
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 four distinct algorithms:

**Item-Item Collaborative Filtering,
Graph-based Hybrid Recommendation System (GHR),
Content-Based Filtering,
Popularity-Based Recommendation.**

The system accurately predicts personalized movie recommendations by analyzing users' historical ratings, movie attributes, and interaction patterns. An interactive, user-friendly web application was built using Flask for backend APIs and React for the frontend, providing seamless and personalized user experiences. Major challenges addressed included efficiently managing large-scale datasets, optimizing recommendation performance, and ensuring low-latency responses for a smooth user interaction.

1 Introduction

Recommendation systems play a crucial role in enhancing user experience on digital platforms by providing personalized content, thereby increasing user engagement and satisfaction. This project presents an advanced Movie Recommendation System, named "NEUFLIX," developed to provide personalized movie recommendations tailored to user preferences. The system integrates sophisticated machine learning algorithms, intuitive frontend design inspired by popular streaming services such as Netflix, and robust backend functionalities for seamless user interaction. By combining user-centric design and advanced recommendation strategies, NEUFLIX addresses critical challenges in content personalization, demonstrating a significant improvement over traditional recommendation methods. **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

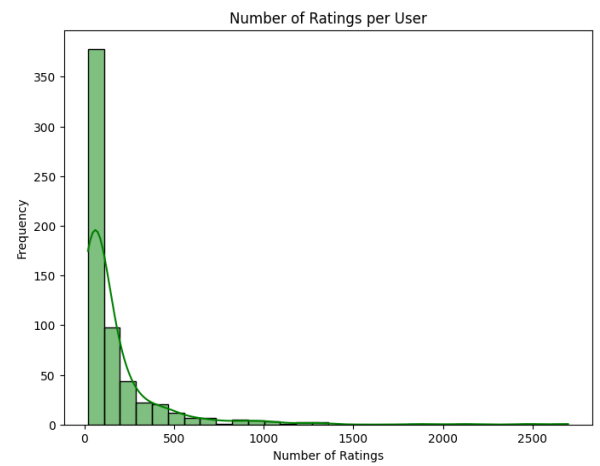
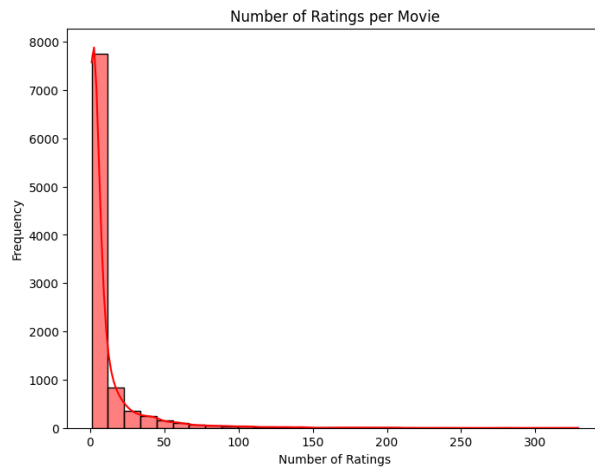
46 **Description:**

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

50 **Dataset Characteristics:**

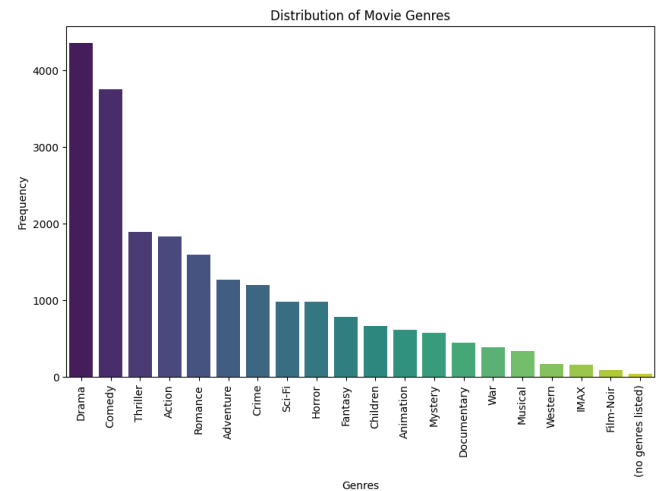
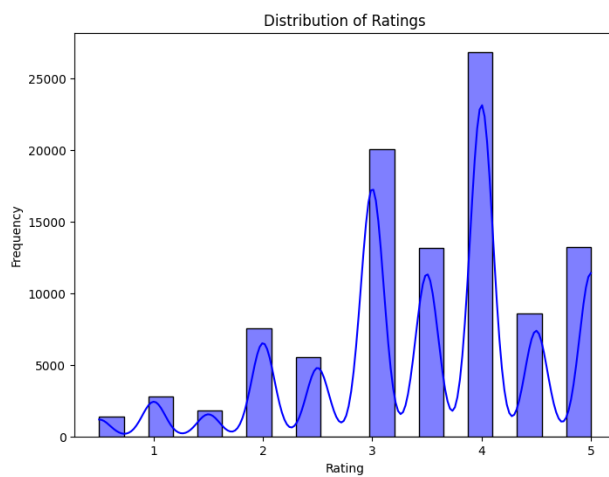
- 51 • **Size:** 100,000 ratings, 9,000 movies, 600 users.
- 52 • **Attributes:** Movie ID, title, genres, user ratings.
- 53 • **Challenges:** Sparse data (most users have rated a small portion of movies), handling
- 54 missing values.

