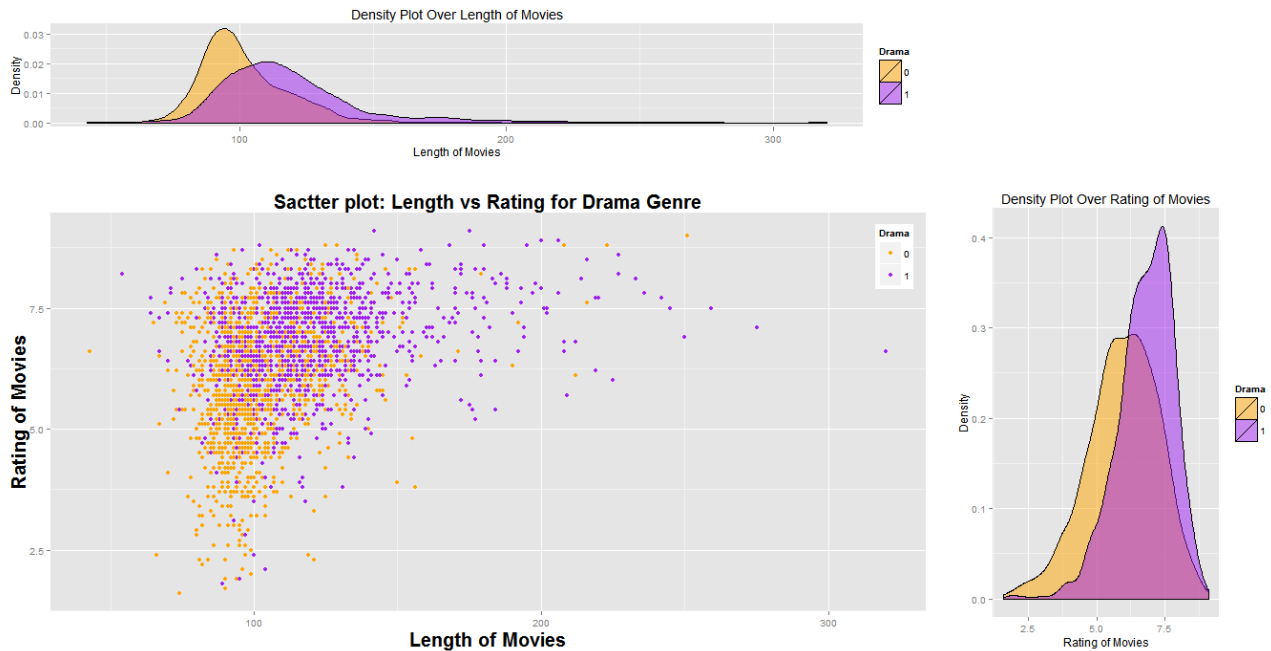
The box plot below is a representation of genre vs rating and genre vs budget. It can be seen that, drama genre tends to have the maximum outliers while documentary having the least. This could be attributed to the number of films made in the respective genres. Drama tends to have the maximum number of movies made in a decade whereas documentary tends to have the least.



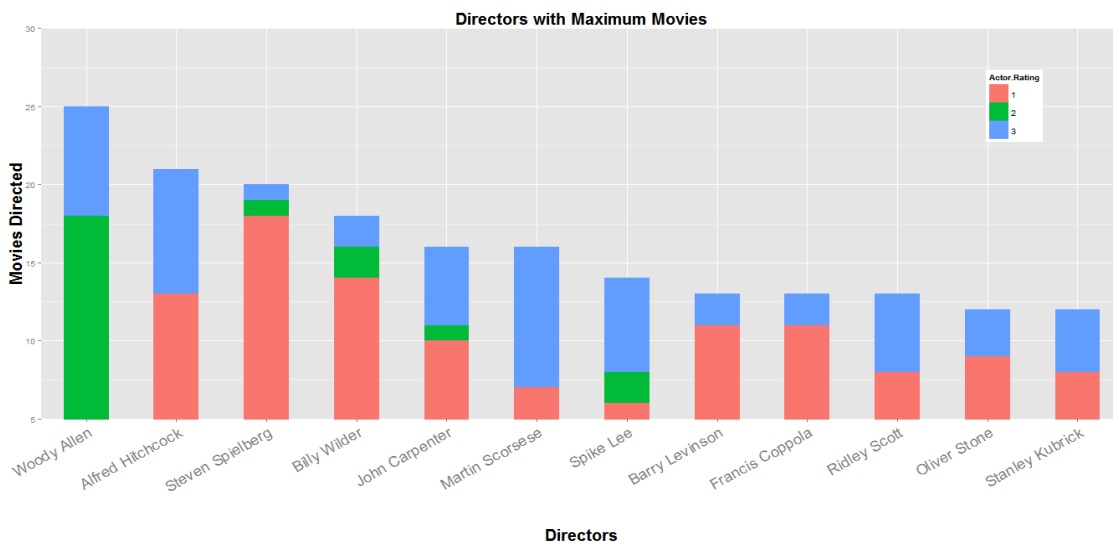
The scatter and density plots depicted below indicate the relation between length of the movies and the ratings. The scatter plot has been divided on the basis of Director Rating. The excellent directors are showcased in orange whereas the good and average directors are showcased in purple and green respectively. It can be observed that, as the length of the movie increases, its rating is likely to be higher. It is also evident that, director rating impacts the rating of the movie directly. A large fraction of excellent directors are in the higher side of the ratings. From here it can be inferred that, length of the movie and the director rating are important predictors directly influencing the ratings.



The scatter and density plot below depicts the relation between the length of movies and their rating. Purple indicates drama genre and orange indicates non – drama movies. It can be observed that, longer movies are likely to be in drama genre and tend to have higher ratings. This reiterates the fact that long drama movies tend to do better on an average.



Below is a list of 12 directors in the decreasing order of the number of movies made by them. The three colors represent the type of actors, these directors have worked with. It can be said that, a vast majority of these directors prefer to work with excellent and good actors as opposed to average or new actors. However, from the graph of actors and the type of directors they worked with [1], the same is not true. Implies, actors are more liberal in choosing the directors they work with as opposed to the directors choosing the actors. This indicates that, director rating can prove to be more important than actor rating during prediction.



