

Movielens project

Predicting movie ratings by machine learning modeling

Contents

1	Introduction	2
1.1	Overview	2
1.2	Scope of the project	2
2	Method of analysis	2
2.1	Importing and checking movielens data set	2
2.2	Choosing the parameters for modeling	2
2.3	Additional parameters	3
2.4	Creating training and test sets	3
2.5	Deciding on modeling approach	3
2.6	Naive model	4
2.7	Regularized model	5
2.8	Results on edX test set	10
3	Results	10
4	Conclusions	10
5	Appendix	10
5.1	Session info	10
5.2	Full code	10

1 Introduction

1.1 Overview

This report presents a summary of **Movielens project**, which is part of Data Science Professional Certificate led by HarvardX on edX platform. The Movielens project was part of final graded assessments.

In the following sections some parts of code have been highlighted, but for reference full code has been included in the Appendix. This report, as well as a script file, can be found on GitHub repository.

1.2 Scope of the project

The scope of this project is to find a way to predict user ratings for a movie, based on chosen set of predictors from provided dataset, using possible machine learning techniques. The calculations will be based on GroupLens research lab Movielens 10M dataset, which consists of 10 million ratings on 10 thousand movies, made by 72 thousand users. Each rating has a different set of predictor. The goal of the project is to find a prediction method that would generate residual mean squared error (RMSE) lower than 0.86490.

2 Method of analysis

2.1 Importing and checking movielens data set

Starting the project I've had to import the Movielens data and create useful data set out of it. As a first step I download the zip file from web and transform the two dat files inside it (movies and ratings) using `fread` and `str_split_fixed` functions from `data.table` and `stringr` packages (part of tidyverse pack). Transformed data frames are saved to *movies* and *ratings* respectively. Both data frames are then joined together by *movieId* to *movielens* data frame that will be used onward. All code for those actions was provided in the project description.

After creating *movielens* data set, I inspected it using `glimpse` function. Our movielens data set consists of 6 parameters:

- *userId* (integer)
- *movieId* (double class)
- *rating* (double class)
- *timestamp* (integer)
- *title* (character)
- *genres* (character)

2.2 Choosing the parameters for modeling

Establishing that *rating* is our outcome, I've considered the influence of the rest of the five parameters on the expected outcome. The PCA analysis couldn't be used to help with cumulative variance explanation, as *movielens* dataset is too big for machine calculations possibilities (and another transformation on character parameters would be needed).

The parameters that have been chosen as primary ones were *userId*, *movieId* and *genres*. *timestamp* and *title* parameters have been used to create additional parameters in different form (see next section).

2.3 Additional parameters

Additional two parameters were created and included in *movielens* data set, which we suspect might be influencing the rating gave by users.

First one is *year_of_release*, indicating in which year the movie has been released. It has been created by taking the years from *title* parameter, as those years are included at the end of character string in the brackets. It has been transformed into double-class using *str_sub* function from *stringr* package, and *as.numeric* function.

Another support parameter created to *movielens* data frame is *rateday*, which is the day of the week the rating has been done (assuming it's happening on the same day the movie is watched), in scale 1-7 where 1 is Monday. This parameter is included based on assumption that users rate differently if they are relaxed over the weekend, or it's middle of stressful week. *rateday* has been added to *movielens* data set as transformation of *timestamp* parameter by *as_datetime* and *wday* functions from *lubridate* package.

```
movielens <- mutate(movielens,
  year_of_release = as.numeric(str_sub(title, start=-5, end=-2)),
  rateday=wday(as_datetime(timestamp), week_start = 1))
```

2.4 Creating training and test sets

For the modeling purposes I've created a training (*edx*) and test (*validation*) sets out of *movielens* data frame. Test set *validation* is consisting of 10% of original *movielens* set. This part was done with *createDataPartition* function from the *caret* package. Code for creating *edx* and *validations* sets was provided as part of project description.

As test set *validation* has to be used only for final check of prepared prediction model, I've created additional test and training sub-sets of training set *edx*, with same approach as used before, thus creating *edx_train* and *edx_test* variables. The test set was created using 50% of *edx* data set. From now on, until final prediction model is created, all trainings and tests would be done on those data frames.