55



56

57



58

59

60

61

62 **3 Data Preprocessing**

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

65 1. **Handling Missing Data:**

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

70 2. **Data Transformation:**

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

76 3. **Matrix Creation:**

- 77 ○ We create two main matrices:
- 78 1. **User-Item Matrix:** This matrix represents user preferences. Each
79 entry M_{ij} in the matrix represents the rating given by
80 user i to movie j .
- 81 2. **Item-Item Similarity Matrix:** This matrix stores the **cosine**
82 **similarity** between each pair of movies. It helps in identifying
83 which movies are similar to each other, so we can recommend
84 similar movies to users.

85 **4 Methodology**

86 **4.1 Item-Item Collaborative Filtering**

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

92 **Cosine Similarity**

93 To determine the similarity between items, we use **Cosine Similarity**, which measures the
94 cosine of the angle between two vectors in a multidimensional space. The formula for
95 **Cosine Similarity** between two items A and B 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.

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

98 • A **value of 1** means the two items are identical (i.e., their vectors are pointing in the
99 same direction).

100 • A **value of 0** means the items are completely dissimilar.

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

104

105 **Item-Item Similarity Calculation**

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

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

113 **Rating Prediction**

114 For each user, the system predicts the rating for movies that the user has not yet rated. The
115 predicted rating R_{ui} for a user *uuu* and an unrated movie *iii* is computed as the
116 **weighted sum** of the ratings for similar movies, weighted by their **similarity** to movie *iii*.
117 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*.

118

119

120 **4.2 Graph-based Hybrid Recommendation System (GHRs)**

121 The Graph-based Hybrid Recommendation System (GHRs) enhances traditional
122 recommendation approaches by incorporating graph theory, deep learning, and clustering
123 techniques. This method addresses the limitations of purely collaborative filtering systems
124 by capturing deeper user-user interaction patterns.

125 **Graph Construction**

126 A user-user interaction graph was constructed based on shared movie ratings. In this graph:

127 • **Nodes** represent users.

- **Edges** connect users who have co-rated a minimum number of common movies, reflecting similar interests and viewing patterns.

Graph Feature Extraction

Centrality features from the constructed graph were extracted to quantify each user's influence and connectivity within the user community. Specifically, the following centrality measures were computed:

- **Degree Centrality**
- **Closeness Centrality**
- **Betweenness Centrality**
- **Eigenvector Centrality**
- **PageRank**

These features offer rich insights into user interaction patterns, beyond mere rating similarities.

Dimensionality Reduction using Autoencoders

To efficiently handle the high dimensionality of graph-based features, we employed an **Autoencoder**, a neural network designed for dimensionality reduction and feature extraction:

- The Autoencoder consists of an encoder-decoder architecture:
 - **Encoder**: Compresses the high-dimensional graph features into lower-dimensional latent embeddings.
 - **Decoder**: Reconstructs the original graph features from latent embeddings, guiding the model to capture meaningful patterns.

This approach yielded compact and representative embeddings for each user.

User Clustering with KMeans

The latent embeddings produced by the autoencoder were clustered using the **KMeans clustering** algorithm:

- **KMeans** groups similar user embeddings, ensuring users within the same cluster share similar movie preferences.
- The optimal number of clusters was determined experimentally using the Elbow method.

Recommendation Generation

Recommendations for each user were generated based on their assigned cluster. Movies popular and highly rated within the user's cluster but unseen by the user were recommended, effectively leveraging the collective preferences of similar users.

