

Movielens
Movie Recommendation System
A Harvard Capstone Project

Manoj Bijoor

March 15, 2021

Abstract

...this is the abstract text...

Contents

List of tables	iv
List of figures	v
List of Equations	vi
1 Project Overview: MovieLens - A Harvard Capstone Project	1
1.1 Create Train and Final Hold-out Test Sets	1
1.1.1 Important: Data sets usage	1
1.2 Final Product	1
1.2.1 My submission for this project is three files:	1
2 Exploratory Data Analysis	2
2.1 Data Wrangling	2
2.1.1 Initial data Exploration & Visualization	2
2.1.1.1 A Very Sparse Matrix	4
2.1.1.2 General Properties of the data	5
2.1.2 Further data Exploration, Visualization & Modification	7
2.1.2.1 Modify edx	7
2.1.2.2 Modify validation, repeat above steps	8
2.1.2.3 Movie Release Date - a closer look	9
2.1.2.4 Movie Rating Date-Time - a closer look	14
2.1.2.5 Genres combinations per movie - a closer look	16
3 Analysis - Model Building and Evaluation	17
3.1 Split the edx data into separate training and test sets	17
3.1.1 Loss function	17
3.2 Model 1: A first naive “mean” model	18
3.2.1 Results Table	19

List of tables

1	Movielens data	2
2	Unique Users, Movies and Genres	2
3	Matrix of seven users and four movies	3
4	Movielens edx data with rating date-time	7
5	Movielens edx data with movie release date	7
6	Movielens validation data with rating date-time	8
7	Movielens validation data with movie release date	8
8	Movielens edx data with average rating due to rating time effect	14
9	Movielens validation data with average rating due to rating time effect	15
10	RMSE Results	19

List of figures

1	A Very Sparse Matrix	4
2	Movies getting rated distribution	5
3	Users rating movies distribution	6
4	Ratings Movie Release Date - All dates	9
5	25 Movies with the most ratings per year and their average rating post 1991	10
6	25 Movies with the most ratings per year and their average rating post 1993	11
7	Movies average ratings versus ratings per year post 1993	12
8	Movies average ratings versus ratings per year pre 1993	13
9	Movies average ratings for each week versus day	14
10	Movies genres error bar plots	16

List of Equations

1	Equation 1	14
2	Equation 2	17
3	Equation 3	18

1 Project Overview: MovieLens - A Harvard Capstone Project

A movie recommendation system using the MovieLens dataset.

For this project, I will be creating a movie recommendation system using the MovieLens dataset, provided by GroupLens Research¹, a research lab in the Department of Computer Science and Engineering at the University of Minnesota, Twin Cities specializing in recommender systems, online communities, mobile and ubiquitous technologies, digital libraries, and local geographic information systems.

GroupLens Research² has collected and made available rating data sets from the MovieLens web site³. The data sets were collected over various periods of time, depending on the size of the set.

I will use the 10M version of the MovieLens dataset⁴ to make the computation a little easier.

First, I will download the MovieLens data and run code provided to generate my datasets.

Second, I will train a machine learning algorithm using the inputs in one subset to predict movie ratings in the validation set.

1.1 Create Train and Final Hold-out Test Sets

I will develop my algorithm using the edx set. For a final test of my final algorithm, I predict movie ratings in the validation set (the final hold-out test set) as if they were unknown. RMSE⁵ will be used to evaluate how close my predictions are to the true values in the validation set (the final hold-out test set).

1.1.1 Important: Data sets usage

The validation data (the final hold-out test set) will NOT be used for training, developing, or selecting my algorithm and it will ONLY be used for evaluating the RMSE of my final algorithm. The final hold-out test set will only be used at the end of my project with my final model. It will not be used to test the RMSE of multiple models during model development. I will split the edx data into separate training and test sets to design and test my algorithm.

1.2 Final Product

1.2.1 My submission for this project is three files:

1. My report in Rmd format
2. My report in PDF format (knit from my Rmd file)
3. A script in R format that generates my predicted movie ratings and RMSE score (contains all code and comments for my project)

The report documents the analysis and presents the findings, along with supporting statistics and figures. The report assumes that the reader is not familiar with the project or the data. The report includes the RMSE generated and the following sections:

1. an introduction/overview/executive summary section that describes the dataset and summarizes the goal of the project and key steps that were performed
2. a methods/analysis section that explains the process and techniques used, including data cleaning, data exploration and visualization, insights gained, and my modeling approach
3. a results section that presents the modeling results and discusses the model performance
4. a conclusion section that gives a brief summary of the report, its limitations and future work

¹<https://grouplens.org/>

²<https://grouplens.org/datasets/movielens/>

³<https://movielens.org>

⁴<https://grouplens.org/datasets/movielens/10m/>

⁵https://en.wikipedia.org/wiki/Root-mean-square_deviation

2 Exploratory Data Analysis

2.1 Data Wrangling

Data wrangling⁶, sometimes referred to as data munging, is the process of transforming and mapping data from one “raw” data form into another format with the intent of making it more appropriate and valuable for a variety of downstream purposes such as analytics.

The main steps could be described as follows:

1. Discovering
2. Structuring
3. Cleaning
4. Enriching
5. Validating
6. Publishing

Let’s perform the steps or combinations thereof starting with Initial data Exploration & Visualization in the next few subsections.

2.1.1 Initial data Exploration & Visualization

The first ten rows out of 9000055 rows of the Movielens data can be found in Table 1.

Table 1: Movielens data

userId	movieId	rating	timestamp	title	genres
1	122	5	838985046	Boomerang (1992)	Comedy Romance
1	185	5	838983525	Net, The (1995)	Action Crime Thriller
1	292	5	838983421	Outbreak (1995)	Action Drama Sci-Fi Thriller
1	316	5	838983392	Stargate (1994)	Action Adventure Sci-Fi
1	329	5	838983392	Star Trek: Generations (1994)	Action Adventure Drama Sci-Fi
1	355	5	838984474	Flintstones, The (1994)	Children Comedy Fantasy
1	356	5	838983653	Forrest Gump (1994)	Comedy Drama Romance War
1	362	5	838984885	Jungle Book, The (1994)	Adventure Children Romance
1	364	5	838983707	Lion King, The (1994)	Adventure Animation Children Drama Musical
1	370	5	838984596	Naked Gun 33 1/3: The Final Insult (1994)	Action Comedy

Each row represents a rating given by one user to one movie.

We can see the number of unique users that provided ratings and how many unique movies were rated. The unique genres here is based on a column called “genres” that includes every genre that applies to the movie. Some movies fall under several genres. Defining a category as whatever combination appears in this column, we are going to refer to this category as “unique genres” or simply “genres” from now on.

The number of unique users, movies and genres can be found in Table 2.

