

# Harvard PH125.9x Capstone: MovieLens Analysis

Mauro Berlanda

April 2020

## Introduction/Overview/Executive Summary

The goal of this project is to develop an algorithm to predict movie ratings.

This analysis is performed on ([MovieLens](#)) dataset. It contains 10 millions ratings and 100 000 tags applications applied to 10 000 movies by 72 000 users.

After downloading the entire dataset, two dataframes are created:

- **edx** which is used to perform the analysis and develop the algorithm to predict ratings
- **validation** which contains the true values of the predictions

The original datasets provided by the exercise are:

Table 1: Default summary for edx dataframe

	Length	Class	Mode
userId	9000055	-none-	numeric
movieId	9000055	-none-	numeric
rating	9000055	-none-	numeric
timestamp	9000055	-none-	numeric
title	9000055	-none-	character
genres	9000055	-none-	character

Table 2: Default summary for validation dataframe

	Length	Class	Mode
userId	999999	-none-	numeric
movieId	999999	-none-	numeric
rating	999999	-none-	numeric
timestamp	999999	-none-	numeric
title	999999	-none-	character
genres	999999	-none-	character

Table 3: Head of edx dataframe

	userId	movieId	rating	timestamp	title	genres
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
7	1	355	5	838984474	Flintstones, The (1994)	Children Comedy Fantasy

Table 4: Summary of edx dataframe

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

The greatest challenge in analysing these data is their volume. The most of the built-in models distributed via [caret](#) will simply run out of memory.

This dataset was already used for a test regarding date parsing in the module [HarvardX: PH125.6x Data Science: Wrangling](#). The module [HarvardX: PH125.8x Data Science: Machine Learning](#) illustrated two important concept using this dataset:

- [Recommendation systems](#): use ratings that users have already given to some items in order to make specific recommendations
- [Regularization](#): add information in order to solve an ill-posed problem or to prevent overfitting

The key steps required to identify an optimal solution are:

- prepare the data for the analysis (normalize some columns, extract new features to prepare the analysis)
- create a train and a test set partition out of `edx` dataset to start and validate the analysis
- explore the content of `train_set` and visualize the distribution of the outcome (`rating`) by every feature
- use `train_set` implement several prediction models (see Method/Analysis section for details)
- run models and ensembles against `test_set`
- evaluate the accuracy and the RMSE of each solution based on `test_set` (see Result for the output)
- report the RMSE against the true values to complete the exercise (RMSE < 0.86490 required)

The quality of the prediction is evaluated using the root mean squared error (RMSE):

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_i (\hat{y}_i - y_i)^2}$$

```
RMSE <- function(true_ratings, predicted_ratings){
  sqrt(mean((true_ratings - predicted_ratings)^2))
}
```

## Method/Analysis

The process of this analysis includes data wrangling, data exploration and the description of the modeling approach.

### Data Wrangling/Cleaning

The `edx` and `validation` are the result of some manipulation on a downloaded dataset from `grouplens.org`.

1. download the zip file in a temporary directory and unzip it into the workspace directory
2. read the raw data from two `.dat` files creating `ratings` and `movies` dataframes
3. normalize column names and cast value as `numeric` and `character`, then join both datasets into the `movielens` dataframe
4. partition the data in `movielens` to create a train test used for developing the model (`edx`, 90%) and the true values used to calculate the RMSE of the final model (`validation`, 10%)
5. ensure that `userId` and `movieId` in validation set are also in `edx` set

```
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")

# 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")
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)
```

Before starting the analysis, we are going to prepare the data for the analysis:

```
wrangle_data <-function(ds) {  
  ds %>%  
    mutate(  
      userId = factor(userId), # int to factor  
      movieId = factor(movieId), # num to factor  
      datetime = as_datetime(timestamp), # parse int into date  
      movieYear = str_sub(title,-5,-2), # extract movieYear from title (xxxx),  
      reviewYear = factor(year(datetime)),  
      reviewWeek = factor(week(datetime))  
    )  
}
```

Instead of using all the known observations from `edx` dataset, we create a `train_set` and a `test_set` to continue this analysis as follows:

```
set.seed(17, sample.kind="Rounding")  
test_index <- createDataPartition(y = edx$rating, times = 1, p = 0.1, list = FALSE)  
train_set <- edx[-test_index,]  
temp <- edx[test_index,]  
  
# Ensure that userId and movieId in the test_set are included in the train_set  
test_set <- temp %>%  
  semi_join(train_set, by = "movieId") %>%  
  semi_join(train_set, by = "userId")  
  
removed <- anti_join(temp, test_set)
```

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

```
train_set <- rbind(train_set, removed)
```

```
rm(test_index, temp, removed)
```

We perform our analysis using `train_set` and we validate it with `test_set`.

## Data Exploration and Insights

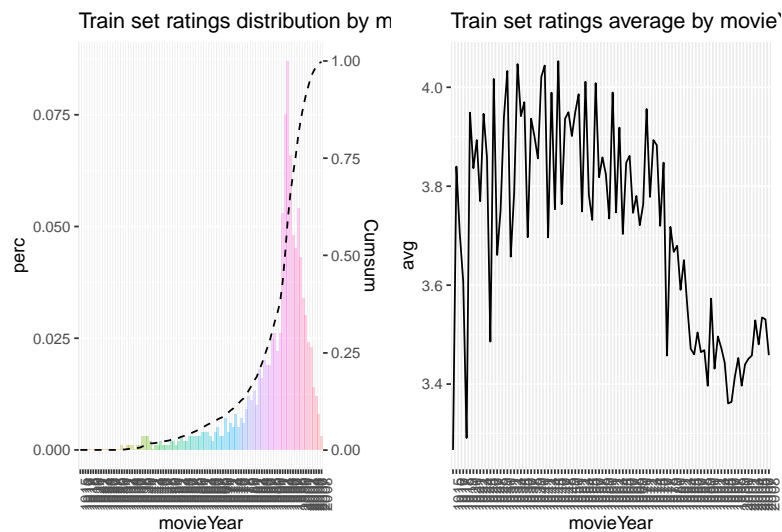
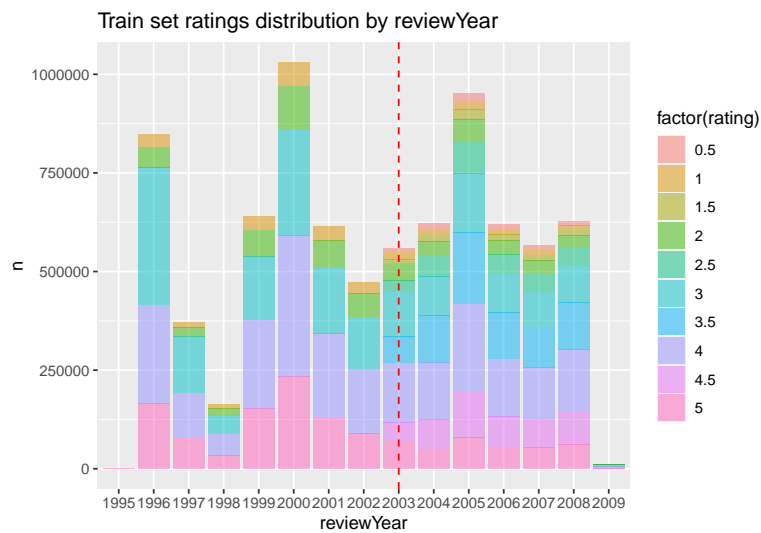
The features we may use in our analysis:

```
setdiff(colnames(edx), c("rating", "title", "timestamp"))
```

```
## [1] "userId"      "movieId"     "genres"      "datetime"    "movieYear"  
## [6] "reviewYear" "reviewWeek"
```

Since our goal is to predict ratings, we start looking at the overall ratings distribution in the `train_set` dataset.

