# HarvardX: PH125.9x Data Science - EDX Rating Prediction Project

Matias Ezequiel Maurig

2024-11-06

## Contents

# 1 Introduction and Project Objective

## 1.1 Project Overview

This report delves into the development of a personalized movie recommendation system, crafted as part of the MovieLens project for the HarvardX: PH125.9x Data Science Capstone course. The project's aim is to build a prediction model capable of estimating movie ratings based on user preferences, leveraging machine learning techniques and predictive algorithms.

The dataset employed is the 10M version of MovieLens, created by GroupLens Research, known for its comprehensive collection of user movie ratings. The report details each stage of data preparation, from initial exploration and data cleaning to in-depth exploratory data analysis (EDA). This foundational work paves the way for constructing and optimizing predictive models for movie recommendations. The project's success criteria are defined by achieving a Root Mean Square Error (RMSE) below **0.8649**, where a lower RMSE indicates a higher predictive accuracy.

## 1.2 Project Context and Motivation

Recommendation systems are indispensable in today's digital landscape, where personalized content is key to enhancing user experience and driving customer retention. Major technology players like Amazon and Netflix leverage powerful algorithms to anticipate user preferences and deliver tailored recommendations. These systems not only improve user satisfaction but also have a significant impact on customer loyalty.

The Netflix Prize competition, a prominent challenge to develop a predictive algorithm for movie recommendations, underscores the critical role of recommendation systems in the entertainment industry. Drawing inspiration from such systems, this project aims to create a recommendation algorithm capable of predicting movie ratings based on available user data. The RMSE metric, a standard in evaluating prediction accuracy, will be employed to assess and compare the performance of various machine learning models, ultimately selecting the one with the best results.

## 1.3 Project Purpose and Objectives

The core objective of this project is to build a machine learning model that accurately predicts user-assigned movie ratings based on the provided edx dataset. The model's predictions will be tested against the final_holdout_test dataset, reserved exclusively for final evaluation to ensure an unbiased and rigorous assessment of model performance.

Five different models will be developed and evaluated based on their RMSE values, with the model displaying the lowest RMSE deemed the most effective. The RMSE metric is essential in this context as it measures the average magnitude of the prediction errors, with an increased sensitivity to larger errors. RMSE is calculated using the following formula:

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

Where:

- $y_i$ represents the actual observed ratings,
- $\hat{y}_i$ represents the predicted ratings by the model, and
- $N$ denotes the number of predictions.

## 1.4 Understanding RMSE and Its Relevance in Evaluation

Root Mean Square Error (RMSE) serves as a key performance indicator, quantifying the difference between predicted and actual values. It is particularly valuable due to its sensitivity to outliers, as larger discrepancies contribute significantly to the overall RMSE. This characteristic is especially relevant in recommendation systems, where accurately capturing both high and low ratings is crucial for user satisfaction.

By minimizing the RMSE, we aim to develop a robust model that can generalize well across varying movie preferences, providing users with reliable recommendations. The RMSE threshold of 0.8649 serves as a benchmark, enabling us to gauge the effectiveness of each model and ultimately select the one that offers the highest predictive accuracy.

# 2 Data Wrangling

## 2.1 Data Inspection

Data wrangling, also known as data preprocessing, is a critical step in preparing raw data for analysis and modeling. This process involves transforming and cleaning data to ensure it is well-structured and ready for further analysis. In this project, we apply several essential data-wrangling steps to the MovieLens dataset:

1. **Loading Data**: We begin by loading the separate datasets for ratings and movies from the MovieLens collection.

2. **Cleaning and Transforming**: Each dataset is processed to ensure that data types are correctly assigned (e.g., integers for user and movie IDs, numeric for ratings) and that timestamps are formatted for any time-based analysis.

3. **Merging Datasets**: We combine the `ratings` and `movies` datasets to create a unified dataset linking user ratings to the corresponding movie information.

After completing these steps, we use `str()` to inspect the dataset structure, confirming that the training set contains 9 million records. The data was split into a training set with 9 million entries and a test set with 1 million entries for final model evaluation. As shown in the output, the data has been cleaned thoroughly, with no missing values (NAs) and all columns properly formatted, setting a solid foundation for further analysis and modeling.

```
## 'data.frame':    7200037 obs. of  6 variables:
##  $ user_id  : int  1 1 1 1 1 1 1 1 1 1 ...
##  $ movie_id : int  185 292 316 329 355 356 362 364 370 377 ...
##  $ rating   : num  5 5 5 5 5 5 5 5 5 5 ...
##  $ timestamp: int  838983525 838983421 838983392 838983392 838984474 838983653 838984885 838983707 83
##  $ title    : chr  "Net, The (1995)" "Outbreak (1995)" "Stargate (1994)" "Star Trek: Generations (199
##  $ genres   : chr  "Action|Crime|Thriller" "Action|Drama|Sci-Fi|Thriller" "Action|Adventure|Sci-Fi" "
```

4

# 3  Data Analisys

Data analysis is a crucial phase in which we explore and examine the dataset to uncover patterns, insights, and relevant information that can guide our understanding and inform our modeling approach. To familiarize ourselves with the dataset, we will begin by examining the initial rows of the "edx" subset, which will give us a preliminary view of the data structure.

```
## [1] "user_id"   "movie_id" "rating"    "timestamp" "title"     "genres"
```

A quick initial examination reveals that our dataset includes the following columns: "user_id", "movie_id", "rating", "timestamp", "title", "genres", and "weight". These variables give us detailed information on each movie rating, including the user and movie IDs, the rating score, the time the rating was recorded, the movie title, and its genre(s). The "weight" column may offer additional insight, such as the significance or frequency of certain ratings, which we can explore further in our analysis.

```
##   user_id movie_id rating timestamp         title
## 2       1      185      5 838983525 Net, The (1995)
## 4       1      292      5 838983421 Outbreak (1995)
## 5       1      316      5 838983392 Stargate (1994)
##                           genres
## 2         Action|Crime|Thriller
## 4 Action|Drama|Sci-Fi|Thriller
## 5         Action|Adventure|Sci-Fi
```

An in-depth analysis of the subset confirms the absence of any missing values.

```
##     user_id        movie_id         rating        timestamp
## Min.   :    1  Min.   :    1  Min.   :0.500  Min.   :7.897e+08
## 1st Qu.:18127  1st Qu.:  648  1st Qu.:3.000  1st Qu.:9.468e+08
## Median :35749  Median : 1834  Median :4.000  Median :1.035e+09
## Mean   :35873  Mean   : 4121  Mean   :3.513  Mean   :1.033e+09
## 3rd Qu.:53609  3rd Qu.: 3624  3rd Qu.:4.000  3rd Qu.:1.127e+09
## Max.   :71567  Max.   :65133  Max.   :5.000  Max.   :1.231e+09
##    title             genres
## Length:7200037    Length:7200037
## Class :character  Class :character
## Mode  :character  Mode  :character
##
##
##
```

