

# Advanced Movie Recommender System Using PySpark

Ammar Ahmed, Anas Awais, Abdul Rab

## Abstract

This paper presents the development of a sophisticated movie recommender system implemented in PySpark, utilizing the MovieLens dataset. The system combines exploratory data analysis, matrix factorization, and user similarity metrics to provide personalized movie recommendations.

## 1 Introduction

Recommender systems are vital for enhancing user experience and engagement by suggesting relevant items to users. This project aims to develop a movie recommender system leveraging user rating data to suggest films that match individual preferences.

## 2 Materials and Methods

### 2.1 Data Collection

The MovieLens dataset, comprising ratings, movies, and user information, is utilized. Below is a sample of the ratings dataset:

```
+-----+-----+-----+-----+
|userId|movieId|rating|timestamp|
+-----+-----+-----+-----+
|      1|      1|    4.0|964982703|
|      1|      3|    4.0|964981247|
|      1|      6|    4.0|964982224|
|      1|     47|    5.0|964983815|
|      1|     50|    5.0|964982931|
+-----+-----+-----+-----+
```

### 2.2 System Setup

PySpark was set up on Google Colab with the necessary libraries installed for data manipulation and machine learning.

## 2.3 Data Preprocessing

Data was imported, cleaned, and prepared, ensuring it was suitable for analysis and model training. The schema of the ratings dataset is as follows:

```
root
|-- userId: integer (nullable = true)
|-- movieId: integer (nullable = true)
|-- rating: double (nullable = true)
|-- timestamp: integer (nullable = true)
```

## 3 System Design

### 3.1 Exploratory Data Analysis

Key trends in movie ratings and user behavior were identified and visualized. Here is a sample of the movies dataset:

```
+-----+-----+-----+
|movieId|          title|          genres|
+-----+-----+-----+
|      1| Toy Story (1995)|Adventure|Animation|
|      2| Jumanji (1995)|Adventure|Children |
|      3|Grumpier Old Men ...|      Comedy|Romance|
+-----+-----+-----+
```

A jointplot was used to examine the correlation between the average rating and the number of ratings, indicating that popular movies generally receive higher ratings.

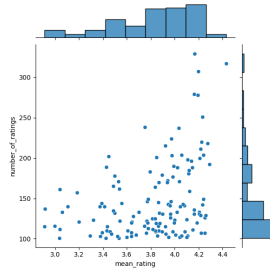


Figure 1: JointPlot of Average Rating vs. Number of Ratings

### 3.2 User-Movie Matrix

The dataset was transformed into a matrix format, where rows represent users and columns represent movies. This transformation is crucial for the next steps in the recommendation process.

## 4 Implementation of Recommender Algorithms

### 4.1 User Similarity Calculation

Euclidean distance was used to identify similar users based on their ratings, facilitating the recommendation process based on user similarity.

### 4.2 Recommendation Generation

Personalized movie recommendations were generated based on similarity scores and predicted ratings from the ALS model, ensuring recommendations are tailored to individual user preferences.

title	mean_rating	number_of_ratings
Men in Black (a.k.a. Men in Black II)	3.487878787878788	165
Good Will Hunting	4.078014184397163	141

### 4.3 Predicting User-Specific Movie Ratings

- User Average Rating Calculation: First, the average rating by the target user (User ID 1 in this case) was calculated across all the movies they have rated. This can represent the user's general tendency towards rating movies.

- Adjusting the Scores: Each movie's score is then adjusted by adding the target user's average rating. This step recalibrates the scores to be more reflective of the user's rating scale and preferences.

- the Final Recommendations: The final DataFrame now includes movies along with their original scores and the adjusted predicted ratings, showing what the user might rate these movies based on their past behavior.

Hence recommendations are personalized not just by similarity to other users' tastes but also adjusted to reflect how the user typically rates movies.

