

# Modeling and prediction for movies

## Setup

### Load packages

```
library(ggplot2)

## Warning: package 'ggplot2' was built under R version 4.0.3

library(dplyr)

## Warning: package 'dplyr' was built under R version 4.0.3

library(stats)

## Warning: package 'stats' was built under R version 4.0.3

library(GGally)

## Warning: package 'GGally' was built under R version 4.0.3

library(gridExtra)

## Warning: package 'gridExtra' was built under R version 4.0.3
```

### Load data

```
load("movies.Rdata")
```

## Part 1: Data

The movie data-set gathers data on movies and a few attributes about its structure, like movie length in minutes, release day, month and year, MPAA rating, and so on. Besides these variables, there are a few columns of interest regarding the movie evaluation, such as the Internet Movie Database - IMDb rating, critics rating and audience ratings. The goal of this project is to bring a perspective in one of these classifications, the IMDb ratings, and its relationship with the other variables in the data-set, excluding the other variables of interest.

For this project, it is essential to state the following points:-

- i. The original data-set about each movie's score was taken from both the IMDb and Rotten Tomatoes database. Characteristics of this site user community influence trends in specific movie scores, thus this results can be generalized to the whole population.
- ii. The article was an observational study. There is no causation can be established because random assignment is not used in this study. In this study, I will assume some voting results about a film on website A could be used to predict the same movie's score on website B. However, in the real world the score for the same movie occurs at the same time on two sites A and B.

## Part 2: Research question

A few features regarding the movie can influence how people perceive the movie. These features not only include genre and run-time, but also consider facts like whether the movie was in the Top 200 Box Office list on BoxOfficeMojo.

What are the factors/variables that influence the IMDb ratings received by a movie from the audience? To what extent are these factors significant in predicting the IMDb score of a movie?

This question is important for us because it is important for online streaming services to predict ratings based on these factors. So, they can accordingly invest in contents meeting certain conditions.