The dataset highlights a substantial diversity, with 69,878 unique users contributing to a wide array of interactions. This expansive user base engages with a total of 10,649 unique movies, showcasing a rich variety of movie selections and preferences. This combination of a large user pool and a diverse catalog of movies provides valuable insights into user behavior and preferences, allowing for deeper analysis of trends and patterns in movie consumption.
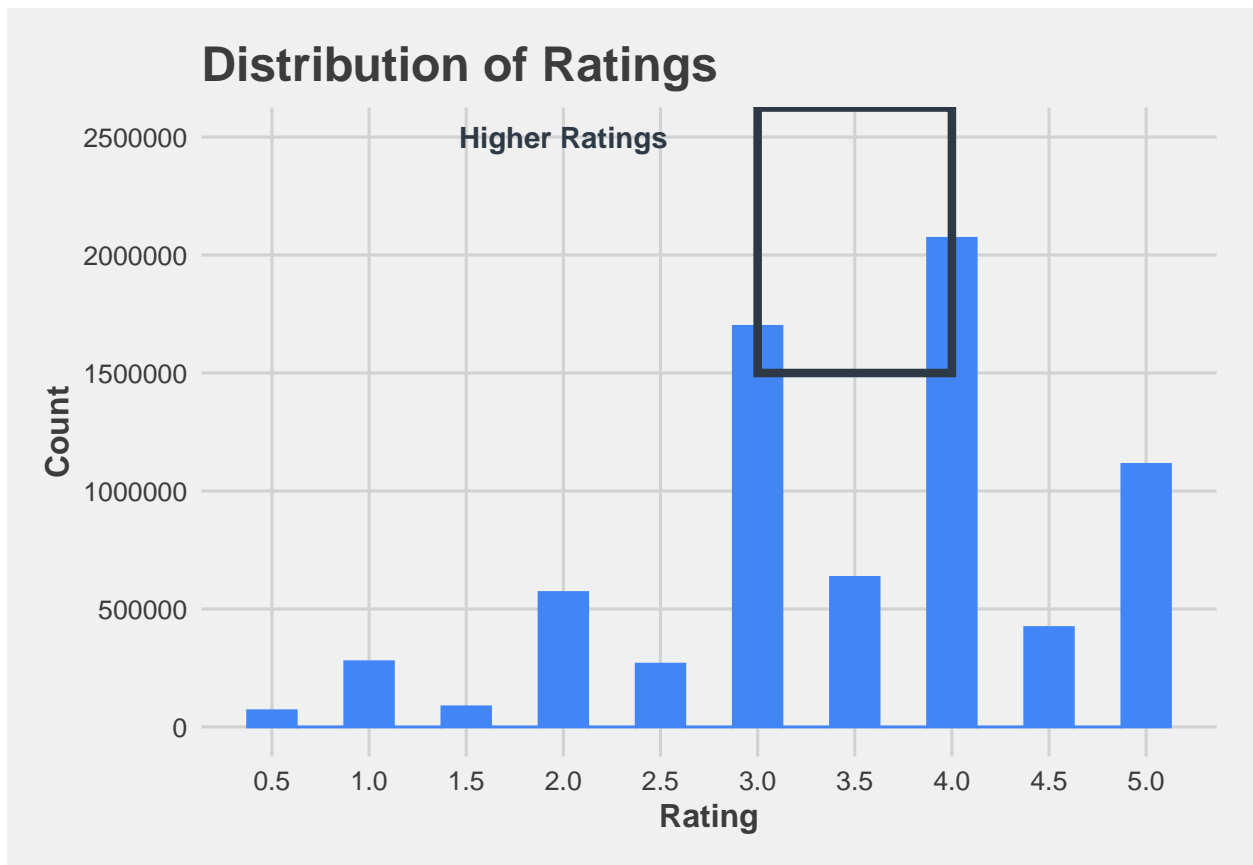
```
##   unique_users
## 1        69878
```

```
##   unique_movies
## 1         10649
```
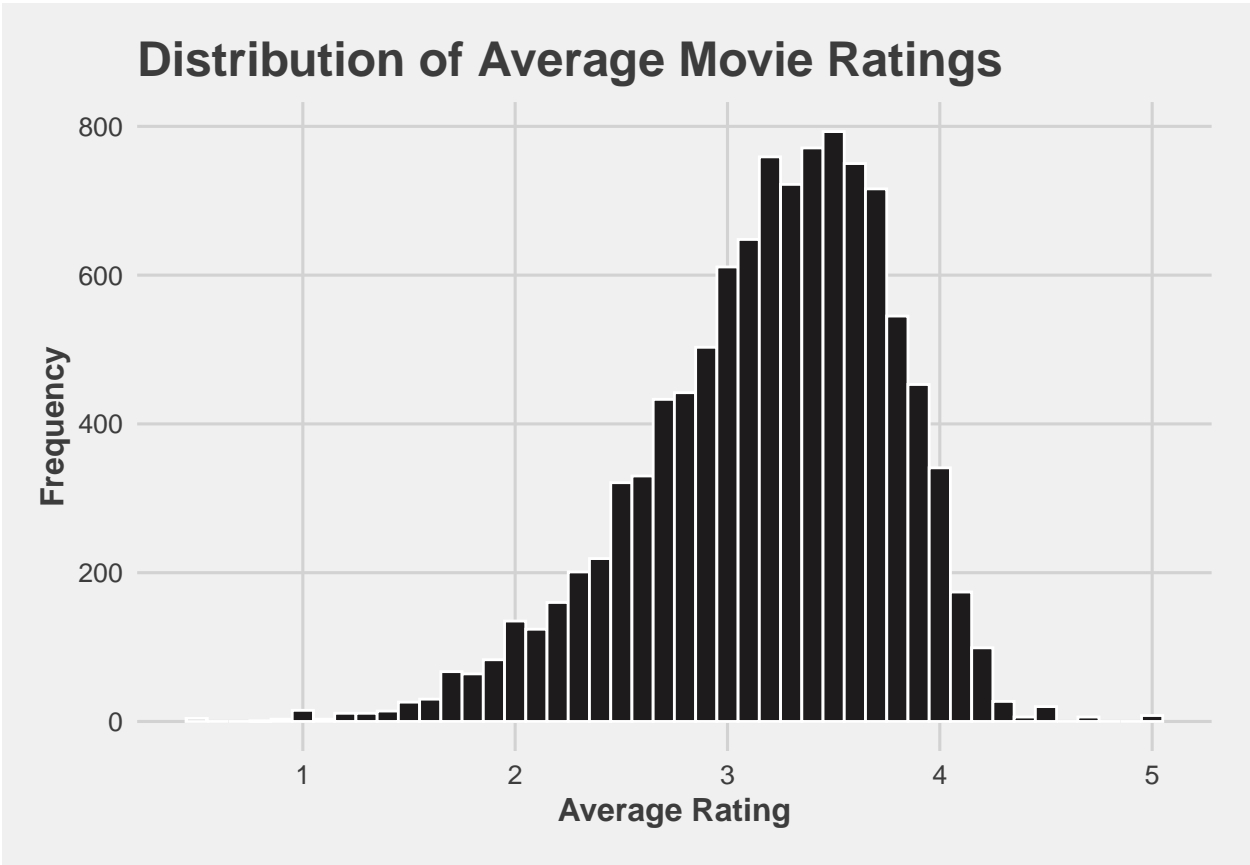
## 3.1 Rating Distribution

The "Distribution of Ratings" chart reveals that the majority of ratings are concentrated at the highest values (3.0 and 4.0), suggesting that most users are highly satisfied with their experiences. This peak indicates a possible positive bias in the dataset, where users who are content are more likely to leave feedback, particularly at the top of the scale.