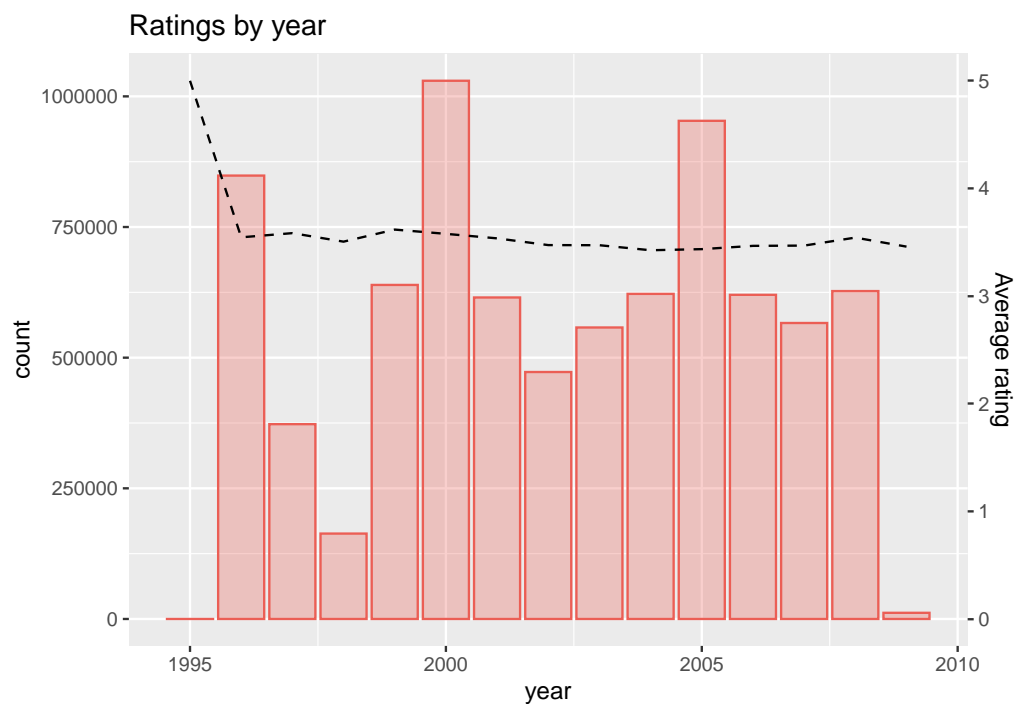
##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.500	3.000	4.000	3.513	4.000	5.000



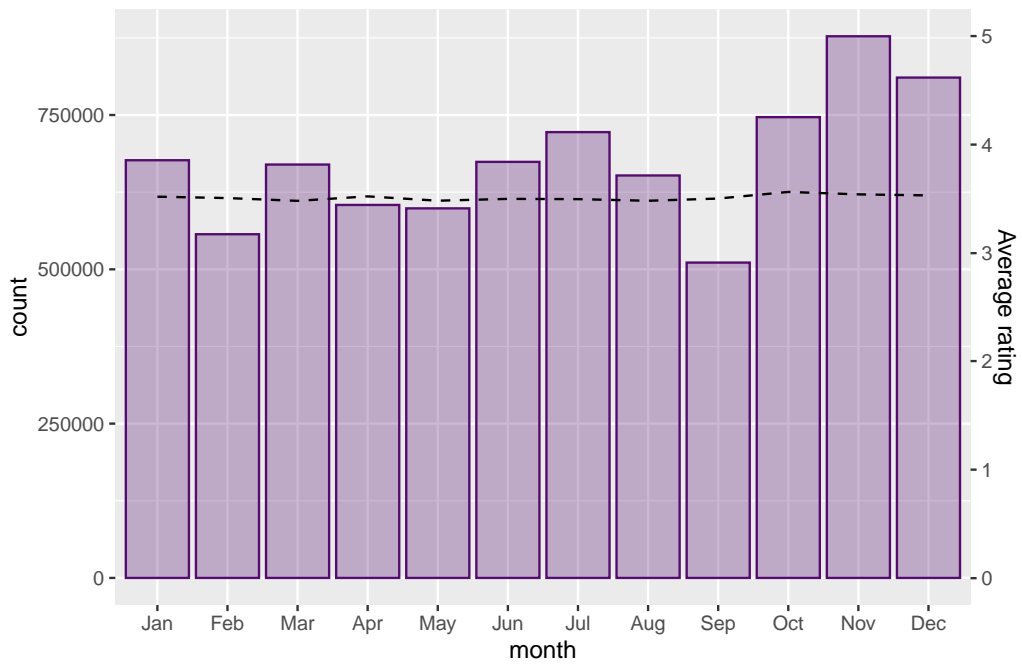
```
## movieYear      avg
## Length:94      Min.   :3.267
## Class :character 1st Qu.:3.511
## Mode  :character Median :3.748
##                  Mean   :3.721
##                  3rd Qu.:3.898
##                  Max.   :4.053
```

We notice that whole star ratings are predominant compares to half star rating. One reason to explain the importance of this difference is that half star ratings have been introduced in 2003. We can also see that a movie is more likely to receive a rating if it was shot between the early 90s and 2000s. On average, movies before the 80s tend to have an higher rating.

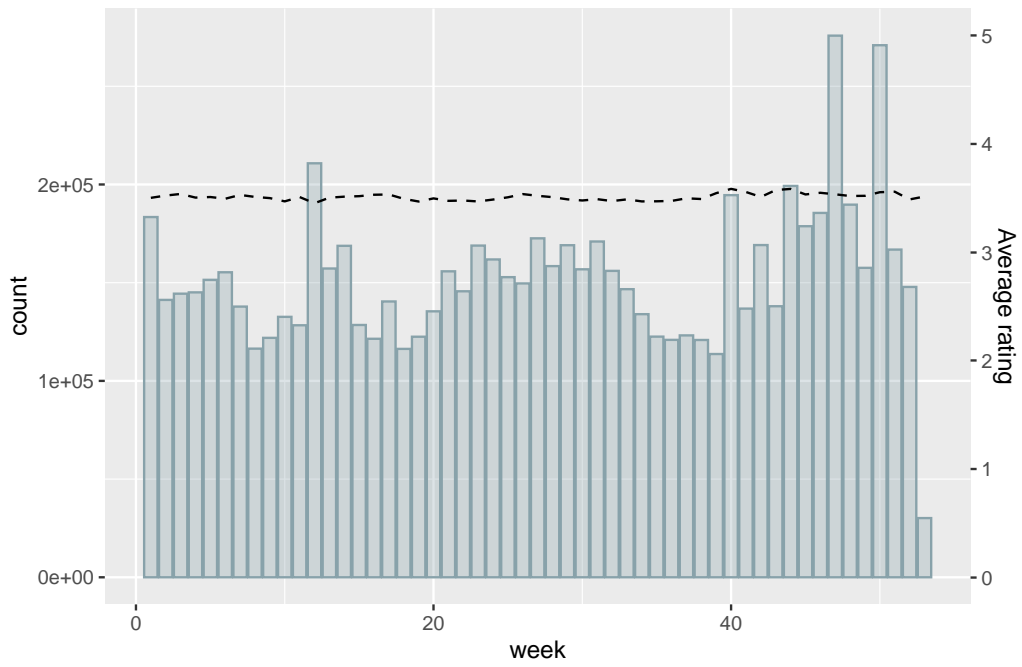
Breaking down the distribution by the review **timestamp**, even if the rating count varies somehow over the time, the average rating seems quite constant.