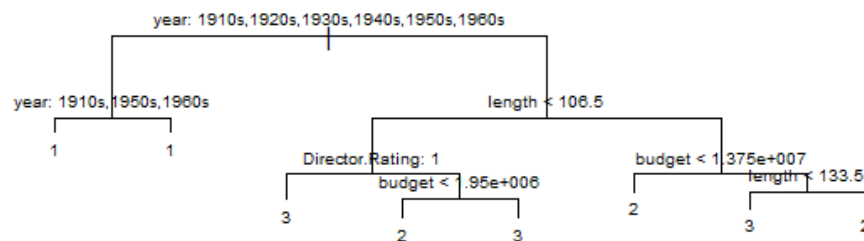
Remark: From the Exploratory analysis it can be said that some of the important predictors that can affect the response are year, length of the movie, director rating and budget.

Tree Based Methods

As this is a classification problem, one of the approaches that was considered was Tree Based Methods. Due to the presence of many predictors and many possible interactions between them, trees were the first choice. Trees tend to perform variable screening or feature selection automatically during the process of classification. Secondly, non-linear relationships between parameters do not affect trees significantly. Lastly, it is easier to interpret the results of decision trees.

Analysis with single trees

From the below plot it can be seen that only certain predictors are considered for the model. The main predictors the tree has chosen for classification are Year of release, Length of the movie, Director rating and Budget.



From the above plot, the following inferences can be drawn.

It can be said that, all movies during and before the 1960's are classified as excellent movies with rating of 1. This is probably due to the following two reasons

1. During this period, the exposure to technology was minimal and hence people's expectations of movies were very less as well. Also, it is a possibility that user ratings for old movies is likely to be higher due to a general bias towards old items.
2. Secondly, it is possible that this particular dataset could be biased towards old movies, having chosen only the good ones from the past.

It can be observed that, movies after the 1960's are classified based on length and budget primarily. If the length of the movie is between 100 and 133 minutes and the budget is less than 13 Million dollars, it is classified as a good movie with a rating of 2.

Similarly, if the length of the movie is between 100 and 133 minutes and the budget is greater than 13 Million dollars, it is classified as an average movie with a rating of 3. However, all movies with a budget greater than 13 Million and length greater than 133 minutes are classified as good movies with a rating of 2.

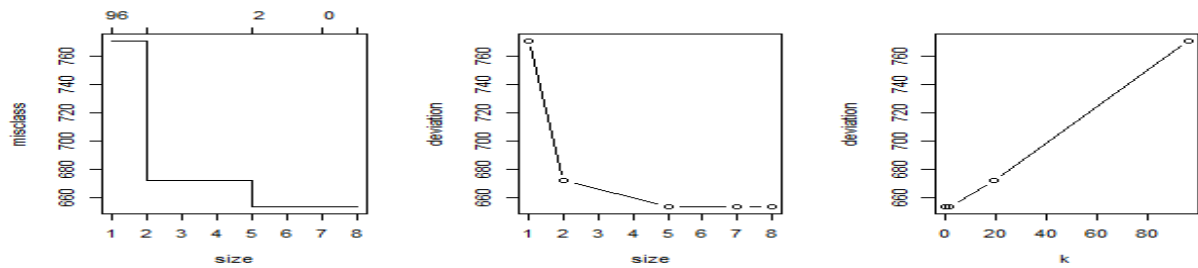
Remark: This is likely due to the fact that, movies within a certain budget and certain time limit, tend to explore ideas fully, whereas short movies with a very high budget tend to leave the audience with unexplored ideas or concepts thus reducing the movie experience. This is corroborated by the fact that movies greater than a certain length and certain budget are classified as good movies.

Accuracy and Error Rate in Model

	Training Dataset	Testing Dataset
Accuracy	59.13	54.86

Remark: It can be seen that, the classification accuracy is very poor. It is just slightly above chance. Also, since the above model uses only a single tree on one particular training dataset, there is a very high variance in the model. Thus to reduce variance and improve the accuracy of prediction, other methods such as Pruning, Bagging and Random Forests were carried out.

Pruning



From the above plots it can be seen that as the size of the tree continues to grow, the misclassification error rate continuously reduces, becoming a constant at size 5. Since, it is better to have a smaller tree with lesser number of splits to reduce variance, pruning is done with size = 5.

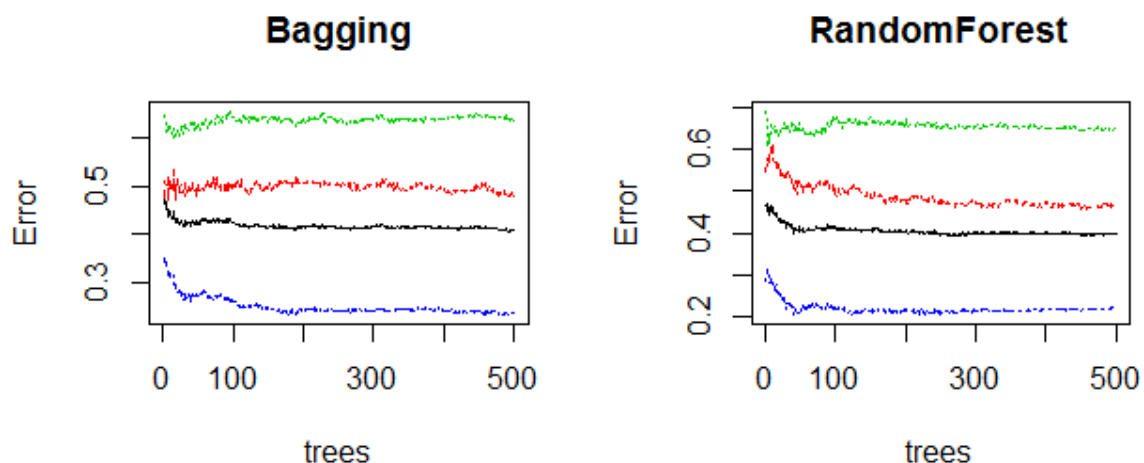
Accuracy and Error rate in Pruned tree

	Training Dataset	Testing Dataset
Accuracy	58.87	55.35

Remark: It can be seen that, there is a slight dip in the training accuracy and an equivalent increase in the testing accuracy. This is mainly because of the reduction in the total number of terminal nodes thus reducing the variance in the model. However, the accuracy is still low and cannot be used for real time prediction.

