# Capstone: MovieLens

Saketha Male

8/21/2021

## Introduction

Recommendation system is a Machine Learning technique that provides users with suggestions of relevant products or services through the use of algorithms. Recommendation Systems help businesses and companies profoundly. In today's world, they exist in almost everything we use; for example, Netflix, Amazon, Spotify, Instagram, etc. These companies generate algorithms by collecting data and processing information on users' history, such as purchases and viewing activity.

The MovieLens 10M Dataset will be used in this report. This report will contain predictions based on user ratings, inferencing/modeling, exploratory analysis, comparisons to the validation set, and results derived from the MovieLens dataset.

## Methods

- Cleaning and Tidying Data

  - A segment of the code was provided by the course instructor. The code separates and cleans the data.

- Inferencing and modeling

  - The dimensions and different aspects of the movielens dataset are first introduced and made familiar through fragments of code.
  - The data analysis was depicted and supported with use of graphs, plots, and other visual features.

- Exploratory Analysis

  - Exploring the relationships between variables in the dataset through visuals such as graphs/plots helped outline the conclusions.

- Validation

  - An RMSE function was defined and used to compare the true ratings in the validation set to the predictions and models generated in this report.

## Overview of the dataset

```
# Number of rows
nrow(edx)
```

```
## [1] 9000061
```

```
# Number of columns
ncol(edx)
```

```
## [1] 6
```

```
# Column names
colnames(edx)
```

```
## [1] "userId"    "movieId"   "rating"    "timestamp" "title"     "genres"
```

```
# First 6 rows
head(edx)
```

```
##    userId movieId rating timestamp                         title
## 1:      1     122      5 838985046               Boomerang (1992)
## 2:      1     185      5 838983525               Net, The (1995)
## 3:      1     231      5 838983392           Dumb & Dumber (1994)
## 4:      1     292      5 838983421                Outbreak (1995)
## 5:      1     316      5 838983392                Stargate (1994)
## 6:      1     329      5 838983392 Star Trek: Generations (1994)
##                            genres
## 1:                 Comedy|Romance
## 2:           Action|Crime|Thriller
## 3:                         Comedy
## 4:   Action|Drama|Sci-Fi|Thriller
## 5:         Action|Adventure|Sci-Fi
## 6: Action|Adventure|Drama|Sci-Fi
```

```
# Summary of the dataset
summary(edx)
```

```
##      userId         movieId          rating        timestamp
##  Min.   :    1   Min.   :    1   Min.   :0.500   Min.   :7.897e+08
##  1st Qu.:18122   1st Qu.:  648   1st Qu.:3.000   1st Qu.:9.468e+08
##  Median :35743   Median : 1834   Median :4.000   Median :1.035e+09
##  Mean   :35869   Mean   : 4120   Mean   :3.512   Mean   :1.033e+09
##  3rd Qu.:53602   3rd Qu.: 3624   3rd Qu.:4.000   3rd Qu.:1.127e+09
##  Max.   :71567   Max.   :65133   Max.   :5.000   Max.   :1.231e+09
##     title             genres
##  Length:9000061     Length:9000061
```

```
##  Class :character   Class :character
##  Mode  :character   Mode  :character
##
##
##
```

