

MovieLens

JP Bradley

2025-06-22

Introduction

In today’s digital landscape, from e-commerce platforms like Amazon to streaming services like Netflix, recommender systems are a critical component for enhancing user experience and driving engagement. They seek to solve the problem of information overload by intelligently filtering and predicting a user’s potential interest in items they have not yet encountered. This project directly engages with this challenge by building and evaluating a movie recommendation system based on the principles of collaborative filtering.

To accomplish this, we leverage the well-established **MovieLens** 10M dataset, a public benchmark dataset provided by the GroupLens research group. This rich dataset contains over 10 million ratings applied by approximately 71,000 users to nearly 10,000 unique movies. Each record provides crucial data points for analysis: unique identifiers for users and movies, a numerical rating on a scale of 0.5 to 5, a UNIX timestamp indicating when the rating was given, and associated metadata such as movie titles and genres. The fundamental challenge lies in the inherent sparsity of the data—most users have rated only a tiny fraction of the available movies, leaving a vast matrix of unknown preferences that our model must learn to predict.

The primary methodology employed in this analysis is collaborative filtering, a technique that makes automatic predictions about the interests of a user by collecting preferences from many other users. Unlike content-based methods, which would analyze the properties of a movie such as its genre or director, collaborative filtering relies solely on historical user-item interaction data. The underlying assumption is that users who have agreed in the past (e.g., gave similar ratings to the same movies) are likely to agree in the future. The model identifies these patterns to predict a user’s rating for a new item by finding a “neighborhood” of similar users or, more commonly, by modeling the intrinsic biases associated with each user (e.g., a user who tends to give high ratings) and each movie (e.g., a movie that is universally acclaimed). The structure of the **MovieLens** dataset, being a large matrix of user-movie ratings, makes it an ideal candidate for this approach.

Project Objective:

The principal objective of this project is to develop and validate a machine learning model that accurately predicts movie ratings. The success of the model is quantitatively measured by the Root Mean Squared Error (**RMSE**), and the primary goal is to minimize this error metric on a final, unseen hold-out validation set, denoted as `final_holdout_test`.

The **RMSE** is a standard metric for regression tasks that calculates the square root of the average of the squared differences between predicted and actual ratings. It is chosen for two key reasons: it heavily penalizes larger errors, making it sensitive to significant prediction misses, and its value is interpretable in the same units as the rating itself (i.e., “stars”). Therefore, a lower **RMSE** signifies a model whose predictions are, on average, closer to the true user ratings. Our methodology will follow an incremental approach, starting with a simple baseline model and systematically incorporating more complex factors—such as movie-specific and user-specific biases—to progressively improve predictive accuracy and achieve the lowest possible **RMSE** on the final hold-out set.

Project Overview

1. Data Acquisition & Preprocessing

1.1 Setting Up the R Environment

```
# Check if the 'tidyverse' package is installed; if not, install it from CRAN
if (!require(tidyverse))
  install.packages("tidyverse", repos = "http://cran.us.r-project.org")

# Check if the 'caret' package is installed; if not, install it from CRAN
if (!require(caret))
  install.packages("caret", repos = "http://cran.us.r-project.org")

# Check if the 'knitr' package is installed; if not, install it from CRAN
if(!require(knitr))
  install.packages("knitr", repos = "http://cran.us.r-project.org")

# Check if the 'skimr' package is installed; if not, install it from CRAN
if(!require(skimr)) install.packages("skimr", repos = "http://cran.us.r-project.org")
```

1.1.2 Installing Required Packages To ensure a fully reproducible and self-contained analysis, the script begins by programmatically managing its core dependencies, `tidyverse` and `caret`. For each package, the code first attempts to load it into the R session. If a package is not already installed on the system, the `require()` function fails, which in turn triggers an automatic installation from a specified Comprehensive R Archive Network (CRAN) repository. This conditional logic automates the environment setup, guaranteeing that the script can run on any machine with R installed without requiring manual pre-installation of these essential packages.

```
# Load tidyverse for data wrangling, visualization, and functional programming workflows
library(tidyverse)

# Load caret to support end-to-end machine learning pipelines including preprocessing,
↪ model tuning, and validation
library(caret)

# Load knitr for dynamic report generation and integration of R output in reproducible
↪ research documents
library(knitr)

# Load skimr to produce compact, well-formatted summaries of data frames for exploratory
↪ analysis and diagnostics
library(skimr)
```