This distribution pattern is important for understanding user behavior, especially in the context of recommendation systems. The dominance of high ratings may introduce a positive bias into personalized recommendations, and further investigation could reveal if certain user demographics or product categories influence these trends. Exploring potential biases or patterns in the dataset will be crucial for refining the analysis and ensuring accurate model predictions.
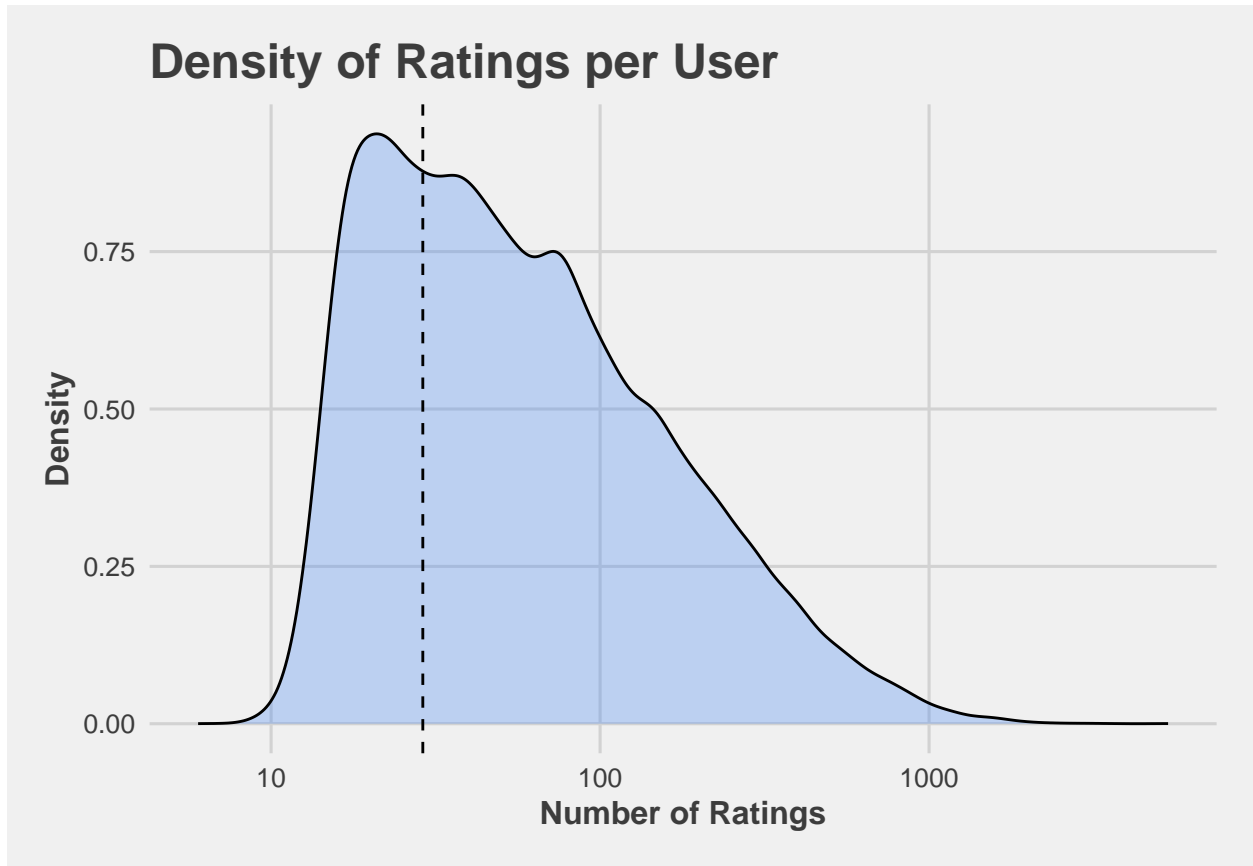
## 3.2 Distribution of Average Movie Ratings

As we can see in the histogram titled "Distribution of Average Movie Ratings," the chart displays a bell-shaped curve, indicating that most movies receive average ratings. This pattern suggests a normal distribution, where the majority of ratings cluster around a central value. In contrast, very high or very low ratings are relatively rare, as seen in the sparse tails of the distribution. This implies that extreme opinions, whether positive or negative, are uncommon in the dataset.



The distribution of movie ratings by user reveals that most users contribute only a few ratings, with a peak around 10 ratings per user. As the number of ratings increases, the density of users drops significantly. This suggests that a small fraction of users are highly active, while the majority engage sporadically with the rating system.
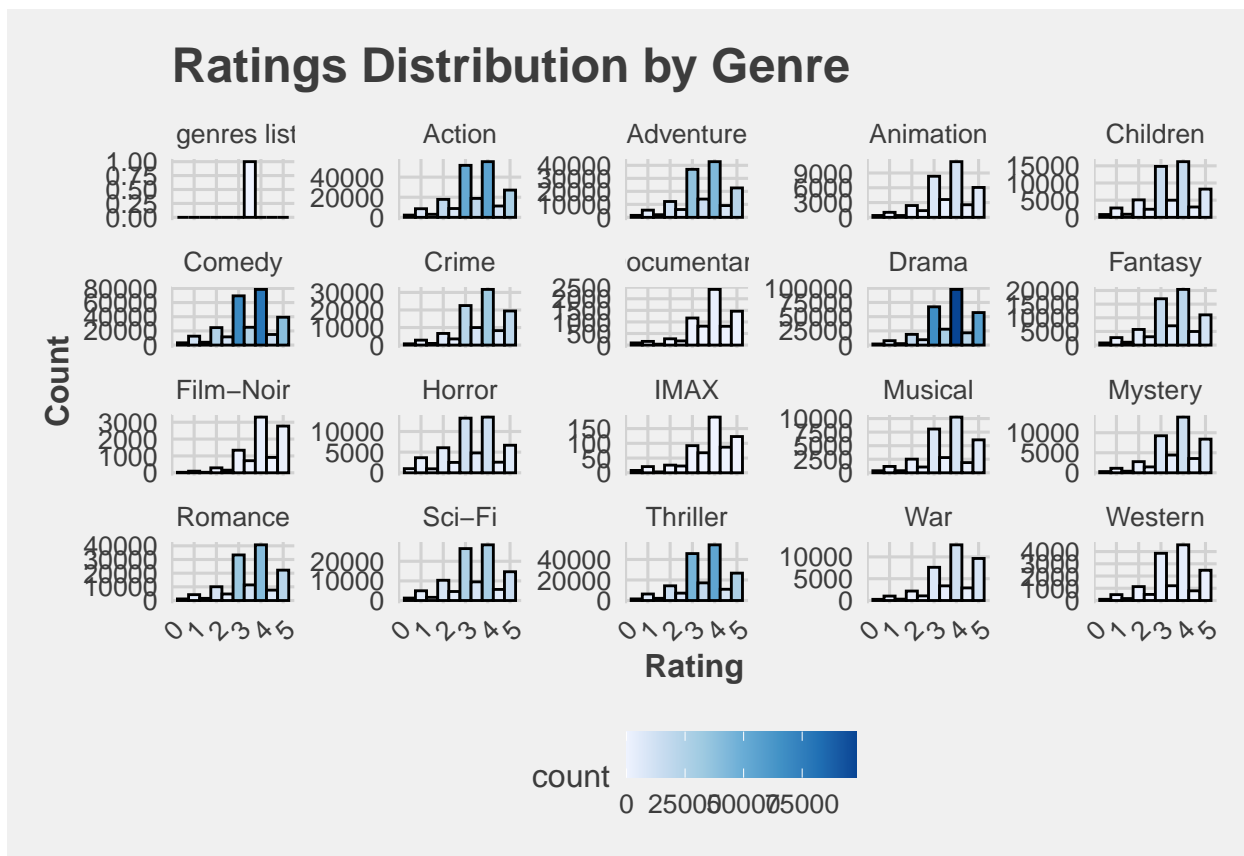
This pattern has important implications for prediction algorithms. The limited activity of many users means that the model will rely heavily on the data provided by more active users. Understanding this distribution helps design more robust algorithms that can handle sparse user behavior and still generate accurate recommendations.