This overview of the data helps us understand the dimensions and content of te dataset. We now know there are 9000055 rows and 6 columns in the edx dataset. Furthermore, we were also able to see the different columns. Finally, we are also able to see the range, quartiles, median, and mean of the different columns in the dataset.
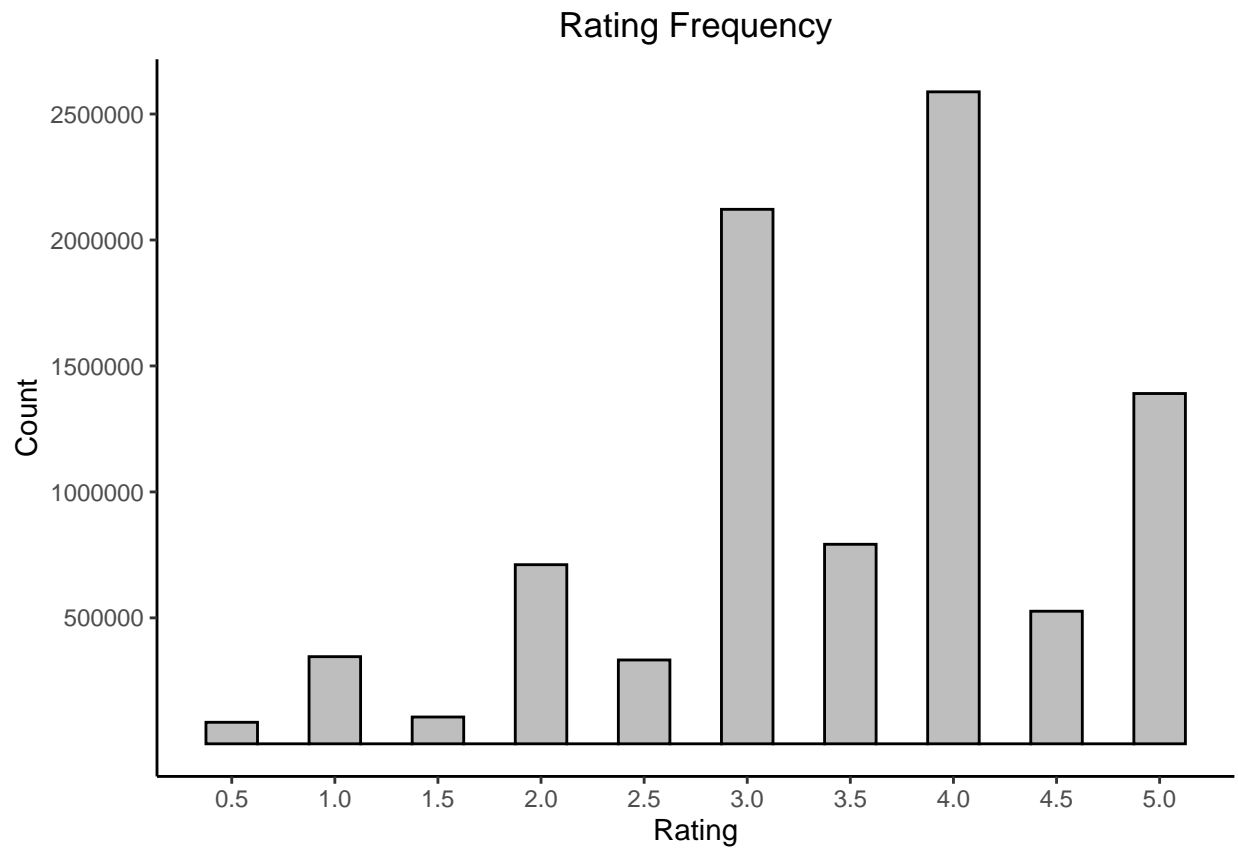
## Exploratory Analysis

Now we will explore the relationships between the variables in the edx dataset to recognize and identify relationships and trends.
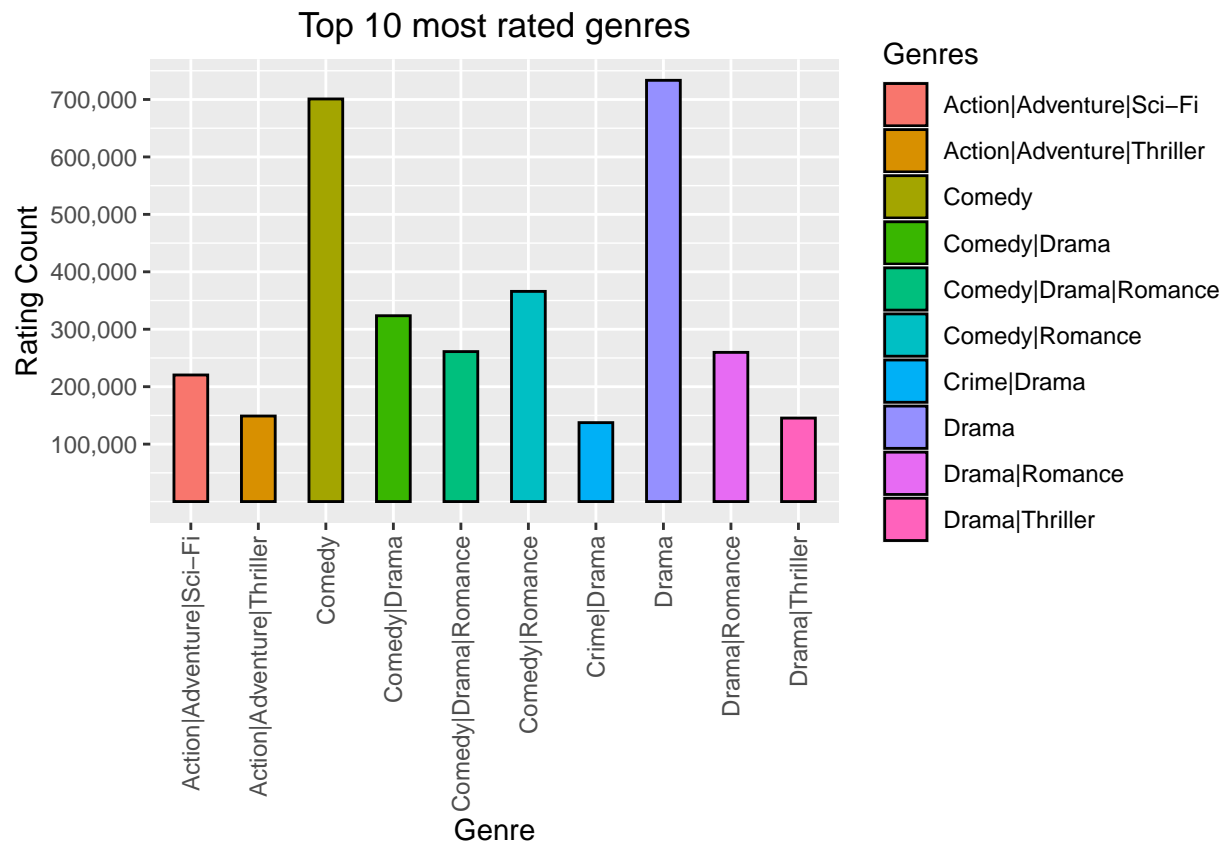
```r
# Loading necessary packages

library(ggplot2)
library(tidyr)
library(dplyr)
library(lattice)
library(grid)
library(lubridate)
library(ggthemes)
library(gridExtra)
library(scales)
```