Bagging and Random Forest

Taking repeated samples from the single training set and generating multiple bootstrapped training sets and fitting a tree model in each of the sets ^[1]. During prediction, the class with a majority across all samples is the predicted class for a data point. This largely reduces variance and can significantly improve accuracy. The main difference between Bagging and Random Forest lies in the number of predictors randomly sampled. In case of bagging all predictors are used and in random forest square root of the total number of predictors is used. Hence, this reduces the correlation between predictors and tends to decrease variance further.

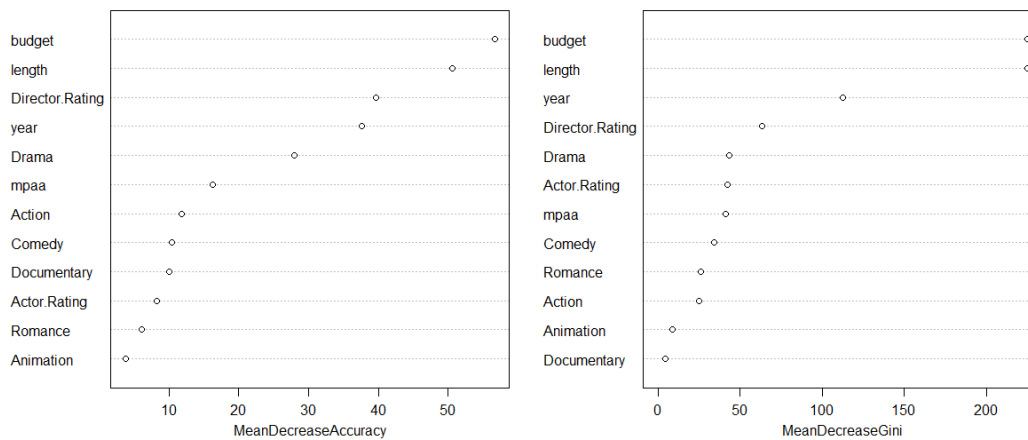


As can be seen from the plot, it is evident that increasing the number of trees is not significantly decreasing the error. However, there is a slight improvement in accuracy as noted below.

	Bagging	RandomForest
Accuracy	58.82%	61.45%

Important Variables

From the below plot, it can be said that the important variables are budget, length of the movie, director rating and year of release. This is true using both accuracy and purity index. Models built using only the top four important variables does not improve results significantly. However, it is very useful for understanding the nature of predictors and the overall dataset.



Model built using top 4 important predictors

	Bagging	Random Forest
Accuracy	59.2%	62.7%

Remark: The model does not significantly improve when only the top four important predictors are used. However, the results don't drop when the rest of the predictors are not used. This shows that these are the four most important predictors driving the model but is not sufficient to adequately predict the ratings. This is likely due to the fact that trees are oversimplifying the prediction process. Hence, to verify, the data was modelled using Logistic regression and SVM.

Multinomial Logistic Regression

A standard logistic regression model is used for 2 class response variable. In case of additional classes, other modelling techniques such as Multinomial Logistic regression is preferred.

Dataset Preprocessing:

Dataset is divided in the ratio 3:2 with three parts for training and two parts for testing. The model is built on the training dataset and is finally tested for accuracy and error rate on the testing dataset.

The package used for this project is 'mlogit' in R. Before 'mlogit' function can be applied the dataset originally in the Wide format needs to be converted to a Long format. The long format consists of each row being repeated the same number of times as there are classes in the response variable. For instance, if there are three classes, each row gets repeated thrice, one row for every class. The output would be the probability of a data point being in a specific class. The class with the highest probability is chosen. The regression coefficients are in the log odds scale.

Model Selection:

For the purpose of building a model with suitable number of predictors, forward selection was the chosen approach. Predictors that gave the highest increase in log likelihood and McFadden R Square values were added one at a time until a stopping criteria was met, in this case the log likelihood value stagnating to a constant.

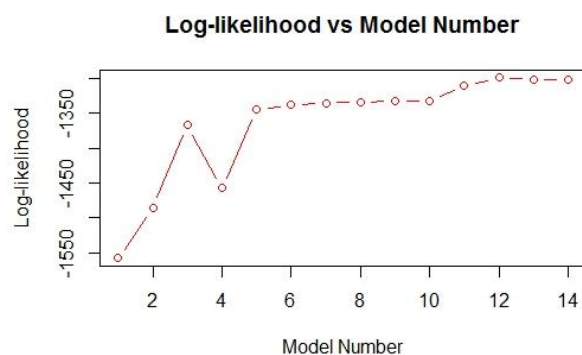
Analysis:

The figure below depicts the graph of log-likelihood vs Model Number. It can be seen that, model 12, gives the highest values of log-likelihood. This model corresponds to the following predictors namely, actor rating, director rating, length, budget, MPAA rating, animation and drama.

From the graph it can also be observed that Model 4 tends to reduce the Log – Likelihood significantly. The model 4 corresponds to Actor rating, director rating and length with length being the newly added term. This indicates the presence of a hidden relationship between length, actor rating and director rating, that yet needs to be explored.

To improve the overall accuracy, possible interactions between the predictors chosen in model 12 were explored. Interactions between director rating and length; budget and MPAA rating; animation and drama improved the overall model in terms of log-likelihood values. Thus the final proposed model consists of actor rating as an independent variable and the above three interactions as other variables.

As expected from model 4, interaction between length and director rating tends to increase log-likelihood whereas length as an independent variable decreases the log likelihood. A possible explanation would be, long movies are usually not preferred unless they are made by good actors and directors.



	Multinomial Logistic Regression
Accuracy (%)	60.07
Error Rate	39.93

Using multinomial logistic regression technique on IMDB data set to predict the rating of any movie, it can be concluded that predictors only affect the movie rating to a certain extent. The model proposed here has an accuracy of about 60% which is reasonable but there is a significant scope for improvement using other classification techniques.

