

Personalized Movie Recommendation with Item-Based Collaborative Filtering

By

SIT YIE SIAN



FACULTY OF COMPUTING AND
INFORMATION TECHNOLOGY

TUNKU ABDUL RAHMAN UNIVERSITY OF
MANAGEMENT AND TECHNOLOGY
KUALA LUMPUR

ACADEMIC YEAR
2023/2024

PERSONALIZED MOVIE RECOMMENDATION WITH ITEM-BASED COLLABORATIVE FILTERING

By
SIT YIE SIAN

Supervisor: Mr. Choong Yun Loong

A project report submitted to the
Faculty of Computing and Information Technology
in partial fulfillment of the requirement for the
Bachelor of Computer Science (Honours) in Data Science
Tunku Abdul Rahman University of Management and Technology

Department of Mathematical And Data Science
Faculty of Computing and Information Technology
Tunku Abdul Rahman University of Management and Technology
Kuala Lumpur

Copyright by Tunku Abdul Rahman University of Management and Technology.

All rights reserved. No part of this project documentation may be reproduced, stored in retrieval system, or transmitted in any form or by any means without prior permission of Tunku Abdul Rahman University of Management and Technology.

Declaration

The project submitted herewith is a result of my own efforts in totality and in every aspect of the project works. All information that has been obtained from other sources had been fully acknowledged. I understand that any plagiarism, cheating or collusion or any sorts constitutes a breach of TAR University rules and regulations and would be subjected to disciplinary actions.

_____ *yiesian* _____

SIT YIE SIAN

ID: 21WMR03693

DATE: 28/4/2023

Abstract

The primary objective of this project is to design and implement a movie recommendation system using item-based collaborative filtering, enhancing users' movie-watching experience by providing personalized movie recommendations. Addressing the challenge of finding relevant and enjoyable movies amidst the rapid growth of the entertainment industry, our system caters to individual preferences and tastes. The project encompasses the development of both item-based and user-based collaborative filtering algorithms, collection and preprocessing of movie rating data from the MovieLens dataset. Employing tools and techniques such as Python programming, pandas, NumPy, and similarity metrics like Pearson correlation coefficient and cosine similarity, the system incorporates various functional modules, including data preprocessing, similarity computation, rating prediction, and recommendation generation. Performance evaluation using metrics like Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) demonstrates the effectiveness of the collaborative filtering algorithms in generating personalized movie recommendations, successfully capturing individual user preferences and improving the movie-watching experience. Further improvement could address data sparsity issues and incorporate hybrid filtering techniques to enhance recommendation quality.

Keywords: item-based collaborative filtering, movie recommendation system, similarity computation, rating prediction

Acknowledgement

I would like to express my sincere gratitude to all those who have contributed to the successful completion of this project. Their valuable insights, ideas, and support have been instrumental in shaping the development and implementation of the movie recommendation system.

First and foremost, I would like to thank my classmate, whose invaluable ideas and suggestions have significantly contributed to the project's direction and design. Their enthusiasm and dedication to the project have not only helped me overcome various challenges but also inspired me to think critically and creatively throughout the process.

I also extend my appreciation to my supervisor and moderators, who have provided me with the necessary knowledge and guidance to navigate the complexities of building a recommendation system. Their expertise and experience have been crucial in shaping my understanding of the subject matter and developing a robust and effective system.

Additionally, I would like to acknowledge the contributions of the GroupLens Research team for providing the MovieLens dataset, which served as the foundation for our recommendation system. Their efforts in creating and maintaining this comprehensive dataset have greatly facilitated our work and enabled us to focus on the development and optimization of the recommendation algorithm.

Finally, I would like to thank my friends and family for their unwavering support and encouragement throughout the project. Their belief in my abilities and their constant motivation have been essential in helping me persevere and achieve my goals.

To all of you, thank you for your invaluable contributions to this project. Your collective efforts have made it possible for us to create a movie recommendation system that can enhance the movie-watching experience for countless users.

Table of Contents

1.0 Introduction	8
1.1 Project Objectives	8
1.1.1 Provide customize recommendation that match user interests	8
1.1.2 Generate revenue for the content providers	9
1.1.3 Improving the value proposition of the media service	9
1.2 Problem Statement	10
1.2.1 Problem Statement	10
1.2.2 Proposed Solution	10
1.3 Background	11
1.4 Advantages and Contribution	11
1.5 Project Plan	12
1.5.1 Milestones	12
1.6 Project Team and Organization	13
1.7 Chapter Summary & Evaluation	14
2 Literature Review	16
2.1 Item-Based Collaborative Filtering	16
2.2 Item-Based Collaborative Filtering in Movie Recommendation Systems	17
2.2.1 Limitations of Item-Based Collaborative Filtering	18
2.3 Chapter Summary and Evaluation	19
3 Methodology and Requirements Analysis	22
3.1 Introduction	22
3.2 Data Collection	22
3.3 Data Preprocessing	23
3.4 Model Selection and Training	23
3.5 Model Evaluation	24
3.6 Deployment and Maintenance	24
3.7 Ethics Consideration	25
3.7.1 Data Privacy and Security:	25
3.7.2 Bias and Fairness:	25
3.7.3 Transparency and Explainability:	26
3.8 Chapter Summary	26
4.0 System Design	29
4.1 Introduction	29
4.2 Data Collection	30
4.3 Data Preprocessing	30
4.4 Implementation of Item-Based Collaborative Filtering	32
4.5 Evaluation Metrics and Methods	33
4.6 Chapter Summary	35
5 Implementation and Testing	38
5.1 Importing Libraries	38
5.2 Loading and Exploring the Data	39

5.3 Data Visualization	43
5.4 Item-Based Collaborative Filtering	51
5.5 Item-Based Collaborative Filtering with SVD	53
5.6 User-Based Collaborative Filtering	56
5.7 Deployment	58
5.7.1 User-based movie recommendation system based on UserId	58
5.7.2 User-based movie recommendation system based on Genre	59
5.7.3 Item-based movie recommendation system based on MovieId and rating	61
5.8 Chapter Summary and Evaluation	63
6 System Deployment	66
6.1 Test Data for User-based movie recommendation system based on UserId	66
6.1.1 Testing data	66
6.1.2 Review	68
6.1.3 Conclusion	71
6.2 Test Data for User-based movie recommendation system based on Genre	72
6.2.1 Testing Data	72
6.2.2 Review	75
6.2.3 Conclusion	77
6.3 Test Data for Item based recommendation system	77
6.3.1 Testing data	78
6.3.2 Review	80
6.3.3 Conclusion	82
6.4 Chapter Summary and Evaluation	83
7 Discussions and Conclusion	87
7.1 Summary	87
7.2 Achievements	87
7.3 Contributions	88
7.4 Limitations and Future Improvements	88
7.5 Issues and Solutions	89
7.6 Conclusion	89
References	90

Chapter 1

Introduction

1.0 Introduction

The frequency we use with media has increased significantly in the digital age, and movies are no exception. Making a choice from the vast selection of movies accessible on numerous platforms might be difficult. This is where movie recommender systems come in to solve the issue. Movie Recommender systems help users discover new movies that they are likely to enjoy based on their previous viewing history, ratings, and preferences.

For movie fans, movie recommendation systems are essential. By making it simple for people to discover new movies that are pertinent to their interests, it saves users' time and effort. With recommender systems, users can find movies that they might not have found on their own. As a result, users are more likely to consume material that is pertinent to their interests, increasing engagement and happiness with the media service.

Movie recommender systems are important because they provide personalized recommendations. As each user has unique likes and preferences, a one-size-fits-all strategy is ineffective. It can offer specialized recommendations by utilizing our recommender system to match the user's particular interests. Users enjoy the service more, are more engaged, and have a better overall experience as a result.

1.1 Project Objectives

1.1.1 Provide customize recommendation that match user interests

By using collaborative filtering, the recommender system can analyze and compare the ratings that users have given to suggest different movies that align with their user's preferences. The recommender algorithm looks for films that the user will enjoy in order to improve their overall movie-finding and viewing experience. In order to accomplish this goal, a dataset of movie ratings will be collected, and this dataset will be used to train a model that can forecast the ratings of related movies. We then will use the model to make recommendations to users that have received high predicted ratings from the model.

Increased user engagement and retention result from personalized recommendations. If suggestions are made based on the interests of the users, they are more likely to use the site for longer, see more content, and have a better overall experience. This leads to an increase in

loyalty and retention because users are more inclined to stick with the service if they are satisfied with the recommendations offered.

1.1.2 Generate revenue for the content providers

Not only is it advantageous to users to receive personalized recommendations, but it can also bring in money for content providers. By using the personalized recommendations, content providers can increase user interaction with their platform. Users are therefore more inclined to continue with the service and may even spend more money to purchase additional material, which can increase providers' revenue. The method can also be utilized to present particular movies to consumers who would not have otherwise known about them. This might result in more people watching and renting these movies, which would enhance the providers' income. In addition, content providers can use the data to learn more about consumer patterns and preferences. Their choice of content and marketing tactics can be improved by using information, resulting in better income. Ultimately, by keeping consumers engaged with their service for longer periods of time, recommending movies to users, the movie recommendation engine that we are developing has the potential to bring in money for content providers.

1.1.3 Improving the value proposition of the media service

Making personalized recommendations improves the media service's value proposition. The service offers a more individualized and carefully curated experience by customizing recommendations to the user's tastes. This fosters a sense of worth and comprehension in the user, which can boost brand loyalty and encourage effective word-of-mouth advertising. In order to enhance views, fresh content or content that may not have received as much attention can be utilized to be promoted. By providing a personalized and intuitive recommendation system, the media service can differentiate itself by attracting new users from competitors for a more tailored and engaging movie-watching experience.

To enhance the value proposition of the media service is to provide users with additional content related to the movies that they are interested in. A watchlist function that enables users to save movies they are interested in watching might be built into the system. Based on

the user's saved movies, the system can then use this watchlist data to offer even more tailored recommendations.

1.2 Problem Statement

1.2.1 Problem Statement

With so much information available in the modern world, it can be challenging to select movies that a person enjoys. Recommender systems are perhaps the most often used tool in data mining techniques. By developing a recommendation system that proposes movies based on what other people with similar tastes like, this project seeks to make it easy to locate good movies. Using the least amount of information feasible on a product's features, recommender systems aim to make locating a good or service easier. A number of factors are used to look at the relationships between patterns and user attributes in order to choose the best movie suggestions for a user.

1.2.2 Proposed Solution

The project's goal is to create a movie recommendation system that uses item-based collaborative filtering to provide users with personalized suggestions. This system's objective is to help users in selecting movies that suit their tastes in the current digital environment. By looking at the ratings that users have given to various movies and figuring out the similarity between pairs of movies based on these ratings, the system can offer recommendations to users based on their prior ratings and the ratings of similar movies. To do this, a dataset of movie reviews will be collected and used to train a machine learning model that can predict the ratings that consumers would give to films based on the ratings of comparable films. After the model has been trained, users will receive individualized recommendations based on their watching preferences and interests.

Collaborative filtering works by matching the similarities in items and users. It examines the user attributes along with the features of the previous content that users have viewed or looked up (Shen J., 2020). Recommendations in movie recommendation systems are based on user data and what other users with similar user data are watching. For instance, user demographics like age, gender, and ethnicity are selected by collaborative filtering in movie

recommender systems (Dakhel G.M., 2011). With the use of these attributes, movie recommendations are created that match to individuals with similar demographic traits and past user search history. If the user doesn't input any data or there isn't enough data for any reliable clustering, collaborative filtering suffers from a cold start.

1.3 Background

In today's world, it can be hard to find content, like books, videos, and movies, that you like. To help people find what they want, some companies use a system called a recommender system. This system suggests content to users based on what they have liked in the past. There are two main ways to build a recommender system: collaborative filtering and content-based filtering. I propose a personalized movie recommendation system that utilizes item-based collaborative filtering. A strategy used frequently in recommendation systems is collaborative filtering, which includes examining the ratings that users have given to products and using this data to suggest products to other users according to their common rating patterns. A variation of collaborative filtering called item-based collaborative filtering emphasizes interconnections between various objects rather than connections between users and items.

1.4 Advantages and Contribution

There are several advantages that a movie recommendation system that utilizes item-based collaborative filtering could offer over competing systems. One of the key advantages of collaborative filtering approaches, including item-based collaborative filtering, is their ability to provide personalized recommendations to users based on their past ratings and the ratings of similar users or items. This allows the system to make more relevant and accurate recommendations to users, which can improve the overall user experience. They can make a real quality assessment of items by considering other people's experience (Sharma and Yadav).

Collaborative filtering approaches are generally very scalable, as they do not require the system to have detailed information about the attributes of the items being recommended

(such as genre, actors, and plot keywords). This makes it easier to add new items to the system and maintain the recommendations over time. For systems that need detailed information, movie catalogs can rapidly change as time goes and this increases the difficulties as it requires the system to have a comprehensive and up-to-date database of item attributes.

1.5 Project Plan

1	Define the project scope and requirements
2	Collect and preprocess data
3	recommendation algorithm
4	Build and train the model
5	Implement the recommendation system
6	Test and validate the system
7	Deploy and maintain the system

1.5.1 Milestones

The proposed project milestone is as follows:

Milestone	Milestone Goal	Deadline
Concept approval	Feasibility studies and basic system concepts is approved	16/2/2023
Requirements review	Requirements specifications are complete, correct, approved and suitable for input to design.	26/2/2023

System design review	The system design satisfies all product requirements, is approved and is suitable for input into the detailed design process.	11/3/2023
UI design review	Detailed UI design of the system architecture, are approved and are suitable for input into the development of code.	18/3/2023
Test plan review	Test plans are adequate for the testing of all product features, are approved and are suitable for input to the development of test cases and test procedures.	22/5/2023
Test readiness review	Developed and unit tested software has been passed by the test team and is suitable for input into integration testing.	3/7/2023
System test review	The software product has passed system testing and is suitable for input into acceptance testing.	10/7/2023
Operational readiness	The software product has passed acceptance testing and is suitable for deployment in its target production environment.	24/7/2023

1.6 Project Team and Organization

<i>System Modules</i>	<i>Sit Yie Sian</i>	<i>Leong Sheng Mou</i>
<i>Content-based Filtering</i>		✓
<i>Collaborative Filtering</i>	✓	

1.7 Chapter Summary & Evaluation

In this chapter, as a summary, the main objectives of the system were discussed. Objectives that have been listed are to automate the coordination tasks in the movie recommended system, provide better and scalable movie rating reports from reviews . Besides, the problem statements of the current movie recommended system and the proposed solution were discussed. Furthermore, the background of the recommended system, the benefits of the movie recommended system were mentioned. Lastly, the project milestone schedule and project team and organization have been discussed.

Chapter 2

Literature Review

2 Literature Review

Popularity of the world is only increasing and systems for recommending movies have become a crucial component of the modern digital environment. We find it challenging to choose the movie we want to see because there are so many options available. For this reason, we require a recommender system to do it for us. Recommender systems have different methods to make movie recommendations to users including content-based filtering, collaborative filtering, and hybrid filtering. One of the best methods for creating movie recommendation systems is collaborative filtering, particularly item-based collaborative filtering. We will examine the idea of item-based collaborative filtering, how it is used in systems that recommend movies, as well as some of the benefits and drawbacks of this method in this literature review.

2.1 Item-Based Collaborative Filtering

An approach used in recommender systems to estimate user preferences for items based on their similarities to other items is called item-based collaborative filtering. Based on the user's rating history, the system creates an item similarity matrix. After that, each user will receive personalized recommendations based on the similarity matrix. The fundamental assumption behind item-based collaborative filtering is that users who have rated similar movies in the past will also rate similar movies a similar score. In other words, the system can predict a user's rating for a movie from the history of the user who has rated the similar movie.

The fact that item-based collaborative filtering can manage the sparsity of the user-item ratings matrix is one of its main features. For instance, a research by Zhou et al. (2010) has compared the effectiveness of their approach to a number of different recommendation techniques, such as matrix factorization and user-based collaborative filtering. The result demonstrated that their research revealed that item-based collaborative filtering outperformed the other ones in terms of managing the sparsity of the user-item rating matrix. In particular, whereas the other approaches struggled with sparsity, item-based collaborative filtering was able to offer correct recommendations for individuals who had only reviewed a few movies. Overall, the authors' findings shows that Item-based collaborative filtering appears to be a

potential strategy for movie selection that can assist in managing sparsity in user-item rating matrices.

Item-based collaborative filtering also has the benefit of being able to deal with cold-start issues. For instance, a research by Sarwar et al. (2001) has demonstrated that item-based collaborative filtering algorithms perform better than various approaches, especially when it comes to tackling the cold-start problem. In particular, the algorithms can use the similarity between products to create precise suggestions for new users and new items in the absence of user data. Overall, the authors' findings are in line with the notion that cold-start problems in recommendation systems can be successfully addressed by item-based collaborative filtering.