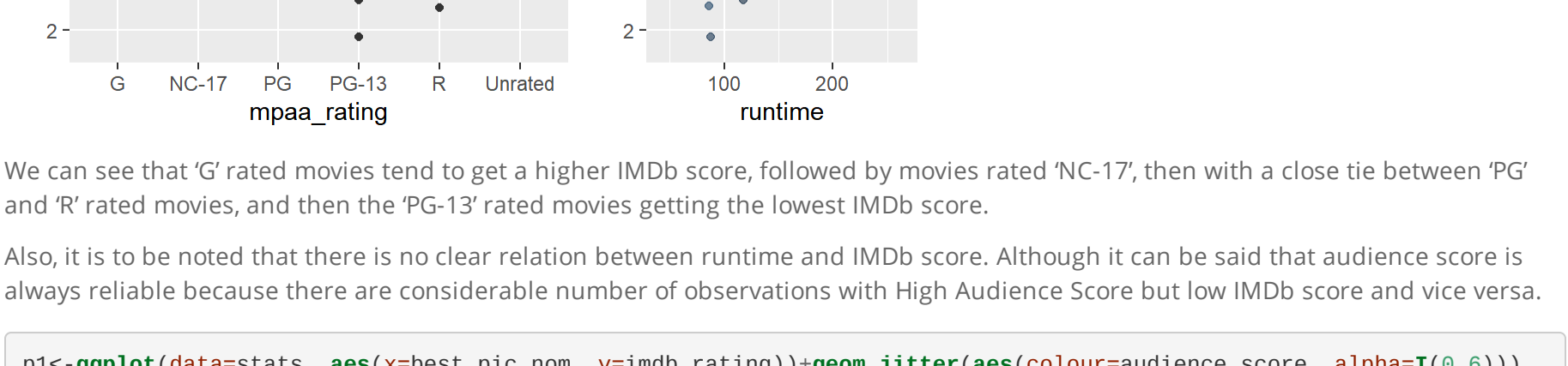
## Part 3: Exploratory data analysis

Let us consider the following variables for constructing a multiple linear regression (MLR) model,

```
genre: Genre of the movie
runtime: Runtime of movie (in movies)
mpaa_rating: MPAA rating of the movie (G, PG, PG-13, R)
critics_score: Critics score on Rotten Tomatoes
audience_score: Audience score on Rotten Tomatoes
best_pic_nom: Whether or not the movie was nominated for a best picture
best_pic_win: Whether or not the movie won a best picture Oscar
top200_box: Whether or not the movie is in the Top 200 Box Office list on BoxOfficeMojo
best_actor_win: Whether or not one of the main actors in the movie ever won an Oscar
best_actress_win: Whether or not one of the main actresses in the movie ever won an Oscar

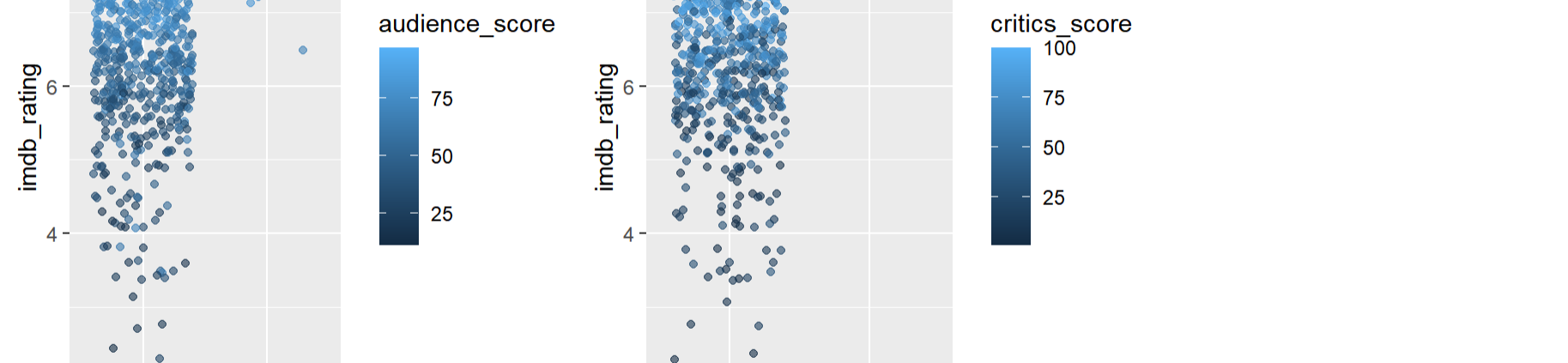
We will use the above variables to try and predict the variable imdb_rating using Multiple Linear Regression model.
First let us clean the model of observations having missing values
```

Note that there is a single observation which has a missing value of one of our variables of our interest. Now, we move on to visualizing the data of IMDb scores based on these factors.

