# MovieLens Project Report

*Dana Saeed*

*23/11/2020*

## Introduction

AI is vital to a profoundly innovative present day business. It depicts information about client bases all in all or customizes an encounter for a solitary client. The point of AI is to handle information into accommodating data and normally natural arrangements. Compelling AI calculations have pulled in a great deal of consideration as of late: for instance, in 2006, Netflix put a seven-figure abundance on a checked improvemnent to their film proposal framework.

We work off the Netflix challenge reason and, all the more explicitly, anticipate film appraisals for clients in a huge dataset. We train a direct model to produce anticipated film appraisals and compute the Root Mean Square Error (RMSE) of the anticipated evaluations versus the real evaluations.

This report is made out of four sections: the presentation has introduced the issue, the synopsis portrays the dataset and creates starter requests, the techniques segment builds up the model and actualizes in with the going with .R record, and the end segment shares the outcomes.

## Summary

We use the MovieLens 10M dataset that consists of 10 million ratings and 100,000 tag applications applied to 10,000 movies by 72,000 users. This leads to wildly varying numbers of ratings for each movie, with the most rated movie being *Pulp Fiction* (1994) with over 31,000 ratings and over 100 titles with a single rating.

```
# Most rated films edx %>%
group_by(title) %>%
summarize(n_ratings = n()) %>%
arrange(desc(n_ratings))
```

```
## # A tibble: 10,676 x 2
##       title                                                    n_ratings
##       <chr>     <int> ## 1 Pulp Fiction (1994)        31362 ## 2 Forrest Gump (1994)   31079
## 3 Silence of the Lambs, The (1991)                               30382
## 4 Jurassic Park (1993)                                           29360
## 5 Shawshank Redemption, The (1994)                               28015
## 6 Braveheart (1995)     26212 ## 7 Fugitive, The (1993)     25998
## 8 Terminator 2: Judgment Day (1991)                              25984
## 9 Star Wars: Episode IV - A New Hope (a.k.a. Star Wars) (1977)   25672
## 10 Apollo 13 (1995)                                              24284
## # ... with 10,666 more rows
```

```
# Number of movies rated once edx %>%
group_by(title) %>%
summarize(n_ratings = n()) %>%
filter(n_ratings==1) %>% count() %>%
pull()
```

## [1] 126

Because of the size of the dataset, worked in AI calculations utilizing R bundles, for example, caret would require an excessive amount of assets for a PC to run in an ideal way if by any means. Along these lines, we accomplish our objective utilizing a straight model created in the techniques segment of this report. The RMSE work is utilized from the DescTools bundle as the proportion of exactness.

The dataset is part 90-10 on train and test sets individually. This is finished in the initial steps of the content. The preparation set (edx) has 9,000,055 sections with 6 segments. Likewise, the test set (approval) has 999,999 passages and 6 segments. The segment data is appeared beneath for the approval dataset.**glimpse**(validation)

## Observations: 999,999 ##
Variables: 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, 8... ## $ 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.... ## $ timestamp <int> 838983392, 838983653, 838984068, 868246450, 868245645, 86... ## $ title      <chr> "Dumb & Dumber (1994)", "Jurassic Park (1993)", "Home Alo... ## $ genres        <chr> "Comedy", "Action|Adventure|Sci-Fi|Thriller", "Children|C...

## Methods

The simplest model is to use the average across every user and every movie as all of our predicted ratings. This model follows the simple equation,

$$Y_{u,i} = \mu, \tag{1}$$

where $Y_{u,i}$ is the predicted rating of user $u$ and movie $i$ and $\mu$ is the average rating across all entries. This is computed as 3.512 (mean(edx$rating)).

mu <- **mean**(edx$rating) **RMSE**(validation$rating, mu)

## [1] 1.061202

To improve our model, we add a free mistake term $b_{u,i}$ that communicates rating contrasts for clients and films. We will add the client predisposition term later, yet for the time being we add the film inclination term

2

$b\_i$. This term midpoints the rankings for any film $i$ in light of the fact that some are loved or abhorred more than others. The new model is:

$$Y_{u,i} = \mu + b_i. \tag{2}$$

```r
# add movie bias term, b_i b_i <- edx
%>%
   group_by(movieId) %>% summarize(b_i =
   mean(rating - mu))

# predict all unknown ratings with mu and b_i predicted_ratings <- validation

%>%
   left_join(b_i, by='movieId') %>%
   mutate(pred = mu + b_i) %>% pull(pred)

# calculate RMSE of movie ranking effect

RMSE(validation$rating, predicted_ratings)
```

## [1] 0.9439087

Now we introduce the user bias term $b_u$ in order to further improve our model. This term minimizes the effect of extreme ratings made by users that love or hate every movie. Each user $u$ is given a bias term that sways their predicted movies. Our updated model is:

