

Movielens

Carlo Cadei

22/02/2022

HarvardX PH125.9x Data Science: Capstone

Part 1)

This part of the project consists of course material and some technical integration to make the outcomes easy to reproduce. No comments are provided.

```
#####  
# 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")  
if(!require(caret)) install.packages("caret", repos = "http://cran.us.r-project.org")  
if(!require(data.table)) install.packages("data.table", repos = "http://cran.us.r-project.org")  
  
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")  
  
# if using R 3.6 or earlier:  
# movies <- as.data.frame(movies) %>% mutate(movieId = as.numeric(levels(movieId))[movieId],  
#                                           title = as.character(title),  
#                                           genres = as.character(genres))  
# if using R 4.0 or later:  
movies <- as.data.frame(movies) %>% mutate(movieId = as.numeric(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") # if using R 3.5 or earlier, use `set.seed(1)`
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)
edx <- rbind(edx, removed)

rm(dl, ratings, movies, test_index, temp, movielens, removed)

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

#####
# Technical integration
#####

save(edx, validation, file="ProjMovie.rda")
load(file="ProjMovie.rda")

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

library(lubridate)
library(knitr)
library(xfun)

# Technical integration end

```

Part 2)

Overview: This project is inspired on October 2006 “Netflix challenge” where Netflix offered a prize to the data science community to improve its recommendation algorithm. A large dataset of movies and ratings has been provided with a train dataset: edX and a test dataset: validation. The goal of the project is to minimize the RMSE, the typical error made when predicting a movie rating, to under 0.86490 in a scale rate up to 5. After some preliminary analysis of data, 4 models have been built, from the simplest case based on the mean rating to the more complicated one based on movies, users and regularization.

```

#####
# Modelling
#####

```

```
# Preliminary Analysis
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
```

```
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
##
##
##
```

```
glimpse(edx)
```

```
## Rows: 9,000,055
## Columns: 6
## $ userId      <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, ~
## $ movieId     <dbl> 122, 185, 292, 316, 329, 355, 356, 362, 364, 370, 377, 420, ~
## $ rating      <dbl> 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, ~
## $ timestamp   <int> 838985046, 838983525, 838983421, 838983392, 838983392, 83898~
## $ title       <chr> "Boomerang (1992)", "Net, The (1995)", "Outbreak (1995)", "S~
## $ genres      <chr> "Comedy|Romance", "Action|Crime|Thriller", "Action|Drama|Sci~
```

```
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
```

```
summary(validation)
```

```
##      userId      movieId      rating      timestamp
## Min.   :    1   Min.   :    1   Min.   :0.500   Min.   :7.897e+08
## 1st Qu.:18096   1st Qu.:   648   1st Qu.:3.000   1st Qu.:9.467e+08
## Median :35768   Median :  1827   Median :4.000   Median :1.035e+09
## Mean   :35870   Mean   :  4108   Mean   :3.512   Mean   :1.033e+09
## 3rd Qu.:53621   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:999999   Length:999999
## Class :character Class :character
## Mode  :character Mode  :character
##
##
##
```

```
glimpse(validation)
```

```
## Rows: 999,999
## Columns: 6
## $ userId      <int> 1, 1, 1, 2, 2, 2, 3, 3, 4, 4, 4, 5, 5, 5, 5, 5, 5, 5, 5, ~
## $ movieId     <dbl> 231, 480, 586, 151, 858, 1544, 590, 4995, 34, 432, 434, 85, ~
## $ rating      <dbl> 5.0, 5.0, 5.0, 3.0, 2.0, 3.0, 3.5, 4.5, 5.0, 3.0, 3.0, 3.0, ~
## $ timestamp   <int> 838983392, 838983653, 838984068, 868246450, 868245645, 86824~
## $ title       <chr> "Dumb & Dumber (1994)", "Jurassic Park (1993)", "Home Alone ~
## $ genres      <chr> "Comedy", "Action|Adventure|Sci-Fi|Thriller", "Children|Come~
```

```
# Definition of RMSE
```

```
RMSE <- function(true_ratings, predicted_ratings){
  sqrt(mean((true_ratings-predicted_ratings)^2,na.rm=T))
}
```

