
Machine Learning Report

Akash Bahri
College of Engineering
Northeastern University
Toronto, ON
Bahri.a@northeastern.edu

Abstract

In this project, a **movie recommendation system** was developed using the **Item-Item Collaborative Filtering** algorithm based on **Cosine Similarity**. The system predicts the most relevant movies for users by analyzing their past ratings and movie similarities. An **interactive web application** was built to provide personalized movie recommendations to users. The backend was implemented using **Flask**, while the frontend was created with **React**. Challenges included handling large datasets and ensuring efficient performance for fast recommendations.

1 Introduction

A recommendation system is a software tool that helps suggest relevant items to users based on their preferences or past behavior. In this project, a movie recommendation system is built using Item-Item Collaborative Filtering, which analyzes the similarity between movies based on user ratings. The system predicts and recommends movies that a user may like, improving their experience by suggesting movies they might not have encountered otherwise. The project also includes building an interactive web application to allow users to input their preferences and receive personalized movie recommendations.

Objective: The goal of this project is to develop a recommendation system using Item-Item Collaborative Filtering to suggest movies to users based on their ratings and the similarity between movies. (similar to Netflix)

Scope: This system includes:

- **Jupyter notebook** for data processing and building the model
- A **backend** for processing recommendations using **Flask**.
- A **frontend** developed with **React** to display recommendations in an interactive manner.

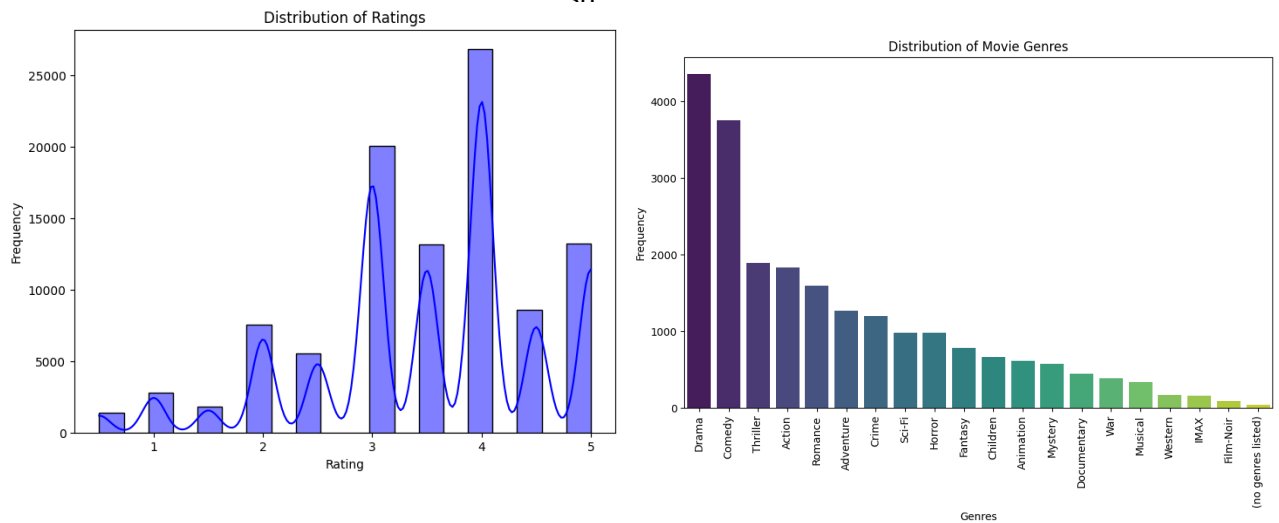
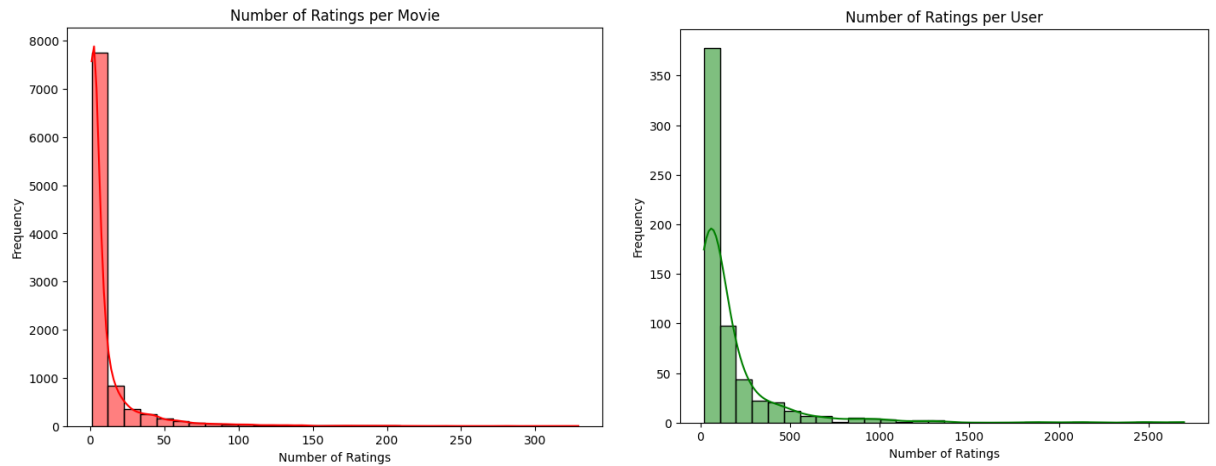
2 Datasets

