# MovieLens Recommendation System

Dr. Luis Satch

## Table of Contents

## Introduction

In this project, I build a movie recommendation system using the MovieLens dataset, which includes user ratings and movie metadata. My goal is to develop a model for predicting movie ratings by incorporating user and movie-specific characteristics. I start by preparing the data, merging ratings and movie information, and splitting it into training and test sets. I then engineer features such as user and movie rating averages to improve predictive accuracy. I explore the data to extract meaningful insights and evaluate various modelling approaches, including baseline models and regularisation techniques to account for both movie and user effects. I validate the model using Root Mean Squared Error (RMSE) to assess its performance, and finally, I apply the best-performing model to a holdout test set to demonstrate its predictive power and accuracy.

## 1. Data Preparation

```r
# Load necessary libraries
if (!require(tidyverse)) install.packages("tidyverse", repos = "http://cran.us.r-project.org")
if (!require(caret)) install.packages("caret", repos = "http://cran.us.r-project.org")

library(tidyverse)
library(caret)

# Clean up from previous runs and clear the console
rm(list = ls())
cat("\f")
```

```r
# Check if datasets already exist in memory, load them if available
if (exists("movielens")) {
    message("Dataset already loaded. Skipping download and loading steps.")
} else {

    # MovieLens 10M dataset: URLs for dataset download
    dl <- "ml-10M100K.zip"
    ratings_file <- "ml-10M100K/ratings.dat"
    movies_file <- "ml-10M100K/movies.dat"

    # Download and unzip the dataset if it doesn't exist locally
    if (!file.exists(dl)) {
        download.file("https://files.grouplens.org/datasets/movielens/ml-10m.zip",
            dl)
    }

    # Unzip ratings and movies data
    if (!file.exists(ratings_file)) {
        unzip(dl, files = ratings_file)
    }
    if (!file.exists(movies_file)) {
        unzip(dl, files = movies_file)
    }

    # Load and process ratings data
    ratings_raw <- read_lines(ratings_file)
    ratings <- as.data.frame(str_split(ratings_raw, fixed("::"), simplify = TRUE),
        stringsAsFactors = FALSE)
    colnames(ratings) <- c("userId", "movieId", "rating", "timestamp")

    # Load and process movies data
    movies_raw <- read_lines(movies_file)
    movies <- as.data.frame(str_split(movies_raw, fixed("::"), simplify = TRUE),
        stringsAsFactors = FALSE)
    colnames(movies) <- c("movieId", "title", "genres")

    # Convert appropriate columns to integer and numeric types
    ratings <- ratings %>%
        mutate(userId = as.integer(userId), movieId = as.integer(movieId), rating = as.numeric(rating),
            timestamp = as.integer(timestamp))

    movies <- movies %>%
        mutate(movieId = as.integer(movieId))

    # Merge ratings and movie data into a single dataframe
    movielens <- left_join(ratings, movies, by = "movieId")

    # Remove temporary variables
    rm(dl, ratings_raw, movies_raw, ratings, movies)

    message("Dataset loaded and processed.")
}
```

## 2. Data Splitting

```r
# Split the data into edx and final_holdout_test sets Ensure that both movieId
# and userId in the final_holdout_test set exist in the edx set

# Set a seed for reproducibility
set.seed(1, sample.kind = "Rounding")
```

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

```r
# Split the data: 90% for edx (training set) and 10% for final_holdout_test
# (test set)
test_index <- createDataPartition(y = movielens$rating, times = 1, p = 0.1, list = FALSE)
edx <- movielens[-test_index, ]  # Training set
temp <- movielens[test_index, ]  # Temporary test set

# Ensure that movieId and userId in final_holdout_test are also in edx
final_holdout_test <- temp %>%
    semi_join(edx, by = "movieId") %>%
    semi_join(edx, by = "userId")

# Add rows removed from final_holdout_test back into edx set
removed <- anti_join(temp, final_holdout_test)
edx <- bind_rows(edx, removed)

# Clean up temporary variables
rm(test_index, temp, removed)
```

