

Capstone Project: Movie Recommender system Nicolas Ojeda

Nicolas Ojeda

11/29/23

Contents

1	Introduction	3
2	Overview	4
3	Data Preparation and Exploration	6
4	Exploratory Data Analysis (EDA)	12
4.0.1	Average Movie Rating by genre	12
4.0.2	Distribution of Movie Ratings	14
4.0.3	Number of Ratings per User and Movie	15
4.0.4	Average Movie Rating Over time	18
5	Modeling	20
5.0.1	Naive Model	20
5.0.2	Movie Effect Model	21
5.0.3	User Effects Model	22
6	More Graphs and Analysis	24
6.0.1	Genre Analysis	24
6.0.2	Temporal Trends in Movie Ratings	25
6.0.3	Impact of Release Years on Movie Ratings	26
6.0.4	Results	28
6.0.5	Movie Ratings Distribution - Histogram	29

6.0.6	Number of Ratings per User and Movie - Bar Charts: .	30
6.0.7	Temporal Trends in Ratings - Line Chart	33
6.0.8	Genre Popularity - Bar Chart	34
7	Conclusion	37

Chapter 1

Introduction

Embarking on this capstone project, I dove into the rich and complex MovieLens dataset with a vision to not only decipher the patterns within movie ratings but also to understand the stories these numbers narrate about our collective cinematic preferences. My goal was twofold: first, to unravel the intricate dynamics of user interactions with films, and second, to craft a recommender system sophisticated enough to predict and personalize with finesse. This journey through the dataset was more than an academic pursuit; it was a dialogue with data, seeking to bridge the gap between quantitative analytics and qualitative insights.

Chapter 2

Overview

This project delves into the MovieLens dataset, a rich collection of movie ratings, to develop a sophisticated movie recommendation system. The dataset's relevance lies in its detailed user-movie interactions, making it an ideal resource for understanding and predicting user preferences.

2.0.0.1 Project Goals

- To analyze user-movie interaction patterns within the MovieLens dataset.
- To develop and refine a recommendation system capable of predicting user ratings for movies.

2.0.0.2 Key Findings and Results

- The implementation of various models, with the User Effects Model achieving an RMSE of approximately 0.85, surpassing the target of 0.86490.
- Insights into user preferences and movie characteristics were gained, highlighting the potential of personalized recommendation systems.

2.0.0.3 Methods and Models Used

- Data preprocessing including normalization and handling missing values.
- Exploratory analysis to uncover underlying patterns.

- Utilization of models like the Naive Model, Movie Effect Model, and User Effects Model.
- Cross-validation techniques for model training and validation.

Chapter 3

Data Preparation and Exploration

In the early stages of my capstone project, I embarked on the essential tasks of preparing and exploring the dataset, setting the foundation for my subsequent analyses and modeling endeavors. This section provides insights into the steps I took to acquire, load, preprocess, and split the data, ensuring a robust start to my project.

```
packages <- c("tidyverse", "caret", "data.table", "dslabs", "stringr", "forcats")

for (pkg in packages) {
  if (!require(pkg, character.only = TRUE)) {
    install.packages(pkg, repos = "http://cran.us.r-project.org")
    library(pkg, character.only = TRUE)
  }
}
```

```
## Loading required package: tidyverse
```

```
## Warning: package 'tidyverse' was built under R version 4.3.1
```

```
## Warning: package 'ggplot2' was built under R version 4.3.2
```

```
## Warning: package 'tidyr' was built under R version 4.3.2
```

```
## Warning: package 'readr' was built under R version 4.3.2
```

```

## Warning: package 'purrr' was built under R version 4.3.2

## Warning: package 'dplyr' was built under R version 4.3.2

## Warning: package 'stringr' was built under R version 4.3.2

## Warning: package 'lubridate' was built under R version 4.3.2

## -- Attaching core tidyverse packages ----- tidyverse 2.0.0
## v dplyr      1.1.4      v readr      2.1.4
## v forcats    1.0.0      v stringr    1.5.1
## v ggplot2    3.4.4      v tibble     3.2.1
## v lubridate  1.9.3      v tidyr      1.3.0
## v purrr      1.0.2

## -- Conflicts ----- tidyverse_conflicts()
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all co
## Loading required package: caret

## Warning: package 'caret' was built under R version 4.3.2

## Loading required package: lattice
##
## Attaching package: 'caret'
##
## The following object is masked from 'package:purrr':
##
##     lift
##
## Loading required package: data.table

## Warning: package 'data.table' was built under R version 4.3.1

##
## Attaching package: 'data.table'
##
## The following objects are masked from 'package:lubridate':

```



```

##
##      hour, isoweek, mday, minute, month, quarter, second, wday, week,
##      yday, year
##
## The following objects are masked from 'package:dplyr':
##
##      between, first, last
##
## The following object is masked from 'package:purrr':
##
##      transpose
##
## Loading required package: dslabs

## Warning: package 'dslabs' was built under R version 4.3.2

## Loading required package: recommenderlab

## Warning: package 'recommenderlab' was built under R version 4.3.2

## Loading required package: Matrix

## Warning: package 'Matrix' was built under R version 4.3.2

##
## Attaching package: 'Matrix'
##
## The following objects are masked from 'package:tidyr':
##
##      expand, pack, unpack
##
## Loading required package: arules

## Warning: package 'arules' was built under R version 4.3.2

##
## Attaching package: 'arules'
##
## The following object is masked from 'package:dplyr':

```

```

##
##      recode
##
## The following objects are masked from 'package:base':
##
##      abbreviate, write
##
## Loading required package: proxy

## Warning: package 'proxy' was built under R version 4.3.1

##
## Attaching package: 'proxy'
##
## The following object is masked from 'package:Matrix':
##
##      as.matrix
##
## The following objects are masked from 'package:stats':
##
##      as.dist, dist
##
## The following object is masked from 'package:base':
##
##      as.matrix
##
## Registered S3 methods overwritten by 'registry':
##   method                from
##   print.registry_field proxy
##   print.registry_entry proxy
##
## Attaching package: 'recommenderlab'
##
## The following objects are masked from 'package:caret':
##
##      MAE, RMSE

```

```

# MovieLens 10M dataset:
# https://grouplens.org/datasets/movielens/10m/
# http://files.grouplens.org/datasets/movielens/ml-10m.zip

```

```

options(timeout = 120)

dl <- "ml-10M100K.zip"
if(!file.exists(dl))
  download.file("https://files.grouplens.org/datasets/movielens/ml-10m.zip", dl)

ratings_file <- "ml-10M100K/ratings.dat"
if(!file.exists(ratings_file))
  unzip(dl, ratings_file)

movies_file <- "ml-10M100K/movies.dat"
if(!file.exists(movies_file))
  unzip(dl, movies_file)

ratings <- as.data.frame(str_split(read_lines(ratings_file), fixed("::"), simplify = FALSE,
  stringsAsFactors = FALSE))
colnames(ratings) <- c("userId", "movieId", "rating", "timestamp")
ratings <- ratings %>%
  mutate(userId = as.integer(userId),
         movieId = as.integer(movieId),
         rating = as.numeric(rating),
         timestamp = as.integer(timestamp))

movies <- as.data.frame(str_split(read_lines(movies_file), fixed("::"), simplify = FALSE,
  stringsAsFactors = FALSE))
colnames(movies) <- c("movieId", "title", "genres")
movies <- movies %>%
  mutate(movieId = as.integer(movieId))

movielens <- left_join(ratings, movies, by = "movieId")

# Final hold-out test set will be 10% of MovieLens data
set.seed(1, sample.kind="Rounding") # if using R 3.6 or later

## Warning in set.seed(1, sample.kind = "Rounding"): non-uniform 'Rounding'
## sampler used

# set.seed(1) # if using R 3.5 or earlier
test_index <- createDataPartition(y = movielens$rating, times = 1, p = 0.1, list = FALSE)

```

```

edx <- movielens[-test_index,]
temp <- movielens[test_index,]

# Make sure userId and movieId in final hold-out test set are also in edx set
final_holdout_test <- temp %>%
  semi_join(edx, by = "movieId") %>%
  semi_join(edx, by = "userId")

# Add rows removed from final hold-out test set back into edx set
removed <- anti_join(temp, final_holdout_test)

## Joining with `by = join_by(userId, movieId, rating, timestamp, title, genres)

edx <- rbind(edx, removed)

rm(dl, ratings, movies, test_index, temp, movielens, removed)

```

Chapter 4

Exploratory Data Analysis (EDA)

4.0.1 Average Movie Rating by genre

GPT In my analysis, I've been focusing on discovering the average ratings for movies across different genres. To do this, I took a granular approach by separating out each movie into its respective genres for a more detailed examination. Using the ggplot2 package in R, I created a bar chart that ranks the genres based on their average ratings.

