# Title: Movielens

## Author: Monica Bustamante

## Date: May/27/2020

The Movie Lens is a project to develop and train algoritmo the analysis the customers preferences in an overviwe of the data, analysis, results and conclusions.

The Methods to used for analysis consist of preparing the data: Cleaning, exploration, visualization, a

# 1 INSTALL PACKAGES, LIBRARIES

```
1 #INSTALL PACKAGES AND LIBRARIES
2 list.of.packages <- c("lubridate","stringi",
3                       "lattice", "tidyverse", "caret",
4                       "tidyr","stringr","ggplot2",
5                       "readr")
6 new.packages <- list.of.packages[!(list.of.packages %in%
7                              installed.packages()[,"Package"])]
8 if(length(new.packages)) install.packages(new.packages)
```

```
1 #Install packages
2 install.packages("rmarkdown")
```

```
1 #Install packages
2 install.packages("tidyverse", repos = "http://cran.us.r-project.org")
```

```
1 #Install packages caret
2 install.packages("caret", repos = "http://cran.us.r-project.org")
```

```
1 #Install packages
2 install.packages("data.table", repos = "http://cran.us.r-project.org")
```

```
1 #Install libraries
2 library(ggplot2)
3 library(readr)
4 library(lubridate)
5 library(stringi)
```

```
5 library(stringi)
6 library(tidyverse)
7 library(caret)
8 library(tidyr)
9 library(stringr)
```

## ▼ 2 DOWNLOAD DATA SET, SPLIT AND MUTATE.

```
1 #download data set Movielens
2 dl <- tempfile()
3  download.file("http://files.grouplens.org/datasets/movielens/ml-10m.zip", dl)
```

```
1 #Read table
2 ratings <- read.table(text = gsub("::", "\t", readLines(unzip(dl, "ml-10M100K/ratings.dat"
3                       col.names = c("userId", "movieId", "rating", "timestamp"))
```

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

```
1  #Mutate, rename title
2  movies <- as.data.frame(movies) %>% mutate(movieId = as.numeric(levels(movieId))[movieId]
3                                             title = as.character(title),
4                                             genres = as.character(genres))
```

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

## ▼ 3 VALIDATION AND TRAIN DATA SET

```
1 # Validation set will be 10% of MovieLens data
2 set.seed(1, sample.kind="Rounding")
3
4 # if using R 3.5 or earlier, use `set.seed(1)` instead
5 test_index <- createDataPartition(y = movielens$rating, times = 1, p = 0.1, list = FALSE)
6  edx <- movielens[-test_index,]
7  temp <- movielens[test_index,]
```

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

```
1 # Add rows removed from validation set back into edx set
```

```
2 removed <- anti_join(temp, validation)
3  edx <- rbind(edx, removed)
4
5 rm(dl, ratings, movies, test_index, temp, movielens, removed)
6
```

```
1 #validation dataset
2 validation  <- validation %>% select(-rating)
```

# GENERAL QUESTIONS

## ▾ How many rows and columns are there in the edx dataset?

```
1 #To see more information about the dataset
2 head(edx, 5)
```

A data.frame: 5 × 6

| | userId | movieId | rating | timestamp | title | genres |
|---|---|---|---|---|---|---|
| | <int> | <dbl> | <dbl> | <int> | <chr> | <chr> |
| **1** | 1 | 122 | 5 | 838985046 | Boomerang (1992) | Comedy\|Romance |
| **2** | 1 | 185 | 5 | 838983525 | Net, The (1995) | Action\|Crime\|Thriller |
| **4** | 1 | 292 | 5 | 838983421 | Outbreak (1995) | Action\|Drama\|Sci-Fi\|Thriller |
| **5** | 1 | 316 | 5 | 838983392 | Stargate (1994) | Action\|Adventure\|Sci-Fi |
| **6** | 1 | 329 | 5 | 838983392 | Star Trek: Generations (1994) | Action\|Adventure\|Drama\|Sci-Fi |

```
1 #Dimension Dataset
2 dim(edx)
```

