

BOOK RECOMMENDATION SYSTEM PROJECT

Reza Ameri, Ph.D.

Contents

1	Abstract	1
2	Introduction	2
3	Book-Crossing Dataset	2
4	Data Processing	3
5	Methodology and Procedure	10
5.1	Method I: Average Book Rating	11
5.2	Method II: Book Effect Model	12
5.3	Method III: Book & User Effect Model	12
5.4	Method IV: Regularized Book & User Effect Model	13
5.5	Method V: Parallel Matrix Factorization Model	14
6	Results	16
7	Summary and Conclusion	16
8	Acknowledgments	17
9	References	17

1 Abstract

Recommender systems have become prevalent in recent years as they tackle the problem of information overload by suggesting the most relevant products to end users. In fact, recommender systems are information filtering tools that aspire to predict the rating for users and items, predominantly from big data to recommend their likes. Specifically, book recommendation systems provide a mechanism to assist users in classifying users with similar interests. In this project, exploratory data analysis is used in order to develop various machine learning algorithms that predict book ratings with reasonable accuracy. The project is part of the capstone for the professional certificate in data science program at Harvard University. The Book-Crossing Dataset that includes 1,149,780 ratings (explicit/implicit) about 271,379 books is used for creating a book recommendation system.

2 Introduction

Fast development of technology and the internet has resulted in a rapid growth of available information over the past few decades. Recommendation systems, as one of the most successful information filtering applications, have become an efficient way to solve the information overload problem [1]. The main goal of a recommendation system is to automatically generate suggested items for end users based on their historical preferences.

Recommendation systems have become an indispensable component in various e-commerce applications [2]. Recommender systems collect information about the user's preferences of different items (e.g. books, movies, shopping, tourism, TV, taxi) by two ways, either implicitly or explicitly [3-7]. Explicit feedback, such as rating scales, provides users with a mechanism to unequivocally express their interests in items. On the other hand, implicit feedback is generated by the recommendation system itself, through inferences it makes about the user's behavior. What constitutes implicit feedback depends on the application domain: Typically, it will be one or multiple observable and measurable parameters that arise out of the user's interactions with the recommender system [2, 14].

In this project, a book recommendation system is developed using the Book-Crossing Dataset [8] which includes 1,149,780 ratings (explicit/implicit) about 271,379 books, Collected by Cai-Nicolas Ziegler [9]. Considering five different models, the Mean Square Error (RMSE) is used to evaluate how close the predictions are to the true values in the final hold-out test set. The model with the lowest RMSE is reported as the best model. It should be noted that, in the machine learning models, the final hold-out test set (i.e. the validation data) has not been used for training, developing, or selecting the algorithm, and has been only used for evaluating the RMSE of the final algorithm. Moreover, only explicit ratings are considered in developing the recommendation systems, as most of the research in the recommender systems field has focused on using one or the other type of feedback, and only few have combined both implicit and explicit feedback [14].

3 Book-Crossing Dataset

The Book-Crossing dataset can be downloaded from ref. [8]. Then, the following code can be run in order to generate the datasets:

```
# Create train and validation sets

# Note: this process could take a couple of minutes

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(recosystem)) install.packages("recosystem", repos = "http://cran.us.r-project.org")
if(!require(data.table)) install.packages("data.table", repos = "http://cran.us.r-project.org")

library(tidyverse)
library(caret)
library(data.table)

# Book-Crossing Dataset
# http://www2.informatik.uni-freiburg.de/~cziegler/BX/

dl <- tempfile()
download.file("http://www2.informatik.uni-freiburg.de/~cziegler/BX/BX-CSV-Dump.zip", dl)

ratings <- fread(text = gsub(":", "\t", readLines(unzip(dl, "BX-Book-Ratings.csv"))),
                  col.names = c("userId", "ISBN", "rating"))
```

```
books <- fread(text = gsub("::", "\t", readLines(unzip(dl, "BX-Books.csv"))),
  col.names = c("ISBN", "BookTitle", "BookAuthor", "YearOfPublication",
    "Publisher", "ImageURLS", "ImageURLM", "ImageURLL"))

bookcrossing <- left_join(ratings, books, by = "ISBN")
bookcrossing <- bookcrossing[, 1:(length(bookcrossing)-3)] # Remove the URLs
bookcrossing <- bookcrossing[complete.cases(bookcrossing), ] # Consider complete cases only
bookcrossing <- bookcrossing[apply(bookcrossing!=0, 1, all),] # Consider explicit ratings only

# Validation set will be 10% of bookcrossing data
set.seed(1, sample.kind="Rounding")
test_index <- createDataPartition(y = bookcrossing$rating, times = 1, p = 0.1, list = FALSE)
edx <- bookcrossing[-test_index,]
temp <- bookcrossing[test_index,]

# Make sure userId and ISBN in validation set are also in edx set
validation <- temp %>%
  semi_join(edx, by = "ISBN") %>%
  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, books, test_index, temp, bookcrossing, removed)
```

As previously mentioned, the algorithms will be developed using the train (i.e. edx) set. Book ratings will be then predicted in the validation set, for a final test of the final algorithm.

4 Data Processing

In order to get familiar with the dataset, the structure of the data as well as its summary are evaluated as shown below:

```
# Structure of the dataset
str(edx)
```

```
## Classes 'data.table' and 'data.frame': 353570 obs. of 7 variables:
## $ userId      : int 276726 276729 276729 276744 276747 276747 276747 276748 276751 ...
## $ ISBN        : chr  "0155061224" "052165615X" "0521795028" "038550120X" ...
## $ rating      : int  5 3 6 7 9 9 7 7 6 8 ...
## $ BookTitle    : chr  "Rites of Passage" "Help!: Level 1" "The Amsterdam Connection : Level 4 (
## $ BookAuthor   : chr  "Judith Rae" "Philip Prowse" "Sue Leather" "JOHN GRISHAM" ...
## $ YearOfPublication: int 2001 1999 2001 2001 2003 1995 1995 1998 2002 1998 ...
## $ Publisher    : chr  "Heinle" "Cambridge University Press" "Cambridge University Press" "Doub
## - attr(*, ".internal.selfref")=<externalptr>
```

```
# Headers of the dataset
head(edx) %>%
  print.data.frame()
```

```
##   userId      ISBN rating
## 1 276726 0155061224      5
## 2 276729 052165615X      3
## 3 276729 0521795028      6
## 4 276744 038550120X      7
## 5 276747 0060517794      9
## 6 276747 0671537458      9
##
##                                     BookTitle      BookAuthor
## 1                                     Rites of Passage      Judith Rae
## 2                                     Help!: Level 1      Philip Prowse
## 3 The Amsterdam Connection : Level 4 (Cambridge English Readers)      Sue Leather
## 4                                     A Painted House      JOHN GRISHAM
## 5                                     Little Altars Everywhere      Rebecca Wells
## 6                                     Waiting to Exhale      Terry McMillan
##   YearOfPublication      Publisher
## 1                2001      Heinle
## 2                1999 Cambridge University Press
## 3                2001 Cambridge University Press
## 4                2001      Doubleday
## 5                2003      HarperTorch
## 6                1995      Pocket
```

```
# Summary of the dataset
summary(edx)
```