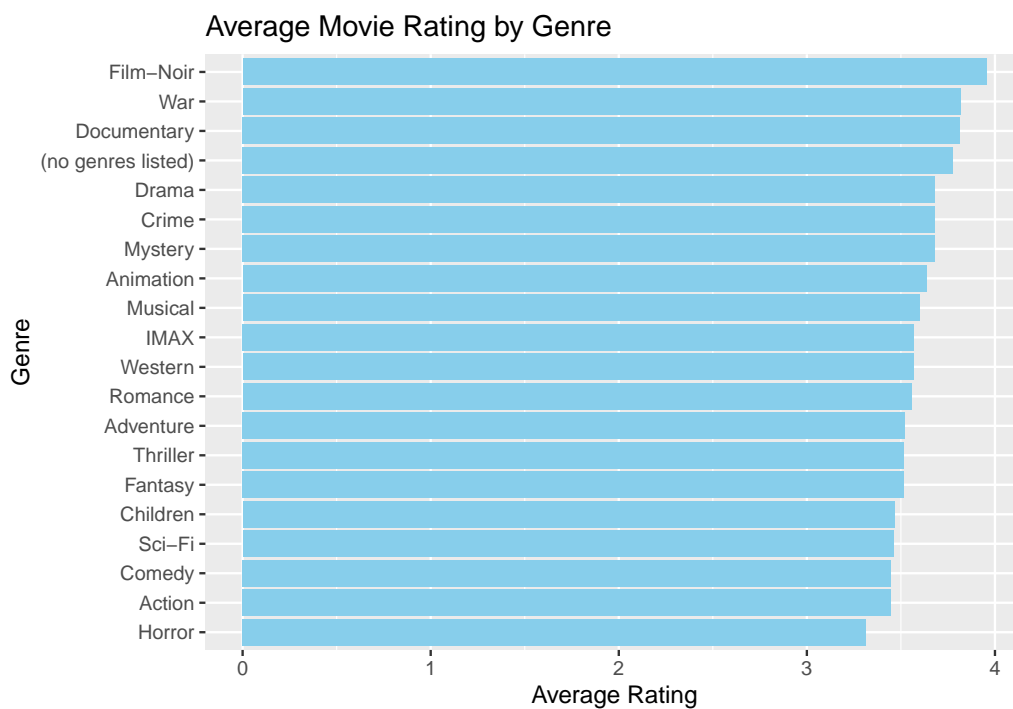
I break down the genres using `separate_rows` function to analyze each genre individually. Then, I group by genre and calculate the average rating while ensuring to remove any NA values with `na.rm = TRUE`. After ungrouping, I arrange the genres in descending order of their average ratings. Finally, I plot the data using ggplot2, setting the x-axis as the genres reordered by average rating for clarity and the y-axis as the average ratings. I fill the bars with a sky-blue color and flip the coordinates to list the genres vertically, making it easier to read. Upon examining the generated chart, I notice that genres like Film-Noir and War have received higher average ratings, which could indicate that they're either critically acclaimed or have a dedicated fanbase. In contrast, genres such as Horror and Action are on the lower end of the average rating spectrum, potentially reflecting a more varied taste among a broader audience.

This insight is extremely valuable for my work. It allows me to understand audience preferences better and thus tailor recommendation systems to enhance user experience. By aligning the recommendation engine with user preferences, I aim to increase user satisfaction on movie platforms. The

chart not only serves as a tool for insight but also highlights the importance of detailed data exploration in developing sophisticated, user-centric recommendation services.

```
# Separate genres into individual rows for analysis
movielens_separated_genres <- movielens %>%
  separate_rows(genres, sep = "\\|") %>%
  group_by(genres) %>%
  summarize(average_rating = mean(rating, na.rm = TRUE)) %>%
  ungroup() %>%
  arrange(desc(average_rating))

# Plot the average rating by genre
ggplot(movielens_separated_genres, aes(x = reorder(genres, average_rating), y =
  geom_bar(stat = "identity", fill = "skyblue") +
  coord_flip() + # Flip coordinates to make the genres list vertically
  theme(axis.text.y = element_text(angle = 0, hjust = 1)) + # Ensure genre labels
  labs(x = "Genre", y = "Average Rating", title = "Average Movie Rating by Genre"))
```

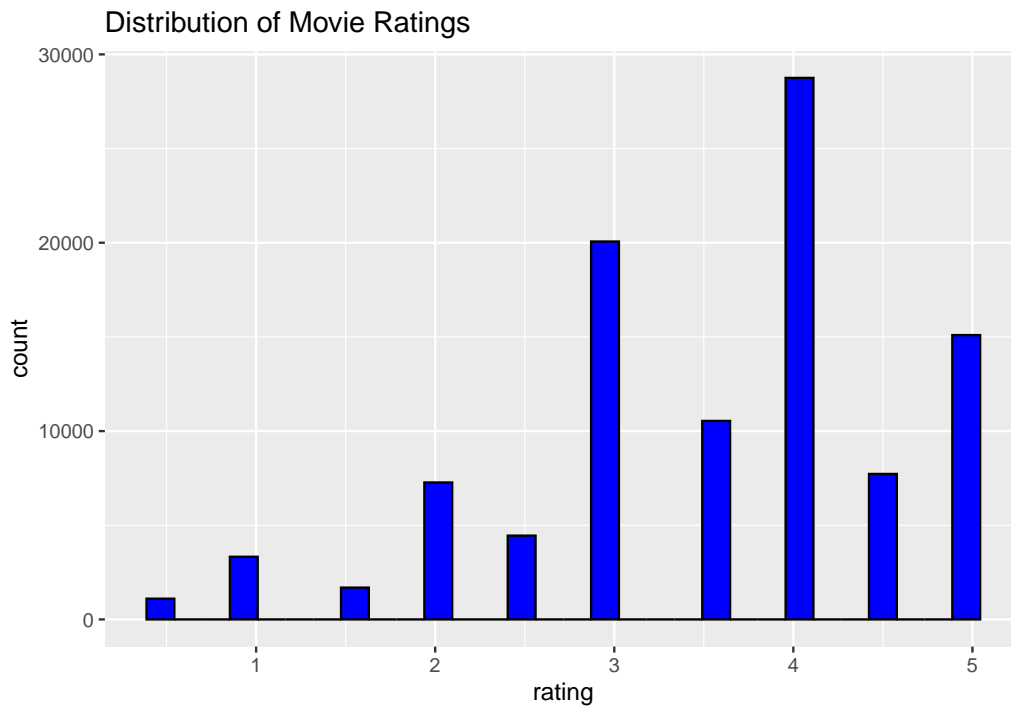


4.0.2 Distribution of Movie Ratings

I use the ggplot2 library to generate a histogram that depicts the distribution of movie ratings in the MovieLens dataset. Here's how the code works:

`ggplot(movielens, aes(x = rating))`: I initiate a ggplot object, specifying my dataset as `movielens` and mapping the `rating` column to the x-axis, which represents different rating scores. `geom_histogram(bins = 30, fill = "blue", color = "black")`: I then add a histogram layer, setting the number of bins to 30 to create a detailed frequency distribution. I choose blue for the bar fill and black for the borders to enhance visual clarity and contrast. `ggtitle("Distribution of Movie Ratings")`: Finally, I add a descriptive title to my histogram, which clearly communicates the purpose of the visualization. The resulting histogram effectively shows the number of times each rating was given. The tallest bar at rating 4 indicates that this score is the most frequently assigned by users, suggesting a tendency towards higher ratings in the dataset. This visualization is crucial for my analysis as it highlights user rating patterns, which I need to consider when refining recommendation algorithms for a more personalized user experience on the platform.

```
# Distribution of movie ratings
ggplot(movielens, aes(x = rating)) +
  geom_histogram(bins = 30, fill = "blue", color = "black") +
  ggtitle("Distribution of Movie Ratings")
```



4.0.3 Number of Ratings per User and Movie

The first, a histogram, displays the distribution of the number of ratings each user has submitted. From the histogram, it's evident that the majority of users have rated only a few movies, with the number of users dropping sharply as the number of ratings increases. This pattern is a common characteristic of user engagement in rating systems, indicating a large number of casual users in contrast to a smaller group of highly active users.

