

# Report

Abrham Birru

12/24/2020

## MovieLens Rating Project Code

### Introduction

This data analysis report is prepared as a part of HarvardX Data Science Capstone (HarvardX: PH125.9x) Project. The given dataset (movielens) has been analysed by using different tools and techniques.

In this project, I will combine several machine learning strategies to construct a movie recommendation system based on the “MovieLens” dataset.

### Goal

The goal is to predict movie ratings, and evaluate the accuracy of the predicted model from the given code the dataset called “edx” was split into the training and validation sets.

### Data Loading

The code is provided in the edx capstone project module: The following lines of code will create training and validation sets

### Create edx set, validation set (final hold-out test set)

Note: this process could take a couple of minutes

```
if(!require(tidyverse)) install.packages("tidyverse", repos = "http://cran.us.r-project.org")
```

```
## Loading required package: tidyverse
```

```
## -- Attaching packages -----
```

```
## v ggplot2 3.3.2      v purrr   0.3.4
## v tibble  3.0.3      v dplyr   1.0.2
## v tidyr   1.1.1      v stringr 1.4.0
## v readr   1.3.1      v forcats 0.5.0
```

```
## -- Conflicts -----
```

```
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```

if(!require(caret)) install.packages("caret", repos = "http://cran.us.r-project.org")

## Loading required package: caret

## Loading required package: lattice

##
## Attaching package: 'caret'

## The following object is masked from 'package:purrr':
##
## lift

if(!require(data.table)) install.packages("data.table", repos = "http://cran.us.r-project.org")

## Loading required package: data.table

##
## Attaching package: 'data.table'

## The following objects are masked from 'package:dplyr':
##
## between, first, last

## The following object is masked from 'package:purrr':
##
## transpose

library(tidyverse)
library(caret)
library(data.table)

# MovieLens 10M dataset:
# https://grouplens.org/datasets/movielens/10m/
# http://files.grouplens.org/datasets/movielens/ml-10m.zip

dl <- tempfile()
download.file("http://files.grouplens.org/datasets/movielens/ml-10m.zip", dl)

ratings <- fread(text = gsub(":", "\t", readLines(unzip(dl, "ml-10M100K/ratings.dat"))),
  col.names = c("userId", "movieId", "rating", "timestamp"))

movies <- str_split_fixed(readLines(unzip(dl, "ml-10M100K/movies.dat")), "\\:", 3)
colnames(movies) <- c("movieId", "title", "genres")

movies <- as.data.frame(movies) %>% mutate(movieId = as.numeric(levels(movieId))[movieId],
  title = as.character(title),
  genres = as.character(genres))

```

```

movielens <- left_join(ratings, movies, by = "movieId")

# Validation set will be 10% of MovieLens data
set.seed(1, sample.kind="Rounding")

## Warning in set.seed(1, sample.kind = "Rounding"): non-uniform 'Rounding' sampler
## used

test_index <- createDataPartition(y = movielens$rating, times = 1, p = 0.1, list = FALSE)
edx <- movielens[-test_index,]
temp <- movielens[test_index,]

# Make sure userId and movieId in validation set are also in edx set
validation <- temp %>%
  semi_join(edx, by = "movieId") %>%
  semi_join(edx, by = "userId")

# Add rows removed from validation set back into edx set
removed <- anti_join(temp, validation)

## Joining, by = c("userId", "movieId", "rating", "timestamp", "title", "genres")

edx <- rbind(edx, removed)
rm(dl, ratings, movies, test_index, temp, movielens, removed)

```

## Exploring The Data

```
head(edx)
```

```

##      userId movieId rating timestamp                title
## 1:      1      122      5 838985046      Boomerang (1992)
## 2:      1      185      5 838983525      Net, The (1995)
## 3:      1      292      5 838983421      Outbreak (1995)
## 4:      1      316      5 838983392      Stargate (1994)
## 5:      1      329      5 838983392 Star Trek: Generations (1994)
## 6:      1      355      5 838984474      Flintstones, The (1994)
##                                     genres
## 1:                                Comedy|Romance
## 2:                        Action|Crime|Thriller
## 3: Action|Drama|Sci-Fi|Thriller
## 4:                        Action|Adventure|Sci-Fi
## 5: Action|Adventure|Drama|Sci-Fi
## 6:                   Children|Comedy|Fantasy

```

```
head(validation)
```

```

##      userId movieId rating timestamp
## 1:      1      231      5 838983392
## 2:      1      480      5 838983653