**Frequency Of Each Rating**

## Rating Frequency



The most common rating given is a 4.0. We see that whole number ratings (1, 2, etc.) are more likely to be given than decimals (.5, 1.5, etc.). Furthermore, users are less likely to give a very high or low rating.
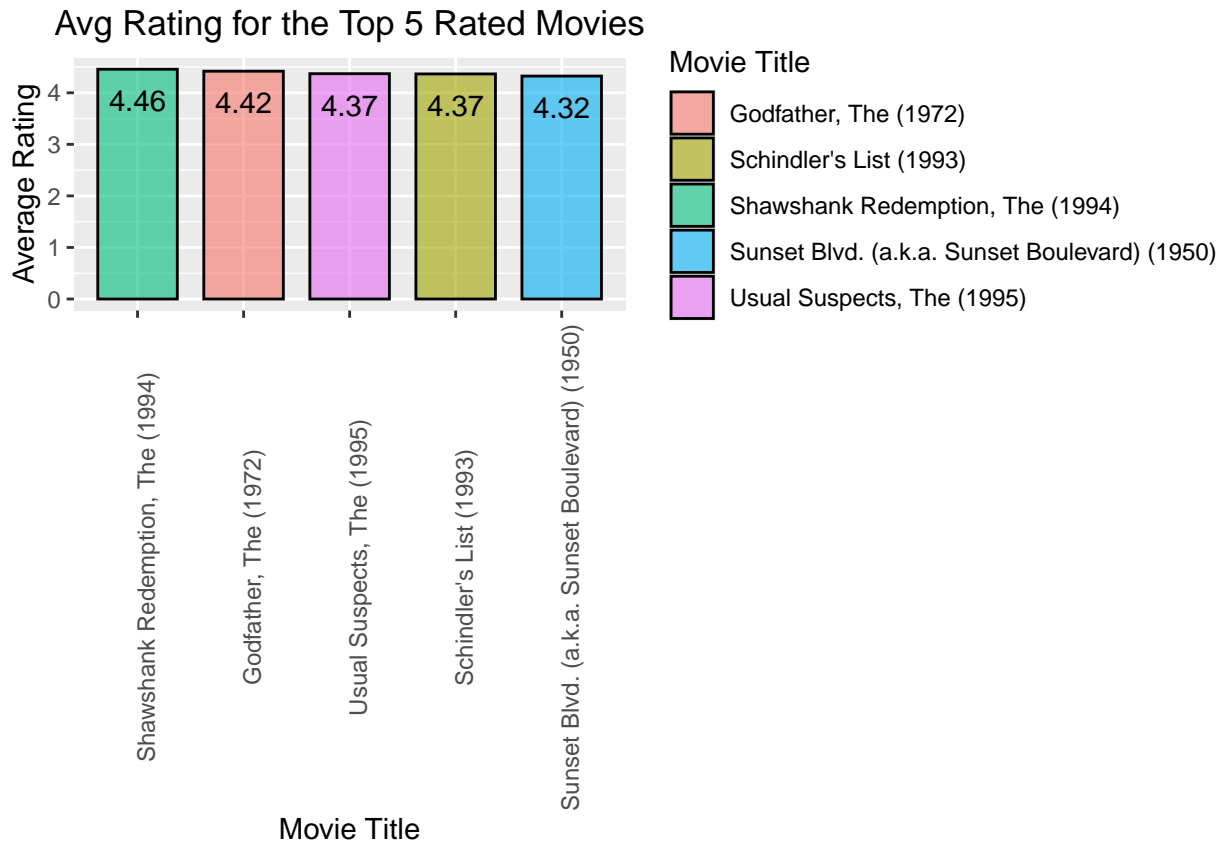
**The Top 10 Most Rated Genres**

## Top 10 most rated genres



The most popular genres are Drama and Comedy. Both of them exceed approximately 700,000 ratings. However, because many genre entires also consist of multiple other genres, the rating numbers are most likely higher.