Support Vector Machines

Support vector machines are generalization of a simple and intuitive classifier called the maximal margin classifier, which separates data using a hyper-plane, of the form

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p = 0, \text{ Where, } p \text{ is the number of dimensions}$$

This hyperplane is chosen to be one farthest away from the training observations by computing the (perpendicular) distance from each training observation to a given separating hyper-plane; the smallest such distance is the minimal distance from the observations to the hyperplane, and is known as the margin. Support vectors are the training points, which are closest/on the margin i.e. these points will directly influence the shape of the hyper-plane ^[3]. The maximal margin classifier is a very natural way to perform classification, if a separating hyperplane exists; else there are techniques like soft margins, kernel trick etc. to help us in that endeavor.

In building our model for the IMDB dataset, which is a classification problem, we noticed complex relationships between the data points and the predictors. This coupled with the fact that our response

is multiclass, we decided to apply SVM, to see whether we get improved results compared to the methods we used earlier.

We observed a large number of variables in our dataset, which suggests that we should start with a linear kernel. We then proceed to building the model, on a sampled training dataset with all the predictor variables at once. Since, our response is multiclass, the SVM model will be generated using a 'one-against-one' approach in which $k*(k-1)/2$, k is the number of classes in the response, binary classifiers are trained and the predicted class is determined by a voting scheme ^[4].

	Training Dataset	Testing Dataset
Accuracy %	95.3 %	87.8 %

After generating the model we looked at tuning the SVM, to get optimal '*cost*' and '*gamma*' value in order to minimize the error rate across a 10-fold cross validation. And thereafter using, said values were able to generate the optimal SVM model.

	Training Dataset	Testing Dataset
Accuracy % with cost = 0.1	95.7 %	89.12 %

It is difficult to visualize the classification of the data points in an n -dimension hyperspace, but using the predicted values from the model and the confusion matrix generate thereafter, we are able to make reasonable predictions.

Results

From the analysis above it can be said that, Data Exploration, Tree Based Methods and Multinomial Logistic Regression gave concurrent results indicating certain predictors to be more important than others. The predictors are Actor Rating, Director rating, year, length and budget primarily. This implies the model gives the best results with these predictors and the model either deteriorates or does not significantly improve on the addition of other predictors. However in the case of Support Vector Machines, the best results were obtained when all of the predictors were used in the model building process. Also, the results of SVM tends to be significantly better than the results of Tree based Methods and Logistic Regression.

This could be attributed to the fact that, there are certain complex dependencies between predictors that was overlooked by Trees and the Logistic Model. Trees tend to make splits based on the strongest predictors, thus choosing predictors that could be highly correlated among themselves (bagging) or ignoring a complex interaction between two predictors. Also trees tend to oversimplify the classification process. In the logistic model however, it is required to manually understand various possible interactions thus missing out on some vital interaction. SVM's on the other hand does not have both of the above constraints thus emerging as a suitable model for such complex scenarios.

It can be concluded that prediction of movie ratings is possible with a certain degree of accuracy with Support Vector Machines.

Discussion

Performing the analysis and observing the results, we felt to be able to build better models, it is important to participate in the model building process. In the dataset used in the above analysis, we are unsure of how the data was collected and the biases it is likely to have. The model built was not able to capture these fallacies. We also noticed, handling the missing values in a different way could drastically change the results. An interesting future work could involve building models on an independently collected dataset and comparing the results with this. That way the true prediction accuracy can be estimated.

References

1. http://www.rmdk.ca/boosting_forests_bagging.html
2. https://en.wikipedia.org/wiki/Support_vector_machine
3. An Introduction to Statistical Learning with Applications in R by : Gareth James, Daniela Witten, Trevor Hastie and Robert Tibshirani
4. <http://www.inside-r.org/node/57517>

Appendix

- Code to extract and clean the dataset from IMDB is in the following Github link below. The final dataset used in the analysis is saved as 'movies.csv' in the same link.

