

Capstone Movielens Report

Azamat Kurbanaev

2025-05-08

Contents

Introduction / Overview / Executive Summary	1
Methods / Analysis	12
Conclusion	34

Introduction / Overview / Executive Summary

The goal of the project is to build a Recommendation System using a [10M version of the MovieLens dataset](#). Following the [Netflix Grand Prize Contest](#) requirements, we will evaluate the *Root Mean Square Error (RMSE)* score, which, as shown in [Section 23.2 Loss function](#) of the *Course Textbook*, is defined as:

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

with N being the number of user/movie combinations for which we make predictions and the sum occurring over all these combinations[1].

Our goal is to achieve a value of less than 0.86490 (compare with the *Netflix Grand Prize* requirement: of at least 0.8563[2]).

Datasets Overview

To start with we have to generate two datasets derived from the *MovieLens* one mentioned above:

- **edx:** we use it to develop and train our algorithms;
- **final_holdout_test:** according to the course requirements, we use it exclusively to evaluate the *RMSE* of our final algorithm.

For this purpose the following package has been developed by the author of this report: `edx.capstone.movielens.data`. The source code of the package is available [on GitHub](#)[3].

Let's install the development version of this package from the GitHub repository and attach the correspondent library to the global environment:

```
if(!require(edx.capstone.movielens.data)) pak::pak("AzKurban-edX-DS/edx.capstone.movielens.data")

library(edx.capstone.movielens.data)
edx <- edx.capstone.movielens.data::edx
final_holdout_test <- edx.capstone.movielens.data::final_holdout_test
```

Now, we have the datasets listed above:

```
summary(edx)
```

```
##      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
##  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
```

```
summary(final_holdout_test)
```

```
##      userId      movieId      rating      timestamp      title      genres
##  Min.   :    1   Min.   :    1   Min.   :0.500   Min.   :7.897e+08  Length:999999   Length:999999
##  1st Qu.:18096  1st Qu.:  648  1st Qu.:3.000   1st Qu.:9.467e+08  Class  :character  Class  :character
##  Median :35768  Median : 1827  Median :4.000   Median :1.035e+09  Mode   :character  Mode   :character
##  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
```

edx Dataset

Let's look into the details of the edx dataset:

```
str(edx)
```

```
## 'data.frame': 9000055 obs. of 6 variables:
## $ userId   : int 1 1 1 1 1 1 1 1 1 ...
## $ movieId  : int 122 185 292 316 329 355 356 362 364 370 ...
## $ rating   : num 5 5 5 5 5 5 5 5 5 ...
## $ timestamp: int 838985046 838983525 838983421 838983392 838983392 838984474 838983653 838984885 838984885 ...
## $ title    : chr "Boomerang (1992)" "Net, The (1995)" "Outbreak (1995)" "Stargate (1994)" ...
## $ genres   : chr "Comedy|Romance" "Action|Crime|Thriller" "Action|Drama|Sci-Fi|Thriller" "Action|A
```

Note that we have 9000055 rows and six columns in there:

```
dim_edx <- dim(edx)
print(dim_edx)
```

```
## [1] 9000055      6
```

First, let's note that we have 10677 different movies:

```
n_movies <- n_distinct(edx$movieId)
print(n_movies)
```

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

and 69878 different users in the dataset:

```
n_users <- n_distinct(edx$userId)
print(n_users)
```

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

Now, note the expressions below which confirm the fact explained in [Section 23.1.1 MovieLens data](#) of the *Course Textbook*[4] that not every user rated every movie:

```
max_possible_ratings <- n_movies*n_users
sprintf("Maximum possible ratings: %s", max_possible_ratings)
```

```
## [1] "Maximum possible ratings: 746087406"
```

```
sprintf("Rows in `edx` dataset: %s", dim_edx[1])
```

```
## [1] "Rows in 'edx' dataset: 9000055"
```

```
sprintf("Not every movie was rated: %s", max_possible_ratings > dim_edx[1])
```

```
## [1] "Not every movie was rated: TRUE"
```

As also explained in that section, we can think of these data as a very large matrix, with users on the rows and movies on the columns, with many empty cells. Therefore, we can think of a recommendation system as filling in the NAs in the dataset for the movies that some or all the users do not rate. A sample from the edx data below illustrates this idea[5]:

```
keep <- edx |>
  dplyr::count(movieId) |>
  top_n(4, n) |>
  pull(movieId)

tab <- edx |>
  filter(movieId %in% keep) |>
  filter(userId %in% c(13:20)) |>
  select(userId, title, rating) |>
  mutate(title = str_remove(title, ", The"),
        title = str_remove(title, ":.*")) |>
  pivot_wider(names_from = "title", values_from = "rating")

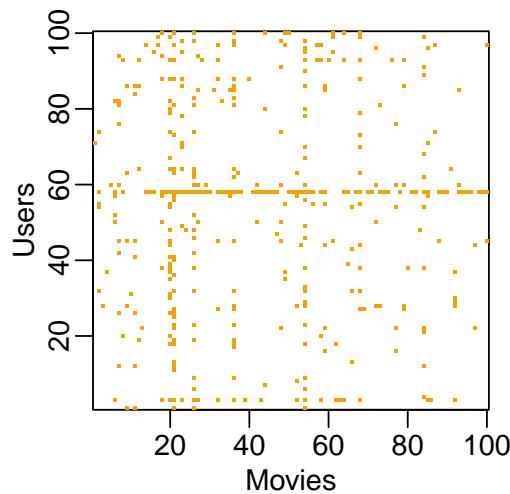
print(tab)

## # A tibble: 5 x 5
##   userId 'Pulp Fiction (1994)' 'Jurassic Park (1993)' 'Silence of the Lambs (1991)' 'Forrest Gump (1994)'
##   <int>          <dbl>           <dbl>            <dbl>           <dbl>
## 1    13             4              NA              NA
## 2    16             NA             3               NA
## 3    17             NA              NA              5
## 4    18             5              3               5
## 5    19             NA              1               NA
```

The following plot of the matrix for a random sample of 100 movies and 100 users with yellow indicating a user/movie combination for which we have a rating shows how *sparse* the matrix is:

```
users <- sample(unique(edx$userId), 100)

rafaelib::mypar()
edx |>
  filter(userId %in% users) |>
  select(userId, movieId, rating) |>
  mutate(rating = 1) |>
  pivot_wider(names_from = movieId, values_from = rating) |>
  (\(mat) mat[, sample(ncol(mat), 100)]())() |>
  as.matrix() |>
  t() |>
  image(1:100, 1:100, z = _, xlab = "Movies", ylab = "Users")
```



Further observations highlighted there that, as we can see from the distributions the author presented, some movies get rated more than others, and some users are more active than others in rating movies:

```
p1 <- edx |>
  count(movieId) |>
  ggplot(aes(n)) +
  geom_histogram(bins = 30, color = "black") +
  scale_x_log10() +
  ggtitle("Movies")

p2 <- edx |>
  count(userId) |>
  ggplot(aes(n)) +
  geom_histogram(bins = 30, color = "black") +
  scale_x_log10() +
  ggtitle("Users")

gridExtra::grid.arrange(p2, p1, ncol = 2)
```



Finally, we can see that no movies have a rating of 0. Movies are rated from 0.5 to 5.0 in 0.5 increments:

```
#library(dplyr)
s <- edx |> group_by(rating) |>
  summarise(n = n())
print(s)

## # A tibble: 10 x 2
##   rating     n
##   <dbl>   <int>
## 1 0.5     85374
## 2 1       345679
## 3 1.5     106426
## 4 2       711422
## 5 2.5     333010
## 6 3       2121240
## 7 3.5     791624
## 8 4       2588430
## 9 4.5     526736
## 10 5      1390114
```

Further analysis of the `edx` dataset have been also inspired by the article mentioned above[6], from which the code and explanatory notes below were cited.

Rating distribution plot[6]

The code below demonstrates another way of visualizing the rating distribution:

```
edx |>
  group_by(rating) |>
  summarise(count = n()) |>
  ggplot(aes(x = rating, y = count)) +
  geom_bar(stat = "identity", fill = "#8888ff") +
  ggtitle("Rating Distribution") +
  xlab("Rating") +
  ylab("Occurrences Count") +
```

```

scale_y_continuous(labels = comma) +
scale_x_continuous(n.breaks = 10) +
theme_economist() +
theme(axis.title.x = element_text(vjust = -5, face = "bold"),
axis.title.y = element_text(vjust = 10, face = "bold"),
plot.margin = margin(0.7, 0.5, 1, 1.2, "cm"))

```



This graph is another confirmation of what we found out above: rounded ratings occur more often than half-stared ones. The upward trend previously discussed is now perfectly clear, although it seems to top right between the 3 and 4-star ratings lowering the occurrences count afterward. That might be due to users being more hesitant to rate with the highest mark for whichever reasons they might hold[6].