```

```
## 3:      1      586      5 838984068
## 4:      2      151      3 868246450
## 5:      2      858      2 868245645
## 6:      2     1544      3 868245920
##
##                                     title
## 1:                                     Dumb & Dumber (1994)
## 2:                                     Jurassic Park (1993)
## 3:                                     Home Alone (1990)
## 4:                                     Rob Roy (1995)
## 5:                                     Godfather, The (1972)
## 6: Lost World: Jurassic Park, The (Jurassic Park 2) (1997)
##                                     genres
## 1:                                     Comedy
## 2:      Action|Adventure|Sci-Fi|Thriller
## 3:                                     Children|Comedy
## 4:      Action|Drama|Romance|War
## 5:                                     Crime|Drama
## 6: Action|Adventure|Horror|Sci-Fi|Thriller
```

```
str(edx)
```

```
## Classes 'data.table' and 'data.frame':  9000055 obs. of  6 variables:
## $ userId   : int  1 1 1 1 1 1 1 1 1 1 ...
## $ movieId  : num  122 185 292 316 329 355 356 362 364 370 ...
## $ rating   : num  5 5 5 5 5 5 5 5 5 5 ...
## $ timestamp: int  838985046 838983525 838983421 838983392 838983392 838984474 838983653 838984885 8...
## $ title    : chr  "Boomerang (1992)" "Net, The (1995)" "Outbreak (1995)" "Stargate (1994)" ...
## $ genres   : chr  "Comedy|Romance" "Action|Crime|Thriller" "Action|Drama|Sci-Fi|Thriller" "Action|A...
## - attr(*, ".internal.selfref")=<externalptr>
```

```
str(validation)
```

```
## Classes 'data.table' and 'data.frame':  999999 obs. of  6 variables:
## $ userId   : int  1 1 1 2 2 2 3 3 4 4 ...
## $ movieId  : num  231 480 586 151 858 ...
## $ rating   : num  5 5 5 3 2 3 3.5 4.5 5 3 ...
## $ timestamp: int  838983392 838983653 838984068 868246450 868245645 868245920 1136075494 1133571200...
## $ title    : chr  "Dumb & Dumber (1994)" "Jurassic Park (1993)" "Home Alone (1990)" "Rob Roy (1995)"...
## $ genres   : chr  "Comedy" "Action|Adventure|Sci-Fi|Thriller" "Children|Comedy" "Action|Drama|Roman...
## - attr(*, ".internal.selfref")=<externalptr>
```

```
names(edx)
```

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

```
names(validation)
```

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

## summary of unique movies and users

```
summary(edx)
```

```
##      userId      movieId      rating      timestamp
## Min.   :    1   Min.   :    1   Min.   :0.500   Min.   :7.897e+08
## 1st Qu.:18124   1st Qu.:   648   1st Qu.:3.000   1st Qu.:9.468e+08
## Median :35738   Median :  1834   Median :4.000   Median :1.035e+09
## Mean   :35870   Mean   :   4122   Mean   :3.512   Mean   :1.033e+09
## 3rd Qu.:53607   3rd Qu.:  3626   3rd Qu.:4.000   3rd Qu.:1.127e+09
## Max.   :71567   Max.   :65133   Max.   :5.000   Max.   :1.231e+09
##      title      genres
## Length:9000055   Length:9000055
## Class :character Class :character
## Mode  :character Mode  :character
##
##
##
```

## Number of unique movies and users in dataset

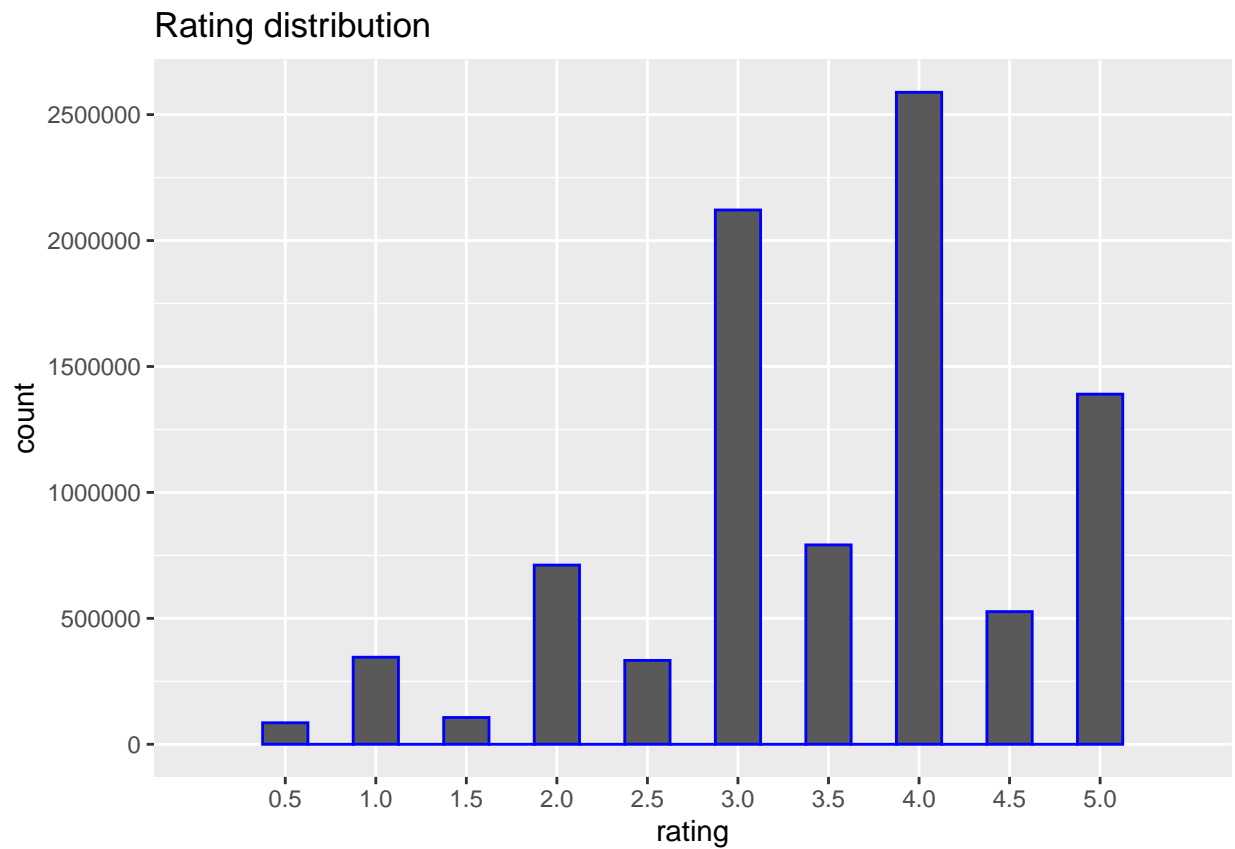
```
edx %>%
  summarize(number_of_users = n_distinct(userId),
            number_of_movies = n_distinct(movieId))
```

```
##   number_of_users number_of_movies
## 1             69878             10677
```