1.1.3 Loading Required Packages The analysis and modeling for this project are conducted in R, leveraging the `tidyverse` and `caret` packages to establish a robust and reproducible data science pipeline. The `tidyverse` suite is instrumental in the data wrangling and exploratory analysis phases. Its packages, particularly `dplyr` and `ggplot2`, enable efficient data manipulation—such as joining the ratings and movies datasets and engineering new features—and the creation of insightful visualizations to understand underlying data distributions. For the modeling phase, the `caret` package provides a unified framework for the predictive modeling workflow. We utilize `caret` to perform a stratified split of the data into training and validation sets, to train various machine learning models with a consistent syntax, and to rigorously evaluate their performance through cross-validation, using the Root Mean Squared Error (RMSE) as the primary metric. This combination of tools provides a seamless transition from raw data exploration to the development and

validation of our final recommendation model.

1.2 Downloading and Extracting the Dataset

```
# Define the filename for the MovieLens 10M dataset zip file
dl <- "ml-10M100K.zip"

# Check if the file already exists in the current working directory
# If it doesn't, download the dataset zip file from the specified URL and save it with
↪ the given filename
if (!file.exists(dl))
  download.file("https://files.grouplens.org/datasets/movielens/ml-10m.zip", dl)
```

1.2.1 Downloading the Dataset To ensure a self-contained and reproducible workflow, the R script programmatically manages the data acquisition process. The script first checks if the dataset's compressed zip archive, `ml-10M.zip`, already exists in the local working directory. The download from the specified URL is initiated only if the file is not found. This conditional logic ensures that the large dataset is not downloaded redundantly upon subsequent executions of the script.

```
# Define the file paths for the ratings and movies data within the extracted folder
ratings_file <- "ml-10M100K/ratings.dat"
movies_file <- "ml-10M100K/movies.dat"

# Check if the ratings file exists; if not, extract it from the downloaded zip archive
if (!file.exists(ratings_file)) unzip(dl, ratings_file)

# Check if the movies file exists; if not, extract it from the downloaded zip archive
if (!file.exists(movies_file)) unzip(dl, movies_file)
```

1.2.2 Extracting the Dataset Following the data acquisition, the script programmatically prepares the raw data for loading into the R environment. Our analysis specifically requires the `ratings.dat` and `movies.dat` files, which are contained within the downloaded zip archive. To manage these files efficiently, the code first checks for the existence of each required file in the designated project directory. The `unzip` function is then called to selectively extract a specific file from the archive only if it is not already present. This conditional logic streamlines the data preparation process by avoiding redundant file extraction during subsequent script executions, ensuring the necessary data is readily available for analysis without unnecessary overhead.

1.3 Parsing the Ratings and Movies Data

```
# Read the ratings file line by line, splitting each line at '::' into a matrix-like
↪ structure
ratings <- as.data.frame(str_split(read_lines(ratings_file),
                                   fixed("::"), simplify = TRUE),
                  stringsAsFactors = FALSE)

# Rename the columns to meaningful variable names
colnames(ratings) <- c("userId", "movieId", "rating", "timestamp")

# Convert each column to its appropriate data type for analysis
ratings <- ratings %>%
  mutate(userId = as.integer(userId),
         movieId = as.integer(movieId),
         rating = as.numeric(rating),
```

```
timestamp = as.integer(timestamp))
```