Ratings per movie

Movie popularity count[6]

```

print(edx |>
  group_by(movieId) |>
  summarize(count = n()) |>
  slice_head(n = 10)
)

```

```

## # A tibble: 10 x 2
##   movieId count
##   <int> <int>
## 1       1 23790
## 2       2 10779
## 3       3  7028
## 4       4  1577
## 5       5  6400
## 6       6 12346
## 7       7  7259
## 8       8   821
## 9       9  2278
## 10      10 15187

summary(edx |> group_by(movieId) |> summarize(count = n()) |> select(count))

```

```

## #> #> #> count
## #> #> Min.    : 1.0
## #> #> 1st Qu.: 30.0
## #> #> Median  : 122.0
## #> #> Mean    : 842.9
## #> #> 3rd Qu.: 565.0
## #> #> Max.    :31362.0

```

Ratings per movie plot[6]

```

edx |>
  group_by(movieId) |>
  summarize(count = n()) |>
  ggplot(aes(x = movieId, y = count)) +
  geom_point(alpha = 0.2, color = "#4020dd") +
  geom_smooth(color = "red") +
  ggtitle("Ratings per movie") +
  xlab("Movies") +
  ylab("Number of ratings") +
  scale_y_continuous(labels = comma) +
  scale_x_continuous(n.breaks = 10) +
  theme_economist() +
  theme(axis.title.x = element_text(vjust = -5, face = "bold"),
        axis.title.y = element_text(vjust = 10, face = "bold"),
        plot.margin = margin(0.7, 0.5, 1, 1.2, "cm"))

## `geom_smooth()` using method = 'gam' and formula = 'y ~ s(x, bs = "cs")'

```



Movies' rating histogram[6]

```
edx |>
  group_by(movieId) |>
  summarize(count = n()) |>
  ggplot(aes(x = count)) +
  geom_histogram(fill = "#8888ff", color = "#4020dd") +
  ggtitle("Movies' rating histogram") +
  xlab("Rating count") +
  ylab("Number of movies") +
  scale_y_continuous(labels = comma) +
  scale_x_log10(n.breaks = 10) +
  theme_economist() +
  theme(axis.title.x = element_text(vjust = -5, face = "bold"),
        axis.title.y = element_text(vjust = 10, face = "bold"),
        plot.margin = margin(0.7, 0.5, 1, 1.2, "cm"))

## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Ratings per user[6]

User rating count (activity measure)

```
print(edx |>
  group_by(userId) |>
  summarize(count = n()) |>
  slice_head(n = 10)
)
```

```
## # A tibble: 10 x 2
##   userId count
##     <int> <int>
## 1      1    19
## 2      2    17
## 3      3    31
## 4      4    35
## 5      5    74
## 6      6    39
## 7      7    96
## 8      8   727
## 9      9    21
```

```
## 10      10     112
```

User rating summary

```
summary(edx |> group_by(userId) |> summarize(count = n()) |> select(count))
```

```
##      count
##  Min.   : 10.0
##  1st Qu.: 32.0
##  Median : 62.0
##  Mean   : 128.8
##  3rd Qu.: 141.0
##  Max.   :6616.0
```

Ratings per user plot

```
edx |>
  group_by(userId) |>
  summarize(count = n()) |>
  ggplot(aes(x = userId, y = count)) +
  geom_point(alpha = 0.2, color = "#4020dd") +
  geom_smooth(color = "red") +
  ggtitle("Ratings per user") +
  xlab("Users") +
  ylab("Number of ratings") +
  scale_y_continuous(labels = comma) +
  scale_x_continuous(n.breaks = 10) +
  theme_economist() +
  theme(axis.title.x = element_text(vjust = -5, face = "bold"),
        axis.title.y = element_text(vjust = 10, face = "bold"),
        plot.margin = margin(0.7, 0.5, 1, 1.2, "cm"))
```

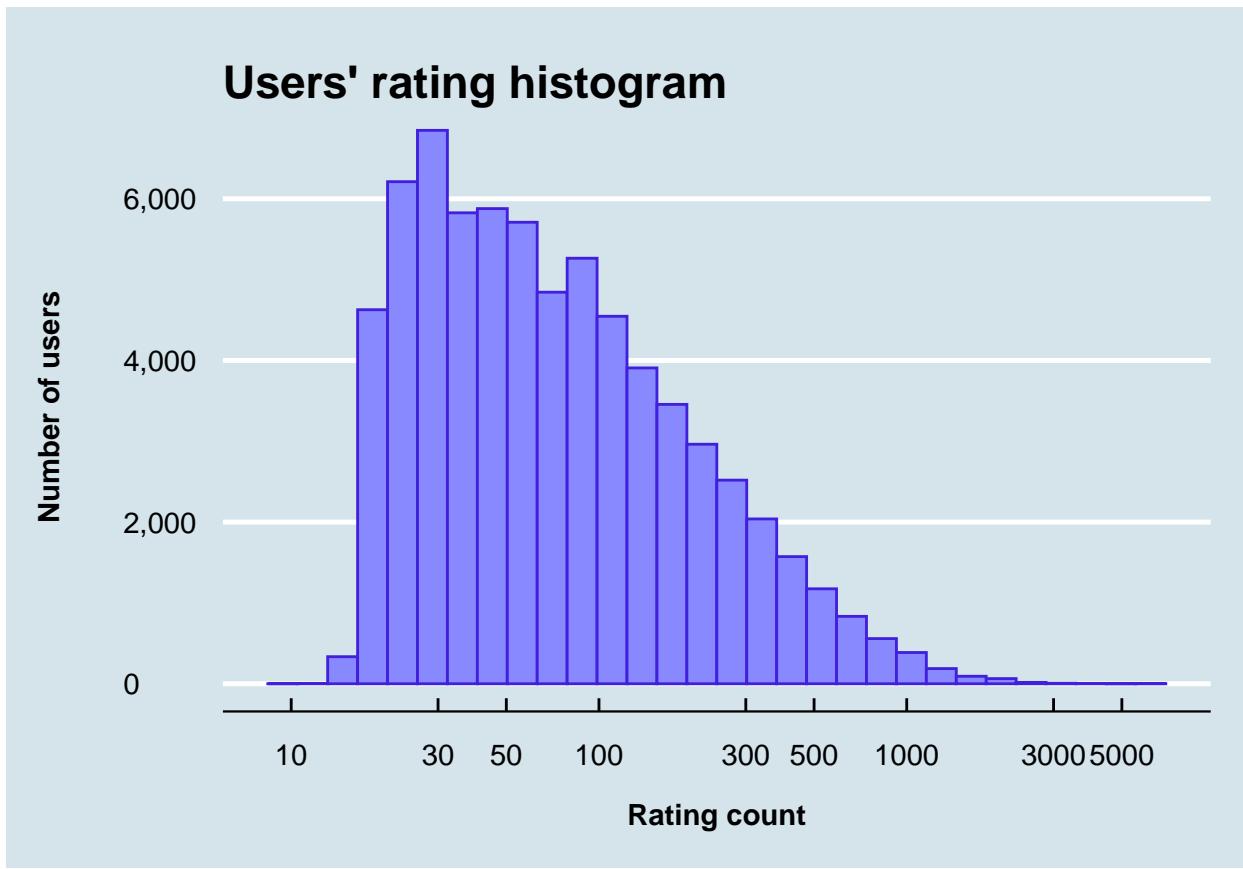
```
## `geom_smooth()` using method = 'gam' and formula = 'y ~ s(x, bs = "cs")'
```



Users' rating histogram

```
edx |>
  group_by(userId) |>
  summarize(count = n()) |>
  ggplot(aes(x = count)) +
  geom_histogram(fill = "#8888ff", color = "#4020dd") +
  ggtitle("Users' rating histogram") +
  xlab("Rating count") +
  ylab("Number of users") +
  scale_y_continuous(labels = comma) +
  scale_x_log10(n.breaks = 10) +
  theme_economist() +
  theme(axis.title.x = element_text(vjust = -5, face = "bold"),
        axis.title.y = element_text(vjust = 10, face = "bold"),
        plot.margin = margin(0.7, 0.5, 1, 1.2, "cm"))

## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Methods / Analysis



All the source code of the R-scripts is available on the project's [GitHub repository](#)[7].

Defining Logging and Time Measuring Helper Functions

First, let's define some helper functions for logging and time-measuring features that we will use in our R scripts. Some of them are listed below:

```
## Logging Helper functions -----
open_logfile <- function(file_name){
  log_file_name <- as.character(Sys.time()) |>
    str_replace_all(':', '_') |>
    str_replace(' ', 'T') |>
    str_c(file_name)

  log_open(file_name = log_file_name)
}

print_start_date <- function(){
```

```

print(date())
Sys.time()
}
put_start_date <- function(){
  put(date())
  Sys.time()
}
print_end_date <- function(start){
  print(date())
  print(Sys.time() - start)
}
put_end_date <- function(start){
  put(date())
  put(Sys.time() - start)
}

msg.set_arg <- function(msg_template, arg, arg.name = "%1") {
  msg_template |>
    str_replace_all(arg.name, as.character(arg))
}
msg.glue <- function(msg_template, arg, arg.name = "%1"){
  msg_template |>
    msg.set_arg(arg, arg.name) |>
    str_glue()
}

print_log <- function(msg){
  print(str_glue(msg))
}
put_log <- function(msg){
  put(str_glue(msg))
}

get_log1 <- function(msg_template, arg1) {
  str_glue(str_replace_all(msg_template, "%1", as.character(arg1)))
}
print_log1 <- function(msg_template, arg1){
  print(get_log1(msg_template, arg1))
}
put_log1 <- function(msg_template, arg1){
  put(get_log1(msg_template, arg1))
}

get_log2 <- function(msg_template, arg1, arg2) {
  msg_template |>
    str_replace_all("%1", as.character(arg1)) |>
    str_replace_all("%2", as.character(arg2)) |>
    str_glue()
}
print_log2 <- function(msg_template, arg1, arg2){
  print(get_log1(msg_template, arg1, arg2))
}
put_log2 <- function(msg_template, arg1, arg2){

```

```

    put(get_log1(msg_template, arg1, arg2))
}

# ...

```



The full source code of these functions is available in the [Logging Helper function section](#) of the [Capstone MovieLens Main R Script](#).