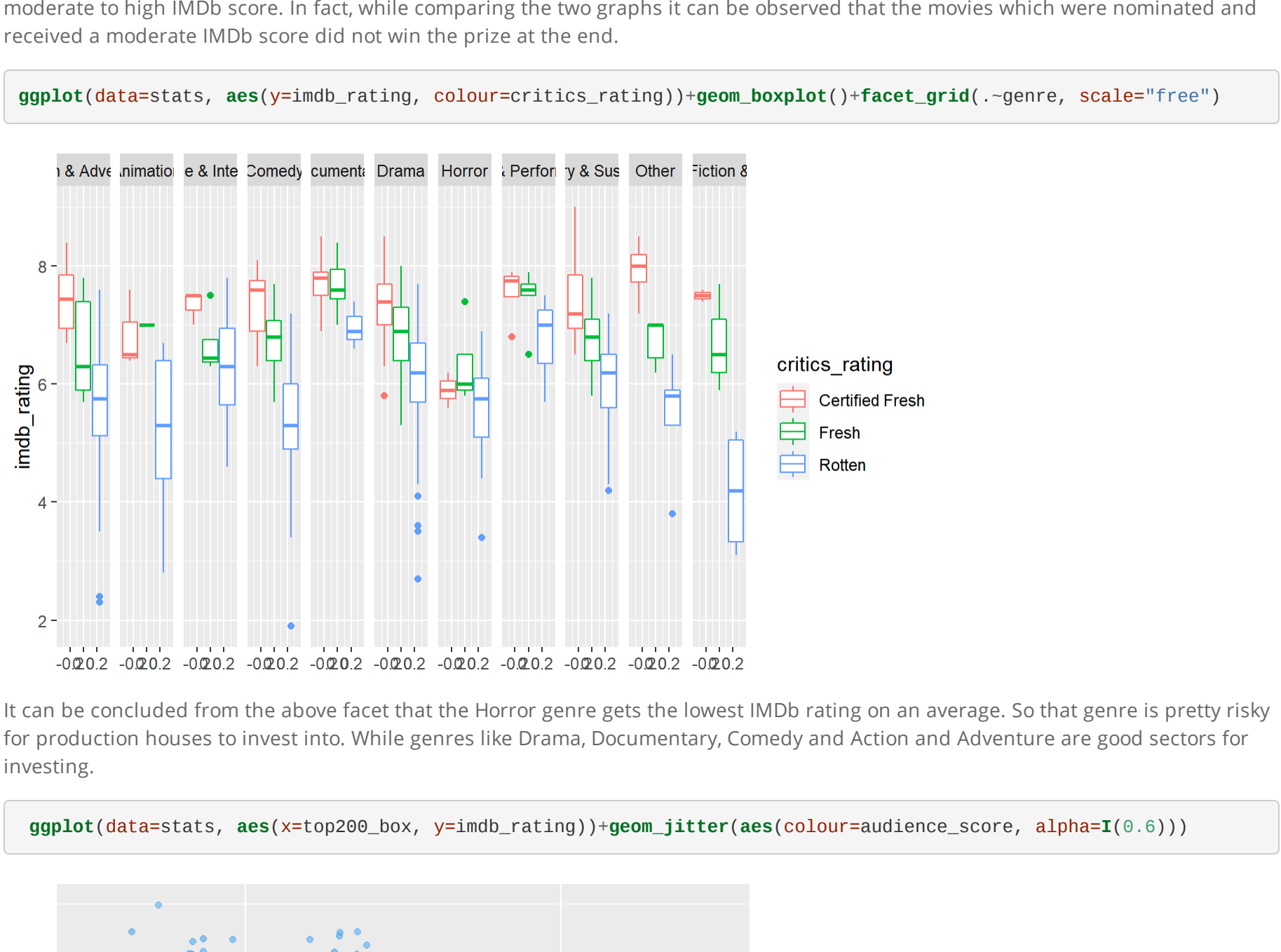


We can see that 'G' rated movies tend to get a higher IMDb score, followed by movies rated 'NC-17'; then with a close tie between 'PG' and 'R' rated movies, and then the 'PG-13' rated movies getting the lowest IMDb score.

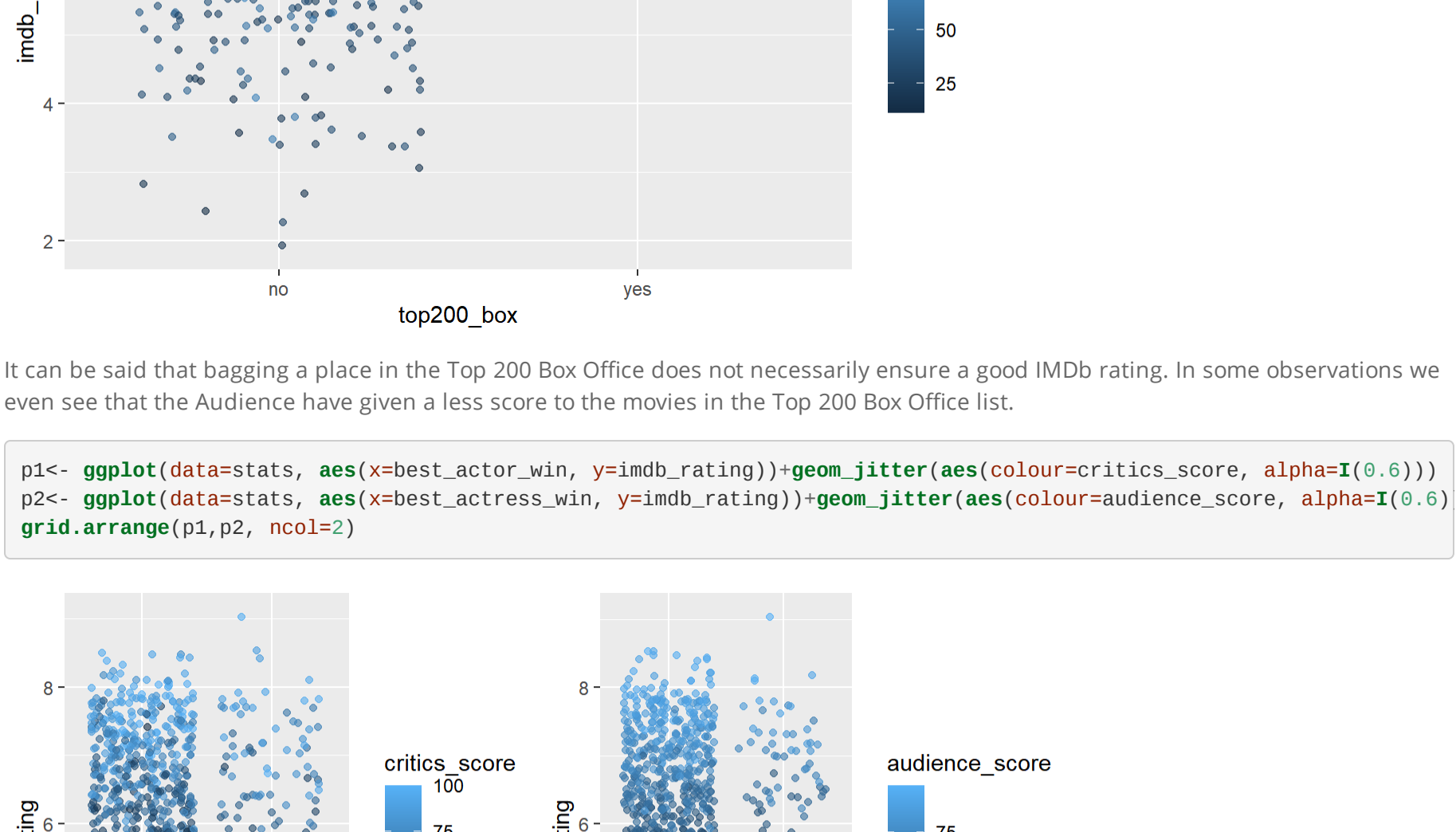
Also, it is to be noted that there is no clear relation between runtime and IMDb score. Although it can be said that audience score is always reliable because there are considerable number of observations with High Audience Score but low IMDb score and vice versa.