Preliminary analysis: This simple data exploration is necessary to gain insights on possible predictors for analysis. Definition of RMSE as a loss function based on the residual mean squared error, the function computes the RMSE for vectors of ratings and their corresponding predictors.

```
# Simple model
mu <- mean(edx$rating)
mu
```

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

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

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

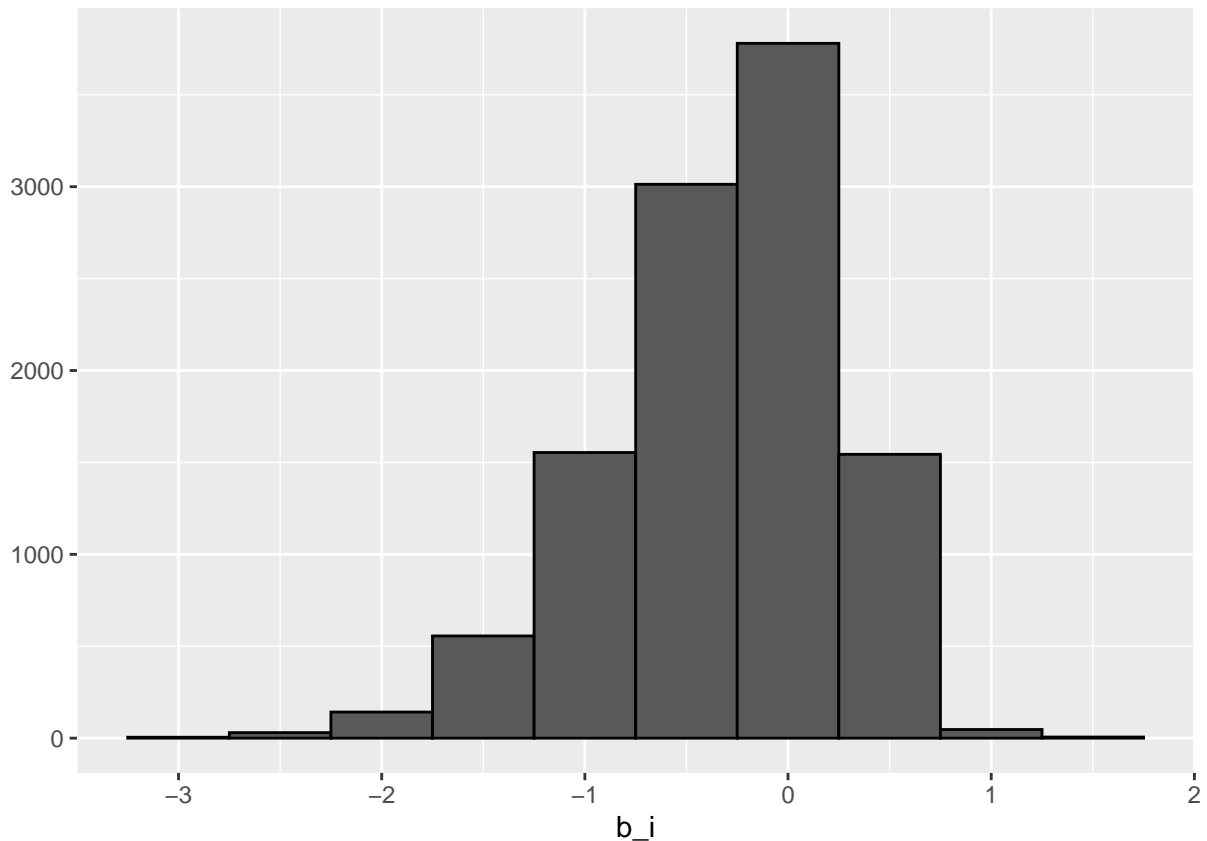
```
rmse_results <- data_frame(method = "Just the average", RMSE = naive_rmse)
```

```
## Warning: `data_frame()` was deprecated in tibble 1.1.0.
## Please use `tibble()` instead.
```