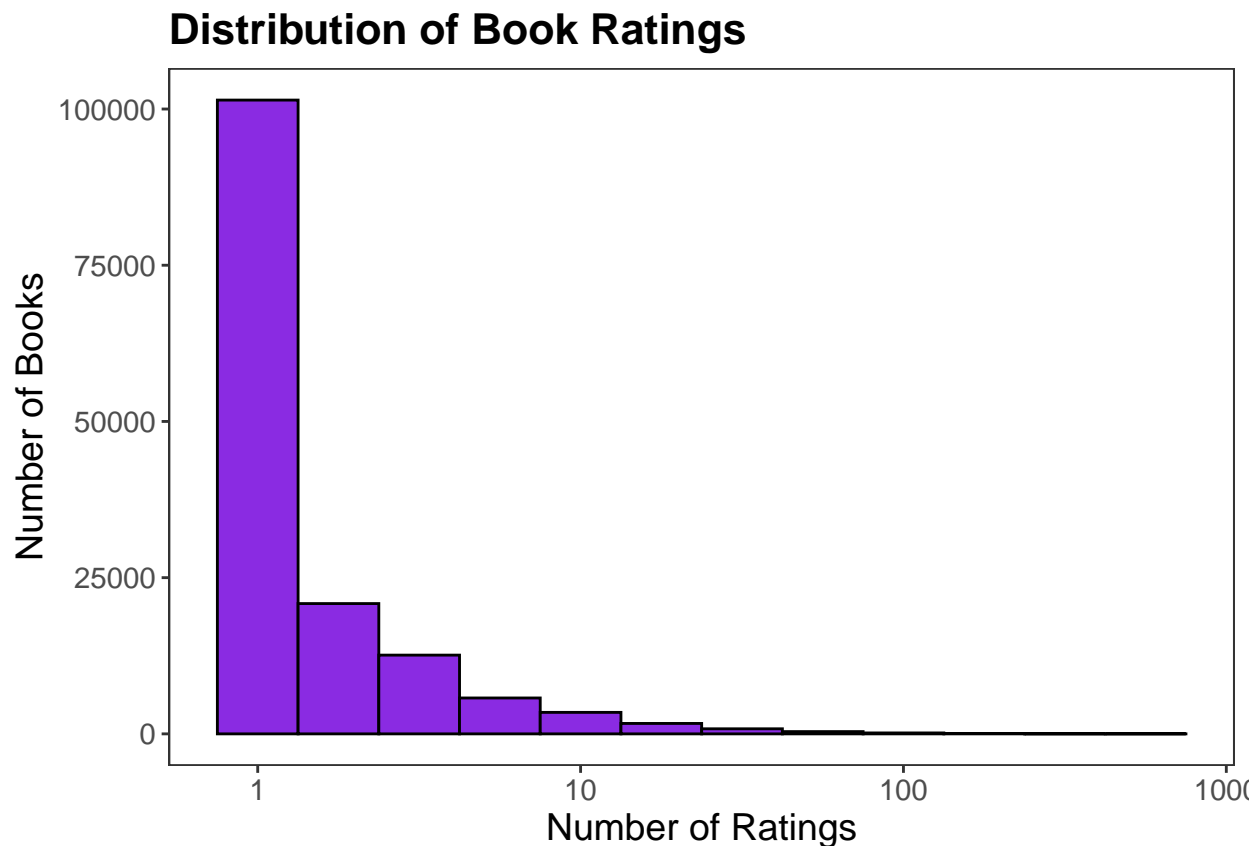
```
##      userId      ISBN      rating      BookTitle
## Min.   :      8   Length:353570   Min.   : 1.000   Length:353570
## 1st Qu.: 67840   Class :character   1st Qu.: 7.000   Class :character
## Median :134005   Mode  :character   Median : 8.000   Mode  :character
## Mean   :136124
## 3rd Qu.:206204
## Max.   :278854
##                                     Max.   :10.000
##   BookAuthor      YearOfPublication      Publisher
## Length:353570   Min.   :1376   Length:353570
## Class :character   1st Qu.:1992   Class :character
## Mode  :character   Median :1998   Mode  :character
##                                     Mean   :1996
##                                     3rd Qu.:2001
##                                     Max.   :2050
```

As can be seen, the dataset contains six variables, including “userId”, “ISBN”, “rating”, “BookTitle”, “BookAuthor”, “YearOfPublication”, and “Publisher”. The dataset is in the tidy format and ready for processing.

We do know, from experience, that some books are generally rated higher than others. We also know that there is substantial variability across users (e.g. some users love every book and others are very critical). There are some other important sources of variation related to the fact that groups of books have similar rating patterns and groups of users have similar rating patterns. In order to have a better understanding of the variability in data, several insightful plots can be made as presented below.

First, we can plot the number of ratings per book using the following code:

```
# Ratings per book distribution
edx %>%
  count(ISBN) %>%
  ggplot(aes(n)) +
  geom_histogram(binwidth = 0.25, color="black", fill="blueviolet") +
  scale_x_log10() +
  xlab("Number of Ratings") +
  ylab("Number of Books") +
  ggtitle("Distribution of Book Ratings") + theme_bw() +
  theme(plot.title = element_text(size = 16, face = "bold"),
        axis.text.y = element_text(size=11),
        axis.text.x = element_text(size=11),
        axis.title.y = element_text(size=14),
        axis.title.x = element_text(size=14),
        panel.grid.major = element_blank(),
        panel.grid.minor = element_blank(),
        panel.background=element_rect(fill = "white"))
```