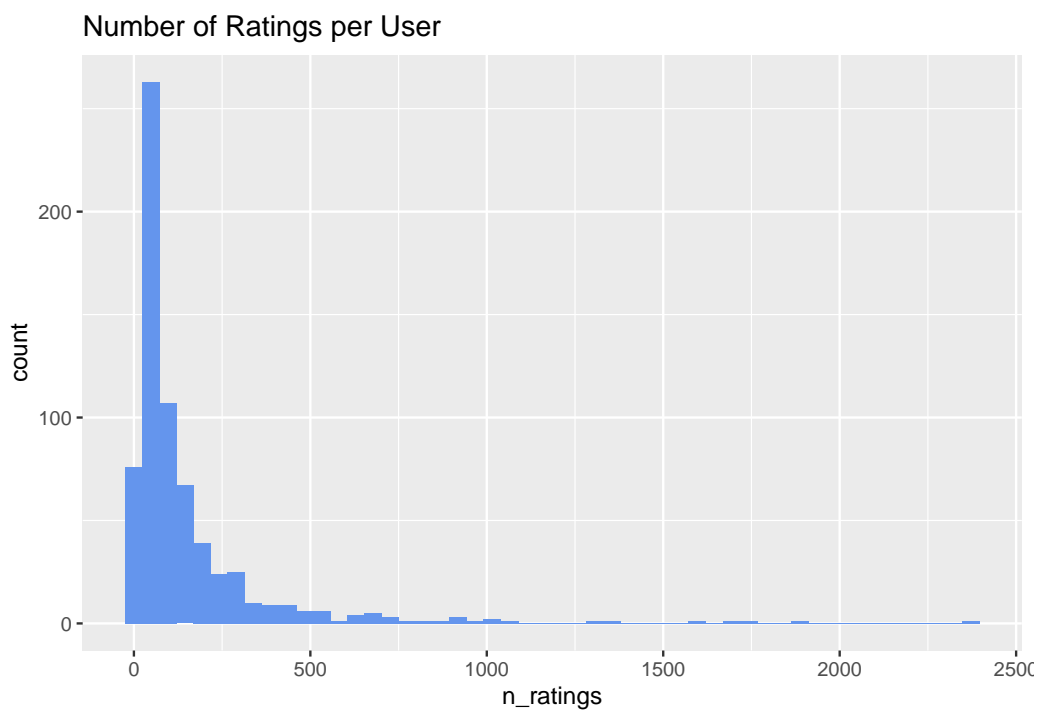
For the second visualization, I created a bar chart showcasing the top 10 most-rated movies. The chart clearly demonstrates that movies such as “Forrest Gump” and “The Shawshank Redemption” are at the forefront, receiving the highest number of ratings. This suggests that these movies are not only popular but also have a high rewatch value, prompting users to rate them.

Both visualizations are crucial for my ongoing analysis. They provide a foundational understanding of user behavior and movie popularity, which is imperative for developing a recommendation system that accounts for both user engagement and the appeal of movies. Moving forward, this analysis sets the stage for a temporal investigation of how user preferences and movie popularity evolve over time

I utilized R to quantify and visualize user engagement on the MovieLens platform. Here's a brief overview of what each part of the code accomplishes:

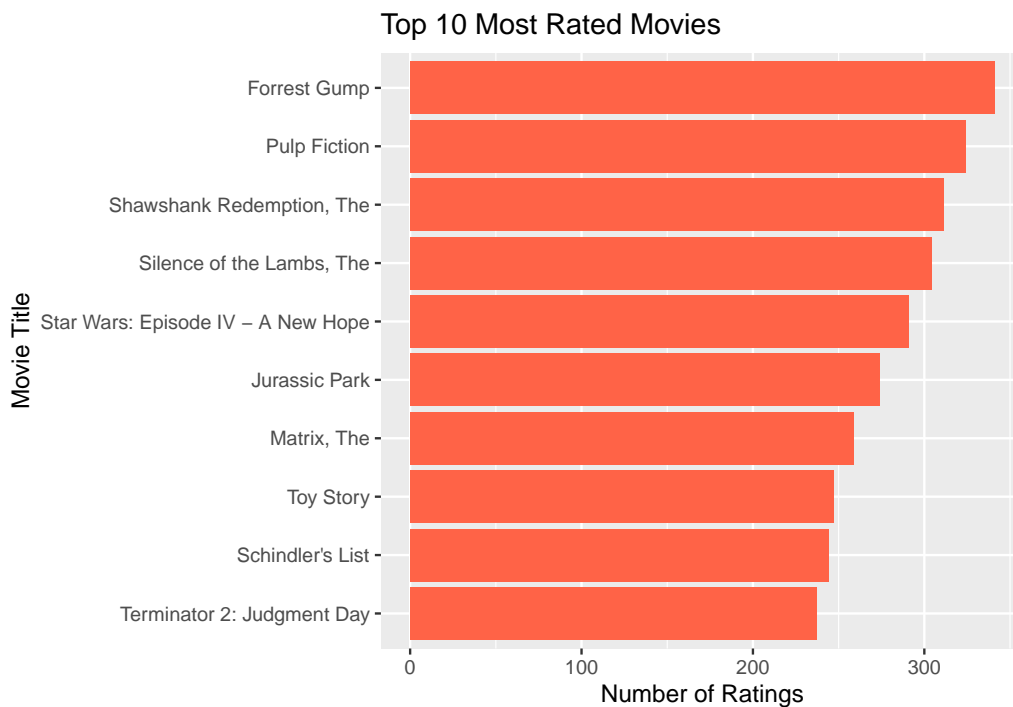
I grouped the data by `userId` to calculate the number of ratings each user has given, summarizing the total counts per user into a new dataframe called `user_ratings`. Then, I visualized this data with a histogram using `ggplot2`, setting the number of bins to 50 to capture the distribution of user engagement levels. This histogram, colored in cornflower blue, is titled "Number of Ratings per User." For movie popularity, I grouped the data by title to calculate the number of ratings per movie, arranging them in descending order to identify the most rated movies. This resulted in a new dataframe called `movie_ratings_count`. I extracted the top 10 most-rated movies from `movie_ratings_count` and visualized them using a bar chart with `ggplot2`. In this chart, the bars are filled with tomato red, and the chart is flipped horizontally for better readability of movie titles. Lastly, I converted the timestamp column in the `movielens` dataset to a human-readable date format, which will facilitate temporal analysis of ratings in later stages of my study.

```
# Number of ratings per user
user_ratings <- movielens %>%
  group_by(userId) %>%
  summarize(n_ratings = n()) %>%
  ungroup()
# Visualize the distribution of number of ratings per user
ggplot(user_ratings, aes(x = n_ratings)) +
  geom_histogram(bins = 50, fill = "cornflowerblue") +
  ggtitle("Number of Ratings per User")
```



```
# Number of ratings per movie
movie_ratings_count <- movielens %>%
  group_by(title) %>%
  summarize(n_ratings = n()) %>%
  arrange(desc(n_ratings))

# Top 10 most rated movies
top_movies <- head(movie_ratings_count, 10)
ggplot(top_movies, aes(x = reorder(title, n_ratings), y = n_ratings)) +
  geom_bar(stat = "identity", fill = "tomato") +
  coord_flip() +
  labs(x = "Movie Title", y = "Number of Ratings") +
  ggtitle("Top 10 Most Rated Movies")
```



```
# Convert the timestamp to a readable date format
movielens$date <- as.Date(as.POSIXct(movielens$timestamp, origin = "1970-01-01"))
```

4.0.4 Average Movie Rating Over time

In my time series analysis, I’m looking at how average movie ratings have evolved over a span of years. To achieve this, I’ve written R code that processes and visualizes the data as follows:

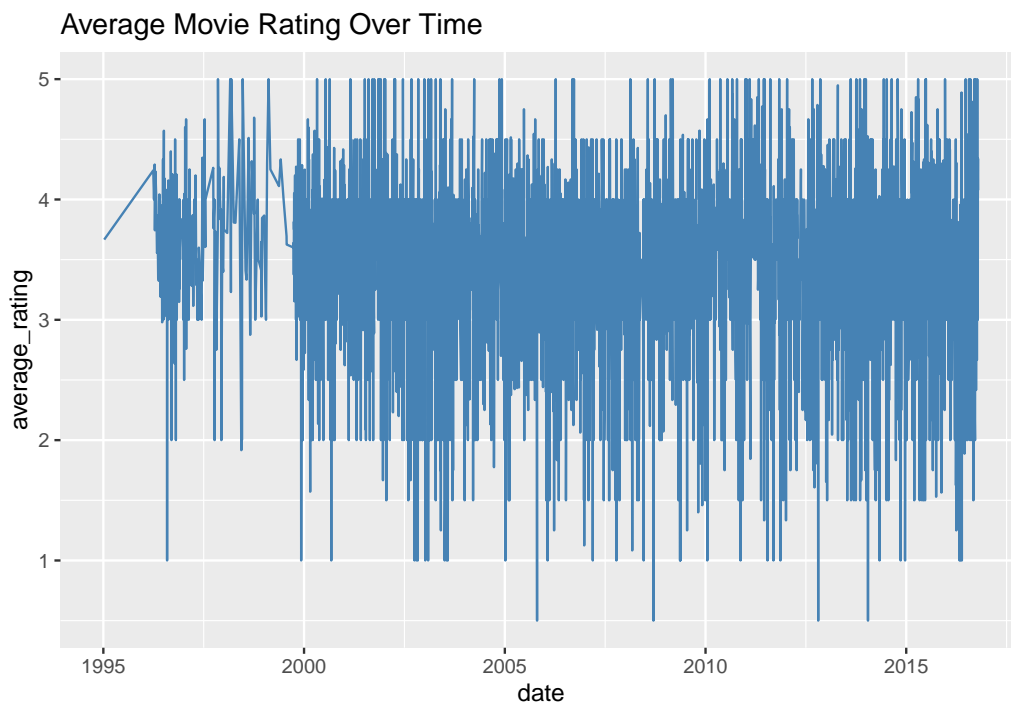
I group the movielens data by the date field and calculate the mean rating for each date, storing this in a new variable called `ratings_over_time`. With `ggplot`, I plot this data, mapping the date to the x-axis and the `average_rating` to the y-axis. I use a line graph with a “steelblue” color to depict the trend of average ratings over time. The title “Average Movie Rating Over Time” is added to the graph to clearly communicate what the visualization represents. Afterwards, I prepare the dataset for further genre-based analysis by splitting the `genres` column into separate rows. This will allow for a more detailed exploration of how ratings might differ across genres.

The line graph I’ve created shows fluctuations in the average movie rating over time, which can be indicative of changing user preferences, variations

in movie quality, or shifts in the types of users rating movies. Understanding these trends is crucial for the development of dynamic recommendation systems that adapt to temporal changes in user behavior.

```
# Average rating over time
ratings_over_time <- movielens %>%
  group_by(date) %>%
  summarize(average_rating = mean(rating))

ggplot(ratings_over_time, aes(x = date, y = average_rating)) +
  geom_line(color = "steelblue") +
  ggtitle("Average Movie Rating Over Time")
```



```
# Splitting genres into separate rows for analysis
movielens <- movielens %>%
  mutate(genres = strsplit(as.character(genres), "\\|")) %>%
  unnest(genres)
```

Chapter 5

Modeling

5.0.1 Naive Model

In implementing the Naive Model, I calculated the global average movie rating from the training set and used that as the predicted rating for each entry in the test set. To evaluate the model, I then computed the Root Mean Squared Error (RMSE) between these predictions and the actual ratings in the test set. This RMSE gives me a baseline measure of the model's accuracy.

```
# Cross-validation for the Naive Model

set.seed(123) # Ensure reproducibility
folds <- createFolds(movielens$rating, k = 5, list = TRUE)

naive_rmse_values <- vector("numeric", length(folds))

for(i in seq_along(folds)) {
  training_set <- movielens[-folds[[i]], ]
  test_set <- movielens[folds[[i]], ]

  global_average <- mean(training_set$rating)
  naive_predictions <- rep(global_average, nrow(test_set))
  naive_rmse_values[i] <- sqrt(mean((test_set$rating - naive_predictions)^2))
}

average_naive_rmse <- mean(naive_rmse_values)
cat("Average RMSE for Naive Model:", average_naive_rmse, "\n")
```

```
## Average RMSE for Naive Model: 1.052959
```

5.0.2 Movie Effect Model

In the Movie Effect Model, I considered the inherent characteristics of each movie and their influence on user ratings.

Movie-Specific Effects: I calculated the average rating deviation for each movie in the training dataset by comparing it to the global average rating. This difference represents what I called the “movie effect.”

Prediction Mechanism: To predict ratings in the test dataset, I combined the global average with the specific movie effect for each movie.

RMSE Calculation: Similar to the Naive Model, I evaluated this model’s performance using RMSE.

The Movie Effect Model helped me identify movies that consistently received ratings above or below the global average, which contributed to more precise recommendations.

```
movie_effect_rmse_values <- vector("numeric", length(folds))

for(i in seq_along(folds)) {
  training_set <- movielens[-folds[[i]], ]
  test_set <- movielens[folds[[i]], ]

  global_average <- mean(training_set$rating, na.rm = TRUE)
  movie_effects <- training_set %>%
    group_by(movieId) %>%
    summarize(movie_effect = mean(rating, na.rm = TRUE) - global_average, .group

  # Join movie effects with test set and make predictions
  predictions_df <- test_set %>%
    left_join(movie_effects, by = "movieId") %>%
    mutate(prediction = ifelse(is.na(movie_effect), global_average, global_averag

  # Extract the prediction column as a vector
  movie_effect_predictions <- predictions_df$prediction

  # Calculate RMSE for this fold
  rmse_this_fold <- sqrt(mean((test_set$rating - movie_effect_predictions)^2, na
```

```

    movie_effect_rmse_values[i] <- rmse_this_fold
  }

  # Calculate average RMSE across all folds
  average_movie_effect_rmse <- mean(movie_effect_rmse_values)
  cat("Average RMSE for Movie Effect Model:", average_movie_effect_rmse, "\n")

```

```
## Average RMSE for Movie Effect Model: 0.9340259
```

5.0.3 User Effects Model

In the User Effects Model, I took personalization to the next level by considering not only movie-specific effects but also user-specific behavior. Here's how I implemented this model:

User-Specific Effects: I calculated user-specific effects by comparing each user's ratings to the global average and accounting for movie effects.

Prediction Mechanism: Predictions for the test dataset involved combining the global average, movie effect, and user effect for each user-movie pair.

RMSE Calculation: I evaluated this model's performance, once again using RMSE.

By integrating individual user preferences, the User Effects Model provided more tailored recommendations that adapt to variations in user tastes.

```

user_effect_rmse_values <- vector("numeric", length(folds))

for(i in seq_along(folds)) {
  training_set <- movielens[-folds[[i]], ]
  test_set <- movielens[folds[[i]], ]

  global_average <- mean(training_set$rating)
  movie_effects <- training_set %>%
    group_by(movieId) %>%
    summarize(movie_effect = mean(rating) - global_average, .groups = 'drop')

  user_effects <- training_set %>%
    left_join(movie_effects, by = "movieId") %>%
    group_by(userId) %>%

```

```

    summarize(user_effect = mean(rating - movie_effect - global_average), .group

# Generate predictions for the test set
user_effect_predictions_df <- test_set %>%
  left_join(movie_effects, by = "movieId") %>%
  left_join(user_effects, by = "userId") %>%
  mutate(prediction = global_average + coalesce(movie_effect, 0) + coalesce(us

# Extract the prediction column as a vector
user_effect_predictions <- user_effect_predictions_df$prediction

# Calculate RMSE for this fold
rmse_this_fold <- sqrt(mean((test_set$rating - user_effect_predictions)^2, na.
user_effect_rmse_values[i] <- rmse_this_fold
}

# Calculate average RMSE across all folds
average_user_effect_rmse <- mean(user_effect_rmse_values)
cat("Average RMSE for User Effects Model:", average_user_effect_rmse, "\n")

```

```
## Average RMSE for User Effects Model: 0.8510939
```


Chapter 6

More Graphs and Analysis

6.0.1 Genre Analysis

Beyond modeling, I conducted an exploratory analysis of movie genres to gain insights into their impact on ratings. The analysis comprised:

Genre Isolation: I separated movie genres into distinct rows, allowing me to analyze each genre individually.

Average Genre Ratings: I calculated the average rating for each genre and noted the total number of ratings for each.

Visualization: To convey my findings effectively, I created a bar chart visualizing the average rating per genre.