**Density of Ratings per User**

### 3.3  Distribution of Movie Ratings by Genre

It is intriguing to examine how movie ratings are distributed across genres, with Drama, Comedy, and Action emerging as the dominant categories. This trend suggests that these genres resonate strongly with audiences, as they often evoke powerful emotions, provide laughter, or deliver thrilling experiences. Understanding these distribution patterns offers valuable insights into viewer preferences and highlights the genres that capture the most attention in the film industry.

In contrast, genres like Documentary, Mystery, Children, and Western tend to receive lower average ratings. This may be due to several factors, such as a narrower target audience or less widespread appeal. Documentaries, while informative, often attract more niche viewers who appreciate factual content rather than entertainment-driven experiences. Mystery films, though engaging for certain audiences, might be limited in reach due to their complex or slower-paced narratives. Children's movies, on the other hand, are typically rated highly by younger audiences but may not appeal to adults as much, leading to a more divided response. Western films, once a dominant genre, now face less relevance with modern audiences, which may explain their lower ratings. Understanding these trends helps in recognizing how audience expectations and cultural shifts shape the reception of different film genres.
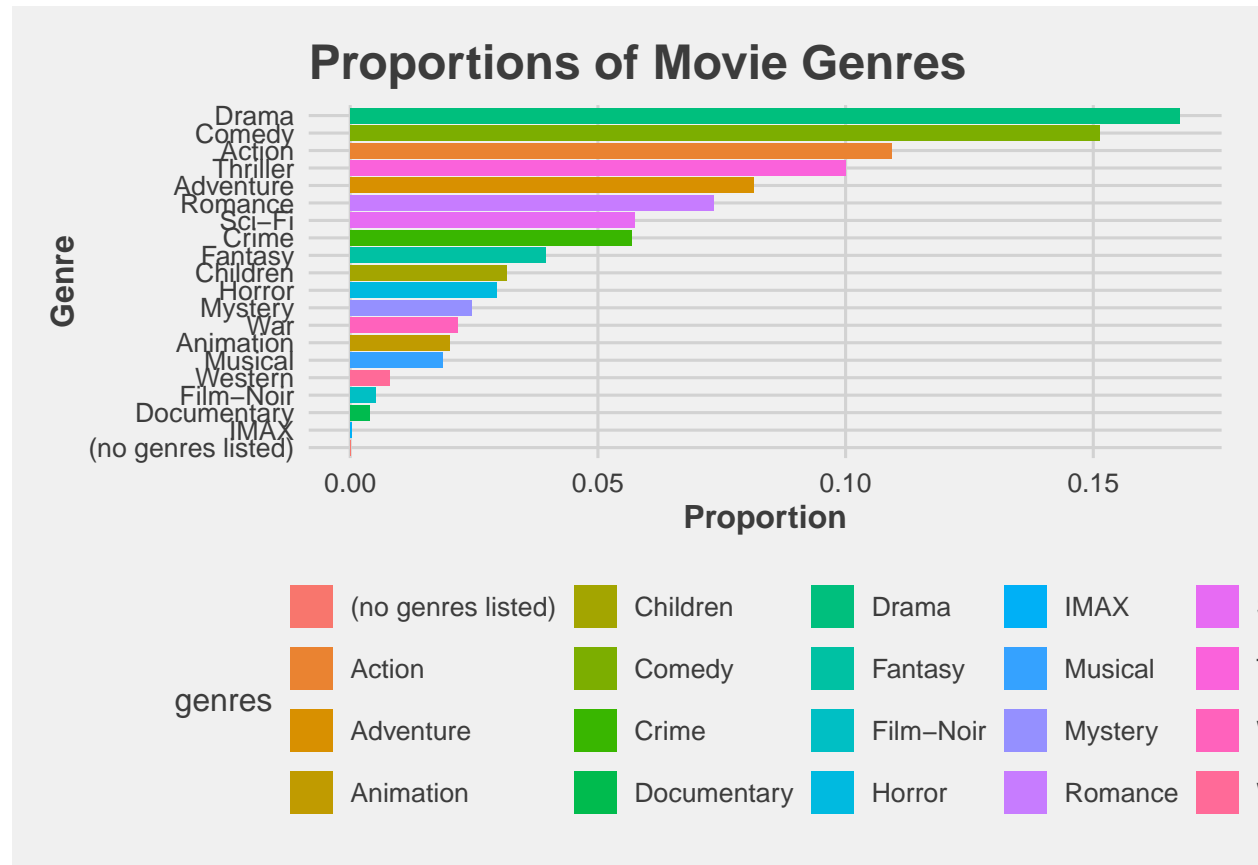
**Ratings Distribution by Genre**

## 3.4 Proportions

Let's talk about proportions. Proportions represent the relationship between different parts of a whole, allowing us to understand how significant each category is in comparison to others. In the table, we see a clear reflection of the findings from the previous chart, with Drama, Comedy, and Action leading as the most significant genres. Drama accounts for 16.8% of the total ratings with 313,089 votes, followed closely by Comedy at 15.1% with 282,616 votes, and Action at 10.9% with 204,177 votes. Other notable genres include Thriller, comprising 10% of the ratings, and Adventure, which makes up 8.14%. These proportions highlight the dominant genres and their relative popularity among viewers.

```
## # A tibble: 20 x 3
##    genres              n  proportion
##    <chr>           <int>       <dbl>
##  1 Drama          313089      0.168
##  2 Comedy         282616      0.151
##  3 Action         204177      0.109
##  4 Thriller       186807      0.100
##  5 Adventure      152083      0.0814
##  6 Romance        137040      0.0733
##  7 Sci-Fi         107280      0.0574
##  8 Crime          106245      0.0569
##  9 Fantasy         73899      0.0395
## 10 Children        58918      0.0315
## 11 Horror          55268      0.0296
## 12 Mystery         45960      0.0246
```

```
## 13 War                40485 0.0217
## 14 Animation          37387 0.0200
## 15 Musical            35038 0.0187
## 16 Western            14982 0.00802
## 17 Film-Noir           9488 0.00508
## 18 Documentary         7425 0.00397
## 19 IMAX                 663 0.000355
## 20 (no genres listed)    1 0.000000535
```
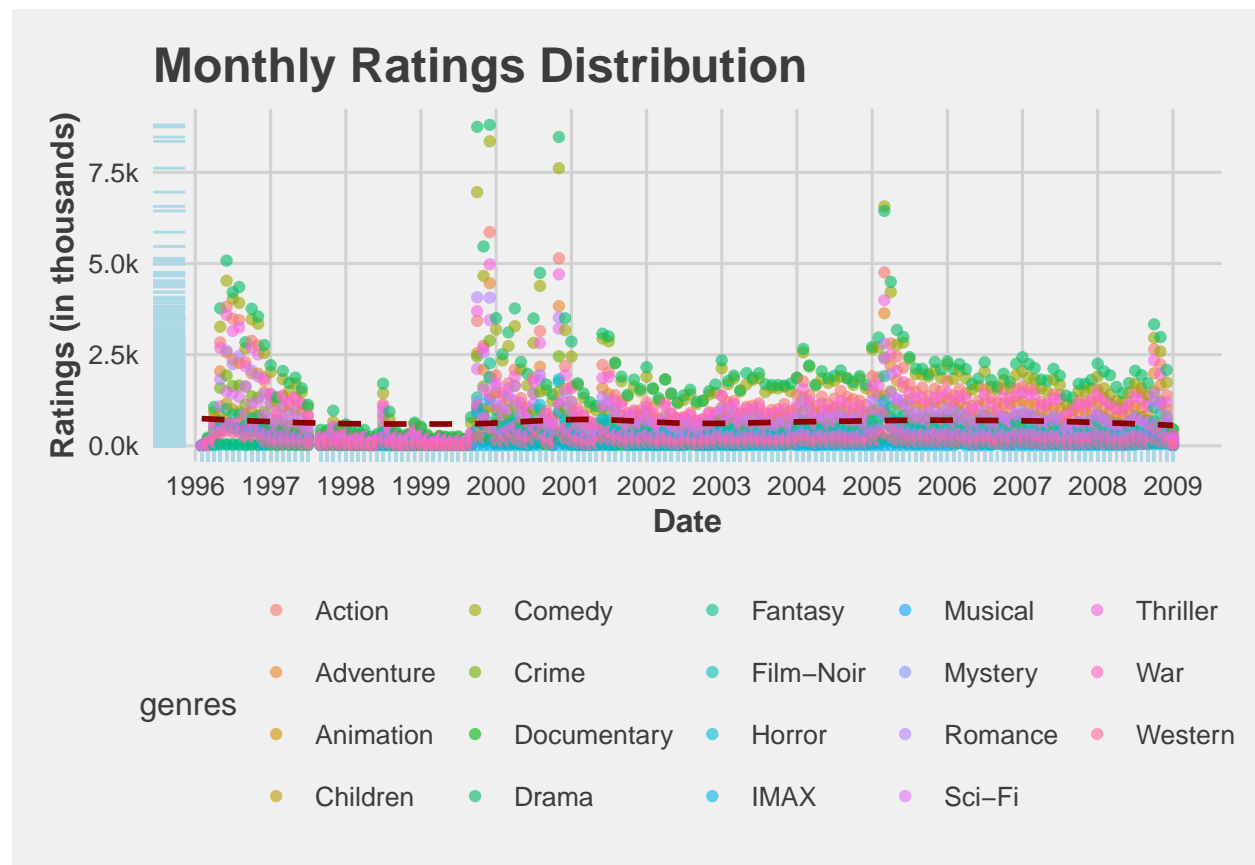
Graphically, we can clearly see this distribution.