## Plot Rating distribution

```
edx %>%
  ggplot(aes(rating)) +
  geom_histogram(binwidth = 0.25, color = "blue") +
  scale_x_discrete(limits = c(seq(0.5,5,0.5))) +
  scale_y_continuous(breaks = c(seq(0, 3000000, 500000))) +
  ggtitle("Rating distribution")
```

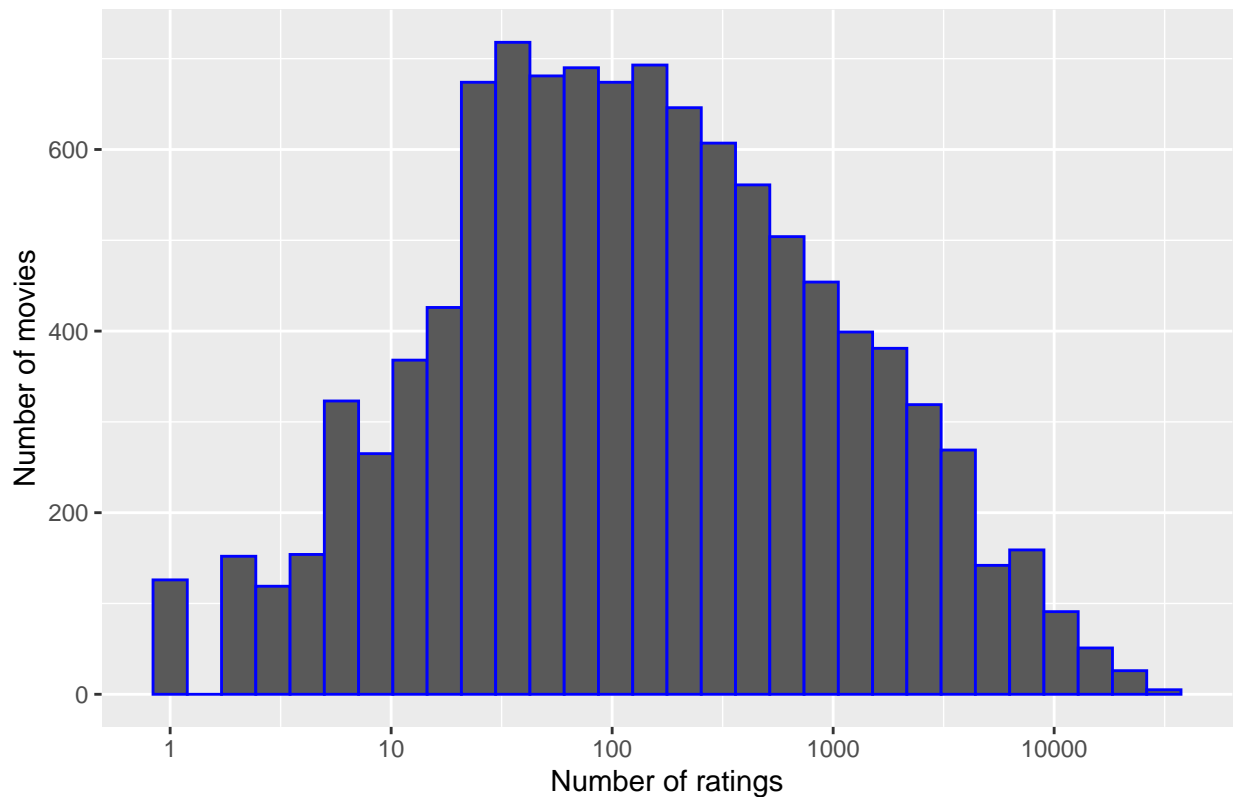
```
## Warning: Continuous limits supplied to discrete scale.
## Did you mean 'limits = factor(...)' or 'scale*_continuous()'?
```