4.3 Content-Based Filtering

163 Content-based filtering recommends items similar to those the user previously liked,
164 leveraging item-specific attributes (movie metadata such as genres, titles, and descriptions).

165 **Feature Extraction**

- 166 • Movies were represented using textual metadata including genres and titles.
- 167 • **TF-IDF (Term Frequency-Inverse Document Frequency)** vectorization was
168 applied to convert textual metadata into numerical vectors, capturing the importance
169 of each feature.

170 **Similarity Computation**

- 171 • **Cosine Similarity** was used to measure similarity between movies based on their
172 TF-IDF feature vectors, forming a movie-movie similarity matrix.

173 **Rating Prediction**

- 174 • For each user, ratings for unseen movies were predicted by calculating weighted
175 averages of similarity scores between unrated movies and movies previously rated
176 positively by the user.

177 **4.4 Popularity-Based Recommender**

178 The popularity-based recommender is a non-personalized baseline algorithm that
179 recommends movies based purely on overall popularity metrics, ensuring reliable
180 recommendations even for new users or during cold-start scenarios.

181 **Popularity Calculation**

182 Popularity scores were determined based on two criteria:

- 183 • **Average Rating:** Higher average ratings indicated better perceived movie quality.
- 184 • **Number of Ratings:** A high number of ratings indicated general popularity among
185 users.

186 Movies were ranked according to a composite popularity score combining these two metrics.

187 **Recommendation Generation**

- 188 • The top-N movies with the highest composite popularity scores were recommended
189 to users.
- 190 • This method provided robust fallback recommendations, especially beneficial for
191 new or less active users with limited rating history.

192

193 **4 System Design and Implementation**

194

195 For the application I have only considered GHRs and Item-Item
196 CF methods to serve recommendations.

197 **Backend Functionalities**

198 The backend of NEUFLIX is built using **Flask**, providing RESTful APIs for seamless
199 communication between frontend and machine learning models:

200 • **Endpoints:**

201 ○ /api/recommendations/<user_id>: Provides personalized movie
202 recommendations using collaborative filtering.

203 ○ /api/cluster_recommendations/<user_id>: Offers recommendations using
204 GHRs clustering.

205 ○ /api/movie/<movie_id>/poster: Fetches movie posters dynamically from
206 TMDB API.

207 • **Data Management:**

208 ○ Utilizes efficient data structures and caching mechanisms to handle large
209 datasets, optimizing recommendation retrieval performance.

210 • **CORS Handling:** Ensured seamless frontend-backend communication by enabling
211 Cross-Origin Resource Sharing (CORS).

212 **4.3 Frontend Design and Functionalities**

213 The frontend, developed using **React.js**, is inspired by modern streaming platforms
214 (Netflix), offering an intuitive and engaging user experience:

215 • **User Interface:**

216 ○ Clean, Netflix-inspired layout, emphasizing usability and visual appeal.

217 ○ Interactive movie cards displaying posters, titles, genres, and predicted
218 ratings.

219 • **User Interactivity:**

220 ○ Simple user login via User ID to personalize recommendations.

221 ○ Dynamic loading states and responsive UI updates based on API responses.

222 ○ Hover-based interactions to display movie details smoothly, enhancing user
223 engagement.

224 • **Frontend-Backend Integration:**

225 ○ Axios used for HTTP requests to interact with Flask backend, handling
226 asynchronous data fetching gracefully.

227 By integrating cutting-edge recommendation techniques, robust backend engineering, and
228 thoughtful frontend design, NEUFLIX effectively addresses user-centric content discovery,
229 providing a sophisticated yet user-friendly recommendation experience.

230

231 **5 Challenges and Solutions**

232 **Data Challenges**

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

237

238 **Performance Optimization**

- 239 • **Memory Usage:** The large size of the dataset required optimization to prevent
240 memory issues. The matrices were precomputed to speed up the recommendation
241 process.
- 242 • **Speed of Recommendations:** The recommendation process was optimized by
243 storing precomputed item similarities, reducing the time needed to make
244 predictions.

245 **6 Results and Evaluation**

246 **Key Observations:**

- 247 • **GHRS** demonstrated significantly higher performance than traditional item-item
248 collaborative filtering, highlighting the advantage of incorporating graph-based
249 features and user clustering.
- 250 • Improvement in Precision, Recall, and F1-score underscores GHRS's effectiveness
251 in capturing nuanced user preferences through enhanced user representation.