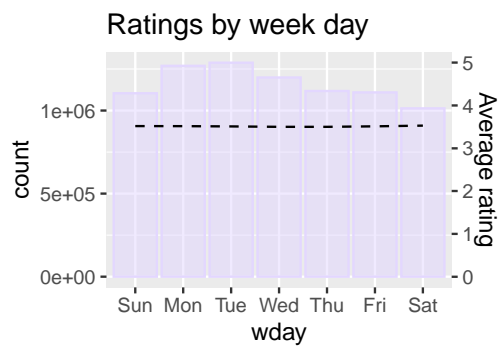
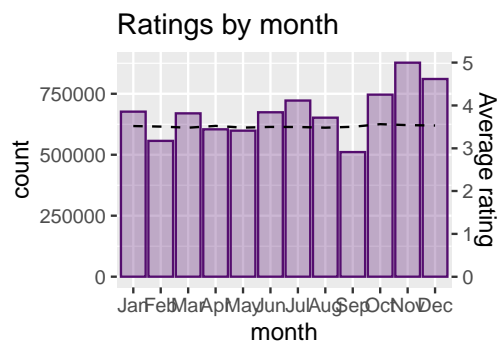
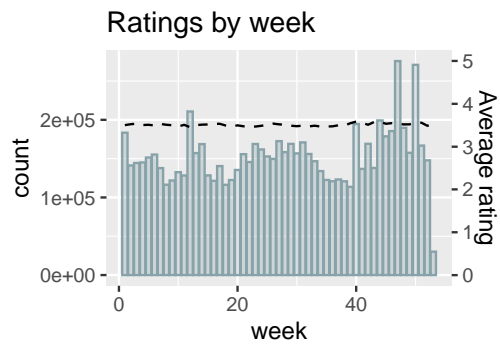
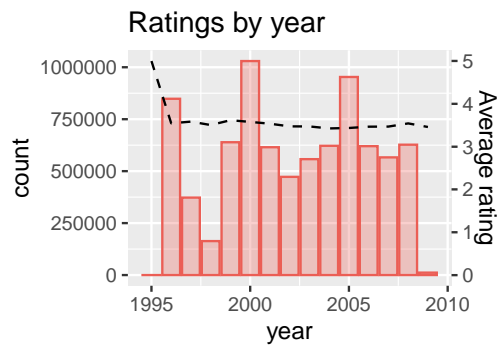
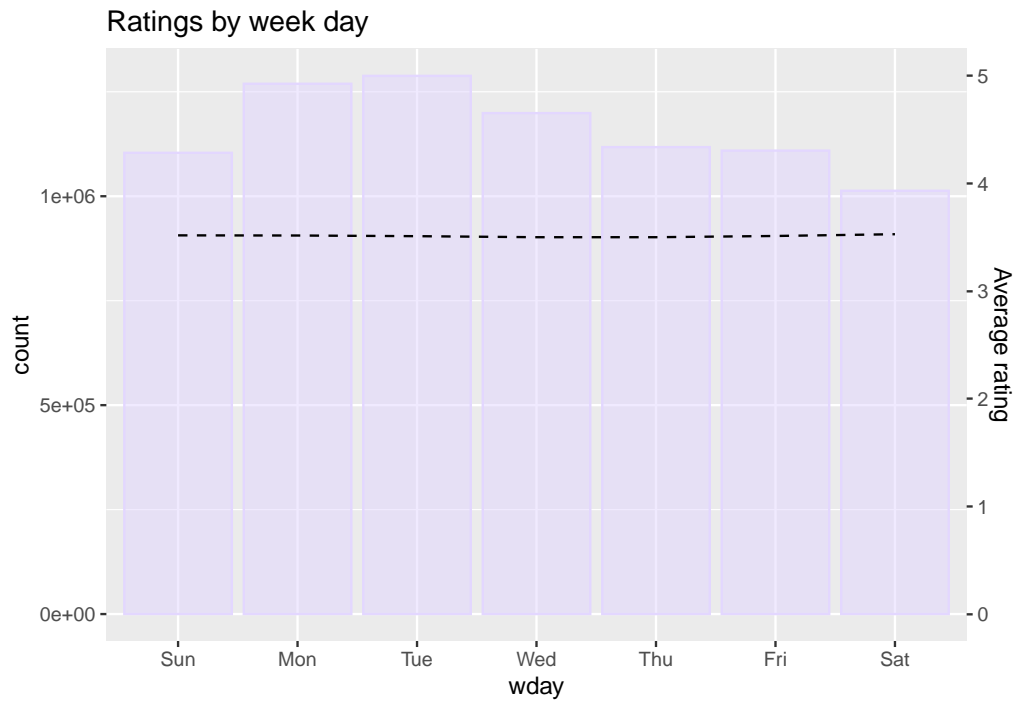


Ratings by month



Ratings by week





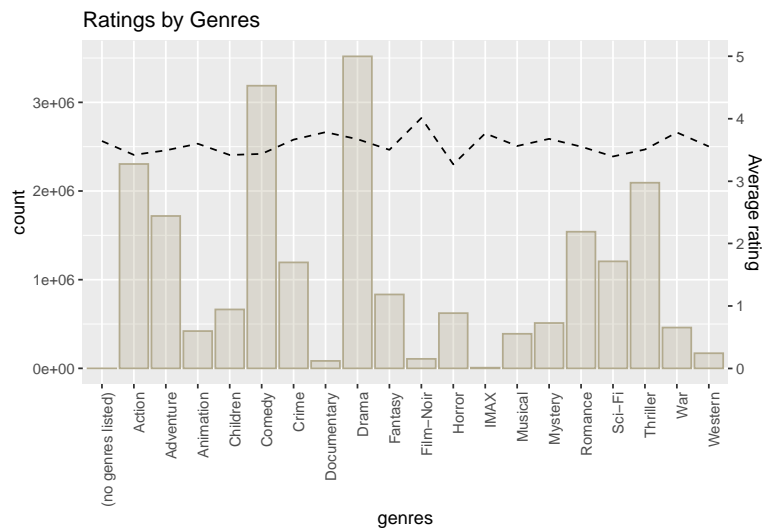
We can clearly see that on average ratings is constant across the **month** and the **wday**. There is some seasonality in the reviews distribution, with a majority of ratings given in the last quarter of the year.