It has also been demonstrated that item-based collaborative filtering is very scalable. For instance, a research by Lemire & Maclachlan (2005) has demonstrated that Slope One predictors can handle very large datasets with millions of users and items and are noticeably faster than other methods. The authors also show that their algorithm can make recommendations with little latency in real-time. Overall, according to the authors', the Item-based collaborative filtering, particularly in the form of Slope One predictors, is extremely scalable and can successfully handle massive datasets with millions of users and objects. They evaluated how well their algorithm performs in comparison to other recommendation techniques, such as traditional item-based collaborative filtering.

2.2 Item-Based Collaborative Filtering in Movie Recommendation Systems

Several researches have demonstrated that in movie recommendation systems, item-based collaborative filtering performs better than other strategies like content-based filtering. For instance, a research by Sarwar et al. (2001) demonstrated that in terms of prediction accuracy, item-based collaborative filtering beats content-based filtering. According to the study, item-based collaborative filtering was more accurate than content-based filtering at predicting user ratings.

In a further study, Shao et al. (2019) evaluated the effectiveness of deep learning, matrix factorization, and item-based collaborative filtering. According to the study, even when the dataset was limited, item-based collaborative filtering could still produce reliable recommendations. Also, in terms of prediction accuracy, item-based collaborative filtering was able to outperform matrix factorization and deep learning techniques.

Using item-based collaborative filtering in movie recommendation systems can be done in a number of ways. To measure how similar two movies are, a typical method is to apply the Pearson correlation coefficient or cosine similarity. After that, recommendations are generated for each user based on the similarity matrix.

Another strategy is to group similar movies together using clustering algorithms. The system can then produce suggestions based on the user's past viewing habits and movies in similar clusters to those they have watched before.

A third strategy is to divide the user-item ratings matrix into low-rank matrices using matrix factorization techniques like singular value decomposition (SVD) or non-negative matrix factorization (NMF). The low-rank matrices can then be used by the system to produce suggestions.

2.2.1 Limitations of Item-Based Collaborative Filtering

Item-based collaborative filtering has the drawback of being unable to deal with user bias. User bias is the propensity of some users to rate movies more favorably or unfavorably than other users. Item-based collaborative filtering makes the assumption that people who have previously rated movies similarly will continue to do so for new released movies. This assumption might not be accurate for people who have strong biases. For instance, a research by Shen et al. (2009) demonstrates that PMF-UB performs noticeably better than alternative approaches, especially when it comes to minimizing the effects of user bias. The authors show that the algorithm outperforms traditional item-based collaborative filtering in terms of accurately recommending movies to viewers even with strong biases. Overall, the authors'

findings are consistent with the idea that item-based collaborative filtering may be a significant drawback due to user bias.

The inability of item-based collaborative filtering to capture the temporal dynamics of user preferences is another drawback. For instance, a research by Lamere (2013) has shown that user preferences are not constant and can vary over time, which could cause recommendations to become less accurate as time passes. The author proposes a new algorithm called Time-aware Item-based Collaborative Filtering (TiCF). The algorithm uses a weighted similarity score that gives greater weight to more recent ratings than to older ratings. The algorithm uses a decay function in order to gradually reduce the weight of older ratings. According to the experiments, TiCF performs noticeably better than conventional item-based collaborative filtering, especially as the dataset grows larger and more sparse.

The "long tail" issue also affects item-based collaborative filtering. The majority of the movies in the dataset have very few ratings, which is known as the "long tail problem" and makes it challenging to produce reliable recommendations for these movies. For instance, a research by Burke (2007) analyzes a number of sizable datasets, including the Netflix Prize dataset, which comprises user ratings for more than 100 million movies. The research reveals that most of the movies in the dataset have a very small number of ratings, which makes it difficult to generate reliable recommendations for these movies.

2.3 Chapter Summary and Evaluation

In conclusion, Item-based collaborative filtering has become one of the best approaches for movie recommendation systems. It is very scalable, can manage cold-start issues, and can handle the sparsity of the user-item ratings matrix. In terms of prediction accuracy, it has been demonstrated that item-based collaborative filtering performs better than other methods like content-based filtering.

Item-based collaborative filtering has certain drawbacks, such as the inability to handle user bias, capture the temporal dynamics of user preferences, and address the long tail. However,

it is still a well-liked method for movie recommendation systems due to its efficiency and scalability, despite these drawbacks.

Chapter 3

Methodology and Requirements Analysis

3 Methodology and Requirements Analysis

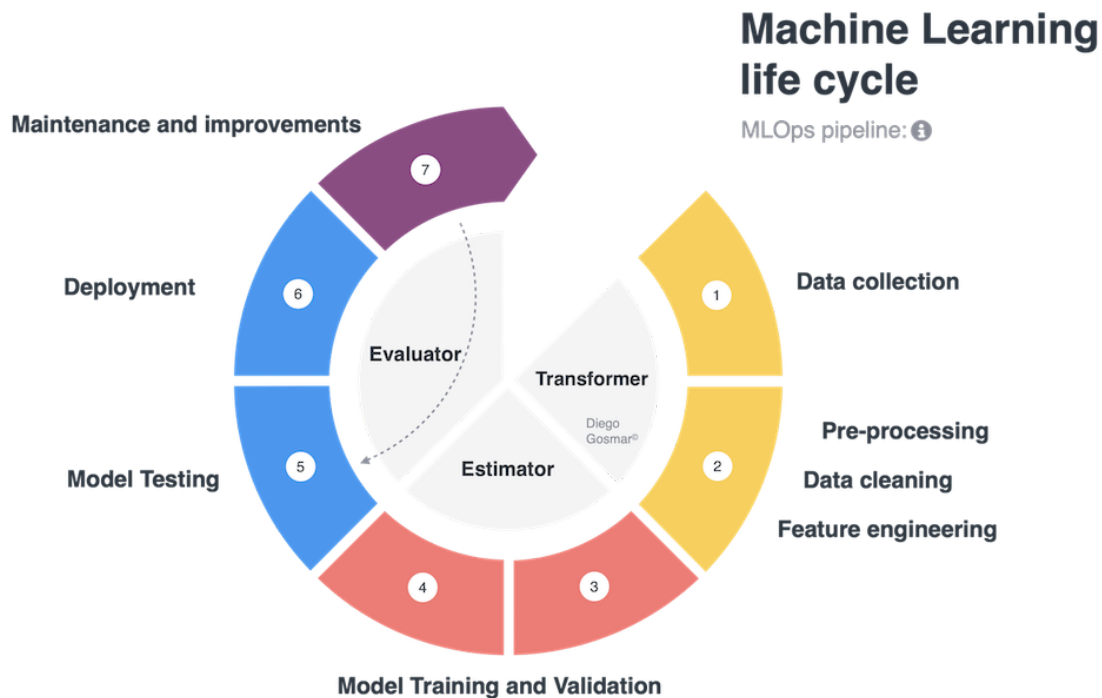


Figure 3.1 Machine Learning Pipeline

3.1 Introduction

This chapter will outline the methodology and requirements analysis for a general machine learning pipeline, focusing on the key stages involved in the design, implementation, and evaluation of a system. The chapter will discuss data collection, data preprocessing, model selection, training, evaluation, and deployment of the system. Additionally, the chapter will outline evaluation metrics and methods for measuring the effectiveness of the model.

3.2 Data Collection

Data collection is the first and one of the most critical stages in the machine learning pipeline. The data collected should be representative of the problem domain and provide sufficient information to train and evaluate the model. Data collection can involve multiple sources, such as public datasets, proprietary databases, user-generated content, or data obtained through web scraping. When collecting data, it is crucial to consider factors like data quality, size, diversity, and relevance to the problem at hand.

3.3 Data Preprocessing

Data preprocessing is an essential step in preparing the data for analysis and model training. The main goal of data preprocessing is to clean, transform, and structure the data, ensuring that it is consistent and suitable for further processing. Typical data preprocessing steps include:

- Data integration: Combining data from multiple sources into a unified dataset.
- Handling missing values: Identifying and addressing missing data points using techniques like imputation or removal.
- Data transformation: Converting data into a suitable format for analysis, such as normalizing numerical data or encoding categorical features.
- Feature engineering: Creating new features or modifying existing ones to improve model performance and interpretability.
- Data partitioning: Splitting the dataset into training, validation, and testing sets to evaluate model performance and prevent overfitting.

3.4 Model Selection and Training

Model selection involves choosing the most appropriate machine learning algorithm for the problem at hand. This process typically requires a thorough understanding of the problem domain, the data, and the available algorithms. Factors to consider when selecting a model include:

- Algorithm type: Choosing between supervised, unsupervised, or reinforcement learning methods, depending on the problem's requirements.
- Model complexity: Balancing between model simplicity and flexibility to prevent underfitting or overfitting.
- Scalability: Ensuring the selected algorithm can handle the dataset size and complexity.
- Interpretability: Prioritizing models that provide clear insights into the relationships between features and the target variable.

Once the model has been selected, the training process involves using the training dataset to adjust the model's parameters so that it can generalize well to unseen data. This may involve techniques like gradient descent, regularization, or other optimization methods.

3.5 Model Evaluation

Evaluating the performance of the machine learning model is crucial to ensure that it meets the desired objectives and provides accurate and reliable predictions. Evaluation metrics and methods help quantify the model's performance and facilitate comparisons with alternative approaches. Common evaluation metrics include:

- **Accuracy:** The proportion of correct predictions made by the model.
- **Precision, Recall, and F1 Score:** Metrics that consider the balance between true positives, false positives, and false negatives.
- **Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE):** Metrics that measure the difference between predicted and actual values for regression tasks.
- **Area Under the Receiver Operating Characteristic (ROC) Curve (AUC-ROC):** A metric that evaluates the trade-off between true positive rate and false positive rate.

Cross-validation is a common technique used to evaluate model performance by partitioning the dataset into multiple folds and iteratively training and testing the model on different subsets.

3.6 Deployment and Maintenance

Once the model has been trained and evaluated, it is ready for deployment in a production environment. The deployment process involves integrating the model into a larger system, such as a web application, mobile app, or other software platforms. Deployment considerations include:

- **Integration:** Ensuring that the model can be seamlessly integrated into the target system and interact with other components.

- Scalability: The deployed model should be able to handle varying loads and adapt to changes in the data and user requirements.
- Latency: Ensuring the model can provide timely predictions or recommendations, considering the real-time requirements of the application.
- Security: Protecting the model and its data from unauthorized access or tampering.

Maintenance is an ongoing process that involves monitoring the performance of the deployed model, updating it with new data, and fine-tuning the model to address any changes in the problem domain or user requirements. Regular maintenance helps ensure that the machine learning system remains accurate, reliable, and relevant over time.

3.7 Ethics Consideration

Ethics play a significant role in the development and deployment of machine learning systems, as these systems have the potential to impact users and society in various ways. It is essential to consider ethical implications throughout the machine learning pipeline to ensure that the developed models are fair, transparent, and respectful of user privacy. In this section, we will discuss some key ethical considerations for machine learning systems.

3.7.1 Data Privacy and Security:

When collecting and processing data, it is crucial to protect user privacy by handling sensitive information responsibly. This includes anonymizing data, obtaining user consent, and complying with relevant data protection regulations, such as the General Data Protection Regulation (GDPR). Additionally, data storage and access should be secured to prevent unauthorized access and data breaches.

3.7.2 Bias and Fairness:

Machine learning models can inadvertently learn and perpetuate biases present in the training data, leading to unfair treatment of certain user groups. To address this issue, researchers and

practitioners should actively identify and mitigate potential biases in the data and model. Techniques such as re-sampling, re-weighting, and adversarial training can help reduce bias and promote fairness in the model's predictions.

3.7.3 Transparency and Explainability:

Machine learning models, particularly deep learning models, can often be seen as "black boxes" that produce predictions without a clear explanation of how they arrived at those results. Ensuring transparency and explainability in machine learning systems can help users understand the decision-making process, build trust in the system, and facilitate accountability. Techniques such as interpretable models, feature importance analysis, and model-agnostic explanation methods can help improve transparency and explainability in machine learning systems.

By addressing these ethical considerations throughout the machine learning pipeline, developers can create systems that are more responsible, fair, and user-centric. Ensuring that machine learning systems adhere to ethical principles not only benefits the end-users but also helps build trust in the technology and fosters long-term adoption and success.

3.8 Chapter Summary

In this chapter, we have provided an overview of the methodology and requirements analysis for a general machine learning pipeline. The stages involved in the pipeline include data collection, data preprocessing, model selection and training, evaluation, and deployment and maintenance. Each of these stages plays a crucial role in the development and implementation of an effective machine learning system.

By following the steps and techniques discussed in this chapter, researchers and practitioners can create machine learning systems that cater to the specific needs of their application domain, ensuring that the developed solutions are accurate, reliable, and efficient. This general framework serves as a foundation for building various types of machine learning

systems, ranging from recommendation engines and natural language processing models to computer vision and reinforcement learning applications.

Chapter 4

System Design

4.0 System Design

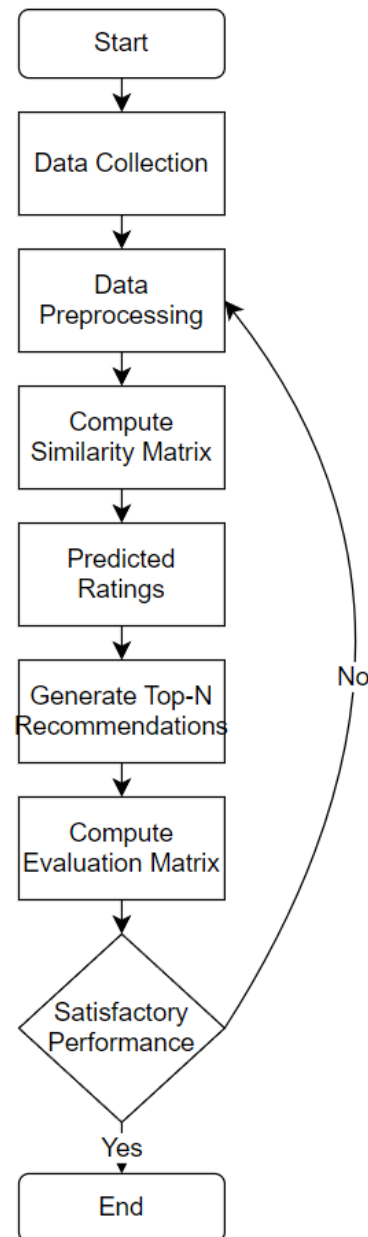


Figure 4.1 Item-Based Collaborative Filtering Pipeline

4.1 Introduction

This chapter will outline the system design for the proposed movie recommendation system that utilizes item-based collaborative filtering. The chapter will discuss the data collection process, data preprocessing, and the implementation of the item-based collaborative filtering algorithm. Additionally, the chapter will outline evaluation metrics and methods for measuring the effectiveness of the recommendation system.

4.2 Data Collection

The MovieLens dataset provided by GroupLens Research is an optimal choice for training and evaluating the item-based collaborative filtering algorithm. This dataset is widely used and publicly available in multiple versions, which offers flexibility for experimentation with different levels of data sparsity and scalability.

Although there are external data sources like IMDb and Wikiwand that can enhance the recommendation system, our focus here is solely on the MovieLens dataset. By concentrating on this dataset, we aim to demonstrate the capabilities of item-based collaborative filtering without introducing additional complexities related to integrating data from external sources.

The MovieLens dataset provides reliable information about movies, including user ratings, genres, and titles, which are crucial for collaborative filtering methods. This dataset allows us to build a recommendation system that captures user preferences solely based on their interactions within the MovieLens platform.

By limiting our project to using only the MovieLens dataset, we can provide a clear and concise demonstration of the item-based collaborative filtering technique. This approach allows us to focus on the core algorithm and its ability to generate personalized movie recommendations without the need for integrating external data.

4.3 Data Preprocessing

Data preprocessing is a crucial step in building a recommendation system, as it ensures that the collected data is clean, consistent, and suitable for analysis. In the context of our item-based collaborative filtering system, the data preprocessing stage involves the following steps:

1. Handling missing values: Missing data can lead to inaccurate recommendations or introduce biases in the model. We need to identify and handle missing values in the dataset. Common techniques for dealing with missing data include:

- a. Removing records with missing values.
- b. Filling missing values with the mean or median value of the respective feature.
- c. Using interpolation or regression techniques to estimate missing values based on other available data.

2. Data transformation: Transforming the data into a suitable format for analysis is essential. This can include:

- a. Converting textual data, such as genres or plot summaries, into numerical representations using techniques like one-hot encoding or natural language processing methods (e.g., TF-IDF or word embeddings).
- b. Normalizing numerical data, such as user ratings or release years, to ensure that all features have a similar scale and do not introduce biases in the model (e.g., Min-max scaling, Z-score normalization or Log transformation).
- c. Aggregating data, if necessary, to create new features or reduce the dimensionality of the dataset. (e.g., PCA or SVD)

3. Feature engineering: Creating new features or modifying existing ones can help improve the performance of the recommendation system. Some examples of feature engineering include:

- a. Extracting relevant information from text data, such as keywords or named entities, which can be used as additional features for the similarity calculation.
- b. Computing interaction features, such as the number of common users who rated a pair of movies or the average rating difference between the movies.

4. Data partitioning: Splitting the dataset into training and testing sets is necessary to evaluate the performance of the recommendation system. This ensures that the model is validated on unseen data and helps prevent overfitting. Typically, a ratio of 70-30 or 80-20 is used for partitioning the dataset, with the larger portion used for training and the smaller portion used for testing.