Description:

- The dataset contains movie ratings, user ratings, and movie metadata (including movie titles and genres).
- MovieLens dataset was used, which includes ratings and metadata for movies.

Dataset Characteristics:

- **Size:** 100,000 ratings, 9,000 movies, 600 users.
- **Attributes:** Movie ID, title, genres, user ratings.
- **Challenges:** Sparse data (most users have rated a small portion of movies), handling missing values.



3 Methodology

Item-Item Collaborative Filtering

The recommendation system implemented in this project uses **Item-Item Collaborative Filtering**, a popular algorithm in recommendation systems. This approach works by recommending items that are similar to those the user has previously rated or interacted with. The main idea is to find the **similarity** between items (movies, in this case) and suggest those that are most similar to the movies the user has rated highly.

61 Cosine Similarity

62 To determine the similarity between items, we use **Cosine Similarity**, which measures the
63 cosine of the angle between two vectors in a multidimensional space. The formula for
64 **Cosine Similarity** between two items AAA and BBB is given by:

$$\text{Cosine Similarity} = \frac{A \cdot B}{\|A\| \|B\|}$$

Where:

- $A \cdot B$ is the dot product of vectors A and B ,
- $\|A\|$ and $\|B\|$ are the magnitudes (norms) of the vectors.

66 This metric produces values between **0** and **1**:

- 67 • A **value of 1** means the two items are identical (i.e., their vectors are pointing in the
68 same direction).
- 69 • A **value of 0** means the items are completely dissimilar.

70 In our system, **item-item similarity** is calculated for all movie pairs using their **ratings** from
71 all users. Once the similarity between items is computed, the next step is to predict ratings
72 for items that a user hasn't rated yet.

73 Data Preprocessing

74 The dataset used in this project is the **MovieLens dataset**, which includes ratings from users
75 for various movies. The data preprocessing steps include:

76 1. Handling Missing Data:

- 77 ○ The **user-item interaction matrix** contains missing values where users
78 haven't rated some movies. These missing values are handled by using
79 **zeros** in the matrix, which means unrated movies are considered as "not
80 rated yet."

81 2. Data Transformation:

- 82 ○ The **MovieLens dataset** is in CSV format. We use **pandas** to load the
83 dataset and transform it into the required **user-item matrix** format, where
84 each row represents a user, and each column represents a movie. The entries
85 in the matrix are the ratings given by users to movies, with **zeros** indicating
86 unrated movies.

87 3. Matrix Creation:

- 88 ○ We create two main matrices:
 - 89 1. **User-Item Matrix:** This matrix represents user preferences. Each
90 entry M_{ij} in the matrix represents the rating given by
91 user i to movie j .
 - 92 2. **Item-Item Similarity Matrix:** This matrix stores the **cosine**
93 **similarity** between each pair of movies. It helps in identifying
94 which movies are similar to each other, so we can recommend

95 similar movies to users.

96 Item-Item Similarity Calculation

97 The similarity between items is calculated using **Cosine Similarity**. For each movie *jjj*, the
98 system computes the similarity between it and all other movies based on user ratings. The
99 resulting similarity values are stored in the **item-item similarity matrix**, where each entry
100 S_{ij} represents the similarity between movie *iii* and movie *jjj*.

101 Once the similarity matrix is computed, the **top-N most similar movies** are identified for
102 each movie. These similar movies are then used to predict the rating a user would give to
103 unrated movies based on their ratings for the similar movies.

104 Rating Prediction

105 For each user, the system predicts the rating for movies that the user has not yet rated. The
106 predicted rating R_{ui} for a user *uuu* and an unrated movie *iii* is computed as the
107 **weighted sum** of the ratings for similar movies, weighted by their **similarity** to movie *iii*.
108 Specifically, the prediction is calculated as:

$$R_{ui} = \frac{\sum_{j \in \text{similar to } i} S_{ij} \cdot R_{uj}}{\sum_{j \in \text{similar to } i} |S_{ij}|}$$

Where:

- R_{uj} is the rating given by user *u* to movie *j*,
- S_{ij} is the similarity between movies *i* and *j*,
- The sum is taken over all movies *j* that are similar to movie *i*.