Simple model: This model simply predicts the same rating for all movies, the average of all ratings. The model returns an RMSE larger than 1 meaning that this first model typical error is larger than one star. Not so good!

```
# First model
movie_avgs <- edx %>%
  group_by(movieId) %>%
  summarize(b_i = mean(rating - mu))

movie_avgs %>% qplot(b_i, geom = "histogram", bins = 10, data = ., color = I("black"))
```



```
predicted_ratings <- mu + validation %>%
  left_join(movie_avgs, by='movieId') %>%
  .$b_i

model_1_rmse <- RMSE(predicted_ratings, validation$rating)
rmse_results <- bind_rows(rmse_results,
  data_frame(method="Movie Effect Model",
    RMSE = model_1_rmse ))

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

method	RMSE
Just the average	1.0612018
Movie Effect Model	0.9439087

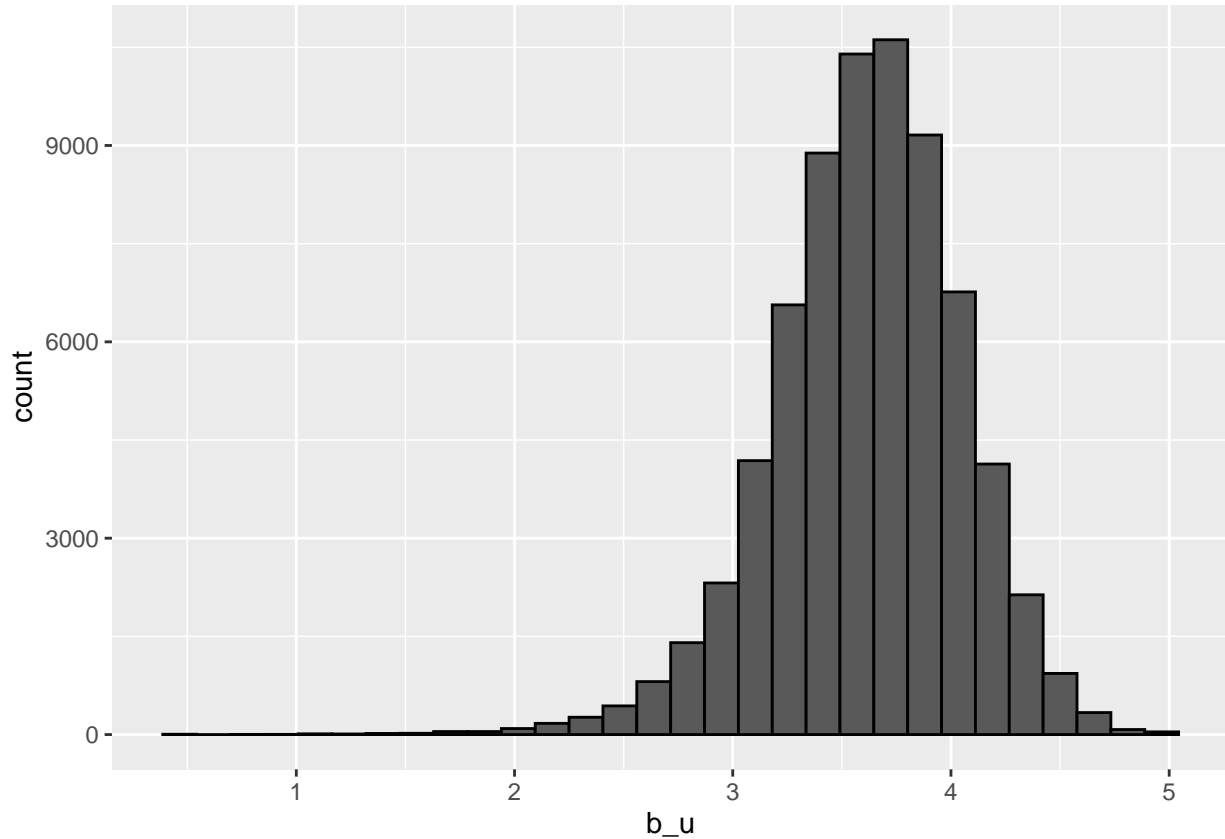
First model: This model is based on the fact that some movies are rated higher than others. We introduce the term b_i to represent average ranking for movie i . Plotting the effect b_i , we can appreciate that these estimates vary substantially. The model returns an RMSE of 0.9439, a much better datum than the simple model one.