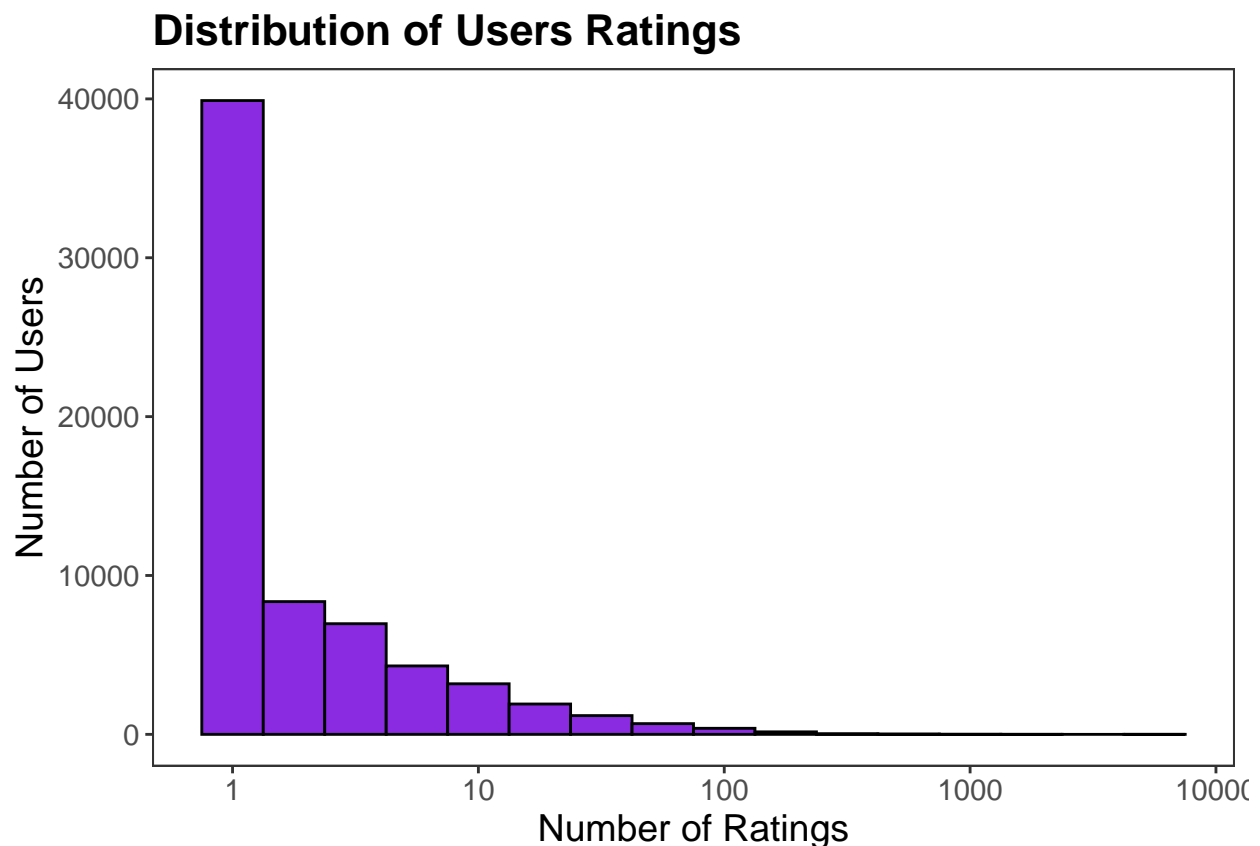
It can be seen that majority of books have been rated approximately between 1 to 10 times. The number of users' ratings frequency can also be plotted using the following code:

```
# User ratings distribution
edx %>%
  count(userId) %>%
  ggplot(aes(n)) +
  geom_histogram(binwidth = 0.25, color="black", fill="blueviolet") +
```

```

scale_x_log10() +
xlab("Number of Ratings") +
ylab("Number of Users") +
ggtitle("Distribution of Users Ratings") + theme_bw() +
theme(plot.title = element_text(size = 16, face = "bold"),
      axis.text.y = element_text(size=11),
      axis.text.x = element_text(size=11),
      axis.title.y = element_text(size=14),
      axis.title.x = element_text(size=14),
      panel.grid.major = element_blank(),
      panel.grid.minor = element_blank(),
      panel.background=element_rect(fill = "white"))

```



The plot represents that majority of users have rated approximately between 1 to 100 books. We can use the following code to plot the distribution of star ratings given to books by users:

```

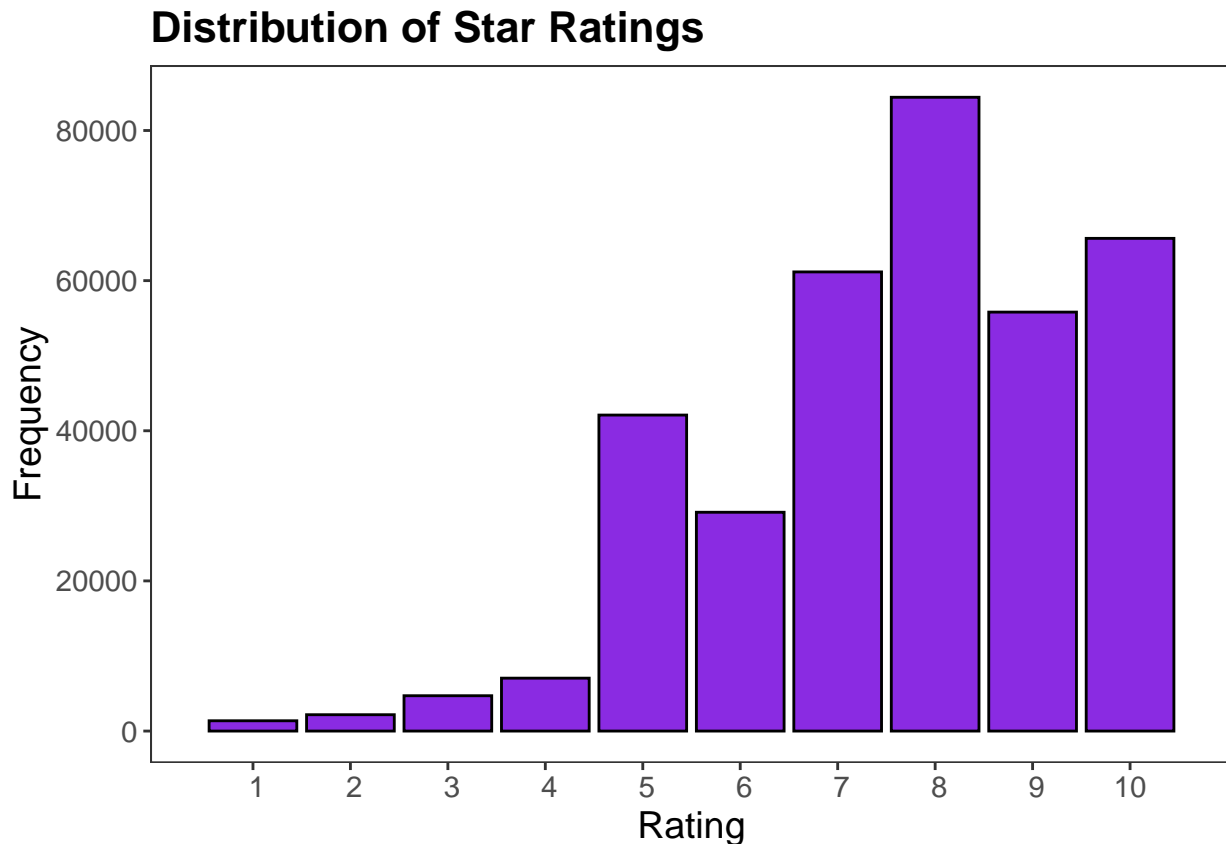
# Star ratings distribution
edx %>%
  ggplot(aes(rating)) +
  geom_bar(color="black", fill="blueviolet") +
  scale_x_discrete(limits = c(seq(1,10))) +
  xlab("Rating") +
  ylab("Frequency") +
  ggtitle("Distribution of Star Ratings") + theme_bw() +
  theme(plot.title = element_text(size = 16, face = "bold"),
        axis.text.y = element_text(size=11),