109

110

111 4 System Design and Implementation

112 Backend Implementation (Flask)

113 The backend of the system is built using **Flask**, a lightweight web framework for Python.
114 The backend handles the following:

115 1. API for Movie Recommendations:

- 116 ○ This API receives a **user ID** as input and returns a list of **top-N**
117 **recommended movies**. The recommendation is based on the **item-item**
118 **similarity** and the **predicted ratings**.

119 2. API for Movie Posters:

- 120 ○ The system fetches **movie posters** from the **TMDB API** using the **movie**
121 **ID**. It retrieves the **poster URL** and returns it along with the movie title and
122 genres.

123 3. Data Storage:

- 124 ○ The **user-item matrix** and **item similarity matrix** are stored in **Pickle files**

125 for fast retrieval during the recommendation process. This minimizes the
126 need for recalculating the similarity matrix every time a recommendation
127 request is made.

128 Frontend Development (React)

129 The frontend is built using **React** to create an interactive user interface:

130 1. User Input:

- 131 ○ Users can input their **user ID**, which is sent to the backend to retrieve
132 personalized movie recommendations.

133 2. Movie Display:

- 134 ○ The frontend displays the movie title, genres, predicted rating, and the
135 movie poster for each recommended movie. The UI is styled using **CSS** to
136 ensure that it is responsive and aesthetically pleasing, with a 3-column grid
137 layout for movie cards.

138 3. Interaction:

- 139 ○ The user can interact with the app by clicking the "**Get**
140 **Recommendations**" button to fetch movie suggestions based on their user
141 ID. The recommended movies are displayed dynamically, along with
142 relevant metadata.

143 5 Challenges and Solutions

144 Data Challenges

- 145 • **Sparsity:** The user-item matrix was sparse, with many missing ratings. This made it
146 difficult to find enough similar movies for accurate recommendations.
- 147 • **Cold Start Problem:** New users or movies that have few ratings had no data to base
148 predictions on. A fallback strategy was used for handling new users and movies.

149

150 Performance Optimization

- 151 • **Memory Usage:** The large size of the dataset required optimization to prevent
152 memory issues. The matrices were precomputed to speed up the recommendation
153 process.
- 154 • **Speed of Recommendations:** The recommendation process was optimized by
155 storing precomputed item similarities, reducing the time needed to make
156 predictions.

157 6 Results and Evaluation

158

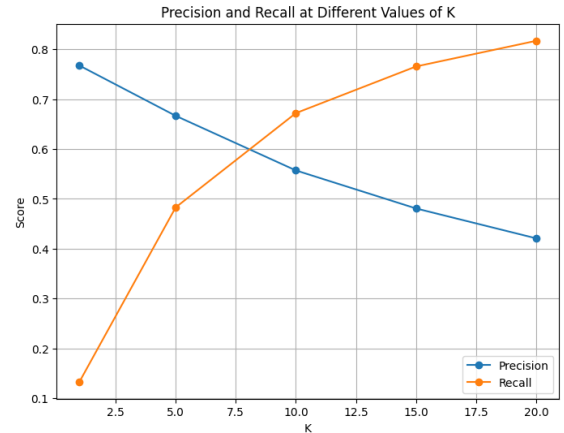
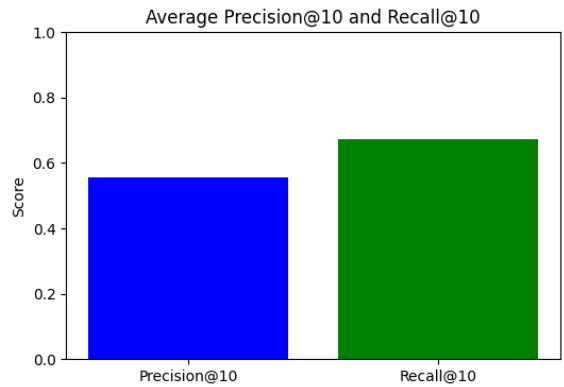
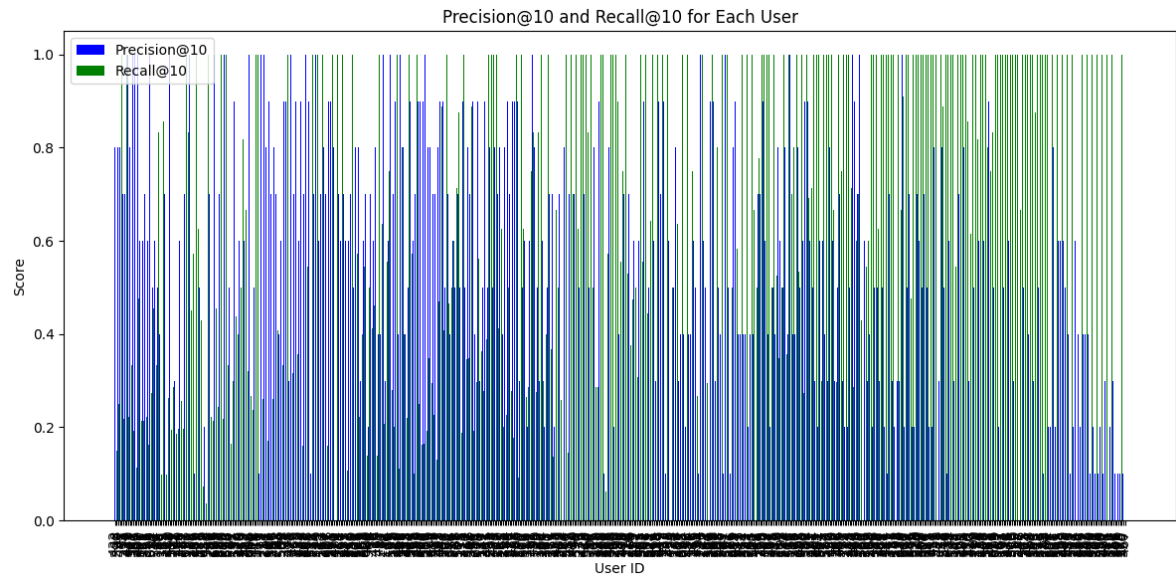
159 Performance Evaluation