1.3.1 Parsing the Ratings Data The script ingests the raw `ratings.dat` file, which uses a non-standard double-colon (`::`) delimiter. To handle this format, the data is first read as a vector of text lines, and each line is then parsed into separate columns using a string-splitting function. The resulting matrix is converted into a standard R data frame, to which we assign the descriptive column headers: `userId`, `movieId`, `rating`, and `timestamp`. In the final and critical step of this process, we perform data type coercion. The `userId`, `movieId`, and `timestamp` columns are converted to integers, and the `rating` column is converted to a numeric type to correctly handle fractional values like 3.5. This procedure transforms the raw text data into a clean, properly typed, and structured format that is essential for all subsequent computational and analytical tasks.

```
# Use kable() for a professional-looking preview of the data
knitr::kable(
  head(ratings),
  caption = "A preview of the first six rows of the ratings data."
)
```

DATA SUMMARY: Cleaned Ratings Data

Table 1: A preview of the first six rows of the ratings data.

userId	movieId	rating	timestamp
1	122	5	838985046
1	185	5	838983525
1	231	5	838983392
1	292	5	838983421
1	316	5	838983392
1	329	5	838983392

```
# Use skim() for a rich and clean statistical summary
skim(ratings)
```

Table 2: Data summary

Name	ratings
Number of rows	10000054
Number of columns	4
Column type frequency:	
numeric	4
Group variables	None

Variable type: numeric

skim_variable	n_missing	n_complete	mean	sd	p0	p25	p50	p75	p100	hist
userId	0	1	3.586986e+00	1.0585.34	1.0	18123	35740.5	53608	71567	
movieId	0	1	4.120290e+00	3.38.40	1.0	648	1834.0	3624	65133	
rating	0	1	3.510000e+00	1.06	0.5	3	4.0	4	5	

skim_variable	n_missing	n_complete	rate	mean	sd	p0	p25	p50	p75	p100	hist
timestamp	0	1	1.032606e+10	1.05963962e+09	65200990167658800354764811026749071231131736						

With ratings distributed across a 0.5–5.0 scale and user IDs spanning a large, heterogeneous population, the dataset supports scalable modeling of rating behavior and temporal shifts in consumption. Its structure provides a stable foundation for advancing hybrid recommender algorithms, behavioral clustering, and fairness-aware personalization strategies in collaborative filtering research

```
# Read each line of the movies file, split on the delimiter ":",
# and assemble into a data frame (strings stay as character type)
movies <- as.data.frame(
  str_split(
    read_lines(movies_file),
    fixed(":"),
    simplify = TRUE
  ),
  stringsAsFactors = FALSE
)

# Assign meaningful column names
colnames(movies) <- c("movieId", "title", "genres")

# Convert movieId from character to integer for proper joins and filtering
movies <- movies %>%
  mutate(movieId = as.integer(movieId))
```

1.3.2 Parsing the Movies Data In parallel with the ratings data, the script loads and processes the `movies.dat` file to create a structured lookup table for movie information. The process mirrors the handling of the ratings file, first ingesting the raw text lines and then parsing each line based on the double-colon (`::`) delimiter to separate the distinct data fields. This parsed data is converted into a data frame, and its columns are programmatically named `movieId`, `title`, and `genres`. To ensure data integrity and to enable correct matching with the ratings data, the `movieId` column is explicitly converted from a character string to an integer type. This structured movies table is essential, as it provides the critical mapping between a `movieId` and its corresponding title and genre information for our analysis.

```
# Use kable() for a professional-looking preview of the data
knitr::kable(
  head(movies),
  caption = "A preview of the first six rows of the movies data."
)
```

DATA SUMMARY: Cleaned Movies Data

Table 4: A preview of the first six rows of the movies data.

movieId	title	genres
1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
2	Jumanji (1995)	Adventure Children Fantasy
3	Grumpier Old Men (1995)	Comedy Romance
4	Waiting to Exhale (1995)	Comedy Drama Romance

movieId	title	genres
5	Father of the Bride Part II (1995)	Comedy
6	Heat (1995)	Action Crime Thriller

```
# Use skim() for a rich and clean statistical summary
skim(movies)
```

Table 5: Data summary

Name	movies
Number of rows	10681
Number of columns	3
Column type frequency:	
character	2
numeric	1
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
title	0	1	8	160	0	10680	0
genres	0	1	3	60	0	797	0

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100	hist
movieId	0	1	13120.52	17808.85	1	2755	5436	8713	65133	

The movies metadata infuses content-based context—via genre ontologies and item-level identifiers—into user behavior modeling. This integration enables the development of hybrid recommender systems that align user preferences with interpretable item features, enhancing algorithmic personalization, fairness, and scalability.