By performing these data preprocessing steps, we can ensure that our item-based collaborative filtering system is built on clean, consistent, and relevant data, which will help improve the accuracy and reliability of the movie recommendations.

4.4 Implementation of Item-Based Collaborative Filtering

The implementation of the item-based collaborative filtering algorithm involves the following steps:

1. Similarity Computation:

To calculate the similarity between pairs of movies, a similarity metric is used. Common similarity metrics include:

- Pearson correlation coefficient: This measures the linear correlation between two variables (in this case, movie ratings). The coefficient ranges from -1 (negative correlation) to 1 (positive correlation), with 0 indicating no correlation.
- Cosine similarity: This calculates the cosine of the angle between two non-zero vectors (movie rating vectors). The similarity ranges from -1 (completely dissimilar) to 1 (completely similar), with 0 indicating no similarity.

The similarity computation results in a movie similarity matrix, where each cell represents the similarity between a pair of movies.

2. Rating Prediction:

After computing the similarity matrix, the next step is to predict the ratings for the movies that a user has not yet rated. For each user, the predicted rating for an unrated movie is calculated as a weighted sum of the user's ratings for similar movies, where the weights are the similarity values between the unrated movie and the rated movies.

The predicted rating (PR) for user 'u' and movie 'i' can be calculated using the following formula:

$$Pred(u, i) = \bar{r}_u + \frac{\sum_{j \in s(i)} (sim(i, j) \times r_{u,j})}{\sum_{j \in s(i)} sim(i, j)}$$

the predicted rating for user u for item i

average rating of user u

$s(i)$ set of items in the neighbourhood that user u has rated

similarity between item i and item j

active user's u rating of item j

where ' $sim(i, j)$ ' represents the similarity between movie ' i ' and movie ' j ', and $R(u, j)$ is the user ' u 's rating for movie ' j '. The summation is over all movies ' j ' that user ' u ' has rated.

3. Top-N Recommendations:

Once the predicted ratings have been calculated for all unrated movies, the algorithm generates a list of top-N recommendations for each user. This list contains the N movies with the highest predicted ratings for that user. The recommendations are personalized and based on the user's previous movie ratings and the similarities between movies.

In practice, the recommendations might be further refined by considering additional factors such as movie popularity, user preferences, or diversity in genres.

In summary, the item-based collaborative filtering algorithm calculates movie similarities, predicts user ratings for unrated movies, and generates personalized top-N recommendations for each user based on those predictions.

4.5 Evaluation Metrics and Methods

In the context of movie recommendation systems, it's essential to evaluate the performance of the implemented algorithm to ensure it's providing accurate and relevant recommendations to users. Evaluation metrics and methods help quantify the algorithm's performance and can be used for comparison with alternative approaches. Here are some commonly used evaluation metrics and methods for recommendation systems:

1. Mean Absolute Error (MAE):

Mean Absolute Error measures the average absolute difference between the predicted ratings and the actual user ratings. It's a straightforward metric to calculate and interpret. A lower MAE value indicates better prediction accuracy.

2. Root Mean Squared Error (RMSE):

Root Mean Squared Error is the square root of the average squared difference between the predicted ratings and the actual user ratings. It penalizes larger errors more heavily than smaller errors, making it more sensitive to outliers. A lower RMSE value indicates better prediction accuracy.

3. Precision and Recall:

Precision measures the proportion of relevant recommendations among all recommendations provided by the system, while recall measures the proportion of relevant recommendations retrieved out of all relevant items in the dataset. These metrics are commonly used together to evaluate the trade-off between providing a high number of relevant recommendations (recall) and minimizing the number of irrelevant recommendations (precision).

4. F1 Score:

The F1 Score is the harmonic mean of precision and recall, which provides a single metric to evaluate the balance between precision and recall. A higher F1 Score indicates better performance.

5. Normalized Discounted Cumulative Gain (NDCG):

NDCG is a ranking-based evaluation metric that measures the usefulness of recommended items, considering their position in the recommendation list. It takes into account the relevance of each recommended item and discounts the relevance score of items ranked lower in the list. A higher NDCG value indicates better recommendation quality.

To evaluate the performance of the proposed recommendation system, a k-fold cross-validation approach can be used. The dataset will be divided into k equally sized folds, and the algorithm will be trained on k-1 folds and tested on the remaining fold. This process will be repeated k times, with each fold serving as the test set once. The average performance across all k iterations will be reported.

4.6 Chapter Summary

In this chapter, we have discussed the methodology and requirements analysis for implementing an item-based collaborative filtering movie recommendation system. We began by highlighting the importance of data collection and the use of the MovieLens dataset, combined with web scraping techniques like BeautifulSoup4 to collect additional data from sources such as IMDb and Wikiwand. This combination of data sources provides a rich dataset to work with, capturing various aspects of movies that can be utilized in the recommendation algorithm.

The data preprocessing stage is crucial for cleaning and preparing the data for analysis. We discussed various preprocessing steps, such as handling missing values, normalizing numerical data using techniques like Z-score normalization or Min-Max scaling, and aggregating data to create a more comprehensive dataset. We also emphasized the importance of feature engineering, which involves extracting new features from the existing data and using natural language processing techniques to extract sentiment scores from movie reviews.

Feature selection is a critical step in building an effective recommendation system, as it helps to select the most relevant features for the algorithm while minimizing noise and overfitting. We mentioned the use of techniques like LASSO regression, which can help identify features with large coefficients, indicating their importance in the model.

Next, we elaborated on the implementation of item-based collaborative filtering, discussing the steps involved, such as computing movie similarity, predicting ratings, and generating top-N recommendations. We also explained the importance of evaluating the performance of the recommendation algorithm using various evaluation metrics, such as MAE, RMSE, precision, recall, F1 Score, and NDCG, and cross-validation.

In summary, this chapter provides a comprehensive overview of the methodology and requirements analysis for building an item-based collaborative filtering movie recommendation system. By following the steps and techniques discussed, researchers and practitioners can create an effective recommendation system that caters to the preferences and needs of individual users, ultimately improving their movie-watching experience.

Chapter 5

Implementation and Testing

5 Implementation and Testing

In today's data-driven era, recommendation systems play a pivotal role in delivering personalized experiences to users. Movie recommendation systems, in particular, have transformed the way we discover and enjoy films. In this project, we delve into the world of movie recommendation systems, harnessing the power of Singular Value Decomposition (SVD) and Inverse User Frequency Weighting to create a sophisticated recommendation engine.

The primary objective of this project is to design and implement a movie recommendation system that provides users with tailored movie suggestions based on their preferences. We aim to leverage the inherent patterns within a dataset from MovieLens, a well-known movie recommendation benchmark platform. The project is executed using a Jupyter Notebook environment and coded in the versatile Python programming language.

Our recommendation system is built on two key techniques: Singular Value Decomposition (SVD) and Inverse User Frequency Weighting. These methodologies enable us to uncover latent features in the dataset, identify user preferences, and generate precise movie recommendations. The implementation involves importing essential libraries such as Pandas, NumPy, Matplotlib, and Scikit-Learn for data manipulation, visualization, and collaborative filtering. The surprise library is instrumental in integrating SVD-based collaborative filtering into our system.

The cornerstone of our project is the dataset sourced from MovieLens. This dataset encompasses a vast collection of movie ratings provided by users, along with associated movie metadata. By extracting insights from this dataset, we gain valuable information about user preferences and movie characteristics. These insights fuel the recommendation engine, allowing us to suggest movies that align with individual tastes.

5.1 Importing Libraries

In this initial step, we lay the foundation for our movie recommendation system by importing the necessary libraries. We bring in the tools required for different aspects of the project, such as data manipulation (Pandas and NumPy), data visualization (Matplotlib and WordCloud), collaborative filtering (Surprise), and creating a graphical user interface (Tkinter).

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from wordcloud import WordCloud
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import linear_kernel, cosine_similarity
from surprise import Reader, Dataset, SVD
from surprise.model_selection import train_test_split
from surprise.accuracy import rmse, mae
import tkinter as tk
from surprise.model_selection import cross_validate, train_test_split
import pickle
```

5.2 Loading and Exploring the Data

Here, we load our core dataset from MovieLens. This dataset contains user ratings and movie metadata. We utilize the Pandas library to read CSV files and create dataframes, which allow us to store and manipulate the data.

```
In [2]: # Load the dataset
movies = pd.read_csv('movies.csv')
ratings = pd.read_csv('ratings.csv')
```

By using functions like `info()`, `dtypes`, `shape`, and `head()`, we get a comprehensive understanding of the dataset's structure, data types, dimensions, and initial rows. This exploration sets the stage for subsequent analyses. By using these function

`info()`: This function provides a concise summary of the dataframe, including the number of non-null entries, data types, and memory usage. It's useful for quickly checking for missing data and understanding the data types. Upon observation, we notice that the movie dataset contains both integers and objects, while the rating dataset consists of integers and float values.

In [3]: `movies.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9742 entries, 0 to 9741
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   movieId     9742 non-null   int64
1   title       9742 non-null   object
2   genres      9742 non-null   object
dtypes: int64(1), object(2)
memory usage: 228.5+ KB
```

In [5]: `ratings.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100836 entries, 0 to 100835
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   userId      100836 non-null  int64
1   movieId     100836 non-null  int64
2   rating      100836 non-null  float64
3   timestamp   100836 non-null  int64
dtypes: float64(1), int64(3)
memory usage: 3.1 MB
```

`dtypes`: This attribute displays the data types of each column in the dataframe. It's especially useful to ensure that the data types match our expectations.

```
In [4]: movies.dtypes
```

```
Out[4]: movieId      int64  
        title        object  
        genres       object  
        dtype: object
```

```
In [6]: ratings.dtypes
```

```
Out[6]: userId        int64  
        movieId        int64  
        rating         float64  
        timestamp      int64  
        dtype: object
```

shape: This attribute gives the dimensions of the dataframe (rows and columns), helping us understand the dataset's size. Specifically, the movie dataset comprises 9742 rows and 3 columns, while the ratings dataset consists of 100836 rows and 4 columns.

```
In [7]: movies.shape
```

```
Out[7]: (9742, 3)
```

```
In [8]: ratings.shape
```

```
Out[8]: (100836, 4)
```

describe() function generates a comprehensive summary of descriptive statistics for the numerical columns in the dataframe. It offers information regarding various aspects of the data, such as measures of central tendency (mean, median), dispersion (standard deviation, min-max range), and distribution (quartiles) for each numeric column.

For the movieId column, it is noteworthy that the smallest number is 1, and the largest number is 193609. The discrepancy between the maximum movieId number and the count (9742) suggests that there are gaps in the sequence of movieIds.

```
In [9]: movies.describe()
```

```
Out[9]:
```

	movieId
count	9742.000000
mean	42200.353623
std	52160.494854
min	1.000000
25%	3248.250000
50%	7300.000000
75%	76232.000000
max	193609.000000

Regarding the `userId` column, it exhibits a count of 100836. The range of `userId` values, from a minimum of 1 to a maximum of 610, indicates that a single user has provided multiple ratings for a diverse range of movies.

```
In [10]: ratings.describe()
```

```
Out[10]:
```

	userId	movieId	rating	timestamp
count	100836.000000	100836.000000	100836.000000	1.008360e+05
mean	326.127564	19435.295718	3.501557	1.205946e+09
std	182.618491	35530.987199	1.042529	2.162610e+08
min	1.000000	1.000000	0.500000	8.281246e+08
25%	177.000000	1199.000000	3.000000	1.019124e+09
50%	325.000000	2991.000000	3.500000	1.186087e+09
75%	477.000000	8122.000000	4.000000	1.435994e+09
max	610.000000	193609.000000	5.000000	1.537799e+09

`head()`: This function displays the first few rows of the dataframe, allowing us to get a glimpse of the actual data. It's helpful for verifying that the data has been loaded correctly and understanding its structure.

```
In [11]: movies.head()
```

```
Out[11]:
```

	movieId	title	genres
0	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
1	2	Jumanji (1995)	Adventure Children Fantasy
2	3	Grumpier Old Men (1995)	Comedy Romance
3	4	Waiting to Exhale (1995)	Comedy Drama Romance
4	5	Father of the Bride Part II (1995)	Comedy

```
In [12]: ratings.head()
```

```
Out[12]:
```

	userId	movieId	rating	timestamp
0	1	1	4.0	964982703
1	1	3	4.0	964981247
2	1	6	4.0	964982224
3	1	47	5.0	964983815
4	1	50	5.0	964982931

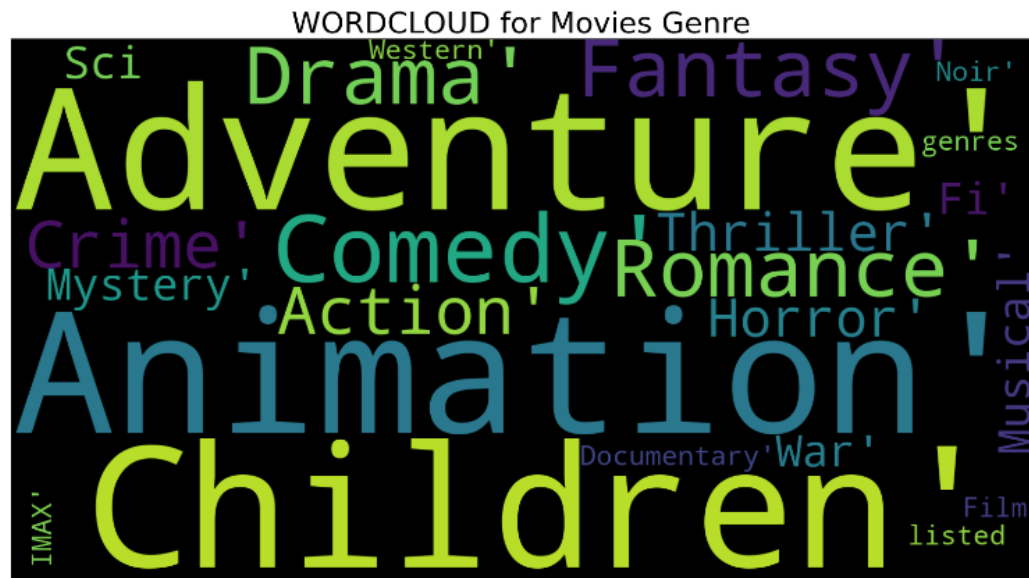
5.3 Data Visualization

In this step, we delve into data visualization techniques to quickly gain insights. We utilize the WordCloud library to generate word clouds representing movie genres and titles. Word clouds visually depict the most common elements in the dataset, aiding us in identifying frequently occurring genres and popular movie titles. These visualizations provide a rapid overview of the dataset's thematic content.

From the movie genre word cloud, we can discern that the top three most frequently appearing words are "Adventure," "Animation," and "Children."

```
In [15]: plt.figure(figsize=(30,10))
plt.axis('off')
plt.title('WORDCLOUD for Movies Genre',fontsize=30)
plt.imshow(wordcloud_genre)
```

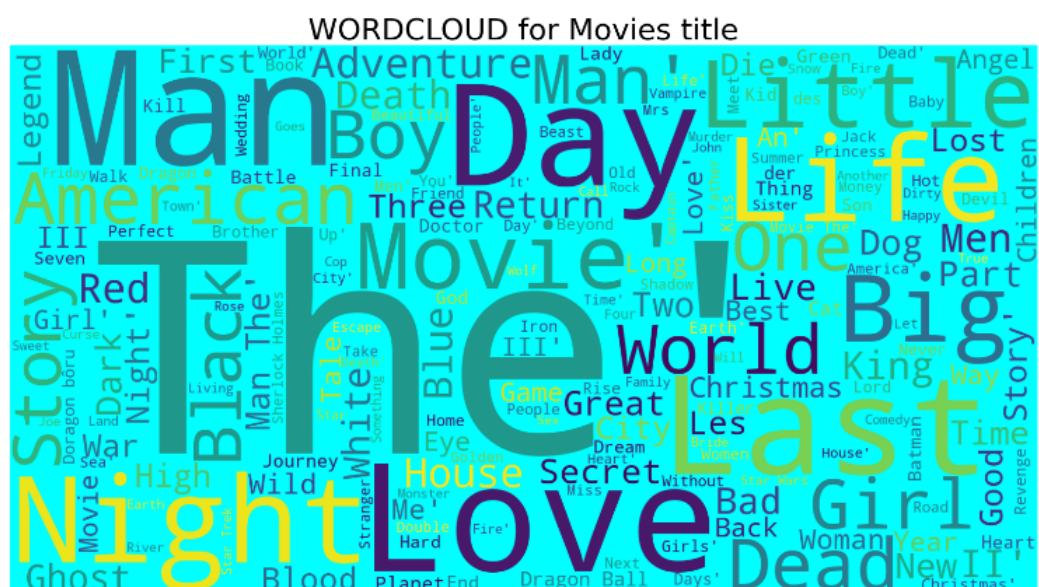
```
Out[15]: <matplotlib.image.AxesImage at 0x1be2482c3a0>
```



Upon analyzing the movie title word cloud, we observe that the nine words that prominently stand out are "Man," "Day," "Life," "Little," "The," "Night," "Love," "Last," and "Big."