Preparing train and set datasets

We will split the `edx` dataset into a training set, which we will use to build and train our models, and a test set in which we will compute the accuracy of our predictions, the way described in [Section 23.1.1 Movielens data](#) of the *Course Textbook* mentioned above[5]. We will also use the *5-Fold Cross Validation* method as described in [Section 29.6 Cross validation](#) of the *Course Textbook*. To prepare datasets for processing, we will use the following functions, specifically designed for these operations:

```

make_source_datasets()
init_source_datasets()

```



The full source code of the function listed above is available in the [Initialize input datasets](#) section of the `data.helper.functions.R` script on [GitHub](#).

The `make_source_datasets` function

Let's take a closer look at the objects we will receive as a result of executing this function.

```

make_source_datasets <- function(){
  # ...
  list(edx_CV = edx_CV,
       edx_mx = edx.mx,
       edx_sgr = edx.sgr,
       tuning_sets = tuning_sets,
       movie_map = movie_map,
       date_days_map = date_days_map)
}

```

`edx.mx` Matrix Object

We will use the array representation described in [Section 17.5 of the Textbook](#), for the training data: we denote ranking for movie j by user i as $y_{i,j}$. To create this matrix, we use `tidyR::pivot_wider` function:

```

put_log("Function: `make_source_datasets`: Creating Rating Matrix from `edx` dataset...")
edx.mx <- edx |>
  mutate(userId = factor(userId),
         movieId = factor(movieId)) |>

```

```

  select(movieId, userId, rating) |>
  pivot_wider(names_from = movieId, values_from = rating) |>
  column_to_rownames("userId") |>
  as.matrix()

  put_log("Function: `make_source_datasets`:
Matrix created: `edx.mx` of the following dimentions:")

```

```
str(edx.mx)
```

```

##  num [1:69878, 1:10677] 5 NA NA NA NA NA NA NA NA NA ...
##  - attr(*, "dimnames")=List of 2
##    ..$ : chr [1:69878] "1" "2" "3" "4" ...
##    ..$ : chr [1:10677] "122" "185" "292" "316" ...

```

edx.sgr Object

To account for the Movie Genre Effect more accurately, we need a dataset with split rows for movies belonging to multiple genres:

```

put_log("Function: `make_source_datasets`:
To account for the Movie Genre Effect, we need a dataset with split rows
for movies belonging to multiple genres.")
edx.sgr <- splitGenreRows(edx)

```

```
str(edx.sgr)
```

```

## tibble [23,371,423 x 6] (S3: tbl_df/tbl/data.frame)
## $ userId    : int [1:23371423] 1 1 1 1 1 1 1 1 1 ...
## $ movieId   : int [1:23371423] 122 122 185 185 185 292 292 292 292 316 ...
## $ rating    : num [1:23371423] 5 5 5 5 5 5 5 5 5 5 ...
## $ timestamp: int [1:23371423] 838985046 838985046 838983525 838983525 838983525 838983421 838983421
## $ title     : chr [1:23371423] "Boomerang (1992)" "Boomerang (1992)" "Net, The (1995)" "Net, The (1995"
## $ genres    : chr [1:23371423] "Comedy" "Romance" "Action" "Crime" ...

```

```
summary(edx.sgr)
```

	userId	movieId	rating	timestamp	title	genres
## Min.	: 1	Min. : 1	Min. : 0.500	Min. : 7.897e+08	Length:23371423	Length:23371423
## 1st Qu.	:18140	1st Qu.: 616	1st Qu.: 3.000	1st Qu.: 9.472e+08	Class :character	Class :character
## Median	:35784	Median : 1748	Median : 4.000	Median : 1.042e+09	Mode :character	Mode :character
## Mean	:35886	Mean : 4277	Mean : 3.527	Mean : 1.035e+09		
## 3rd Qu.	:53638	3rd Qu.: 3635	3rd Qu.: 4.000	3rd Qu.: 1.131e+09		
## Max.	:71567	Max. : 65133	Max. : 5.000	Max. : 1.231e+09		

Note that we use the `splitGenreRows` function to split rows of the original dataset:

```

splitGenreRows <- function(data){
  put("Splitting dataset rows related to multiple genres...")
  start <- put_start_date()
  gs_splitted <- data |>
    separate_rows(genres, sep = "\\|")
  put("Dataset rows related to multiple genres have been splitted to have single genre per row.")
  put_end_date(start)
  gs_splitted
}

```



The source code of the function mentioned above is also available in the [Initialize input datasets](#) section of the `data.helper.functions.R` script on *GitHub*.

`movie_map` Object

To be able to map movie IDs to titles we create the following lookup table:

```

movie_map <- edx |> select(movieId, title, genres) |>
  distinct(movieId, .keep_all = TRUE)

  put_log("Function: `make_source_datasets`: Dataset created: movie_map")

str(movie_map)

## 'data.frame': 10677 obs. of 3 variables:
## $ movieId: int 122 185 292 316 329 355 356 362 364 370 ...
## $ title  : chr "Boomerang (1992)" "Net, The (1995)" "Outbreak (1995)" "Stargate (1994)" ...
## $ genres : chr "Comedy|Romance" "Action|Crime|Thriller" "Action|Drama|Sci-Fi|Thriller" "Action|Adv

summary(movie_map)

##      movieId          title           genres
##  Min.   :    1  Length:10677  Length:10677
##  1st Qu.: 2754  Class :character  Class :character
##  Median : 5434  Mode  :character  Mode  :character
##  Mean   :13105
##  3rd Qu.: 8710
##  Max.   :65133

```

Note that titles cannot be considered unique, so we can't use them as IDs[5].

`date_days_map` Object

We have a `timestamp` field in the `edx` dataset. To be able to map the date, year, and number of days since the earliest record in the `edx` dataset with the corresponding value in this field, we create the following lookup table:

```

put_log("Function: `make_source_datasets`: Creating Date-Days Map dataset...")
date_days_map <- edx |>
  mutate(date_time = as_datetime(timestamp)) |>
  mutate(date = as_date(date_time)) |>
  mutate(year = year(date_time)) |>
  mutate(days = as.integer(date - min(date))) |>
  select(timestamp, date_time, date, year, days) |>
  distinct(timestamp, .keep_all = TRUE)

put_log("Function: `make_source_datasets`: Dataset created: date_days_map")

```

```
str(date_days_map)
```

```

## 'data.frame':    6519590 obs. of  5 variables:
## $ timestamp: int  838985046 838983525 838983421 838983392 838984474 ...
## $ date_time: POSIXct, format: "1996-08-02 11:24:06" "1996-08-02 10:58:45" ...
## $ date      : Date, format: "1996-08-02" "1996-08-02" "1996-08-02" ...
## $ year      : num  1996 1996 1996 1996 1996 ...
## $ days       : int  571 571 571 571 571 571 571 571 571 571 ...

```

```
summary(date_days_map)
```

	timestamp	date_time	date	year	days
## Min.	:7.897e+08	Min. :1995-01-09 11:46:49.00	Min. :1995-01-09	Min. :1995	Min. :
## 1st Qu.	:9.783e+08	1st Qu.:2001-01-01 05:05:01.75	1st Qu.:2001-01-01	1st Qu.:2001	1st Qu.:21...
## Median	:1.091e+09	Median :2004-08-03 01:08:18.50	Median :2004-08-03	Median :2004	Median :34...
## Mean	:1.066e+09	Mean :2003-10-10 23:15:02.07	Mean :2003-10-10	Mean :2003	Mean :31...
## 3rd Qu.	:1.152e+09	3rd Qu.:2006-07-04 20:41:57.50	3rd Qu.:2006-07-04	3rd Qu.:2006	3rd Qu.:41...
## Max.	:1.231e+09	Max. :2009-01-05 05:02:16.00	Max. :2009-01-05	Max. :2009	Max. :51...

edx_CV Object

Here we have a list of sample objects we need to perform the *5-Fold Cross Validation* as explained in Section 29.6.1 K-fold cross validation of the *Course Textbook*:

```

start <- put_start_date()
edx_CV <- lapply(kfold_index,  function(fold_i){

  put_log1("Method `make_source_datasets`:
Creating K-Fold Cross Validation Datasets, Fold %1", fold_i)

  #> We split the initial datasets into training sets, which we will use to build
  #> and train our models, and validation sets in which we will compute the accuracy
  #> of our predictions, the way described in the `Section 23.1.1 MovieLens data`
  #> (https://rafaelab.dfci.harvard.edu/dsbook-part-2/highdim/regularization.html#movielens-data)
  #> of the Course Textbook.

  split_sets <- edx |>
    sample_train_validation_sets(fold_i*1000)
}

```

```

train_set <- split_sets$train_set
validation_set <- split_sets$validation_set

put_log("Function: `make_source_datasets`:
Sampling 20% from the split-row version of the `edx` dataset...")
split_sets.gs <- edx.sgr |>
  sample_train_validation_sets(fold_i*2000)

train.sgr <- split_sets.gs$train_set
validation.sgr <- split_sets.gs$validation_set

# put_log("Function: `make_source_datasets`: Dataset created: validation.sgr")
# put(summary(validation.sgr))

##> We will use the array representation described in `Section 17.5 of the Textbook` 
##> (https://rafalab.dfcii.harvard.edu/dsbook-part-2/linear-models/treatment-effect-models.html#sec-a)
##> for the training data.
##> To create this matrix, we use `tidy::pivot_wider` function:

put_log("Function: `make_source_datasets`: Creating Rating Matrix from Train Set...")
train_mx <- train_set |>
  mutate(userId = factor(userId),
         movieId = factor(movieId)) |>
  select(movieId, userId, rating) |>
  pivot_wider(names_from = movieId, values_from = rating) |>
  column_to_rownames("userId") |>
  as.matrix()

put_log("Function: `make_source_datasets`:
Matrix created: `train_mx` of the following dimentions:")
put(dim(train_mx))

list(train_set = train_set,
     train_mx = train_mx,
     train.sgr = train.sgr,
     validation_set = validation_set)
}

put_end_date(start)
put_log("Function: `make_source_datasets`:
Set of K-Fold Cross Validation datasets created: edx_CV")

```

```
str(edx_CV)
```

```

## List of 5
## $ :List of 4
##   ..$ train_set      :'data.frame': 7172311 obs. of 6 variables:
##   ...$ userId       : int [1:7172311] 1 1 1 1 1 1 1 1 1 ...
##   ...$ movieId     : int [1:7172311] 122 185 292 329 356 362 364 370 420 466 ...
##   ...$ rating      : num [1:7172311] 5 5 5 5 5 5 5 5 5 ...
##   ...$ timestamp   : int [1:7172311] 838985046 838983525 838983421 838983392 838983653 ...
##   ...$ title       : chr [1:7172311] "Boomerang (1992)" "Net, The (1995)" "Outbreak (1995)" "Star Trek (1995)" ...
##   ...$ genres      : chr [1:7172311] "Comedy|Romance" "Action|Crime|Thriller" "Action|Drama|Sci-Fi|Thriller" ...
##   ..$ train_mx     : num [1:69878, 1:10677] 5 NA NA NA NA NA NA NA NA ...