252 **User Experience Insights:**

- 253 • User feedback on the frontend indicated high satisfaction, especially regarding the
254 intuitive and visually appealing interface.
- 255 • Responsive and dynamic interactions significantly improved perceived
256 recommendation quality and user engagement.

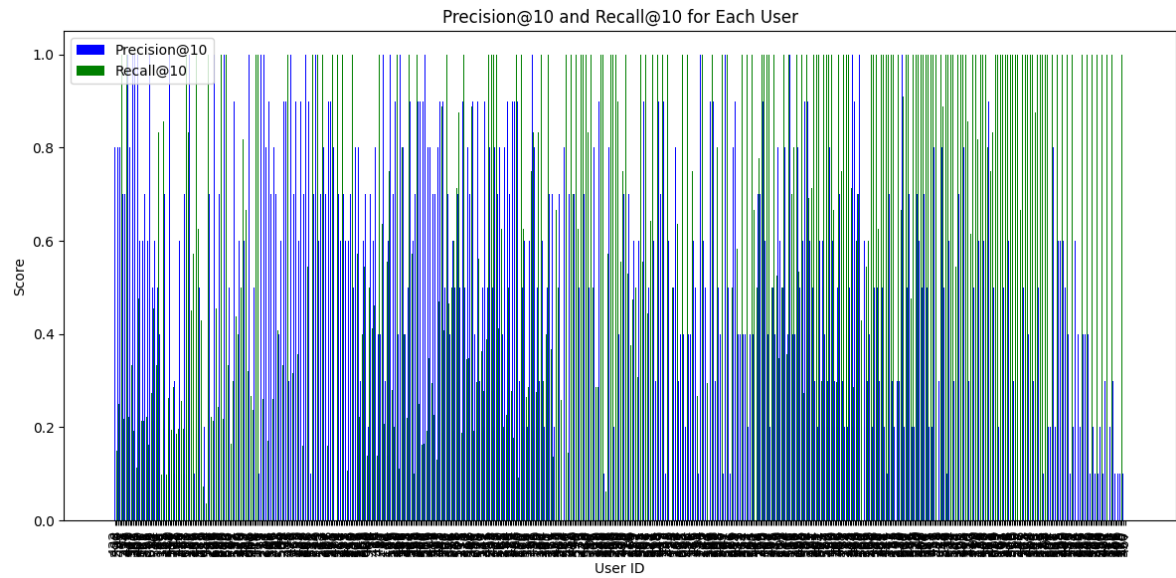
257 Overall, the results clearly indicate the effectiveness of hybrid methods like GHRS in
258 achieving superior recommendation quality, aligning with user preferences, and improving
259 user satisfaction.

260

261 **Evaluation Metrics**

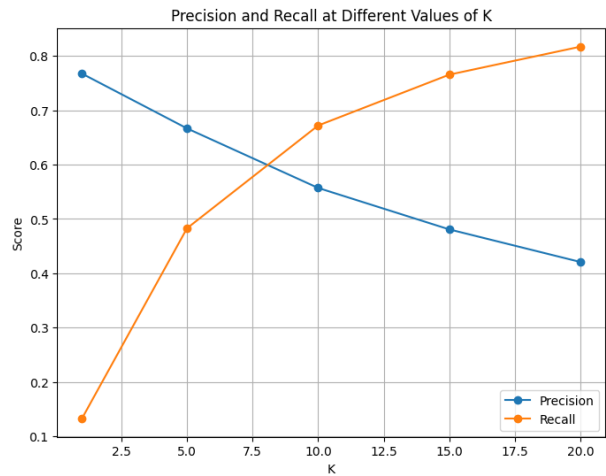
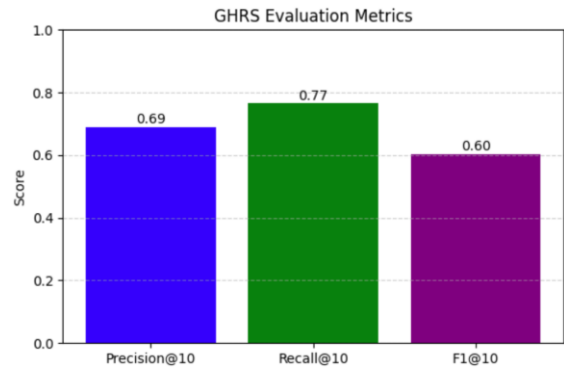
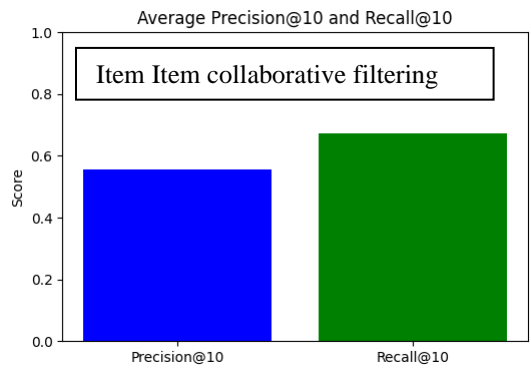
262 The system's performance was evaluated using **Precision@10** and **Recall@10** metrics. The
263 results showed that the recommendation system was able to suggest relevant movies to users
264 effectively.