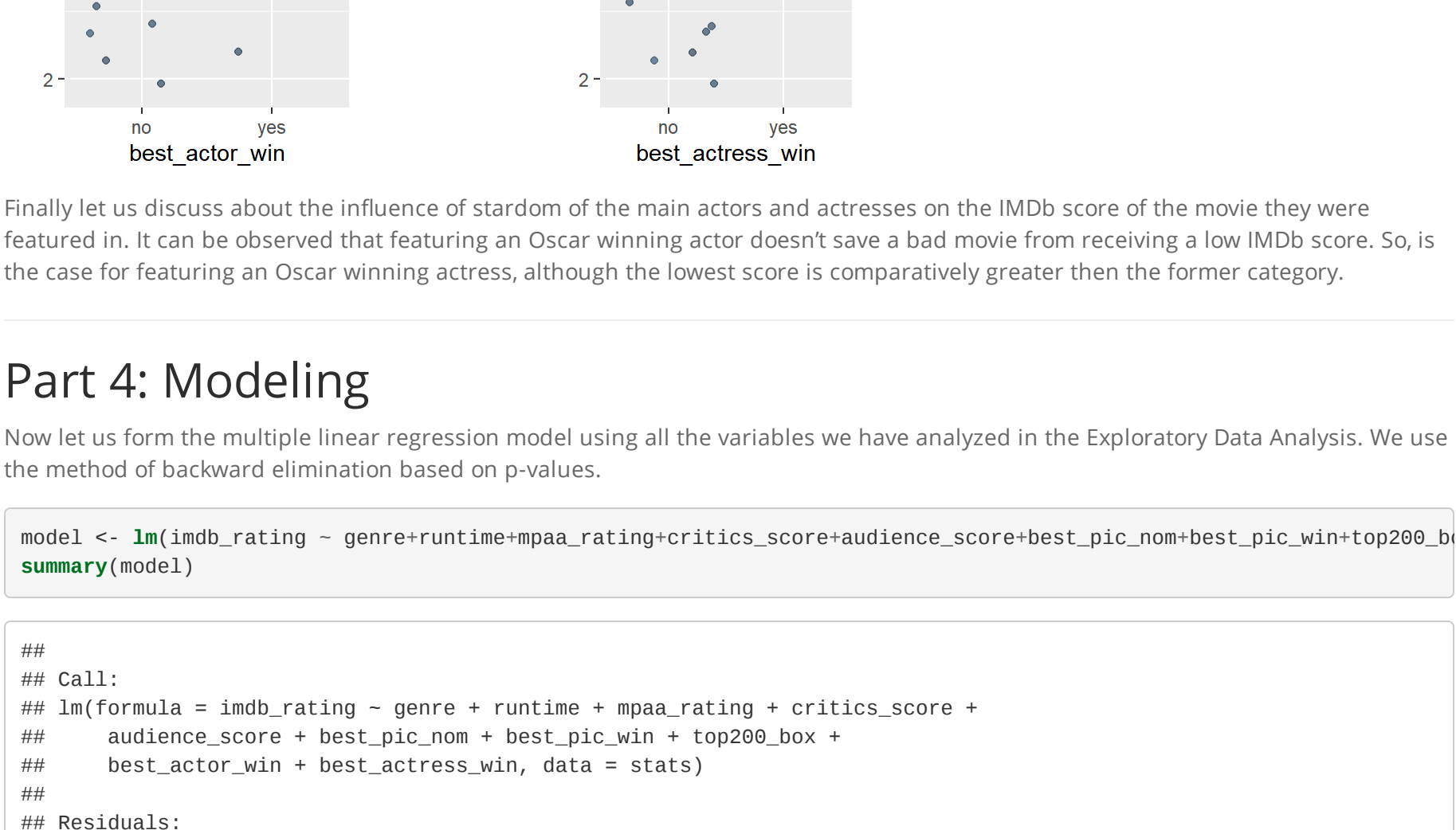
From the above graph we can say that the movies which have been nominated for best picture or won the award for best picture have moderate to high IMDb score. In fact, while comparing the two graphs it can be observed that the movies which were nominated and received a moderate IMDb score did not win the prize at the end.