```

```

## ... - attr(*, "dimnames")=List of 2
## ... .$. : chr [1:69878] "1" "2" "3" "4" ...
## ... .$. : chr [1:10677] "122" "185" "292" "329" ...
## ... $.train.sgr : tibble [18,669,190 x 6] (S3: tbl_df/tbl/data.frame)
## ... $.userId : int [1:18669190] 1 1 1 1 1 1 1 1 1 1 ...
## ... $.movieId : int [1:18669190] 122 122 185 185 292 292 292 292 316 316 ...
## ... $.rating : num [1:18669190] 5 5 5 5 5 5 5 5 5 5 ...
## ... $.timestamp: int [1:18669190] 838985046 838985046 838983525 838983525 838983421 838983421 838983421 838983421 838983421 838983421 ...
## ... $.title : chr [1:18669190] "Boomerang (1992)" "Boomerang (1992)" "Net, The (1995)" "Net, The (1995)" ...
## ... $.genres : chr [1:18669190] "Comedy" "Romance" "Action" "Crime" ...
## ... $.validation_set:'data.frame': 1827744 obs. of 6 variables:
## ... $.userId : int [1:1827744] 1 1 1 1 2 2 2 2 3 3 ...
## ... $.movieId : int [1:1827744] 316 355 377 588 260 376 648 1049 110 1252 ...
## ... $.rating : num [1:1827744] 5 5 5 5 5 3 2 3 4.5 4 ...
## ... $.timestamp: int [1:1827744] 838983392 838984474 838983834 838983339 868244562 868245920 868245920 868245920 868245920 868245920 ...
## ... $.title : chr [1:1827744] "Stargate (1994)" "Flintstones, The (1994)" "Speed (1994)" "Aladdin (1992)" ...
## ... $.genres : chr [1:1827744] "Action|Adventure|Sci-Fi" "Children|Comedy|Fantasy" "Action|Romantic" ...
## $ :List of 4
## ... $.train_set : 'data.frame': 7172306 obs. of 6 variables:
## ... $.userId : int [1:7172306] 1 1 1 1 1 1 1 1 1 ...
## ... $.movieId : int [1:7172306] 122 185 292 316 329 355 356 364 370 377 ...
## ... $.rating : num [1:7172306] 5 5 5 5 5 5 5 5 5 5 ...
## ... $.timestamp: int [1:7172306] 838985046 838983525 838983421 838983392 838983392 838984474 838984474 838984474 838984474 838984474 ...
## ... $.title : chr [1:7172306] "Boomerang (1992)" "Net, The (1995)" "Outbreak (1995)" "Stargate (1994)" ...
## ... $.genres : chr [1:7172306] "Comedy|Romance" "Action|Crime|Thriller" "Action|Drama|Sci-Fi|Thriller" ...
## ... $.train_mx : num [1:69878, 1:10677] 5 NA NA NA NA NA NA NA NA ...
## ... - attr(*, "dimnames")=List of 2
## ... .$. : chr [1:69878] "1" "2" "3" "4" ...
## ... .$. : chr [1:10677] "122" "185" "292" "316" ...
## ... $.train.sgr : tibble [18,669,201 x 6] (S3: tbl_df/tbl/data.frame)
## ... $.userId : int [1:18669201] 1 1 1 1 1 1 1 1 1 ...
## ... $.movieId : int [1:18669201] 122 122 185 185 185 292 292 316 316 329 ...
## ... $.rating : num [1:18669201] 5 5 5 5 5 5 5 5 5 5 ...
## ... $.timestamp: int [1:18669201] 838985046 838985046 838983525 838983525 838983525 838983525 838983525 838983525 838983525 838983525 ...
## ... $.title : chr [1:18669201] "Boomerang (1992)" "Boomerang (1992)" "Net, The (1995)" "Net, The (1995)" ...
## ... $.genres : chr [1:18669201] "Comedy" "Romance" "Action" "Crime" ...
## ... $.validation_set:'data.frame': 1827749 obs. of 6 variables:
## ... $.userId : int [1:1827749] 1 1 1 1 2 2 2 2 3 3 ...
## ... $.movieId : int [1:1827749] 362 520 539 594 539 590 733 1210 1252 1408 ...
## ... $.rating : num [1:1827749] 5 5 5 5 3 5 3 4 4 3.5 ...
## ... $.timestamp: int [1:1827749] 838984885 838984679 838984068 838984679 868246262 868245608 868245608 868245608 868245608 868245608 ...
## ... $.title : chr [1:1827749] "Jungle Book, The (1994)" "Robin Hood: Men in Tights (1993)" "Sleeping Beauty (1959)" ...
## ... $.genres : chr [1:1827749] "Adventure|Children|Romance" "Comedy" "Comedy|Drama|Romance" "Animation|Family|Romance" ...
## $ :List of 4
## ... $.train_set : 'data.frame': 7172307 obs. of 6 variables:
## ... $.userId : int [1:7172307] 1 1 1 1 1 1 1 1 1 ...
## ... $.movieId : int [1:7172307] 122 185 292 316 329 355 362 370 377 420 ...
## ... $.rating : num [1:7172307] 5 5 5 5 5 5 5 5 5 5 ...
## ... $.timestamp: int [1:7172307] 838985046 838983525 838983421 838983392 838983392 838984474 838984474 838984474 838984474 838984474 ...
## ... $.title : chr [1:7172307] "Boomerang (1992)" "Net, The (1995)" "Outbreak (1995)" "Stargate (1994)" ...
## ... $.genres : chr [1:7172307] "Comedy|Romance" "Action|Crime|Thriller" "Action|Drama|Sci-Fi|Thriller" ...
## ... $.train_mx : num [1:69878, 1:10677] 5 NA NA NA NA NA NA NA NA ...
## ... - attr(*, "dimnames")=List of 2
## ... .$. : chr [1:69878] "1" "2" "3" "4" ...