**Barplots Of The Most & Least Rated Movies**

```
## `summarise()` has grouped output by 'title'. You can override using the `.groups` argument.
## `summarise()` has grouped output by 'title'. You can override using the `.groups` argument.
```
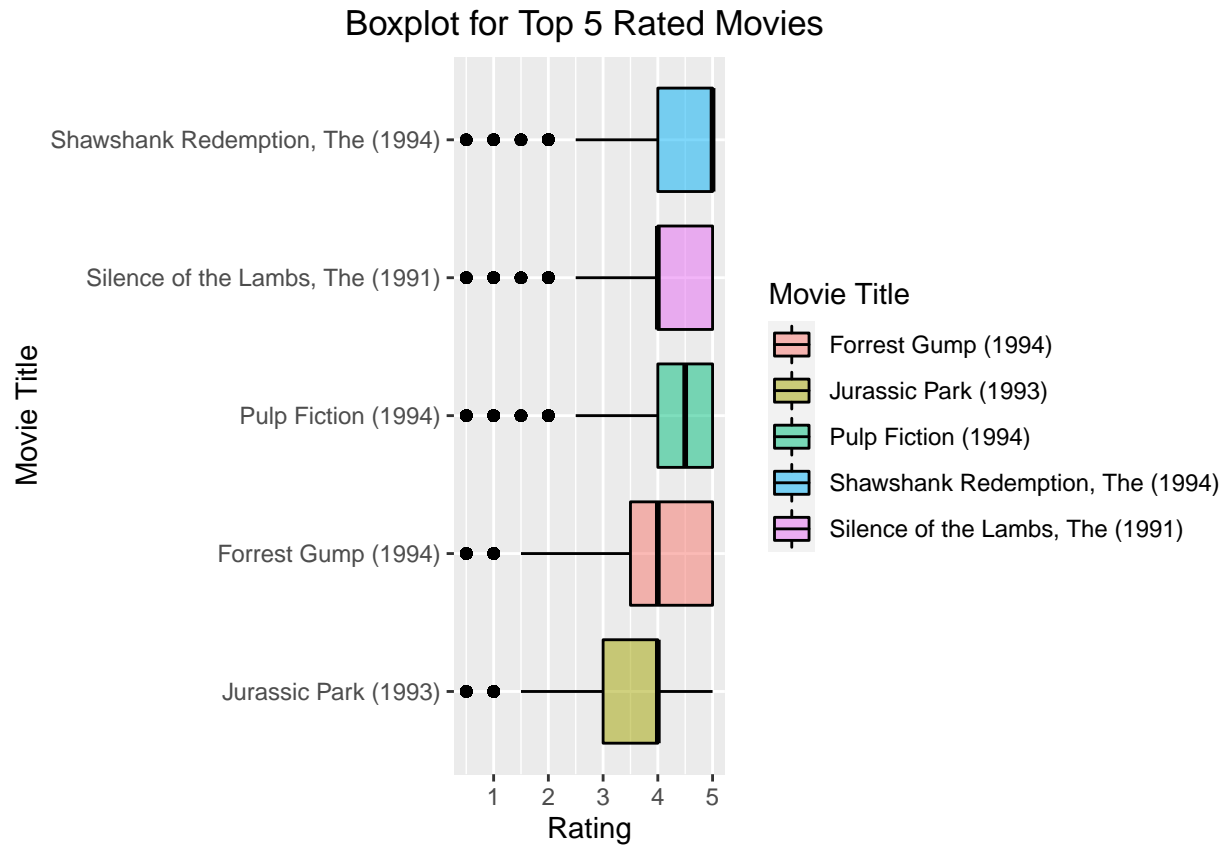
## Avg Rating for the Top 5 Rated Movies



**Movie Title**

- Godfather, The (1972)
- Schindler's List (1993)
- Shawshank Redemption, The (1994)
- Sunset Blvd. (a.k.a. Sunset Boulevard) (1950)
- Usual Suspects, The (1995)

## Avg Rating for the Bottom 5 Rated Movies



Average Rating

1.15    1.18    1.37    1.45    1.58

Glitter (2001)

Gigli (2003)

Turbo: A Power Rangers Movie (1997)

3 Ninjas: High Noon On Mega Mountain (1998)

Battlefield Earth (2000)

Movie Title

**Movie Title**

- 3 Ninjas: High Noon On Mega Mountain (1998)
- Battlefield Earth (2000)
- Gigli (2003)
- Glitter (2001)
- Turbo: A Power Rangers Movie (1997)

These plots show the top 5 and botttom 5 rated movies. The plot for top 5 has been filtered to have atleast 100 ratings, whereas the bottom 5 plot has been filtered to have atleast 300 ratings. This is to avoid external bias. The highest rated movie is "The Shawshank Redemption (1994)" and the least rated movie is "Glitter (2001)".