## 3. Feature Engineering: Add user and movie averages

```r
# Calculate the average rating per user
user_avgs <- edx %>%
    group_by(userId) %>%
    summarise(user_avg = mean(rating))

# Calculate the average rating per movie
movie_avgs <- edx %>%
    group_by(movieId) %>%
    summarise(movie_avg = mean(rating))

# Merge these averages into the edx dataset
edx <- edx %>%
    left_join(user_avgs, by = "userId") %>%
    left_join(movie_avgs, by = "movieId")

cat("User and movie averages added to the edx dataset.\n")
```

```
## User and movie averages added to the edx dataset.
```

## 4. Exploratory Data Analysis (EDA)

```r
# Summarise basic statistics (number of users, movies, ratings)

# Number of unique users
num_users <- edx %>%
    summarise(users = n_distinct(userId)) %>%
    pull(users)

# Number of unique movies
num_movies <- edx %>%
    summarise(movies = n_distinct(movieId)) %>%
    pull(movies)

# Total number of ratings
num_ratings <- nrow(edx)

# Summary output
cat("Number of unique users:", num_users, "\n")
```

```
## Number of unique users: 69878
```

```r
cat("Number of unique movies:", num_movies, "\n")
```

```
## Number of unique movies: 10677
```
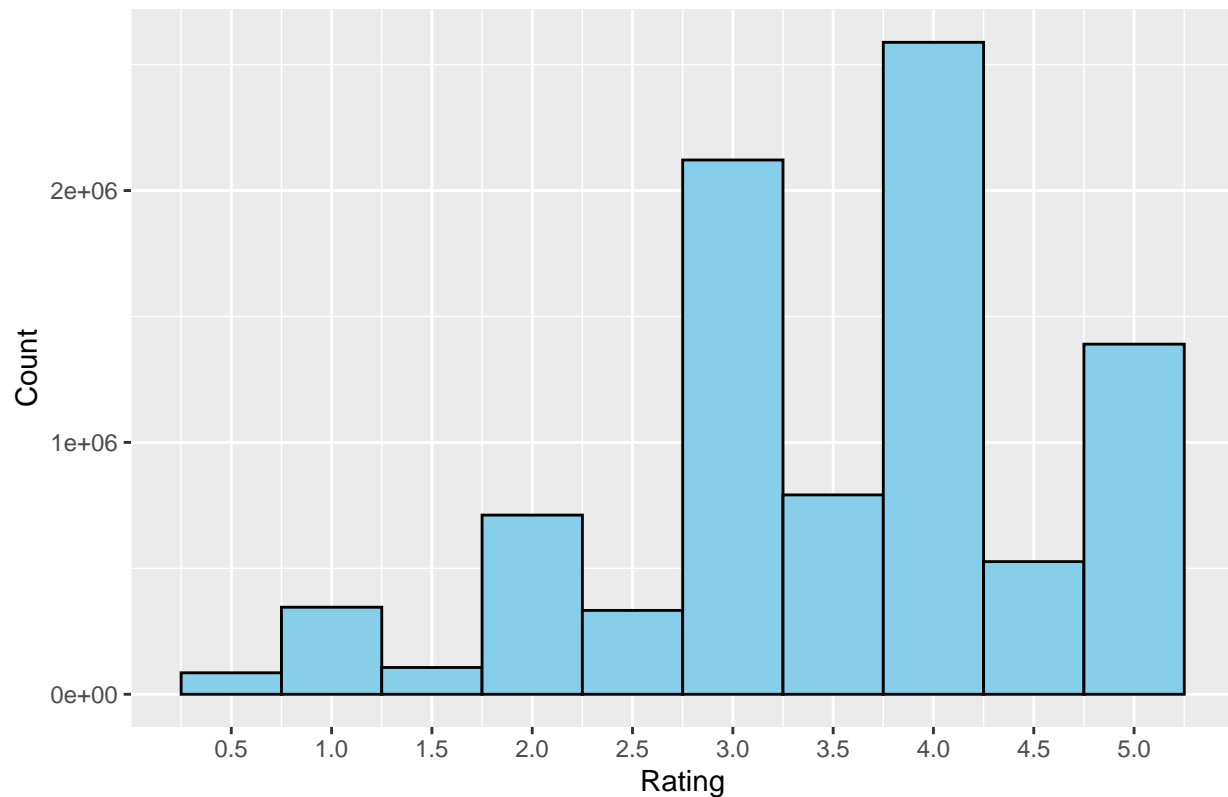
```r
cat("Total number of ratings:", num_ratings, "\n")
```

```
## Total number of ratings: 9000055
```

```r
# Visualise rating distributions
edx %>%
    ggplot(aes(x = rating)) + geom_histogram(binwidth = 0.5, color = "black", fill = "skyblue") +
    scale_x_continuous(breaks = seq(0.5, 5, by = 0.5)) + labs(title = "Distribution of Movie Ratings",
    x = "Rating", y = "Count")
```