```

```

## ... .$. : chr [1:10677] "122" "185" "292" "316" ...
## ..$ train.sgr : tibble [18,669,195 x 6] (S3: tbl_df/tbl/data.frame)
## ... $.userId : int [1:18669195] 1 1 1 1 1 1 1 1 1 1 ...
## ... $.movieId : int [1:18669195] 122 122 185 185 185 292 292 292 316 329 ...
## ... $.rating : num [1:18669195] 5 5 5 5 5 5 5 5 5 5 ...
## ... $.timestamp: int [1:18669195] 838985046 838985046 838983525 838983525 838983525 838983421 838983421 838983421 838983421 838983421 ...
## ... $.title : chr [1:18669195] "Boomerang (1992)" "Boomerang (1992)" "Net, The (1995)" "Net, The (1995)" ...
## ... $.genres : chr [1:18669195] "Comedy" "Romance" "Action" "Crime" ...
## ... $.validation_set:'data.frame': 1827748 obs. of 6 variables:
## ... $.userId : int [1:1827748] 1 1 1 1 2 2 2 2 3 3 ...
## ... $.movieId : int [1:1827748] 356 364 539 616 590 719 780 786 151 213 ...
## ... $.rating : num [1:1827748] 5 5 5 5 5 3 3 3 4.5 5 ...
## ... $.timestamp: int [1:1827748] 838983653 838983707 838984068 838984941 868245608 868246191 868246191 868246191 868246191 ...
## ... $.title : chr [1:1827748] "Forrest Gump (1994)" "Lion King, The (1994)" "Sleepless in Seattle (1995)" ...
## ... $.genres : chr [1:1827748] "Comedy|Drama|Romance|War" "Adventure|Animation|Children|Drama|Mystery" ...
## $ :List of 4
## ... $.train_set : 'data.frame': 7172311 obs. of 6 variables:
## ... $.userId : int [1:7172311] 1 1 1 1 1 1 1 1 1 ...
## ... $.movieId : int [1:7172311] 122 185 292 316 329 355 356 362 364 370 ...
## ... $.rating : num [1:7172311] 5 5 5 5 5 5 5 5 5 5 ...
## ... $.timestamp: int [1:7172311] 838985046 838983525 838983421 838983392 838983392 838984474 838984474 838984474 838984474 838984474 ...
## ... $.title : chr [1:7172311] "Boomerang (1992)" "Net, The (1995)" "Outbreak (1995)" "Stargate (1994)" ...
## ... $.genres : chr [1:7172311] "Comedy|Romance" "Action|Crime|Thriller" "Action|Drama|Sci-Fi|Thriller" ...
## ... $.train_mx : num [1:69878, 1:10677] 5 NA NA NA NA NA NA NA NA ...
## ... -- attr(*, "dimnames")=List of 2
## ... ... $. : chr [1:69878] "1" "2" "3" "4" ...
## ... ... $. : chr [1:10677] "122" "185" "292" "316" ...
## ... $.train.sgr : tibble [18,669,192 x 6] (S3: tbl_df/tbl/data.frame)
## ... $.userId : int [1:18669192] 1 1 1 1 1 1 1 1 1 ...
## ... $.movieId : int [1:18669192] 122 122 185 185 292 292 316 316 329 ...
## ... $.rating : num [1:18669192] 5 5 5 5 5 5 5 5 5 ...
## ... $.timestamp: int [1:18669192] 838985046 838985046 838983525 838983525 838983421 838983421 838983421 838983421 838983421 ...
## ... $.title : chr [1:18669192] "Boomerang (1992)" "Boomerang (1992)" "Net, The (1995)" "Net, The (1995)" ...
## ... $.genres : chr [1:18669192] "Comedy" "Romance" "Action" "Thriller" ...
## ... $.validation_set:'data.frame': 1827744 obs. of 6 variables:
## ... $.userId : int [1:1827744] 1 1 1 1 2 2 2 3 3 ...
## ... $.movieId : int [1:1827744] 377 520 588 616 110 648 1049 1356 1148 1276 ...
## ... $.rating : num [1:1827744] 5 5 5 5 5 2 3 3 4 3.5 ...
## ... $.timestamp: int [1:1827744] 838983834 838984679 838983339 838984941 868245777 868244699 868244699 868244699 868244699 ...
## ... $.title : chr [1:1827744] "Speed (1994)" "Robin Hood: Men in Tights (1993)" "Aladdin (1992)" ...
## ... $.genres : chr [1:1827744] "Action|Romance|Thriller" "Comedy" "Adventure|Animation|Children|Drama|Mystery" ...
## $ :List of 4
## ... $.train_set : 'data.frame': 7172301 obs. of 6 variables:
## ... $.userId : int [1:7172301] 1 1 1 1 1 1 1 1 1 ...
## ... $.movieId : int [1:7172301] 122 185 292 316 355 356 364 370 420 466 ...
## ... $.rating : num [1:7172301] 5 5 5 5 5 5 5 5 5 5 ...
## ... $.timestamp: int [1:7172301] 838985046 838983525 838983421 838983392 838984474 838983653 838983653 838983653 838983653 838983653 ...
## ... $.title : chr [1:7172301] "Boomerang (1992)" "Net, The (1995)" "Outbreak (1995)" "Stargate (1994)" ...
## ... $.genres : chr [1:7172301] "Comedy|Romance" "Action|Crime|Thriller" "Action|Drama|Sci-Fi|Thriller" ...
## ... $.train_mx : num [1:69878, 1:10677] 5 NA NA NA NA NA NA NA NA ...
## ... -- attr(*, "dimnames")=List of 2
## ... ... $. : chr [1:69878] "1" "2" "3" "4" ...
## ... ... $. : chr [1:10677] "122" "185" "292" "316" ...
## ... $.train.sgr : tibble [18,669,194 x 6] (S3: tbl_df/tbl/data.frame)

```

```

## ...$ userId    : int [1:18669194] 1 1 1 1 1 1 1 1 1 ...
## ...$ movieId   : int [1:18669194] 122 122 185 185 292 292 316 329 329 355 ...
## ...$ rating    : num [1:18669194] 5 5 5 5 5 5 5 5 5 ...
## ...$ timestamp: int [1:18669194] 838985046 838985046 838983525 838983525 838983421 838983421 838983421 ...
## ...$ title     : chr [1:18669194] "Boomerang (1992)" "Boomerang (1992)" "Net, The (1995)" "Net, The (1995)" ...
## ...$ genres    : chr [1:18669194] "Comedy" "Romance" "Crime" "Thriller" ...
## ...$ validation_set:'data.frame': 1827754 obs. of 6 variables:
## ...$ userId    : int [1:1827754] 1 1 1 1 2 2 2 2 3 3 ...
## ...$ movieId   : int [1:1827754] 329 362 377 594 110 376 539 736 1252 1408 ...
## ...$ rating    : num [1:1827754] 5 5 5 5 5 3 3 3 4 3.5 ...
## ...$ timestamp: int [1:1827754] 838983392 838984885 838983834 838984679 868245777 868245920 868245920 ...
## ...$ title     : chr [1:1827754] "Star Trek: Generations (1994)" "Jungle Book, The (1994)" "Speed Racer (2008)" ...
## ...$ genres    : chr [1:1827754] "Action|Adventure|Drama|Sci-Fi" "Adventure|Children|Romance" "Action|Thriller" ...

```



This code snippet is a part of the `make_source_datasets` function code described above.

Note that we used the `sample_train_validation_sets` function call to split the original dataset (`edx` in this case):

```

split_sets <- edx |>
  sample_train_validation_sets(fold_i*1000)

```

which returns a pair of train/validation sets:

```

sample_train_validation_sets <- function(data, seed){
  put_log("Function: `sample_train_validation_sets`: Sampling 20% of the `data` data...")
  set.seed(seed)
  validation_ind <-
    sapply(splitByUser(data),
      function(i) sample(i, ceiling(length(i)*.2))) |>
  unlist() |>
  sort()

  put_log("Function: `sample_train_validation_sets`:
Extracting 80% of the original `data` not used for the Validation Set,
excluding data for users who provided no more than a specified number of ratings: {min_nratings}.") 

  train_set <- data[-validation_ind,]

  put_log("Function: `sample_train_validation_sets`: Dataset created: train_set")
  put(summary(train_set))

  put_log("Function: `sample_train_validation_sets`:
To make sure we don't include movies in the Training Set that should not be there,
we exclude entries using the semi_join function from the Validation Set.")
  tmp.data <- data[validation_ind,]

  validation_set <- tmp.data |>
    semi_join(train_set, by = "movieId") |>
    semi_join(train_set, by = "userId") |>

```

```

    as.data.frame()

  # Add rows excluded from `validation_set` into `train_set`
  tmp.excluded <- anti_join(tmp.data, validation_set)
  train_set <- rbind(train_set, tmp.excluded)

  put_log("Function: `sample_train_validation_sets`: Dataset created: validation_set")
  put(summary(validation_set))

  # CV train & test sets Consistency Test
  validation.left_join.Nas <- train_set |>
    mutate(tst.col = rating) |>
    select(userId, movieId, tst.col) |>
    data.consistency.test(validation_set)

  put_log("Function: `sample_train_validation_sets`:
Below are the data consistency verification results")
  put(validation.left_join.Nas)

  # Return result datasets -----
  list(train_set = train_set,
       validation_set = validation_set)
}

```



The `sample_train_validation_sets` function is defined in the same script as the `make_source_datasets` one, from where it is called.

Common Helper Functions

For our further analysis, we are going to use the following *common helper functions*:

`clamp` function

As explained in [Section 24.4 User effects](#) of the *Course Textbook* we know ratings can't be below 0.5 or above 5. For this reason, we will use the `clamp` function described in that section:

```
clamp <- function(x, min = 0.5, max = 5) pmax(pmin(x, max), min)
```

Functions to calculate (*Root*) Mean Squared Error

We will need the following functions to calculate (*R*)MSEs:

```

mse <- function(r) mean(r^2)

mse_cv <- function(r_list) {
  mses <- sapply(r_list, mse(r))
  mean(mses)
}

```

```

rmse <- function(r) sqrt(mse(r))
# rmse_cv <- function(r_list) sqrt(mse_cv(r_list))

rmse2 <- function(true_ratings, predicted_ratings) {
  rmse(true_ratings - predicted_ratings)
}

```



All the *common helper functions*, including those described above, are defined in the [common-helper.functions.R](#) script on *GitHub*.

Overall Mean Rating (Naive) Model

Let's begin our analysis by evaluating the simplest model described in Section 23.3 *The First Model* of the [Course Textbook](#), and then gradually refine it through further research. It is about a model that assumes the same rating for all movies and users with all the differences explained by random variation would look as follows:

$$Y_{i,j} = \mu + \varepsilon_{i,j}$$

with $\varepsilon_{i,j}$ independent errors sampled from the same distribution centered at 0 and μ the *true* rating for all movies.

We know that the estimate that minimizes the RMSE is the least squares estimate of μ and, in this case, is the average of all ratings:

```

mu <- mean(edx$rating)
print(mu)

```

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

If we predict all unknown ratings with $\hat{\mu}$, we obtain the following RMSE:

```

mu.MSEs <- naive_model_MSEs(mu)
data.frame(fold_No = 1:5, MSE = mu.MSEs) |>
  data.plot(title = "MSE results of the 5-fold CV method applied to the Overall Mean Rating Model",
            xname = "fold_No",
            yname = "MSE")

```

MSE results of the 5-fold CV method applied to the Overall Mean Rating



```
mu.RMSE <- sqrt(mean(mu.MSEs))
mu.RMSE
```

```
## [1] 1.060346
```



For the *Mean Squared Error* data visualization we used `data.plot` function] defined in the [Data Visualization](#) section of the `data.helper.function.R` script.

```
data.plot <- function(data,
                      title,
                      xname,
                      yname,
                      xlabel = NULL,
                      ylabel = NULL,
                      line_col = "blue",
                      # scale = 1,
                      normalize = FALSE) {
  y <- data[, yname]

  if (normalize) {
    y <- y - min(y)
  }

  if (is.null(xlabel)) {
```

```

    xlabel = xname
}
if (is.null(ylabel)) {
  ylabel = yname
}

aes_mapping <- aes(x = data[, xname], y = y)

data |>
  ggplot(mapping = aes_mapping) +
  ggtitle(title) +
  xlab(xlabel) +
  ylab(ylabel) +
  geom_point() +
  geom_line(color=line_col)
}