It can be concluded from the above facet that the Horror genre gets the lowest IMDb rating on an average. So that genre is pretty risky for production houses to invest into. While genres like Drama, Documentary, Comedy and Action and Adventure are good sectors for investing.



It can be said that bagging a place in the Top 200 Box Office does not necessarily ensure a good IMDb rating. In some observations we even see that the Audience have given a less score to the movies in the Top 200 Box Office list.



Finally let us discuss about the influence of stardom of the main actors and actresses on the IMDb score of the movie they were featured in. It can be observed that featuring an Oscar winning actor doesn't save a bad movie from receiving a low IMDb score. So, is the case for featuring an Oscar winning actress, although the lowest score is comparatively greater than the former category.

## Part 4: Modeling

Now let us form the multiple linear regression model using all the variables we have analyzed in the Exploratory Data Analysis. We use the method of backward elimination based on p-values.

```
model <- lm(imdb_rating ~ genre+runtime+mpaa_rating+critics_score+audience_score+best_pic_nom+best_pic_win+top200_box,
summary(model))

## Call:
## lm(formula = imdb_rating ~ genre + runtime + mpaa_rating + critics_score +
## audience_score + best_pic_nom + best_pic_win + top200_box +
## best_actor_win + best_actress_win, data = stats)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.35880  -0.19831  0.03439  0.28111  1.19759
##
## Coefficients:
## (Intercept)                3.293966  0.172484 19.235  < 2e-16 ***
## genreAnimation             -0.430195  0.183393  -2.348  0.01926 *
## genreArt_House & International 0.2113828 0.1415625  1.493  0.13589
## genreComedy                -0.1485145 0.0783843  -1.895  0.05859
## genreDocumentary           0.3127179 0.1072217  2.963  0.00382 **
## genreDrama                 0.0393476  0.0684997  0.574  0.56589
## genreError                 0.0828272  0.1171954  0.787  0.47999
## genreMusical & Performing Arts 0.0222248 0.1512875  0.214  0.83688
## genreMystery & Suspense      0.2245986  0.0889947  2.548  0.01114 *
## genreOther                 -0.0466777 0.1327961  -0.351  0.72533
## genreScience Fiction & Fantasy -0.1943448 0.1607645  -1.105  0.24431
## runtime                   0.0050390  0.0011668  4.159  6.46e-06 ***
## mpaa_ratingNC-17          -0.1712827  0.3551833  -0.482  0.62973
## mpaa_ratingPG             -0.1564002  0.1293984  -1.210  0.22692
## mpaa_ratingPG-13         -0.1842124  0.1202650  -0.761  0.44361
## mpaa_ratingR              -0.0716892  0.1286623  -0.557  0.57765
## mpaa_ratingUnrated        -0.1952550  0.1468137  -1.339  0.18492
## critics_score              0.0164261  0.0089597 10.864  < 2e-16 ***
## audience_score             0.0336862  0.0013594 25.268  < 2e-16 ***
## best_pic_nomies            -0.0366553 0.1218988  -0.251  0.80165
## mpaa_ratingR              0.1089759 0.2044895  -0.549  0.58733
## top200_boxies              0.0132653 0.1273853  0.104  0.91718
## best_actor_winyes          0.0371077 0.0556715  0.677  0.49842
## best_actress_winyes        0.0615546 0.0617252  0.997  0.31904
##
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.487 on 626 degrees of freedom
## Multiple R-squared:  0.821, Adjusted R-squared:  0.8147
## F-statistic: 125.1 on 23 and 626 DF, p-value: < 2.2e-16
```