9000055 · 6

```
1 str(edx)
```

```
'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
 $ title    : chr  "Boomerang (1992)" "Net, The (1995)" "Outbreak (1995)" "Stargate (199
 $ genres   : chr  "Comedy|Romance" "Action|Crime|Thriller" "Action|Drama|Sci-Fi|Thrille
```

```
1 #General information about dataset
2 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
```

```
1 #How many rows and columns are there in the edx dataset
2 paste('The dataset has',nrow(edx),'rows and',ncol(edx),'columns.')
```

'The dataset has 9000055 rows and 6 columns.'

```
1 #To see more information about dataset
2 edx %>% summarise(
3   uniq_movies = n_distinct(movieId),
4   uniq_users = n_distinct(userId),
5   uniq_genres = n_distinct(genres))
```

A data.frame: 1 × 3

| uniq_movies | uniq_users | uniq_genres |
|---|---|---|
| <int> | <int> | <int> |
| 10677 | 69878 | 797 |

```
1 #Mean of rating dataset
2 rating_mean <- mean(edx$rating)
3 rating_mean
```

3.51246520160155

## ▾ How many zeros were given as ratings in the edx dataset?

```
1 #How many zeros were given as ratings in the edx dataset.
2 paste(sum(edx$rating == 0), 'ratings and',
3       sum(edx$rating == 3),'ratings with 3')
```

'0 ratings and 2121238 ratings with 3'

```
1 edx %>% filter(rating == 3) %>% tally()
```

> A
> data.frame:
>   1 × 1
>     **n**
>   **<int>**

## How many different movies are in the edx dataset?

```
1 #How many different movies are in the edx dataset
2 n_distinct(edx$movieId)
```

⌐→   10677

```
1 edx %>% summarize(n_movies = n_distinct(movieId))
```

⌐→        A
      data.frame:
        1 × 1
      **n_movies**
        **<int>**
        10677

## How many different users are in the edx dataset?

```
1 #How many different users are in the edx dataset. n_distinct or lenght
2 n_distinct(edx$userId)
```

⌐→   69878

```
1 edx %>% summarize(n_users = n_distinct(userId))
```

⌐→        A
      data.frame:
        1 × 1
      **n_users**
        **<int>**
        69878

## How many movie ratings are in each of the following genres in the edx

```
1 # str_detect
2 genres = c("Drama", "Comedy", "Thriller", "Romance")
3 sapply(genres, function(g) {
4     sum(str_detect(edx$genres, g))
```

```
 5 })
 6
 7 # separate_rows, much slower!
 8 edx %>% separate_rows(genres, sep = "\\|") %>%
 9     group_by(genres) %>%
10     summarize(count = n()) %>%
11     arrange(desc(count))
```

Drama:         3910127 Comedy:        3540930 Thriller:        2325899 Romance:         1712100

```
1 #Movie ratings by Drama. str_detect Detect The Presence Or Absence Of A Pattern In A Strin
2 drama <- edx %>% filter(str_detect(genres,"Drama"))
3 paste('Drama has',nrow(drama),'movies')
```

'Drama has 3910127 movies'

```
1 #Movie ratings by Comedy
2 comedy <- edx %>% filter(str_detect(genres,"Comedy"))
3 paste('Comedy has',nrow(comedy),'movies')
```

'Comedy has 3540930 movies'

```
1 ##Movie ratings by Thriller
2 thriller <- edx %>% filter(str_detect(genres,"Thriller"))
3 paste('Thriller has',nrow(thriller),'movies')
```

'Thriller has 2325899 movies'

```
1 #Movie ratings by Romance
2 romance <- edx %>% filter(str_detect(genres,"Romance"))
3 paste('Romance has',nrow(romance),'movies')
```

'Romance has 1712100 movies'

## ▾ Which movie has the greatest number of ratings?

```
1 #Greatest number of ratings. Arrange rows by variables
2 edx %>% group_by(title) %>%
3 summarise(number = n()) %>%
4 arrange(desc(number))
```