```

```

axis.text.x = element_text(size=11),
axis.title.y = element_text(size=14),
axis.title.x = element_text(size=14),
panel.grid.major = element_blank(),
panel.grid.minor = element_blank(),
panel.background=element_rect(fill = "white")

```



As can be seen, users who have rated books tend to assign higher scores in general (most users have rated a 5-star or higher, with an 8-star being the most common rating). Moreover, as previously mentioned users differ from one another regarding how critical they are, and while some users give higher star rates to majority of books they rate, others leave lower star ratings. This can be shown via the following code:

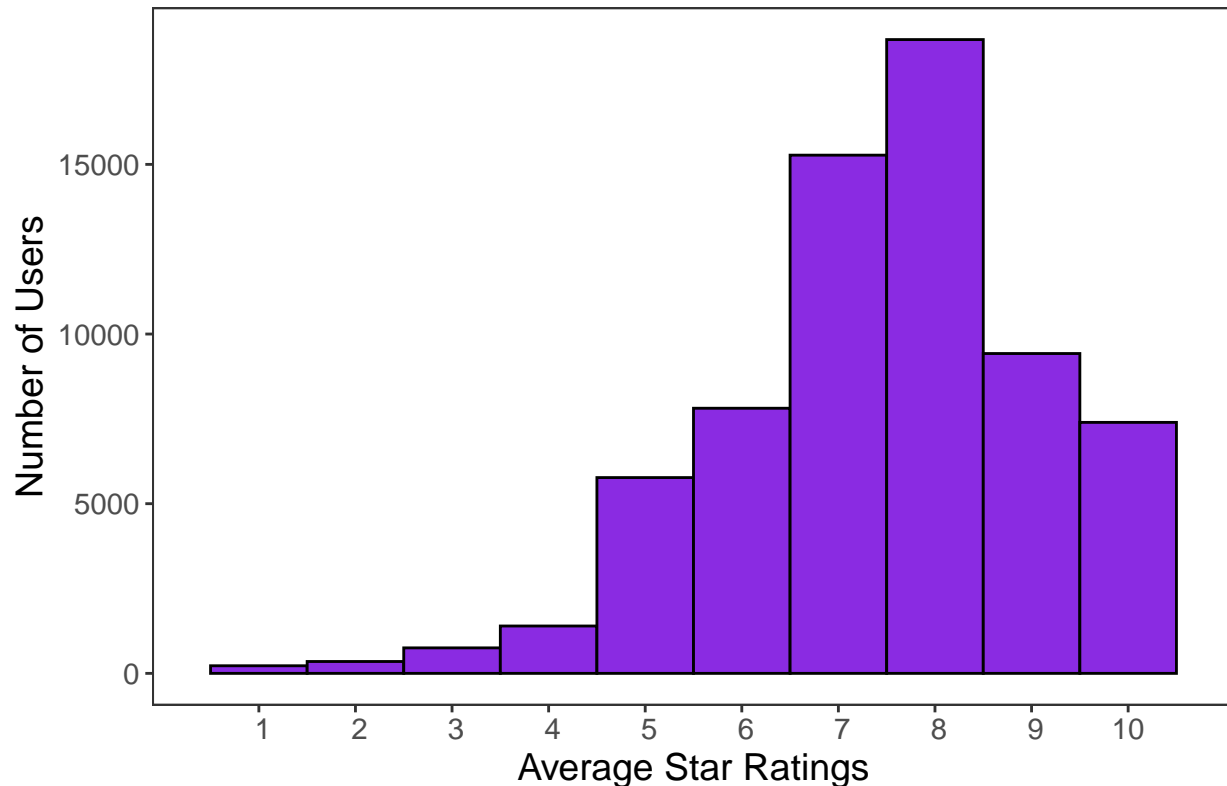
```

# Average book ratings by users
edx %>%
  group_by(userId) %>%
  summarize(b_u = mean(rating)) %>%
  ggplot(aes(b_u)) +
  geom_histogram(binwidth = 1, color="black", fill="blueviolet") +
  xlab("Average Star Ratings") +
  ylab("Number of Users") +
  ggtitle("Average Book Ratings by Users") +
  scale_x_discrete(limits = c(seq(1,10))) + theme_bw() +
  theme(plot.title = element_text(size = 16, face = "bold"),
        axis.text.y = element_text(size=11),
        axis.text.x = element_text(size=11),
        axis.title.y = element_text(size=14),

```

```
axis.title.x = element_text(size=14),
panel.grid.major = element_blank(),
panel.grid.minor = element_blank(),
panel.background=element_rect(fill = "white"))
```

Average Book Ratings by Users



Based on the plots shown above, users and book effects are apparent. Therefore, it is important to account for these effects in the machine learning models. It is noteworthy that, if one digs more into the data, they could potentially observe some other effects that could be included in the machine learning algorithms. Based upon common sense, for example, that the older the book, the higher its number of ratings will be. This can be shown in the following plot:

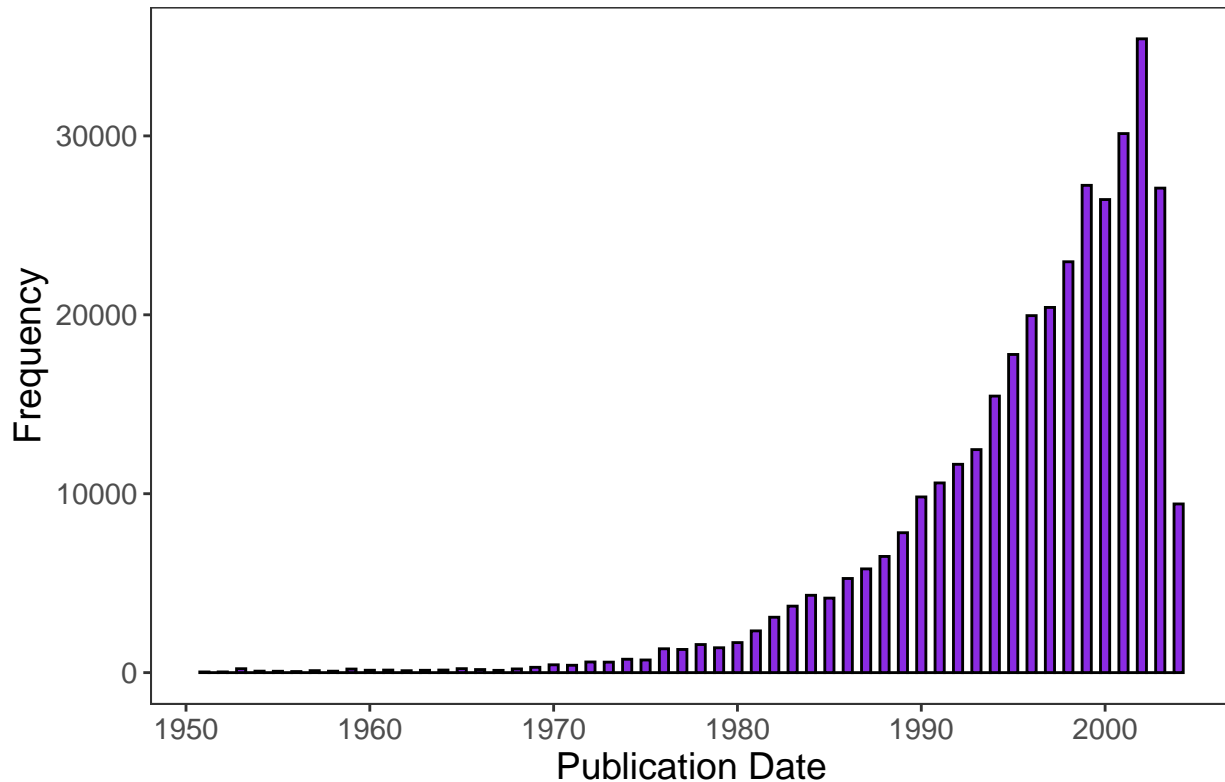
```
# Distribution of rated books based on the publication date
dat <- edx %>% filter(edx$YearOfPublication > 1950 & edx$YearOfPublication < 2005)

dat %>%
  ggplot(aes(dat$YearOfPublication)) +
  geom_histogram(binwidth = 0.5, color="black", fill="blueviolet") +
  xlab("Publication Date") +
  ylab("Frequency") +
  ggtitle("Average Rating Based on Publication Date") + theme_bw() +
  theme(plot.title = element_text(size = 16, face = "bold"),
        axis.text.y = element_text(size=11),
        axis.text.x = element_text(size=11),
        axis.title.y = element_text(size=14),
        axis.title.x = element_text(size=14),
        panel.grid.major = element_blank(),
```



```
panel.grid.minor = element_blank(),
panel.background=element_rect(fill = "white"))
```

Average Rating Based on Publication Date



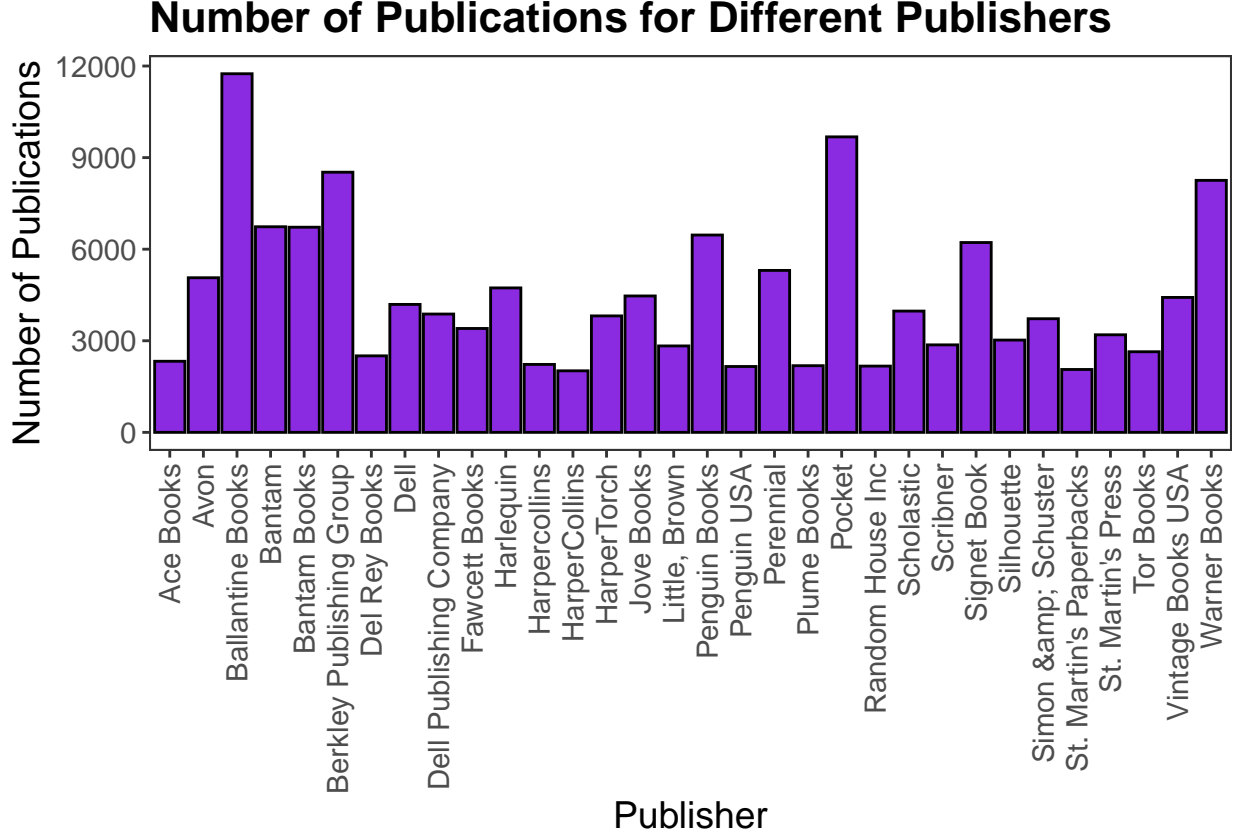
Similarly and again, based on common sense, different publishers might be more popular among users, and thus, books published by certain publishers could be purchased (and thus, rated) more frequently than others. The distribution of number of publications for each publisher can be plotted, using the following code:

```
# Number of publications for different publishers
mydata <- edx %>%
  group_by(Publisher) %>%
  summarise(count = n())

mydata <- mydata[order(-mydata$count),]

mydata %>%
  filter(count > 2000) %>%
  ggplot(aes(Publisher)) +
  geom_bar(aes(weights=count), color="black", fill="blueviolet") +
  xlab("Publisher") +
  ylab("Number of Publications") +
  ggtitle("Number of Publications for Different Publishers") + theme_bw() +
  theme(plot.title = element_text(size = 16, face = "bold"),
        axis.text.y = element_text(size=11),
        axis.text.x = element_text(angle = 90, vjust = 0.5, hjust=1, size=11),
        axis.title.y = element_text(size=14),
```

```
axis.title.x = element_text(size=14),
panel.grid.major = element_blank(),
panel.grid.minor = element_blank(),
panel.background=element_rect(fill = "white"))
```