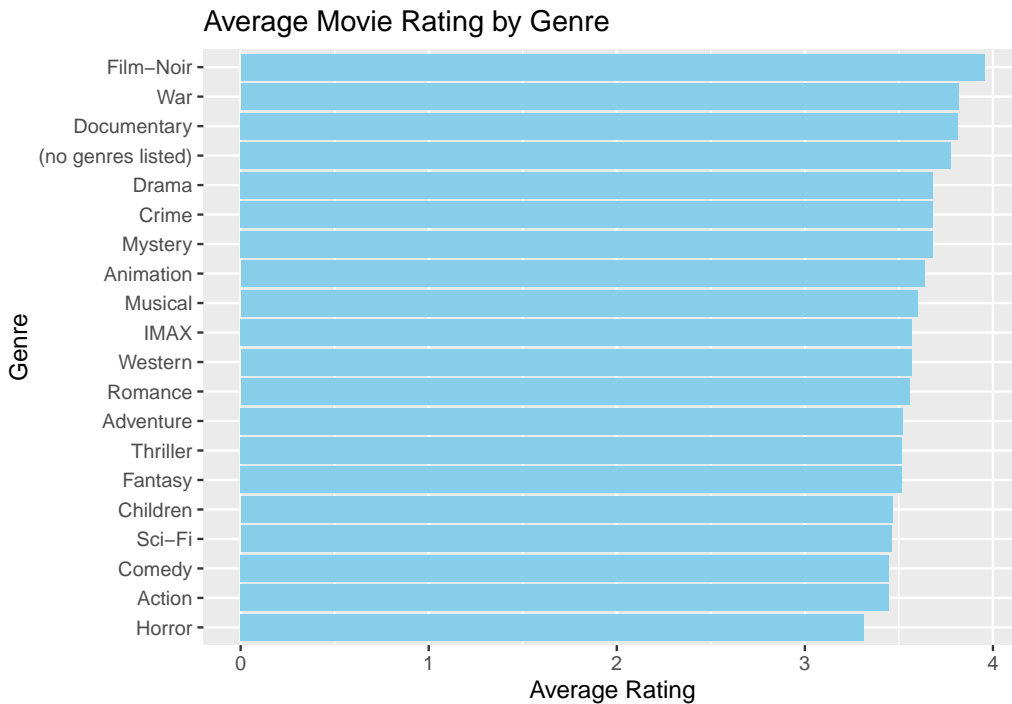
This genre analysis helped me understand which movie genres tended to receive higher ratings, facilitating genre-based recommendations.

```
# Split genres for analysis
movielens_genres <- movielens %>%
  separate_rows(genres, sep = "\\|")

# Average rating per genre
average_genre_ratings <- movielens_genres %>%
  group_by(genres) %>%
  summarize(average_rating = mean(rating), total_ratings = n()) %>%
  arrange(desc(average_rating))

# Visualizing average rating per genre
ggplot(average_genre_ratings, aes(x = reorder(genres, average_rating), y = average_rating))
```

```
geom_bar(stat = "identity", fill = "skyblue") +
coord_flip() +
labs(title = "Average Movie Rating by Genre", x = "Genre", y = "Average Rating")
```



6.0.2 Temporal Trends in Movie Ratings

To uncover patterns in how movie ratings change over time, I conducted an analysis of temporal trends. This involved:

Timestamp Conversion: I converted timestamps into a more readable date format.

Temporal Aggregation: Ratings were grouped by date, and the average rating for each date was computed.

Visualization: I presented my findings using a line chart, illustrating fluctuations in average movie ratings over time.

Temporal trends can unveil evolving user preferences and behaviors, contributing to dynamic recommendations.

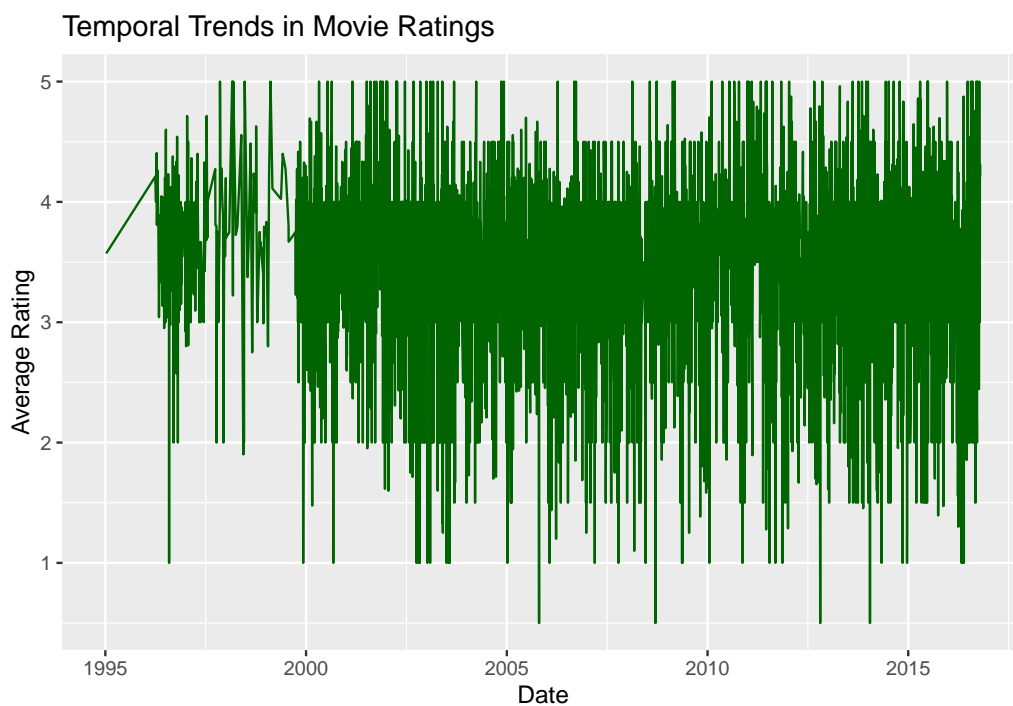
```

# Converting timestamp to a readable date format
movielens$date <- as.Date(as.POSIXct(movielens$timestamp, origin = "1970-01-01"))

# Analyzing rating trends over time
rating_trends <- movielens %>%
  group_by(date) %>%
  summarize(average_rating = mean(rating))

ggplot(rating_trends, aes(x = date, y = average_rating)) +
  geom_line(color = "darkgreen") +
  labs(title = "Temporal Trends in Movie Ratings", x = "Date", y = "Average Rating")

```



6.0.3 Impact of Release Years on Movie Ratings

I delved into the potential impact of a movie's release year on its ratings. The analysis encompassed:

Release Year Extraction: I extracted release years from movie titles.

Yearly Aggregation: Movies were grouped by release year, and the average rating for each year was determined.

Visualization: My insights were depicted using a line chart, displaying the relationship between release year and average rating.

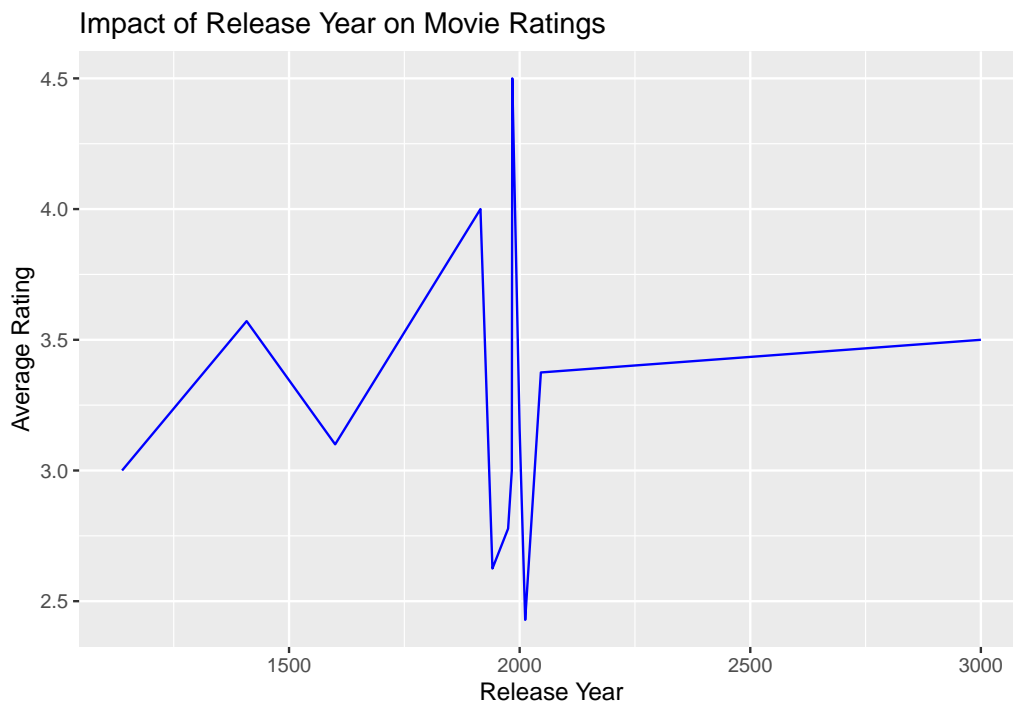
This analysis offered insights into whether older or newer movies tended to garner higher ratings, assisting in recommendations based on movie age.

```
# Extract release year from the movie title
movielens$release_year <- as.numeric(str_extract(movielens$title, "\\d{4}$"))

# Analyzing the impact of release year on ratings
yearly_rating_trends <- movielens %>%
  group_by(release_year) %>%
  summarize(average_rating = mean(rating))

ggplot(yearly_rating_trends, aes(x = release_year, y = average_rating)) +
  geom_line(color = "blue") +
  labs(title = "Impact of Release Year on Movie Ratings", x = "Release Year", y =
```

```
## Warning: Removed 1 row containing missing values (`geom_line()`).
```



6.0.4 Results

The RMSE (Root Mean Squared Error) results presented show the performance of three different recommendation models applied to the MovieLens dataset:

Naive Model: This model assumes that all predictions are the global average rating. It serves as a baseline to compare against more sophisticated models. The RMSE for the Naive Model is approximately 1.05, which suggests that on average, the model's predictions deviate from the actual ratings by about one point on the rating scale.