<https://github.com/Qartks89/DataMiningIMDBParsing>

- The R code for the rest of the analysis is as below.

```
1. ##### Exploratory Data Analysis #####
2.
3. require(plyr)
4. require(ggplot2)
5. require(reshape2)
6. require(gridExtra)
7.
8.
9. dat <- read.csv("movies.csv", header = T)
10.
11. dat$Director.Rating <- as.factor(dat$Director.Rating)
12. dat$Actor.Rating <- as.factor(dat$Actor.Rating)
13.
14.
15.
16. # Average Budget of Movies in Each Decade
17.
18. meanbudget <- tapply(dat$budget, dat$year, mean)
19. meanbudget <- as.data.frame(meanbudget)
20. meanbudget$years1 <- c("1910s", "1920s", "1930s", "1940s", "1950s", "1960s", "1970s",
21.   "1980s", "1990s", "2000s")
22. plot1 <- ggplot(data=meanbudget, aes(x = years1, y = meanbudget)) +
23.   geom_histogram(stat="identity", fill = "blue", alpha = 0.5) +
24.   coord_cartesian(ylim=c(1000, 4000000)) +
25.   xlab("Decades") + ylab("Average budget") + ggtitle("Average Budget of Movies in Each Decade") +
26.   theme(axis.text.x = element_text(angle = 30, hjust = 1, size = 14))
27.
28.
29. # Average rating of movies in each decade
30.
31. meanrating <- tapply(dat$rating, dat$year, mean)
32. meanrating <- as.data.frame(meanrating)
33. meanrating$years1 <- c("1910s", "1920s", "1930s", "1940s", "1950s", "1960s", "1970s",
34.   "1980s", "1990s", "2000s")
35. plot2 <- ggplot(data=meanrating, aes(x = years1, y = meanrating)) +
36.   geom_histogram(stat = "identity", colour = "blue", fill = "red", alpha = 0.5) +
37.   coord_cartesian(ylim = c(5, 8.5)) +
38.   xlab("Decades") + ylab("Average rating of movies") + ggtitle("Average Rating of Movies in Each Decade")
39.
40.
41. grid.arrange(plot1, plot2, ncol=2, nrow=1)
42.
43.
44. # Type of genres over the decades
45.
46. dat$Drama <- as.numeric(dat$Drama)
47. movie_data_sub <- dat[, c(1,7,8,9,10,11,23,24,25,26,27,28)]
48. movie_data_sub <- melt(movie_data_sub, c(1,2,3,4,5,6))
49. names(movie_data_sub)[7] <- c("Genre")
50. movie_data_sub <- subset(movie_data_sub, value == 1)
51.
```

```

52. plot3 <- ggplot(movie_data_sub, aes(year)) + geom_bar(fill = "red", alpha = 0.4) +
53.   facet_wrap(~ Genre) + theme(axis.text.x = element_text(angle = 45, hjust = 1, size = 18),
54.     plot.title = element_text(size = 20, face = "bold"),
55.     axis.title.x = element_text(size = 20, face = "bold"),
56.     axis.title.y = element_text(size = 20, face = "bold")
57.   ) +
57.   xlab("Decades") + ylab("Number of Movies") + ggtitle("Distribution of Each Genre over Decades")
58. plot3
59.
60.
61. # Box plot rating and genre
62.
63.
64. plot4 <- ggplot(data = movie_data_sub, aes(x = Genre, y = rating, fill = Genre)) +
   geom_boxplot() +
65.   xlab("Genre") + ylab("Rating") + ggtitle("Distribution of ratings for various genre") +
66.   theme(axis.text.x = element_text(angle = 30, hjust = 1, size = 14))
67.
68.
69. # Distribution of budget for various genre boxplot
70.
71. plot5 <- ggplot(data = movie_data_sub, aes(x = Genre, y = budget, fill = Genre)) +
72.   geom_boxplot() + xlab("Genre") + ylab("Budget") + ggtitle("Distribution of budget for various genre") +
73.   theme(axis.text.x = element_text(angle = 30, hjust = 1, size = 14))
74.
75.
76. grid.arrange(plot4, plot5, ncol=2, nrow=1)
77.
78.
79.
80. # Scatter plot between length and rating over director rating
81.
82. empty <- ggplot()+geom_point(aes(1,1), colour="white") +theme(
83.   plot.background = element_blank(),
84.   panel.grid.major = element_blank(),
85.   panel.grid.minor = element_blank(),
86.   panel.border = element_blank(),
87.   panel.background = element_blank(),
88.   axis.title.x = element_blank(),
89.   axis.title.y = element_blank(),
90.   axis.text.x = element_blank(),
91.   axis.text.y = element_blank(),
92.   axis.ticks = element_blank())
93.
94. # Scatterplot of x and y variables
95.
96. scatter <- ggplot(data = dat,aes(x = length, y = rating)) + geom_point(aes(color = Director.Rating)) +
97.   scale_color_manual(values = c("orange", "purple", "green")) +
98.   ggtitle("Scatter Plot") + xlab("Length of Movies") +
99.   ylab("Rating of Movies") +
100.   theme(legend.position=c(1,1),legend.justification=c(1,1),
101.     plot.title = element_text(size = 20, face = "bold"),
102.     axis.title.x = element_text(size=20, face = "bold"),
103.     axis.title.y = element_text(size=20, face = "bold"))
104.
105.   # Marginal density of x - plot on top
106.
107.   plot_top <- ggplot(data = dat, aes(x = length, fill = Director.Rating)) +
108.     geom_density(alpha=.5) +
109.     scale_fill_manual(values = c("orange", "purple", "green")) +

```

```

110.     theme(legend.position = "none") + ggtitle("Length Distribution") + xlab("
Length of Movies") +
111.     ylab("Density")
112.
113.     # Marginal density of y - plot on the right
114.
115.     plot_right <- ggplot(data = dat, aes(x = rating, fill = Director.Rating)) +
116.     geom_density(alpha=.5) +
117.     coord_flip() +
118.     scale_fill_manual(values = c("orange", "purple", "green")) +
119.     theme(legend.position = "none") + ggtitle("Rating Distribution") + xlab("R
ating of Movies") +
120.     ylab("Density")
121.
122.
123.     # The above three plots combined
124.
125.     plot6 <- grid.arrange(plot_top, empty, scatter, plot_right, ncol=2, nrow=2,
widths=c(4, 1), heights=c(1.5, 4))
126.
127.
128.
129.     # Drama v/s All Genre Scatter
130.
131.     dat$Drama <- as.factor(dat$Drama)
132.     scatter2 <- ggplot(data = dat, aes(x = length, y = rating)) + geom_point(aes(
color = Drama)) +
133.     scale_color_manual(values = c("orange", "purple")) +
134.     ggtitle("Sactter plot: Length vs Rating for Drama Genre") +
135.     xlab("Length of Movies") + ylab("Rating of Movies") +
136.     theme(legend.position=c(1,1), legend.justification=c(1,1),
137.     plot.title = element_text(size = 20, face = "bold"),
138.     axis.title.x = element_text(size=20, face = "bold"),
139.     axis.title.y = element_text(size=20, face = "bold"))
140.
141.     # Drama and all other genre length density plot
142.
143.     plot_top2 <- ggplot(data = dat, aes(x = length, fill = Drama)) + geom_densit
y(alpha=.5) +
144.     scale_fill_manual(values = c("orange", "purple")) +
145.     ggtitle("Density Plot Over Length of Movies") + xlab("Length of Movies") +
ylab("Density")
146.
147.     # Drama and all other genre rating density plot
148.
149.     plot_right2 <- ggplot(data = dat, aes(x = rating, fill = Drama)) + geom_dens
ity(alpha=.5) +
150.     scale_fill_manual(values = c("orange", "purple")) +
151.     ggtitle("Density Plot Over Rating of Movies") + xlab("Rating of Movies") +
ylab("Density")
152.
153.
154.     # The above three plots combined
155.
156.     plot7 <- grid.arrange(plot_top2, empty, scatter2, plot_right2, ncol=2, nrow=
2, widths=c(4, 1.5), heights=c(1.5, 4))
157.
158.
159.
160.     # Director names and number of movies
161.
162.     dat2 <- dat[(dat$Director != " "), ]
163.     dat2 <- dat2[(dat2$Director != "De"), ]
164.     dat2 <- dat2[(dat2$Director != " De"), ]
165.     dat2 <- within(dat2, Director <- factor(Director, levels=names(sort(table(Di
rector), decreasing=TRUE))))
166.
167.     plot8 <- ggplot(data=dat2, aes(x = Director, fill = Actor.Rating)) + geom_hi
stogram(stat="bin", width = .5) +