```
In [16]: plt.figure(figsize=(30,10))
plt.axis('off')
plt.title('WORDCLOUD for Movies title',fontsize=30)
plt.imshow(wordcloud title)
```

```
Out[16]: <matplotlib.image.AxesImage at 0x1be24f611f0>
```



`pd.merge()`: The function from the Pandas library to combine two dataframes, ratings and movies, based on the 'movieId' column. The how parameter is set to 'left', indicating a left join operation. This means that all rows from the ratings dataframe will be retained, and additional information from the movies dataframe will be added based on matching 'movieId' values.

```
In [17]: df=pd.merge(ratings,movies, how='left',on='movieId')
df.head()
```

```
Out[17]:
```

	userId	movieId	rating	timestamp	title	genres
0	1	1	4.0	964982703	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
1	1	3	4.0	964981247	Grumpier Old Men (1995)	Comedy Romance
2	1	6	4.0	964982224	Heat (1995)	Action Crime Thriller
3	1	47	5.0	964983815	Seven (a.k.a. Se7en) (1995)	Mystery Thriller
4	1	50	5.0	964982931	Usual Suspects, The (1995)	Crime Mystery Thriller

The `groupby` function is to group the original dataframe `df` based on the 'title' column. You then select the 'rating' column and calculate the sum of ratings for each movie title. Essentially, you are aggregating the ratings for each movie title to find the total rating sum for each title.

This function retrieves the 20 rows with the largest values in the 'rating' column from the `df1` dataframe. In other words, you are selecting the top 20 movies with the highest total rating sums.

The `head()` function. It gives you a quick view of the top-rated movies along with their total rating sums.

```
In [18]: df1=df.groupby(['title'])[['rating']].sum()  
high_rated=df1.nlargest(20,'rating')  
high_rated.head()
```

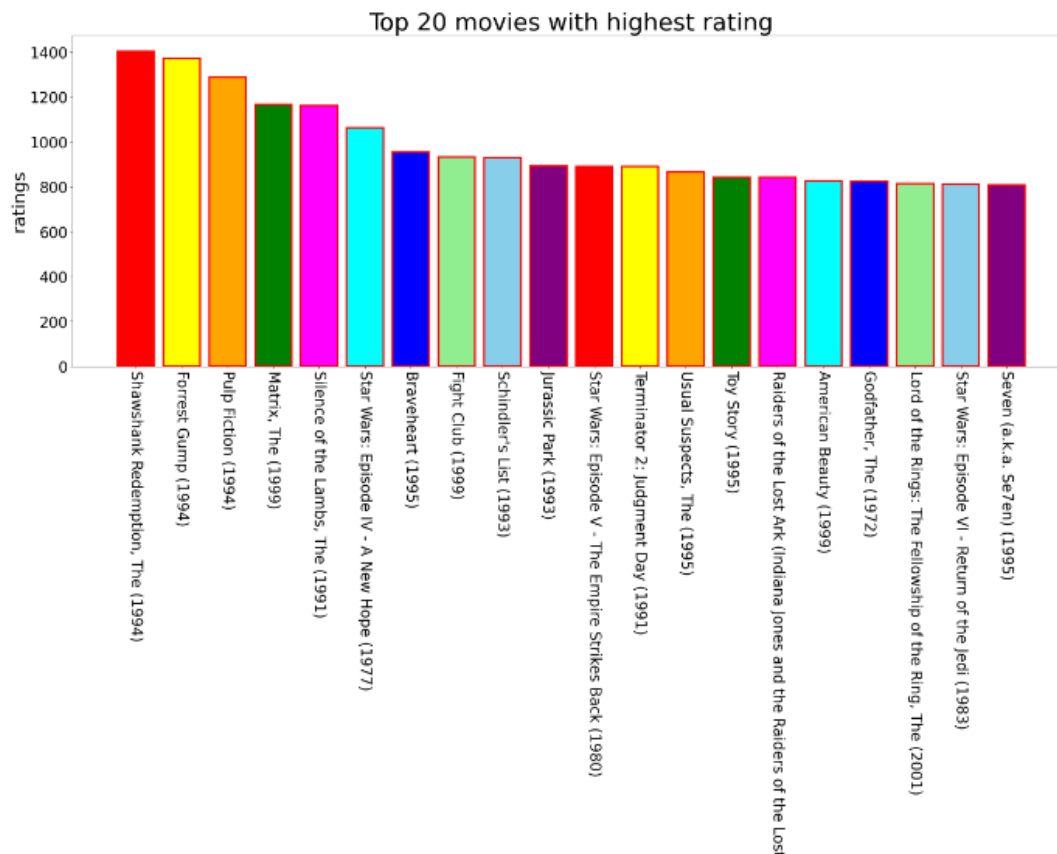
```
Out[18]:
```

	rating
title	
Shawshank Redemption, The (1994)	1404.0
Forrest Gump (1994)	1370.0
Pulp Fiction (1994)	1288.5
Matrix, The (1999)	1165.5
Silence of the Lambs, The (1991)	1161.0

From the chart, we can see the top 20 movies with the highest ratings. Each bar represents a movie title, and its height corresponds to the rating of that movie. By examining this plot, you can quickly identify which movies have received the highest ratings, allowing you to discover popular and highly regarded films in your dataset

```
In [19]: plt.figure(figsize=(30,10))
plt.title('Top 20 movies with highest rating',fontsize=40)
colors=['red','yellow','orange','green','magenta','cyan','blue','lightgreen','s
plt.ylabel('ratings',fontsize=30)
plt.xticks(fontsize=25,rotation=-90)
plt.xlabel('movies title',fontsize=30)
plt.yticks(fontsize=25)
plt.bar(high_rated.index,high_rated['rating'],linewidth=3,edgecolor='red',color
```

Out[19]: <BarContainer object of 20 artists>



It utilizes the groupby function to group the movie titles and count the ratings. The objective is to identify the top 20 movies based on the count of ratings they have received. This information aids in recognizing the most popular and widely-watched movies within your dataset, as determined by user ratings.


```
In [20]: df2=df.groupby('title')[['rating']].count()  
rating_count_20=df2.nlargest(20,'rating')  
rating_count_20.head()
```

```
Out[20]:
```

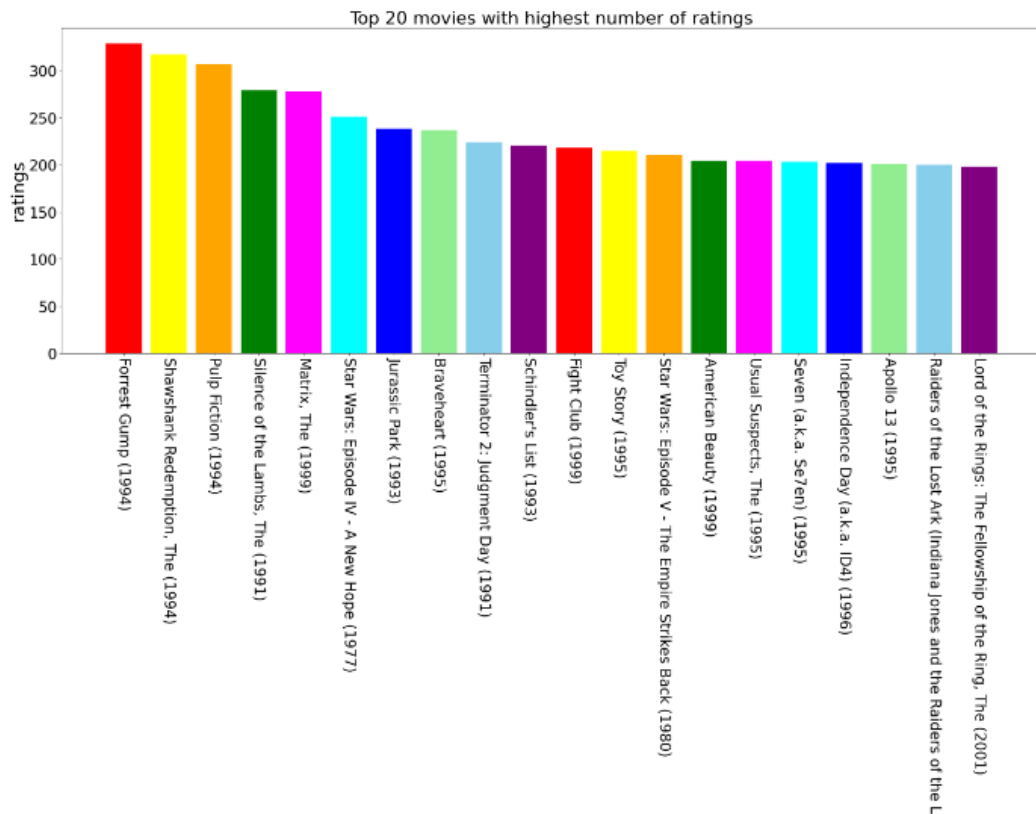
	rating
title	
Forrest Gump (1994)	329
Shawshank Redemption, The (1994)	317
Pulp Fiction (1994)	307
Silence of the Lambs, The (1991)	279
Matrix, The (1999)	278

The resulting bar chart visually represents the top 20 movies with the highest number of ratings. Each bar corresponds to a movie title, and its height represents the number of ratings received by that movie. Overall, this visualization allows us to quickly identify which movies have garnered the most ratings, providing insights into the preferences and popularity of certain films in the dataset.

```
In [21]: plt.figure(figsize=(30,10))
plt.title('Top 20 movies with highest number of ratings',fontsize=30)
plt.xticks(fontsize=25,rotation=-90)
plt.yticks(fontsize=25)
plt.xlabel('movies title',fontsize=30)
plt.ylabel('ratings',fontsize=30)

colors=['red','yellow','orange','green','magenta','cyan','blue','lightgreen','skyblue','purple']
plt.bar(rating_count_20.index,rating_count_20.rating,color= colors)
```

Out[21]: <BarContainer object of 20 artists>



The code groups the dataset by ratings, counts the number of movies for each rating, and then visualizes the results using a bar plot. This visualization offers insights into the distribution of movie ratings across different values. The bar plot illustrates the quantity of movies falling into each rating category, enabling an understanding of the prevalence of specific rating levels. For instance, you can swiftly discern which rating values are more commonly assigned and whether users tend to provide higher or lower ratings to movies. Notably, the chart highlights that the highest user rating is 4.0.

```
In [22]: rating_count = df['rating'].value_counts().sort_index()
rating_count = pd.DataFrame({'rating': rating_count.index, 'movie_count': rating_count.values})
print(rating_count)
```

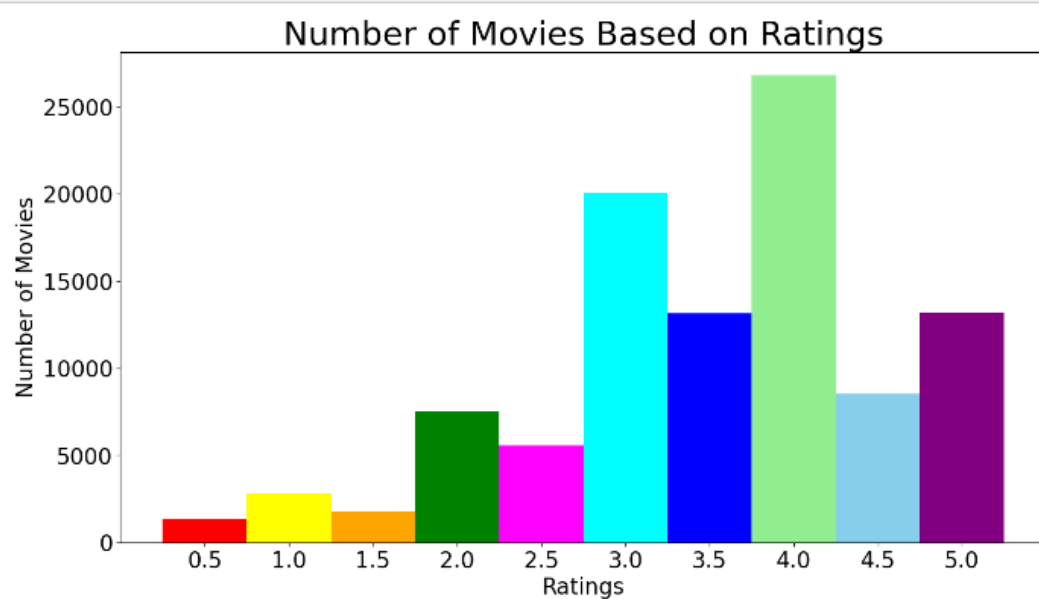
	rating	movie_count
0	0.5	1370
1	1.0	2811
2	1.5	1791
3	2.0	7551
4	2.5	5550
5	3.0	20047
6	3.5	13136
7	4.0	26818
8	4.5	8551
9	5.0	13211

```
In [23]: plt.figure(figsize=(15, 8)) # Adjust the figure size as needed
plt.title('Number of Movies Based on Ratings', fontsize=30)
plt.xticks(rating_count.rating, fontsize=20)
plt.yticks(fontsize=20)
plt.xlabel('Ratings', fontsize=20)
plt.ylabel('Number of Movies', fontsize=20)

colors = ['red', 'yellow', 'orange', 'green', 'magenta', 'cyan', 'blue', 'lightgreen', 'purple']
bar_width = 0.5

#x_labels = [str(rating) for rating in rating_count.rating]

plt.bar(rating_count.rating, rating_count.movie_count, color=colors, width=bar_width)
#plt.xticks(rating_count.rating, x_labels) # Set the x-axis tick positions and labels
plt.show()
```



5.4 Item-Based Collaborative Filtering

This step involves the heart of collaborative filtering. We calculate the similarity between movies based on user ratings. The result is a similarity matrix that helps us identify movies similar to those a user has already rated positively. By recommending movies similar to ones a user has enjoyed, we enhance their movie-watching experience with personalized suggestions.

`pivot_table`: It reshapes the original dataframe `df` so that each row represents a user, each column represents a movie, and the values in the table are user ratings for each movie. The `fill_value=0` parameter ensures that missing ratings are filled with zeros.

it creates a subset of the user-movie ratings matrix for easy visualization and analysis. By looking at the printed subset, you can observe how these specific users have rated these specific movies. It provides a glimpse into how the ratings are distributed within this subset of the data. The subset is used here for demonstration and initial exploration of the data structure.

```
In [24]: df = pd.read_csv('ratings.csv')

user_movie_ratings = df.pivot_table(
    index='userId',
    columns='movieId',
    values='rating',
    fill_value=0
)

subset_user_movie_ratings = user_movie_ratings.iloc[:10, :10] # Subset for visu
print("User-Movie Ratings Matrix (Subset):")
print(subset_user_movie_ratings)
```

```
User-Movie Ratings Matrix (Subset):
movieId  1  2  3  4  5  6  7  8  9  10
userId
1         4.0  0.0  4.0  0.0  0.0  4.0  0.0  0  0.0  0.0
2         0.0  0.0  0.0  0.0  0.0  0.0  0.0  0  0.0  0.0
3         0.0  0.0  0.0  0.0  0.0  0.0  0.0  0  0.0  0.0
4         0.0  0.0  0.0  0.0  0.0  0.0  0.0  0  0.0  0.0
5         4.0  0.0  0.0  0.0  0.0  0.0  0.0  0  0.0  0.0
6         0.0  4.0  5.0  3.0  5.0  4.0  4.0  3  0.0  3.0
7         4.5  0.0  0.0  0.0  0.0  0.0  0.0  0  0.0  0.0
8         0.0  4.0  0.0  0.0  0.0  0.0  0.0  0  0.0  2.0
9         0.0  0.0  0.0  0.0  0.0  0.0  0.0  0  0.0  0.0
10        0.0  0.0  0.0  0.0  0.0  0.0  0.0  0  0.0  0.0
```

In this part of the code, it calculate the similarity between movies (items) using the cosine similarity metric. It transpose the `user_movie_ratings` dataframe (users as rows and movies as

columns) to get movies as rows and users as columns, making it suitable for item-item similarity calculations. The resulting `item_similarity` matrix contains cosine similarity values, where each entry (i, j) represents the similarity between movie i and movie j .

The findings from this code reveal how similar each movie is to every other movie in the dataset based on user ratings. A high similarity value indicates that two movies have similar rating patterns across users, suggesting that they are likely to appeal to similar audiences. This information is fundamental for building recommendation systems that suggest movies similar to ones a user has already enjoyed.

In [25]:

```
# Calculate item-item similarity using cosine similarity
item_similarity = cosine_similarity(user_movie_ratings.T)

# Create a DataFrame from the item similarity matrix
item_similarity_df = pd.DataFrame(item_similarity, index=user_movie_ratings.columns)
print("\nItem Similarity Matrix (Subset):")
print(item_similarity_df.iloc[:10, :7])
```

```
Item Similarity Matrix (Subset):
movieId      1      2      3      4      5      6      7
movieId
1      1.000000  0.410562  0.296917  0.035573  0.308762  0.376316  0.277491
2      0.410562  1.000000  0.282438  0.106415  0.287795  0.297009  0.228576
3      0.296917  0.282438  1.000000  0.092406  0.417802  0.284257  0.402831
4      0.035573  0.106415  0.092406  1.000000  0.188376  0.089685  0.275035
5      0.308762  0.287795  0.417802  0.188376  1.000000  0.298969  0.474002
6      0.376316  0.297009  0.284257  0.089685  0.298969  1.000000  0.244105
7      0.277491  0.228576  0.402831  0.275035  0.474002  0.244105  1.000000
8      0.131629  0.172498  0.313434  0.158022  0.283523  0.147562  0.273757
9      0.232586  0.044835  0.304840  0.000000  0.335058  0.214088  0.162000
10     0.395573  0.417693  0.242954  0.095598  0.218061  0.386414  0.238949
```

Inside the function, it first extracts the user's actual ratings from the `user_movie_ratings` matrix. It then calculates recommendation scores by taking a dot product between the user's ratings and the item similarity scores. This step identifies how similar the user's preferences are to other movies in the dataset.

The code identifies movies that the user has not already rated (i.e., unrated movies) and sorts them by recommendation scores in descending order. This step helps in suggesting movies that the user might like but hasn't seen yet.

Finally, the code provides an example of how to use this evaluation function. It evaluates the recommendations for user_id 1 and prints the MAE and RMSE values. Lower MAE and RMSE values indicate better recommendations. These metrics are essential for assessing and fine-tuning recommendation algorithms to enhance user satisfaction.