**Box Plots For The Most Rated Movies**
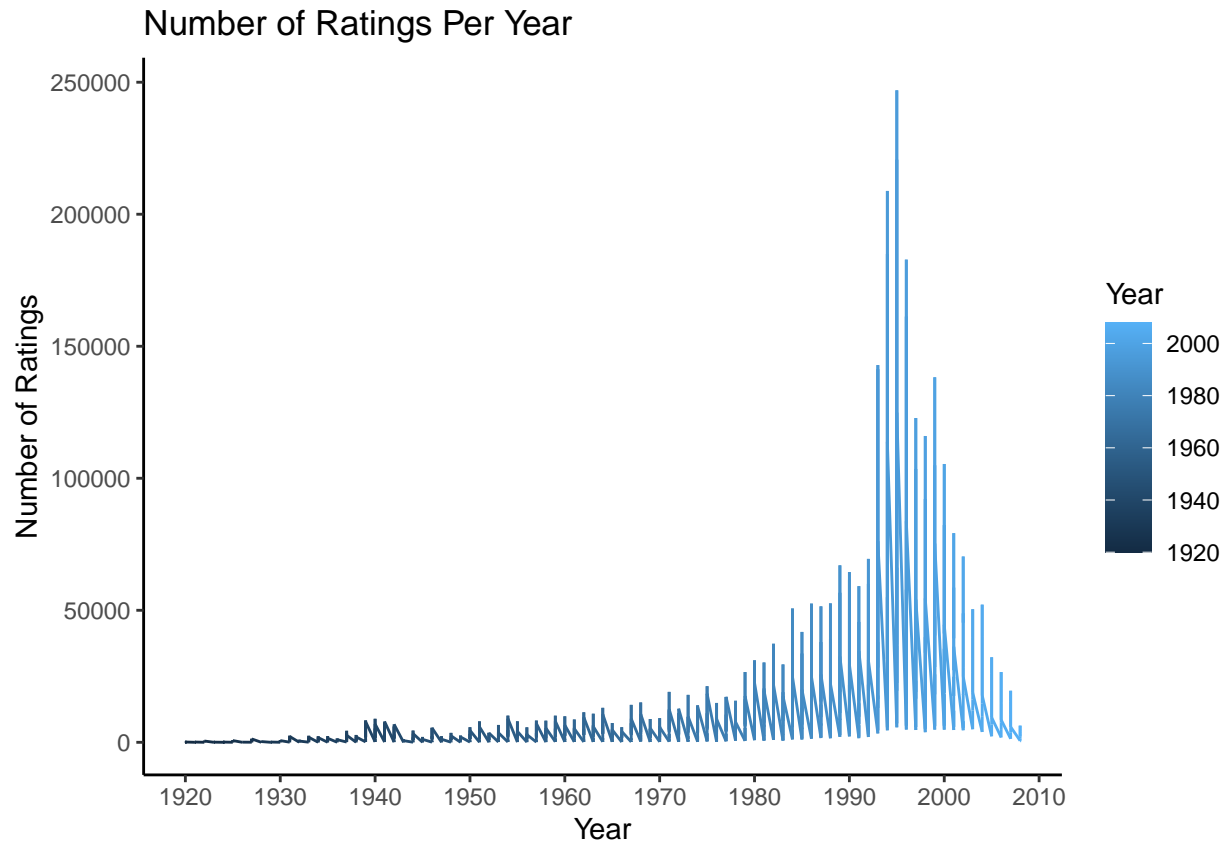
## Boxplot for Top 5 Rated Movies



Box plots give us a very good overview of data. These box plots depict how the most rated movies vary in averages, medians, ranges, and quartiles.