1.4 Merging Ratings and Movies Data

```
# Combine user ratings with movie metadata into one comprehensive table:
# - left_join preserves every rating record (even if some movies lack metadata)
# - appends title and genres columns to each rating
# Result: movielens has userId, movieId, rating, timestamp, title, and genres
movielens <- left_join(ratings, movies, by = "movieId")
```

To create a single, comprehensive dataset for analysis, the ratings and movies data frames are merged using a `left_join` operation. This function links the two tables by matching rows that share a common `movieId`. The result is a new, unified data frame named `movielens`, where each rating record is now enriched with the corresponding movie’s title and genres. This combined dataset is the foundational data structure for all subsequent exploratory analysis, feature engineering, and model training, as it contains all the necessary user, movie, and rating information in a single, tidy format.

```
# Use kable() for a professional-looking preview of the data
knitr::kable(
  head(movielens),
  caption = "A preview of the first six rows of the combined data."
)
```

DATA SUMMARY: Merged Data

Table 8: A preview of the first six rows of the combined data.

userId	movieId	rating	timestamp	title	genres
1	122	5	838985046	Boomerang (1992)	Comedy Romance
1	185	5	838983525	Net, The (1995)	Action Crime Thriller
1	231	5	838983392	Dumb & Dumber (1994)	Comedy
1	292	5	838983421	Outbreak (1995)	Action Drama Sci-Fi Thriller
1	316	5	838983392	Stargate (1994)	Action Adventure Sci-Fi
1	329	5	838983392	Star Trek: Generations (1994)	Action Adventure Drama Sci-Fi

```
# Use skim() for a rich and clean statistical summary
skim(movielens)
```

Table 9: Data summary

Name	movielens
Number of rows	10000054
Number of columns	6
Column type frequency:	
character	2
numeric	4
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
title	0	1	8	160	0	10676	0
genres	0	1	3	60	0	797	0

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100	hist
userId	0	1	3.586986e+00	1.0585.34	1.0	18123	35740.5	53608	71567	
movieId	0	1	4.120290e+00	1.0338.40	1.0	648	1834.0	3624	65133	
rating	0	1	3.510000e+00	1.06	0.5	3	4.0	4	5	
timestamp	0	1	1.032606e+10	1.05963962.78965200991	1.0	16765880	1003547648110	10267490712	1231131736	

The movielens dataset constitutes a temporally indexed, user-item interaction matrix enriched with item-level covariates (title, year, genres) and user-specific behavioral data (userId, rating, timestamp). Its schema supports advanced modeling tasks including matrix factorization with side information, dynamic collaborative filtering, and hierarchical genre-aware embeddings. By unifying observational feedback with content-based taxonomies, it facilitates research into high-dimensional preference estimation, temporal drift in item relevance, and fine-grained regularization strategies for recommender architectures.

1.5 Final Hold-out Test Set Preparation

```
# -----
# STEP 1: Create a reproducible 10% hold-out sample from the full movielens data
# -----
# Set a random seed so that anyone running this code will get the same split.
# For R 3.6, you need sample.kind = "Rounding" to replicate older behavior;
# if you're on R 3.5, just use set.seed(1).
set.seed(1, sample.kind = "Rounding")

# createDataPartition() returns row indices for a stratified sample of ratings.
# Here p = 0.1 means "take 10% of the rows" while preserving the rating distribution.
test_index <- createDataPartition(
  y      = movielens$rating, # the outcome we want to stratify by
  times  = 1,               # only one partition
  p      = 0.1,             # proportion in the hold-out set
  list   = FALSE           # return a vector, not a list
)

# Split the data:
# - edx      : 90% of the data, to train and tune our models
# - temp     : initial 10% hold-out, which we'll prune next
edx <- movielens[-test_index, ]
temp <- movielens[test_index, ]

# -----
# STEP 2: Guarantee that the final test set only contains users and movies seen
#         during training (so we aren't forced to predict on unseen data).
# -----
# semi_join(x, y, by) keeps only rows in x that have matching keys in y.
final_holdout_test <- temp %>%
  semi_join(edx, by = "movieId") %>% # drop any ratings for brand-new movies
  semi_join(edx, by = "userId")      # drop any ratings from brand-new users

# Any rows from temp that got dropped because of unseen user/movie:
removed <- anti_join(temp, final_holdout_test)

# Put those removed rows back into the edx (training) set to preserve all data.
edx <- bind_rows(edx, removed)

# -----
# STEP 3: Clean up workspace
# -----
# Remove files and objects we no longer need to free up memory
rm(dl, ratings, movies, test_index, temp, movielens, removed)
```