## Distribution of Movie Ratings



```r
# Most rated movies (top 10)
top_movies <- edx %>%
    group_by(title) %>%
    summarise(count = n()) %>%
    arrange(desc(count)) %>%
    top_n(10, count)

# Display the top 10 most rated movies
cat("\nTop 10 most rated movies:\n")
```

```
##
## Top 10 most rated movies:
```

```r
print(top_movies)
```

```
## # A tibble: 10 x 2
##    title                               count
##    <chr>                               <int>
##  1 Pulp Fiction (1994)                 31362
##  2 Forrest Gump (1994)                 31079
##  3 Silence of the Lambs, The (1991)    30382
##  4 Jurassic Park (1993)                29360
##  5 Shawshank Redemption, The (1994)    28015
##  6 Braveheart (1995)                   26212
##  7 Fugitive, The (1993)                25998
```
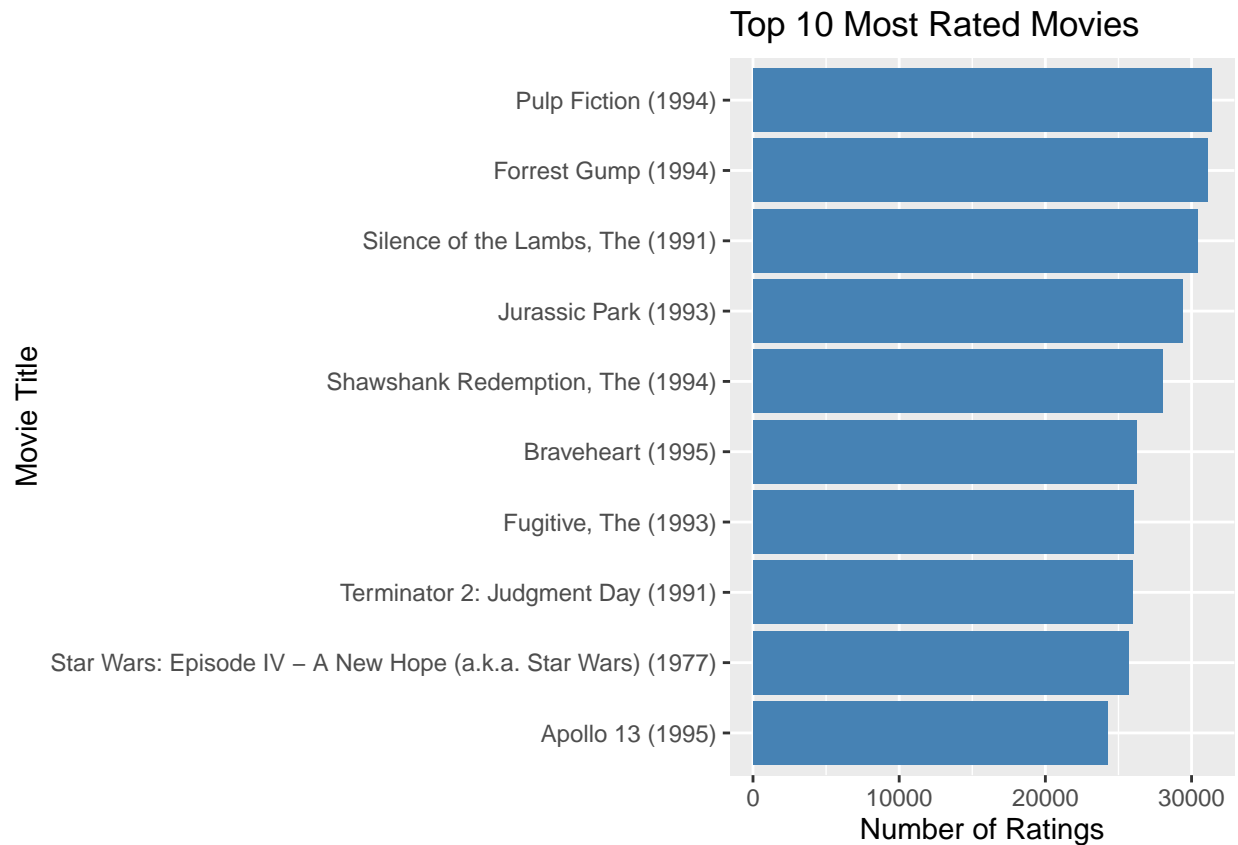
```
## 8 Terminator 2: Judgment Day (1991)                              25984
## 9 Star Wars: Episode IV - A New Hope (a.k.a. Star Wars) (1977) 25672
## 10 Apollo 13 (1995)                                              24284
```

```
# Visualise the top 10 most rated movies
top_movies %>%
    ggplot(aes(x = reorder(title, count), y = count)) + geom_bar(stat = "identity",
    fill = "steelblue") + coord_flip() + labs(title = "Top 10 Most Rated Movies",
    x = "Movie Title", y = "Number of Ratings")
```