Note that we are only looking at publishers with more than 2000 publications in the plot above. This is while some publishers have only one publications! Of course, there are other observations that could be made from the available dataset. Note that, however, the focus in this project will be only on the books and users effects. Nonetheless, in order to investigate the results from a more advanced machine learning algorithm, the recommendation system developed by Chin et al. [10] is used for the sake of comparison. This method is explained as the “Method V: Parallel Matrix Factorization Model” in the next section.

5 Methodology and Procedure

In this project, five different methods are implemented in order to develop the recommendation systems. In all the models, the Mean Square Error (RMSE) has been used as a metric to evaluate the performance of each machine learning (ML) model. The RMSE is a frequently used measure of the differences between values predicted by a model and the values actually observed, and can be defined as follows:

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

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

where N is the number of user/book combinations (i.e. samples), $\hat{y}_{u,i}$ are the predicted values, and $y_{u,i}$ are the observed values. Of course, the lower this metric is, the more accurate the model is.

5.1 Method I: Average Book Rating

This model is a simple model that predicts the same rating for all books, and ignores all the variability within the data, explained in the previous section. The mathematical formulation for this model is as follows:

$$Y_{u,i} = \mu + \epsilon_{u,i}$$

with $\epsilon_{u,i}$ being the independent errors sampled from the same distribution centered at 0 and μ the true rating for all books. This should be apparent that for this specific model, the value that minimizes the RMSE is the average of all ratings:

The average of the all ratings is:

```
mu <- mean(edx$rating)
cat("mu = ", mu)
```

```
## mu = 7.623704
```

If we predict all unknown ratings with μ , we obtain the first naive RMSE:

```
naive_rmse <- RMSE(validation$rating, mu)
cat("RMSE = ", naive_rmse)
```

```
## RMSE = 1.823992
```

In order to keep track of the results, we make a results table and save the RMSE values in this table for the sake of comparison:

```
rmse_results <- data_frame(Method = "Method I: Average Book Rating", RMSE = naive_rmse)
rmse_results %>% knitr::kable()
```

Method	RMSE
Method I: Average Book Rating	1.823992

This RMSE can be regarded as the loss value for the most basic ML model, and can be used as a baseline to evaluate the results from our next ML models.

In the next three models, we will be using some of the insights discussed in the previous section, in order to see how the results will change in reference to our base ML model.

5.2 Method II: Book Effect Model

In this model, we will specifically add the book effects to our base model as follows:

$$Y_{u,i} = \mu + b_i + \epsilon_{u,i}$$

where b_i represents average ranking for book i . As previously discussed, some books are generally rated higher than others, which could be due to their more popularity among users or the book age. This effect has been incorporated in our basic model using the term b_i .

```
book_avgs <- edx %>%
  group_by(ISBN) %>%
  summarize(b_i = mean(rating - mu))

predicted_ratings <- mu + validation %>%
  left_join(book_avgs, by='ISBN') %>%
  pull(b_i)
model_2_rmse <- RMSE(predicted_ratings, validation$rating)

rmse_results <- bind_rows(rmse_results,
  data_frame(Method = "Method II: Book Effect Model",
    RMSE = model_2_rmse ))
rmse_results %>% knitr::kable()
```

Method	RMSE
Method I: Average Book Rating	1.823992
Method II: Book Effect Model	1.969008

In this context, if a book is on average rated below the average rating for all books μ , the model predicts that it will be rated lower than μ by b_i (i.e. the difference of the individual book average from the total average). Hypothetically, we would expect the results to get better, but that is not the case here!

5.3 Method III: Book & User Effect Model

In this model, we will add both the book and user effects to our base model as follows:

$$Y_{u,i} = \mu + b_i + b_u + \epsilon_{u,i}$$

where b_u is a user-specific effect. As previously mentioned, it makes sense that some users are more cranky while others love every book they read. Thus, in this model, if a cranky user (negative b_u) rates a high-ranked book (positive b_i), the effects counter each other and we may be able to correctly predict that this cranky user may have given this good book a 3-star instead of 5-star.

Using this model which accounts for both user and book effects, the accuracy of the model will change, as can be seen below:

```
user_avgs <- edx %>%
  left_join(book_avgs, by='ISBN') %>%
  group_by(userId) %>%
  summarize(b_u = mean(rating - mu - b_i))
```

```

predicted_ratings <- validation%>%
  left_join(book_avgs, by='ISBN') %>%
  left_join(user_avgs, by='userId') %>%
  mutate(pred = mu + b_i + b_u) %>%
  pull(pred)
model_3_rmse <- RMSE(predicted_ratings, validation$rating)

rmse_results <- bind_rows(rmse_results,
  data_frame(Method = "Method III: Book & User Effect Model",
    RMSE = model_3_rmse))
rmse_results %>% knitr::kable()

```

Method	RMSE
Method I: Average Book Rating	1.823992
Method II: Book Effect Model	1.969008
Method III: Book & User Effect Model	1.867419