## Number Of Ratings Given Throughout The Years

## `summarise()` has grouped output by 'yearId'. You can override using the `.groups` argument.
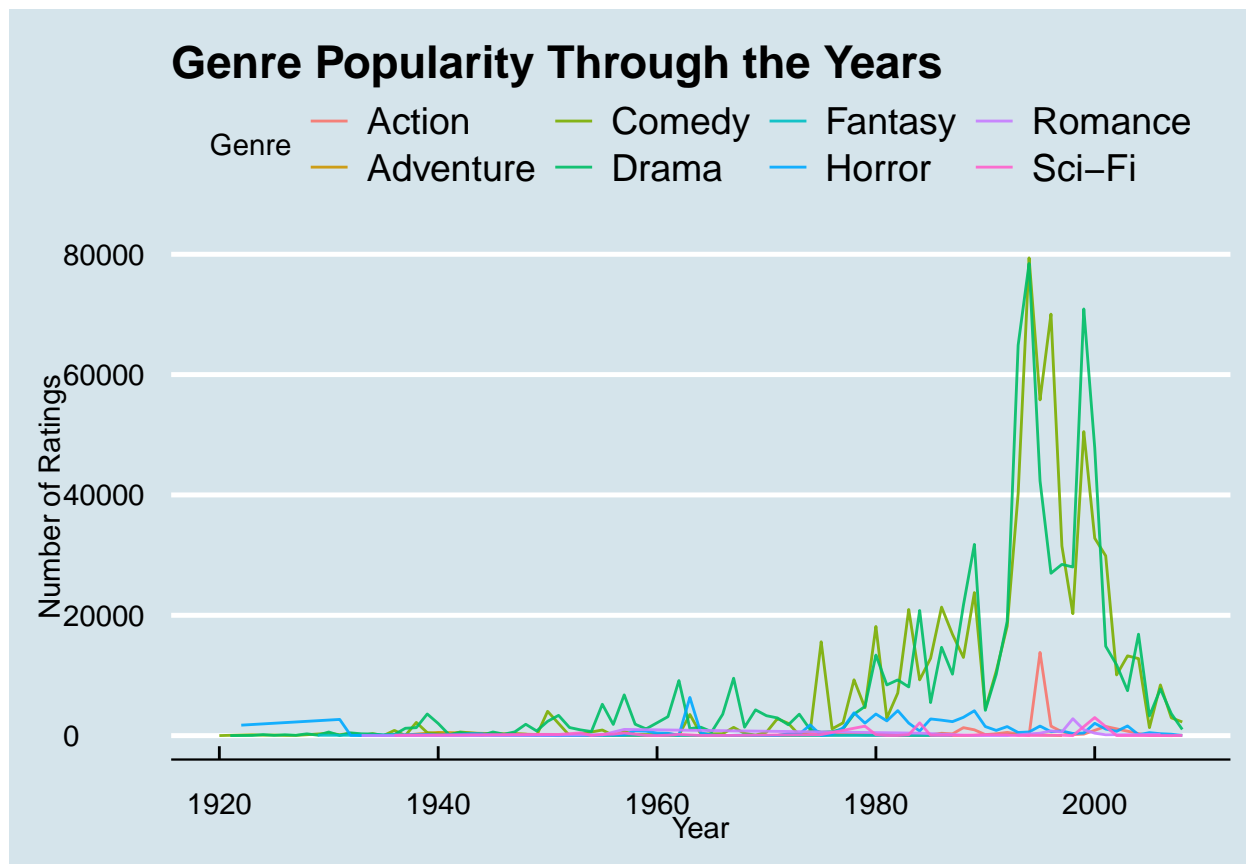
## Number of Ratings Per Year



For this graph, we had to specifically add another column by altering the title column values to extract the date each movie was released. We see that the ratings were given the most often between 1990 and 2000 but then declined dramatically.

## Genre Popularity Throughout The Years

```
## [1] "The 8 genres are Adventure, Action, Comedy, Drama, Fantatsy, Horror, Romance, Sci-Fi"
```

```
## `summarise()` has grouped output by 'yearId'. You can override using the `.groups` argument.
```

**Genre Popularity Through the Years**

For this plot, we used the sample function to choose 8 distinct genres from the dataset. The 8 Genres were then plotted with a different color representing each one. We see that throughout most of the years, comedy and drama continued to be the most popular. All genres mostly peaked to their highest point between 1990 and 2005, except for horror which peaked in the 1960s.

### Recommender systems simulation

Using the sample function once again, we have selected 3 randomly chosen users and are on a task of giving them a genre recommendation. This simualation is a very basic and uncomplex example compared to those other comapnies use. We do this by analyzing their past viewing preferences and the ratings they gave to each genre/movie.

```
## [1] "The 3 randomly chosen userIds are 28291, 58469, 21149"
```

- Analysis for User 28291:

```
## # A tibble: 54 x 3
##   genres               NumberOfRatings AvgRating
##   <chr>                          <int>     <dbl>
## 1 Drama                             31      2.87
## 2 Drama|Romance                     18      3.17
## 3 Drama|Thriller                     6      2.67
## 4 Comedy                             4      3
## 5 Action|Drama|Thriller              3      3
## 6 Action|Thriller                    3      3.33
```

```
##  7 Comedy|Drama                                3       2
##  8 Comedy|Drama|Romance                        3       2.67
##  9 Drama|War                                   3       4
## 10 Action|Romance|Thriller                     2       3
## # ... with 44 more rows


##       movieId                        title rating
##  1:      110          Braveheart (1995)      5
##  2:      527    Schindler's List (1993)      5
##  3:      587                 Ghost (1990)      5
##  4:      590 Dances with Wolves (1990)      5
##  5:     1286   Somewhere in Time (1980)      5
##  6:     1411                Hamlet (1996)      5
##  7:     1721                Titanic (1997)      5
##  8:       62 Mr. Holland's Opus (1995)      4
##  9:      150          Apollo 13 (1995)      4
## 10:      168        First Knight (1995)      4
##                                          genres
##  1:                          Action|Drama|War
##  2:                                 Drama|War
##  3: Comedy|Drama|Fantasy|Romance|Thriller
##  4:              Adventure|Drama|Western
##  5:                          Drama|Romance
##  6:                    Crime|Drama|Romance
##  7:                          Drama|Romance
##  8:                                     Drama
##  9:                         Adventure|Drama
## 10:                  Action|Drama|Romance
```