# Make a plot of number of ratings per movie in edx dataset

```
edx %>%  
  count(movieId) %>%  
  ggplot(aes(n)) +  
  geom_histogram(bins = 30, color = "blue") +  
  scale_x_log10() +  
  xlab("Number of ratings") +  
  ylab("Number of movies") +  
  ggtitle("Number of ratings per movie")
```

Number of ratings per movie



# only once rated 20 movies table

```
edx %>%
  group_by(movieId) %>%
  summarize(count = n()) %>%
  filter(count == 1) %>%
  left_join(edx, by = "movieId") %>%
  group_by(title) %>%
  summarize(rating = rating, n_rating = count) %>%
  slice(1:20) %>%
  knitr::kable()
```

```
## 'summarise()' ungrouping output (override with '.groups' argument)
## 'summarise()' ungrouping output (override with '.groups' argument)
```

title	rating	n_rating
1, 2, 3, Sun (Un, deuz, trois, soleil) (1993)	2.0	1
100 Feet (2008)	2.0	1
4 (2005)	2.5	1
Accused (Anklaget) (2005)	0.5	1
Ace of Hearts (2008)	2.0	1
Ace of Hearts, The (1921)	3.5	1
Adios, Sabata (Indio Black, sai che ti dico: Sei un gran figlio di...) (1971)	1.5	1
Africa addio (1966)	3.0	1
Aleksandra (2007)	3.0	1

title	rating	n_rating
Bad Blood (Mauvais sang) (1986)	4.5	1
Battle of Russia, The (Why We Fight, 5) (1943)	3.5	1
Bellissima (1951)	4.0	1
Big Fella (1937)	3.0	1
Black Tights (1-2-3-4 ou Les Collants noirs) (1960)	3.0	1
Blind Shaft (Mang jing) (2003)	2.5	1
Blue Light, The (Das Blaue Licht) (1932)	5.0	1
Borderline (1950)	3.0	1
Brothers of the Head (2005)	2.5	1
Chapayev (1934)	1.5	1
Cold Sweat (De la part des copains) (1970)	2.5	1

## Make a plot of mean movie ratings given by users

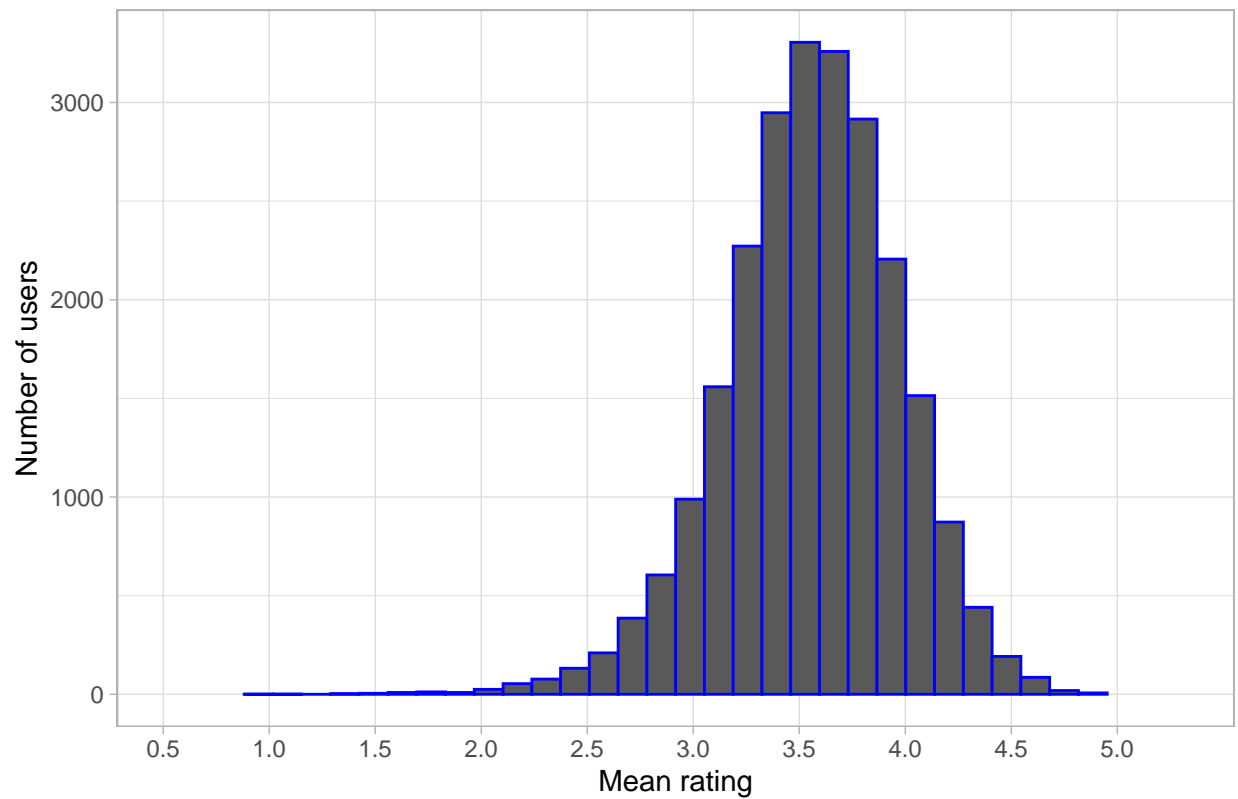
```
edx %>%
  group_by(userId) %>%
  filter(n() >= 100) %>%
  summarize(b_u = mean(rating)) %>%
  ggplot(aes(b_u)) +
  geom_histogram(bins = 30, color = "blue") +
  xlab("Mean rating") +
  ylab("Number of users") +
  ggtitle("Mean movie ratings given by users") +
  scale_x_discrete(limits = c(seq(0.5,5,0.5))) +
  theme_light()

## 'summarise()' ungrouping output (override with '.groups' argument)

## Warning: Continuous limits supplied to discrete scale.
## Did you mean 'limits = factor(...)' or 'scale*_continuous()'?
```