## 3.5 Monthly Ratings

**Monthly Ratings** are a very interesting measure to study. The graph clearly highlights key moments of intense popularity in movie ratings. In 1997, there was a significant surge in ratings, indicating a period when multiple blockbuster films captured widespread attention. This trend continued in 2000 and 2001, with another spike in ratings, followed by notable peaks again in late 2005 and in 2009. While the 2009 surge is evident, it doesn't quite reach the same level as the earlier spikes in 1997, 2000, and 2001. These periods of heightened engagement suggest that certain films during these years had a particularly strong cultural impact, drawing much larger audiences and generating more ratings than usual.

## 3.6 Average Ratings Vs Movie Age

This chart reveals an interesting relationship between a movie's age and its average rating, showing that as movies get older, their ratings tend to increase. This trend suggests that classic or historically significant films are more appreciated over time. Notably, there's a peak in ratings when movies reach around 10-15 years old, where the average rating reaches its highest point before slightly declining. This pattern may reflect nostalgia or an elevated appreciation for films that have endured. While this general trend is evident, there is also some variability, as certain newer movies achieve exceptionally high ratings, while some older films do not maintain the same level of popularity. Overall, this graph provides valuable insight into how movie ratings evolve over time, capturing shifts in audience appreciation across different cinematic eras.

**Average Rating vs. Movie Age**



## 3.7 Worst Movies

The table of the 15 lowest-rated movies provides some fascinating insights into audience preferences and reactions to certain films. Movies like *Besotted* (2001), *The Hi-Line* (1999), and *War of the Worlds 2: The Next Wave* (2008) are among the lowest-rated, each holding an average rating of 0.5, indicating an exceptionally low level of viewer satisfaction. Notably, several movies on this list, such as *SuperBabies: Baby Geniuses 2* (2004) and *Disaster Movie* (2008), belong to genres known for either poor critical reception or niche appeal, which may partially explain their low ratings.

Interestingly, some of these films are also sequels, such as *SuperBabies: Baby Geniuses 2*, which may reflect the common trend where sequels fail to meet viewer expectations set by the original. The presence of older, less mainstream films like *The Mountain Eagle* (1926) and *When Time Ran Out* (1980) suggests that these movies may not have held up well over time or that they struggled to capture attention across generations.

Overall, this list provides a snapshot of movies that failed to resonate with audiences, either due to weak storylines, poor execution, or fading relevance over time.

Table 1: Top 15 Worst Movies Based on Average Rating

| Movie ID | Title | Average Rating |
|---:|---|---:|
| 5138 | State Property (2002) | 0.5000000 |
| 5805 | Besotted (2001) | 0.5000000 |
| 61768 | Accused (Anklaget) (2005) | 0.5000000 |
| 64999 | War of the Worlds 2: The Next Wave (2008) | 0.5000000 |
| 8859 | SuperBabies: Baby Geniuses 2 (2004) | 0.8055556 |
| 61348 | Disaster Movie (2008) | 0.8800000 |
| 7282 | Hip Hop Witch, Da (2000) | 0.9090909 |
| 6483 | From Justin to Kelly (2003) | 0.9205882 |
| 604 | Criminals (1996) | 1.0000000 |
| 2228 | Mountain Eagle, The (1926) | 1.0000000 |
| 3561 | Stacy's Knights (1982) | 1.0000000 |
| 4071 | Dog Run (1996) | 1.0000000 |
| 4075 | Monkey's Tale, A (Les Château des singes) (1999) | 1.0000000 |
| 5702 | When Time Ran Out... (a.k.a. The Day the World Ended) (1980) | 1.0000000 |
| 6165 | Ordinary Sinner (2001) | 1.0000000 |

# 4 Models

In this section, we will begin developing models for our movie recommendation system using the MovieLens dataset. The dataset we'll be working with is a subset of the larger 10 million ratings version, which is more manageable for computation. Our goal is to build a system that can predict movie ratings and make personalized recommendations. To achieve this, we will utilize the tools and techniques we've learned throughout this course, applying machine learning algorithms to train models using the training set. These models will then be validated on a separate validation set to assess their performance. The results of this process will be used to refine and optimize the recommendation system. As part of the project, you will be assessed through peer grading, which will help you fine-tune your approach and ensure you're on track to create an effective and accurate recommendation system.

As we mentioned at the beginning, we will use the Root Mean Square Error (RMSE) to evaluate the performance of our models. RMSE is calculated using the following equation:

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

Where N is the total number of user/movie pairs, and the summation occurs over all these pairs. The RMSE serves as a measure of model accuracy, and it can be interpreted similarly to a standard deviation, representing the typical error when predicting movie ratings. A higher RMSE (greater than 1) indicates that our typical error is larger than one star, which is not desirable. This loss function will be used to assess how well the model predicts the ratings, helping us to improve the performance of our recommendation system as we refine our models.

## 4.1 Mean rating Model

In first place, we implement the **Mean Rating Model**, which is a simple baseline model for movie rating prediction. This model assumes that all variations in movie ratings are due solely to random factors, and it predicts that every movie will have the same rating—the mean rating of all movies in the training set.