In both cases I've made sure that *userId*, *movieId*, *year_of_release*, *rateday* and *genres* appearing in test sets are also in training sets. This has been done by using *semi_join* function on those two sets, by this selected parameters.

2.5 Deciding on modeling approach

Given the size of dataset, more complex algorithms for predicting *ratings* outcome, like random forest, K-nearest neighbors (KNN) or any regression forms couldn't be used (train function was taking too much time on given machine to calculate). I've decided to follow step by step Naive Bayes approach on all 5 chosen features. This could be described by following equation

$$Y_{i,u,y,d,g} = \mu + b_i + b_u + b_y + b_d + b_g + \epsilon_{i,u,y,d,g}$$

where:

- $Y_{i,u,y,d,g}$ - predicted movie rating
- μ - actual rating for all movies
- b_i - *movieId* effect
- b_u - *userId* effect
- b_y - *year_of_release* effect
- b_d - *rateday* effect
- b_g - *genres* effect
- $\epsilon_{i,u,y,d,g}$ - independent errors

2.6 Naive model

First step to creating a naive prediction model is to estimate μ . This can be done by assuming it's close to the average of all outcomes (*rating*). Effect of *movieId* feature, as well as *userId*, *year_of_release*, *rateday* and *genres* effect has been calculated one by one as the difference of the predicted rating and all effects calculated before (and the independence error). Therefore the equation for *movieId* effect b_i is as follow

$$b_i = Y_{i,u,y,d,g} - \mu - \epsilon_{i,u,y,d,g}$$

and equation for *userId* b_u effect being

$$b_u = Y_{i,u,y,d,g} - \mu - b_i - \epsilon_{i,u,y,d,g}$$

and all other effects of b_y , b_d and b_g created similarly.

All of those effects are calculated on `edx_train` train set, with code as followed:

```
mu <- mean(edx_train$rating)

b_i <- edx_train %>%
  group_by(movieId) %>%
  summarise(b_i = mean(rating-mu))

b_u <- edx_train %>%
  left_join(b_i, by="movieId") %>%
  group_by(userId) %>%
  summarise(b_u=mean(rating-mu-b_i))

b_y <- edx_train %>%
  left_join(b_i, by="movieId") %>%
  left_join(b_u, by="userId") %>%
  group_by(year_of_release) %>%
  summarise(b_y=mean(rating-mu-b_i-b_u))

b_d <- edx_train %>%
  left_join(b_i, by="movieId") %>%
  left_join(b_u, by="userId") %>%
  left_join(b_y, by="year_of_release") %>%
  group_by(rateday) %>%
  summarise(b_d=mean(rating-mu-b_i-b_u-b_y))

b_g <- edx_train %>%
  left_join(b_i, by="movieId") %>%
  left_join(b_u, by="userId") %>%
  left_join(b_y, by="year_of_release") %>%
  left_join(b_d, by="rateday") %>%
  group_by(genres) %>%
  summarise(b_g=mean(rating-mu-b_i-b_u-b_y-b_d))
```

The effects calculated that way are then used to calculate predictions on movie rate $Y_{i,u,y,d,g}$. For this purpose, I've created *predictions* function, which will be used through all project, also for final testing on *validations* test set.

```

predictions<- function(x,i,u,y,d,g){
  x %>%
    left_join(i, by="movieId") %>%
    left_join(u, by="userId") %>%
    left_join(y, by="year_of_release") %>%
    left_join(d, by="rateday") %>%
    left_join(g, by="genres") %>%
    mutate(pred=mu+b_i+b_u+b_y+b_d+b_g) %>%
    pull(pred)
}

pred <- predictions(edx_test, b_i, b_u, b_y, b_d, b_g)

```

Creating prediction vector, I was able to test out what would be the residual mean squared error (RMSE) of current model. This could be calculated by following

$$RMSE = \sqrt{\frac{1}{N} \sum_{i,u,y,d,g} (\hat{y}_{i,u,y,d,g} - y_{i,u,y,d,g})^2}$$

where:

- $y_{i,u,y,d,g}$ - actual rating for a movie with selected features
- $\hat{y}_{i,u,y,d,g}$ - predicted rating for a movie with selected features

equation, which in this project can be translated to code:

```
rmse <- sqrt(mean((pred-edx_test$rating)^2))
```

The RMSE obtained from naive model, although calculated based on train and test subsets and to proper whole *edx* set, is equal 0.86907, so above 0.86490 threshold required by the project. Therefore current model is not enough, and further modification of it is required.