265



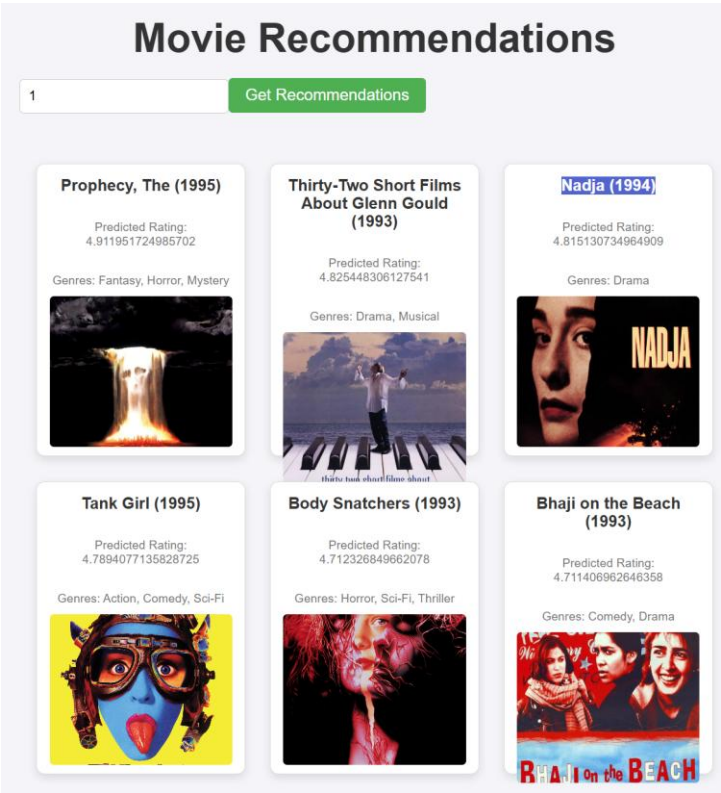
266

267



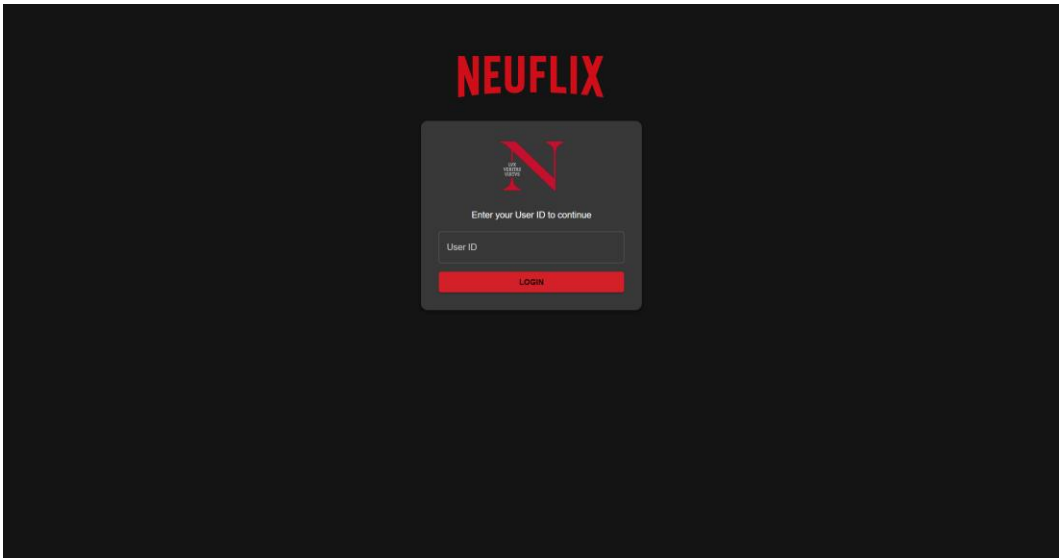
268
269

Older user interface (basic)