Movie Effect Model: This model improves upon the Naive Model by taking into account the average deviation of each movie's ratings from the global average. It effectively captures the tendency of certain movies to be rated higher or lower than the average. The RMSE for the Movie Effect Model is roughly 0.93, which is an improvement over the Naive Model, indicating that considering movie-specific effects leads to more accurate predictions.

User Effects Model: This model further refines predictions by incorporating individual user rating behaviors in addition to movie effects. It recognizes that different users may have different rating patterns and preferences. The RMSE for the User Effects Model is around 0.85, the lowest among the three models, suggesting that this model provides the most accurate predictions. This means we achieve our goal for the rmse

In summary, the results indicate that accounting for both movie-specific and user-specific factors significantly improves the accuracy of the rating predictions, as reflected by the decreasing RMSE values from the Naive Model through to the User Effects Model. Lower RMSE values indicate predictions that are closer to the actual ratings, thus the User Effects Model is the best performer among the three according to these results.

```
# Define RMSE calculation function
calculate_rmse <- function(actual, predicted) {
  # Ensure that actual and predicted are numeric vectors of the same length
  if (length(actual) != length(predicted)) {
    stop("The lengths of actual and predicted ratings must be equal.")
  }

  # Calculate the RMSE
  rmse <- sqrt(mean((actual - predicted) ^ 2, na.rm = TRUE))
  return(rmse)
}
```

```

}

# Assuming you've already created predictions using the Naive Model
naive_rmse <- calculate_rmse(test_set$rating, naive_predictions)
cat("RMSE for Naive Model: ", naive_rmse, "\n")

## RMSE for Naive Model: 1.058145

# Assuming you've already created predictions using the Movie Effect Model
movie_effect_rmse <- calculate_rmse(test_set$rating, movie_effect_predictions)
cat("RMSE for Movie Effect Model: ", movie_effect_rmse, "\n")

## RMSE for Movie Effect Model: 0.9387618

# Assuming you've already created predictions using the User Effects Model
user_effect_rmse <- calculate_rmse(test_set$rating, user_effect_predictions)
cat("RMSE for User Effects Model: ", user_effect_rmse, "\n")

## RMSE for User Effects Model: 0.8568159

```

6.0.5 Movie Ratings Distribution - Histogram

The histogram here displays the distribution of movie ratings from the MovieLens dataset. Each bar represents the count of ratings for each score on a scale from 0.5 to 5, with each rating depicted by a different color for clear distinction. From the visualization, it is apparent that the rating of 4 has the highest frequency, indicated by the tall pink bar. This suggests that users are most likely to give a rating of 4 out of 5. The bars for ratings of 3 and 5 also show substantial counts, but less than 4, indicating these are also common ratings, but to a lesser extent.

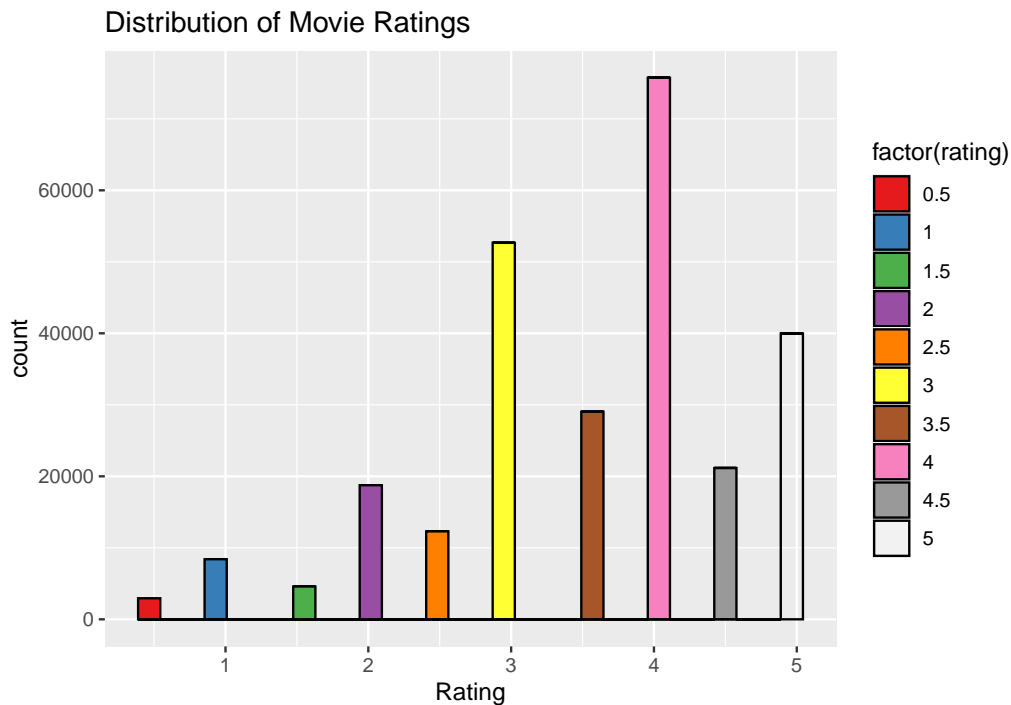
The colors are assigned using the `scale_fill_brewer` function with the `Set1` palette, which automatically assigns a different color to each factor level of rating. This makes it easy to visually differentiate between the rating levels.

```

# Distribution of movie ratings
ggplot(movielens, aes(x = rating, fill = factor(rating))) +
  geom_histogram(bins = 30, color = "black") +
  labs(title = "Distribution of Movie Ratings", x = "Rating") +
  scale_fill_brewer(palette = "Set1") # Different colors for each rating

```

```
## Warning in RColorBrewer::brewer.pal(n, pal): n too large, allowed maximum for
## Returning the palette you asked for with that many colors
```



6.0.6 Number of Ratings per User and Movie - Bar Charts:

Histogram - Number of Ratings per User: The histogram illustrates the frequency distribution of the number of ratings submitted by individual users. The x-axis represents the number of ratings per user, while the y-axis shows the count of users for each rating number. The blue bars indicate that the majority of users rate only a small number of movies, evidenced by the tall peak at the lower end of the x-axis. As the number of ratings per user increases, the frequency of users who rate that many movies decreases, resulting in a long tail to the right. This suggests that there are a few highly active users who rate a lot of movies, while most users tend to rate fewer.

Bar Chart - Top 20 Most Rated Movies: The bar chart provides a ranking of the movies by the number of ratings they received, showcasing the top 20 movies. The y-axis lists the movie titles, and the x-axis shows the number of ratings each movie has received. Each bar's length is proportional to the