Table 2: Unique Users, Movies and Genres

unique_users	unique_movies	unique_genres
69878	10677	797

If we multiply the number of unique users by number of unique movies, we get a very large number, actually 746087406, yet our data table has 9000055 rows. This implies that not every user rated every movie. So

⁶https://en.wikipedia.org/wiki/Data_wrangling

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. The ‘gather’ function permits us to convert it to this format, but if we try it for the entire matrix, it will crash R.

Let’s show the matrix of seven users ie: userId’s 13-20 and four movies in Table 3.

Table 3: Matrix of seven users and four movies

userId	Forrest Gump (1994)	Jurassic Park (1993)	Pulp Fiction (1994)	Silence of the Lambs, The (1991)
13	NA	NA	4	NA
16	NA	3	NA	NA
17	NA	NA	NA	5
18	NA	3	5	5
19	4	1	NA	NA

2.1.1.1 A Very Sparse Matrix You can think of the task of a recommendation system as filling in the ‘NAs’ in the table above. To see how sparse the matrix is, here is the matrix in Figure 1 for a random sample of 100 movies and 100 users with yellow indicating a user/movie combination for which we have a rating.

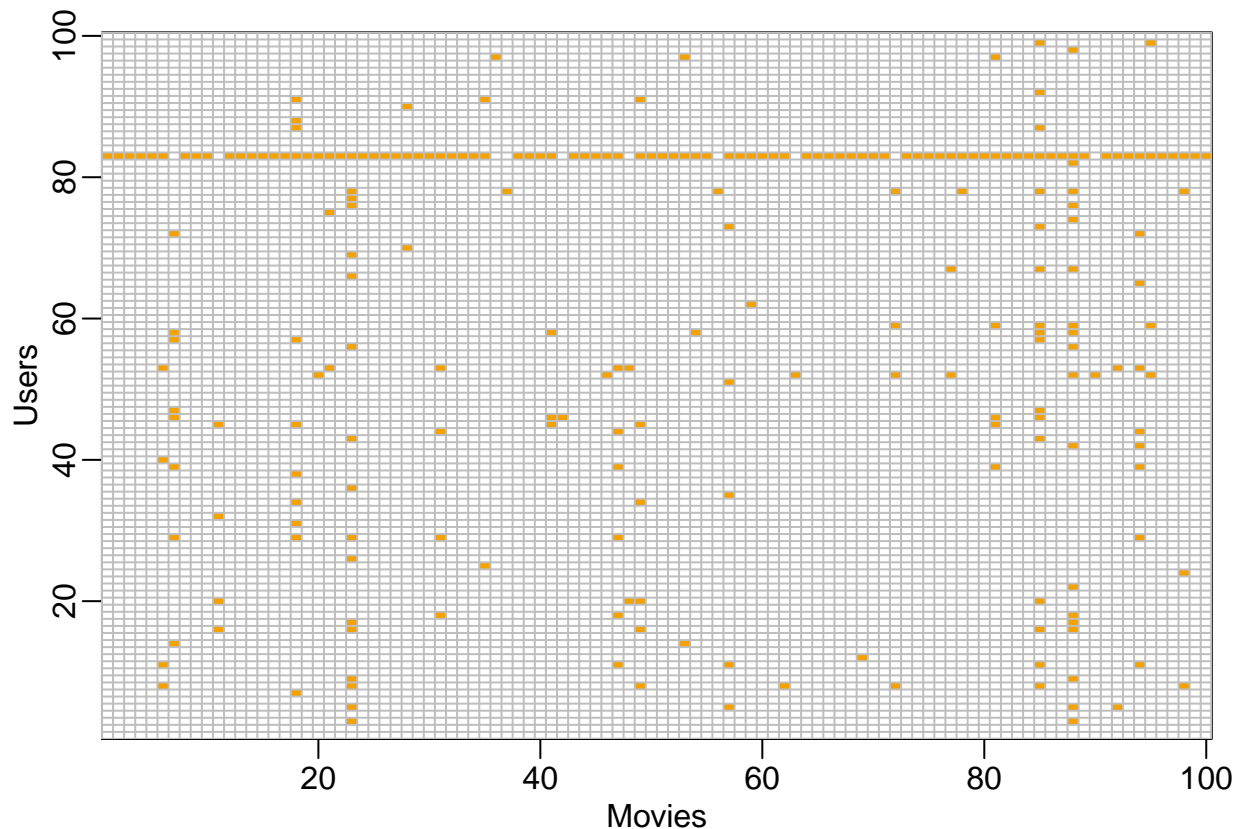


Figure 1: A Very Sparse Matrix

This machine learning challenge is quite complicated, because each outcome Y has a different set of predictors. To see this, note that if we are predicting the rating for movie i by user u , in principle, all other ratings related to movie i and by user u may be used as predictors, but different users rate different movies and a different number of movies. Furthermore, we may be able to use information from other movies that we have determined are similar to movie i or from users determined to be similar to user u . In essence, the entire matrix can be used as predictors for each cell.

2.1.1.2 General Properties of the data Let's look at some of the *general properties* of the data to better understand the challenges.

The *first thing* we notice is that some movies get rated more than others. Figure 2 shows the Movies getting rated distribution. This should not surprise us given that there are blockbuster movies watched by millions and artsy, independent movies watched by just a few:

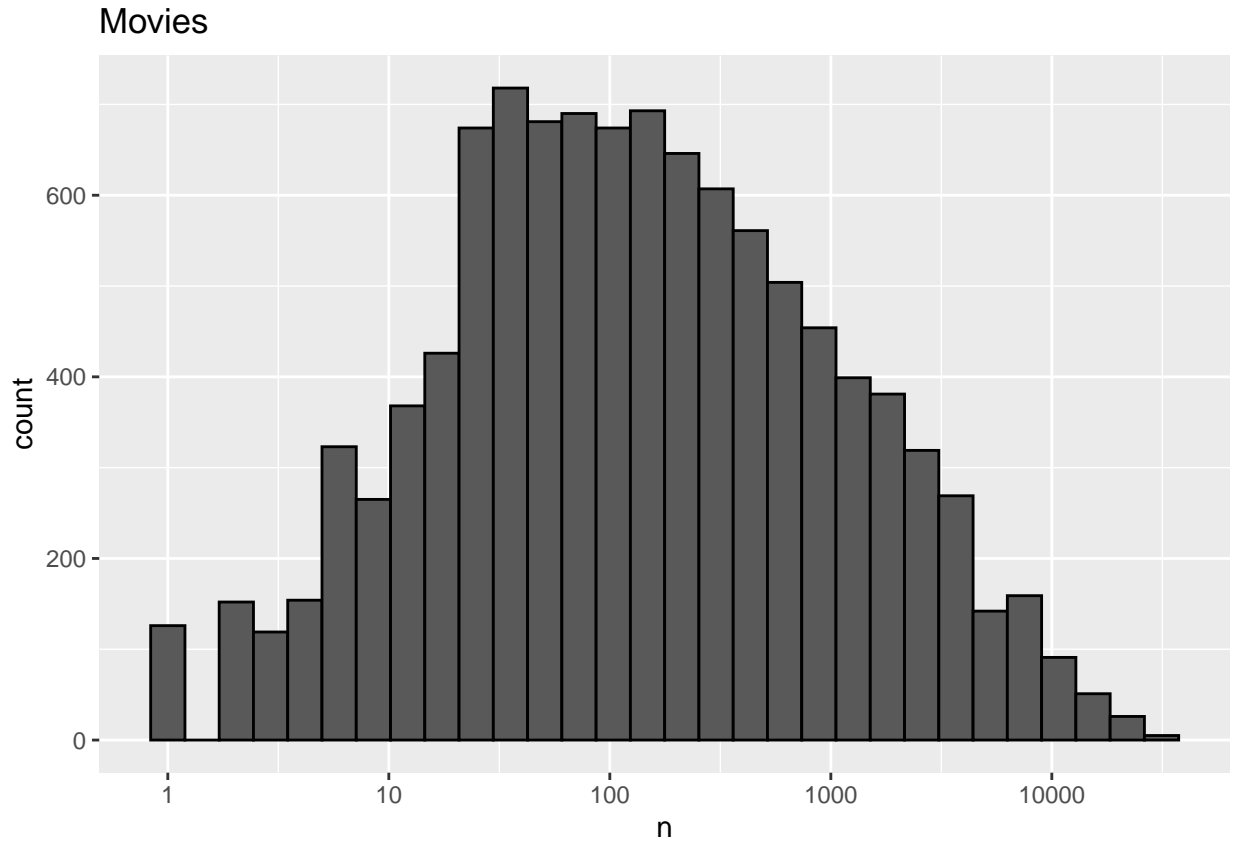


Figure 2: Movies getting rated distribution

Our *second observation* is that some users are more active than others at rating movies. Figure 3 shows Users rating movies distribution:

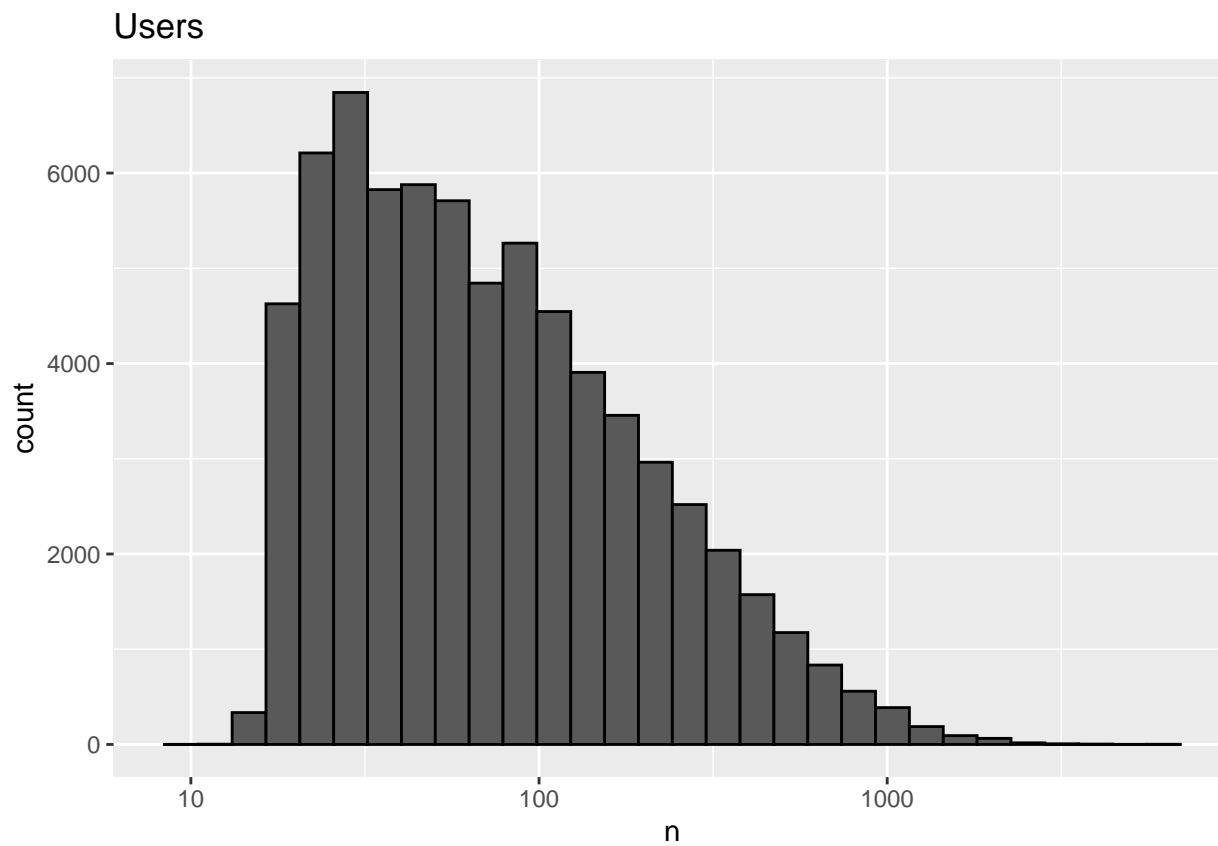


Figure 3: Users rating movies distribution

2.1.2 Further data Exploration, Visualization & Modification

2.1.2.1 Modify edx Convert **timestamp** in Movielens edx data Table 1 into date-time, a more readable and useful format named *rating_date* in Table 4 below.

Table 4: Movielens edx data with rating date-time

userId	movieId	rating	title	genres	rating_date
1	122	5	Boomerang (1992)	Comedy Romance	1996-08-02 11:24:06
1	185	5	Net, The (1995)	Action Crime Thriller	1996-08-02 10:58:45
1	292	5	Outbreak (1995)	Action Drama Sci-Fi Thriller	1996-08-02 10:57:01
1	316	5	Stargate (1994)	Action Adventure Sci-Fi	1996-08-02 10:56:32
1	329	5	Star Trek: Generations (1994)	Action Adventure Drama Sci-Fi	1996-08-02 10:56:32

Split title in Movielens edx data Table 1 into title and year movie released, a more useful format named *movie_dt* in Table 5 below.