To rigorously evaluate the final model's performance, we partition the complete `movielens` dataset into a primary training set (`edx`) and a final hold-out validation set (`final_holdout_test`). To ensure this process is reproducible, we first set a random seed. Using the `createDataPartition` function, we perform an initial

90/10 stratified split based on the rating outcome, which preserves the rating distribution in both sets. However, a critical step is then taken to address a common challenge in recommendation systems: ensuring the validation set only contains users and movies that are also present in the training set. We achieve this by filtering the hold-out set using `semi_join` operations, which guarantees that every `userId` and `movieId` in the final `final_holdout_test` set also exists in the `edx` set. To maintain data integrity, any rows removed from the hold-out set during this process are identified and added back to the `edx` training set. Finally, all intermediate data objects are removed from the environment to ensure a clean workspace, leaving only the final, well-defined training and validation sets for the modeling phase.

Methodologies

1. **Data Wrangling and Cleaning** To begin with, key variables such as `userId`, `movieId`, `rating`, `timestamp`, and `genres` must be thoroughly parsed to ensure data integrity. The dataset should be examined for duplicates, missing values, and potentially corrupt entries. Anomalous ratings—those falling outside the valid range—should be identified and addressed. Using regular expressions, the movie title field can be parsed to extract the release year, while timestamps should be converted into standard date-time objects to facilitate the creation of time-based features (e.g., year, month, day). Furthermore, genre information ought to be transformed into a format suitable for modeling, such as multi-hot encoding or a normalized long-table format.
2. **Auditing for Bias** A comprehensive bias audit involves evaluating the distribution of rating counts at both the user and movie levels to assess potential participation imbalances. Comparative analyses should be conducted across genres, time periods, and popularity strata to identify disparities in coverage. The relationship between the number of ratings and average rating per item can reveal the presence of popularity bias. Additionally, one should assess representation gaps, focusing on users or movies with limited interaction history, such as obscure genres or newly introduced items. Differences in rating behavior across user frequency tiers—high-frequency versus low-frequency users—should also be analyzed to expose behavioral heterogeneity.
3. **Exploratory Trend Analysis** This phase entails a visual exploration of global rating distributions, with stratification by genre, temporal dimension, and user activity level. Temporal trends in rating behaviors can be captured via line plots segmented by rating year, or heatmaps mapped by genre over time. Analysis of user and movie activity trajectories can uncover patterns such as user churn or content consumption surges. To probe cultural or temporal shifts, average ratings can be compared across genres and release decades. Scatter plots serve to illustrate relationships between rating and movie age, as well as to explore the interaction between rating variance and sample size.
4. **Feature Engineering Preparation** For effective predictive modeling, several derived features should be constructed. Time-sensitive features such as `movie_age`, `rating_year`, and `user_tenure` enrich the temporal resolution of the dataset. Normalization techniques should be employed to adjust ratings by user or movie means to reduce idiosyncratic bias. Genre information may be encoded through one-hot vectors, multi-hot matrices, or embedded representations suitable for downstream algorithms. Users and movies may be binned according to their activity levels, enabling the generation of aggregate statistics per bin. Finally, statistical outliers—users with anomalous rating variance or consistent directional bias—should be flagged to inform data preprocessing or modeling strategies.