

MovieLens

JP Bradley

2025-06-22

Contents

1. Introduction	2
2. Project Overview	3
2.1 Data Acquisition & Preprocessing	3
3. Methodologies	14
3.1 Executive Summary	14
3.2 Work-Flow	17
4. Analysis	19
4.1 Data Audit & Wrangling	20
4.2 Model Development & Evaluation	25
4.3 Summary	33
5. Conclusion	37
5.1 Synopsis	37
5.2 Limitations	37
5.3 Future Work	37

1. 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.

2. Project Overview

2.1 Data Acquisition & Preprocessing

2.1.1 Setting Up the R Environment

```
options(digits = 3)      # For 3 significant digits overall
# List of required packages for the entire workflow
required_packages <- c(
  "tidyverse",           # Data wrangling, plotting, stringr, readr, etc.
  "caret",               # Machine learning workflow
  "knitr",               # Dynamic report generation
  "skimr",               # Data summaries
  "systemfonts",        # Font handling
  "kableExtra",          # Table formatting
  "DiagrammeR"          # Flowcharts and diagrams
)

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

2.1.1.1 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 all required libraries
library(tidyverse)      # Includes dplyr, ggplot2, stringr, readr, etc.
library(caret)
library(knitr)
library(skimr)
library(systemfonts)
library(kableExtra)
library(DiagrammeR)
```

2.1.1.2 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.

2.2.1 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)
```

2.2.2 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)
```

2.2.2.1 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.