The ratings breakdown by top 10 genres is:



Table 5: Top rated movie genres

genres	count	avg
Drama	3518593	3.673084
Comedy	3187624	3.436912
Action	2304788	3.421162
Thriller	2093364	3.507701
Adventure	1718289	3.493975
Romance	1540803	3.553572
Sci-Fi	1206731	3.395655
Crime	1194688	3.665547
Fantasy	833106	3.502179
Children	664169	3.419442



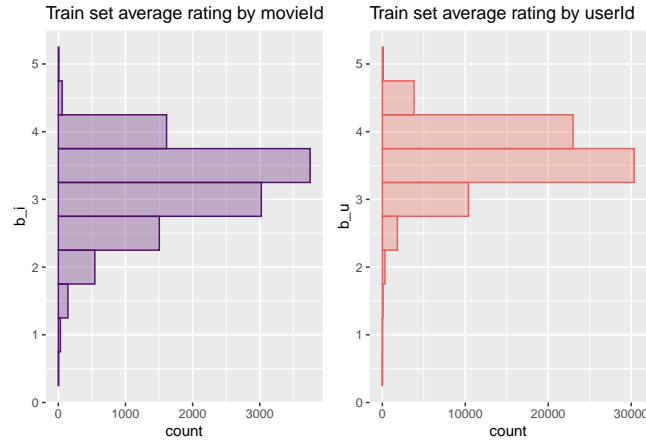
Comparing the distribution of rating by `movieId` and by `userId`

```
n_distinct(train_set$movieId)
```

```
## [1] 10677
```

```
n_distinct(train_set$userId)
```

```
## [1] 69878
```



Out of curiosity, the 10 most rated movies are:

Table 6: Top rated movies

movieId	title	count	avg_rating
296	Pulp Fiction (1994)	28201	4.158611
356	Forrest Gump (1994)	27912	4.014528
593	Silence of the Lambs, The (1991)	27386	4.202585
480	Jurassic Park (1993)	26434	3.662821
318	Shawshank Redemption, The (1994)	25207	4.452731
110	Braveheart (1995)	23631	4.081715
589	Terminator 2: Judgment Day (1991)	23423	3.930432
457	Fugitive, The (1993)	23390	4.008743
260	Star Wars: Episode IV - A New Hope (a.k.a. Star Wars) (1977)	23075	4.221105
150	Apollo 13 (1995)	21842	3.886205

This data exploration has been done within the limits of my machine capacity. More powerful machines can help to find out more advanced patterns.

## Modeling

In the first data exploration we could not find identify a single feature allowing to predict clearly.

As I first step we tried to solve the problem as a **classification** problem with the following outcomes:

```
seq(0.5, 5, 0.5)
```

```
## [1] 0.5 1.0 1.5 2.0 2.5 3.0 3.5 4.0 4.5 5.0
```

With this result normalization, every result could have been converted to the real values levels:

```
n <- 2.7
ceiling(n*2)/2
```

In *HarvardX: PH125.8x* we have seen that if we want to handle more than two features we can use regression trees or random forest.

Running a regression tree on all these dimensions was not really helpful:

```
# Regression tree model
rt_model <- rpart(rating ~ ., data = train_set)
plot(rt_model, margin = 0.1)
text(rt_model, cex = 0.75)

# Warning message:
# In labels.rpart(x, minlength = minlength) :
#   more than 52 levels in a predicting factor, truncated for printout
```

Other approaches like randomForrest will just run out of memory.

Models such as logistic\_regression, lda, qda, loess, knn perform better with at maximum two features. We may just elect two features or use principal component analysis to define two features.

In this report we are modeling the solution taking inspiration from the approach suggested by the [team](#) winning the [Netflix \\$1 million award](#) in 2009.

The solution suggested was based on two models:

*a. Normalization of global effects*

- consider a baseline rating (overall mean)
- a user specific effect based on the previous ratings
- a movie specific effect based on the previous ratings
- a specific interaction (e.g. I like this movie because a given actor is playing)

*b. Neighborhood Models*

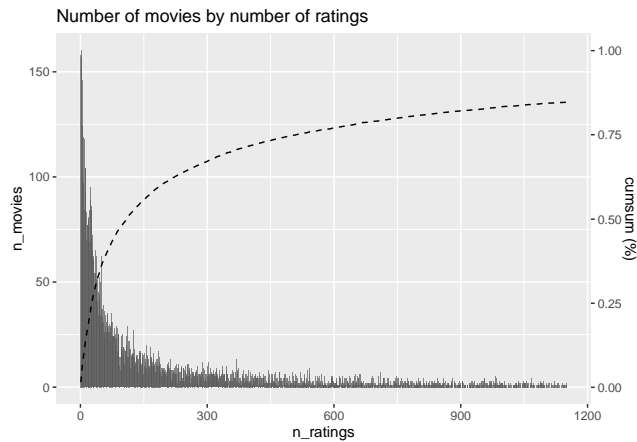
- item-item approach: similar items may receive similar ratings
- user-user approach: similar users may give similar ratings

In our model we are going to regularize the data using *penalized least square* to constraint the total variability of the effect sizes.

As we can see most of the items have very few ratings:

```
movies_reviews_count <- train_set %>% group_by(movieId) %>% count() %>% as_tibble() %>%
  group_by(n) %>% count() %>% rename(n_ratings=n, n_movies=nn)
movies_reviews_count <- movies_reviews_count %>%
  add_column(cumsum=cumsum(movies_reviews_count$n_movies/sum(movies_reviews_count$n_movies)))

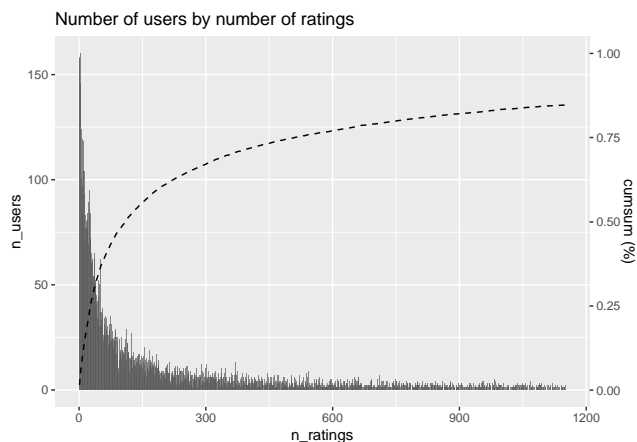
movies_reviews_count %>% head(1000) %>%
  ggplot(aes(x=n_ratings)) +
  geom_histogram(aes(y = n_movies), stat = "identity", alpha = 0.7) +
  geom_line(aes(y=cumsum * max(movies_reviews_count$n_movies)), linetype = "dashed") +
  scale_y_continuous(sec.axis = sec_axis(~./max(movies_reviews_count$n_movies), name = "cumsum (%)")) +
  ggtitle("Number of movies by number of ratings")
```



```
rm(movies_reviews_count)
```

```
users_reviews_count <- train_set %>% group_by(movieId) %>% count() %>% as_tibble() %>%
  group_by(n) %>% count() %>% rename(n_ratings=n, n_users=nn)
users_reviews_count <- users_reviews_count %>%
  add_column(cumsum=cumsum(users_reviews_count$n_users/sum(users_reviews_count$n_users)))
```

```
users_reviews_count %>% head(1000) %>%
  ggplot(aes(x=n_ratings)) +
  geom_histogram(aes(y = n_users), stat = "identity", alpha = 0.7) +
  geom_line(aes(y=cumsum * max(users_reviews_count$n_users)), linetype = "dashed") +
  scale_y_continuous(sec.axis = sec_axis(~./max(users_reviews_count$n_users), name = "cumsum (%)")) +
  ggtitle("Number of users by number of ratings")
```



```
rm(users_reviews_count)
```

In order to identify the `lambda` parameter (penalty term) to apply to our model, we are going to use the `test_set` we extracted earlier.