Based on the results, we recommend for user 28291: Movies in the Drama genre since the user watches that
genre the most, has the highest average rating in its sub category (Drama|Romance) when filtered for genres
that have been rated atleast 5 times, and was a genre in 7/7 of the movies the user rated the highest.

- Analysis for User 58469:

```
## # A tibble: 65 x 3
##    genres                        NumberOfRatings AvgRating
##    <chr>                                   <int>     <dbl>
##  1 Comedy                                     10       4.2
##  2 Drama                                       5       3.8
##  3 Action|Adventure|Sci-Fi                     3       3.33
##  4 Action|Crime|Thriller                       3       4
##  5 Drama|Thriller                              3       3
##  6 Action|Adventure|Thriller                   2       3.5
##  7 Action|Sci-Fi                               2       4.5
##  8 Action|Thriller                             2       4
##  9 Comedy|Crime|Drama                          2       5
## 10 Comedy|Fantasy|Romance|Sci-Fi               2       4
## # ... with 55 more rows


##     movieId                        title rating               genres
## 1:       62        Mr. Holland's Opus (1995)      5                Drama
## 2:      110              Braveheart (1995)      5    Action|Drama|War
```

```
## 3:     296                    Pulp Fiction (1994)    5 Comedy|Crime|Drama
## 4:     318  Shawshank Redemption, The (1994)    5              Drama
## 5:     333                       Tommy Boy (1995)    5             Comedy
## 6:     344  Ace Ventura: Pet Detective (1994)    5             Comedy
```

Based on the results, we recommend for user 58469: Movies in the Comedy genre since the user watches that genre the most, has the highest average rating when filtered for genres that have been rated atleast 3 times, and was a genre for 7/13 movies the user gave a 5 rating.

- Analysis for User 21149:

```
## # A tibble: 237 x 3
##    genres              NumberOfRatings AvgRating
##    <chr>                         <int>     <dbl>
##  1 Drama                           116      3.82
##  2 Comedy                           77      3.24
##  3 Comedy|Drama                     38      3.5
##  4 Comedy|Romance                   37      3.30
##  5 Drama|Thriller                   28      3.66
##  6 Drama|Romance                    24      3.71
##  7 Crime|Drama                      21      3.62
##  8 Crime|Drama|Thriller             20      3.85
##  9 Comedy|Drama|Romance             18      3.47
## 10 Drama|War                        17      3.85
## # ... with 227 more rows
```

Based solely on the table for genres and average rating, we recommend for user 21149: movies in the Drama and Comedy genre since those are the genres watched most by the user and those are the genres with the highest average rating.

Additionally, to provide better recommendations, we can look for movies that have the highest ratings in their respective genres and recommend those movies to users that enjoy that genre the most.

# Validation

Defining the RMSE function

```
rmse <- function(validation, y_hat){
  sqrt(mean((validation-y_hat)^2))
}
```

## Analysis/Accuracy of the report's models will now be assesed.

- Model 1: Ratings predicted/approximated by the mean ratings in the edx dataset

```
## [1] 3.512464
```

```
## [1] 1.060651
```

- Model 2: Ratings predicted/approximated based on individual genres; some genres are rated higher/lower than others.

```
## [1] 1.018114
```

- Model 3: Ratings predicted/approximated based on individual movies; some movies are rated higher/lower than others.

```
## [1] 0.9437046
```

- Model 4: Ratings predicted/approximated based on individual users (adding onto Model 3); some users give higher/lower ratings for specific movies than do others.

```
## [1] 0.8655329
```