### Top 10 Most Rated Movies



```
#
```

## 5. Modeling Approach

```
# Baseline Model: Predict the average rating for all movies
mu <- mean(edx$rating)  # Global average rating

# Display the baseline model prediction (global average rating)
cat("Baseline Model (Global Average Rating):", mu, "\n")
```

```
## Baseline Model (Global Average Rating): 3.512465
```

```r
# Movie Effect Model: Adjust the rating based on each movie's average rating
movie_effects <- edx %>%
    group_by(movieId) %>%
    summarise(b_i = mean(rating - mu))  # Movie-specific effect (deviation from global average)

# Display a few movie effects
cat("Movie Effect Model: Displaying a few movie effects\n")
```

## Movie Effect Model: Displaying a few movie effects

```r
print(head(movie_effects))
```

```
## # A tibble: 6 x 2
##   movieId    b_i
##     <int>  <dbl>
## 1       1  0.415
## 2       2 -0.307
## 3       3 -0.365
## 4       4 -0.648
## 5       5 -0.444
## 6       6  0.303
```

```r
# Movie + User Effect Model: Adjust for both movie and user-specific effects
user_effects <- edx %>%
    left_join(movie_effects, by = "movieId") %>%
    group_by(userId) %>%
    summarise(b_u = mean(rating - mu - b_i))  # User-specific effect

# Display a few user effects
cat("User Effect Model: Displaying a few user effects\n")
```

## User Effect Model: Displaying a few user effects

```r
print(head(user_effects))
```

```
## # A tibble: 6 x 2
##   userId     b_u
##    <int>   <dbl>
## 1      1  1.68
## 2      2 -0.236
## 3      3  0.264
## 4      4  0.652
## 5      5  0.0853
## 6      6  0.346
```

```r
# Regularisation: Penalise complexity by shrinking movie and user effects
lambda <- 5  # Regularisation parameter
movie_effects_reg <- edx %>%
    group_by(movieId) %>%
    summarise(b_i = sum(rating - mu)/(n() + lambda))  # Regularised movie effect
```

```r
user_effects_reg <- edx %>%
    left_join(movie_effects_reg, by = "movieId") %>%
    group_by(userId) %>%
    summarise(b_u = sum(rating - mu - b_i)/(n() + lambda))  # Regularised user effect

# Display regularised movie and user effects
cat("Regularised Movie Effect Model: Displaying a few regularised movie effects\n")
```