```
# Second model
edx %>%
  group_by(userId) %>%
```

```

summarize(b_u = mean(rating)) %>%
filter(n()>=100) %>%
ggplot(aes(b_u)) +
geom_histogram(bins = 30, color = "black")

```



```

user_avgs <- edx %>%
  left_join(movie_avgs, by='movieId') %>%
  group_by(userId) %>%
  summarize(b_u = mean(rating - mu - b_i))

predicted_ratings <- validation %>%
  left_join(movie_avgs, by='movieId') %>%
  left_join(user_avgs, by='userId') %>%
  mutate(pred = mu + b_i + b_u) %>%
  .$pred

model_2_rmse <- RMSE(predicted_ratings, validation$rating)
rmse_results <- bind_rows(rmse_results,
  data_frame(method="Movie + User Effects Model",
    RMSE = model_2_rmse ))

rmse_results %>% knitr::kable()

```

method	RMSE
Just the average	1.0612018
Movie Effect Model	0.9439087
Movie + User Effects Model	0.8653488

Second model: This model is based on the fact that there is variability across users too. We introduce the term `b_u` to represent the user effect. Plotting the effect `b_u`, we see that these estimates vary substantially as well. The model returns an RMSE of 0.8653 that is significantly better than the first one.

```
# Third model
movie_titles <- edx %>%
  select(movieId, title) %>%
  distinct()

movie_avgs %>% left_join(movie_titles, by="movieId") %>%
  arrange(desc(b_i)) %>%
  select(title, b_i) %>%
  slice(1:10) %>%
  knitr::kable()
```

title	b_i
Hellhounds on My Trail (1999)	1.487535
Satan's Tango (SĀġtĀġtangĀ ³) (1994)	1.487535
Shadows of Forgotten Ancestors (1964)	1.487535
Fighting Elegy (Kenka erejii) (1966)	1.487535
Sun Alley (Sonnenallee) (1999)	1.487535
Blue Light, The (Das Blaue Licht) (1932)	1.487535
Who's Singin' Over There? (a.k.a. Who Sings Over There) (Ko to tamo peva) (1980)	1.237535
Human Condition II, The (Ningen no joken II) (1959)	1.237535
Human Condition III, The (Ningen no joken III) (1961)	1.237535
Constantine's Sword (2007)	1.237535

```
movie_avgs %>% left_join(movie_titles, by="movieId") %>%
  arrange(b_i) %>%
  select(title, b_i) %>%
  slice(1:10) %>%
  knitr::kable()
```

title	b_i
Besotted (2001)	-3.012465
Hi-Line, The (1999)	-3.012465
Accused (Anklaget) (2005)	-3.012465
Confessions of a Superhero (2007)	-3.012465
War of the Worlds 2: The Next Wave (2008)	-3.012465
SuperBabies: Baby Geniuses 2 (2004)	-2.717822
Hip Hop Witch, Da (2000)	-2.691037
Disaster Movie (2008)	-2.653090
From Justin to Kelly (2003)	-2.610455

title	b_i
Criminals (1996)	-2.512465