2.7 Regularized model

To make RMSE of the model smaller, all five features were adjusted by separate regularization parameters. Using regularization parameters (named *alpha*, *lambda delta*, *kappa* and *omega*) allowed to influence total variability of effect sizes. This can be done by transforming RMSE equation with penalized regression (for controlling variety), with focus on minimizing paramter effect.

For example for b_i it would be

$$\frac{1}{N} \sum_{i,u,y,d,g} (\hat{y}_{i,u,y,d,g} - \mu - b_i)^2 + \lambda \sum_i b_i^2$$

transformed into

$$\hat{b}_i(\lambda) = \frac{1}{\lambda + n_i} \sum_{u,y,d,g=1}^{n_i} (Y_{i,u,y,d,g} - \hat{\mu})$$

2.7.1 Regularized movie effect

For regularized movie effect (b_i), regularization has been done with α parameter. The best value of α has been chosen using supply function on α values from 0:10 (in the code it has been narrowed down for the sake of calculation speed).

Value of b_i has been calculated using following code

```
b_i <- edx_train %>%  
  group_by(movieId) %>%  
  summarise(b_i = sum(rating-mu)/(n()+alpha))
```

and the rest of effects has been calculated using basic naive approach.

The best value of α has been determined as 3.85.

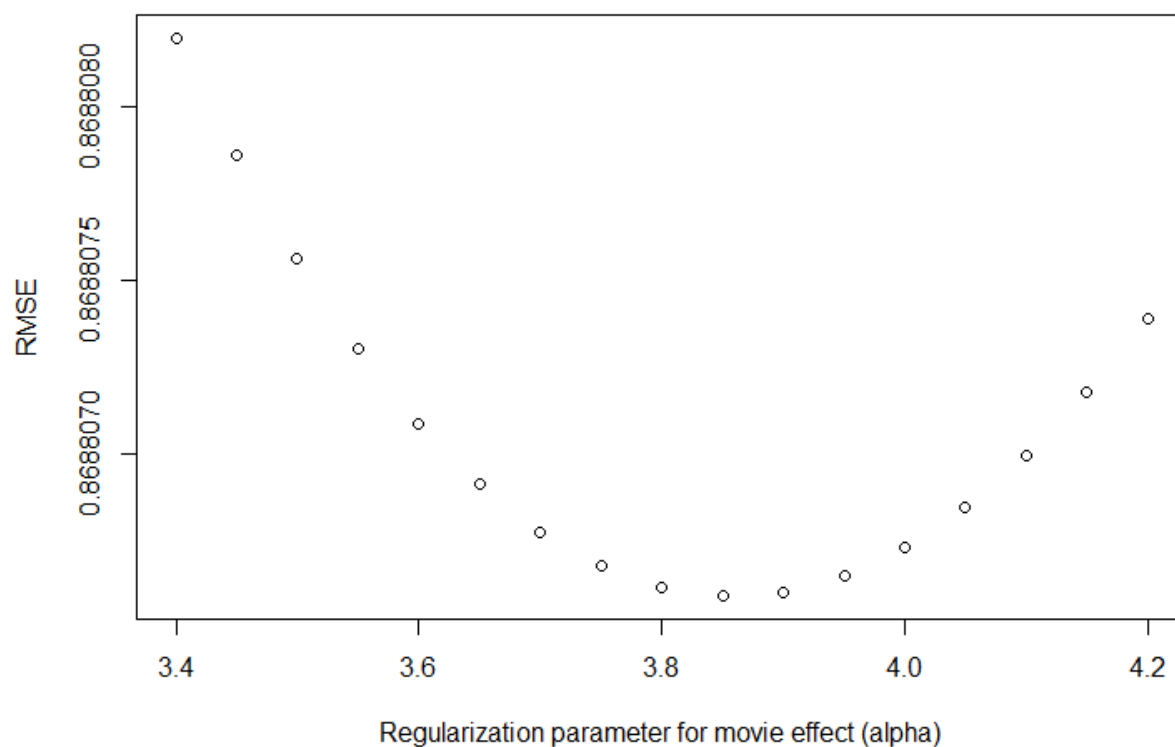


Figure 1: RMSE vs. Alpha plot

For this value of α RMSE has been calculated as 0.86881, which is 0.03% improvement from fully naive approach.

2.7.2 Regularized user effect

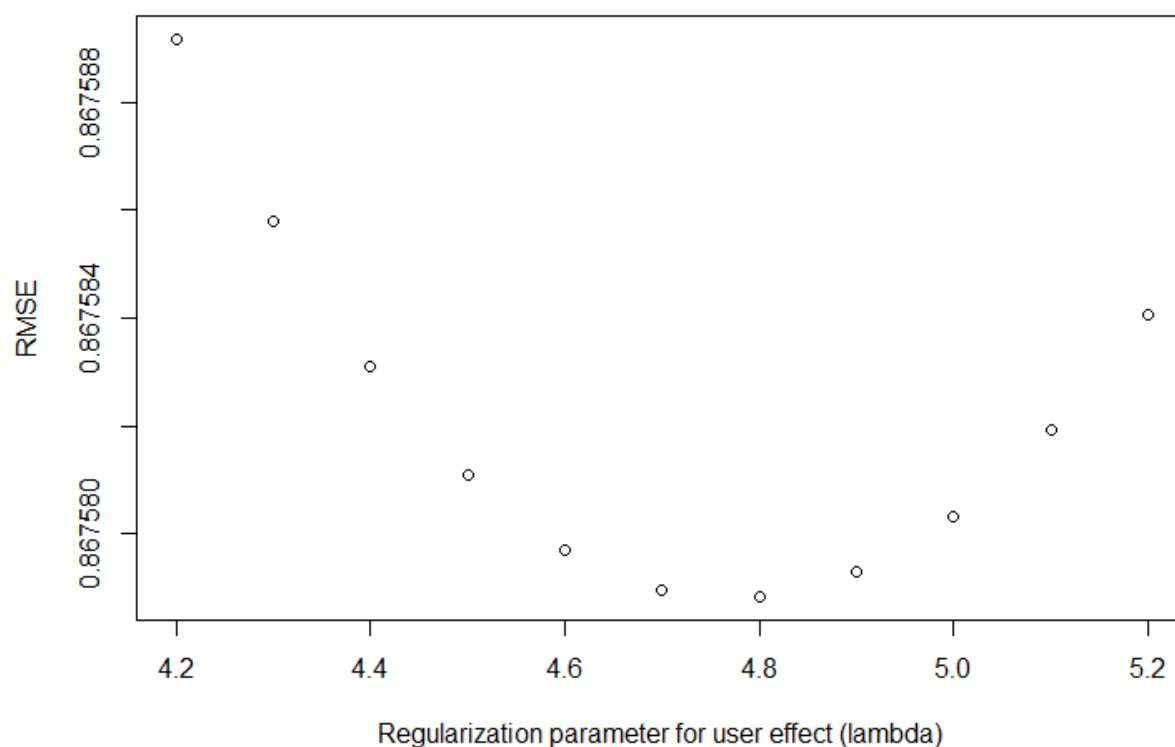
For regularized user effect (b_u), regularization has been done with λ parameter. The best value of λ has been chosen using supply function on λ values from 0:10 (in the code it has been narrowed down for the sake of calculation speed).

Value of b_u has been calculated using following code

```
b_u <- edx_train %>%
  left_join(b_i, by="movieId") %>%
  group_by(userId) %>%
  summarise(b_u=sum(rating-mu-b_i)/(n()+lambda))
```

and the rest of effects (besides movie effect b_i) has been calculated using basic naive approach. Movie effect used for this regularization has been calculated with previously defined α regularization parameter of 3.85.