It is important to note that all users and movies both in `test_set` and in `validation` are included by design in the `train_set`. This allows us to stretch the impact of the section b. of the winning algorithm introduced above.

To identify an acceptable solution ( $RMSE < 0.86490$ ), we first define the RMSE function:

```
RMSE <- function(true_ratings, predicted_ratings){  
  sqrt(mean((true_ratings - predicted_ratings)^2))  
}  
  
mu <- mean(train_set$rating)
```

1. naive approach: use the average

A first naive approach could be just to take the average rating and predict all the results with it:

```
predict_model_0 <- function(df) {  
  rep(mu, nrow(df))  
}  
naive_rmse <- RMSE(test_set$rating, predict_model_0(test_set))  
rmse_results <- tibble(method = "Model 0: Average", RMSE = naive_rmse)
```

2. Introduce the movie effect

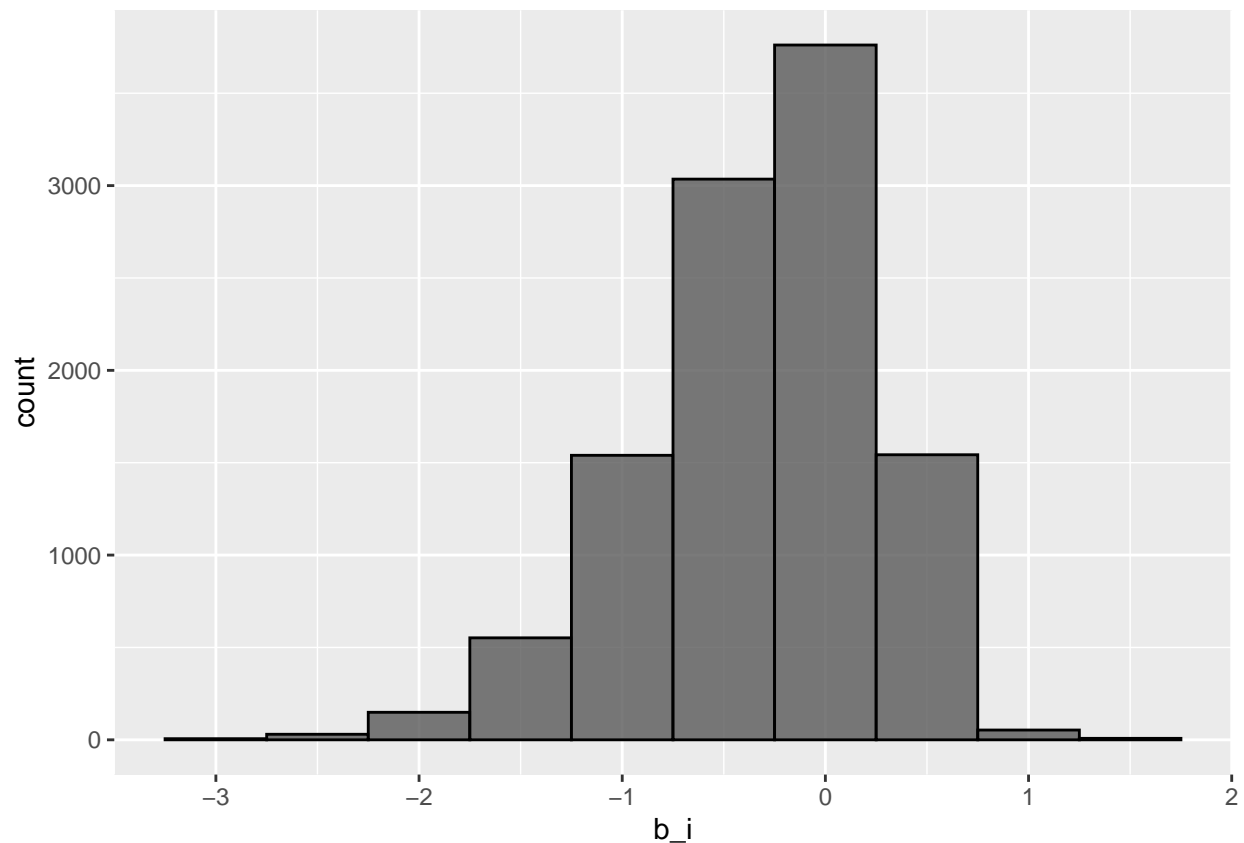
```
# fit <- lm(rating ~ movieId, data = train_set)  
# Error: vector memory exhausted (limit reached?)
```

On a regular machine is not possible to train a linear model on by movieId. There are two options:

- reduce the size of the `train_set`
- avoid using some vendor library and write the model by ourselves