A tibble: 10676 × 2

| title | number |
|-------|--------|
| <chr> | <int> |
| Pulp Fiction (1994) | 31362 |
| Forrest Gump (1994) | 31079 |
| Silence of the Lambs, The (1991) | 30382 |
| Jurassic Park (1993) | 29360 |
| Shawshank Redemption, The (1994) | 28015 |
| Braveheart (1995) | 26212 |
| Fugitive, The (1993) | 25998 |
| Terminator 2: Judgment Day (1991) | 25984 |
| Star Wars: Episode IV - A New Hope (a.k.a. Star Wars) (1977) | 25672 |
| Apollo 13 (1995) | 24284 |
| Batman (1989) | 24277 |
| Toy Story (1995) | 23790 |
| Independence Day (a.k.a. ID4) (1996) | 23449 |
| Dances with Wolves (1990) | 23367 |
| Schindler's List (1993) | 23193 |
| True Lies (1994) | 22823 |
| Star Wars: Episode VI - Return of the Jedi (1983) | 22584 |
| 12 Monkeys (Twelve Monkeys) (1995) | 21891 |
| Usual Suspects, The (1995) | 21648 |
| Fargo (1996) | 21395 |
| Speed (1994) | 21361 |
| Aladdin (1992) | 21173 |
| Matrix, The (1999) | 20908 |
| Star Wars: Episode V - The Empire Strikes Back (1980) | 20729 |
| Seven (a.k.a. Se7en) (1995) | 20311 |
| American Beauty (1999) | 19950 |
| Raiders of the Lost Ark (Indiana Jones and the Raiders of the Lost Ark) (1981) | 19678 |
| Back to the Future (1985) | 19034 |
| Mission: Impossible (1996) | 18992 |
| Ace Ventura: Pet Detective (1994) | 18959 |
| ⋮ | ⋮ |
| Please Vote for Me (2007) | 1 |
| Quarry, The (1998) | 1 |
| Quiet City (2007) | 1 |
| Relative Strangers (2006) | 1 |
| Ring of Darkness (2004) | 1 |
| Rockin' in the Rockies (1945) | 1 |
| Säg att du älskar mig (2006) | 1 |
| Shadows of Forgotten Ancestors (1964) | 1 |
| Small Cuts (Petites coupures) (2003) | 1 |
| Splinter (2008) | 1 |
| Stacy's Knights (1982) | 1 |
| Stone Angel, The (2007) | 1 |
| Strange Planet (1999) | 1 |

Strange Planet (1999)　　　　　　　　　　　　　　　　　1

Sun Alley (Sonnenallee) (1999)　　　　　　　　　　　　　1

Sun Shines Bright, The (1953)　　　　　　　　　　　　　1

Symbiopsychotaxiplasm: Take One (1968)　　　　　　　1

Tattooed Life (Irezumi ichidai) (1965)　　　　　　　　　1

Testament of Orpheus, The (Testament d'Orphée) (1960)　　1

Tokyo! (2008)　　　　　　　　　　　　　　　　　　　1

Train Ride to Hollywood (1978)　　　　　　　　　　　　1

Twice Upon a Time (1983)　　　　　　　　　　　　　　1

Uncle Nino (2003)　　　　　　　　　　　　　　　　　1

## ▾ What are the five most given ratings in order from most to least?

When Time Ran Out... (a.k.a. The Day the World Ended) (1980)　　　　1

```
1 #Sort a variable in descending order.
2 edx %>% group_by(rating) %>%
3 summarize(count = n()) %>%
4 top_n(5) %>%
5   arrange(desc(count))
```

⇨　Selecting by count

A tibble: 5 × 2

| rating | count |
| --- | --- |
| <dbl> | <int> |
| 4.0 | 2588430 |
| 3.0 | 2121240 |
| 5.0 | 1390114 |
| 3.5 | 791624 |
| 2.0 | 711422 |

```
1 head(sort(-table(edx$rating)),5)
```

⇨
```
        4         3         5       3.5         2
 -2588430  -2121240  -1390114   -791624   -711422
```
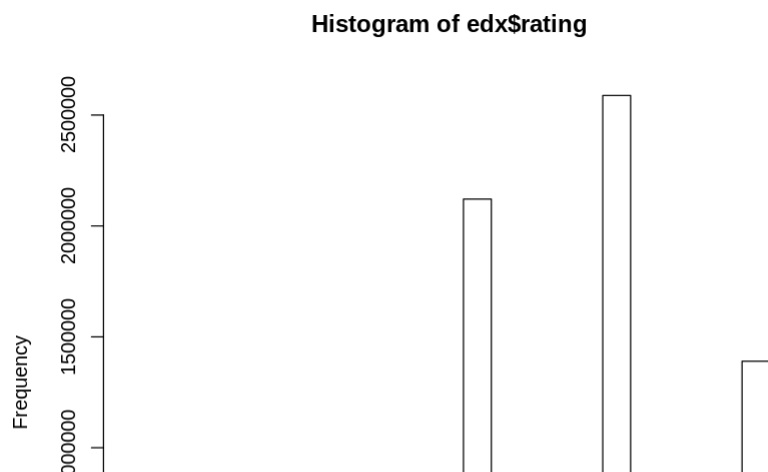
```
1 hist(edx$rating)
2 summary(edx$rating)
```