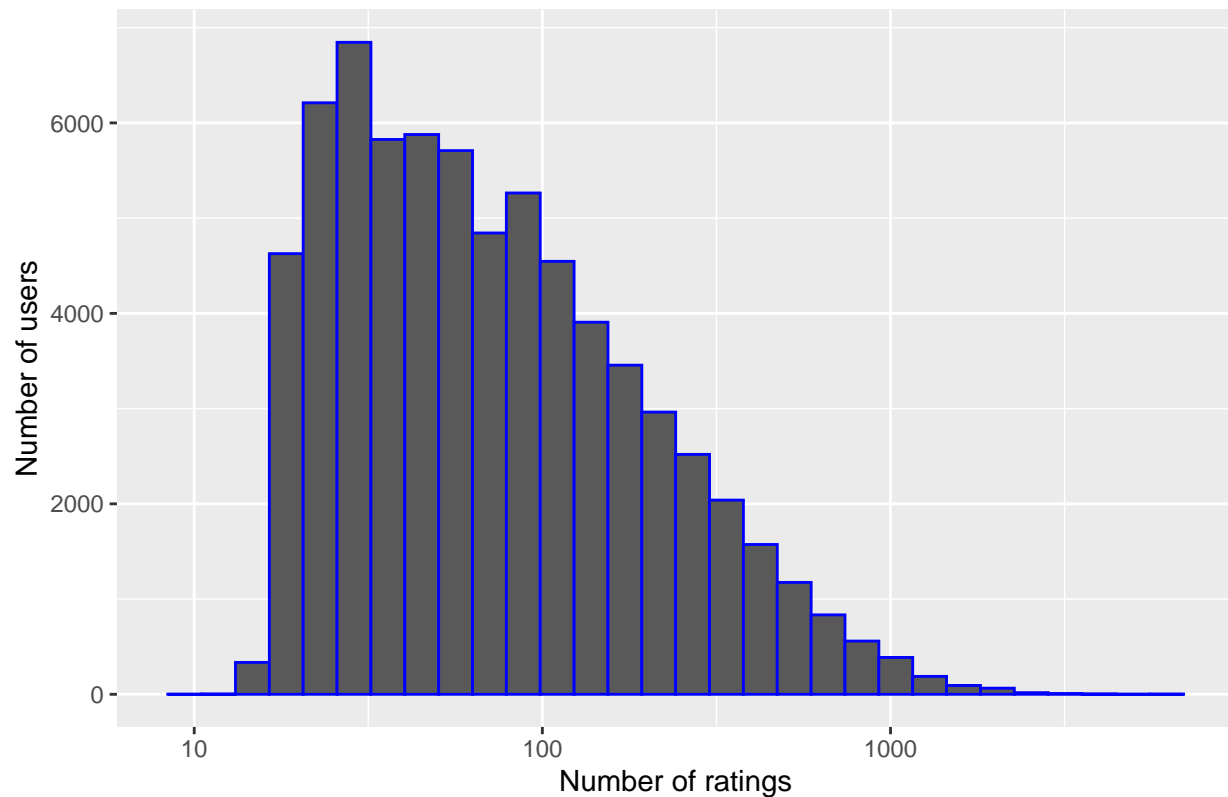
Mean movie ratings given by users



# Make a plot of number of ratings given by users

```
edx %>%  
  count(userId) %>%  
  ggplot(aes(n)) +  
  geom_histogram(bins = 30, color = "blue") +  
  scale_x_log10() +  
  xlab("Number of ratings") +  
  ylab("Number of users") +  
  ggtitle("Number of ratings given by users")
```

Number of ratings given by users



```
## Modelling
```

## Model-1

### Let's take the average(mean)

This model is developed by using average rating (mean) to train the model and predict the movie rating in the validation model.

```
average <- mean(edx$rating)
average
```

```
## [1] 3.512465
```