Following the second approach:

```
# regularized movie averages  
reg_movie_avgs <- train_set %>%  
  group_by(movieId) %>%  
  summarize(b_i = mean(rating - mu))  
  
ggplot(aes(x=b_i), data=reg_movie_avgs) + geom_histogram(bins=10, alpha=0.8, color="black")
```



```
predict_model_1 <- function(df) {
  mu + df %>% left_join(reg_movie_avgs, by='movieId') %>% .$b_i
}

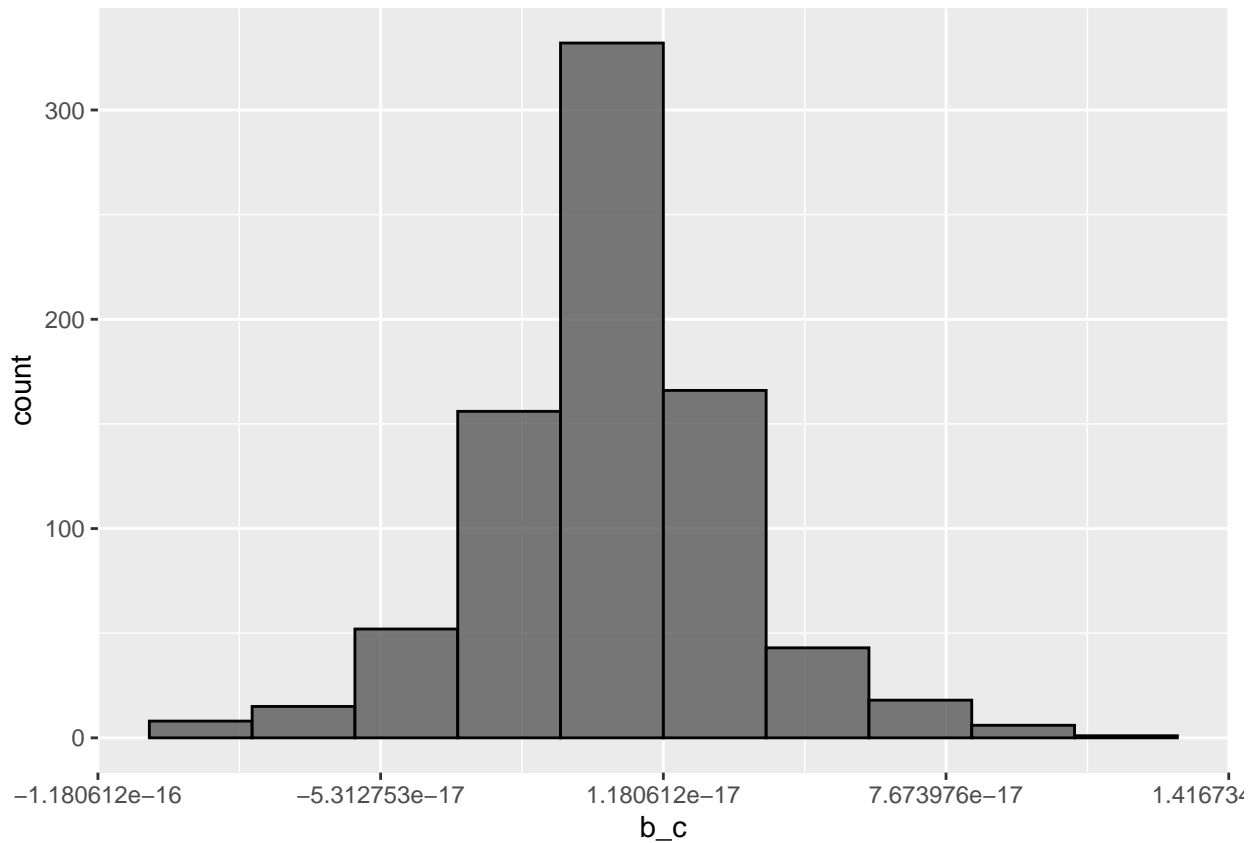
model_1_rmse <- RMSE(predict_model_1(test_set), test_set$rating)
rmse_results <- bind_rows(rmse_results,
  data_frame(method="Model 1: Movie Effect",
    RMSE = model_1_rmse ))
```

3. Introduce the genres effect -> don't

In the model b was described an item-item approach. In order to group similar movies we may used the `genres` column provided in the datasets.

```
## regularized genres averages
reg_genres_avgs <- train_set %>%
  left_join(reg_movie_avgs, by='movieId') %>%
  group_by(genres) %>%
  summarize(b_c = mean(rating - mu - b_i))

ggplot(aes(x=b_c), data=reg_genres_avgs) + geom_histogram(bins=10, alpha=0.8, color="black")
```



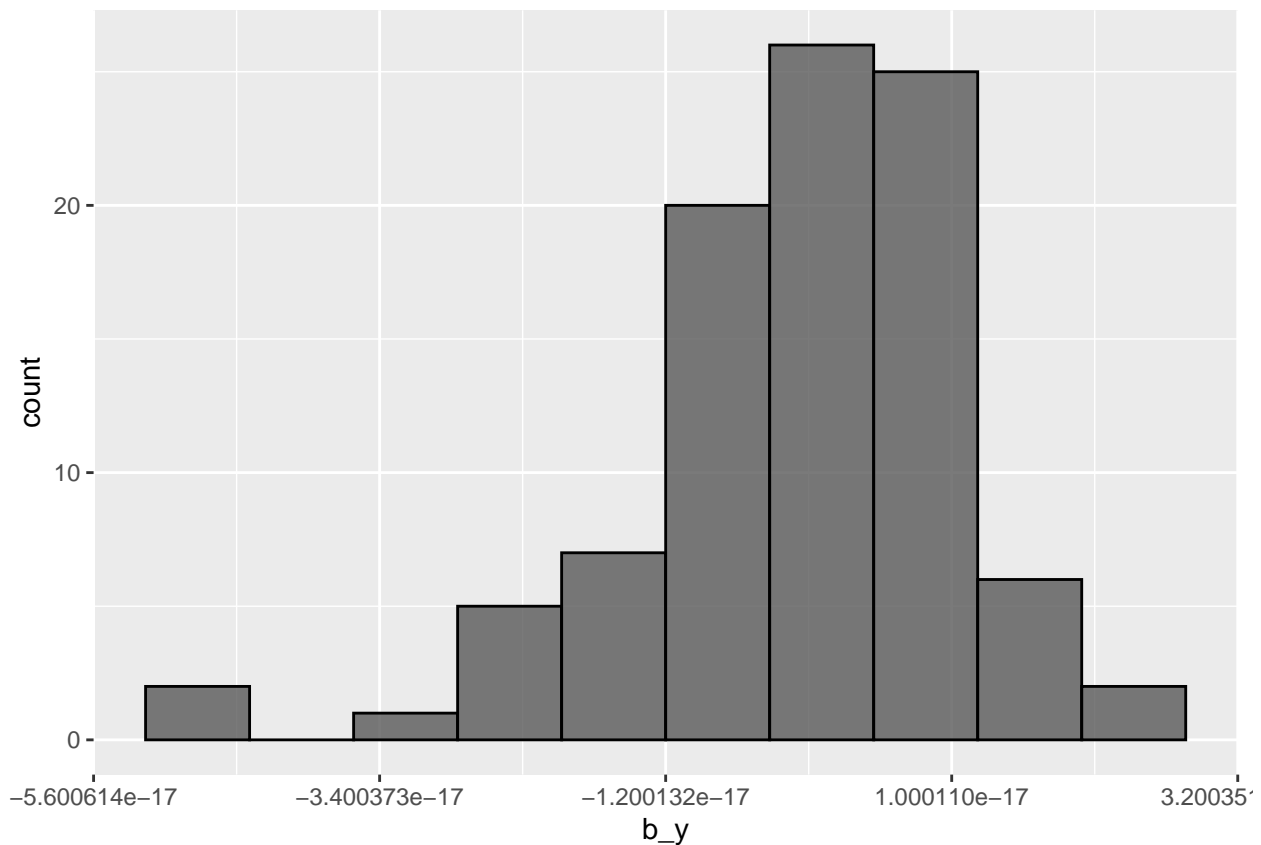
As we can see in the plot, the **genres** effect after the movie effect does not seem to be relevant (affecting the rating only by  $10^{-17}$ ). It is not worth to include it in our model.