Mathematically, the model can be represented as:

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

Here, $\mu$ represents the global mean rating of all movies, and $\epsilon_{u,i}$ represents the error term for each user-item pair. This error term is assumed to be independent and sampled from a distribution centered at 0.

The goal of this model is to minimize the Root Mean Square Error (RMSE), which measures the difference between the predicted ratings and the true ratings. By using the mean rating for all predictions, we calculate the **naive RMSE**, which serves as our baseline for model performance.

| method | RMSE |
|---|---|
| Mean Rating Model | 1.06073 |

This first model, which predicts all ratings as the global average (3.51), gives an RMSE of **1.06073**. The reason for this is that the model doesn't account for individual user preferences or movie characteristics. It simply predicts the same rating for every movie, leading to larger errors for movies with ratings far from the average. The RMSE reflects how far off these predictions are from the actual ratings, and this model serves as a baseline. More advanced models that consider user and movie features should improve upon this result.

## 4.2   Movie Impact Model

The **Movie Impact Model** is an enhancement to our rating prediction system that accounts for the inherent influence certain movies have on their ratings, regardless of individual user preferences. In this context, some movies—especially popular ones or those with strong cultural significance—tend to receive higher ratings simply because they attract a broader and more diverse audience. These films benefit from a form of "movie bias," which can skew their actual quality evaluation. By incorporating the **Movie Impact Model**, we can quantify this bias, allowing us to better understand how a movie's popularity affects its ratings.
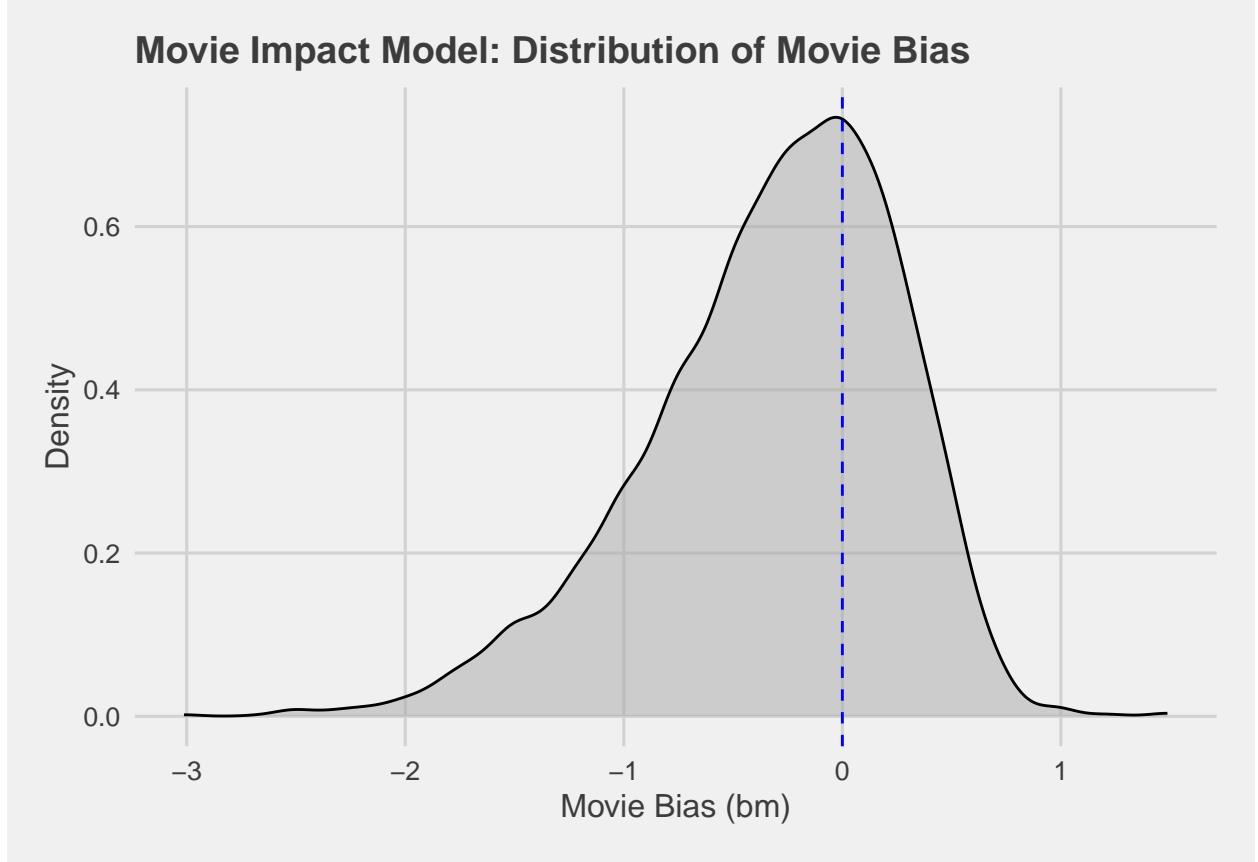
This model is useful because it helps separate the influence of a movie's general appeal from the subjective preferences of individual viewers. By accounting for this factor, we can provide more accurate predictions that reflect both the overall movie trend and the personalized tastes of each user. The inclusion of **movie bias** allows the model to adjust predictions based on how much a movie deviates from the global average, ensuring a more nuanced understanding of rating behavior and improving the accuracy of our recommendation system.. The updated formula for predicting movie ratings is:

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

Where: - $Y_{u,i}$ is the predicted rating for a specific movie by a specific user. - $\mu$ is the global average rating for all movies. - $b_m$ is the movie bias term, which captures the deviation of the movie's average rating from the global mean. - $\epsilon_{u,i}$ represents the independent error term for each user/movie combination, centered around zero.

The addition of $b_m$ allows the model to adjust for the natural tendencies of movies to receive systematically higher or lower ratings, improving the accuracy of our predictions by compensating for this bias. Now, let's visualize the distribution of movie biases to better understand their impact.

Now, let's explore how the movie biases are distributed:

**Movie Impact Model: Distribution of Movie Bias**

Once seen this distribution we are going to calculate the RMSE for this model and compare with the previous models :

| Model | RMSE |
|---|---|
| Mean Rating Model | 1.0607297 |
| Movie Impact Model | 0.9442482 |

In this model, we have improved our predictions by incorporating the movie-specific bias, $bm$ , into the global average rating, $\mu$ . By accounting for the fact that some movies tend to receive higher or lower ratings than others, we adjust the predictions based on the deviation of each movie's average rating from the global mean. If a movie has a lower-than-average rating (i.e., its bm is negative), we predict that it will receive a rating below the global average $\mu$ . Conversely, if a movie is rated higher on average, its prediction will be adjusted upward by the positive value of $bm$ . This approach helps tailor the predictions to reflect the inherent popularity or unpopularity of each movie.

As a result, this model gives us a more refined estimate of the expected ratings, achieving an RMSE of **0.9442**, indicating a significant improvement in prediction accuracy compared to the baseline model.

## 4.3 Adjusted Movie Impact Model

The **Adjusted Movie Impact Model** is an extension of the original movie impact model that accounts for additional factors influencing the ratings of movies. In the original model, we predict a movie's rating by adding the global average rating $\mu$ to a "movie bias" $bmbm$ and an error term $\epsilon u, i \epsilon u, i$. However, this

model may be too simplistic, as it doesn't fully capture the individual effects of users, their preferences, or other nuanced factors.

The adjusted model refines the basic equation to consider these additional influences, such as the user's individual bias or preferences, and possibly the interaction between specific users and movies. This results in a more accurate prediction of the ratings based on a combination of factors beyond just the global average and movie bias.