The results reveal that the 4th model is the best predictor of the 4 models. However, to achieve a better result, we must account for the fact that some movies are rated more or less often than others are. Similarly, some users rate movies more or less often than others do and more or less critically than others do.. In order to produce more accurate results, we need to put a penalty to limit the impact such occurrences have.

```r
# Refining model 4

# Adding a number of ratings per movie (numberofratingspermovie = norpm)

lambdas <- seq(0, 10, 0.25)

ratingspermovie <- edx %>% group_by(movieId) %>% summarize(total = sum(rating - mu), norpm = n())

finalrmsevalues <- sapply(lambdas, function(l) {
  model5 <- validation %>% left_join(ratingspermovie, by = "movieId") %>% left_join(userratingavgs, by =
  return(rmse(validation$rating, model5$prediction))

})
```
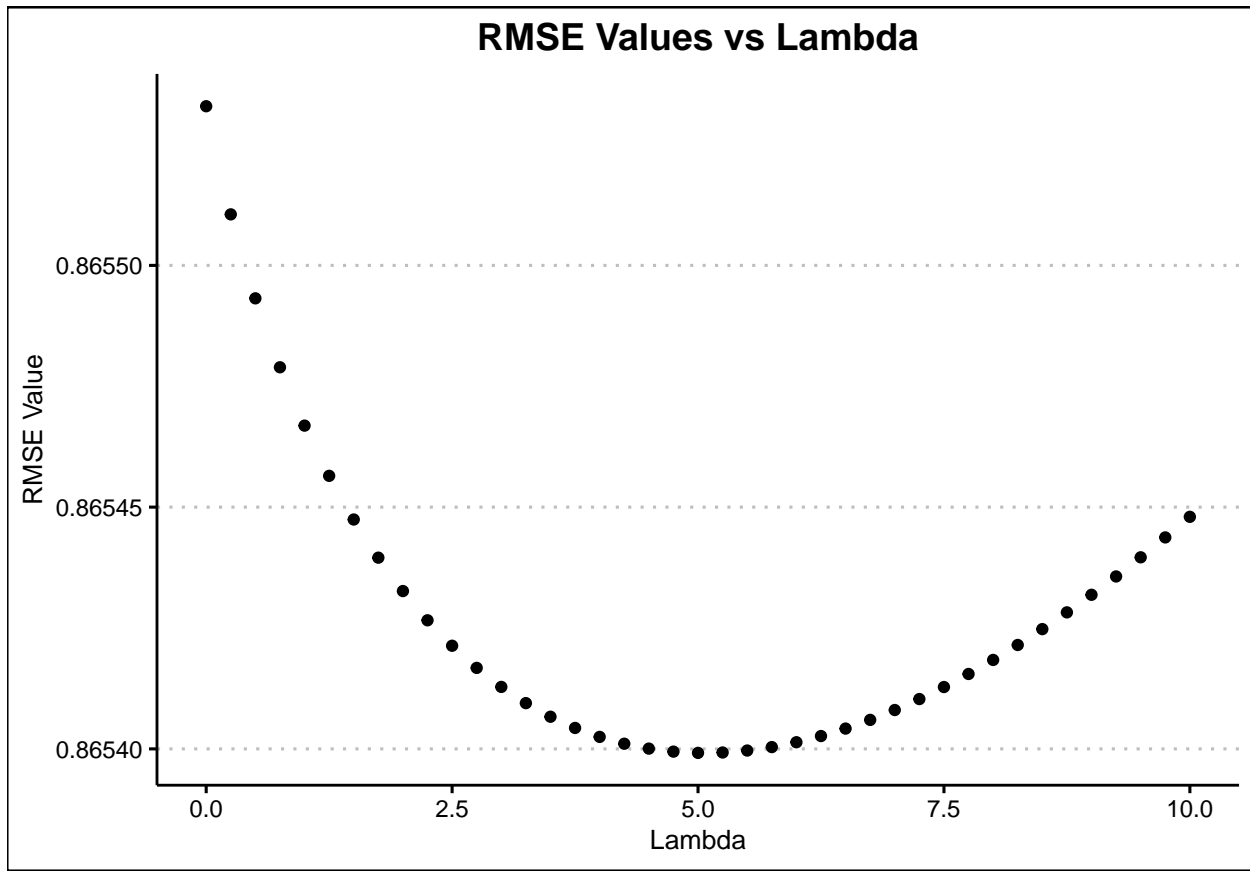
## RMSE Values Plot

```r
rmseplot <- qplot(lambdas, finalrmsevalues, main = "RMSE Values vs Lambda", xlab = "Lambda", ylab = "RMS

rmseplot
```

**RMSE Values vs Lambda**



```
# Final RMSE value

min(finalrmsevalues)
```

```
## [1] 0.8653992
```

```
# Lambda Value at which it occurs:

lambdas[which.min(finalrmsevalues)]
```

```
## [1] 5
```

## Concluding Remarks

In this report, we have profoundly explored the dataset and identified relationships among variables and columns in the datatset. In addition, we have also added a year variable to identify trends throughout time. We have simulated a basic recommender system for 3 randomly chosen users. FInally, we chose model 4, which predicted ratings based on users and the movies they've watched, to see how accurate or true our models are to the true ratings in the validation set. We have arrived at an RMSE value of .8652069 that occurs at lambda 4.75. Although not ideal, of the 4 models this was the minimized value. In future reports, stricter limits and parameters can be added to bring this value even lower.

The projects goal was to analyze the MovieLens dataset, present the analysis through tables and graphs, and minimize the RMSE loss function of the true ratings in the validation dataset.