⇨

```
     Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
    0.500   3.000   4.000   3.512   4.000   5.000
```

**Histogram of edx$rating**



True or False: In general, half star ratings are less common than whole fewer ratings of 3.5 than there are ratings of 3 or 4, etc.).

```
 1 #Rating movies
 2 rating4 <- table(edx$rating)["4"]
 3 rating35 <- table(edx$rating)["3.5"]
 4 rating3 <- table(edx$rating)["3"]
 5
 6 Result <- (rating35 < rating3 && rating35 < rating4)
 7
 8 print(Result)
 9
10 rm(rating3, rating35,  rating4, Result)
```
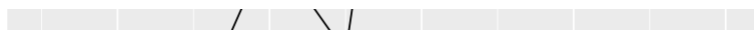
```
[1] TRUE
```

```
 1 #Graphic Rating movies
 2 edx %>%
 3   group_by(rating) %>%
 4   summarize(count = n()) %>%
 5   ggplot(aes(x = rating, y = count)) +
 6   geom_line()
```

# 4 MODELING



## ▾ Predicted movie ratings and calculates RMSE.

Movie rating predictions will be compared to the true ratings in the validation set using RMSE

```
1 data <- movies %>% separate_rows(genres, sep ="\\|")
2 DAT.aggregate <- aggregate(formula = cbind(n = 1:nrow(dat)) ~ genres, data = data, FUN = l
```
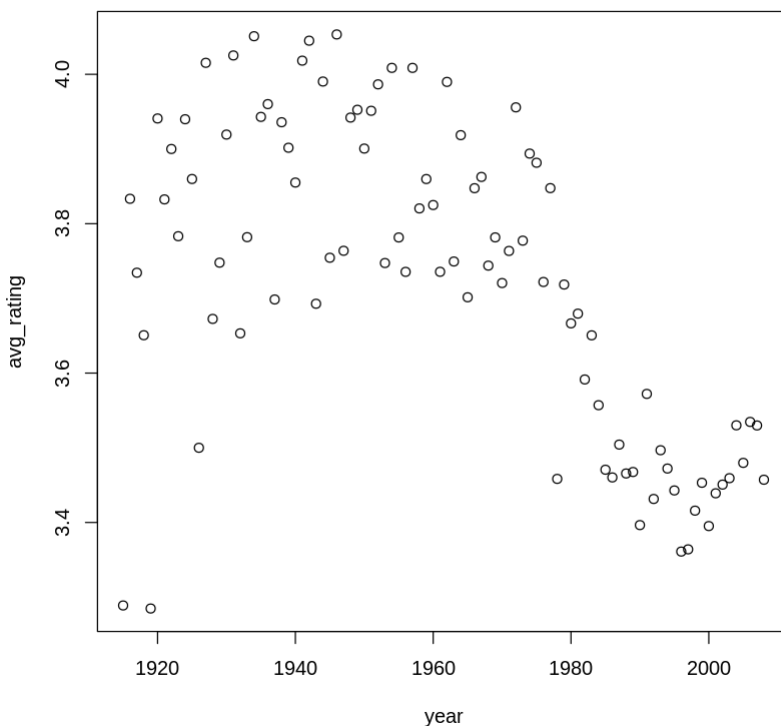
```
1 #Size of dataset
2 movielens <- left_join(ratings,
3                     movies, by = "movieId")
4 nrow(movielens)
```