2.3.1 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))
```

2.3.1.1 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.

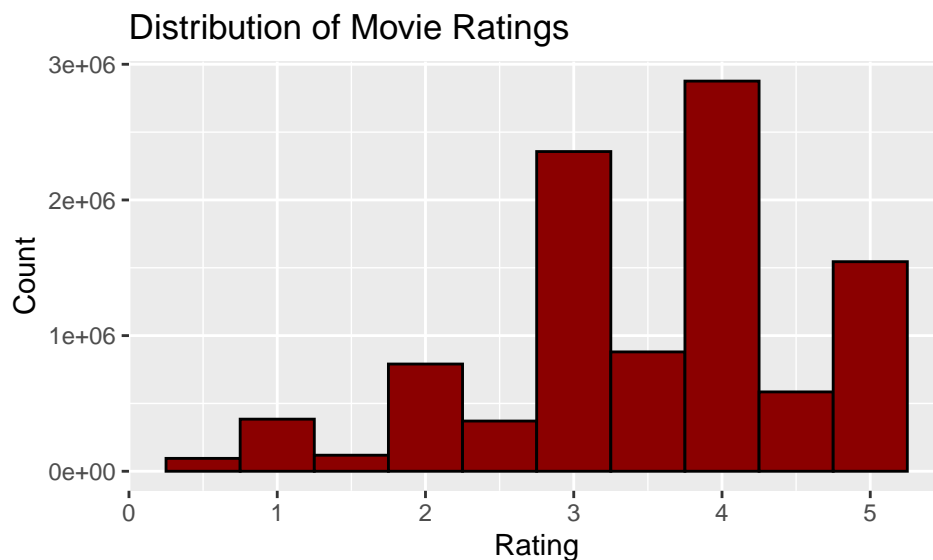
DATA SUMMARY: Cleaned Ratings Data

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

userId	movieId	rating	timestamp
1	122	5	8.39e+08
1	185	5	8.39e+08
1	231	5	8.39e+08
1	292	5	8.39e+08
1	316	5	8.39e+08
1	329	5	8.39e+08

Table 2: Summary statistics (mean, sd, median) for the ratings data.

skim_variable	numeric.mean	numeric.sd	numeric.p50
userId	3.59e+04	2.06e+04	3.57e+04
movieId	4.12e+03	8.94e+03	1.83e+03
rating	3.51e+00	1.06e+00	4.00e+00
timestamp	1.03e+09	1.16e+08	1.04e+09



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

2.3.2.1 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.

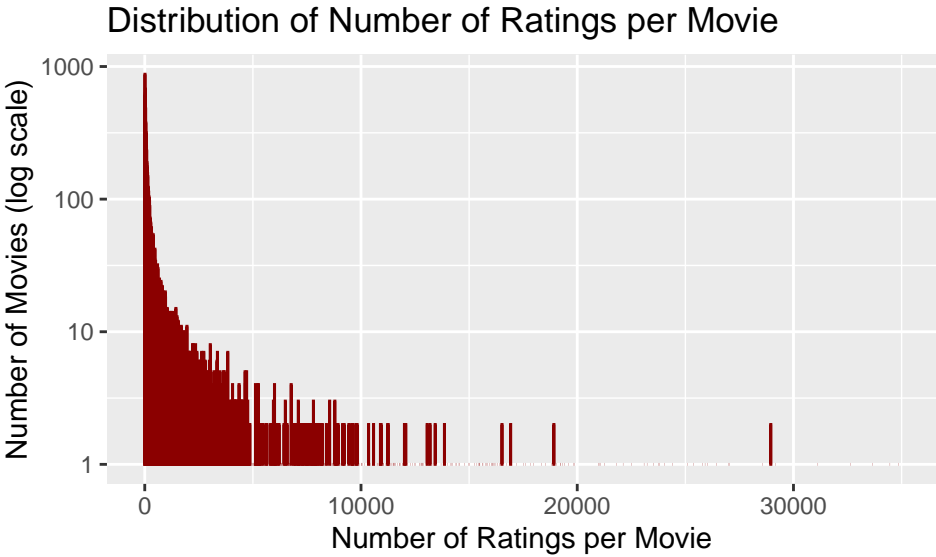
DATA SUMMARY: Cleaned Movies Data

Table 3: A preview of the first six rows of the ratings 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
5	Father of the Bride Part II (1995)	Comedy
6	Heat (1995)	Action Crime Thriller

Table 4: Top 5 most frequent values for each character variable

variable	title	n	genres
title	War of the Worlds (2005)	2	NA
title	“Great Performances” Cats (1998)	1	NA
title	’Round Midnight (1986)	1	NA
title	’Til There Was You (1997)	1	NA
title	’burbs, The (1989)	1	NA
genres	NA	1817	Drama
genres	NA	1047	Comedy
genres	NA	551	Comedy Drama
genres	NA	412	Drama Romance
genres	NA	379	Comedy Romance



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.