```
edx %>% dplyr::count(movieId) %>%
  left_join(movie_avgs) %>%
  left_join(movie_titles, by="movieId") %>%
  arrange(desc(b_i)) %>%
  slice(1:10) %>%
  knitr::kable()
```

Joining, by = "movieId"

movieId	n	b_i	title
3226	1	1.487535	Hellhounds on My Trail (1999)
33264	2	1.487535	Satan's Tango (SÄ;itÄ;ntangÄ³) (1994)
42783	1	1.487535	Shadows of Forgotten Ancestors (1964)
51209	1	1.487535	Fighting Elegy (Kenka erejii) (1966)
53355	1	1.487535	Sun Alley (Sonnenallee) (1999)
64275	1	1.487535	Blue Light, The (Das Blaue Licht) (1932)
5194	4	1.237535	Who's Singin' Over There? (a.k.a. Who Sings Over There) (Ko to tamo peva) (1980)
26048	4	1.237535	Human Condition II, The (Ningen no joken II) (1959)
26073	4	1.237535	Human Condition III, The (Ningen no joken III) (1961)
65001	2	1.237535	Constantine's Sword (2007)

```
edx %>% dplyr::count(movieId) %>%
  left_join(movie_avgs) %>%
  left_join(movie_titles, by="movieId") %>%
  arrange(b_i) %>%
  select(title, b_i, n) %>%
  slice(1:10) %>%
  knitr::kable()
```

Joining, by = "movieId"

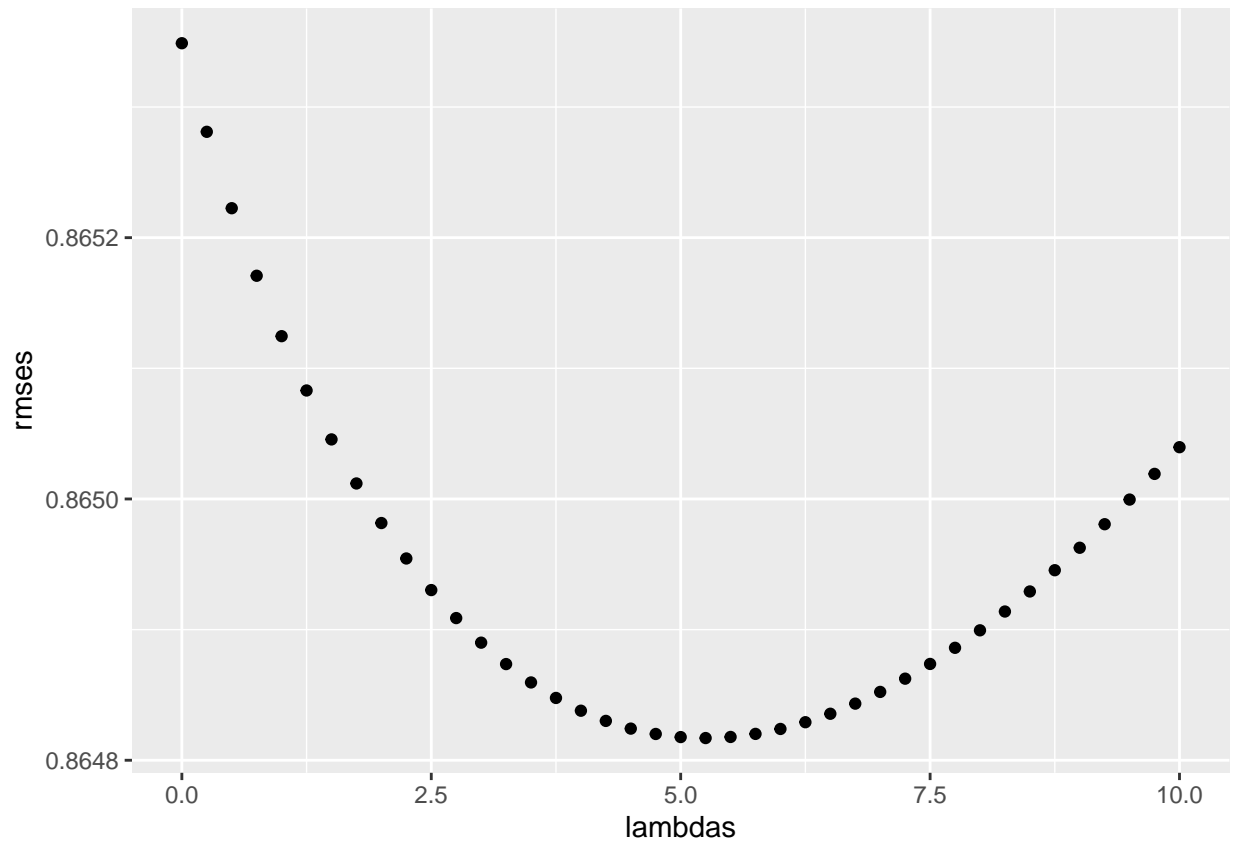
title	b_i	n
Besotted (2001)	-3.012465	2
Hi-Line, The (1999)	-3.012465	1
Accused (Anklaget) (2005)	-3.012465	1
Confessions of a Superhero (2007)	-3.012465	1
War of the Worlds 2: The Next Wave (2008)	-3.012465	2
SuperBabies: Baby Geniuses 2 (2004)	-2.717822	56
Hip Hop Witch, Da (2000)	-2.691037	14
Disaster Movie (2008)	-2.653090	32
From Justin to Kelly (2003)	-2.610455	199
Criminals (1996)	-2.512465	2