```

```

168.     coord_cartesian(xlim=c(0.5, 12.5), ylim=c(5, 30)) + ggtitle("Directors wi
th Maximum Movies") +
169.     xlab("Directors") + ylab("Movies Directed") +
170.     theme(legend.position = c(.9, .8), axis.text.x = element_text(angle = 30,
hjust = 1, size=19),
171.           plot.title = element_text(size = 20, face = "bold"),
172.           axis.title.x = element_text(size=20, face = "bold"),
173.           axis.title.y = element_text(size=20, face = "bold"))
174.     plot8
175.
176.
177.
178.     ##### Trees #####
#####
179.
180.
181.     #Use final movies dataset
182.     movies <- read.csv("movies.csv",header = TRUE)
183.
184.     #Splitting the rating into three categories
185.     #Rating below 6.5 is category 3
186.     #Rating below 7.3 and above 6.5 is category 2
187.     #Rating below 9.1 and above 7.3 is category 1
188.
189.     movies$rating <- as.factor(ifelse(movies$rating < 6.5, 3,ifelse(movies$rating
g <= 7.3, 2,ifelse(movies$rating >=9.1,1,NA))))
190.
191.     #MPAA rating
192.     #PG 13, NC- 17, PG - category 1
193.     #R -category 2
194.     #None - category 3
195.
196.     movies$mpaa <- as.factor(ifelse(movies$mpaa == "PG-
13",1,ifelse(movies$mpaa == "PG",1,ifelse(movies$mpaa == "NC-
17",1,ifelse(movies$mpaa == "R",2,3)))))
197.
198.     #Adding an ID column
199.     movies$ID <- seq.int(nrow(movies))
200.
201.     #Converting actor and director ratings to factor variables
202.
203.     movies$Actor.Rating <- as.factor(movies$Actor.Rating)
204.     movies$Director.Rating <- as.factor(movies$Director.Rating)
205.
206.     #Removing all columns that are not predictors
207.     movies <- movies[,c(29,1,5:9,22:28,10)]
208.
209.     #Understanding non -factor variables(Length and budget)
210.     #Length
211.
212.     hist(movies$length, xlab = "length of movie", col = "grey", prob = TRUE)
213.     lines(density(movies$length), col = "blue", lwd = 2)
214.     lines(density(movies$length, adjust = 2),lty = "dotted",col = "darkgreen", 1
wd = 2)
215.     plot(density(movies$length), xlab = "length of movie")
216.     polygon(density(movies$length), col = "Purple",border = "black")
217.
218.     #Testing for normality - Shapiro -Wilk Test
219.
220.     shapiro.test(movies$length)
221.
222.     #Comments - P value is very less. This implies we can reject the null hypoth
esis
223.     #and go with the alternate. Hence, it can be said that the sample does not
224.     #typically follow a normal distribution.
225.
226.     #QQplots
227.
228.     qqnorm(movies$length, col = "grey")
229.     qqline(movies$length, col = 2)

```

```

230.
231.      #From the above graph, it can be said that, there is a deviation from the
232.      #normal distribution at the tails. Also, since the rest of the sample follow
233.      #a typical order and lies largely within a bell shape(as seen above), the
234.      #lack of normality can be overlooked
235.
236.      #Testing for budget
237.
238.      hist(movies$budget, xlab = "Budget of movie", col = "grey", prob = TRUE)
239.      lines(density(movies$budget), col = "blue", lwd = 2)
240.      lines(density(movies$budget, adjust = 2),lty = "dotted",col = "darkgreen", l
241.      wd = 2)
242.      plot(density(movies$budget), xlab = "length of movie")
243.      polygon(density(movies$budget), col = "Purple",border = "black")
244.
245.      #Shapiro- Wilk normality test
246.      shapiro.test(movies$budget)
247.
248.      #QQ plots
249.      qqnorm(movies$budget, col = "grey")
250.      qqline(movies$budget, col = 2)
251.
252.      #Comments: same as the predictor length
253.
254.      #####
255.      ##
256.      #Classification Trees
257.
258.      library(tree)
259.
260.      #Sampling the dataset into train and test
261.
262.      set.seed(10)
263.      s <- sample(movies$ID, 1500, replace = FALSE)
264.      train <- movies[s,]
265.      test <- movies[-s,]
266.
267.      movies$ID <- NULL
268.      train$ID <- NULL
269.      test$ID <- NULL
270.      train$Title <- NULL
271.      test$Title <- NULL
272.
273.      #Tree object with all predictors
274.
275.      tree.obj <- tree(train$rating ~ ., data = train)
276.      summary(tree.obj)
277.      plot(tree.obj)
278.      text(tree.obj, pretty = 0, cex = 0.6)
279.
280.      #Misclassification error rate = 43%
281.
282.      #Testing error
283.
284.      pred <- predict(tree.obj, test, type = "class")
285.      t <- table(pred,test$rating)
286.      t
287.
288.      accuracy <- sum(diag(t))/sum(t)
289.      accuracy
290.
291.      #Accuracy = 54%
292.      #error rate = 46%
293.
294.      #Pruning the trees
295.
296.      set.seed(10)
297.      cv.obj <- cv.tree(tree.obj, FUN = prune.misclass)

```



```

298. names(cv.obj)
299. cv.obj
300. plot(cv.obj)
301.
302. #Error rates
303.
304. par(mfrow = c(1,3))
305. plot(cv.obj$dev ~ cv.obj$size, type = "b", xlab = "size", ylab = "deviation"
)
306. plot(cv.obj$dev ~ cv.obj$k, type = "b", xlab = "k", ylab = "deviation")
307.
308. #Plotting the pruned tree
309. prune.obj <- prune.misclass(tree.obj, best = 5)
310. plot(prune.obj)
311. text(prune.obj, pretty = 0)
312. summary(prune.obj)
313.
314. #On test dataset
315.
316. pred <- predict(prune.obj, test, type = "class")
317. t <- table(pred, test$rating)
318. t
319.
320. accuracy <- sum(diag(t))/sum(t)
321. accuracy
322.
323. #Accuracy after pruning - 54%
324.
325. #Bagging
326.
327. library(randomForest)
328. set.seed(10)
329. bag.tree <- randomForest(train$rating ~ ., data = train, mtry = 12, importan
ce = TRUE)
330. bag.tree
331. plot(bag.tree, main = "Bagging")
332.
333. #prediction using bagging
334.
335. bag.pred <- predict(bag.tree, newdata = test)
336. t <- table(bag.pred, test$rating)
337. t
338.
339. accuracy <- sum(diag(t))/sum(t)
340. accuracy
341.
342. #accuracy after bagging = 59%
343.
344. #RandomForest
345.
346. set.seed(10)
347. random.tree <- randomForest(train$rating ~ ., data = train, mtry = 4, import
ance = TRUE)
348. random.tree
349. plot(random.tree, main = "RandomForest")
350.
351. #prediction using randomForest
352.
353. random.pred <- predict(random.tree, newdata = test)
354. t <- table(random.pred, test$rating)
355. t
356.
357. accuracy <- sum(diag(t))/sum(t)
358. accuracy
359.
360. #accuracy after random forest = 61%
361.
362. #Finding important predictors
363.
364. importance(random.tree)
365. varImpPlot(random.tree)