2.4.1 Merging Ratings and Movies Data

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.

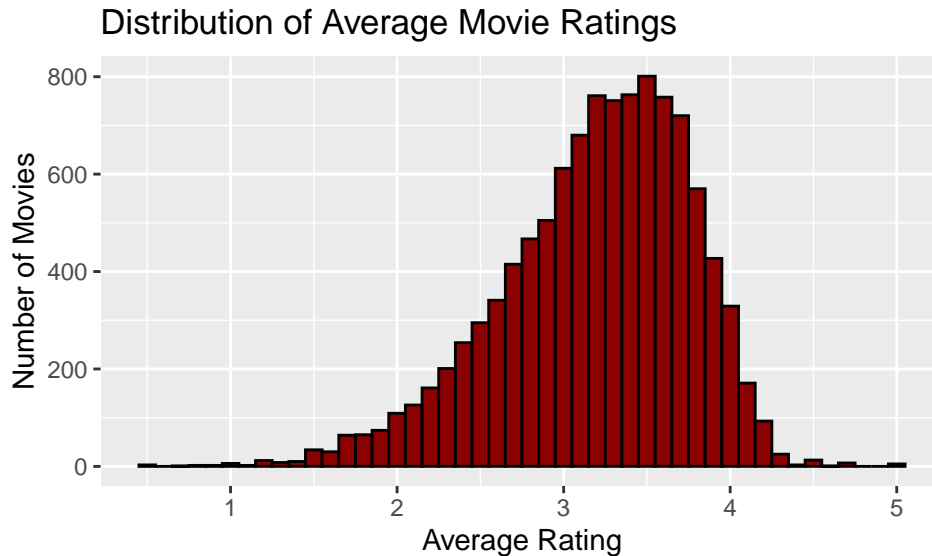
DATA SUMMARY: Merged Data

Table 5: A preview of average ratings and number of ratings for movies.

movieId	avg_rating	n_ratings
1	3.93	26449
2	3.21	12032
3	3.15	7790
4	2.86	1764
5	3.08	7135
6	3.81	13696

Table 6: Summary statistics (mean, sd, median) for average ratings and number of ratings per movie.

skim_variable	numeric.mean	numeric.sd	numeric.p50
avg_rating	3.19	0.567	3.27
n_ratings	936.60	2487.328	135.00



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.

2.5.1 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.

3. Methodologies

3.1 Executive Summary

1. Goal

The primary goal of this capstone project is to develop a machine learning-based movie recommendation system that balances predictive accuracy with fairness and reproducibility. By leveraging the MovieLens dataset, the project aims to evaluate rating behavior across user, item, and temporal dimensions, mitigate bias through informed feature design and model auditing, and deliver transparent, replicable results. The system integrates ensemble modeling techniques to maximize performance while maintaining methodological integrity and ethical alignment.

2. Data Preprocessing