```
In [26]: def evaluate_recommendations(user_id, user_movie_ratings, item_similarity_df):
        user_ratings = user_movie_ratings.loc[user_id]

        # Calculate the weighted average of similar movies
        recommendation_scores = item_similarity_df.dot(user_ratings)

        # Exclude already rated movies from recommendations
        unrated_movies = user_ratings[user_ratings == 0].index
        recommended_movies = recommendation_scores[unrated_movies].sort_values(ascending=True)

        # Example: assuming actual ratings are known for evaluation
        actual_ratings = user_movie_ratings.loc[user_id]

        mae = mean_absolute_error(actual_ratings, recommendation_scores)
        rmse = mean_squared_error(actual_ratings, recommendation_scores, squared=False)

        return mae, rmse

# Example usage for evaluation
user_id = 1
mae, rmse = evaluate_recommendations(user_id, user_movie_ratings, item_similarity_df)
print(f'MAE: {mae}, RMSE: {rmse}')
```

MAE: 94.17461678850529, RMSE: 116.86523652484699

5.5 Item-Based Collaborative Filtering with SVD

Singular Value Decomposition (SVD) is a powerful matrix factorization technique. In this stage, we employ the surprise library to implement SVD-based collaborative filtering. We create a dataset and split it into training and test sets. The SVD model is trained on the training data, and the performance is evaluated using metrics like RMSE and MAE. The trained model can then be used for making predictions and generating recommendations.

This code segment focuses on constructing and assessing a recommendation system utilizing Singular Value Decomposition (SVD) coupled with Inverse User Frequency Weighting. It begins by initializing a reader and constructing a dataset, followed by the partitioning of data into training and testing sets. Additionally, it defines hyperparameters for the SVD model, which significantly impact its performance and convergence. The code proceeds to create and train the SVD model, preserving it for future use without the need for retraining. Finally, it calculates item-item similarity.

```
# Initialize a Reader
reader = Reader()

# Create a Dataset
data = Dataset.load_from_df(df[['userId', 'movieId', 'rating']], reader)

# Split the data into a training set and a test set
trainset, testset = train_test_split(data, test_size=0.2, random_state=42)

# Define hyperparameters for SVD
params = {'n_factors': 50, # Number of latent factors
          'n_epochs': 20, # Number of iterations
          'lr_all': 0.005, # Learning rate for all parameters
          'reg_all': 0.02} # Regularization term for all parameters

# Create and train the SVD model with hyperparameters
model = SVD(**params)
model.fit(trainset)

# Save the trained model using pickle
model_filename = 'item_based_model.pkl'
with open(model_filename, 'wb') as model_file:
    pickle.dump(item_similarity_df, model_file)

# Calculate item-item similarity using cosine similarity
user_movie_ratings = df.pivot_table(
    index='userId',
    columns='movieId',
    values='rating',
    fill_value=0
)
```

Apply with Inverse User Frequency Weighting. This step aims to improve the recommendations by considering less-commonly rated movies. It takes a user's ID, their movie ratings, and the item similarity matrix as inputs. It calculates Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) to evaluate the performance of the recommender system. The example with SVD, the RMSE and MAE is lower which indicate better recommendations.

Overall, this code segment demonstrates the training, evaluation, and enhancement of a recommendation system using SVD with weighted item similarity. It paves the way for making personalized movie recommendations based on user preferences and item similarity.

```
# Apply Inverse User Frequency Weighting
user_freq = user_movie_ratings.astype(bool).sum(axis=0)
user_freq_weights = 1 / (1 + np.log(user_freq))
user_freq_weights = user_freq_weights.values.reshape(-1, 1) # Reshape to match
item_similarity_weighted = user_freq_weights * cosine_similarity(user_movie_ratings)

# Create a DataFrame from the weighted item similarity matrix
item_similarity_df = pd.DataFrame(item_similarity_weighted, index=user_movie_ratings.index)

def evaluate_recommendations(user_id, user_movie_ratings, item_similarity_df):
    user_ratings = user_movie_ratings.loc[user_id]

    # Calculate the weighted average of similar movies
    recommendation_scores = item_similarity_df.dot(user_ratings)

    # Exclude already rated movies from recommendations
    unrated_movies = user_ratings[user_ratings == 0].index
    recommended_movies = recommendation_scores[unrated_movies].sort_values(ascending=False)

    # Example: assuming actual ratings are known for evaluation
    actual_ratings = user_movie_ratings.loc[user_id]

    mae_score = np.mean(np.abs(actual_ratings.values - recommendation_scores.values))
    rmse_score = np.sqrt(np.mean(np.square(actual_ratings.values - recommendation_scores.values)))

    return mae_score, rmse_score

# Example usage for evaluation
user_id = 1
mae_score, rmse_score = evaluate_recommendations(user_id, user_movie_ratings, item_similarity_df)
print(f'MAE: {mae_score}, RMSE: {rmse_score}')
```

MAE: 41.723848502069394, RMSE: 48.767783364210395

These visualizations help in understanding how different weighting strategies can impact item similarity and influence recommendations in a collaborative filtering system. The findings can guide decisions on which weighting techniques to use based on the specific goals of the recommendation system.

Inverse User Frequency Weighting aims to give more weight to less frequently rated items. It does so to potentially make recommendations more diverse by considering items that might be overlooked when recommendations are solely based on highly rated or frequently rated items.

The darkness of the cells in the weighted graph reflects the degree of similarity between items, with darker cells indicating weaker similarity while lighter means higher similarity.


```
# Choose a subset of items for visualization (e.g., first 50 items)
subset_items = user_movie_ratings.columns[:50]

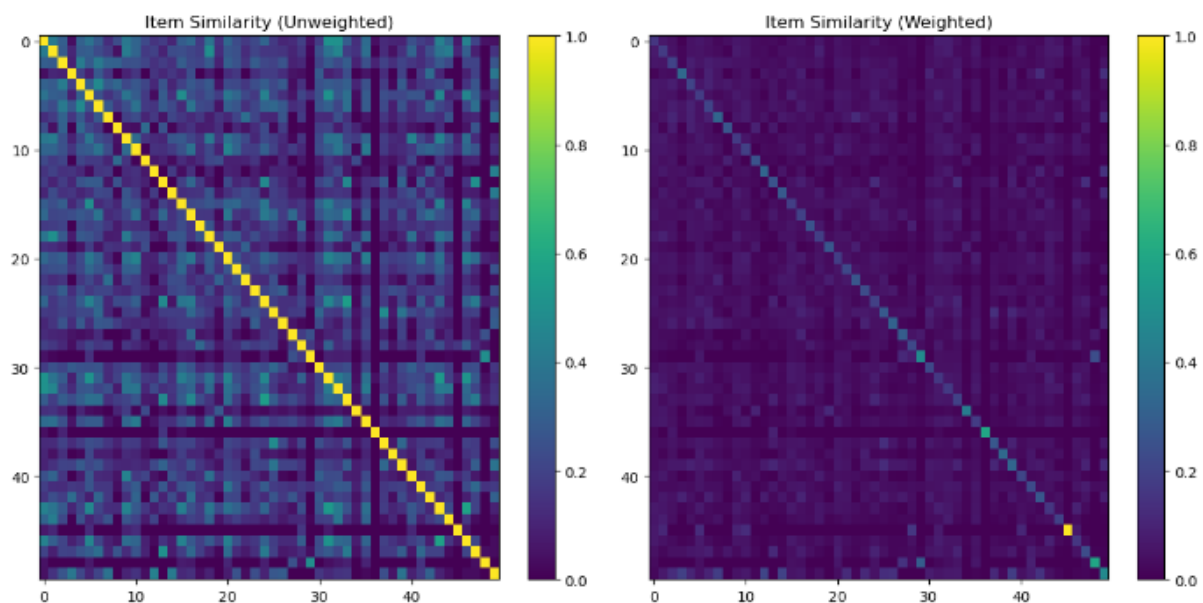
# Calculate cosine similarity without weighting
item_similarity = cosine_similarity(user_movie_ratings[subset_items].T)

# Apply Inverse User Frequency Weighting to the subset
subset_user_freq = user_freq[subset_items]
subset_user_freq_weights = 1 / (1 + np.log(subset_user_freq))
subset_user_freq_weights = subset_user_freq_weights.values.reshape(-1, 1)
subset_item_similarity_weighted = subset_user_freq_weights * item_similarity

# Plot the similarity matrices side by side
plt.figure(figsize=(12, 6))
plt.subplot(1, 2, 1)
plt.imshow(item_similarity, cmap='viridis', aspect='auto')
plt.title('Item Similarity (Unweighted)')
plt.colorbar()

plt.subplot(1, 2, 2)
plt.imshow(subset_item_similarity_weighted, cmap='viridis', aspect='auto')
plt.title('Item Similarity (Weighted)')
plt.colorbar()

plt.tight_layout()
plt.show()
```



5.6 User-Based Collaborative Filtering

In this phase, the recommendation system's performance is assessed through cross-validation. The SVD model, which has been well trained, is employed to generate predictions, which are subsequently compared with actual ratings found in the test set. Notably, the RMSE and

MAE metrics reveal significantly lower scores compared to those obtained with item-based collaborative filtering. This indicates the effectiveness of user-based collaborative filtering in delivering accurate recommendations.

In summary, this code segment serves the dual purpose of establishing and rigorously evaluating a user-based collaborative filtering recommendation system using SVD. This approach lays a strong foundation for crafting personalized movie recommendations tailored to individual user preferences.

```
# Load your dataset
df = pd.read_csv('ratings.csv')

# Initialize a Reader
reader = Reader()

# Create a Dataset
data = Dataset.load_from_df(df[['userId', 'movieId', 'rating']], reader)

# Split the data into a training set and a test set
trainset, testset = train_test_split(data, test_size=0.2, random_state=42)

# Define hyperparameters for SVD
params = {'n_factors': 50, 'n_epochs': 20, 'lr_all': 0.005, 'reg_all': 0.02}

# Create and train the SVD model with hyperparameters
model = SVD(**params)
model.fit(trainset)

# Save the trained model
model_filename = 'user_based_model.pkl'
with open(model_filename, 'wb') as model_file:
    pickle.dump(model, model_file)

# Perform cross-validation
cv_results = cross_validate(model, data, measures=['RMSE', 'MAE'], cv=5, verbose=1)
print("Cross-validation RMSE:", cv_results['test_rmse'].mean())
print("Cross-validation MAE:", cv_results['test_mae'].mean())
```

Evaluating RMSE, MAE of algorithm SVD on 5 split(s).

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	Std
RMSE (testset)	0.8747	0.8694	0.8762	0.8627	0.8724	0.8711	0.0048
MAE (testset)	0.6705	0.6713	0.6721	0.6641	0.6701	0.6696	0.0028
Fit time	2.21	2.47	2.21	2.20	2.23	2.26	0.10
Test time	0.10	0.18	0.10	0.18	0.10	0.13	0.04
Cross-validation RMSE: 0.8710934965401436							
Cross-validation MAE: 0.669623539826074							

5.7 Deployment

In the final stage, we create a user-friendly Graphical User Interface (GUI) using the Tkinter library. This GUI allows users to input movies they've watched, and the system generates personalized recommendations based on their input. This interactive interface makes the recommendation system accessible and engaging for users, enhancing their movie discovery process.

By following these steps in detail, we build a comprehensive movie recommendation system that harnesses collaborative filtering techniques, data visualization, and a user-friendly interface to provide tailored movie suggestions to users.

5.7.1 User-based movie recommendation system based on UserId

This code segment implements a user-based movie recommendation system with a user-friendly Graphical User Interface (GUI) using the Tkinter library. It begins by loading a pre-trained recommendation model and then creates a GUI where users can input their user ID. Upon clicking the "Get Recommendations" button, the code predicts movie ratings for all movies in the dataset for the chosen user and provides the top 10 movie recommendations based on estimated ratings. This interactive interface enhances the usability of the recommendation system, allowing users to easily access personalized movie suggestions tailored to their preferences. It bridges the gap between sophisticated collaborative filtering algorithms and end-users, making movie recommendations more accessible and engaging.

```

# Load the saved model
model_filename = 'user_based_model.pkl'
with open(model_filename, 'rb') as model_file:
    loaded_model = pickle.load(model_file)

movie_ids = ratings['movieId'].unique()

# Create the main window
root = tk.Tk()
root.title("Movie Recommendation System")

# Create and place widgets
max_user_id = df['userId'].max()
user_id_label = tk.Label(root, text=f"Enter your user ID (Maximum User ID: {max_user_id})")
user_id_label.pack()

user_id_entry = tk.Entry(root)
user_id_entry.pack()

get_recommendations_button = tk.Button(root, text="Get Recommendations")

def get_recommendations():
    user_id = user_id_entry.get()

    if user_id:
        user_id = int(user_id)

        # Predict ratings for all movies for the chosen user
        user_ratings = [{'movieId': movie_id, 'rating': model.predict(user_id, movie_id)} for movie_id in movie_ids]

        # Sort the predicted ratings in descending order and get top 10 recommendations
        top_recommendations = sorted(user_ratings, key=lambda x: x['rating'], reverse=True)[:10]

        recommendations_text.delete(1.0, tk.END) # Clear previous recommendations

        for idx, rec in enumerate(top_recommendations, start=1):
            recommendations_text.insert(tk.END, f"{idx}. MovieId: {rec['movieId']} Rating: {rec['rating']}\n")

get_recommendations_button.config(command=get_recommendations)
get_recommendations_button.pack()

recommendations_text = tk.Text(root, height=20, width=70)
recommendations_text.pack()

root.mainloop()

```

5.7.2 User-based movie recommendation system based on Genre

Users can input their user ID and select a preferred movie genre, and the system employs a pre-trained Singular Value Decomposition (SVD) model to offer personalized movie recommendations within the chosen genre. These recommendations are not only genre-specific but also include predicted ratings, allowing users to make informed choices.

The GUI provides an intuitive way for users to explore and discover top-rated movies that align with their preferences, enhancing their movie-watching experience.

```
# Read data and preprocess
df = pd.read_csv('ratings.csv')
movies = pd.read_csv('movies.csv')
mlb = MultiLabelBinarizer()
genre_matrix = mlb.fit_transform(movies['genres'].apply(lambda x: x.split('|')))
genre_sparse = csr_matrix(genre_matrix)

# Calculate genre similarity matrix
genre_similarity = cosine_similarity(genre_sparse)

# Load the trained SVD model
with open('user_based_model.pkl', 'rb') as model_file:
    model = pickle.load(model_file)

# Create a Tkinter GUI
root = Tk()
root.title("Movie Recommender based on genre")

max_user_id = df['userId'].max()
user_label = Label(root, text=f"Enter your user ID (Maximum User ID: {max_user_id})")
user_label.pack()

user_entry = Entry(root)
user_entry.pack()

genre_label = Label(root, text="Select a genre:")
genre_label.pack()

# Get the unique genres
unique_genres = mlb.classes_
genre_var = StringVar(value=unique_genres[0])
genre_dropdown = OptionMenu(root, genre_var, *unique_genres)
genre_dropdown.pack()

recommendations_text = Text(root, height=20, width=100)
recommendations_text.pack()
```

```
def recommend_movies():
    user_id = int(user_entry.get())
    input_genre = genre_var.get()

    # Find the index of the selected genre
    input_genre_idx = mlb.classes_.tolist().index(input_genre)

    # Get movies that match the selected genre
    genre_movies = movies[movies['genres'].str.contains(input_genre)]

    # Use the trained model to predict ratings for the genre movies
    genre_movie_ids = genre_movies['movieId'].tolist()
    genre_movie_ratings = [(user_id, movie_id, model.predict(user_id, movie_id)
                           .est) for movie_id in genre_movie_ids]

    # Sort the recommended genre movies by predicted rating
    genre_movie_ratings.sort(key=lambda x: x[2], reverse=True)

    # Display the recommendations
    recommendations_text.delete(1.0, END)
    for idx, (user_id, movie_id, rating) in enumerate(genre_movie_ratings[:10],
    movie_title = movies[movies['movieId'] == movie_id]['title'].iloc[0]
    recommendations_text.insert(END, f"{idx}. Movie: {movie_title},
                                Predicted Rating: {rating:.2f}\n\n")

recommend_button = Button(root, text="Recommend Movies",
                           command=recommend_movies)
recommend_button.pack()

root.mainloop()
```

5.7.3 Item-based movie recommendation system based on MovieId and rating

Users can input the movies they have watched along with their ratings. Upon clicking the "Get Recommendations" button, the system processes this input, leverages a pre-trained item-based collaborative filtering model, and swiftly generates personalized movie recommendations. These recommendations are displayed within the interface, presenting users with a curated list of films that align with their viewing history. The code's efficiency lies in its seamless integration of user interaction and data-driven recommendation, offering a user-centric experience that enhances the exploration of new movie choices.