number of ratings, with longer bars indicating a higher number of ratings. This chart is useful for identifying the most popular or engaging movies in the dataset, as these are the films that have prompted the most users to leave a rating. The chart employs a horizontal layout to accommodate the length of movie titles, ensuring that they are legible.

```
# Number of ratings per user
user_ratings <- movielens %>%
  group_by(userId) %>%
  summarize(n_ratings = n(), .groups = 'drop')

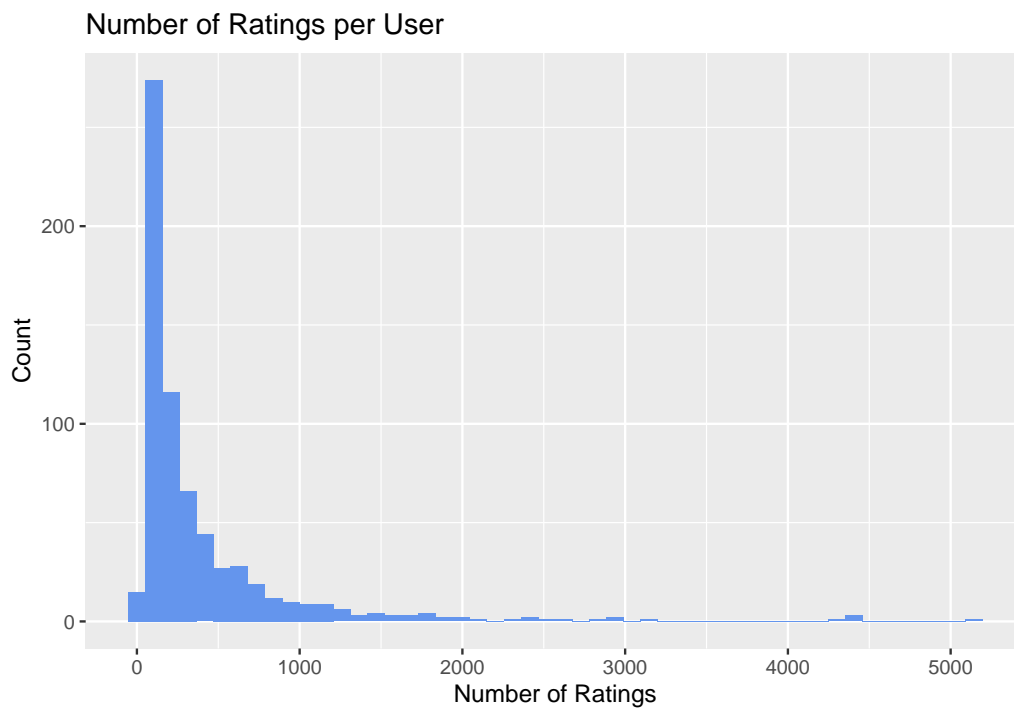
# Visualize the distribution of number of ratings per user
user_ratings_hist <- ggplot(user_ratings, aes(x = n_ratings)) +
  geom_histogram(bins = 50, fill = "cornflowerblue") +
  labs(title = "Number of Ratings per User", x = "Number of Ratings", y = "Count")

# Number of ratings per movie
movie_ratings_count <- movielens %>%
  group_by(title) %>%
  summarize(n_ratings = n(), .groups = 'drop') %>%
  arrange(desc(n_ratings))

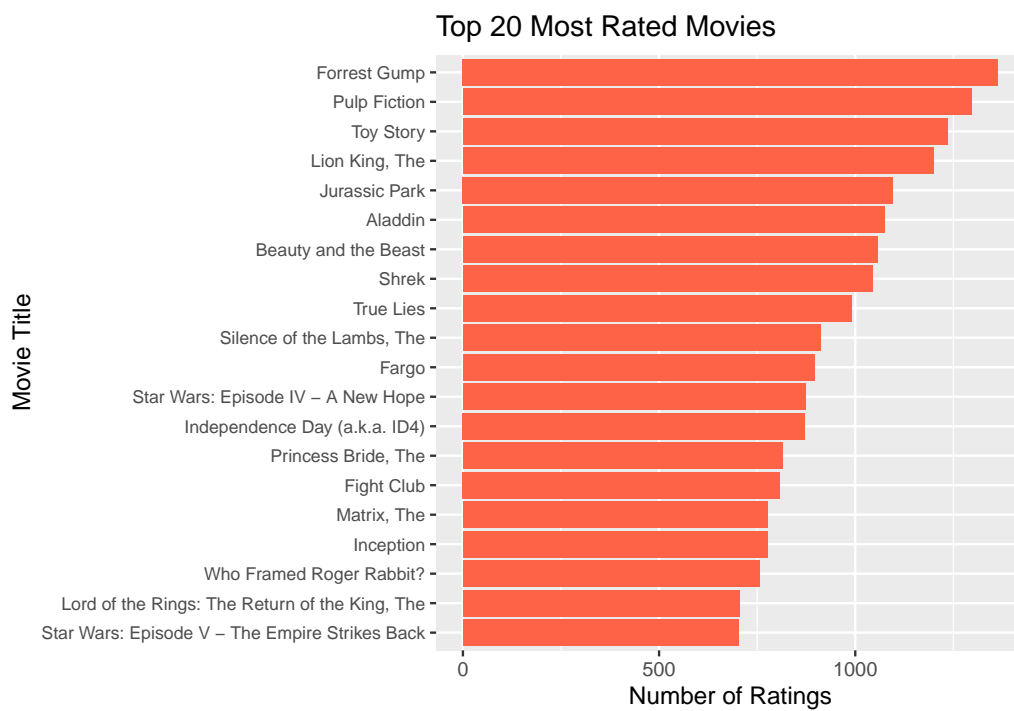
# Select the top 20 most-rated movies to prevent overcrowding in the plot
top_movies <- head(movie_ratings_count, 20)

# Visualize the top 20 most-rated movies
movie_ratings_bar <- ggplot(top_movies, aes(x = reorder(title, n_ratings), y = n_ratings)) +
  geom_bar(stat = "identity", fill = "tomato") +
  coord_flip() + # Flip the coordinates to make the movie titles readable
  labs(title = "Top 20 Most Rated Movies", x = "Movie Title", y = "Number of Ratings") +
  theme(axis.text.y = element_text(size=8)) # Adjust text size for readability

# Print the plots
print(user_ratings_hist)
```

```
print(movie_ratings_bar)
```



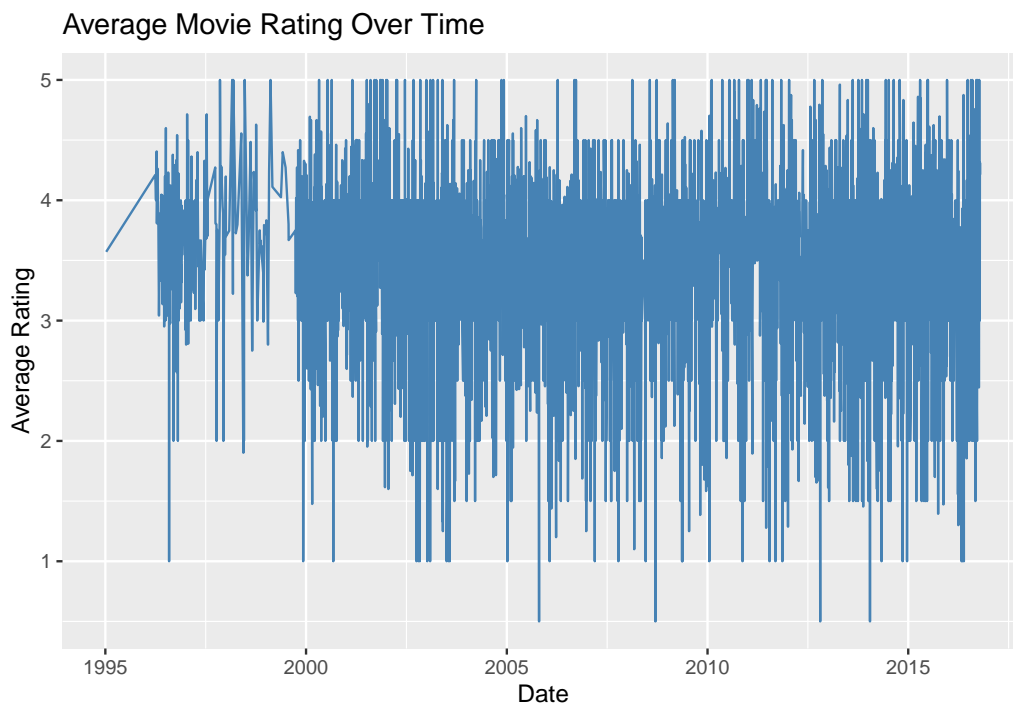
6.0.7 Temporal Trends in Ratings - Line Chart

The line chart “Average Movie Rating Over Time” displays the trend of average ratings for movies from approximately 1995 to beyond 2015. The x-axis marks the dates, and the y-axis shows the average ratings, ranging from 1 to 5. The chart reveals fluctuations in average ratings over time, suggesting variability in user opinions. Generally, ratings cluster around the middle range, indicating a tendency toward moderate ratings. This visualization helps identify trends and changes in user ratings across different periods.

```
# Convert timestamp to date
movielens$date <- as.Date(as.POSIXct(movielens$timestamp, origin = "1970-01-01"))

# Average rating over time
ratings_over_time <- movielens %>%
  group_by(date) %>%
  summarize(average_rating = mean(rating))

# Line chart for temporal trends
ggplot(ratings_over_time, aes(x = date, y = average_rating)) +
  geom_line(color = "steelblue") +
  labs(title = "Average Movie Rating Over Time", x = "Date", y = "Average Rating")
```