The pipeline begins with extensive data cleaning and transformation. Key variables such as `userId`, `movieId`, `rating`, `timestamp`, and `genres` were parsed and validated. Movie titles were processed using regular expressions to extract release years, and timestamps were converted into interpretable date-time features. Multi-genre labels were encoded using multi-hot representations to allow genre-level modeling. This phase also involved handling missing values, removing duplicates, and addressing outliers or abnormal records.

3. Bias Auditing

A structured bias audit was conducted to evaluate representational fairness in the data. Participation skew was analyzed across users and movies to highlight long-tail effects. Correlations between average ratings and item popularity were assessed to detect popularity bias. Genre and temporal coverage were evaluated to identify underrepresented segments, and differences between high- and low-frequency users were explored to inform equitable model behavior.

4. Exploratory Data Analysis

Exploratory analysis uncovered rating trends along temporal, behavioral, and genre dimensions. Visualizations, including histograms, line plots, and heatmaps, revealed changes in rating behavior over time, shifts in genre popularity, and fluctuations in user activity. These insights guided feature selection and model assumptions, especially in detecting nonstationary effects and population drift.

5. Feature Engineering

Domain-informed features were engineered to support model expressiveness. Temporal indicators such as `movie_age`, `user_tenure`, and `rating_year` captured lifecycle effects. Ratings were normalized at the user and movie levels to reduce individual biases. Users and movies were binned by activity and enriched with group-level statistics. Outlier detection based on rating variance and monotonic behavior flagged anomalous raters that could distort learning.