Table 5: Movielens edx data with movie release date

userId	movieId	rating	title	genres	rating_date	movie_dt
1	122	5	Boomerang	Comedy Romance	1996-08-02 11:24:06	1992
1	185	5	Net, The	Action Crime Thriller	1996-08-02 10:58:45	1995
1	292	5	Outbreak	Action Drama Sci-Fi Thriller	1996-08-02 10:57:01	1995
1	316	5	Stargate	Action Adventure Sci-Fi	1996-08-02 10:56:32	1994
1	329	5	Star Trek: Generations	Action Adventure Drama Sci-Fi	1996-08-02 10:56:32	1994

2.1.2.2 Modify validation, repeat above steps Convert **timestamp** in Movielens validation data Table 1 into date-time, a more readable and useful format named *rating_date* in Table 6 below.

Table 6: Movielens validation data with rating date-time

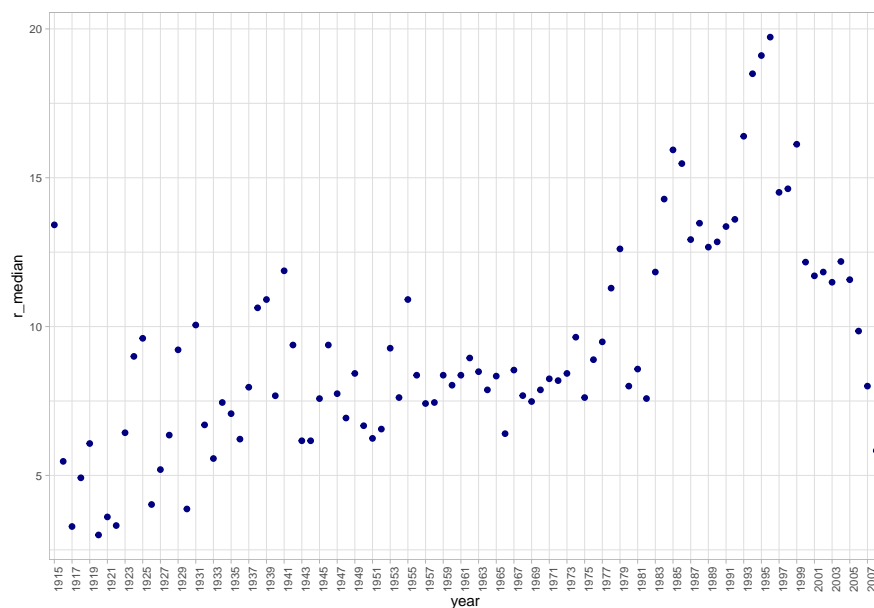
userId	movieId	rating	title	genres	rating_date
1	231	5	Dumb & Dumber (1994)	Comedy	1996-08-02 10:56:32
1	480	5	Jurassic Park (1993)	Action Adventure Sci-Fi Thriller	1996-08-02 11:00:53
1	586	5	Home Alone (1990)	Children Comedy	1996-08-02 11:07:48
2	151	3	Rob Roy (1995)	Action Drama Romance War	1997-07-07 03:34:10
2	858	2	Godfather, The (1972)	Crime Drama	1997-07-07 03:20:45

Split title in Movielens validation data Table 1 into title and year movie released, a more useful format named *movie_dt* in Table 7 below.

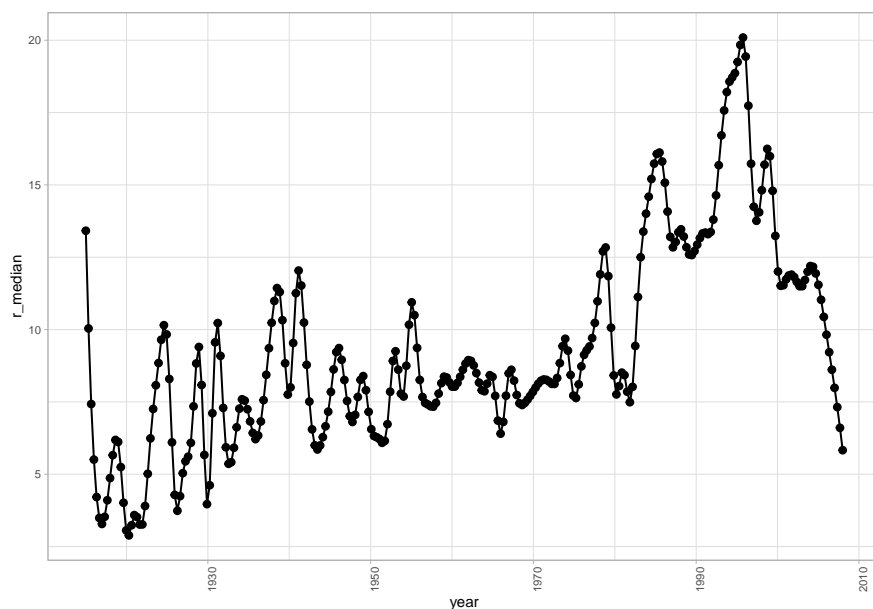
Table 7: Movielens validation data with movie release date

userId	movieId	rating	title	genres	rating_date	movie_dt
1	231	5	Dumb & Dumber	Comedy	1996-08-02 10:56:32	1994
1	480	5	Jurassic Park	Action Adventure Sci-Fi Thriller	1996-08-02 11:00:53	1993
1	586	5	Home Alone	Children Comedy	1996-08-02 11:07:48	1990
2	151	3	Rob Roy	Action Drama Romance War	1997-07-07 03:34:10	1995
2	858	2	Godfather, The	Crime Drama	1997-07-07 03:20:45	1972

2.1.2.3 Movie Release Date - a closer look Computing the number of ratings for each movie and then plotting it against the year the movie came out, that is the release date and using the square root transformation on the counts using Table 5 , we get see Figure 4 :



(a) All data points only



(b) Smooth line through all data points

Figure 4: Ratings Movie Release Date - All dates

we see that, on average, movies that came out after 1993 get more ratings. We also see that with newer movies, starting in 1993, the number of ratings decreases with year: the more recent a movie is, the less time users have had to rate it.