## Regularised Movie Effect Model: Displaying a few regularised movie effects

```r
print(head(movie_effects_reg))
```

```
## # A tibble: 6 x 2
##   movieId    b_i
##     <int>  <dbl>
## 1       1  0.415
## 2       2 -0.307
## 3       3 -0.365
## 4       4 -0.646
## 5       5 -0.443
## 6       6  0.303
```

```r
cat("Regularised User Effect Model: Displaying a few regularised user effects\n")
```

## Regularised User Effect Model: Displaying a few regularised user effects

```r
print(head(user_effects_reg))
```

```
## # A tibble: 6 x 2
##   userId     b_u
##    <int>   <dbl>
## 1      1  1.33
## 2      2 -0.183
## 3      3  0.228
## 4      4  0.571
## 5      5  0.0803
## 6      6  0.306
```

## 6. Model Validation

```r
# - Create a separate training and validation split within the edx dataset for
# model development - Calculate RMSE for each model using cross-validation or a
# validation set

# Load necessary library for calculating RMSE
if (!require(Metrics)) install.packages("Metrics", repos = "http://cran.us.r-project.org")
library(Metrics)

# Split the edx dataset into training (90%) and validation (10%) sets
set.seed(1, sample.kind = "Rounding")  # Ensure reproducibility
```

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

```r
# Create a partition: 90% training, 10% validation
train_index <- createDataPartition(y = edx$rating, times = 1, p = 0.9, list = FALSE)
train_set <- edx[train_index, ]
validation_set <- edx[-train_index, ]

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

# Baseline Model RMSE: Predicting the global average rating
baseline_rmse <- RMSE(validation_set$rating, mu)
cat("Baseline Model RMSE:", as.numeric(baseline_rmse), "\n")
```

```
## Baseline Model RMSE: 1.059799
```

```r
# Movie Effect Model RMSE
predicted_ratings_movie <- validation_set %>%
    left_join(movie_effects, by = "movieId") %>%
    mutate(pred = mu + b_i) %>%
    pull(pred)

movie_effect_rmse <- RMSE(validation_set$rating, predicted_ratings_movie)
cat("Movie Effect Model RMSE:", movie_effect_rmse, "\n")
```

```
## Movie Effect Model RMSE: 0.9427288
```

```r
# Movie + User Effect Model RMSE
predicted_ratings_user <- validation_set %>%
    left_join(movie_effects, by = "movieId") %>%
    left_join(user_effects, by = "userId") %>%
    mutate(pred = mu + b_i + b_u) %>%
    pull(pred)

user_effect_rmse <- RMSE(validation_set$rating, predicted_ratings_user)
cat("Movie + User Effect Model RMSE:", user_effect_rmse, "\n")
```

```
## Movie + User Effect Model RMSE: 0.8561705
```

```r
# Regularised Movie + User Effect Model RMSE
predicted_ratings_reg <- validation_set %>%
    left_join(movie_effects_reg, by = "movieId") %>%
    left_join(user_effects_reg, by = "userId") %>%
    mutate(pred = mu + b_i + b_u) %>%
    pull(pred)

# Calculate RMSE for the regularised model
regularised_rmse <- RMSE(validation_set$rating, predicted_ratings_reg)
```

```r
# Now print the regularised RMSE value
cat("Regularised Movie + User Effect Model RMSE:", regularised_rmse, "\n")
```

## Regularised Movie + User Effect Model RMSE: 0.856558

```r
cat("\nRMSE Summary:\n")
```

```
##
## RMSE Summary:
```

```r
cat("Baseline Model:", baseline_rmse, "\n")
```

## Baseline Model: 1.059799

```r
cat("Movie Effect Model:", movie_effect_rmse, "\n")
```

## Movie Effect Model: 0.9427288

```r
cat("Movie + User Effect Model:", user_effect_rmse, "\n")
```

## Movie + User Effect Model: 0.8561705

```r
cat("Regularised Movie + User Effect Model:", regularised_rmse, "\n")
```

## Regularised Movie + User Effect Model: 0.856558

```r
validation_set <- validation_set %>%
    left_join(user_avgs, by = "userId") %>%
    left_join(movie_avgs, by = "movieId")

final_holdout_test <- final_holdout_test %>%
    left_join(user_avgs, by = "userId") %>%
    left_join(movie_avgs, by = "movieId")
```