The best value of λ has been determined as 4.8.



For this value of λ RMSE has been calculated as 0.86758, which is 0.17% improvement from fully naive approach.

2.7.3 Regularized year of movie release effect

For regularized year of movie release effect (b_y), regularization has been done with δ parameter. The best value of δ has been chosen using supply function on δ values from 0:50 (in the code it has been narrowed down for the sake of calculation speed).

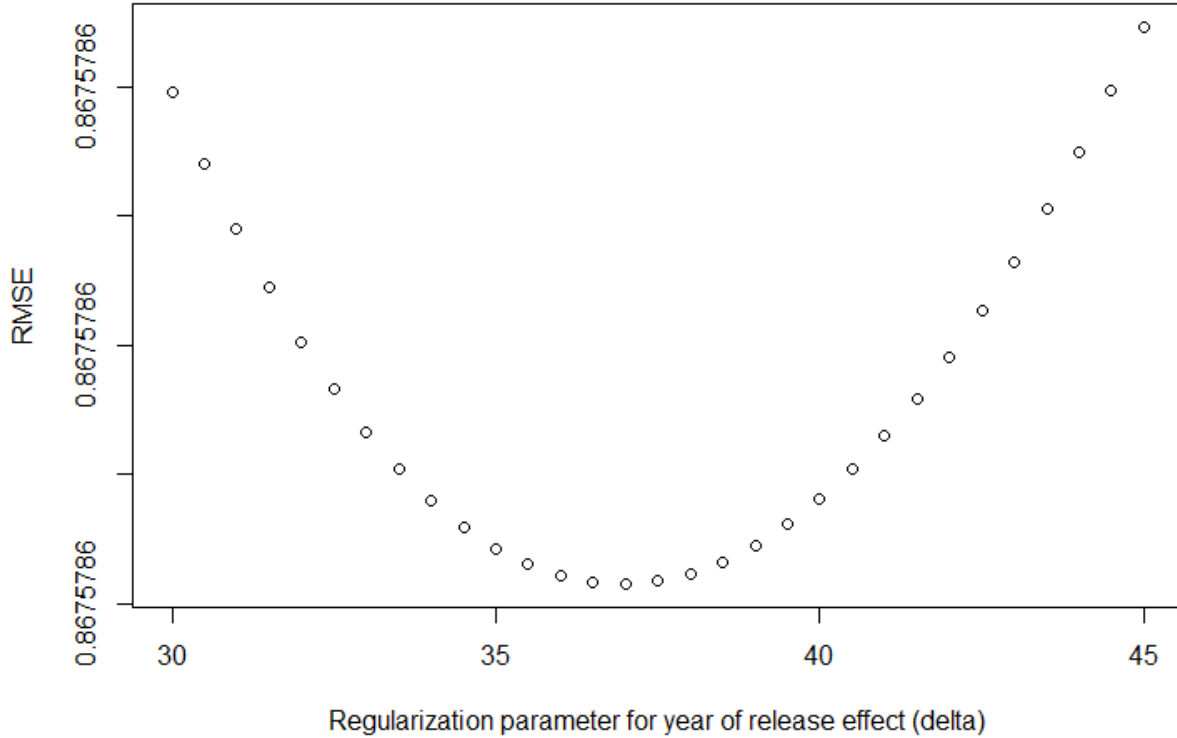
Value of b_y has been calculated using following code

```
b_y <- edx_train %>%
  left_join(b_i, by="movieId") %>%
  left_join(b_u, by="userId") %>%
```

```
group_by(year_of_release) %>%
summarise(b_y=sum(rating-mu-b_i-b_u)/(n()+delta))
```

and the rest of effects (besides movie effect b_i and user effect b_u) has been calculated using basic naive approach. Movie effect used for this regularization has been calculated with previously defined α regularization parameter of 3.85, and user effect used with λ of 4.8.

The best value of δ has been determined as 37.