```

Here we also used `naive_model_MSEs` function defined in the `common-helper.functions.R` script (already mentioned above) to compute *Mean Squared Errors* using *5-Fold Cross Validation* method:

```

naive_model_MSEs <- function(val) {
  sapply(edx_CV, function(cv_item){
    mse(cv_item$validation_set$rating - val)
  })
}

```

One more function, defined in the `same script`, that we will need for further analysis of the current model, is the `naive_model_RMSE` one:

```

naive_model_RMSE <- function(val){
  sqrt(mean(naive_model_MSEs(val)))
}

```

Ensure that `mu.RMSE` value is the best for the current model

If we plug in any other number, we will get a higher RMSE. Let's prove that by the following small investigation:

```

deviation <- seq(0, 6, 0.1) - 3

deviation.RMSE <- sapply(deviation, function(delta){
  naive_model_RMSE(mu + delta)
})

```

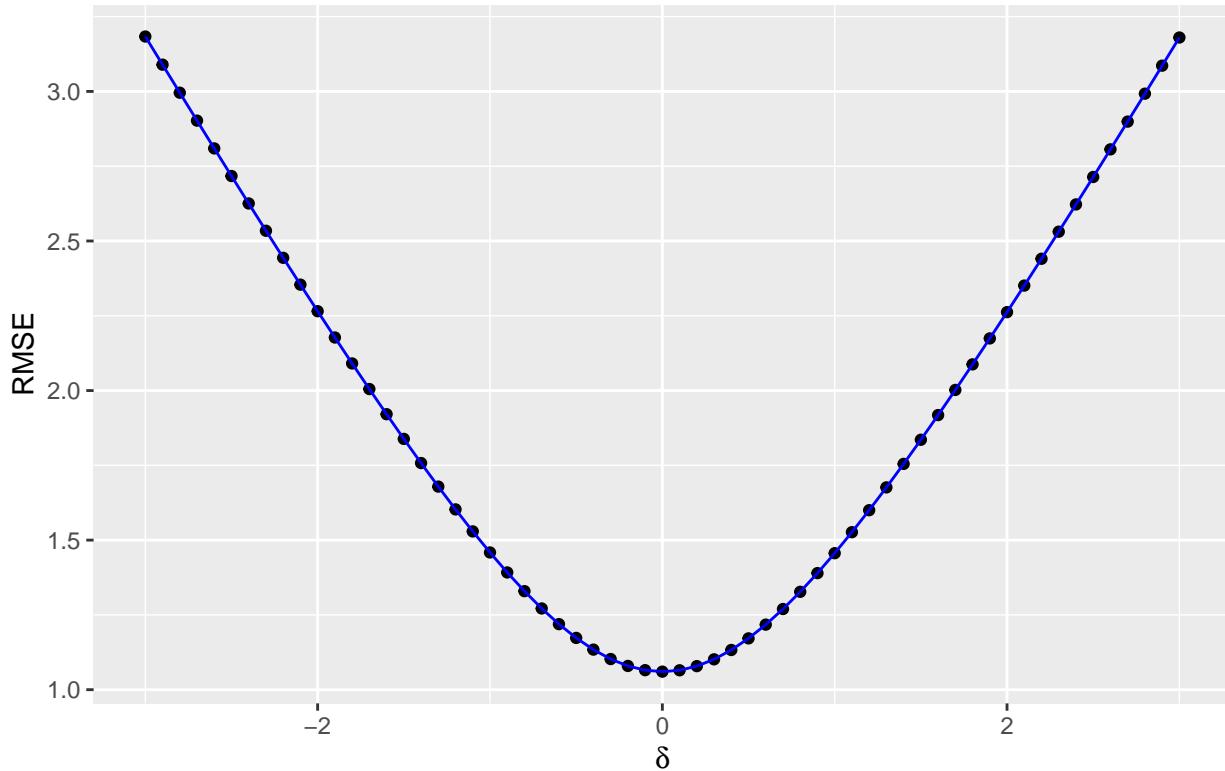
Let's make a quick investigation of the `deviation.RMSE` result we have just got:

```

data.frame(delta = deviation,
           delta.RMSE = deviation.RMSE) |>
data.plot(title = TeX(r' [RMSE as a function of deviation ($\delta$) from the Overall Mean Rating ($\hat{\mu}$) ]'),
          xname = "delta",
          yname = "delta.RMSE",
          xlabel = TeX(r' ['$\delta$']'),
          ylabel = "RMSE")

```

RMSE as a function of deviation (δ) from the Overall Mean Rating ($\hat{\mu}$)



```
which_min_deviation <- deviation[which.min(deviation.RMSE)]
min_rmse = min(deviation.RMSE)

print_log1("Minimum RMSE is achieved when the deviation from the mean is: %1",
          which_min_deviation)
```

```
## Minimum RMSE is achieved when the deviation from the mean is: 0
```

```
print_log1("Is the previously computed RMSE the best for the current model: %1",
          mu.RMSE == min_rmse)
```

```
## Is the previously computed RMSE the best for the current model: TRUE
```

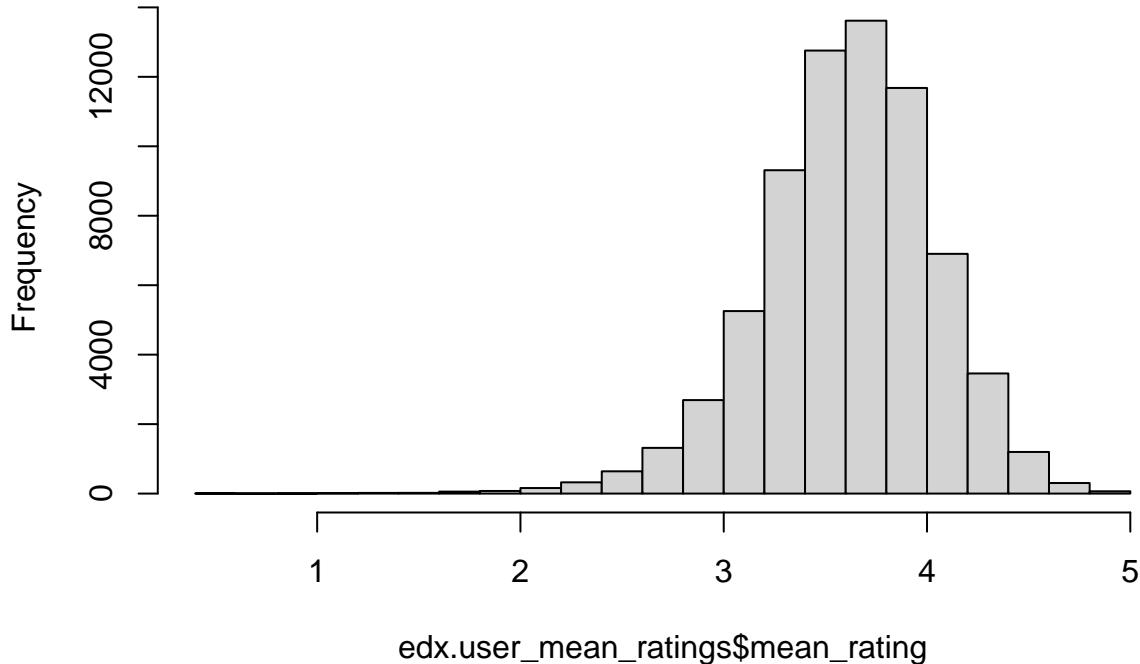
To win the grand prize of \$1,000,000, a participating team had to get an RMSE of at least 0.8563[2]. So we can definitely do better![8]

Taking into account User Effect

To improve our model let's now take into consideration user effects as explained in [Section 23.4 User effects of the Course Textbook](#). If we visualize the average rating for each user the way the [the author](#) shows, we can see that there is substantial variability in the average ratings across users:

```
hist(edx.user_mean_ratings$mean_rating, nclass = 30)
```

Histogram of edx.user_mean_ratings\$mean_rating



```
print("A histogram of the User Mean Rating distribution has been plotted.")
```

```
## [1] "A histogram of the User Mean Rating distribution has been plotted."
```

Following the author's further explanation, to account for this variability, we will use a linear model with a *treatment effect* α_i for each user. The sum $\mu + \alpha_i$ can be interpreted as the typical rating user i gives to movies. So we write the model as follows:

$$Y_{i,j} = \mu + \alpha_i + \varepsilon_{i,j}$$

Statistics textbooks refer to the α s as treatment effects. In the Netflix challenge papers, they refer to them as *bias*[9, 10].