```

```

366.     summary(random.tree)
367.
368.     #Budget > length > year > director rating
369.     #Tree object using important predictors only
370.
371.     tree.imp <- tree(train$rating ~ train$Director.Rating + train$year + train$length + train$budget, data = train)
372.     summary(tree.imp)
373.     plot(tree.imp)
374.     text(tree.imp, pretty = 0)
375.
376.     #Random Forest for important predictors
377.
378.     data1 <- train[,c(2:5,13)]
379.     rand.imp <- randomForest(data1$rating ~ ., data = data1, mtry = 2, importance = TRUE)
380.     rand.imp
381.     plot(rand.imp)
382.
383.     #Accuracy using random forest
384.
385.     data <- test[,c(2:5)]
386.     random.pred.imp <- predict(rand.imp, newdata = data)
387.     t <- table(random.pred.imp, test$rating)
388.     t
389.
390.     accuracy <- sum(diag(t))/sum(t)
391.     accuracy
392.
393.     #accuracy after using important predictors = 58%
394.
395.     #prediction using randomForest
396.
397.     dat <- test[,c(2:5,10,6,1,13)]
398.     random.pred <- predict(random.tree, newdata = dat)
399.     t <- table(random.pred, test$rating)
400.     t
401.
402.     accuracy <- sum(diag(t))/sum(t)
403.     accuracy
404.
405.     #accuracy after random forest = 61%
406.
407.     #Finding important predictors
408.
409.     importance(random.tree)
410.     varImpPlot(random.tree)
411.     summary(random.tree)
412.
413.     #budget > length > year > director.rating > drama > mpaa > actor rating
414.
415.
416.
417.     ##### Multinomial Logistic Regression #####
418.     #####
419.
420.     data <- read.csv("movies.csv")
421.
422.     hist(data$rating)      # Histogram of rating values
423.     summary(data$rating)   # Summary of rating column
424.
425.
426.     # Rating values below median(=6.5) go into category-
427.     # 3, values between median and 3rd quartile(=7.3) go into category-
428.     # 2 and values greater than 3rd quartile go into category-1
429.
430.     data$rating[data$rating<6.5] <- 3
431.     data$rating[data$rating>7.3] <- 1
432.     data$rating[data$rating>=6.5 & data$rating<=7.3] <- 2

```

```

432.      unique(data$rating)          #unique values of rating column
433.
434.      length(data$rating[data$rating==1])      # total= 593
435.      length(data$rating[data$rating==2])      # total= 722
436.      length(data$rating[data$rating==3])      # total= 1222
437.
438.      table(data$rating)
439.
440.      require(mlogit)
441.
442.      data$rating <- as.factor(data$rating)
443.
444.      # The square root of the budget column values has been done to make values s
mall enough to fit in the model
445.
446.      data$budget <- sqrt(data$budget)
447.
448.
449.      # Divide the data set in training and testing data set approximately in 60:4
0 ratio
450.      set.seed(10)
451.      tr <- sample(1:nrow(data), size = 1500, replace=F)
452.      train <- data[tr,]
453.      test <- data[-tr,]
454.
455.      likelihood <- rep(0,14)
456.
457.      #reshaping the data from wide to long format
458.
459.      mydata <- mlogit.data(train,varying=NULL,choice="rating",shape="wide")
460.      head(mydata)
461.
462.
463.      # Forward variable selection method is applied here to get the most signific
ant variables in the model
464.
465.      model1 <- mlogit(rating ~ 1|Actor.Rating,data=mydata,reflevel="1")
466.      summary(model1)
467.
468.      likelihood[1] <- summary(model1)$logLik[1]
469.
470.      model2 <- mlogit(rating ~ 1|Actor.Rating+Director.Rating,data=mydata,refleve
l="1")
471.      summary(model2)
472.
473.      likelihood[2] <- summary(model2)$logLik[1]
474.
475.      model3 <- mlogit(rating ~ 1|Actor.Rating+Director.Rating+year,data=mydata,re
flevel="1")
476.      summary(model3)
477.
478.      likelihood[3] <- summary(model3)$logLik[1]
479.
480.      model4 <- mlogit(rating ~ 1|Actor.Rating+Director.Rating+length,data=mydata,
reflevel="1")
481.      summary(model4)
482.
483.      likelihood[4] <- summary(model4)$logLik[1]
484.
485.      model5 <- mlogit(rating ~ 1|Actor.Rating+Director.Rating+length+budget,data=
mydata,reflevel="1")
486.      summary(model5)
487.
488.      likelihood[5] <- summary(model5)$logLik[1]
489.
490.      model6 <- mlogit(rating ~ 1|Actor.Rating+Director.Rating+length+budget+mpaa,
data=mydata,reflevel="1")
491.      summary(model6)
492.
493.      likelihood[6] <- summary(model6)$logLik[1]
494.