Among movies that came out in 1991 or later, we select the top 25 movies with the highest average number of ratings per year (n/year) and calculate the average rating of each of them. To calculate number of ratings per year, use 2018 as the end year. See Figure 5 :

`'geom_smooth()' using method = 'loess' and formula 'y ~ x'`

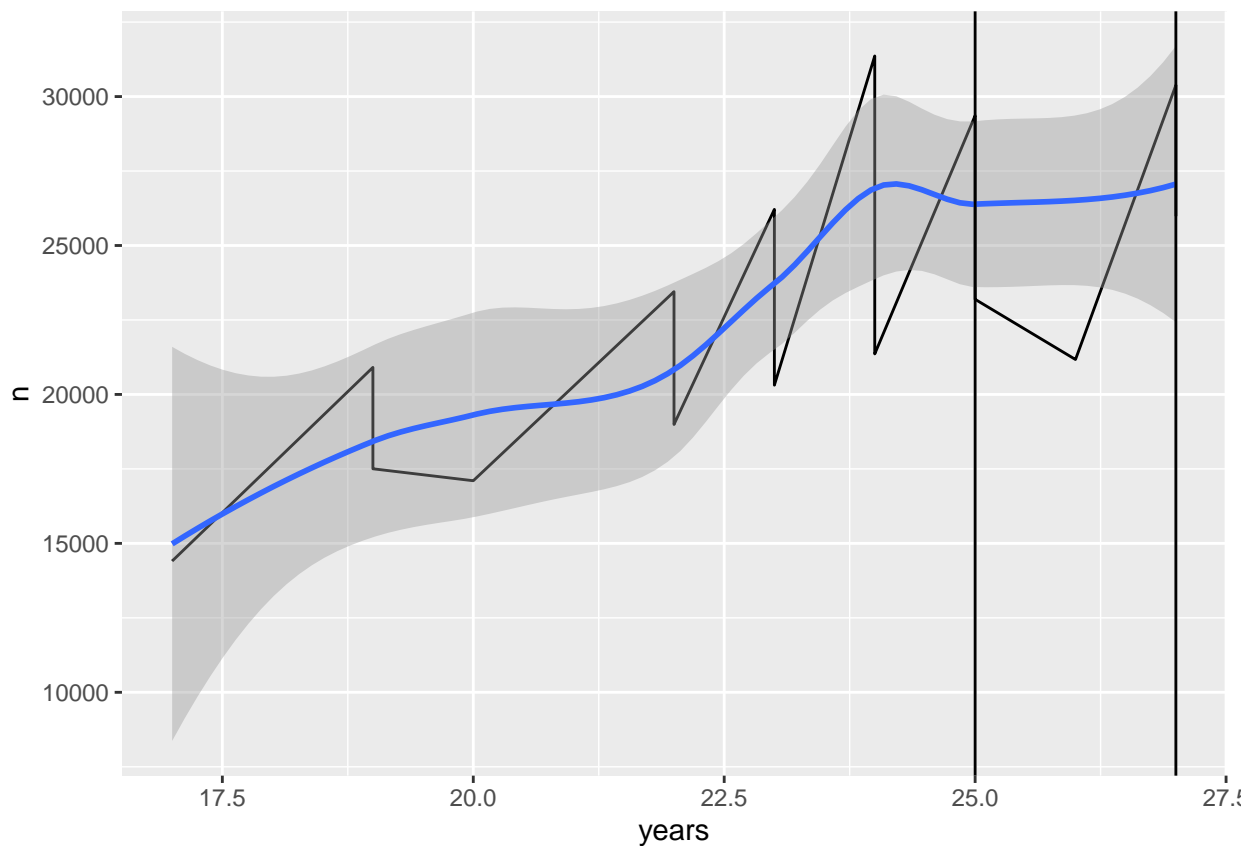


Figure 5: 25 Movies with the most ratings per year and their average rating post 1991

Among movies that came out in 1993 or later, we select the top 25 movies with the highest average number of ratings per year (n/year) and calculate the average rating of each of them. To calculate number of ratings per year, use 2018 as the end year. See Figure 6 :

`'geom_smooth()' using method = 'loess' and formula 'y ~ x'`

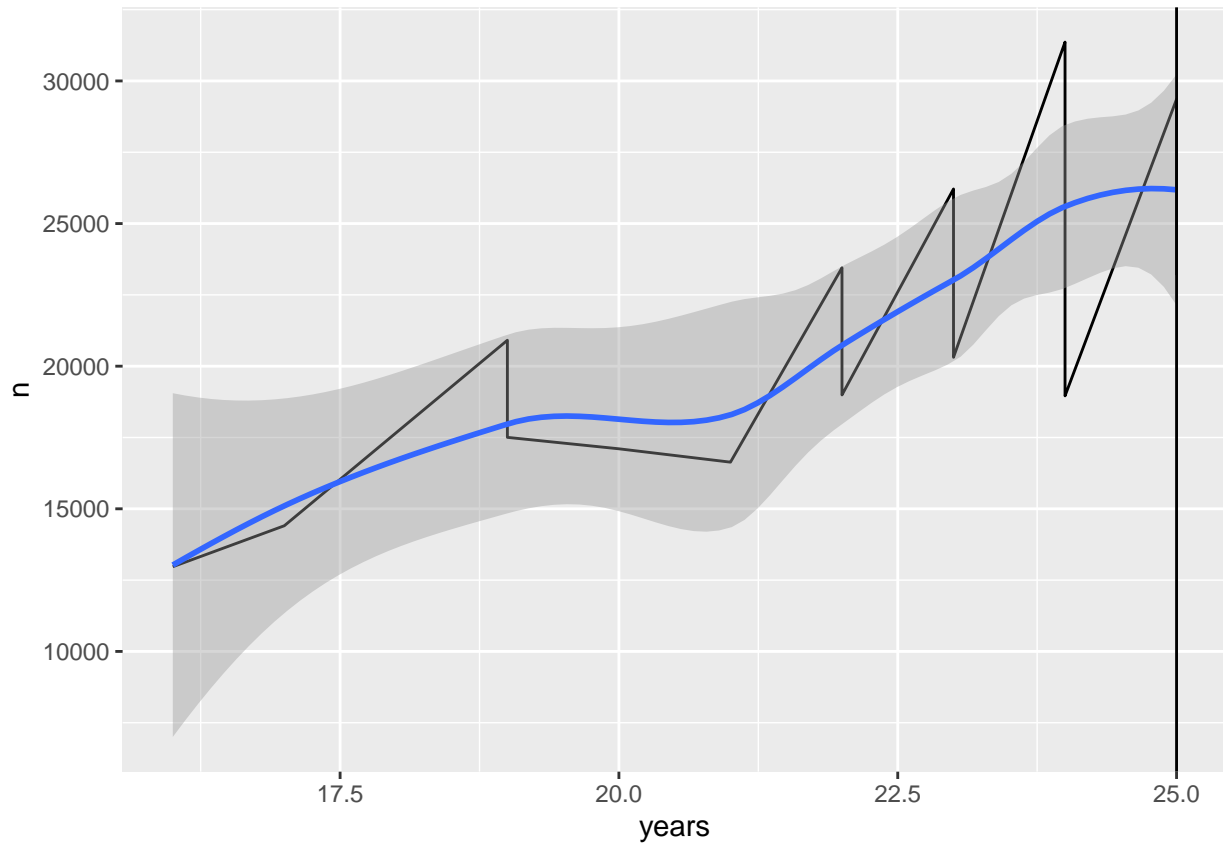


Figure 6: 25 Movies with the most ratings per year and their average rating post 1993

We see that the most rated movies tend to have above average ratings. This is not surprising: more people watch popular movies. To confirm this, we stratify the post 1993 movies by ratings per year and compute their average ratings. Figure 7 is a plot of average ratings versus ratings per year showing an estimate of the trend.