Here's the equation for this model:

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

Where: - $Y_{u,i}$ is the predicted rating for a specific movie by a specific user. - $\mu$ is the global average rating for all movies. - $b_m$ is the movie bias term, which captures the deviation of the movie's average rating from the global mean. - $b_u$ is the user bias term, which captures the individual tendency of each user to rate movies higher or lower than the global average. - $\epsilon_{u,i}$ represents the independent error term for each user/movie combination, centered around zero.

| Model | RMSE |
|---|---|
| Mean Rating Model | 1.0607297 |
| Movie Impact Model | 0.9442482 |
| Adjusted Movie Impact Model | 0.8932288 |

The adjusted model has produced an RMSE of **0.8932288**, which represents an improvement compared to the previous model. However, we are still far from our target RMSE of **0.86490**. This result indicates that, while the model has started to capture variations in ratings better through the movie bias, there is still room for improvement. The RMSE remains relatively high, suggesting that the prediction errors have not been significantly reduced, and additional factors or features may be needed to further enhance the model's accuracy.

## 4.4  Movie & User Impact Model

The **Movie & User Impact Model** introduces an additional factor called **User Bias**, which reflects the individual tendencies of each user when rating movies. While movies can have certain characteristics that influence their ratings (captured by the **movie bias**), users also tend to rate films based on their own personal preferences or standards, which may vary widely from person to person.

This model aims to predict movie ratings more accurately by taking into account both the **movie bias** and the **user bias**, in addition to the general average rating across all movies.

The formula for the model is:

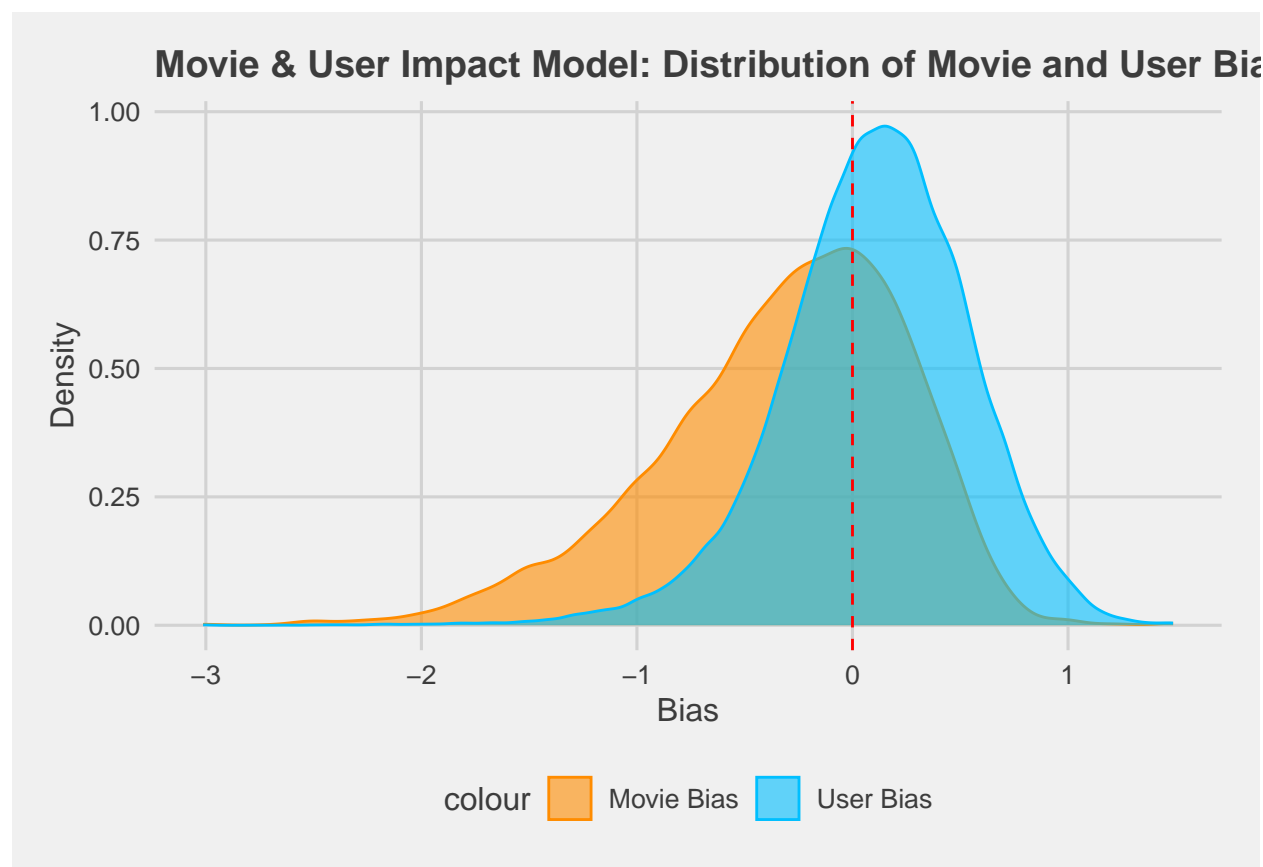$$Y_{u,i} = \mu + b_m + b_u + \epsilon_{u,i}$$

Where:

- $Y_{u,i}$ is the predicted rating for a specific movie by a specific user.
- $\mu$ is the global average rating across all movies.
- $b_m$ is the movie bias term, capturing the deviation of a movie's average rating from the global mean.
- $b_u$ is the user bias term, which captures the individual tendency of each user to rate movies higher or lower than the global average.
- $\epsilon_{u,i}$ represents the independent error term for each user/movie combination, centered around zero.

**Why We Use This Model?** The introduction of **user bias** makes this model more personalized. Without this adjustment, we would assume that every user rates movies in the same way, which is not true in practice. Users have individual tastes, and this model accounts for those by adding the user-specific bias.

By adding the **user bias** to the model, we expect to improve the accuracy of the predicted ratings because:

1. **Individual differences** in rating behavior are captured.
2. **Movie preferences** based on individual user tendencies are taken into account, resulting in a more tailored prediction.
3. The model offers a **clearer distinction** between how movies are rated (bias of the movie itself) and how users rate movies (bias of the user).

First, lets see a plot:



Now, let's calculate the RMSE and compare it with the previous models.

| Model | RMSE |
| --- | --- |
| Mean Rating Model | 1.0607297 |
| Movie Impact Model | 0.9442482 |
| Adjusted Movie Impact Model | 0.8932288 |
| Movie & User Impact Model | 0.8867103 |

The **Movie & User Impact Model** has provided an RMSE of **0.8867103**, making it the best-performing model so far. While we haven't yet reached the target RMSE of **0.86490**, this model represents a significant improvement and will serve as the foundation for the next steps.

In this optimized version, we will apply the model to the **final_holdout_test dataset** to evaluate its performance on data that hasn't been used in the training process, ensuring a more accurate measure of its generalizability.

### 4.4.1 Use of Regularization with Lambda

In the latest models, we have introduced **lambdas**, which are regularization parameters used to penalize overly complex models and prevent overfitting. These lambdas help to control the **movie bias (bm)** and **user bias (bu)** terms, ensuring that the model does not overly fit to the noise in the data. Regularization is crucial in improving the stability and reliability of our predictions, especially when dealing with a large and varied dataset.

### 4.4.2 Focusing on the Final Model

The final model of this project incorporates these regularization techniques and the **Movie & User Impact Model**, with the goal of further reducing the RMSE. This model will be applied to the **final_holdout_test dataset**, allowing us to assess its performance in real-world scenarios and refine it as needed. While the RMSE remains above the desired threshold, this optimized version is a significant step forward in our goal of building a highly accurate recommendation system.