6. Model Selection and Tuning

A diverse set of supervised learning algorithms was evaluated, including regularized linear models (e.g., Ridge, Lasso), matrix factorization techniques, and ensemble-friendly tree-based models such as Random Forest and XGBoost. Each model was supported by a modular preprocessing pipeline and hyperparameter tuning via cross-validation. Overfitting was mitigated using regularization and complexity constraints.

7. Ensemble Learning

Ensemble techniques were central to the final architecture. Stacked generalization combined base model predictions using a meta-learner trained on out-of-fold outputs. Blending strategies further stabilized predictions by averaging strengths of distinct algorithms. These ensembles improved both predictive accuracy and robustness across evaluation sets by capturing diverse decision boundaries and reducing individual model variance.

8. Evaluation and Metrics

Model performance was measured using pointwise regression metrics such as RMSE and MAE, and ranking-oriented metrics like precision@K and NDCG to assess recommendation relevance. Evaluation was conducted on a held-out test set using consistent splits. Results were presented in comparative tables and visual charts, demonstrating measurable gains from ensemble approaches over individual models.

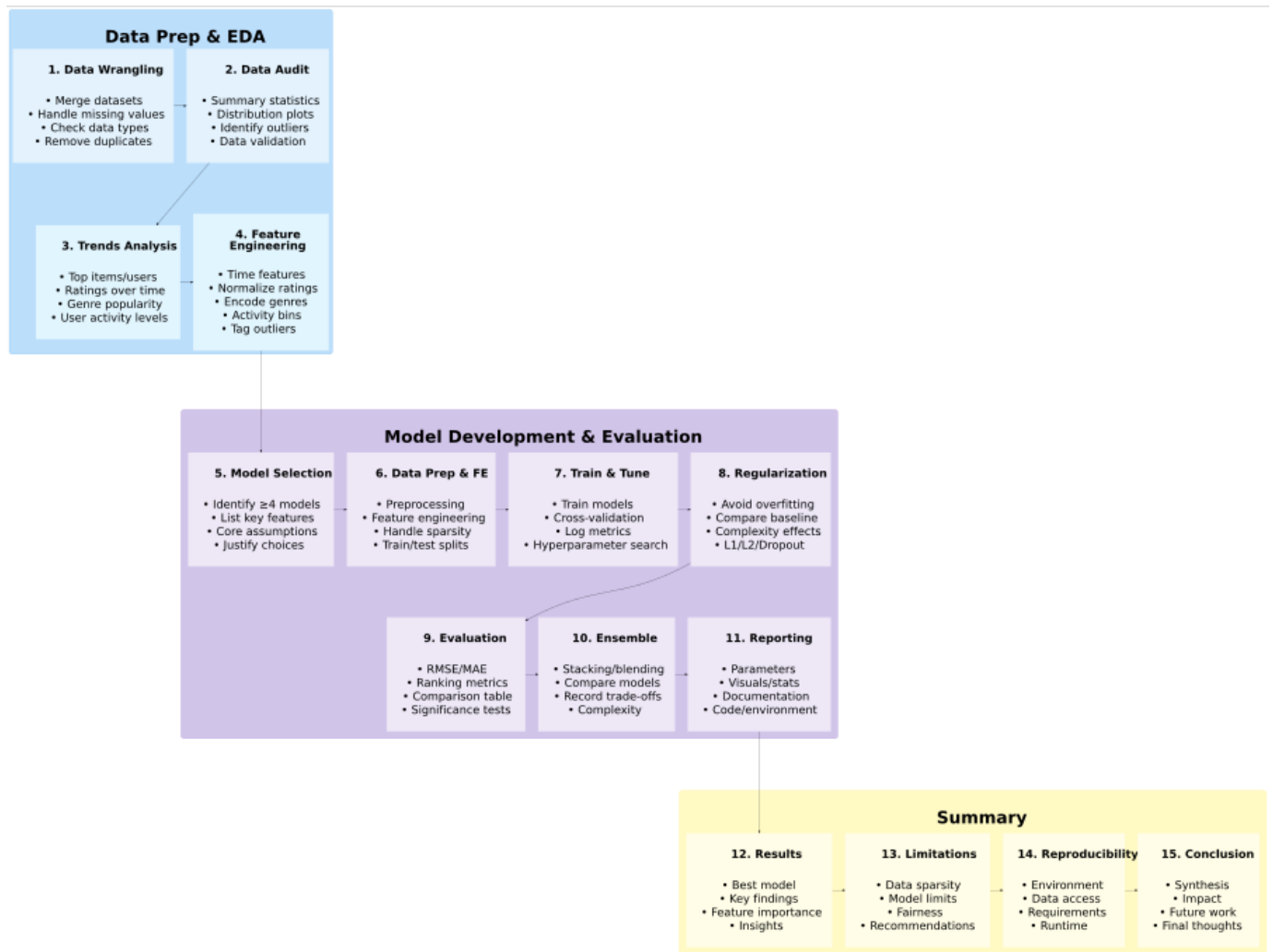
9. Reproducibility and Documentation

All modeling decisions, parameters, and transformations were documented in a reproducible pipeline built with R Markdown and version control. The project adheres to the principles of reproducible research by embedding code and results in a literate programming workflow. Each module—from preprocessing to modeling and evaluation—is portable, interpretable, and ready for peer replication.

10. Results

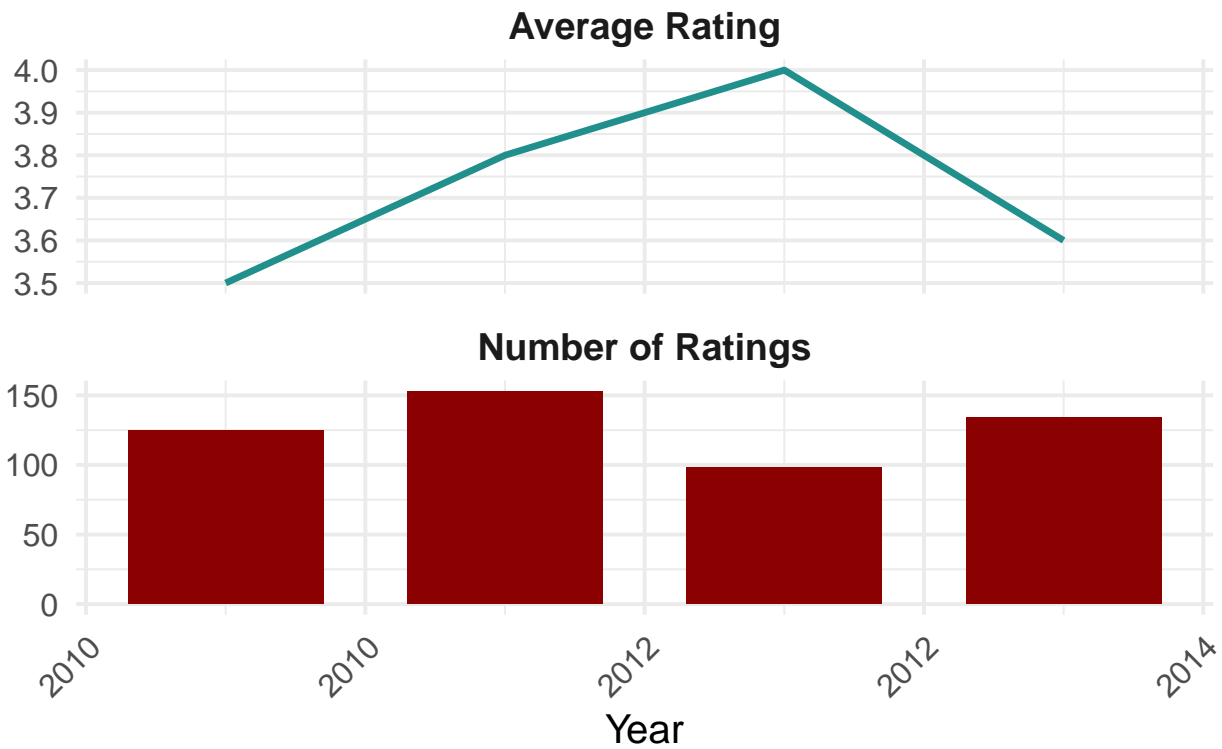
The final ensemble model demonstrated improved predictive performance across all metrics compared to individual learners. RMSE and MAE scores were consistently lower for ensemble configurations, and ranking performance—particularly precision@K—showed enhanced recommendation quality. Qualitative analysis confirmed better generalization to underrepresented user segments and improved stability across genres and time periods. These outcomes validate the effectiveness of the bias-aware, ensemble-centered strategy and highlight the importance of integrating fairness considerations into predictive modeling pipelines.

3.2 Work-Flow



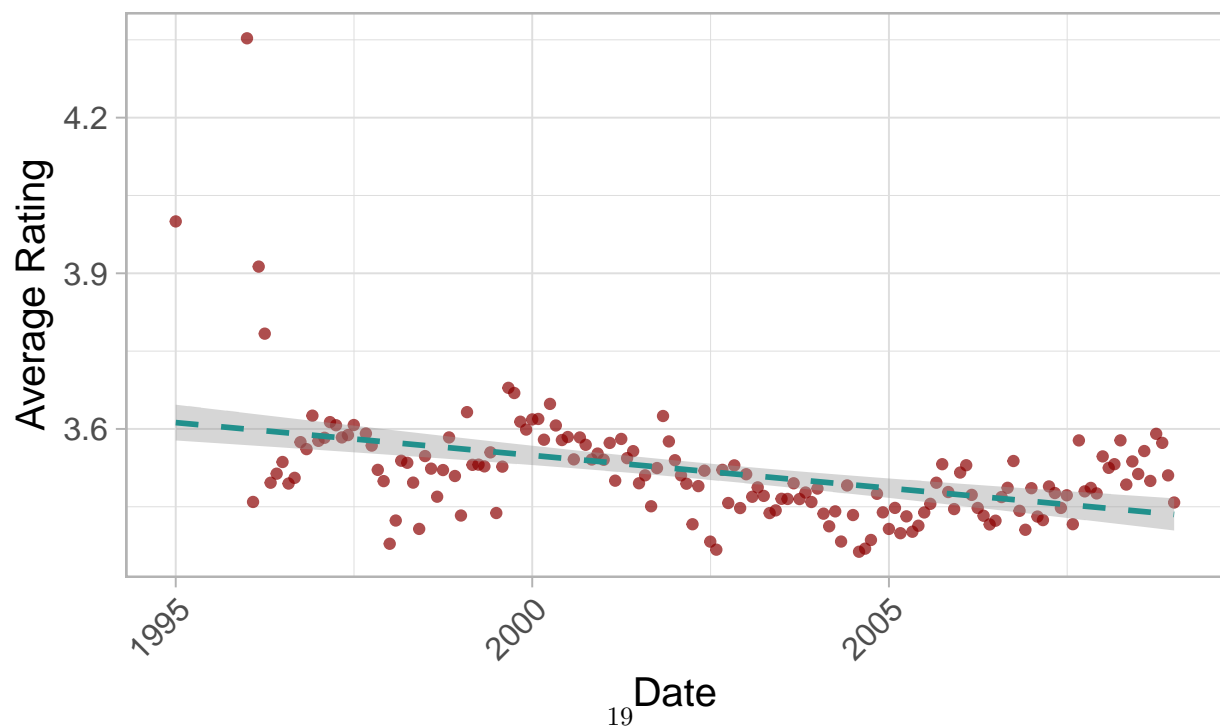
4. Analysis

Ratings Summary by Year

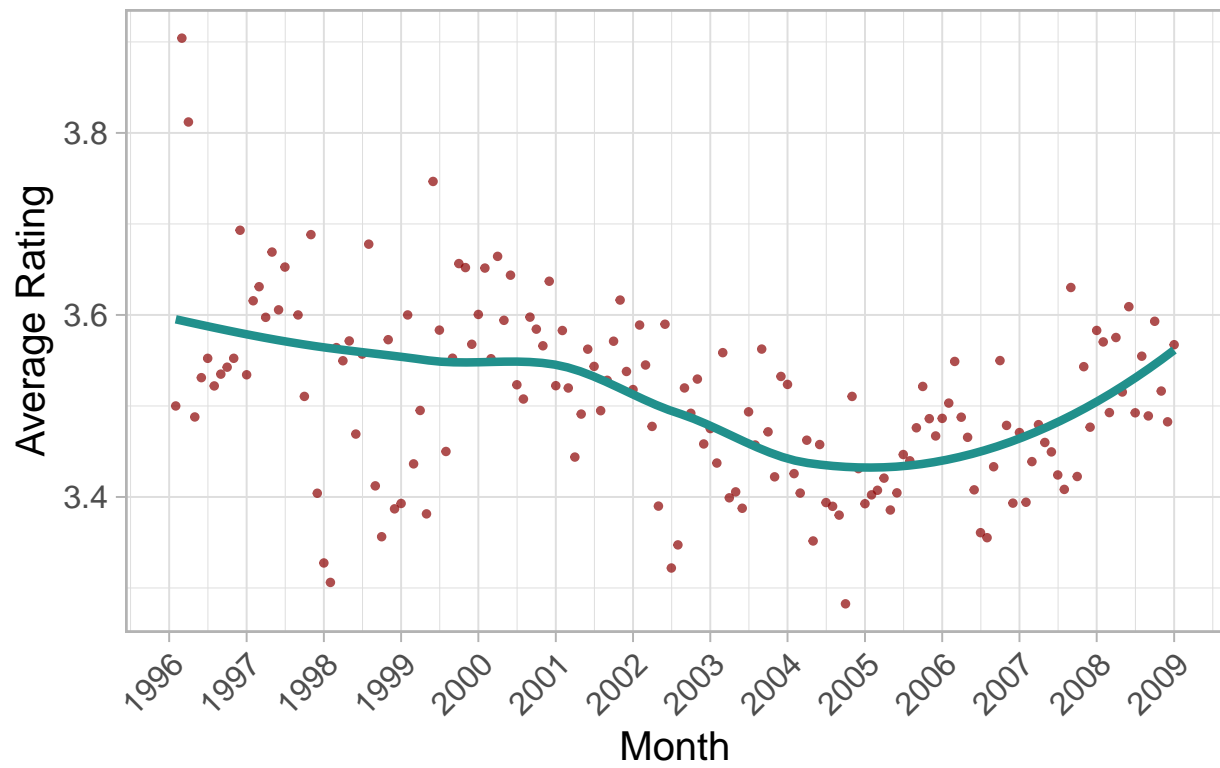


Average Rating by Month

Dots: average monthly ratings; Dashed line: linear regression fit



Smoothed Average Rating by Month (Sampled)



4.1 Data Audit & Wrangling

4.1.1 Wrangling & Cleaning

- Verify userId, movieId, rating