4. Introduce the movieYear effect -> don't

Another item-item approach could be explored using **movieYear**.

```
## regularized movieYear averages
reg_movieYear_avgs <- train_set %>%
  left_join(reg_movie_avgs, by='movieId') %>%
  group_by(movieYear) %>%
  summarize(b_y = mean(rating - mu - b_i))

ggplot(aes(x=b_y), data=reg_movieYear_avgs) + geom_histogram(bins=10, alpha=0.8, color="black")
```



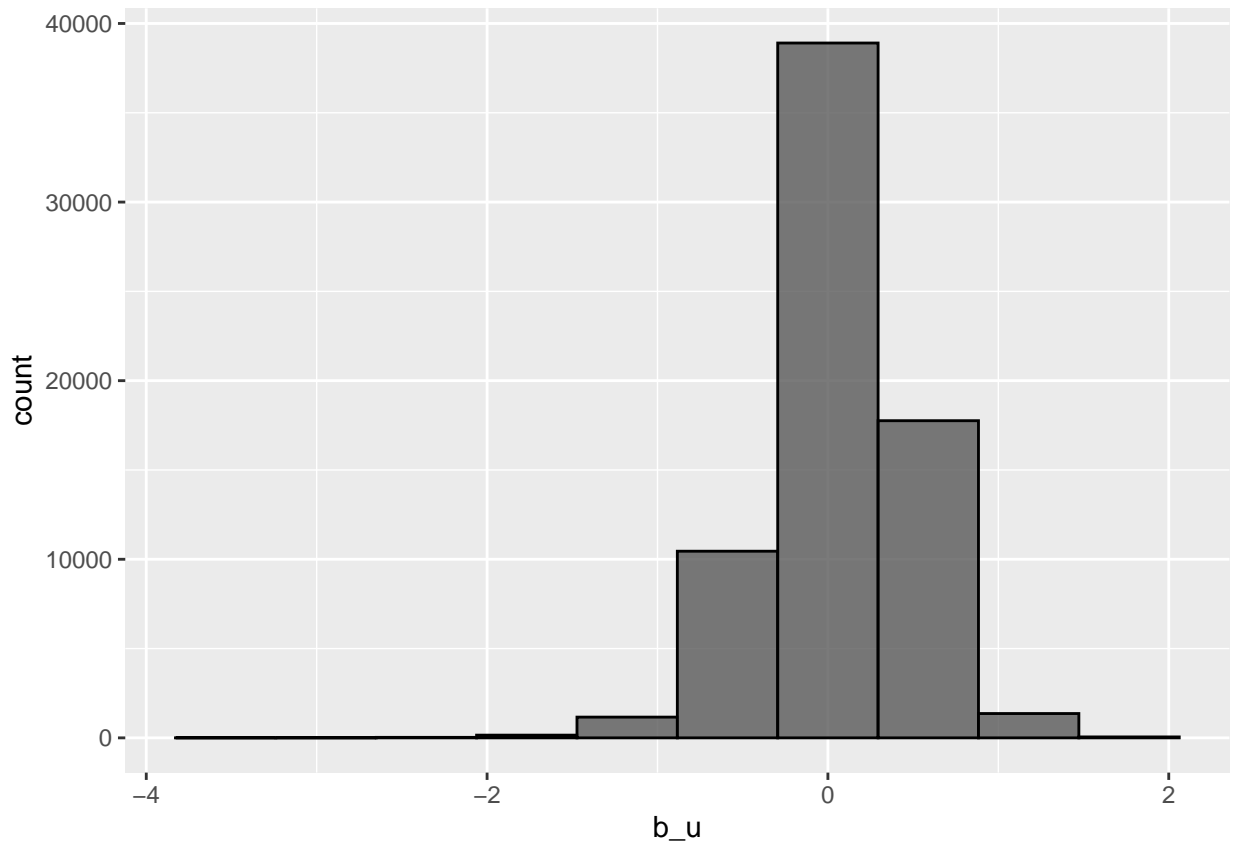
As we can see in the plot, the `movieYear` effect after the movie effect does not seem to be relevant (affecting the rating only by  $10^{-17}$ ). It is not worth to include it in our model.

##### 5. Introduce the `userId` effect

```
## regularized reviewYear averages
reg_user_avgs <- train_set %>%
  left_join(reg_movie_avgs, by='movieId') %>%
  group_by(userId) %>%
  summarize(b_u = mean(rating - mu - b_i))

ggplot(aes(x=b_u), data=reg_user_avgs) + geom_histogram(bins=10, alpha=0.8, color="black")
```





```
predict_model_2 <- function(df) {
  df %>%
    left_join(reg_movie_avgs, by='movieId') %>%
    left_join(reg_user_avgs, by='userId') %>%
    mutate(pred = mu + b_i + b_u) %>% .$pred
}

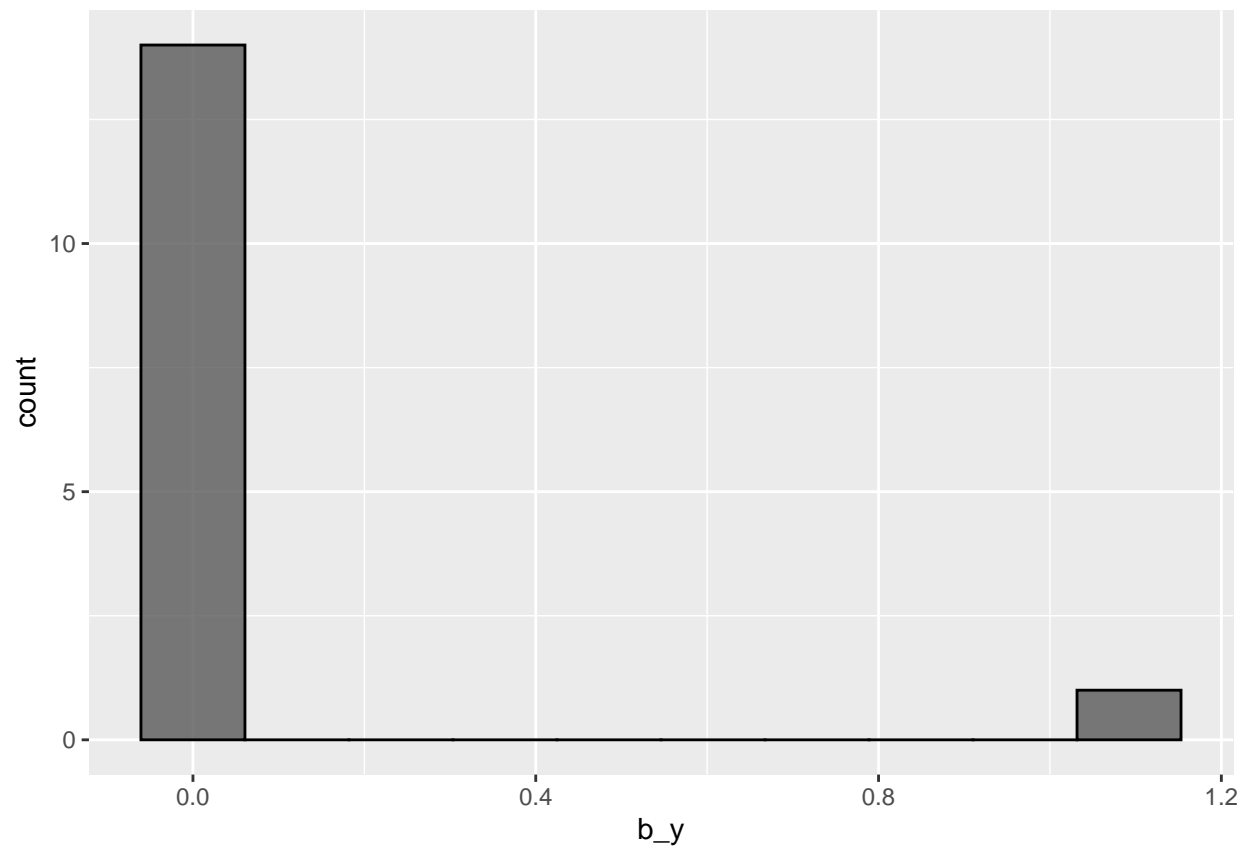
model_2_rmse <- RMSE(predict_model_2(test_set), test_set$rating)
rmse_results <- bind_rows(rmse_results,
  data_frame(method="Model 2: Movie + User Effect",
    RMSE = model_2_rmse )) # 0.8659109
```