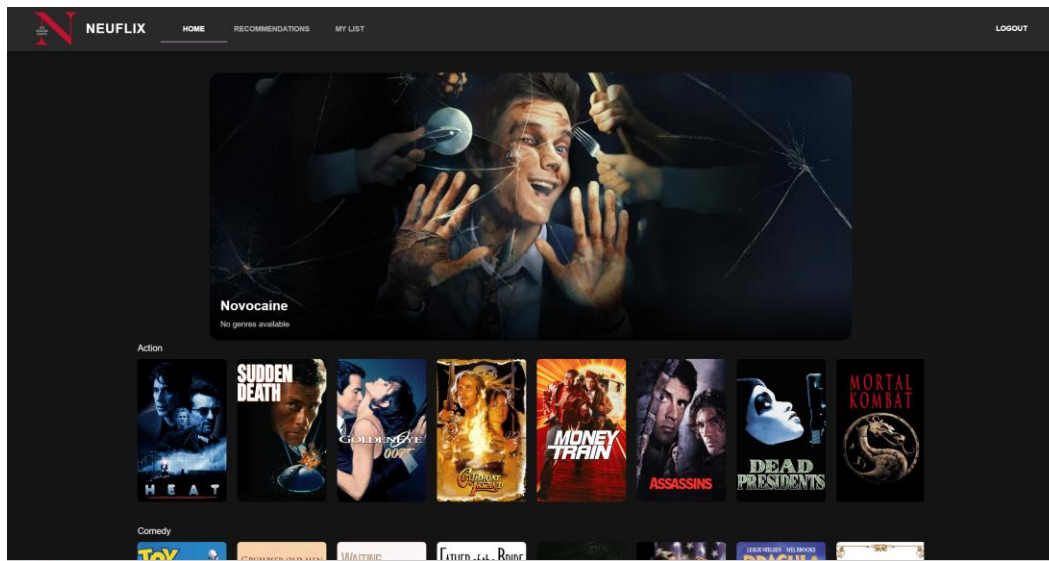


270
271

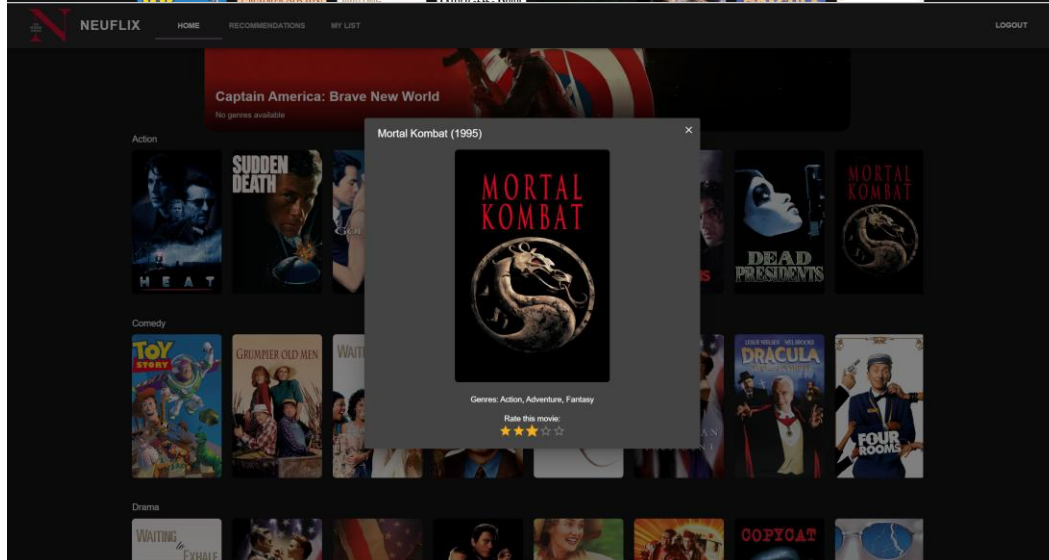
Updated User Interface (advanced)