```

```

495.     model7 <- mlogit(rating ~ 1|Actor.Rating+Director.Rating+budget+length+mpaa+
Action,data=mydata,reflevel="1")
496.     summary(model7)
497.
498.     likelihood[7] <- summary(model7)$logLik[1]
499.
500.     model8 <- mlogit(rating ~ 1|Actor.Rating+Director.Rating+budget+length+mpaa+
Action+Comedy,data=mydata,reflevel="1")
501.     summary(model8)
502.
503.     likelihood[8] <- summary(model8)$logLik[1]
504.
505.     model9 <- mlogit(rating ~ 1|Actor.Rating+Director.Rating+budget+length+mpaa+
Action+Comedy+Documentary,data=mydata,reflevel="1")
506.     summary(model9)
507.
508.     likelihood[9] <- summary(model9)$logLik[1]
509.
510.     model10 <- mlogit(rating ~ 1|Actor.Rating+Director.Rating+budget+length+bud
get+mpaa+Action+Comedy+Documentary+Romance,data=mydata,reflevel="1")
511.     summary(model10)
512.
513.     likelihood[10] <- summary(model10)$logLik[1]
514.
515.     model11 <- mlogit(rating ~ 1|Actor.Rating+Director.Rating+budget+length+mpaa
+Action+Comedy+Documentary+Romance+Animation,data=mydata,reflevel="1")
516.     summary(model11)
517.
518.     likelihood[11] <- summary(model11)$logLik[1]
519.
520.     model12 <- mlogit(rating ~ 1|Actor.Rating+Director.Rating+budget+length+mpaa
+Animation+Drama,data=mydata,reflevel="1")
521.     summary(model12)
522.
523.     likelihood[12] <- summary(model12)$logLik[1]
524.
525.     model13 <- mlogit(rating ~ 1|Director.Rating+budget+length+mpaa+Action+Comed
y+Documentary+Romance+Animation+Drama,data=mydata,reflevel="1")
526.     summary(model13)
527.
528.     likelihood[13] <- summary(model13)$logLik[1]
529.
530.     model14 <- mlogit(rating ~ 1|Director.Rating+length+budget+mpaa+Action+Comed
y+Documentary+Romance+Animation+Drama,data=mydata,reflevel="1")
531.     summary(model14)
532.
533.     likelihood[14] <- summary(model14)$logLik[1]
534.
535.     x <- c(1:14)
536.
537.     # plot of log-likelihood values of all the models
538.
539.     plot(x,likelihood,main="Log-
likelihood vs Model Number",xlab="Model Number",ylab="Log-
likelihood",col="red",type = "b")
540.     likelihood
541.
542.     # based on the log-likelihood and R square values we choose 12th model
543.     # and is selected and analysed further to give better results
544.
545.
546.     # Trying different combinations of interactive predictors, the final model f
or predicting the ratings of movie is the below:
547.
548.     model20 <- mlogit(rating ~ 1|Actor.Rating+Director.Rating*length+budget*mpaa
+Animation*Drama,data=mydata,reflevel="1")
549.     summary(model20)
550.
551.
552.     # Exponential function of logistic coefficients is calculated below:

```

```

553.      # This tells us the log-
      odds of of being in category 2 versus 1 and 3 versus 1.
554.
555.      exp(coef(model20))
556.
557.
558.      #reshaping the test data from wide to long format
559.
560.      testdata <- mlogit.data(test,varying=NULL,choice="rating",shape="wide")
561.
562.
563.      # predict function is used to predict the rating on the basis of proposed mo
      del
564.
565.      p <- predict(model20, newdata = testdata)
566.      p <- data.frame(p)
567.
568.
569.      l <- length(p$X1)
570.      m <- rep(0,l)
571.
572.
573.      # taking out the category from the predicted values which gives the highest
      probability
574.
575.      for (i in 1:l){
576.
577.          m[i] <- which.max(p[i,])
578.
579.      }
580.
581.
582.      #table to form 3 by 3 matrix using predicted and the original values from th
      e data
583.
584.      tab <- table(m,test$rating)
585.      tab
586.
587.      # accuracy is calculated
588.
589.      accuracy <- (tab[1] + tab[5] + tab[9])/sum(tab)
590.      accuracy_percentage <- accuracy*100
591.      accuracy_percentage
592.
593.
594.      # error rate is calculated
595.
596.      error <- 1-accuracy
597.      error_percentage <- error*100
598.      error_percentage
599.
600.
601.
602.      ##### Support Vector Machines #####
      #####
603.
604.
605.      ### Required Library
606.
607.      library(e1071)
608.      library(ROCR)
609.
610.      ### Importing DataSet
611.
612.      movies1 <- read.csv("movies.csv")
613.
614.      movies1$rating <- as.factor(ifelse(movies1$rating < 6.5, 3,ifelse(movies1$ra
      ting <= 7.3, 2,ifelse(movies1$rating <=9.1,1,NA))))
615.
616.
617.      movies1$rating <- as.factor(movies1$rating)

```

```

618. movies1 <- data.frame(movies1)
619.
620. ### Sampling Train and Test Data
621. set.seed(15)
622. s <- sample(nrow(movies1), 1500)
623. train <- movies1[s,]
624. test <- movies1[-s,]
625.
626. ### Checking Number of Classes in Training and Testing
627. table(train$rating)
628. table(test$rating)
629.
630. ### Training Data Frame
631. x = train[,-10]
632. y = train[,10]
633. dat=data.frame(x=x , y=as.factor(y))
634.
635.
636. ### Testing Data Frame
637. xt = test[,-10]
638. yt = test[,10]
639. datt=data.frame(x=xt , y=as.factor(yt))
640.
641. ### Building SVM Model using Linear Kernal
642. svm.model=svm(y ~ ., data=dat , kernel ="linear", cost=10, gamma = 1)
643. summary(svm.model)
644.
645. ### Confusion Matrix for Training Data
646. train.tab = table(svm.model$fitted , dat$y)
647. train.tab
648. ### Checking Accuracy
649. classAgreement(tab)
650.
651. ### Confusion Matrix for Test Data
652. svm.pred <- predict(svm.model, datt)
653. test.tab=table(pred = svm.pred, true = datt$y)
654. ### Checking Accuracy
655. classAgreement(test.tab)
656.
657. ### Tuning The model for Test data
658. tune.out=tune(svm ,y~.,data=datt ,kernel ="linear",
659.               ranges =list(cost=c(0.001, 0.01, 0.1, 1,5,10,100,1000))
660.               , gamma=c(0.5,1,2,3,4))
661. summary(tune.out)

```

Statement of Contributions

	Report Building	Exploratory Analysis	Multinomial Logistic Analysis	Tree Based Method Analysis	SVM Analysis	Dataset Creation
Saurabh	✓		✓			✓
Huzefa	✓	✓				✓
Deepak	✓			✓		✓
Kartikeya	✓				✓	✓