The variable top200\_box has the highest p-value and thus it bears the least statistical significance. So, we refit the MLR model without that variable.

```
model <- lm(imdb_rating ~ genre+runtime+mpaa_rating+critics_score+audience_score+best_pic_nom+best_pic_win+best_actor_win+best_actress_win,
summary(model))

## Call:
## lm(formula = imdb_rating ~ genre + runtime + mpaa_rating + critics_score +
## audience_score + best_pic_nom + best_pic_win + best_actor_win +
## best_actress_win, data = stats)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.35791  -0.19841  0.03438  0.28138  1.19769
##
## Coefficients:
## (Intercept)                3.2934910  0.1749447 19.255  < 2e-16 ***
## genreAnimation             -0.4283592  0.1823424  -2.349  0.01912 *
## genreArt_House & International 0.2122249 0.1412199  1.583  0.11339
## genreComedy                -0.1478841 0.0783253  -1.894  0.05864
## genreDocumentary           0.3124765 0.1065180  2.934  0.00347 **
## genreDrama                 0.0402380 0.0679103  0.593  0.55372
## genreError                 0.0803656  0.1169889  0.713  0.47636
## genreMusical & Performing Arts 0.0336115 0.1506631  0.223  0.82354
## genreMystery & Suspense      0.2256678 0.0876914  2.567  0.01058 *
## genreOther                 -0.0462124 0.1326255  -0.349  0.72745
## genreScience Fiction & Fantasy -0.1947910 0.1665779  -1.169  0.24270
## runtime                   0.0050197  0.0010997  4.565  6.82e-06 ***
## mpaa_ratingPG             -0.1698438 0.1354744  -0.479  0.63241
## mpaa_ratingPG-13         -0.1553253 0.1257840  -0.396  0.69277
## mpaa_ratingR              -0.1028108 0.1325647  -0.776  0.43831
## mpaa_ratingUnrated        -0.0701185 0.1276755  -0.549  0.58311
## critics_score              0.0336864 0.0013592 25.268  < 2e-16 ***
## audience_score             0.0164265 0.0089575 10.883  < 2e-16 ***
## best_pic_nomies            -0.0366619 0.0013370 25.327  < 2e-16 ***
## mpaa_ratingR              0.0336827 0.1228826  -0.251  0.80161
## best_actress_winyes        0.1073802 0.2641933  0.525  0.59943
## best_actor_winyes          0.0376123 0.0556196  0.676  0.49914
## best_actress_winyes        0.0612119 0.0615888  0.994  0.32066
##
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4661 on 627 degrees of freedom
## Multiple R-squared:  0.8213, Adjusted R-squared:  0.815
## F-statistic: 131 on 22 and 627 DF, p-value: < 2.2e-16
```

We notice a slight increase in the adjusted R-square after removing the previous insignificant variable. Now, we remove the next least significant variables best\_pic\_nom, best\_pic\_win, best\_actor\_win, best\_actress\_win, mpaa\_rating, top200\_box with pretty high p-value.

```
model <- lm(imdb_rating ~ genre+ runtime+critics_score+audience_score, data=stats)
summary(model)

## Call:
## lm(formula = imdb_rating ~ genre + runtime + critics_score +
## audience_score, data = stats)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.34430  -0.20098  0.03524  0.27885  1.17364
##
## Coefficients:
## (Intercept)                3.1675348  0.1251241 25.315  < 2e-16 ***
## genreAnimation             -0.3681453 0.1668808  -2.286  0.0277 *
## genreArt_House & International 0.1997289 0.1376430  1.451  0.1473
## genreComedy                -0.1418076 0.0766630  -1.839  0.0663
## genreDocumentary           0.2611971 0.0945446  2.783  0.0059 **
## genreDrama                 0.0572713 0.0655556  0.875  0.3818
## genreError                 0.0953283 0.1116151  0.835  0.4048
## genreMusical & Performing Arts 0.0156689 0.1491699  0.105  0.9184
## genreMystery & Suspense      0.2613679 0.0846405  3.088  0.0021 **
## genreOther                 -0.0599035 0.1211552  -0.457  0.6489
## genreScience Fiction & Fantasy -0.1913924 0.1569902  -1.153  0.2494
## runtime                   0.0052878 0.0013370 25.327  < 2e-16 ***
## critics_score              0.0164097 0.0089876 11.068  < 2e-16 ***
## audience_score             0.1073802 0.2641933  0.525  0.59943
## best_actor_winyes          0.0376123 0.0556196  0.676  0.49914
## best_actress_winyes        0.0612119 0.0615888  0.994  0.32066
##
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4657 on 636 degrees of freedom
## Multiple R-squared:  0.8184, Adjusted R-squared:  0.8157
## F-statistic: 222 on 13 and 636 DF, p-value: < 2.2e-16
```