For this value of δ RMSE has been calculated as 0.86758, which is 0.17% improvement from fully naive approach (and not much different from previous step).

2.7.4 Regularized rating day of the week effect

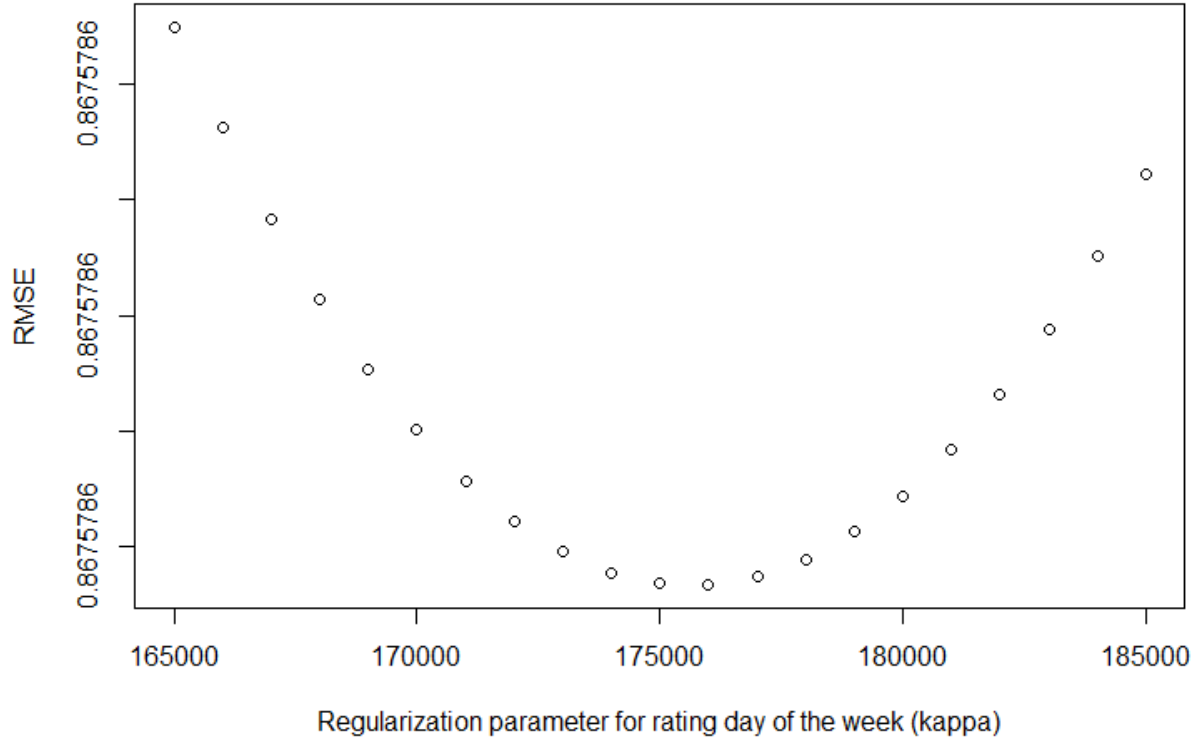
For regularized rating day of the week effect (b_d), regularization has been done with κ parameter. The best value of κ has been chosen using supply function on κ values from 0:1,000,000 (in the code it has been narrowed down for the sake of calculation speed).

Value of b_d has been calculated using following code

```
b_d <- edx_train %>%
left_join(b_i, by="movieId") %>%
left_join(b_u, by="userId") %>%
left_join(b_y, by="year_of_release") %>%
group_by(rateday) %>%
summarise(b_d=sum(rating-mu-b_i-b_u-b_y)/(n()+kappa))
```


and the rest of effects followed previous approach. Genre effect has been calculated using naive approach, and movie, user and year of movie release effects used for this regularization has been calculated with previously defined α (3.85), λ (4.8) and δ (37) regularization parameters.

The best value of κ has been determined as 176,000.



For this value of κ RMSE has been calculated as 0.86758, which is 0.17% improvement from fully naive approach (and only slightly different from two previous steps).