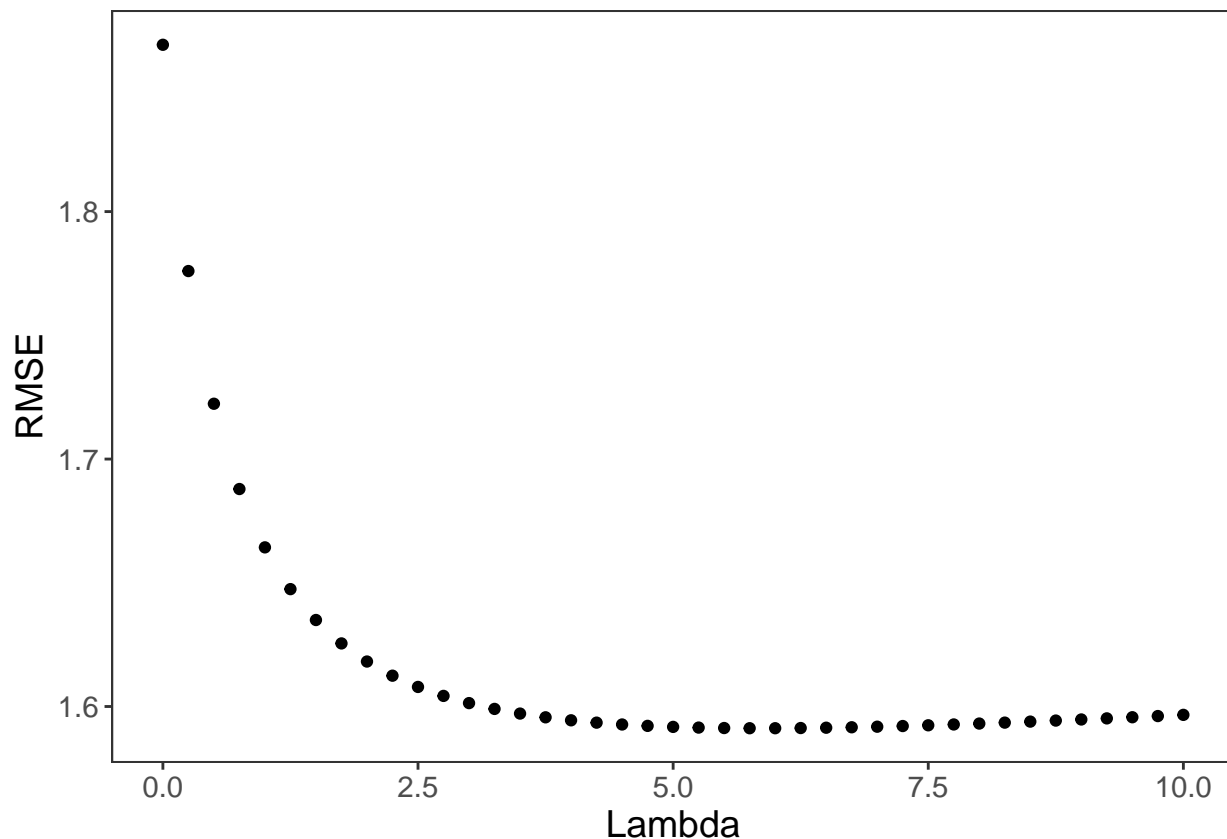
5.4 Method IV: Regularized Book & User Effect Model

In this model, we use regularization in order to make more accurate predictions. Specifically, model regularization is a fundamental technique for improving the generalization performance of a predictive model. Accordingly, many efficient optimization algorithms have been developed for solving various machine learning formulations with different regularizations [11].

For our specific problem, using a simple data analysis we can learn that the supposed best and worst books are rated by very few users. Of course, with few users, we have more uncertainty. Thus, larger estimates of b_i are more likely to happen. Regularization permits to penalize large estimates that are formed using small sample sizes. The general idea is to constrain the total variability of the effect sizes. In this context, instead of minimizing the least squares equation for our previous model, we minimize an equation that adds a penalty, as shown below:

$$\sum (y_{u,i} - \mu - b_i - b_u)^2 + \lambda(\sum b_i^2 + \sum b_u^2)$$

where λ is a tuning parameter and can be calculated using cross-validation. The optimal λ for the full model can be calculated as shown below:



The RMSE calculated using this model can be then calculated, using the following:

```
rmse_results <- bind_rows(rmse_results,
                          data_frame(Method = "Method IV: Regularized Book & User Effect Model",
                                     RMSE = min(rmses)))
rmse_results %>% knitr::kable()
```

Method	RMSE
Method I: Average Book Rating	1.823992
Method II: Book Effect Model	1.969008
Method III: Book & User Effect Model	1.867419
Method IV: Regularized Book & User Effect Model	1.591219

As can be seen, using the concept of regularization we can build a ML model that can predict the results more accurately, in comparison to the previous model explained.

5.5 Method V: Parallel Matrix Factorization Model

In this model, We will use a fast parallel stochastic gradient method for matrix factorization to predict the outcomes. Using the *recosystem* package in R as the wrapper of the *LIBMF* library, we can create a recommender system that will predict unknown entries in the rating matrix based on observed values [10].

A popular technique to solve the recommender system problem is the matrix factorization method. The idea is to approximate the whole rating matrix $R_{m \times n}$ by the product of two matrices of lower dimensions, $P_{k \times m}$

and $Q_{k \times n}$, such that:

$$R \approx P'Q$$

Let p_u be the u -th column of P , and q_v be the v -th column of Q , then the rating given by user u on item v would be predicted as $p_u'q_v$.

A typical solution for P and Q is given by the following optimization problem [12-13]:

$$\min_{P,Q} \sum [f(p_u, q_v; r_{u,v}) + \mu P \|p_u\|_1 + \mu Q \|q_v\|_1 + \frac{\lambda_P}{2} \|p_u\|_2^2 + \frac{\lambda_Q}{2} \|q_v\|_2^2]$$

In this methodology, the data used for training needs to be arranged in sparse matrix triplet form, (i.e. each line in the file should contain three numbers) in order to use *reco*system package. Moreover, the files will need to be saved as tables. Thus, we will create two new matrices (for training and validation) using the features *ISBN*, *userId*, and *rating*, as shown below:

```
edx_factorization <- edx %>% select(ISBN, userId, rating)
validation_factorization <- validation %>% select(ISBN, userId, rating)

edx_factorization <- as.matrix(edx_factorization)
validation_factorization <- as.matrix(validation_factorization)

write.table(edx_factorization, file = "trainingset.txt", sep = " ", row.names = FALSE,
  col.names = FALSE, quote = FALSE)

write.table(validation_factorization, file = "validationset.txt", sep = " ",
  row.names = FALSE, col.names = FALSE, quote = FALSE)

set.seed(1)
training_dataset <- data_file("trainingset.txt")

validation_dataset <- data_file("validationset.txt")

r = Reco() # this will create a model object
```