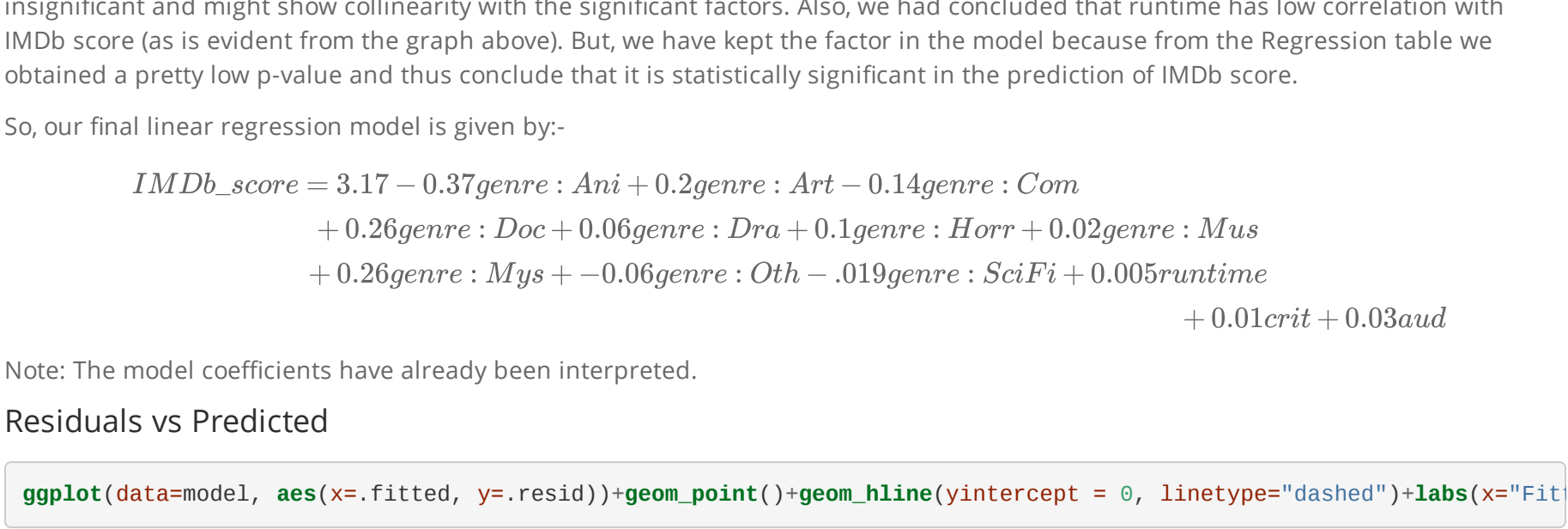
In the above model we see that some of the factor 'Genre' are pretty insignificant, but we retain the factor in our model because it has some levels which are statistically very significant with low p-values. This is our final MLR model.

Note that, if the Genre falls into the categories 'Animation', 'Comedy', 'Science Fiction and Fantasy' and 'Other', the predicted IMDb score tends to decrease (when all other factors are held constant), as the estimated slope coefficient of these factors are negative.

Also, note that the coefficient of determination r-square is pretty high (also the Adjusted r-square), so 81.94% of the variability can be explained.

## Checking conditions and Model Diagnostics

Let us first check for collinearity between the explanatory variables using pairwise plot.



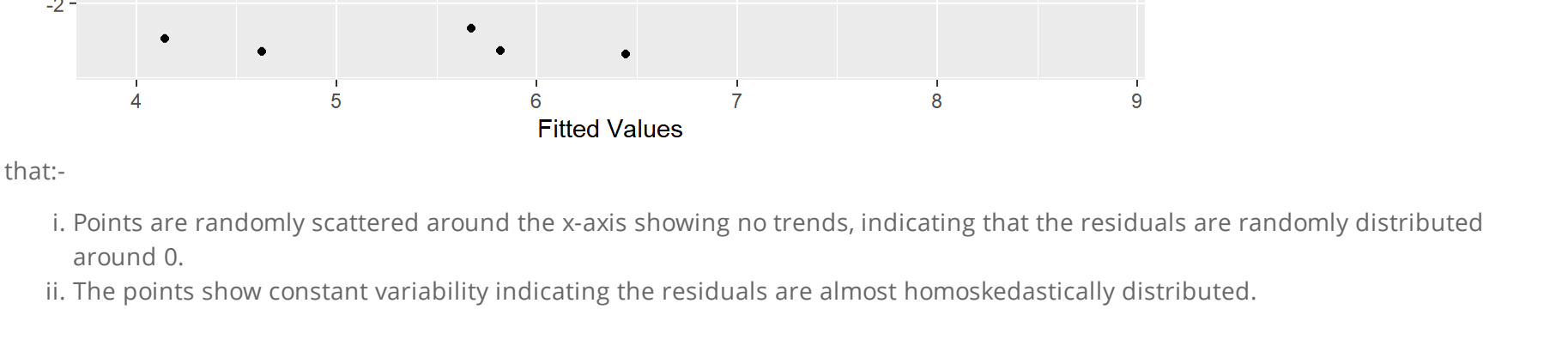
We do not see any collinearity between the explanatory variables because we have already removed the factors which are statistically insignificant and might show collinearity with the significant factors. Also, we had concluded that runtime has low correlation with IMDb score (as is evident from the graph above). But, we have kept the factor in the model because from the Regression table we obtained a pretty low p-value and thus conclude that it is statistically significant in the prediction of IMDb score.