$$Y_{u,i} = \mu + b_i + b_u. \tag{3}$$

```r
# add user bias term, b_u b_u <- edx %>%
left_join(b_i, by='movieId') %>% group_by(userId)
%>% summarize(b_u = mean(rating - mu - b_i))

# predict new ratings with movie and user bias

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

# calculate RMSE of movie and user bias effect RMSE(predicted_ratings,

validation$rating)
```
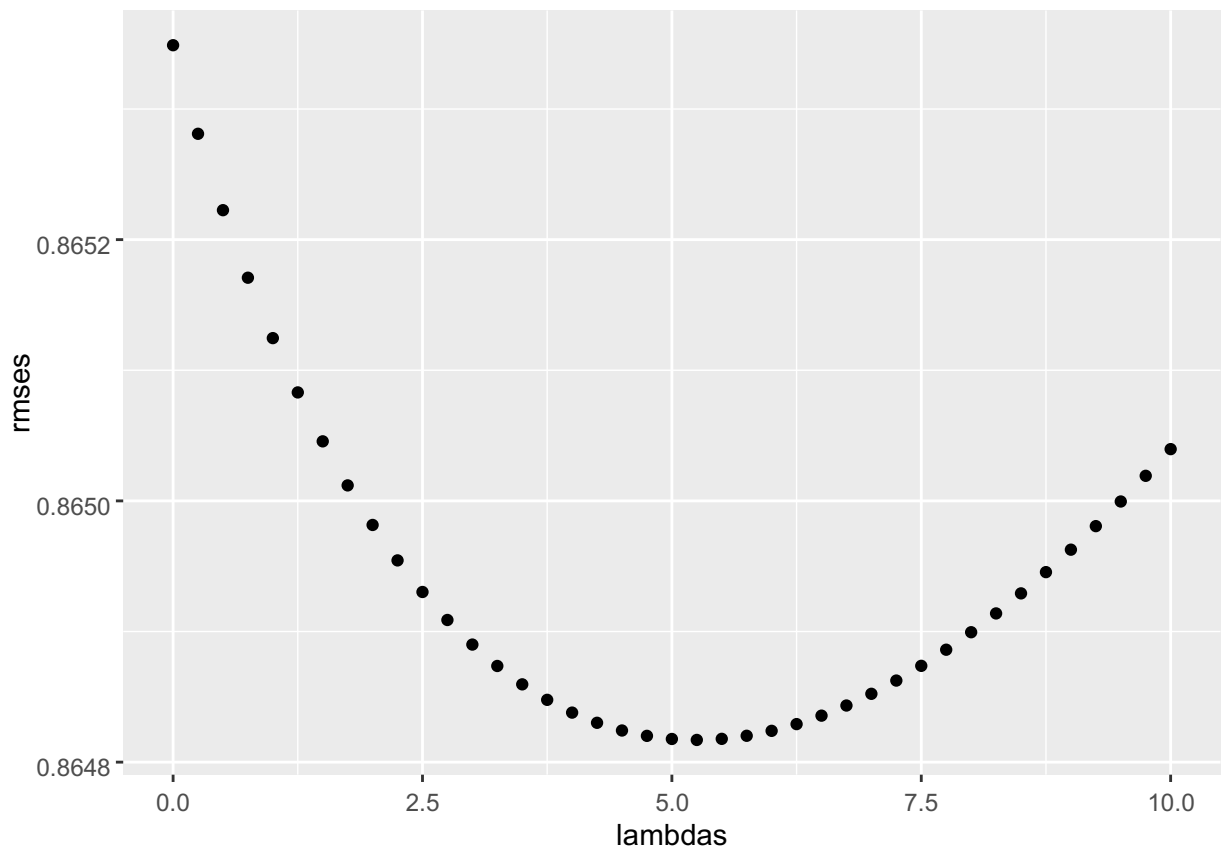
## [1] 0.8653488

At last, we utilize regularization to diminish the impact of huge mistakes in our expectations. Regularization punishes off base assessments on little example sizes. For example, our $b\_i$ term represents the normal deviation on all appraisals of a film, regardless of whether there is 1 or 100 evaluations to the film. We use regularization to diminish the sensational impact that an incredibly outrageous rating will have on our $b\_i$ term. This technique is additionally applied to the client predisposition term $b\_u$ to lessen enormous peculiarities in the evaluations of clients.

Regularization accomplishes similar objective as certainty stretches when you are simply ready to foresee a solitary number, not a span. Our new model is:

$$\frac{1}{N}\sum_{u,i}(Y_{u,i} - \mu - b_i - b_u)^2 + \lambda(\sum_i b_i^2 + \sum_u b_u^2), \tag{4}$$

where the first term is our previous least squares equation and the last term is the penalty with large bias terms. Minimizing the biases using a single $\lambda$ is the goal to our model shown above. We test lamda <seq(from=0, to=10, by=0.25) and plot the results below:

**qplot**(lambdas, rmses)



We see that the minimizing $\lambda$ term is

lambdas[**which.min**(rmses)]

## [1] 5.25

## Results

For completeness, the final model is executed below:

```r
# choose minimized lambda lam <-
lambdas[which.min(rmses)]

# compute regularize movie bias term
b_i <- edx %>%
    group_by(movieId) %>% summarize(b_i = sum(rating -
    mu)/(n()+lam))

# compute regularize user bias term b_u <-
edx %>% left_join(b_i, by="movieId") %>%
group_by(userId) %>%
    summarize(b_u = sum(rating - b_i - mu)/(n()+lam))

# compute predictions on validation set based on these above terms
predicted_ratings <- validation %>% left_join(b_i, by =
    "movieId") %>% left_join(b_u, by = "userId") %>%
    mutate(pred = mu + b_i + b_u) %>% pull(pred)

# output RMSE of our final model
RMSE(predicted_ratings, validation$rating)
```

## [1] 0.864817

We can see incremental improvements to the RMSE as we supplant our model with bias terms and regularization.

| Method | RMSE |
|---|---|
| Average | 1.06120 |
| Movie effect | 0.94391 |
| Movie and user effects | 0.86535 |
| Regularized movie and user effect | 0.86481 |

Our model is significantly more efficient than machine learning algorithms from R packages applied to this large data set. Because of the simplicity of the linear model, we are able to predict movie ratings without a serious toll on the computer resources.