## 6. Introduce the reviewYear and reviewWeek effect

We do not have any element to evaluate a user-user approach. On of the few approaches we can imagine is that at a given moment in the time the people tend to rate in a more similar way.

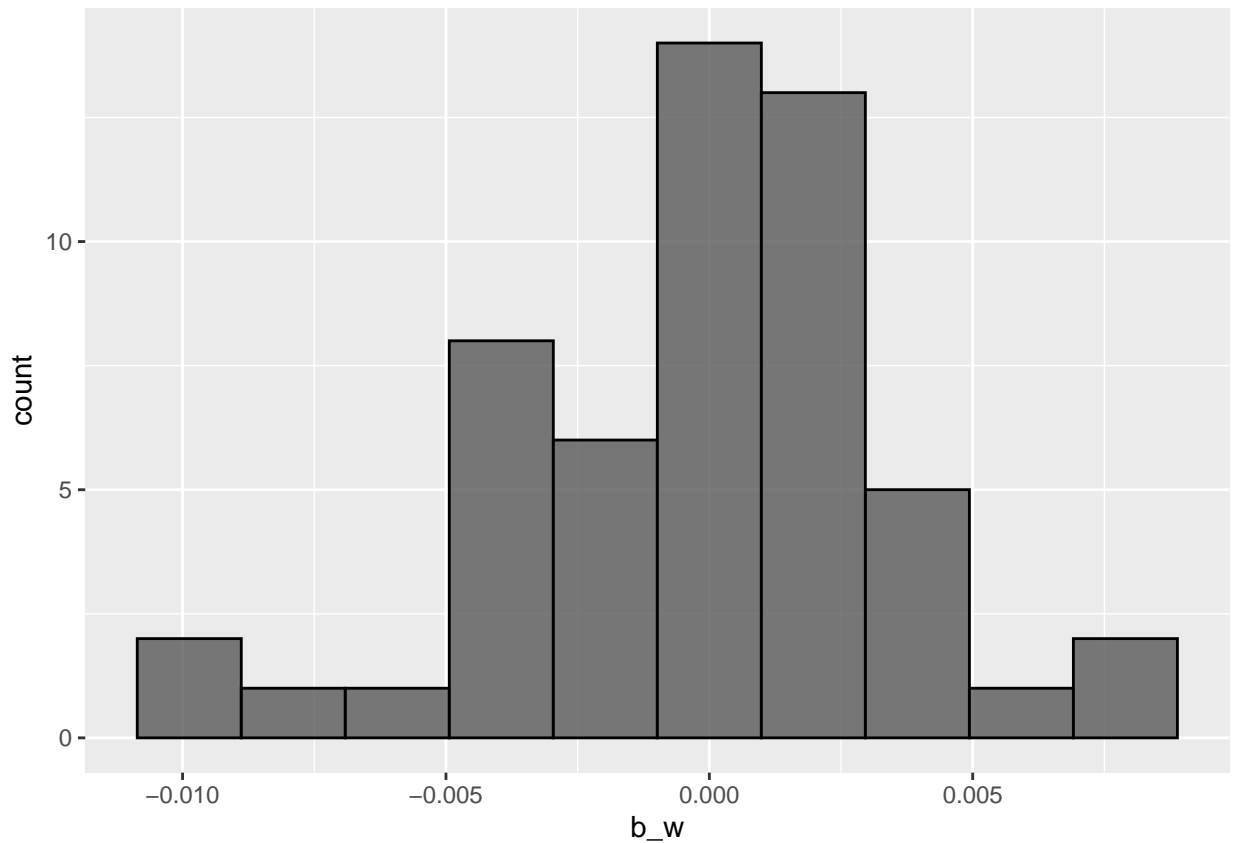
```
## regularized reviewYear averages
reg_reviewYear_avgs <- train_set %>%
  left_join(reg_movie_avgs, by='movieId') %>%
  left_join(reg_user_avgs, by='userId') %>%
  group_by(reviewYear) %>%
  summarize(b_y = mean(rating - mu - b_i - b_u))

ggplot(aes(x=b_y), data=reg_reviewYear_avgs) + geom_histogram(bins=10, alpha=0.8, color="black")
```



```
reg_reviewWeek_avgs <- train_set %>%
  left_join(reg_movie_avgs, by='movieId') %>%
  left_join(reg_user_avgs, by='userId') %>%
  left_join(reg_reviewYear_avgs, by='reviewYear') %>%
  group_by(reviewWeek) %>%
  summarize(b_w = mean(rating - mu - b_i - b_u - b_y))

ggplot(aes(x=b_w), data=reg_reviewWeek_avgs) + geom_histogram(bins=10, alpha=0.8, color="black")
```



```
predict_model_3 <- function(df) {
  df %>%
    left_join(reg_movie_avgs, by='movieId') %>%
    left_join(reg_user_avgs, by='userId') %>%
    left_join(reg_reviewYear_avgs, by='reviewYear') %>%
    # left_join(reg_reviewWeek_avgs, by='reviewWeek') %>%
    mutate(pred = mu + b_i + b_u + b_y) %>% .$pred
}

model_3_rmse <- RMSE(predict_model_3(test_set), test_set$rating) # 0.865905
rmse_results <- bind_rows(rmse_results,
  data_frame(method="Model 3: Movie + User + reviewYearWeek Effect",
    RMSE = model_3_rmse ))
```

Despite the impact is moderate, taking in consideration also `reviewYear` and `reviewWeek` is slightly improving the RMSE.

7. Identify the optimal penalty rate `lambda`

## Results

Let's see the results of the model developed on the `test_set`:

```
#
RMSE(predict_model_3(validation), validation$rating)
```

```
## [1] 0.8658524
```

```
rmse_results %>% kable(caption = "RMSE on test_set")
```

Table 7: RMSE on test\_set

method	RMSE
Model 0: Average	1.0606957
Model 1: Movie Effect	0.9440571
Model 2: Movie + User Effect	0.8662206
Model 3: Movie + User + reviewYearWeek Effect	0.8662164

## Conclusion