160 The recommendation system was evaluated using metrics such as **Precision** and **Recall**. For
161 instance, **Precision@10** was used to measure how accurate the top-10 recommendations
162 were for each user. The system performed well, providing relevant movie recommendations

163 based on user preferences.

164 **Evaluation Metrics**

165 The system's performance was evaluated using **Precision@10** and **Recall@10** metrics. The
166 results showed that the recommendation system was able to suggest relevant movies to users
167 effectively.



168

169

170 **Sample Recommendations**

171 Here are some examples of movie recommendations for user 1:

- 172 • **Movie 1: Prophecy, The (1995)** (Predicted Rating: 4.9)
- 173 • **Movie 2: Thirty-Two Short Films About Glenn Gould (1993)** (Predicted
- 174 Rating: 4.8)
- 175 • **Movie 3: Nadja (1994)** (Predicted Rating: 4.8)


Movie Recommendations

Get Recommendations

Prophecy, The (1995)

Predicted Rating:
4.911951724985702

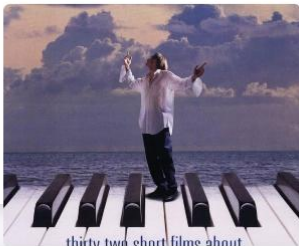
Genres: Fantasy, Horror, Mystery



Thirty-Two Short Films About Glenn Gould (1993)

Predicted Rating:
4.825448306127541


Genres: Drama, Musical



Nadja (1994)

Predicted Rating:
4.815130734964909


Genres: Drama



Tank Girl (1995)

Predicted Rating:
4.7894077135828725


Genres: Action, Comedy, Sci-Fi



Body Snatchers (1993)

Predicted Rating:
4.712326849662078


Genres: Horror, Sci-Fi, Thriller



Bhaji on the Beach (1993)

Predicted Rating:
4.711406962646358

Genres: Comedy, Drama



176

177

178 **7 Conclusion**

179 This project demonstrates the development of a movie recommendation system using **Item-**
180 **Item Collaborative Filtering** and **Cosine Similarity**. The backend was built using **Flask**,
181 and the frontend was developed with **React** to provide an interactive user interface. The
182 system was able to generate personalized movie recommendations based on user ratings and
183 item similarities. Future work could include addressing the **cold start problem**, improving
184 the **speed of recommendations**, and exploring hybrid models that combine multiple
185 recommendation techniques.

186

187 **8 References**

- 188 1. MovieLens Dataset - <https://grouplens.org/datasets/movielens/>
- 189 2. Cosine Similarity Explanation - https://en.wikipedia.org/wiki/Cosine_similarity
- 190 3. TMDb API Documentation - <https://www.themoviedb.org/documentation/api>
- 191 4. Flask Documentation - <https://flask.palletsprojects.com/>
- 192 5. React Documentation - <https://reactjs.org/docs/>
- 193 6. OpenAI. ChatGPT: Language Model 4o for Assistance: <https://www.openai.com/chatgpt>