As it is stated here[9], it can be shown that the least squares estimate $\hat{\alpha}_i$ is just the average of $y_{i,j} - \hat{\mu}$ for each user i . So we compute them this way:

```
a <- rowMeans(y - mu, na.rm = TRUE)
```

Finally, we are ready to compute the RMSE (additionally using the helper function `clamp` we defined above to keep predictions in the proper range):

```
# Compute the RMSE taking into account user effects:  
user_effects_rmse <- test_set |>  
  left_join(data.frame(userId = as.integer(names(a)), a = a), by = "userId") |>  
  mutate(resid = rating - clamp(mu + a)) |>
```

```

filter(!is.na(resid)) |>
  pull(resid) |> rmse()

print(user_effects_rmse)

```

Taking into account Movie effects

In Section 23.5 *Movie effects* of the *Course Textbook* the author draws our attention to the fact that some movies are generally rated higher than others. He also explains that a linear model with a *treatment effect* β_j for each movie can be used in this case, which can be interpreted as movie effect or the difference between the average ranking for movie j and the overall average μ :

$$Y_{i,j} = \mu + \alpha_i + \beta_j + \varepsilon_{i,j}$$

The author then shows how to use an approximation by first computing the least square estimate $\hat{\mu}$ and $\hat{\alpha}_i$, and then estimating $\hat{\beta}_j$ as the average of the residuals $y_{i,j} - \hat{\mu} - \hat{\alpha}_i$:

```
b <- colMeans(y - mu - a, na.rm = TRUE)
```

We can now construct predictors and see how much the RMSE improves[11]:

```

user_and_movie_effects_rmse <- test_set |>
  left_join(data.frame(userId = as.integer(names(a)), a = a), by = "userId") |>
  left_join(data.frame(movieId = as.integer(names(b)), b = b), by = "movieId") |>
  mutate(resid = rating - clamp(mu + a + b)) |>
  filter(!is.na(resid)) |>
  pull(resid) |> rmse()

print(user_and_movie_effects_rmse)

```

Utilizing Penalized least squares

Section 23.6 *Penalized least squares* of the *Course Textbook* explains why and how we should use *Penalized least squares* to improve our predictions. The author also explains that the general idea of penalized regression is to control the total variability of the movie effects: $\sum_{j=1}^n \beta_j^2$. Specifically, instead of minimizing the least squares equation, we minimize an equation that adds a penalty:

$$\sum_{i,j} (y_{u,i} - \mu - \alpha_i - \beta_j)^2 + \lambda \sum_j \beta_j^2$$

The first term is just the sum of squares and the second is a penalty that gets larger when many β_i s are large. Using calculus, we can actually show that the values of β_i that minimize this equation are:

$$\hat{\beta}_j(\lambda) = \frac{1}{\lambda + n_j} \sum_{i=1}^{n_j} (Y_{i,j} - \mu - \alpha_i)$$

where n_j is the number of ratings made for movie j .

This approach will have our desired effect: when our sample size n_j is very large, we obtain a stable estimate and the penalty λ is effectively ignored since $n_j + \lambda \approx n_j$. Yet when the n_j is small, then the estimate $\hat{\beta}_j(\lambda)$ is shrunken towards 0. The larger the λ , the more we shrink[12].

Support function

We will use the following function to calculate $RMSE$ in this section:

```
reg_rmse <- function(b){  
  test_set |>  
    left_join(data.frame(userId = as.integer(names(a)), a = a), by = "userId") |>  
    left_join(data.frame(movieId = as.integer(names(b)), b = b), by = "movieId") |>  
    mutate(resid = rating - clamp(mu + a + b)) |>  
    filter(!is.na(resid)) |>  
    pull(resid) |> rmse()  
}
```

Let's now figure out the λ that minimizes the $RMSE$:

```
# Here we will simply compute the RMSE for different values of `lambda`  
n <- colSums(!is.na(y))  
  
sums <- colSums(y - mu - a, na.rm = TRUE)  
lambdas <- seq(0, 10, 0.1)  
  
rmses <- sapply(lambdas, function(lambda){  
  b <- sums / (n + lambda)  
  reg_rmse(b)  
})  
  
# Here is a plot of the RMSE versus `lambda`:  
plot(lambdas, rmses, type = "l")
```

Now we can determine the minimal $RMSE$:

```
# print(min(rmses))
```

which is achieved for the following λ :

```
lambda <- lambdas[which.min(rmses)]  
print(lambda)
```

Using this λ we can compute the regularized estimates:

```
b_reg <- sums / (n + lambda)  
  
str(b_reg)
```

Finally, let's verify that the penalized estimates $\hat{b}_i(\lambda)$ we have just computed actually result in the minimal $RMSE$ figured out above:

```
reg_rmse(b_reg)
```

Accounting for Date effects

Yearly rating count[6]

```

print(edx |>
  mutate(year = year(as_datetime(timestamp, origin = "1970-01-01"))) |>
  group_by(year) |>
  summarize(count = n())
)

## # A tibble: 15 x 2
##       year   count
##     <dbl>   <int>
## 1  1995      2
## 2  1996  942772
## 3  1997  414101
## 4  1998  181634
## 5  1999  709893
## 6  2000 1144349
## 7  2001  683355
## 8  2002  524959
## 9  2003  619938
## 10 2004  691429
## 11 2005 1059277
## 12 2006  689315
## 13 2007  629168
## 14 2008  696740
## 15 2009  13123

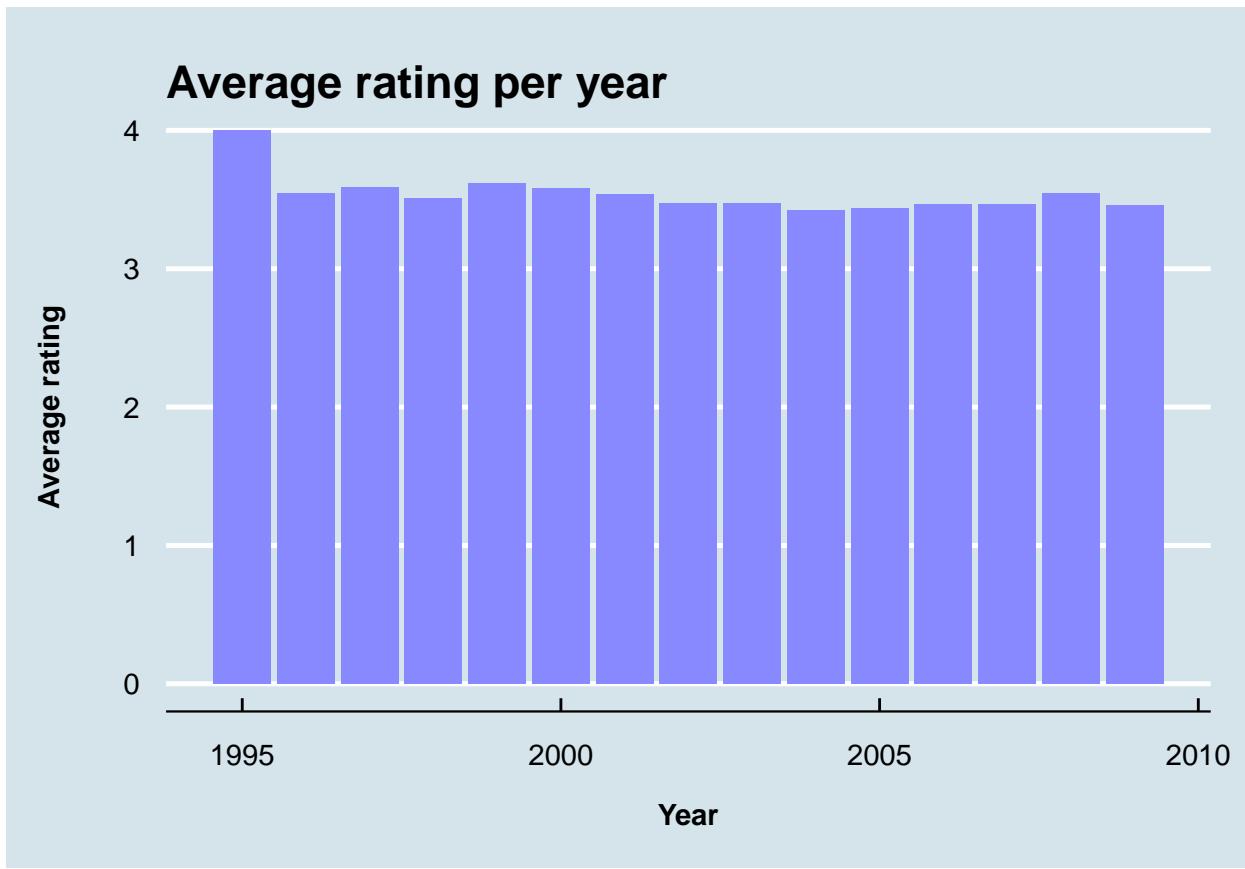
```

Average rating per year plot[6]

```

edx |>
  mutate(year = year(as_datetime(timestamp, origin = "1970-01-01"))) |>
  group_by(year) |>
  summarize(rating_avg = mean(rating)) |>
  ggplot(aes(x = year, y = rating_avg)) +
  geom_bar(stat = "identity", fill = "#8888ff") +
  ggtitle("Average rating per year") +
  xlab("Year") +
  ylab("Average rating") +
  scale_y_continuous(labels = comma) +
  theme_economist() +
  theme(axis.title.x = element_text(vjust = -5, face = "bold"),
        axis.title.y = element_text(vjust = 10, face = "bold"),
        plot.margin = margin(0.7, 0.5, 1, 1.2, "cm"))

```



We use the following models to account for the `date` effect:

$$Y_{i,j} = \mu + \alpha_i + \beta_j + f(d_{i,j}) + \varepsilon_{i,j}$$

Accounting for Genre effect

As mentioned in [Section 23.7: Exercises](#) of the *Chapter “23 Regularization” of the Course Textbook* the MovieLens dataset also has a genres column. This column includes every genre that applies to the movie (some movies fall under several genres)[[13](#)].