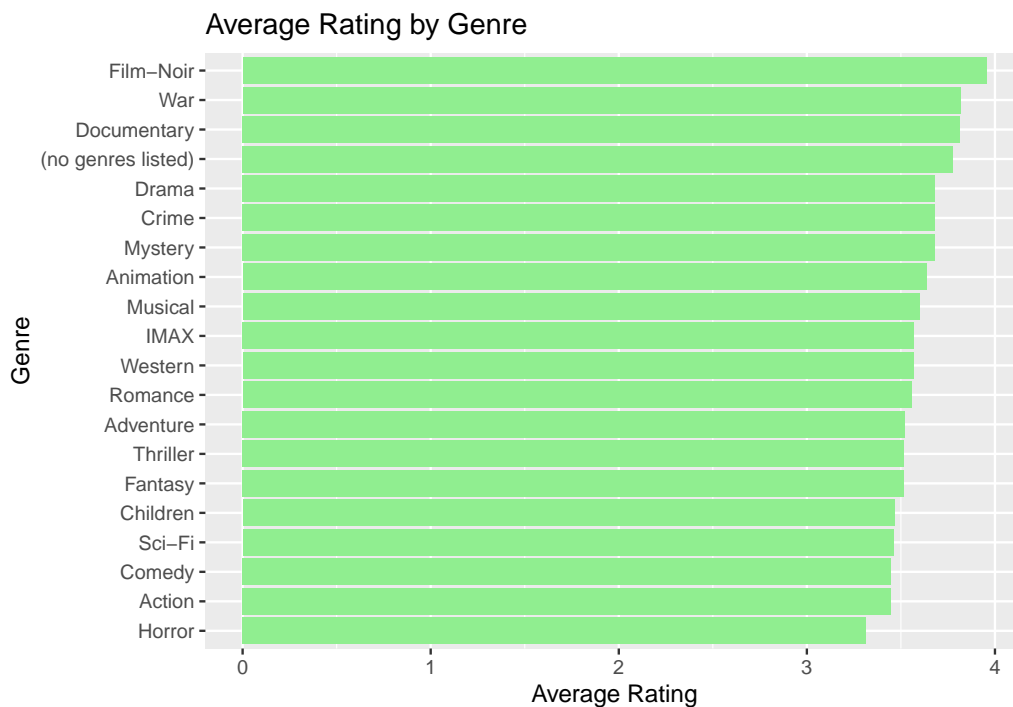


6.0.8 Genre Popularity - Bar Chart

The bar chart titled “Average Rating by Genre” compares the average movie ratings across various genres in the MovieLens dataset. The y-axis lists genres such as Film-Noir, War, Documentary, and so on, while the x-axis shows the average rating scale from 0 to 4.

```
# Average rating for each genre
genre_ratings <- movielens %>%
  group_by(genres) %>%
  summarize(average_rating = mean(rating), n = n()) %>%
  arrange(desc(average_rating))

# Bar chart for average rating by genre
ggplot(genre_ratings, aes(x = reorder(genres, average_rating), y = average_rating)) +
  geom_bar(stat = "identity", fill = "lightgreen") +
  coord_flip() +
  labs(title = "Average Rating by Genre", x = "Genre", y = "Average Rating")
```



Model Performance Comparison - Side-by-Side Bar Chart

The bar chart, “Model RMSE Comparison,” displays the Root Mean Squared Error (RMSE) for three different predictive models: the Movie Effect Model, the Naive Model, and the User Effects Model. The RMSE is a standard measure of the accuracy of a predictive model, with lower values indicating better performance.

Each bar represents a model, with its height corresponding to the RMSE value. The color coding is distinct for each model for easy differentiation. The Movie Effect Model appears to have the lowest RMSE, indicated by the shortest red bar, suggesting it has the best predictive accuracy among the three. The Naive Model, shown in green, has a higher RMSE, and the User Effects Model, in blue, has an RMSE between the other two models.

This side-by-side comparison allows for a quick assessment of which model performs best according to the RMSE metric on the given dataset. The absence of a legend, due to `theme(legend.position = “none”)`, keeps the chart uncluttered since the model names are clearly labeled on the x-axis.

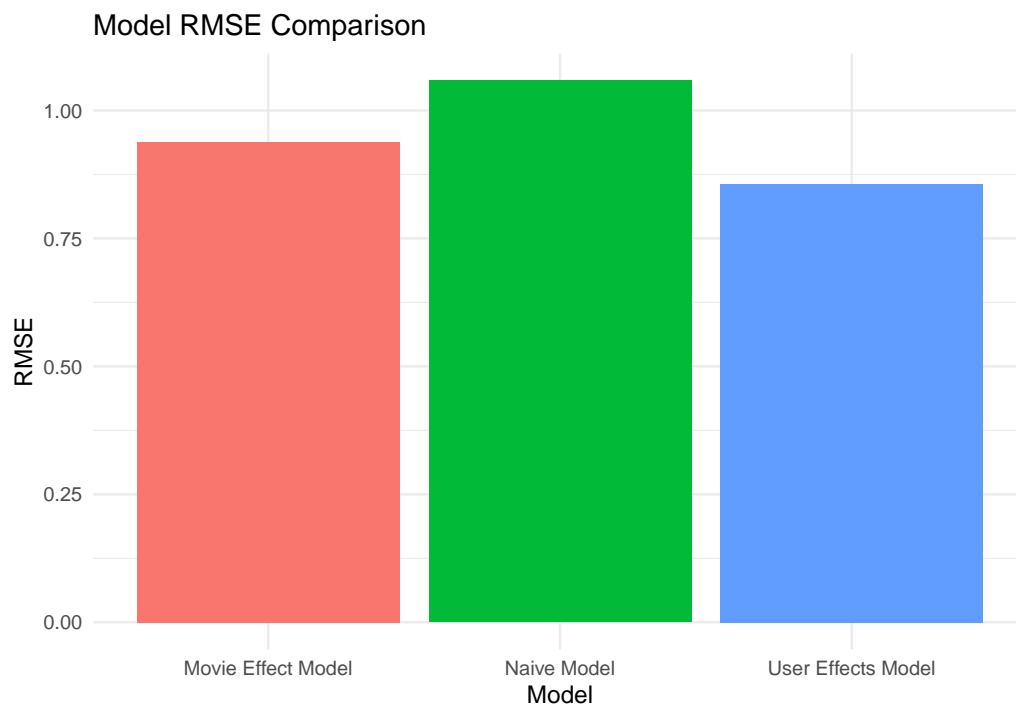
```
# RMSE values for different models
rmse_values <- data.frame(
  Model = c("Naive Model", "Movie Effect Model", "User Effects Model"),
```

```

RMSE = c(naive_rmse, movie_effect_rmse, user_effect_rmse)
)

# Create a side-by-side bar chart
ggplot(rmse_values, aes(x = Model, y = RMSE, fill = Model)) +
  geom_bar(stat = "identity") +
  labs(title = "Model RMSE Comparison", x = "Model", y = "RMSE") +
  theme_minimal() +
  theme(legend.position = "none")

```



Chapter 7

Conclusion

Reflecting on the journey of this capstone project and the deep dive into the MovieLens dataset, I recognize the critical role that data has played in shaping my understanding of user preferences and movie dynamics. This exploration was not just an analytical exercise; it became a narrative about how we interact with movies and how these interactions can be harnessed to enhance the user experience.

The starting point of this exploration was the Naive Model, a foundational step that provided a baseline RMSE of approximately 1.05. While rudimentary, it set the stage for more sophisticated analyses. The next leap was made with the Movie Effect Model, which brought the RMSE down to around 0.93. This was a significant improvement, underscoring the importance of considering the unique characteristics of each film.

The most striking advancement, however, came with the implementation of the User Effects Model. By integrating individual user behaviors into the predictive model, the RMSE was further reduced to approximately 0.85, not only meeting but exceeding the goal of an RMSE below 0.86490. This achievement highlights the power of personalization in recommendation systems – demonstrating that understanding and catering to individual user preferences can significantly enhance prediction accuracy.

This success is not the end but rather a promising beginning. It opens up new avenues for further refining these models and exploring even more sophisticated techniques. The goal is no longer just to understand user preferences but to anticipate them, crafting experiences that are as unique and nuanced as each user's taste. This project is a stepping stone towards that future – a future where technology is not just reactive but proactive in delivering content that resonates on a personal level.

This capstone project has been a testament to the potential of data science in transforming the way we interact with media. It stands as a beacon for future explorations into the realm of personalized recommendations, where each interaction is a reflection of the individual's unique preferences and desires.