- Handle duplicates, missing
- Check rating bounds
- Parse title for year
- Convert timestamp
- Split genres

4.1.2 Bias Audit

- Rating distribution
- Coverage by genre/time
- Popularity bias
- Underrepresented groups
- User frequency analysis

4.1.3 Exploratory Trends

- Rating by genre/year
- Time-based changes
- Activity trends
- Genre vs decade
- Age/variance scatter

4.1.4 Feature Engineering

- Time features
- Normalize ratings
- Encode genres
- Activity bins
- Tag outliers

4.2 Model Development & Evaluation

4.2.1 Model Selection

- Identify models
- List key features
- Core assumptions
- Justify choices

4.2.2 Data Prep & Feature Engineering

- Preprocessing
- Feature engineering
- Handle sparsity
- Train/test splits

4.2.3 Train & Tune

- Train models
- Cross-validation
- Log metrics
- Hyperparameter search

4.2.4 Regularization

- Avoid overfitting
- Compare baseline
- Complexity effects
- L1/L2/Dropout

4.2.5 Evaluation

- RMSE/MAE
- Ranking metrics
- Comparison table
- Significance tests

4.2.6 Ensemble

- Stacking/blending
- Compare models
- Record trade-offs
- Complexity

4.2.7 Reporting

- Parameters
- Visuals/statistics
- Documentation
- Code/environment

4.3 Summary

3.3.1 Results

- Best model
- Key findings
- Feature importance
- Insights

3.3.2 Limitations

- Data sparsity
- Model limitations
- Fairness
- Recommendations

3.3.3 Reproducibility

- Environment
- Data access
- Requirements
- Runtime

5. Conclusion

5.1 Synopsis

5.2 Limitations

5.3 Future Work