### Initiate RMSE and store based on simple prediction

```
naive_rmse <- RMSE(validation$rating, average)
naive_rmse
```

```
## [1] 1.061202
```

```
rmse_data <- data_frame(method = "Average Movie Rating Model", RMSE = naive_rmse)
```

```
## Warning: 'data_frame()' is deprecated as of tibble 1.1.0.
## Please use 'tibble()' instead.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_warnings()' to see where this warning was generated.
```

```
rmse_data %>% knitr::kable()
```

method	RMSE
Average Movie Rating Model	1.061202

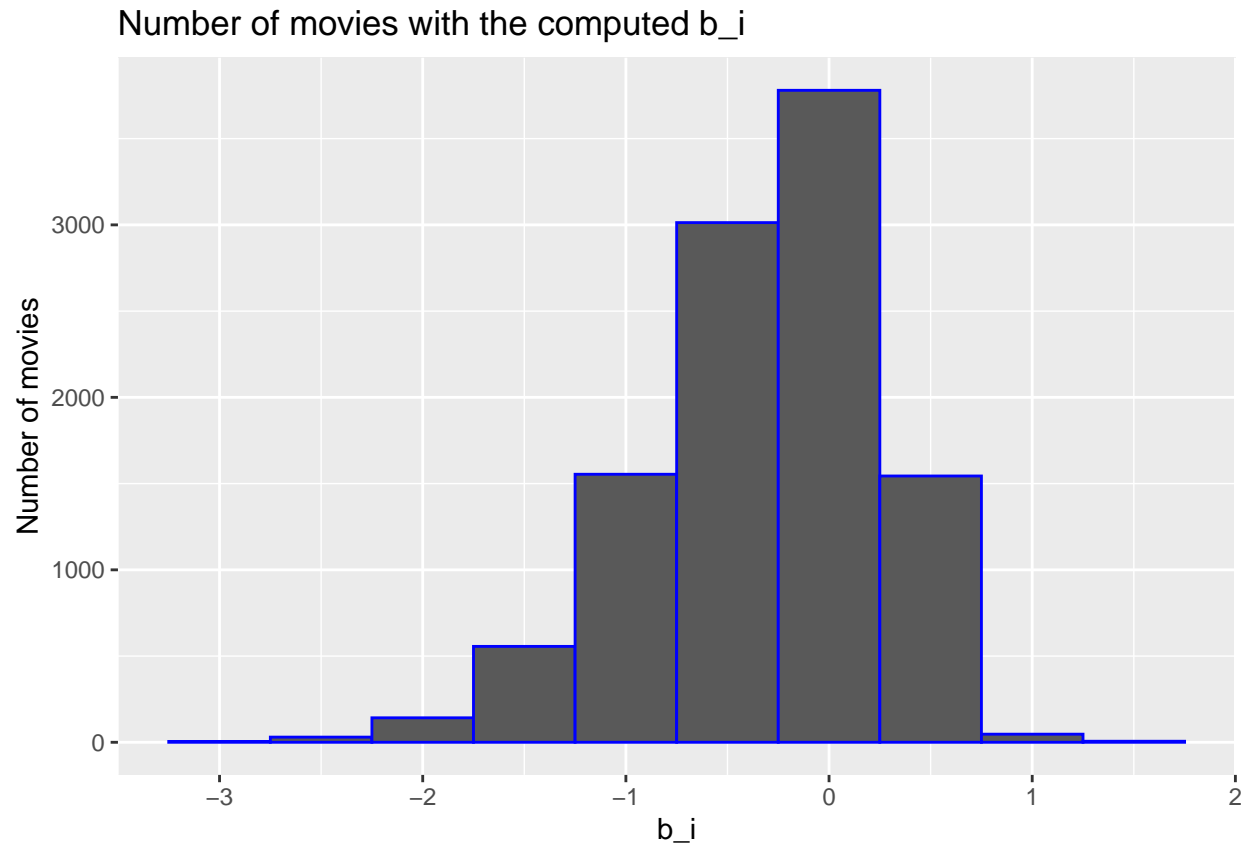
## Model-2

In this model the Movies are used to examine the predictive efficiency Simple model taking into account the movie effect  $b_i$  Subtract the rating minus the mean for each rating the movie received Plot number of movies with the computed  $b_i$

```
movie_avgs <- edx %>%
  group_by(movieId) %>%
  summarize(b_i = mean(rating - average))
```

```
## 'summarise()' ungrouping output (override with '.groups' argument)
```

```
movie_avgs %>% qplot(b_i, geom = "histogram", bins = 10, data = ., color = I("blue"),
  ylab = "Number of movies", main = "Number of movies with the computed b_i")
```



# Store rmse data results

```
predicted_ratings <- average + validation %>%
  left_join(movie_avgs, by='movieId') %>%
  pull(b_i)
model_1_rmse <- RMSE(predicted_ratings, validation$rating)
rmse_data <- bind_rows(rmse_data,
  data_frame(method="Movie Effect Model",
    RMSE = model_1_rmse ))
```