We see that the more a movie is rated, the higher the rating.

Post-1993 movies

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

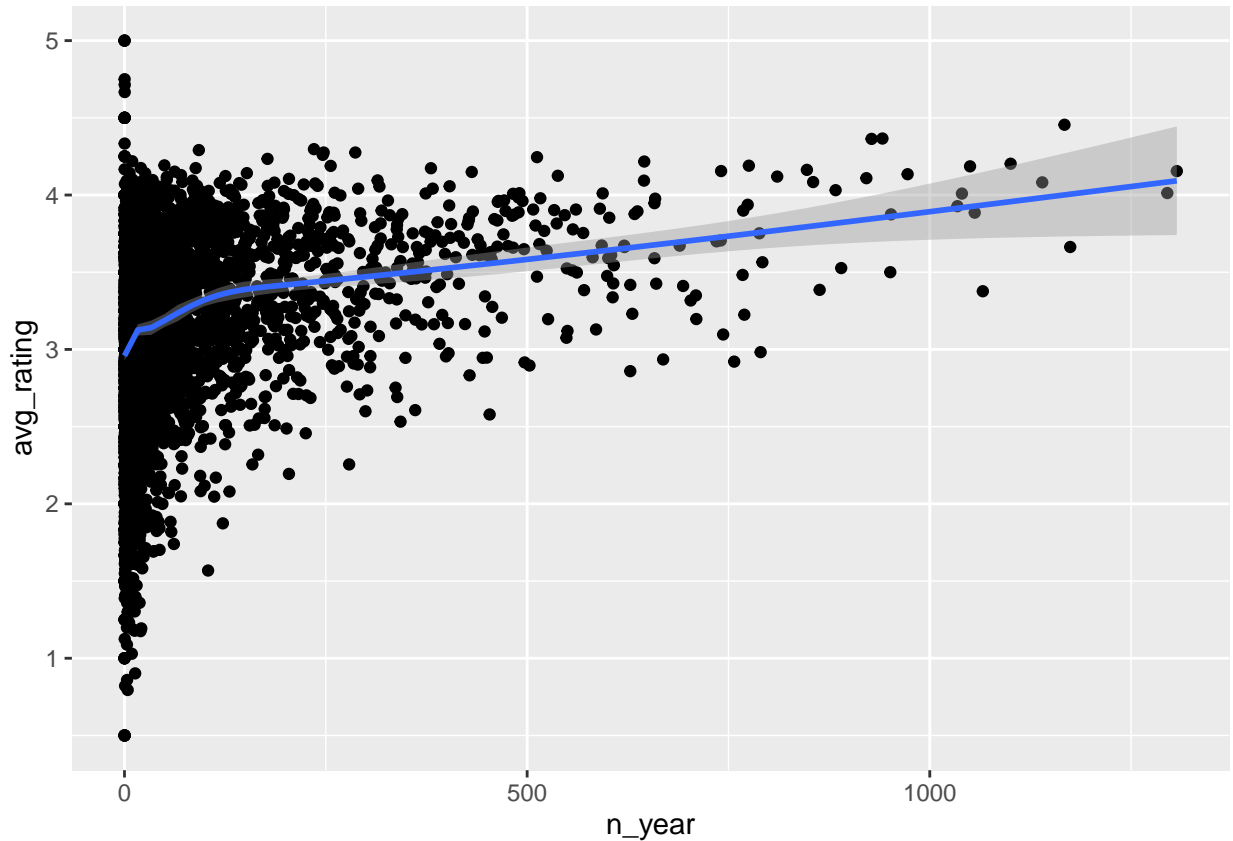


Figure 7: Movies average ratings versus ratings per year post 1993

Pre-1993 movies

Compare Pre-1993 movies trend shown here in Figure 8 Versus Post-1993 movies trend in Figure 7 above.