�localhost   10000054

```
1 #Creates Year column.
2 edx <- edx %>%
3   mutate(title = str_trim(title)) %>%
4   extract(title, c("title_tmp", "year"),
5           regex = "^(.*) \\(([0-9 \\-]*)\\)$",
6           remove = F) %>%
7   mutate(year = if_else(str_length(year) > 4,
8                     as.integer(str_split(year, "-",
9                                         simplify = T)[1]),
10                     as.integer(year))) %>%
11   mutate(title = if_else(is.na(title_tmp), title, title_tmp)) %>%
12   select(-title_tmp)  %>%
13   mutate(genres = if_else(genres == "(No Genres Listed)",
14                         `is.na<-`(genres), genres))
15 validation <- temp %>%
16   semi_join(edx, by = "movieId") %>%
17   semi_join(edx, by = "userId")
```

```
17    semi_join(edx, by = "userId")
```

```
1 avg_ratings <- edx %>%
2 group_by(year) %>%
3 summarise(avg_rating = mean(rating))
4 plot(avg_ratings)
```
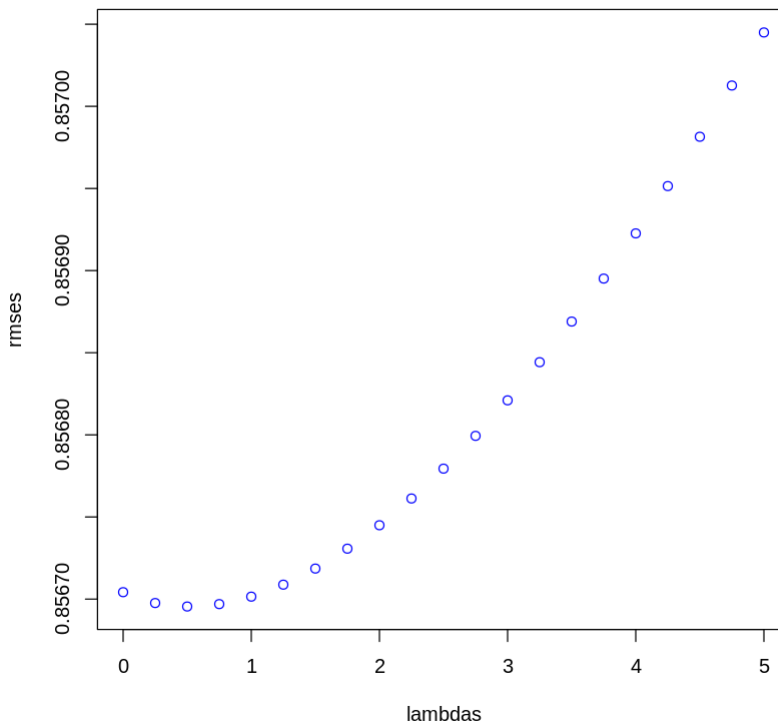


```
1 #Root Mean Square Error
2 RMSE <- function(true_ratings, predicted_ratings){
3       sqrt(mean((true_ratings - predicted_ratings)^2))
4      }
5
6 lambdas <- seq(0, 5, 0.25)
7 rmses <- sapply(lambdas,function(l){
8   mu <- mean(edx$rating) #The mean of ratings from training set
9
10  Movie_effect <- edx %>%  #Adjust mean by movie effect
11    group_by(movieId) %>%
12    summarize(Movie_effect = sum(rating - mu)/(n()+l))
13
14  Movie_user <- edx %>% #Ajdust mean by movie effect and user
15    left_join(Movie_effect, by="movieId") %>%
16    group_by(userId) %>%
17    summarize(Movie_user = sum(rating - Movie_effect - mu)/(n()+l))
18
19  predicted_ratings <-
20    edx %>%
21    left_join(Movie_user, by = "userId") %>%
```

```
21     left_join(movie_user, by = "userId") %>%
22     left_join(Movie_effect, by = "movieId") %>%
23     mutate(pred = mu + Movie_effect + Movie_user) %>%
24     .$pred #Predict ratings
25
26   return(RMSE(predicted_ratings, edx$rating))
27 })
28 plot(lambdas, rmses,
29     col = "blue")
```

⤷



```
1 #Calculate Lambda optimal RMSE
2 lambda <- lambdas[which.min(rmses)]
3 paste('RMSE',min(rmses),'Lambda',lambda)
```

⤷    'RMSE 0.856695492876063 Lambda 0.5'

## ▾ CONCLUSION:

Predict a list of rated movies.

Discovered patterns: as people prefer movies with a medium to high rating. (3 to 5).

The movies preferred by the customers was the end of the 1980 and 1990 periods.