2.7.5 Regularized movie genre effect

For regularized movie genre effect (b_g), regularization has been done with ω parameter. The best value of ω has been chosen using supply function on ω values from 0:100 (in the code it has been narrowed down for the sake of calculation speed).

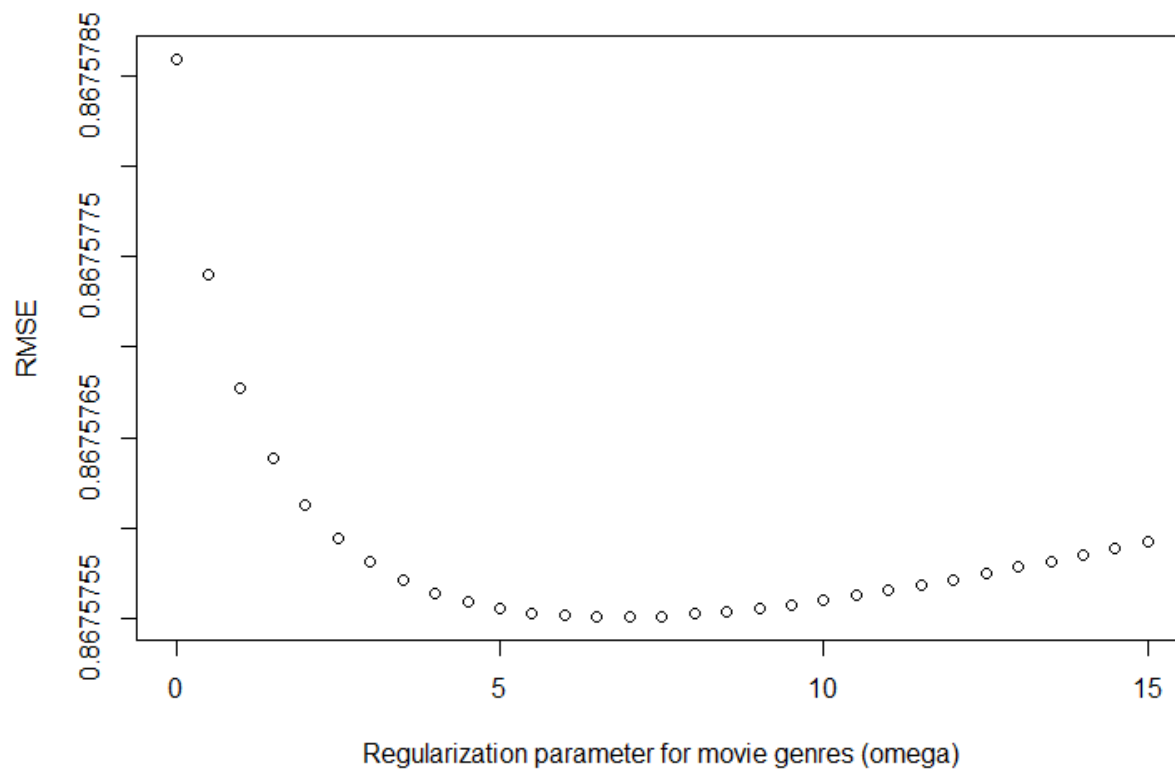
Value of b_g has been calculated using following code

```
b_g <- edx_train %>%
  left_join(b_i, by="movieId") %>%
  left_join(b_u, by="userId") %>%
  left_join(b_y, by="year_of_release") %>%
  left_join(b_d, by="rateday") %>%
  group_by(genres) %>%
  summarise(b_g=sum(rating-mu-b_i-b_u-b_y)/(n()+omega))
```

The rest of effects were calculated using their regularized versions. Movie, user, year of movie release and

rating day of the week effects were using previously defined for them the best regularization parameters values, with α as 3.85, λ as 4.8, δ as 37 and κ as 176,000.

The best value of ω has been determined as 6.5.



For this value of ω RMSE has been calculated as 0.86758, which is 0.17% improvement from fully naive approach (and only slightly different from last three steps).

2.8 Results on edX test set

3 Results

0.864257089475457

4 Conclusions

5 Appendix

5.1 Session info

5.2 Full code