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

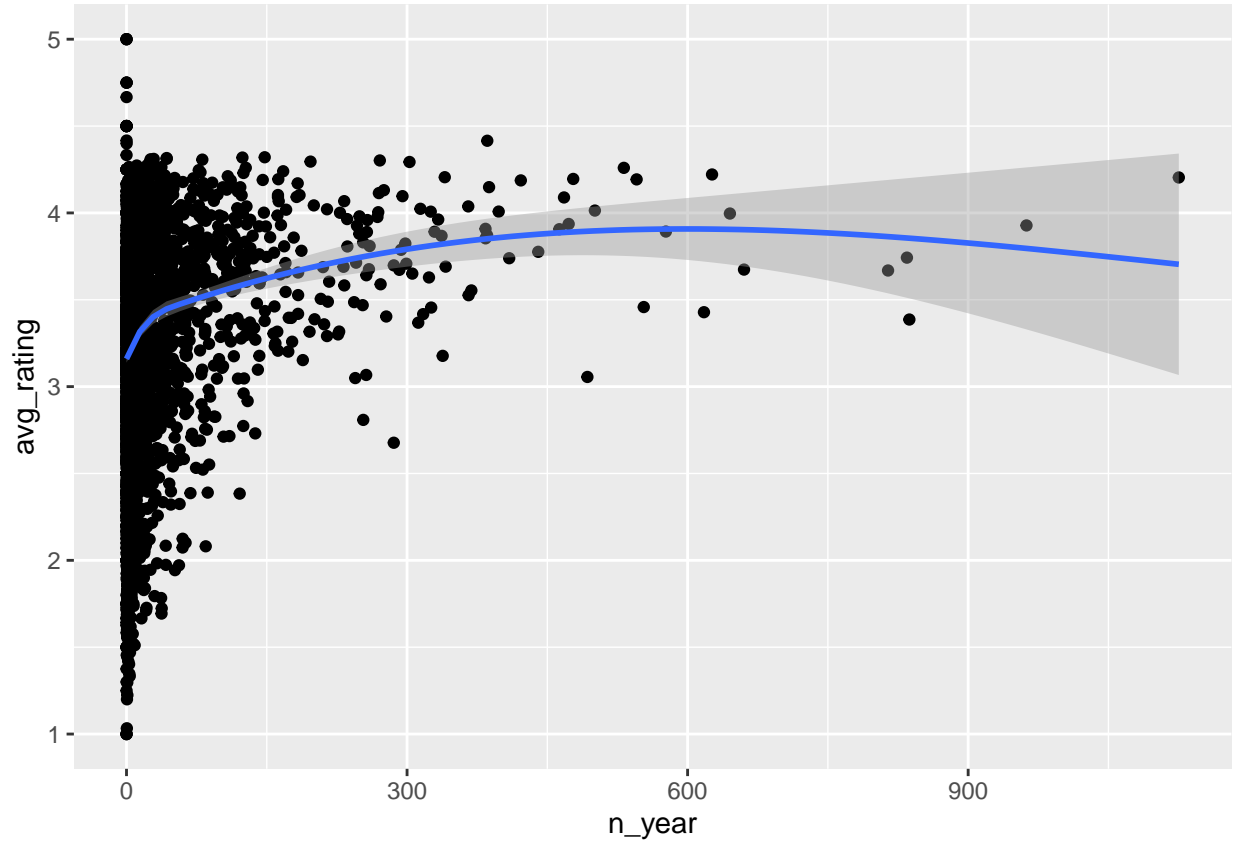


Figure 8: Movies average ratings versus ratings per year pre 1993

2.1.2.4 Movie Rating Date-Time - a closer look The Movielens edx data Table 1 also includes a time stamp. This variable represents the time and date in which the rating was provided. The units are seconds since January 1, 1970. We create a new column date with the date named *rating_date* in subsection [Modify edx](#) to get Table 5 .

We compute the average rating for each week and plot this average against day. See Figure 9 :

`'geom_smooth()' using method = 'loess' and formula 'y ~ x'`

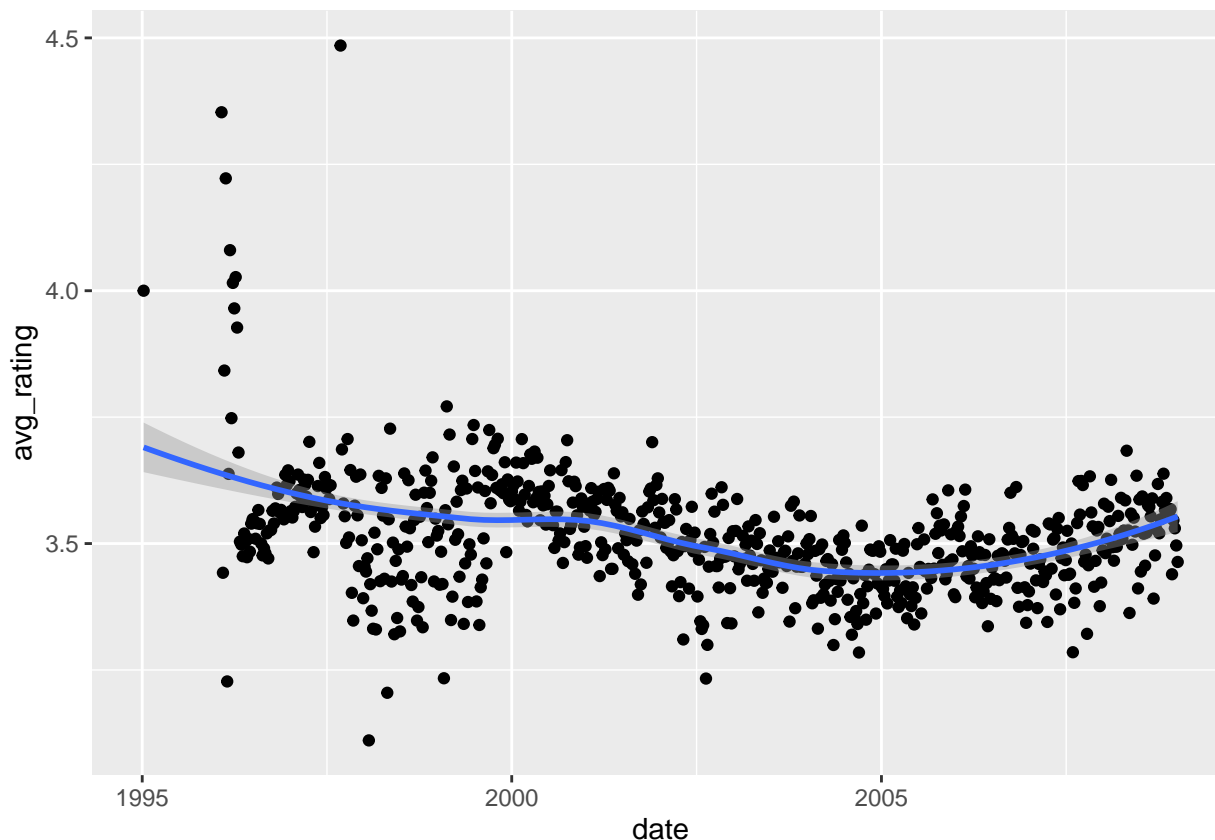


Figure 9: Movies average ratings for each week versus day

The plot shows some evidence of a time effect. If we define $d_{u,i}$ as the day for user's u rating of movie i , then the following model given by Equation 1 is most appropriate:

$$Y_{u,i} = \mu + b_i + b_u + f(d_{u,i}) + \epsilon_{u,i}, \text{ with } f \text{ a smooth function of } d_{u,i} \quad (1)$$

Modify edx Let's update the *edx* table to get Table 8 below:

Table 8: Movielens edx data with average rating due to rating time effect

userId	movieId	rating	title	genres	rating_date	movie_dt	date	avg_rating
1	122	5	Boomerang	Comedy Romance	1996-08-02 11:24:06	1992	1996-08-04	3.538801
1	185	5	Net, The	Action Crime Thriller	1996-08-02 10:58:45	1995	1996-08-04	3.538801
1	292	5	Outbreak	Action Drama Sci-Fi Thriller	1996-08-02 10:57:01	1995	1996-08-04	3.538801
1	316	5	Stargate	Action Adventure Sci-Fi	1996-08-02 10:56:32	1994	1996-08-04	3.538801
1	329	5	Star Trek: Generations	Action Adventure Drama Sci-Fi	1996-08-02 10:56:32	1994	1996-08-04	3.538801

TODO: Repeat above for validation data as well and somehow add this to the modelling section

Modify validation We need to do the above *avg_rating_time_effect* update for the validation data as well. Let's update the *validation* table to get Table 9 below:

Table 9: Movielens validation data with average rating due to rating time effect

userId	movieId	rating	title	genres	rating__date	movie_dt	date	avg__rating
1	231	5	Dumb & Dumber	Comedy	1996-08-02 10:56:32	1994	1996-08-04	3.555820
1	480	5	Jurassic Park	Action Adventure Sci-Fi Thriller	1996-08-02 11:00:53	1993	1996-08-04	3.555820
1	586	5	Home Alone	Children Comedy	1996-08-02 11:07:48	1990	1996-08-04	3.555820
2	151	3	Rob Roy	Action Drama Romance War	1997-07-07 03:34:10	1995	1997-07-06	3.606571
2	858	2	Godfather, The	Crime Drama	1997-07-07 03:20:45	1972	1997-07-06	3.606571

2.1.2.5 Genres combinations per movie - a closer look The movielens data Table 8 also has a genres column. This column includes every genre that applies to the movie. Some movies fall under several genres. We define a category as whatever combination appears in this column.