Genre Data Analysis

Movie Genres Data

The following code computes movie rating summaries by popular genres like Drama, Comedy, Thriller, and Romance:

```
#library(stringr)
genres = c("Drama", "Comedy", "Thriller", "Romance")
sapply(genres, function(g) {
  sum(str_detect(edx$genres, g))
})
```

Further, we can find out the movies that have the greatest number of ratings using the following code:

```
ordered_movie_ratings <- edx |> group_by(movieId, title) |>
  summarize(number_of_ratings = n()) |>
  arrange(desc(number_of_ratings))
print(head(ordered_movie_ratings))
```

and figure out the most given ratings in order from most to least:

```
ratings <- edx |> group_by(rating) |>
  summarise(count = n()) |>
  arrange(desc(count))
print(ratings)
```

The following code allows us to summarize that in general, half-star ratings are less common than whole-star ratings (e.g., there are fewer ratings of 3.5 than there are ratings of 3 or 4, etc.):

```
print(edx |> group_by(rating) |> summarize(count = n()))
```

We can visually see that from the following plot:

```
edx |>
  group_by(rating) |>
  summarize(count = n()) |>
  ggplot(aes(x = rating, y = count)) +
  geom_line()
```

Movie Genres Effect

The plot below shows strong evidence of a genre effect (for illustrative purposes, the plot shows only categories with more than 20, 000 ratings).

```
# Preparing data for plotting:
genre_ratins_grp <- train_set |>
  mutate(genre_categories = as.factor(genres)) |>
  group_by(genre_categories) |>
  summarize(n = n(), rating_avg = mean(rating), se = sd(rating)/sqrt(n())) |>
  filter(n > 20000) |>
  mutate(genres = reorder(genre_categories, rating_avg)) |>
  select(genres, rating_avg, se, n)

dim(genre_ratins_grp)
genre_ratins_grp_sorted <- genre_ratins_grp |> sort_by.data.frame(~ rating_avg)
print(genre_ratins_grp_sorted)

# Creating plot:
genre_ratins_grp |>
  ggplot(aes(x = genres, y = rating_avg, ymin = rating_avg - 2*se, ymax = rating_avg + 2*se)) +
  geom_point() +
  geom_errorbar() +
  ggtitle("Average rating per Genre") +
  ylab("Average rating") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1))
```

Below are worst and best ratings categories:

```
sprintf("The worst ratings are for the genre category: %s",
       genre_ratins_grp$genres[which.min(genre_ratins_grp$genres)])
printf("The best ratings are for the genre category: %s",
       genre_ratins_grp$genres[which.max(genre_ratins_grp$genres)])
```

Another way of visualizing a genre effect is shown in the section [Average rating for each genre](#) of the article “Movie Recommendation System using R - BEST” written by [Amir Moterfaker](#)[6]:

```
# For better visibility, we reduce the data for plotting
# while keeping the worst and best rating rows:
plot_ind <- odd(1:nrow(genre_ratins_grp))
plot_dat <- genre_ratins_grp_sorted[plot_ind,]

plot_dat |>
  ggplot(aes(x = rating_avg, y = genres)) +
  ggtitle("Genre Average Rating") +
  geom_bar(stat = "identity", width = 0.6, fill = "#8888ff") +
  xlab("Average ratings") +
  ylab("Genres") +
  scale_x_continuous(labels = comma, limits = c(0.0, 5.0)) +
  theme_economist() +
  theme(plot.title = element_text(vjust = 3.5),
        axis.title.x = element_text(vjust = -5, face = "bold"),
        axis.title.y = element_text(vjust = 10, face = "bold"),
        axis.text.x = element_text(vjust = 1, hjust = 1, angle = 0),
        axis.text.y = element_text(vjust = 0.25, hjust = 1, size = 8),
        plot.margin = margin(0.7, 0.5, 1, 1.2, "cm"))
```

If we define $g_{i,j}$ as the genre for user’s i rating of movie j , we can use the following models to account for the *genre* effect:

To account for *genre effects* we will use the model suggested in the [Section 23.7: Exercises](#) of the *Chapter “23 Regularization” of the Course Textbook*[13]:

$$Y_{i,j} = \mu + \alpha_i + \beta_j + g_{i,j} + \varepsilon_{i,j}$$

where $g_{i,j}$ is an *aggregation function* which is explained in detail in *Section 22.3: “Review of Aggregation Functions” of “Recommender Systems Handbook”* (*Chapter 22: “Aggregation of Preferences in Recommender Systems”, p. 712*) book[14].

In the formula above $g_{i,j}$ denotes a *genre effect* for user’s i rating of movie j , so that:

$$g_{i,j} = \sum_{k=1}^K x_{i,j}^k \gamma_k$$

with $x_{i,j}^k = 1$ if $g_{i,j}$ includes genre k , and $x_{i,j}^k = 0$ otherwise.

$$Y_{i,j} = \mu + \alpha_i + \beta_j + g_{i,j} + f(d_{i,j})$$

$$\sum_{i=1}^{n_i} (Y_{i,j} - \mu - \alpha_i)$$

Conclusion

Hello Conclusion!

This is a great conclusion, isn't it??!

References

- [1] Rafael A. Irizarry. *Introduction to Data Science, Part II, Chapter 23: Regularization, Section 23.2: Loss function. Statistics and Prediction Algorithms Through Case Studies*. Dec. 27, 2024. URL: <https://rafaelab.dfci.harvard.edu/dsbook-part-2/highdim/regularization.html#sec-netflix-loss-function> (visited on 02/18/2025) (cit. on p. 1).
- [2] Robert M. Bell Andreas Toscher Michael Jahrer. *The BigChaos Solution to the Netflix Grand Prize. commendo research & consulting*. Sept. 5, 2009. URL: https://www.asc.ohio-state.edu/statistics/statgen/joul_aut2009/BigChaos.pdf (visited on 02/18/2025) (cit. on pp. 1, 26).
- [3] Azamat Kurbanayev. *edX Data Science: Capstone, MovieLens Datasets. Package: edx.capstone.movielens.data*. Version 0.0.0.9000. Feb. 5, 2025. URL: <https://github.com/AzKurban-edX-DS/edx.capstone.movielens.data> (visited on 02/05/2025) (cit. on p. 1).
- [4] Rafael A. Irizarry. *Introduction to Data Science, Part II. Statistics and Prediction Algorithms Through Case Studies*. Dec. 27, 2024. URL: <https://rafaelab.dfci.harvard.edu/dsbook-part-2/> (visited on 02/18/2025) (cit. on p. 3).
- [5] Rafael A. Irizarry. *Introduction to Data Science, Part II, Chapter 23: Regularization, Section 23.1.1: Movielens data. Statistics and Prediction Algorithms Through Case Studies*. Dec. 27, 2024. URL: <https://rafaelab.dfci.harvard.edu/dsbook-part-2/highdim/regularization.html#movielens-data> (visited on 02/18/2025) (cit. on pp. 3, 14, 16).
- [6] Amir Motefaker. *Movie Recommendation System using R - BEST*. Version 284. July 18, 2024. URL: <https://www.kaggle.com/code/amirmotefaker/movie-recommendation-system-using-r-best/notebook> (visited on 02/18/2025) (cit. on pp. 5–9, 29, 30, 33).
- [7] Azamat Kurbanayev. *edX Data Science: Capstone-MovieLens Project. A movie recommendation system using the MovieLens dataset*. Version 1.0.0.0. May 5, 2025. URL: <https://github.com/AzKurban-edX-DS/Capstone-MovieLens/tree/main> (visited on 05/05/2025) (cit. on p. 12).
- [8] Rafael A. Irizarry. *Introduction to Data Science, Part II, Chapter 23: Regularization, Section 23.3: A first model. Statistics and Prediction Algorithms Through Case Studies*. Dec. 27, 2024. URL: <https://rafaelab.dfci.harvard.edu/dsbook-part-2/highdim/regularization.html#a-first-model> (visited on 02/18/2025) (cit. on p. 26).
- [9] Rafael A. Irizarry. *Introduction to Data Science, Part II, Chapter 23: Regularization, Section 23.4: User effects. Statistics and Prediction Algorithms Through Case Studies*. Dec. 27, 2024. URL: <https://rafaelab.dfci.harvard.edu/dsbook-part-2/highdim/regularization.html#user-effects> (visited on 02/18/2025) (cit. on p. 27).
- [10] Robert Bell Yehuda Koren Yahoo Research and Chris Volinsky. *Matrix Factorization Techniques for Recommender Systems*. Aug. 1, 2009. URL: [https://datajobs.com/data-science-repo/Recommender-Systems-\[Netflix\].pdf](https://datajobs.com/data-science-repo/Recommender-Systems-[Netflix].pdf) (visited on 02/18/2025) (cit. on p. 27).
- [11] Rafael A. Irizarry. *Introduction to Data Science, Part II, Chapter 23: Regularization, Section 23.5: Movie effects. Statistics and Prediction Algorithms Through Case Studies*. Dec. 27, 2024. URL: <https://rafaelab.dfci.harvard.edu/dsbook-part-2/highdim/regularization.html#movie-effects> (visited on 02/18/2025) (cit. on p. 28).
- [12] Rafael A. Irizarry. *Introduction to Data Science, Part II, Chapter 23: Regularization, Section 23.6: Penalized least squares. Statistics and Prediction Algorithms Through Case Studies*. Dec. 27, 2024. URL: <https://rafaelab.dfci.harvard.edu/dsbook-part-2/highdim/regularization.html#penalized-least-squares> (visited on 02/18/2025) (cit. on p. 28).

- [13] Rafael A. Irizarry. *Introduction to Data Science, Part II, Chapter 23: Regularization, Section 23.7: Exercises. Statistics and Prediction Algorithms Through Case Studies*. Dec. 27, 2024. URL: <https://rafalab.dfci.harvard.edu/dsbook-part-2/highdim/regularization.html#exercises> (visited on 02/18/2025) (cit. on pp. 31, 33).
- [14] Francesco Ricci. *Recommender Systems Handbook*. Ed. by Paul B. Kantor Lior Rokach Bracha Shapira. Springer, New York, 2011. ISBN: ISBN 978-0-387-85819-7. DOI: [10.1007/978-0-387-85820-3](https://doi.org/10.1007/978-0-387-85820-3). URL: https://github.com/vwang0/recommender_system/blob/master/Recommender%20Systems%20Handbook.pdf (cit. on p. 33).