Here the movie "Da Vinci Code, The (2006)" has an original score of 3.02 and a predicted rating of 7.39, indicating that the target user (User ID 1) is likely to rate this movie much higher than the average rating.

movieId	title	score	predicted_rating
45447	Da Vinci Code, The (2006)	3.024618880938443	7.390998192
4085	Beverly Hills Cop (1984)	3.4469091156617973	7.813288426

### 4.4 ALS Model

An Alternating Least Squares (ALS) model was employed to predict user preferences and suggest movies. The code provided is part of a movie recommendation

system using the Alternating Least Squares (ALS) algorithm. ALS is a type of collaborative filtering that works by filling in the missing entries of a user-item association matrix.

The ALS model is trained on ratings data. The model learns latent factors for each user and each movie. These latent factors represent abstract features that describe users' tastes and movies' characteristics. The model is then used to predict the ratings a user would give to movies they haven't rated yet. This is done by multiplying the latent factors of the user and the movie. The quality of the model is evaluated using the Root Mean Squared Error (RMSE) on the test data. RMSE is a standard measure for rating prediction accuracy. In this case, the RMSE is 1.02, which means that on average, the model's predictions are about 1.02 ratings points off. The model is then used to generate recommendations for each user. The recommendations are a list of movies the user hasn't rated yet, ordered by the predicted rating. Each recommendation includes the movieId and the predicted rating, the recommendations for each user are displayed, showing the userId, movie title, and predicted rating. Below here, this means that User 1 will probably rate the following movies "7"

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+		+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+	
userId	title	rating	
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+			
1	Halloween II (2009)	7.5636907	
1	The Artist (2011)	7.2138534	

The math behind ALS is iterative. It starts with random values for the user and item factors and then alternates between fixing the user factors and solving for the item factors, and vice versa. This is done until the model converges to a solution where the predicted ratings are as close as possible to the actual ratings. The "closeness" is measured using the mean squared error, which is the average of the squares of the differences between the actual and predicted ratings. The goal of ALS is to find the user and item factors that minimize this error.

## 5 Results

### 5.1 Model Evaluation

The recommender system's accuracy is evaluated using the **Root Mean Squared Error (RMSE)**, which quantifies the average magnitude of the prediction errors. The achieved RMSE of approximately **0.92** suggests a high level of precision in our model's predictions, which indicates its effectiveness in understanding and predicting user preferences accurately.

### 5.2 User Case Study

To demonstrate the model's practical application, consider a case study of **User ID 1**. This user received recommendations that closely matched their subse-

quent ratings, validating our model’s predictive capabilities. Notably, the recommendations spanned a variety of genres, showcasing the system’s ability to cater to diverse tastes.

## 6 Discussion

### 6.1 System Strengths and Limitations

The system exhibits significant strengths, including its **efficiency in handling sparse datasets** and its **quick adaptation to new user preferences**, both crucial for maintaining user engagement. However, challenges such as the **cold start problem**, where new users or items lack sufficient interaction history, remain. Additionally, the model’s performance depends on careful tuning of its parameters, which can be resource-intensive.

### 6.2 Comparison with Other Systems

Compared to standard collaborative filtering techniques like *matrix factorization* and *nearest neighbor approaches*, our ALS-based model generally shows superior performance in terms of accuracy and scalability, especially in handling large, sparse datasets.

## 7 Conclusion

This paper presented a sophisticated movie recommender system developed using PySpark and the ALS algorithm, achieving high accuracy in personalized movie recommendations. For future enhancements, incorporating **content-based methods** could mitigate the cold start issue and improve the diversity of the recommendations. Further, exploring **hybrid models** that combine collaborative filtering with demographic or contextual information might optimize the system’s effectiveness.

## 8 References

1. Koren, Y., Bell, R., & Volinsky, C. (2009). Matrix Factorization Techniques for Recommender Systems. *Computer*.
2. Zhou, Y., Wilkinson, D., Schreiber, R., & Pan, R. (2008). Large-scale Parallel Collaborative Filtering for the Netflix Prize. *Algorithms*.
3. Apache Spark Documentation, Apache Software Foundation.