## Finally Check rmse data results

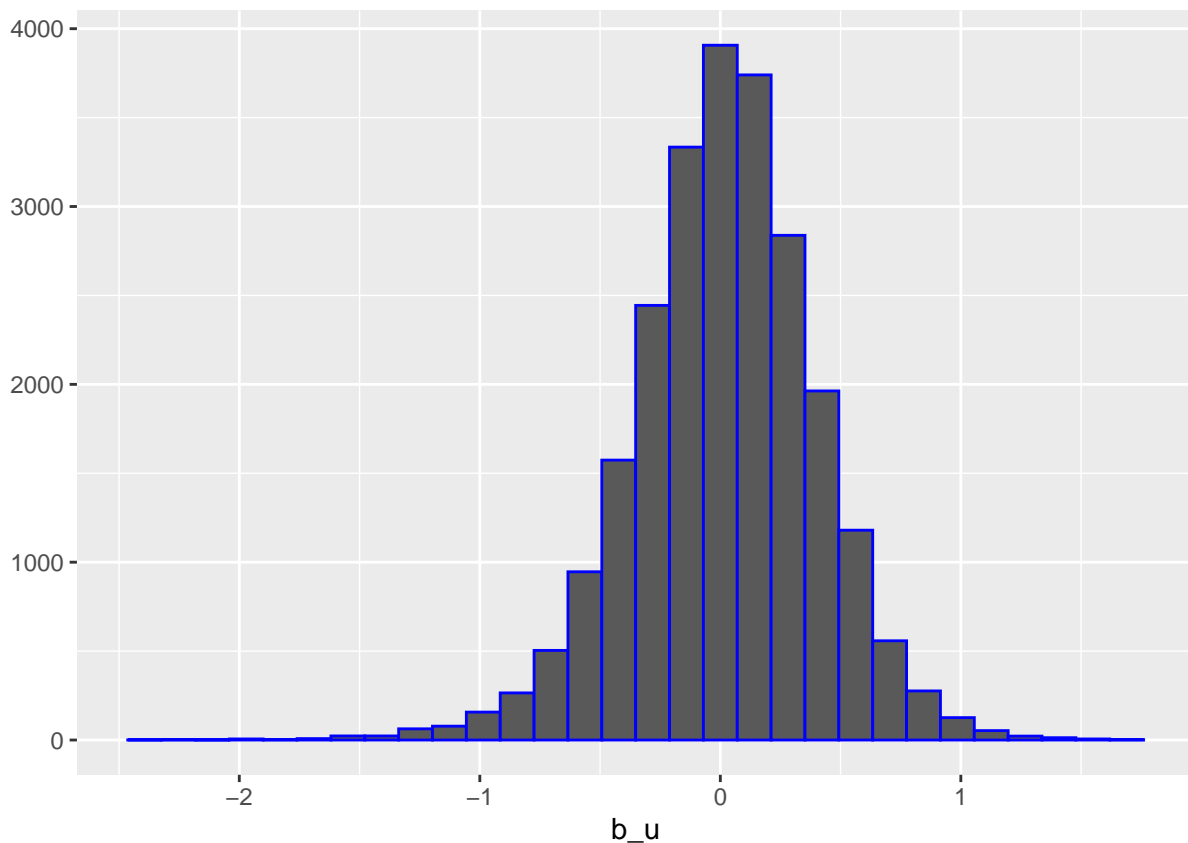
```
rmse_data %>% knitr::kable()
```

method	RMSE
Average Movie Rating Model	1.0612018
Movie Effect Model	0.9439087

## Model-3

Here we will model the Movie-User effects

```
user_avgs<- edx %>%  
  left_join(movie_avgs, by='movieId') %>%  
  group_by(userId) %>%  
  filter(n() >= 100) %>%  
  summarize(b_u = mean(rating - average - b_i))  
  
## 'summarise()' ungrouping output (override with '.groups' argument)  
  
user_avgs%>% qplot(b_u, geom = "histogram", bins = 30, data = ., color = I("blue"))
```



## User averages group by user id

```
user_avgs <- edx %>%  
  left_join(movie_avgs, by='movieId') %>%  
  group_by(userId) %>%  
  summarize(b_u = mean(rating - average - b_i))  
  
## 'summarise()' ungrouping output (override with '.groups' argument)
```

## Test and save rmse data results

```
predicted_ratings <- validation%>%
  left_join(movie_avgs, by='movieId') %>%
  left_join(user_avgs, by='userId') %>%
  mutate(pred = average + b_i + b_u) %>%
  pull(pred)

model_2_rmse <- RMSE(predicted_ratings, validation$rating)
rmse_data <- bind_rows(rmse_data,
  data_frame(method="Movie-User Effect Model",
    RMSE = model_2_rmse))
```

## Check rmse data result

```
rmse_data %>% knitr::kable()
```

method	RMSE
Average Movie Rating Model	1.0612018
Movie Effect Model	0.9439087
Movie-User Effect Model	0.8653488

## Model-4

Here Cross validation is used to select optimum value of lambda. It is a regularized model.

```
lambdas <- seq(0, 10, 0.25)
```

For each lambda, find  $b_i$  &  $b_u$ , and followed by rating prediction and testing

```
rmses <- sapply(lambdas, function(l){
  average <- mean(edx$rating)

  b_i <- edx %>%
    group_by(movieId) %>%
    summarize(b_i = sum(rating - average)/(n()+1))

  b_u <- edx %>%
    left_join(b_i, by="movieId") %>%
```

```
group_by(userId) %>%
  summarize(b_u = sum(rating - b_i - average)/(n()+1))

predicted_ratings <-
  validation %>%
  left_join(b_i, by = "movieId") %>%
  left_join(b_u, by = "userId") %>%
  mutate(pred = average + b_i + b_u) %>%
  pull(pred)

return(RMSE(predicted_ratings, validation$rating))
})
```