Here we keep only categories with more than 1,000 ratings. Then compute the average and standard error for each category, and plot these as error bar plots. See Figure 10 :

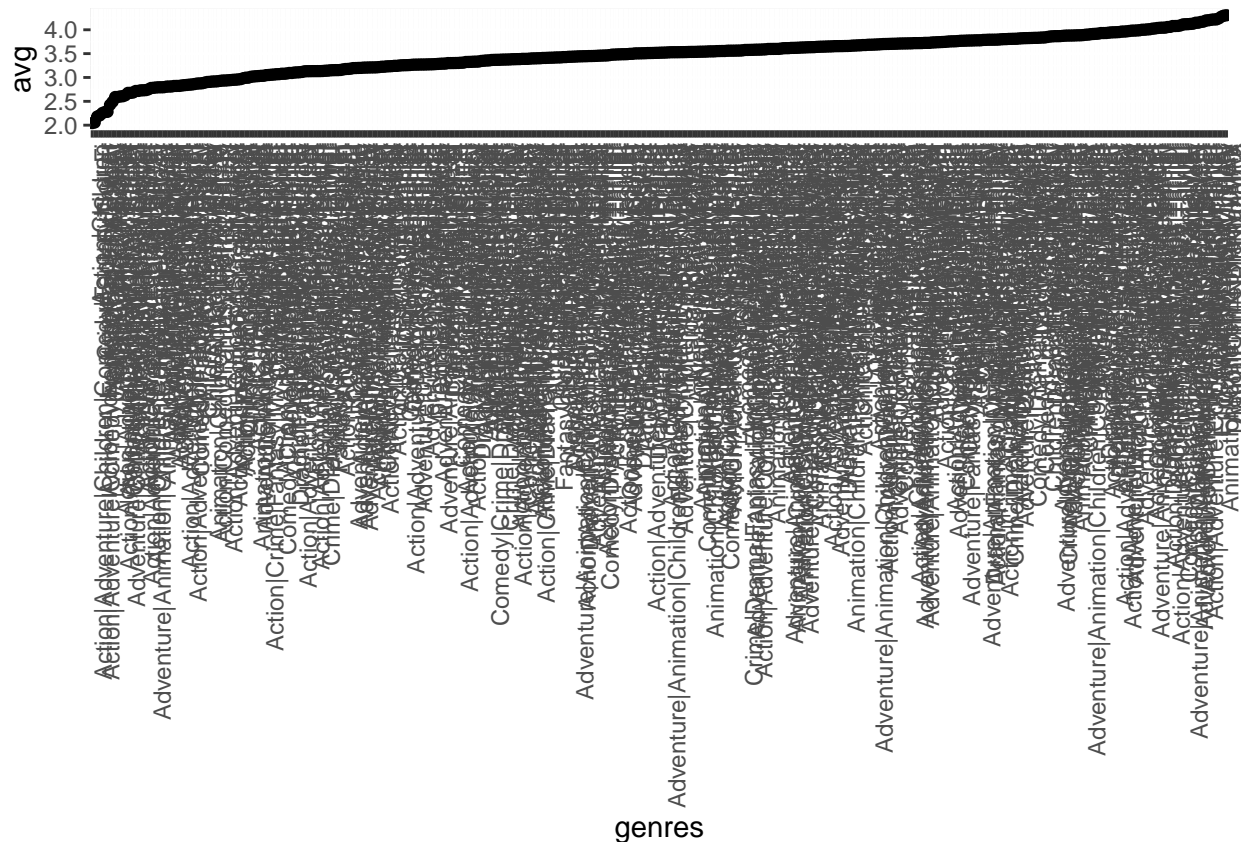


Figure 10: Movies genres error bar plots

The plot shows strong evidence of a genre effect.

3 Analysis - Model Building and Evaluation

3.1 Split the edx data into separate training and test sets

We will develop our algorithm using the edx set only.

We will split the edx data into separate training and test sets to design and test our algorithm, namely `train_set` and `test_set`.

3.1.1 Loss function

For a final test of our algorithm, we predict movie ratings in the test set as if they were unknown. RMSE⁷ (residual mean squared error/root mean square error), the typical error loss, will be used to evaluate how close our predictions are to the true values in the validation set.

We define $y_{u,i}$ as the rating for movie i by user u and denote our prediction with $\hat{y}_{u,i}$.

The RMSE is then defined as Equation 2:

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

with N being the number of user/movie combinations and the sum occurring over all these combinations.

Remember that we can interpret the RMSE similarly to a standard deviation: it is the typical error we make when predicting a movie rating. If this number is larger than 1, it means our typical error is larger than one star, which is not good.

Let's write a function that computes the RMSE for vectors of ratings and their corresponding predictors:

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

⁷https://en.wikipedia.org/wiki/Root-mean-square_deviation

3.2 Model 1: A first naive “mean” model

Let’s start by building the simplest possible recommendation system: we predict the same rating for all movies regardless of user. A model that assumes the same rating for all movies and users with all the differences explained by random variation would look like Equation 3:

$$Y_{i,i} = \mu + \epsilon_{u,i} \quad (3)$$

with $\epsilon_{u,i}$ 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_hat <- mean(train_set$rating))  
[1] 3.512482
```

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

```
(naive_rmse <- RMSE(test_set$rating, mu_hat))  
[1] 1.059904
```

Keep in mind that if we plug in any other number, we get a higher RMSE. For example:

```
predictions <- rep(2.5, nrow(test_set))  
RMSE(test_set$rating, predictions)  
[1] 1.465736
```

```
predictions <- rep(3, nrow(test_set))  
RMSE(test_set$rating, predictions)  
[1] 1.177271
```

```
predictions <- rep(4, nrow(test_set))  
RMSE(test_set$rating, predictions)  
[1] 1.166678
```

From looking at the distribution of ratings, we can visualize that this is the standard deviation of that distribution. We get a RMSE of about 1. Our target is $\text{RMSE} < 0.86490$. So we can definitely do better!

3.2.1 Results Table

As we go along, we will be comparing different approaches. Let's start by creating a results table with this naive approach:

Table 10: RMSE Results

Index	Method	RMSE
1	Just the average	1.059904

```
knitr::knit_exit()
```