```

lambdas <- seq(0, 10, 0.25)
rmsees <- sapply(lambdas, function(l){
  mu <- mean(edx$rating)
  b_i <- edx %>%
    group_by(movieId) %>%
    summarize(b_i = sum(rating - mu)/(n()+1))
  b_u <- edx %>%
    left_join(b_i, by="movieId") %>%
    group_by(userId) %>%
    summarize(b_u = sum(rating - b_i - mu)/(n()+1))
  predicted_ratings <-
    validation %>%
    left_join(b_i, by = "movieId") %>%
    left_join(b_u, by = "userId") %>%
    mutate(pred = mu + b_i + b_u) %>%
    .$pred
  return(RMSE(predicted_ratings, validation$rating))
})

qplot(lambdas, rmsees)

```



```

lambda <- lambdas[which.min(rmsees)]
lambda

```

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

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

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

method	RMSE
Just the average	1.0612018
Movie Effect Model	0.9439087
Movie + User Effects Model	0.8653488
Regularized Movie + User Effect Model	0.8648170

```
# Modelling end
```

Third model: Checking the list of the 10 best and worst movies according to our estimate, we find that most of them have a small number of user ratings. With just a few users we have a lot of uncertainty therefore it is necessary to proceed with regularization, a Bayesian-like approach that penalizes large estimates resulting from the use of small sample sizes. This model is based on a penalty term λ chosen with a cross validation function to minimize RMSE. The model returns an RMSE of 0.864817 that is better than the second model and under 0.86490, the RMSE requested.

Conclusion: The RMSE table shows an improvement with respect to the simplest model that calculates the RMSE more than 1, then incorporating ‘Movie effect’ with a 11% improvement and ‘Movie and user effect’ with a 18.5% improvement. With a few data giving large effect on errors, a regularization model was used to penalize them. The final RMSE is 0.864817 with an improvement above 18.5%. The regularization has not given much improvement but is necessary to reach the under 0.8649 goal. It would be possible to try and improve the model using genres as effect or matrix factorization and principal component analysis but it is not the goal of this exercise.

Reference: Irizarry, R. A. - Introduction to data science & edX data science professional certificate notes