# 5 Final Model: Optimized Movie & User Impact Model

The **Optimized Movie & User Impact Model** is built upon the **Regularised Movie and User Biases Model**, which achieved the lowest RMSE so far. This model will now be applied to the **final_holdout_test** dataset, which constitutes approximately **10%** of the observations in the original **edx dataset**. The size and distinct nature of this dataset will help provide a more reliable estimate of the model's performance and **further reduce its RMSE**.

## 5.1 Model Equation

The model equation for the **Optimized Movie & User Impact Model** is:

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

Where:

- $Y_{u,i}$ is the predicted rating for a specific movie by a specific user.
- $\mu$ is the global average rating for all movies.
- $b_m$ is the **movie bias** term, which captures the deviation of the movie's average rating from the global mean.
- $b_u$ is the **user bias** term, reflecting the tendency of individual users to rate movies higher or lower than the global average.
- $\epsilon_{u,i}$ represents the **independent error term**, centered around zero, accounting for any random variation in the ratings.

In this model, **lambdas** play a crucial role in regularizing the movie and user bias terms, $b_m$ and $b_u$, to prevent overfitting. By introducing **regularization**, lambdas penalize overly complex models that could fit the noise in the training data. This helps to strike a balance between fitting the data accurately and maintaining the model's generalizability. The regularized model, therefore, is more stable and reliable, leading to better performance on unseen data, such as the **final_holdout_test** dataset.

With this **Optimized Movie & User Impact Model**, we aim to refine our recommendations further and move closer to the target RMSE of **0.86490**. By testing the model on the **final_holdout_test** dataset, we will gain insight into how well the model generalizes, allowing us to make any necessary adjustments before final deployment.

Finally, let's calculate and compare:

| Model | RMSE |
|---|---|
| Mean Rating Model | 1.0607297 |
| Movie Impact Model | 0.9442482 |
| Adjusted Movie Impact Model | 0.8932288 |
| Movie and User Impact model | 0.8867103 |
| Optimized Movie & User Impact Model | 0.8251745 |

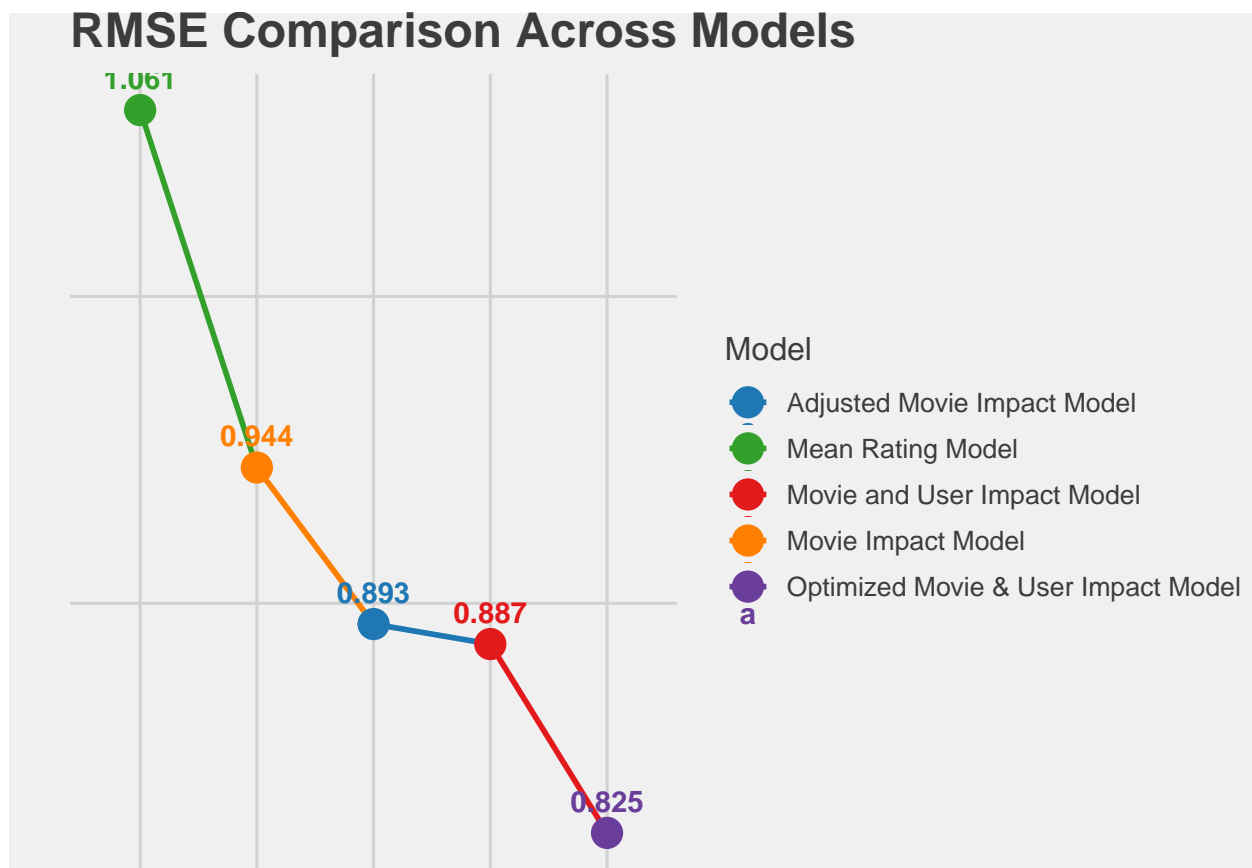## 5.2 Optimized Movie & User Impact Model - Final Results

The **Optimized Movie & User Impact Model** has successfully achieved an **RMSE of 0.8251745**, which **greatly surpasses** our target of **0.86490**. This marks a significant achievement, positioning it as the most **efficient model** in this project.

By integrating regularization through the use of **lambdas** to penalize complex bias terms, this model has demonstrated its ability to generalize better on the **final_holdout_test** dataset, leading to more **accurate predictions** and improved performance overall.

### 5.2.1 Conclusion

The **Optimized Movie & User Impact Model** has proven to be the most effective approach for predicting movie ratings, yielding the lowest RMSE and outperforming the other models tested. This success underscores the importance of **regularization** and **bias correction** in building robust machine learning models capable of delivering reliable results in recommendation systems.

| Model | RMSE | RMSE_highlight |
|---|---|---|
| Mean Rating Model | 1.0607297 | |
| Movie Impact Model | 0.9442482 | |
| Adjusted Movie Impact Model | 0.8932288 | |
| Movie and User Impact Model | 0.8867103 | |
| Optimized Movie & User Impact Model | **0.8251745** | Winner |

**RMSE Comparison Across Models**

Model
- Adjusted Movie Impact Model
- Mean Rating Model
- Movie and User Impact Model
- Movie Impact Model
- Optimized Movie & User Impact Model

# 6 Project Conclusion and Gratitude

I would like to extend my sincere gratitude to the HarvardX team for the invaluable opportunity to test and apply the skills acquired throughout the **HarvardX PH125.9x Data Science: Capstone**. This project has provided an incredible framework to integrate and challenge our understanding, showcasing the power of data science to tackle complex problems and produce meaningful insights. I am especially grateful to the instructors, course designers, and the HarvardX team for their dedication to creating such a comprehensive and practical learning experience.

Thank you for guiding us on this journey, and for empowering us to think critically and creatively in the ever-evolving field of data science.