Similar to the previous model, we tune the algorithm to find the optimal answer:

```
opts = r$tune(training_dataset, opts = list(dim = c(10, 20, 30), lrate = c(0.1,
  0.2), costp_l1 = 0, costq_l1 = 0, nthread = 1, niter = 10))
```

Using the tuned parameters, we can train the model using the training set, as follows:

```
r$train(training_dataset, opts = c(opts$min, nthread = 1, niter = 20))
```

```
## warning: train on an empty training set
```

```
stored_prediction = tempfile()
```

We can then use the *predict()* function to make predictions on the validation set, as follows:

```

r$predict(validation_dataset, out_file(stored_prediction))

## prediction output generated at C:\Users\ADMINI~1\AppData\Local\Temp\2\Rtmp8KsxsL\file1e0459efd03

real_ratings <- read.table("validationset.txt", header = FALSE, sep = " ")$V3
pred_ratings <- scan(stored_prediction)

model_5_rmse <- RMSE(real_ratings, pred_ratings)
rmse_results <- bind_rows(rmse_results,
                          data_frame(Method = "Method V: Parallel Matrix Factorization Model",
                                    RMSE = model_5_rmse ))
rmse_results %>% knitr::kable()

```

Method	RMSE
Method I: Average Book Rating	1.823992
Method II: Book Effect Model	1.969008
Method III: Book & User Effect Model	1.867419
Method IV: Regularized Book & User Effect Model	1.591219
Method V: Parallel Matrix Factorization Model	1.238945

As can be seen, this ML model provides the most accurate predictions, by minimizing the RMSE in comparison to other models.

6 Results

The RMSE values obtained from the models considered in this project are as follows:

As can be seen, the best model is the Parallel Matrix Factorization Model, which results in a RMSE of 1.24. It is noteworthy that there are other ML models that could be used in order to develop an efficient recommendation system. Examples include Neural Networks and Evolutionary Algorithms. Depending on how the parameters are tuned and what variables are considered as predictors, one may achieve better results using the aforementioned methodologies. Implementing such models was beyond the scope of this project, however, and thus, the Parallel Matrix Factorization Model is introduced as the most efficient model in this project.

7 Summary and Conclusion

In this project, exploratory data analysis is used in order to develop various machine learning algorithms that predict book ratings with reasonable accuracy. The project is part of the capstone for the professional certificate in data science program at Harvard University. The Book-Crossing Dataset which includes 1,149,780 ratings (explicit/implicit) about 271,379 books is used for creating a book recommendation system. Considering five different models, the Mean Square Error (RMSE) is used to evaluate how close the predictions are to the true values in the final hold-out test set. The model with the lowest RMSE is reported as the best model. It should be noted that, in the machine learning models, the final hold-out test set (i.e. the validation data) has not been used for training, developing, or selecting the algorithm, and has been only used for evaluating the RMSE of the final algorithm.

Considering the results from the five different models considered in this project, the best model was the Parallel Matrix Factorization Model, which results in a RMSE of 1.24. The model is also reasonably fast,

and therefore, is reported as the most efficient algorithm to make the predictions in this project. It should be noted that, however, there are other ML models that could be implemented in order to develop an efficient recommendation system. Examples include Neural Networks and Evolutionary Algorithms. Depending on how the parameters are tuned for each model and what variables are considered as predictors, better predictions might be achieved using the aforementioned methodologies. Implementing such models was beyond the scope of this project, however, and not considered in the project.

8 Acknowledgments

This project is submitted as part of the capstone project for the professional certificate in data science program at Harvard University. The content expressed in this report are the views of the author and do not necessarily represent the opinions or views of the instructors.

9 References

- [1] Wang, Z., Yu, X., Feng, N., & Wang, Z. (2014). An improved collaborative movie recommendation system using computational intelligence. *Journal of Visual Languages & Computing*, 25(6), 667-675.
- [2] Katarya, R., & Verma, O. P. (2017). An effective collaborative movie recommender system with cuckoo search. *Egyptian Informatics Journal*, 18(2), 105-112.
- [3] Lu, J., Wu, D., Mao, M., Wang, W., & Zhang, G. (2015). Recommender system application developments: a survey. *Decision Support Systems*, 74, 12-32.
- [4] Bobadilla, J., Ortega, F., Hernando, A., & Gutiérrez, A. (2013). Recommender systems survey. *Knowledge-based systems*, 46, 109-132.
- [5] Chen, L., Chen, G., & Wang, F. (2015). Recommender systems based on user reviews: the state of the art. *User Modeling and User-Adapted Interaction*, 25(2), 99-154.
- [6] Konstan, J. A., & Riedl, J. (2012). Recommender systems: from algorithms to user experience. *User modeling and user-adapted interaction*, 22(1-2), 101-123.
- [7] Katarya, R., & Verma, O. P. (2016). Recent developments in affective recommender systems. *Physica A: Statistical Mechanics and its Applications*, 461, 182-190.
- [8] <http://www2.informatik.uni-freiburg.de/~ciegler/BX/>
- [9] Ziegler, C. N., McNee, S. M., Konstan, J. A., & Lausen, G. (2005, May). Improving recommendation lists through topic diversification. In *Proceedings of the 14th international conference on World Wide Web* (pp. 22-32).
- [10] Chin, W. S., Yuan, B. W., Yang, M. Y., Zhuang, Y., Juan, Y. C., & Lin, C. J. (2016). LIBMF: a library for parallel matrix factorization in shared-memory systems. *The Journal of Machine Learning Research*, 17(1), 2971-2975.
- [11] Zhu, D., Cai, C., Yang, T., & Zhou, X. (2018). A machine learning approach for air quality prediction: Model regularization and optimization. *Big data and cognitive computing*, 2(1), 5.
- [12] Chin, W. S., Zhuang, Y., Juan, Y. C., & Lin, C. J. (2015). A fast parallel stochastic gradient method for matrix factorization in shared memory systems. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 6(1), 1-24.
- [13] Chin, W. S., Zhuang, Y., Juan, Y. C., & Lin, C. J. (2015, May). A learning-rate schedule for stochastic gradient methods to matrix factorization. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining* (pp. 442-455). Springer, Cham.

[14] Jawaheer, G., Szomszor, M., & Kostkova, P. (2010, September). Comparison of implicit and explicit feedback from an online music recommendation service. In proceedings of the 1st international workshop on information heterogeneity and fusion in recommender systems (pp. 47-51).