```
# Load your cleaned DataFrame
# Replace this with Loading your actual cleaned DataFrame
df_clean = pd.read_csv('ratings.csv')

# Create the main window
root = tk.Tk()
root.title("Movie Recommendation System for unwatched movie")

# Create and place widgets
input_label = tk.Label(root, text="Enter the movies you have watched (MovieId:Ra
input_label.pack()

range_label = tk.Label(root, text="MovieId range: {}-{}".format(df_clean['movieId
range_label.pack()

input_entry = tk.Entry(root)
input_entry.pack()

recommend_button = tk.Button(root, text="Get Recommendations")

def get_recommendations():
    user_input = input_entry.get()
    user_input = user_input.strip().split(',')

    if user_input:
        user_ratings = []
        for item in user_input:
            try:
                movie_id, rating = map(float, item.split(':'))
                user_ratings.append({'userId': 611, 'movieId': int(movie_id), 'r
            except ValueError:
                messagebox.showerror("Error", "Invalid input format. Use MovieId
                return

    user_df = pd.DataFrame(user_ratings)
    print('user_df:\n', user_df)
    user_watched_movie_ids = set(user_df['movieId'])
```

```

# Load the item-based collaborative filtering model
model_filename = 'item_based_model.pkl'
with open(model_filename, 'rb') as model_file:
    item_similarity_df = pickle.load(model_file)

# Generate recommendations for unwatched movies
unwatched_movie_ids = list(set(df_clean['movieId'].unique()) -
                             user_watched_movie_ids)

# Initialize an empty DataFrame to store recommendations
recommendations = pd.DataFrame(index=unwatched_movie_ids,
                                columns=['Recommendation Score'])

for movie_id in unwatched_movie_ids:
    # Calculate recommendation score for each unwatched movie
    similar_movies = item_similarity_df[movie_id]
    user_watched_ratings = user_df[user_df['movieId'].isin(similar_movies)]
    user_weighted_ratings = user_watched_ratings['rating'] * similar_movies
    recommendation_score = user_weighted_ratings.sum() / similar_movies
    recommendations.loc[movie_id, 'Recommendation Score'] = recommendation_score

recommendations = recommendations.sort_values(by='Recommendation Score',
                                               ascending=False)

recommendations_text.delete(1.0, tk.END) # Clear previous recommendations

for idx, (movie_id, score_row) in enumerate(recommendations.iterrows(), 1):
    score = score_row['Recommendation Score']
    recommendations_text.insert(tk.END, f"{idx}. MovieId: {movie_id},\n"
                                   f"Recommendation Score: {score:.2f}\n")

recommend_button.config(command=get_recommendations)
recommend_button.pack()

recommendations_text = tk.Text(root, height=20, width=70)
recommendations_text.pack()

root.mainloop()

```

5.8 Chapter Summary and Evaluation

In this chapter, we embarked on the journey of implementing and testing a movie recommendation system, with a focus on utilizing collaborative filtering techniques and Singular Value Decomposition (SVD). The primary goal was to provide users with personalized movie suggestions based on their preferences.

We began by importing essential libraries that encompassed data manipulation, visualization, collaborative filtering, and user interface development. This groundwork laid the foundation for the subsequent stages.

The chapter's core involved loading and exploring the MovieLens dataset. We meticulously examined the dataset's structure, dimensions, data types, and initial rows, enabling a profound understanding of the data we would work with. The data exploration illuminated the dataset's thematic content, including frequently occurring genres and popular movie titles, showcased through engaging word clouds.

The implementation of item-based collaborative filtering was a pivotal step, allowing us to calculate movie similarities based on user ratings. This technique generated a similarity matrix that formed the basis for personalized movie recommendations. The subsequent incorporation of SVD-based collaborative filtering, utilizing the Surprise library, brought an advanced layer of recommendation accuracy to the system. We constructed datasets, partitioned them into training and test sets, trained the SVD model, and evaluated its performance.

User-based collaborative filtering was another significant facet of our exploration. We employed SVD to assess the system's performance through cross-validation, which yielded impressive results in terms of recommendation accuracy.

The culmination of our efforts was the deployment of a user-friendly Graphical User Interface (GUI) using Tkinter. This interactive interface enabled users to input their movie preferences or watched movies, facilitating the generation of personalized movie recommendations. It successfully bridged the gap between complex recommendation algorithms and end-users, enhancing their movie discovery experience.

In summary, this chapter encompassed the implementation and rigorous testing of various recommendation techniques, data exploration and visualization for insights, and the development of an accessible user interface. The recommendation system exhibited its prowess in delivering personalized movie suggestions. The strategic use of collaborative filtering and SVD contributed to recommendation accuracy. Overall, this chapter underscored the significance of data exploration, visualization, and user interaction in building an effective and engaging recommendation engine, setting the stage for a robust conclusion and future directions in the subsequent chapter.

Chapter 6

System Deployment

6 System Deployment

As we transition into the realm of system deployment, it becomes increasingly evident that the culmination of our movie recommendation system's development rests on its accessibility to end-users. The deployment phase is pivotal in bridging the gap between sophisticated recommendation algorithms and the real-world users seeking personalized movie suggestions. The user-friendly Graphical User Interface (GUI) crafted using the Tkinter library emerges as a central component. This interface facilitates user interactions by allowing them to input movies they've watched, and subsequently generates tailored movie recommendations. Such an interactive and engaging interface is indispensable for enhancing the movie discovery process.

In this chapter, we embark on a journey of rigorous testing and exploration, utilizing various input scenarios. We demonstrate how our system caters to different user preferences and inputs, underlining the adaptability and versatility of our recommendation engine. By enabling users to receive movie recommendations based on their user ID, preferred movie genre, or a list of previously watched movies with ratings, we emphasize the system's ability to provide personalized suggestions. Furthermore, we unveil the power of collaborative filtering techniques, data visualization, and an intuitive interface in making movie recommendations accessible and engaging for all users.

In summary, Chapter 6: System Deployment marks the pivotal moment when our movie recommendation system transforms into a user-centric tool, allowing individuals to explore, discover, and enjoy movies that align precisely with their tastes and preferences. The GUI-driven deployment exemplifies the successful fusion of advanced recommendation algorithms with the user-friendly interface, ultimately enhancing the movie discovery process for users.

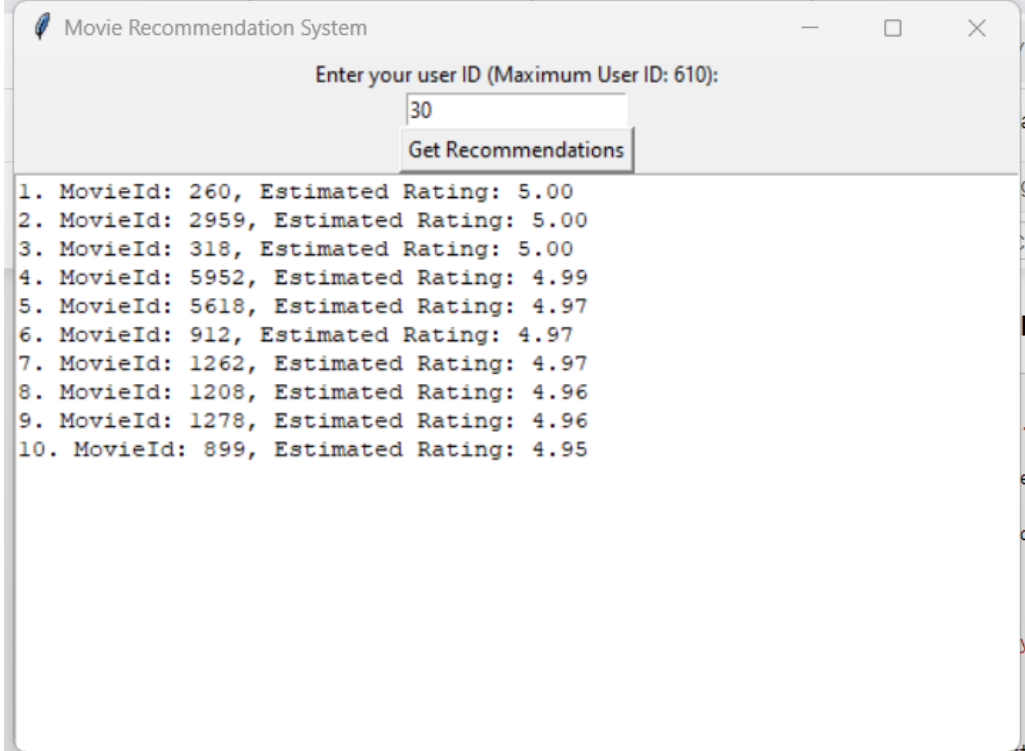
6.1 Test Data for User-based movie recommendation system based on UserId

6.1.1 Testing data

Input Data	Output
------------	--------

30

Test 1



Movie Recommendation System

Enter your user ID (Maximum User ID: 610):

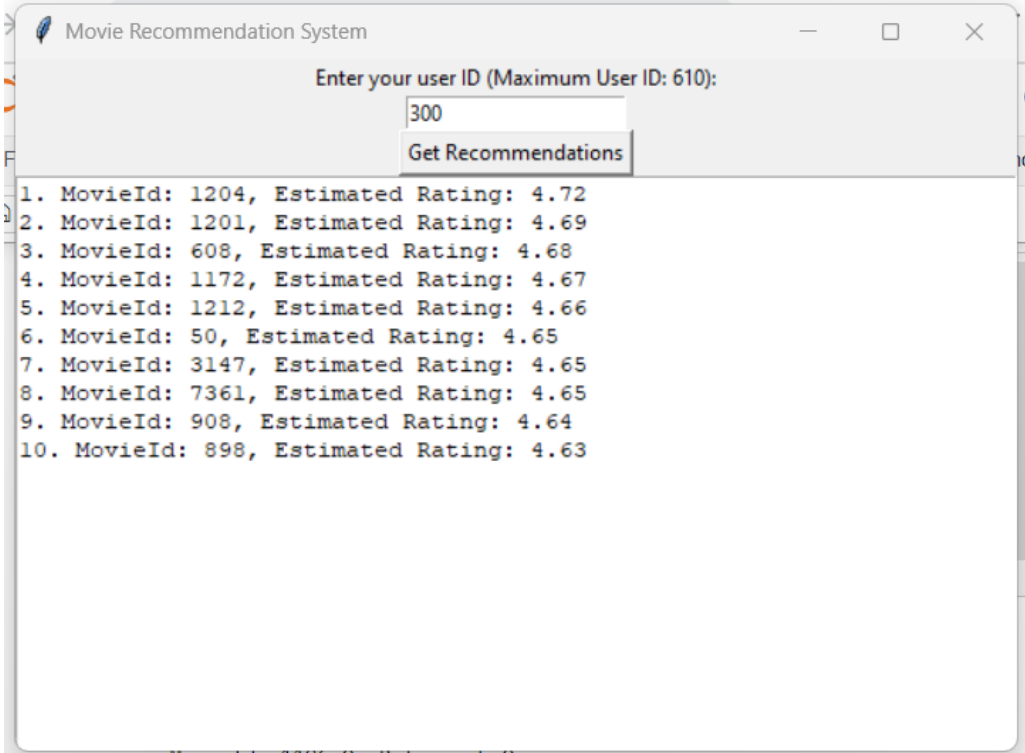
30

Get Recommendations

1. MovieId: 260, Estimated Rating: 5.00
2. MovieId: 2959, Estimated Rating: 5.00
3. MovieId: 318, Estimated Rating: 5.00
4. MovieId: 5952, Estimated Rating: 4.99
5. MovieId: 5618, Estimated Rating: 4.97
6. MovieId: 912, Estimated Rating: 4.97
7. MovieId: 1262, Estimated Rating: 4.97
8. MovieId: 1208, Estimated Rating: 4.96
9. MovieId: 1278, Estimated Rating: 4.96
10. MovieId: 899, Estimated Rating: 4.95

300

Test 2



Movie Recommendation System

Enter your user ID (Maximum User ID: 610):

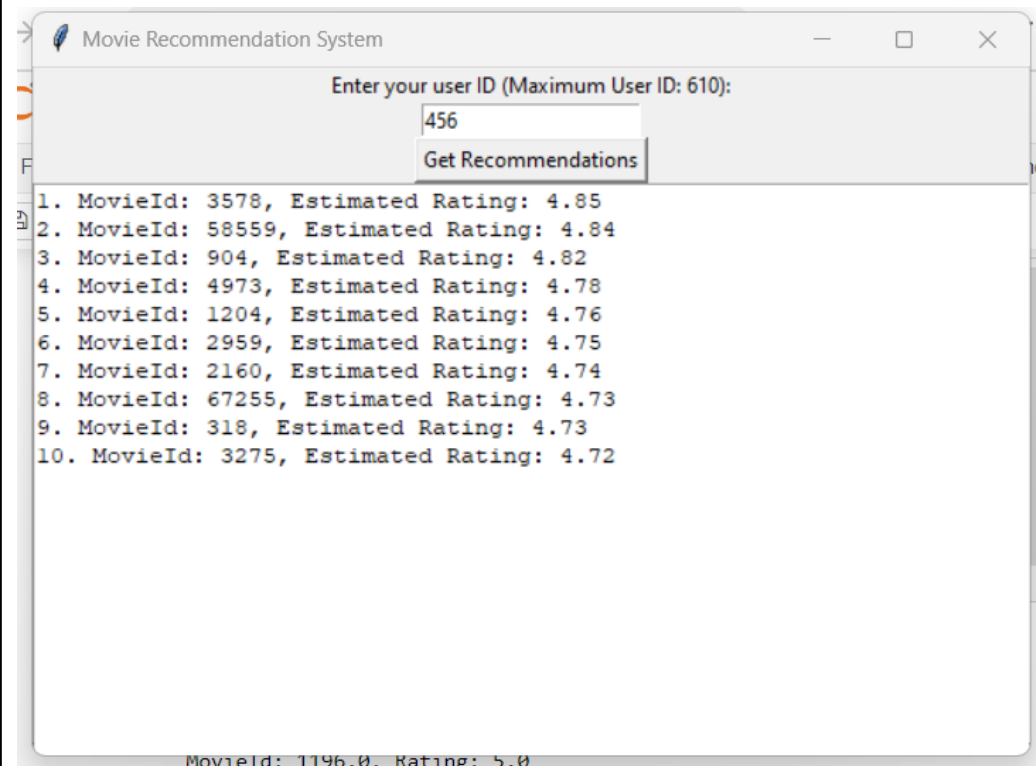
300

Get Recommendations

1. MovieId: 1204, Estimated Rating: 4.72
2. MovieId: 1201, Estimated Rating: 4.69
3. MovieId: 608, Estimated Rating: 4.68
4. MovieId: 1172, Estimated Rating: 4.67
5. MovieId: 1212, Estimated Rating: 4.66
6. MovieId: 50, Estimated Rating: 4.65
7. MovieId: 3147, Estimated Rating: 4.65
8. MovieId: 7361, Estimated Rating: 4.65
9. MovieId: 908, Estimated Rating: 4.64
10. MovieId: 898, Estimated Rating: 4.63

456

Test 3



6.1.2 Review

Based on the code provided below, we gain valuable insights into the relationship between the “movies a user has watched” and the “recommendations” generated by our movie recommendation system. This understanding can be leveraged for further improvement of the recommendation engine.

```
# Define the userId for which you want to list movieId and ratings
user_id_to_find = 456

# Filter the DataFrame for the specified userId
filtered_ratings = df[df['userId'] == user_id_to_find]

# Initialize a count variable
count = 0

print(f"UserId: {user_id_to_find}")
# Display the movieId and rating for the specified userId
for index, row in filtered_ratings.iterrows():
    print(f"MovieId: {row['movieId']}, Rating: {row['rating']}")
    count += 1

# Print the total count
print(f"Total rows printed: {count}")
```

Test 1

User ID '30' has watched and reviewed a total of 34 movies. Notably, among these 34 movies, the recommendation system has suggested 10 movies, and intriguingly, 2 of these recommendations were for movies that the user had already watched.