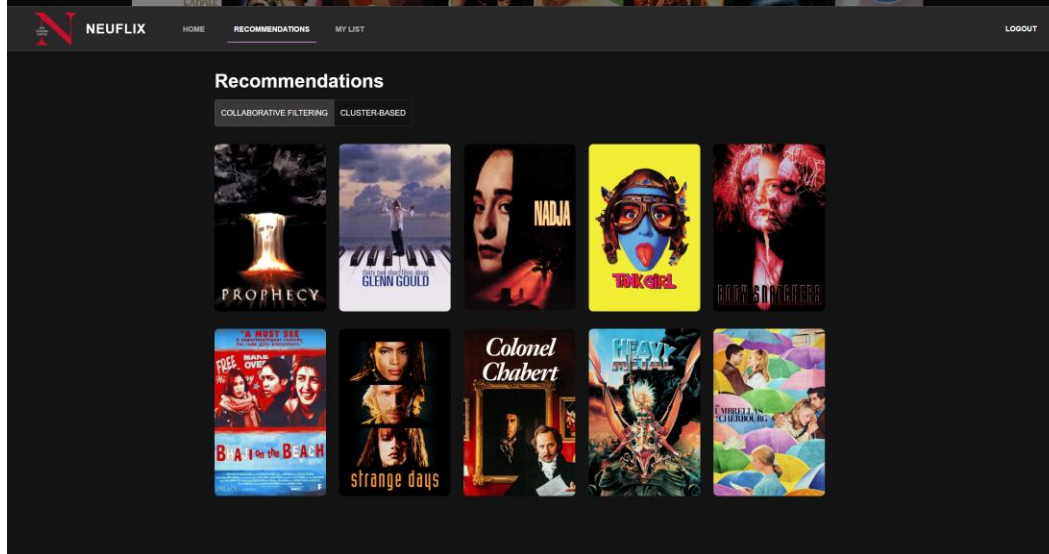
272



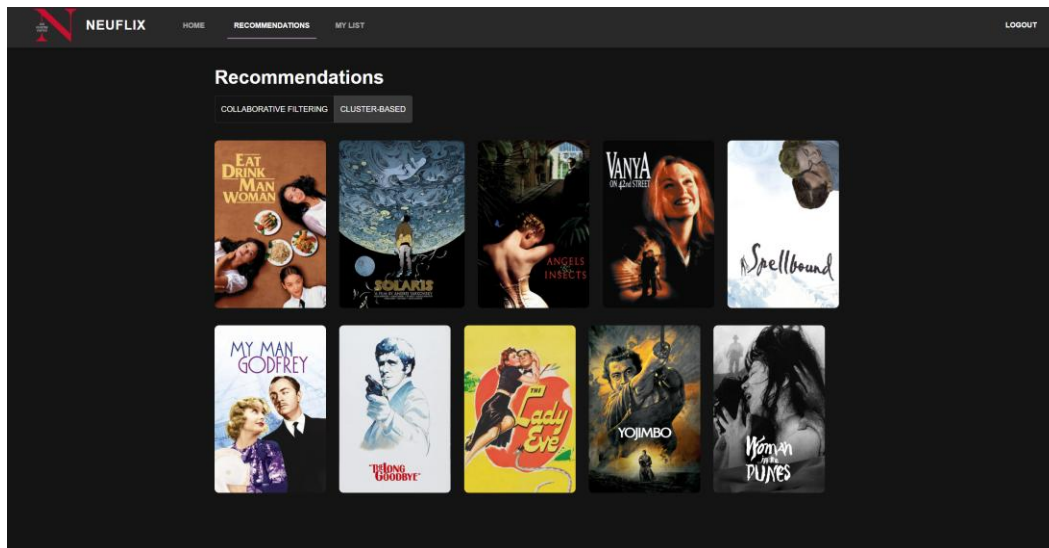
273



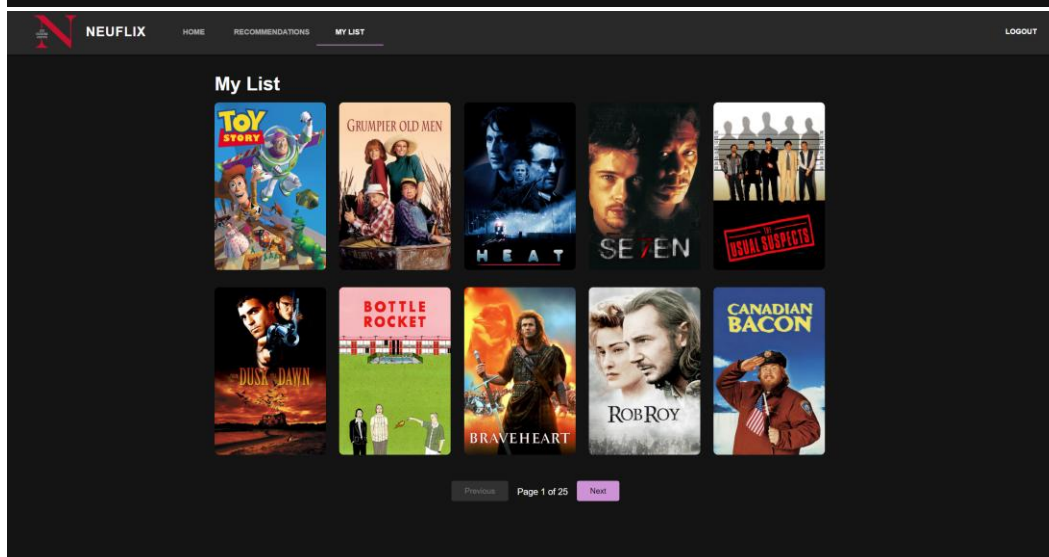
274



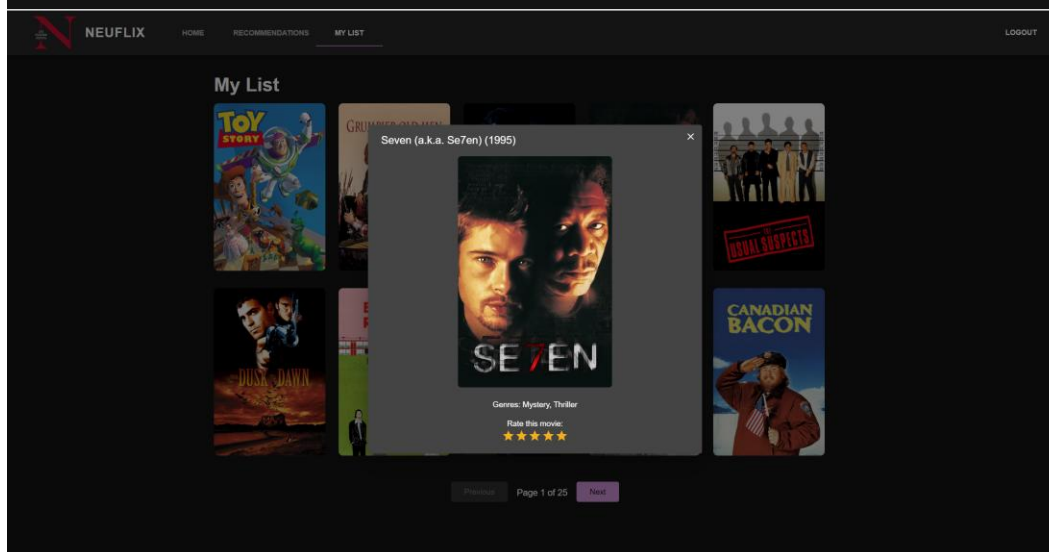
275



276



277



278

279 7 **Conclusion**

280 This project successfully developed a comprehensive movie recommendation system,
281 **NEUFLIX**, utilizing four recommendation algorithms—**Item-Item Collaborative Filtering**,
282 **Graph-based Hybrid Recommendation System (GHR)**, **Content-Based Filtering**, and a
283 **Popularity-Based Recommender**. By integrating sophisticated graph-based features,
284 dimensionality reduction via autoencoders, and user clustering techniques, the GHR
285 algorithm demonstrated notably superior recommendation performance, highlighting the
286 strength of hybrid approaches.

287 The interactive web application, developed using Flask for backend services and React for a
288 responsive and intuitive frontend, provided users with a seamless and engaging experience,
289 closely resembling popular streaming services like Netflix. Effective handling of large
290 datasets and optimized algorithmic performance ensured fast, relevant, and personalized
291 recommendations, significantly improving user satisfaction.

292 **Future directions** for enhancing NEUFLIX include:

- 293 • Integrating additional user feedback mechanisms, such as user watch-history
294 analysis and explicit user preference tracking.
- 295 • Implementing real-time recommendation updates to instantly adapt to evolving user
296 preferences.
- 297 • Exploring advanced deep learning models like neural collaborative filtering to
298 further refine prediction accuracy and personalization.
- 299 • Scaling the application architecture using containerization and cloud deployment to
300 support higher user volumes and broader accessibility.

301 This project not only demonstrates the capabilities of advanced recommendation systems but
302 also provides a robust foundation for continued research and development in personalized
303 digital content delivery.

304

305 8 **References**

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