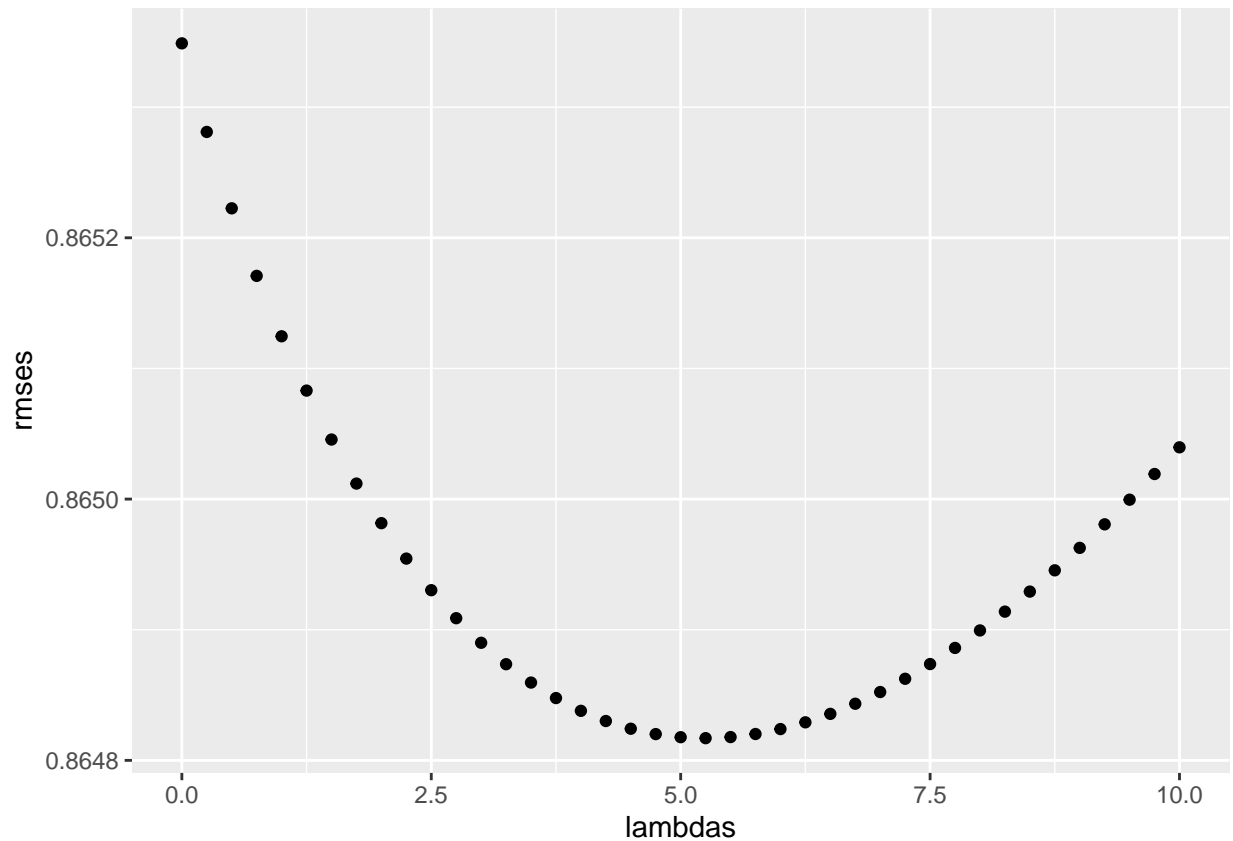
[illegible]

[illegible]

Make a Plot rmse vs lambda to select the optimal lambda

```
qplot(lambdas, rmses)
```





## The optimal lambda

```
lambda <- lambdas[which.min(rmses)]
lambda
```

```
## [1] 5.25
```

## Test and save results

```
rmse_data <- bind_rows(rmse_data,
  data_frame(method="Regularized Movie and User Effect Model",
    RMSE = min(rmses)))
```

## Result

```
rmse_data %>% knitr::kable()
```

method	RMSE
Average Movie Rating Model	1.0612018
Movie Effect Model	0.9439087
Movie-User Effect Model	0.8653488
Regularized Movie and User Effect Model	0.8648170

## Conclusion

We have built a machine learning algorithm to predict movie ratings with a given MovieLens dataset. The regularized model including the effect of user has lower RMSE value(0.8648170) which reduced the error further. We can further improve the RMSE by utilizing age, genre, year and so on.

## Appendix

```
print("Operating System:")
```

```
## [1] "Operating System:"
```

```
version
```

```
##
## platform      x86_64-apple-darwin15.6.0
## arch          x86_64
## os            darwin15.6.0
## system        x86_64, darwin15.6.0
## status
## major         3
## minor         6.3
## year          2020
## month         02
## day           29
## svn rev       77875
## language      R
## version.string R version 3.6.3 (2020-02-29)
## nickname      Holding the Windsock
```