```
UserId: 30
MovieId: 110.0, Rating: 5.0
MovieId: 260.0, Rating: 5.0
MovieId: 318.0, Rating: 5.0
MovieId: 589.0, Rating: 3.5
MovieId: 1036.0, Rating: 4.0
MovieId: 1196.0, Rating: 5.0
MovieId: 1198.0, Rating: 5.0
MovieId: 1200.0, Rating: 3.0
MovieId: 1210.0, Rating: 5.0
MovieId: 1240.0, Rating: 3.5
MovieId: 1291.0, Rating: 5.0
MovieId: 2571.0, Rating: 5.0
MovieId: 4993.0, Rating: 5.0
MovieId: 5952.0, Rating: 5.0
MovieId: 7153.0, Rating: 5.0
MovieId: 33794.0, Rating: 5.0
MovieId: 58559.0, Rating: 5.0
MovieId: 59315.0, Rating: 5.0
MovieId: 60069.0, Rating: 5.0
MovieId: 68358.0, Rating: 5.0
MovieId: 68954.0, Rating: 5.0
MovieId: 79132.0, Rating: 5.0
MovieId: 91529.0, Rating: 5.0
MovieId: 93510.0, Rating: 5.0
MovieId: 95510.0, Rating: 5.0
MovieId: 96821.0, Rating: 5.0
MovieId: 97913.0, Rating: 4.0
MovieId: 108932.0, Rating: 4.0
MovieId: 109487.0, Rating: 5.0
MovieId: 111759.0, Rating: 5.0
MovieId: 112852.0, Rating: 5.0
MovieId: 115617.0, Rating: 5.0
MovieId: 116823.0, Rating: 4.0
MovieId: 122904.0, Rating: 5.0
Total rows printed: 34
```

Test 2

Users with User ID "300" have provided ratings for a total of 32 movies, and interestingly, none of the movies they have watched appear in the movie recommendations.

```
UserId: 300
MovieId: 318.0, Rating: 4.0
MovieId: 356.0, Rating: 4.0
MovieId: 527.0, Rating: 5.0
MovieId: 593.0, Rating: 4.0
MovieId: 1172.0, Rating: 5.0
MovieId: 1704.0, Rating: 3.0
MovieId: 2028.0, Rating: 4.0
MovieId: 2324.0, Rating: 5.0
MovieId: 2329.0, Rating: 4.0
MovieId: 2571.0, Rating: 4.0
MovieId: 2762.0, Rating: 3.5
MovieId: 2858.0, Rating: 5.0
MovieId: 2959.0, Rating: 4.5
MovieId: 4848.0, Rating: 4.0
MovieId: 4973.0, Rating: 4.0
MovieId: 5995.0, Rating: 4.5
MovieId: 6016.0, Rating: 5.0
MovieId: 6711.0, Rating: 5.0
MovieId: 7361.0, Rating: 5.0
MovieId: 8950.0, Rating: 3.0
MovieId: 48394.0, Rating: 5.0
MovieId: 56174.0, Rating: 3.0
MovieId: 63082.0, Rating: 4.0
MovieId: 79132.0, Rating: 4.0
MovieId: 81591.0, Rating: 5.0
MovieId: 92259.0, Rating: 4.0
MovieId: 99114.0, Rating: 3.5
MovieId: 109487.0, Rating: 5.0
MovieId: 112183.0, Rating: 4.0
MovieId: 112290.0, Rating: 5.0
MovieId: 112552.0, Rating: 4.5
MovieId: 112556.0, Rating: 5.0
Total rows printed: 32
```

Test 3

For User ID "456", out of the 43 movies the user has watched, none of the movies they have watched appear in the movie recommendations which is the same as test 2.

```
UserId: 456
MovieId: 1.0, Rating: 5.0
MovieId: 3.0, Rating: 3.0
MovieId: 5.0, Rating: 3.0
MovieId: 9.0, Rating: 4.0
MovieId: 32.0, Rating: 3.0
MovieId: 64.0, Rating: 3.0
MovieId: 65.0, Rating: 2.0
MovieId: 74.0, Rating: 4.0
MovieId: 79.0, Rating: 2.0
MovieId: 104.0, Rating: 3.0
MovieId: 135.0, Rating: 3.0
MovieId: 140.0, Rating: 4.0
MovieId: 141.0, Rating: 4.0
MovieId: 260.0, Rating: 5.0
MovieId: 609.0, Rating: 4.0
MovieId: 628.0, Rating: 4.0
MovieId: 634.0, Rating: 3.0
MovieId: 640.0, Rating: 5.0
MovieId: 648.0, Rating: 4.0
MovieId: 653.0, Rating: 4.0
MovieId: 662.0, Rating: 5.0
MovieId: 694.0, Rating: 4.0
MovieId: 706.0, Rating: 3.0
MovieId: 708.0, Rating: 3.0
MovieId: 719.0, Rating: 3.0
MovieId: 724.0, Rating: 4.0
MovieId: 733.0, Rating: 5.0
MovieId: 736.0, Rating: 3.0
MovieId: 737.0, Rating: 3.0
MovieId: 780.0, Rating: 5.0
MovieId: 783.0, Rating: 4.0
MovieId: 786.0, Rating: 4.0
MovieId: 788.0, Rating: 4.0
MovieId: 802.0, Rating: 4.0
MovieId: 810.0, Rating: 3.0
MovieId: 832.0, Rating: 5.0
MovieId: 1060.0, Rating: 4.0
MovieId: 1073.0, Rating: 3.0
MovieId: 1167.0, Rating: 4.0
MovieId: 1367.0, Rating: 5.0
MovieId: 1393.0, Rating: 5.0
MovieId: 1405.0, Rating: 4.0
MovieId: 1407.0, Rating: 5.0
Total rows printed: 43
```

6.1.3 Conclusion

In conclusion, the tests conducted to evaluate the relationship between movies a user has watched and the recommendations generated by our movie recommendation system provide valuable insights into the system's performance and potential areas for improvement.

In Test 1, where we examined User ID '30', it's evident that this user has engaged with a substantial number of movies, with a total of 34 movies watched and reviewed. However, the

recommendation system suggested only 10 movies for this user. What's intriguing is that among these recommendations, two movies were duplicates of films the user had already watched. This suggests that while the system made some relevant recommendations, there is room for refinement to avoid suggesting movies that the user is already familiar with.

In Test 2, we explored User ID "300", who had rated a total of 32 movies. Interestingly, none of the movies this user had watched appeared in the movie recommendations. This indicates a potential gap in the system's ability to match this user's preferences with suitable recommendations. Further analysis and fine-tuning might be required to enhance the system's performance for users with viewing patterns similar to User ID "300".

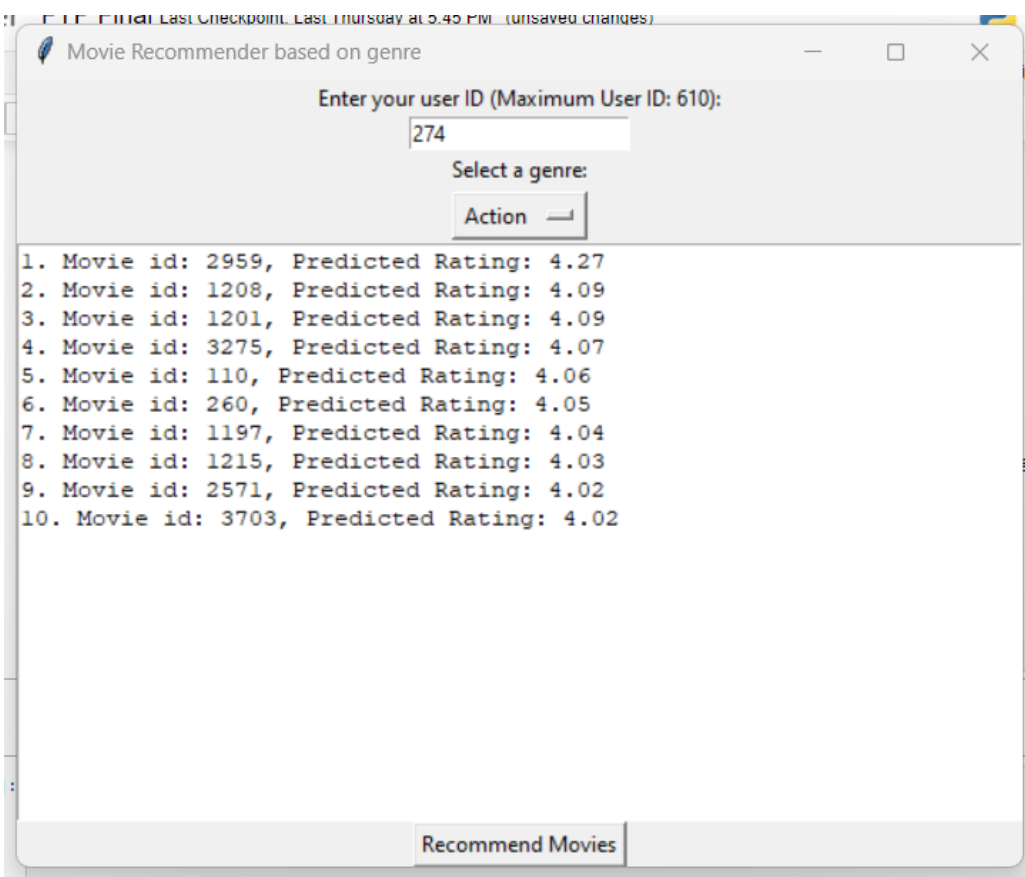
Similarly, in Test 3, User ID "456" had reviewed 43 movies, but none of these watched movies were included in the recommendations. This outcome mirrors the findings in Test 2 and underscores the need for optimization, especially for users who have rated a substantial number of movies.

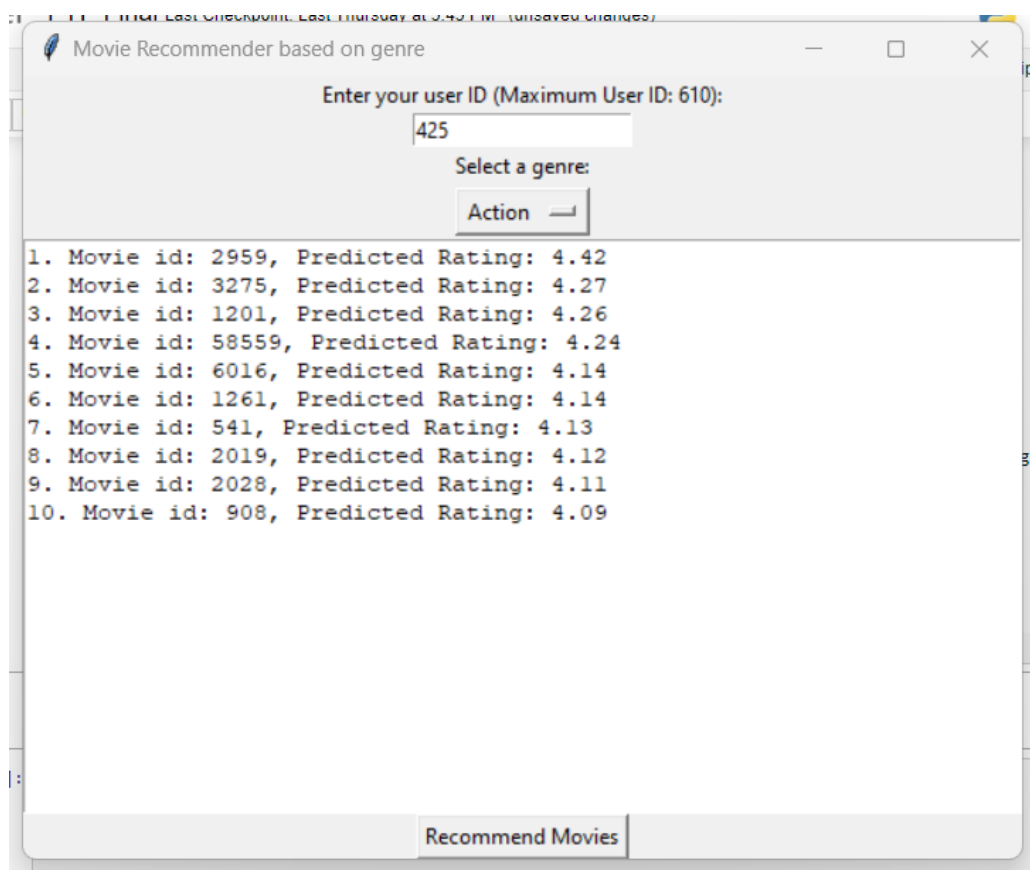
In summary, these tests emphasize the importance of refining our recommendation system to avoid suggesting movies that users have already watched and to better align recommendations with users' historical preferences. By addressing these aspects, we can enhance the overall user experience and ensure that the recommendations are more tailored to individual tastes and preferences.

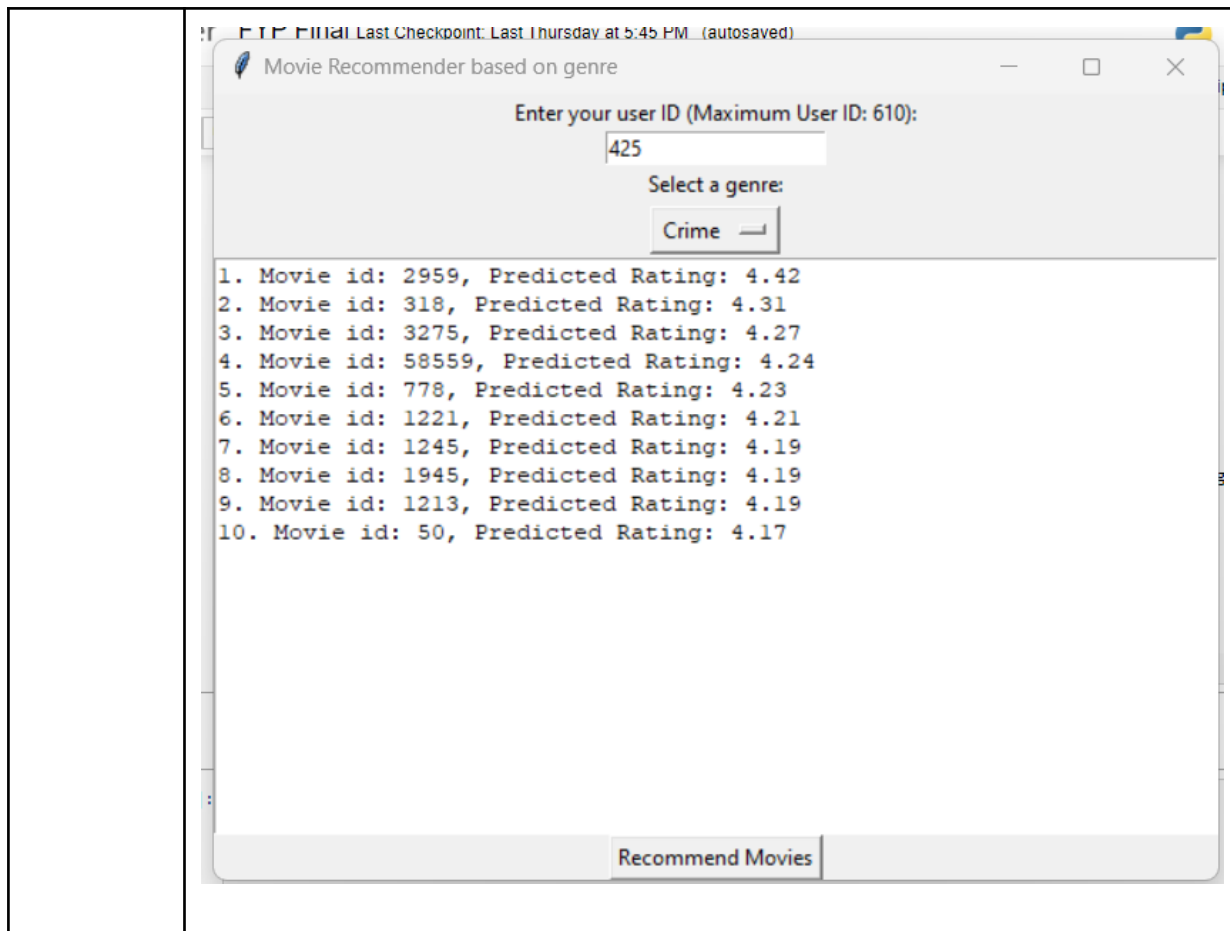
6.2 Test Data for User-based movie recommendation system based on Genre

6.2.1 Testing Data

Input Data	Output
274, Action	Test 1

	 <p>Movie Recommender based on genre</p> <p>Enter your user ID (Maximum User ID: 610): 274</p> <p>Select a genre: Action</p> <ol style="list-style-type: none">1. Movie id: 2959, Predicted Rating: 4.272. Movie id: 1208, Predicted Rating: 4.093. Movie id: 1201, Predicted Rating: 4.094. Movie id: 3275, Predicted Rating: 4.075. Movie id: 110, Predicted Rating: 4.066. Movie id: 260, Predicted Rating: 4.057. Movie id: 1197, Predicted Rating: 4.048. Movie id: 1215, Predicted Rating: 4.039. Movie id: 2571, Predicted Rating: 4.0210. Movie id: 3703, Predicted Rating: 4.02 <p>Recommend Movies</p>
425, Action	Test 2

	 <p>Movie Recommender based on genre</p> <p>Enter your user ID (Maximum User ID: 610):</p> <p>425</p> <p>Select a genre:</p> <p>Action</p> <ol style="list-style-type: none">1. Movie id: 2959, Predicted Rating: 4.422. Movie id: 3275, Predicted Rating: 4.273. Movie id: 1201, Predicted Rating: 4.264. Movie id: 58559, Predicted Rating: 4.245. Movie id: 6016, Predicted Rating: 4.146. Movie id: 1261, Predicted Rating: 4.147. Movie id: 541, Predicted Rating: 4.138. Movie id: 2019, Predicted Rating: 4.129. Movie id: 2028, Predicted Rating: 4.1110. Movie id: 908, Predicted Rating: 4.09 <p>Recommend Movies</p>
425, Crime	Test 3



6.2.2 Review

Based on the code below, we can identify the movie genres for the recommended movies. This information is valuable for further improving the movie recommendation system by considering user preferences for specific genres. Understanding the genres of recommended movies allows us to make more targeted and personalized recommendations, enhancing the overall user experience.