## 7. Final Model

```r
# - Train the final model on the entire edx dataset (using the best-performing
# approach) - Predict movie ratings for the final_holdout_test set

# Based on RMSE results, we will use the Regularised Movie + User Effect Model
# as the final model

# Train the final model on the entire edx dataset (regularised movie + user
# effects)
lambda <- 5  # Regularisation parameter
movie_effects_final <- edx %>%
```

```r
    group_by(movieId) %>%
    summarise(b_i = sum(rating - mu)/(n() + lambda))  # Regularised movie effect

user_effects_final <- edx %>%
    left_join(movie_effects_final, by = "movieId") %>%
    group_by(userId) %>%
    summarise(b_u = sum(rating - mu - b_i)/(n() + lambda))  # Regularised user effect

# Predict movie ratings for the final_holdout_test set
predicted_ratings_final <- final_holdout_test %>%
    left_join(movie_effects_final, by = "movieId") %>%
    left_join(user_effects_final, by = "userId") %>%
    mutate(pred = mu + b_i + b_u) %>%
    pull(pred)

# Calculate RMSE for the final model on the final_holdout_test set
final_rmse <- RMSE(final_holdout_test$rating, predicted_ratings_final)

# Display the final RMSE
cat("Final Model RMSE on final_holdout_test set:", final_rmse, "\n")
```

```
## Final Model RMSE on final_holdout_test set: 0.8648177
```

## 8. Evaluation

```r
# - Compute the RMSE on the final_holdout_test set

# The RMSE for the final model on the final_holdout_test set has already been
# calculated in the previous step
cat("Final Model RMSE on final_holdout_test set:", final_rmse, "\n")
```

```
## Final Model RMSE on final_holdout_test set: 0.8648177
```

```r
# Evaluate the performance of the final model by comparing RMSE with baseline
# and other models
cat("\nModel Evaluation Summary:\n")
```

```
##
## Model Evaluation Summary:
```

```r
cat("Baseline Model RMSE:", baseline_rmse, "\n")
```

```
## Baseline Model RMSE: 1.059799
```

```r
cat("Movie Effect Model RMSE:", movie_effect_rmse, "\n")
```

```
## Movie Effect Model RMSE: 0.9427288
```

```
cat("Movie + User Effect Model RMSE:", user_effect_rmse, "\n")
```

## Movie + User Effect Model RMSE: 0.8561705

```
cat("Regularised Movie + User Effect Model RMSE (Final Model):", regularised_rmse,
    "\n")
```

## Regularised Movie + User Effect Model RMSE (Final Model): 0.856558

```
cat("Final Model RMSE on final_holdout_test set:", final_rmse, "\n")
```

## Final Model RMSE on final_holdout_test set: 0.8648177

# 9. Conclusion

```
##
## ================= Conclusion =================
##
## I developed multiple models to predict movie ratings using the MovieLens dataset.
##
## 1. Baseline Model:
##    The global average rating for all movies provided an RMSE of:  1.059799
##
## 2. Movie Effect Model:
##    Accounted for differences in movie ratings, which reduced the RMSE to:  0.9427288
##
## 3. Movie + User Effect Model:
##    Further reduced the RMSE to:  0.8561705
##
## 4. Regularised Movie + User Effect Model:
##    Best model with an RMSE of:  0.856558
##
## This final model was evaluated on the holdout test set, yielding an RMSE of:  0.8648177
##
## The Regularised Movie + User Effect Model provides a balance between model complexity
##  and generalisation,
##  resulting in the most accurate predictions.
##
## Overall, the final model outperforms the baseline and intermediate models,
##  demonstrating the importance of accounting for both movie and user effects, as well
##  as applying regularisation to prevent overfitting.
##  =================================================
```