So, our final linear regression model is given by:-

$$IMDb\_score = 3.17 - 0.37genre : Ani + 0.29genre : Art - 0.14genre : Com + 0.26genre : Doc + 0.06genre : Dra + 0.19genre : Horr + 0.02genre : Mys + 0.26genre : Mys + -0.06genre : Oth - .0119genre : SciFi + 0.005runtime + 0.01crit + 0.03aud$$

Note: The model coefficients have already been interpreted.

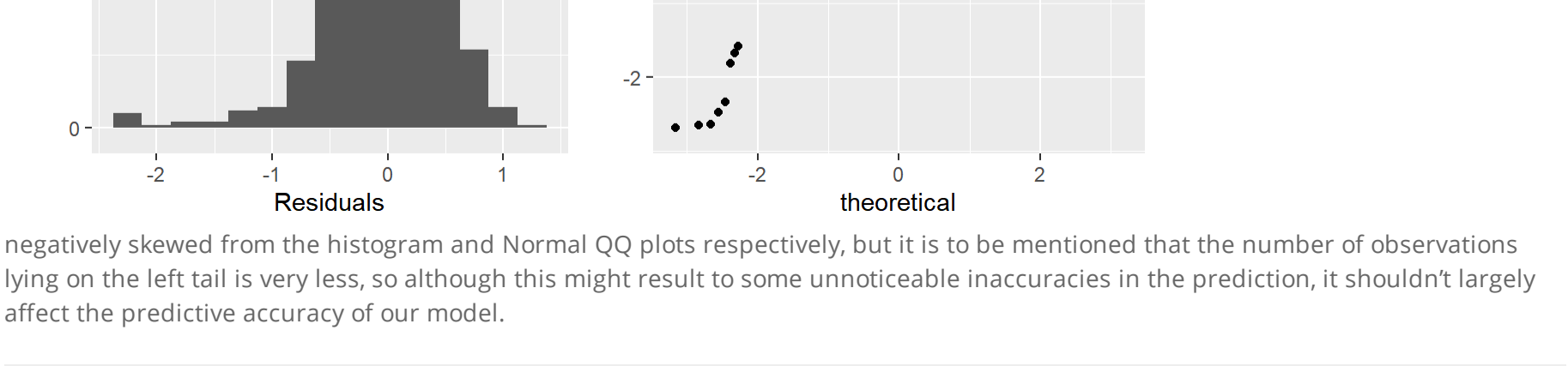
### Residuals vs Predicted



From the above plot we conclude that:-

- i. Points are randomly scattered around the x-axis showing no trends, indicating that the residuals are randomly distributed around 0.
- ii. The points show constant variability indicating the residuals are almost homoskedastically distributed.

## Histogram and QQ Plot of residuals



We see that the residuals are a bit negatively skewed from the histogram and Normal Q-Q plots respectively, but it is to be mentioned that the number of observations lying on the left tail is very less, so although this might result to some unrealistic inaccuracies in the prediction, it shouldn't largely affect the predictive accuracy of our model.

## Part 5: Prediction

For prediction, let us consider the movie "Silence" by Martin Scorsese, released in 2016.

For our model, the required data are (source: Google):-

- i. Genre: Drama
- ii. Runtime: 161 minutes
- iii. Critics Score: 83
- iv. Audience Score: 69

Let us create a data frame for predicting the IMDb score of this movie.

```
Silence <- data.frame(genre="Drama", runtime=161, critics_score=83, audience_score=69)

Now, we predict the IMDb score of this movie and also the interval of precision based on the available data.

predict(model, Silence, interval = "prediction", level=0.95)

##      fitted      lower      upper
## 1 7.278885 6.357171 8.2006
```

Note that, the observed IMDb score is 7.2 which is almost equal to the predicted value of 7.27. Thus, the fit is extremely good. The 95% confidence interval for this estimate is (6.36, 8.20).

So, we can say that the probability that the above interval contains the IMDb score of a movie under the "Drama" genre, running for 161 minutes, with a Rotten Tomatoes Critic Score of 83% and an Audience Score of 69%, is 0.95

## Part 6: Conclusion

We have analyzed different factors so far to predict the IMDb scores and later ended up only 4 of them, because the others were statistically insignificant (judging from the p-values). We have used the backward reduction based on p-values. In some cases, it may be observed that in this process, the coefficient of determination (r-square) has decreased after removing the insignificant factors from the model, but in our model, the adjusted r-square was initially very high and then went up by a small amount after reduction, indicating that the factors removed could have caused collinearity.

During the model selection we saw that the score has a low linear association with the run-time, but we still keep the factor in our model as it is statistically significant and we also see that the final model yields a pretty accurate prediction of the IMDb score, indicating that runtime is an important factor of our model.

Lastly, we can conclude by saying that there is a lot of research that can be done and the one explained in this project shows that we can accurately predict how people will rate (like) the movies from its general characteristics.