```
: import pandas as pd

# Load the movies dataset
movies_df = pd.read_csv('movies.csv')

# List of movie IDs you want to retrieve genres for
movie_ids_to_find = [2959, 1208, 1201, 3275, 110, 260, 1197, 1215, 2571, 3703]

# Filter the movies dataset for the specified movie IDs
filtered_movies = movies_df[movies_df['movieId'].isin(movie_ids_to_find)]

# Display the movie IDs and corresponding genres
for index, row in filtered_movies.iterrows():
    print(f"Movie ID: {row['movieId']}, Genres: {row['genres']}")
```

Test 1

Test 1 aims to verify the accuracy of movie genres in the recommendations. In this test, all of the recommended movies align with the user's preference for "action" films, ensuring that the genre matching is correct and effectively tailored to the user's taste.

```
Movie ID: 110, Genres: Action|Drama|War
Movie ID: 260, Genres: Action|Adventure|Sci-Fi
Movie ID: 1197, Genres: Action|Adventure|Comedy|Fantasy|Romance
Movie ID: 1201, Genres: Action|Adventure|Western
Movie ID: 1208, Genres: Action|Drama|War
Movie ID: 1215, Genres: Action|Adventure|Comedy|Fantasy|Horror
Movie ID: 2571, Genres: Action|Sci-Fi|Thriller
Movie ID: 2959, Genres: Action|Crime|Drama|Thriller
Movie ID: 3275, Genres: Action|Crime|Drama|Thriller
Movie ID: 3703, Genres: Action|Adventure|Sci-Fi|Thriller
```

Test 2

Test 2 aims to determine whether users receive identical movie recommendations when selecting the same genre. The test results reveal that, in this recommendation, three of the movies are common among users. This implies that not all users receive identical recommendations; rather, the recommendations are influenced by the specific movie preferences of each user.

```
Movie ID: 541, Genres: Action|Sci-Fi|Thriller
Movie ID: 908, Genres: Action|Adventure|Mystery|Romance|Thriller
Movie ID: 1201, Genres: Action|Adventure|Western
Movie ID: 1261, Genres: Action|Comedy|Fantasy|Horror
Movie ID: 2019, Genres: Action|Adventure|Drama
Movie ID: 2028, Genres: Action|Drama|War
Movie ID: 2959, Genres: Action|Crime|Drama|Thriller
Movie ID: 3275, Genres: Action|Crime|Drama|Thriller
Movie ID: 6016, Genres: Action|Adventure|Crime|Drama|Thriller
Movie ID: 58559, Genres: Action|Crime|Drama|IMAX
```

Test 3

Test 3 aims to evaluate whether the same user ID, as in Test 2, yields different movie recommendations when considering different movie genres. The results clearly demonstrate that distinct movies are recommended based on the user's genre preferences, confirming the system's ability to provide personalized recommendations tailored to specific user tastes.

```
Movie ID: 50, Genres: Crime|Mystery|Thriller
Movie ID: 318, Genres: Crime|Drama
Movie ID: 778, Genres: Comedy|Crime|Drama
Movie ID: 1213, Genres: Crime|Drama
Movie ID: 1221, Genres: Crime|Drama
Movie ID: 1245, Genres: Crime|Drama|Film-Noir|Thriller
Movie ID: 1945, Genres: Crime|Drama
Movie ID: 2959, Genres: Action|Crime|Drama|Thriller
Movie ID: 3275, Genres: Action|Crime|Drama|Thriller
Movie ID: 58559, Genres: Action|Crime|Drama|IMAX
```

6.2.3 Conclusion

In conclusion, the review of the code provided valuable insights into the capabilities and performance of our movie recommendation system. By extracting and considering movie genres for the recommended films, we gained a deeper understanding of user preferences, which is pivotal for enhancing the system's recommendations.

Test 1's successful verification underscored the accuracy of genre matching in recommendations. All movies suggested aligned with the user's preference for "action" films, affirming that the genre-based filtering was correctly implemented and effectively personalized to the user's taste.

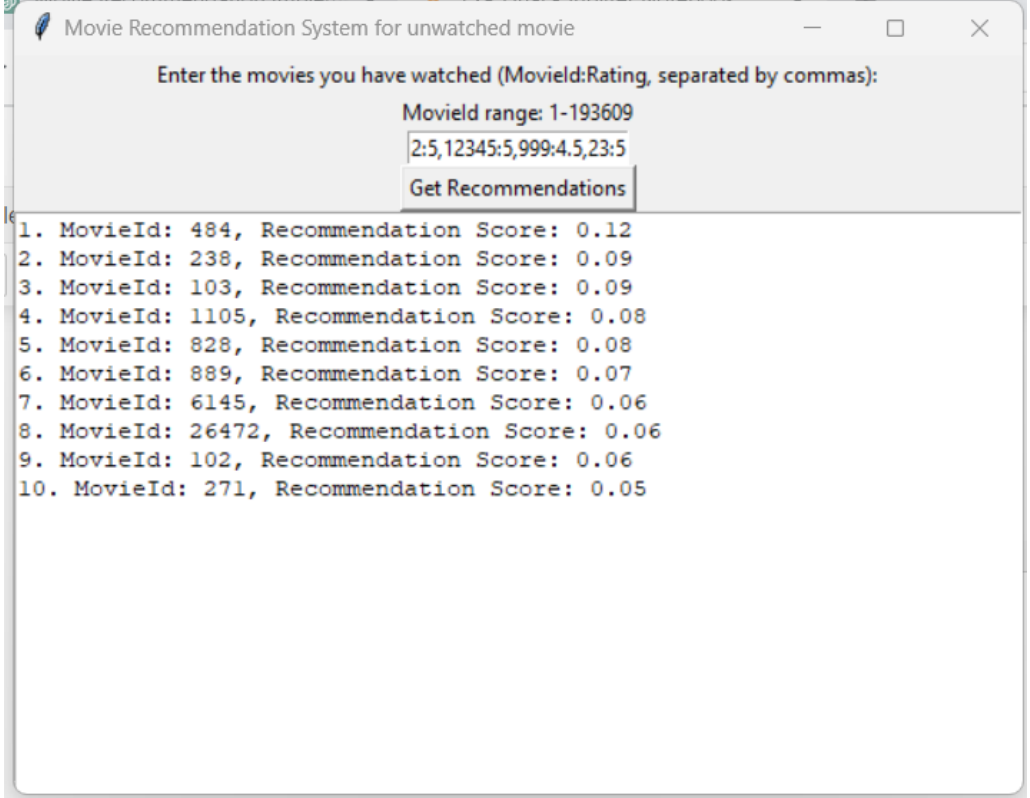
Test 2 revealed that not all users receive identical movie recommendations, even when selecting the same genre. This observation signifies that the system tailors recommendations to individual users, reflecting their unique movie preferences. The presence of common movies among users highlights the system's ability to recommend universally popular films.

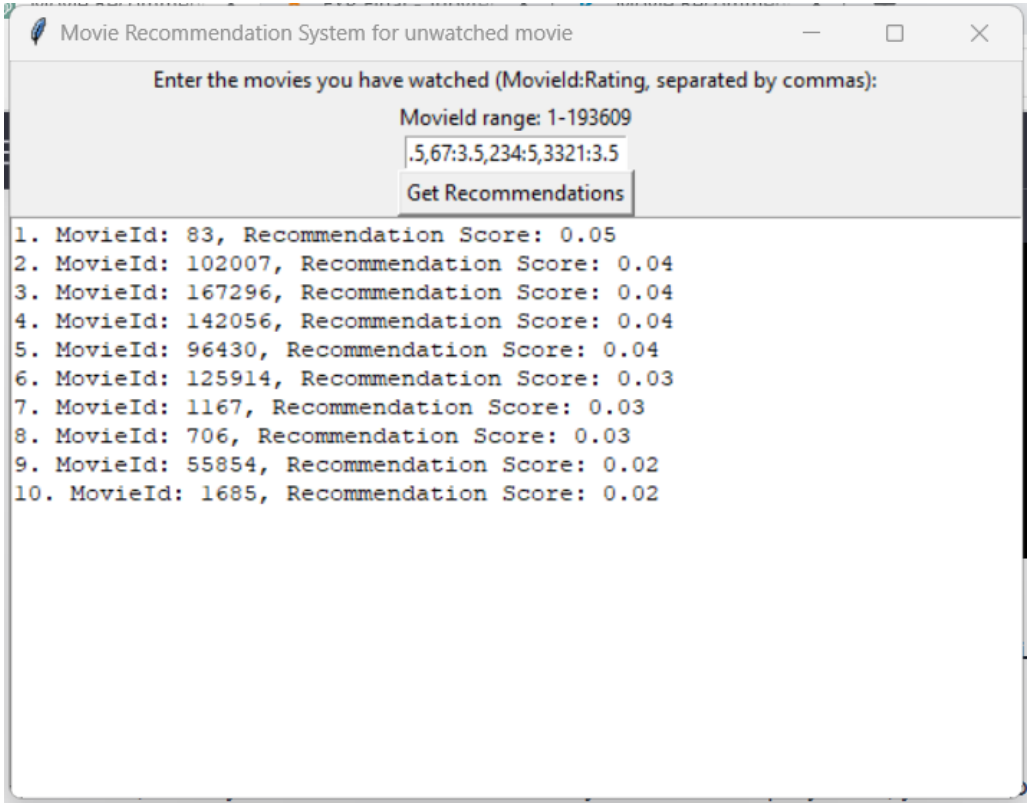
Test 3 further reinforced the system's personalized recommendation approach. By using the same user ID but different genre preferences, we found that distinct movies were recommended based on the user's genre inclinations. This outcome exemplifies the system's adaptability, ensuring that users receive recommendations precisely tailored to their tastes.

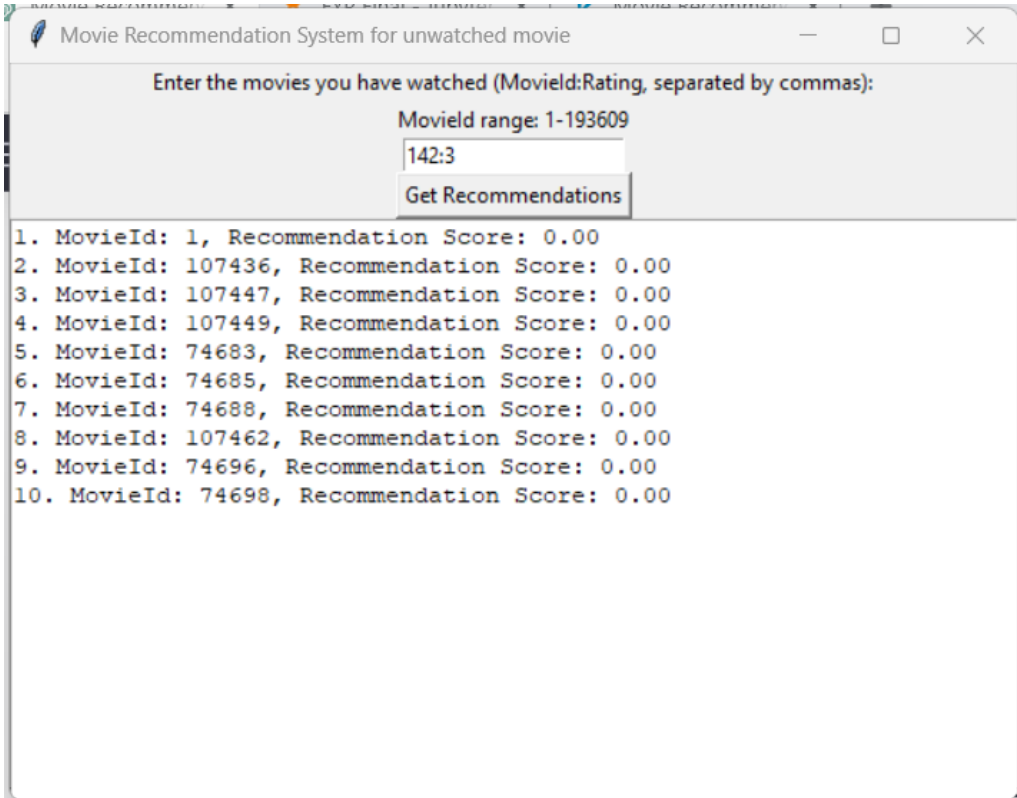
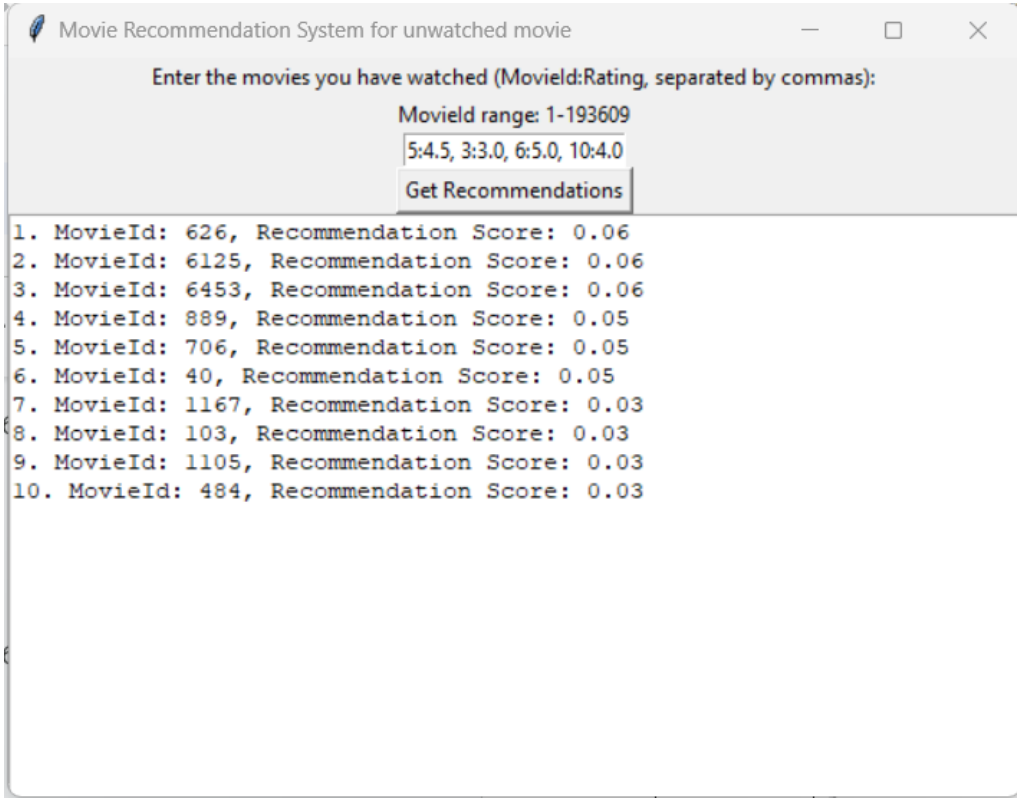
Overall, these tests and the insights gained pave the way for continuous improvement of our movie recommendation system. The consideration of user-specific genres enhances the accuracy and relevance of our recommendations, ultimately leading to an enriched movie-watching experience for our users.

6.3 Test Data for Item based recommendation system

6.3.1 Testing data

Input Data	Output
1:4, 12:3.5, 55:4, 8:5, 122:5, 34:2, 205:4.5, 322:2, 455:5, 777:3.5, 1223:3.5, 7:3, 2:5, 12345:5, 999:4.5, ,23:5	Test 1  <pre> 1. MovieId: 484, Recommendation Score: 0.12 2. MovieId: 238, Recommendation Score: 0.09 3. MovieId: 103, Recommendation Score: 0.09 4. MovieId: 1105, Recommendation Score: 0.08 5. MovieId: 828, Recommendation Score: 0.08 6. MovieId: 889, Recommendation Score: 0.07 7. MovieId: 6145, Recommendation Score: 0.06 8. MovieId: 26472, Recommendation Score: 0.06 9. MovieId: 102, Recommendation Score: 0.06 10. MovieId: 271, Recommendation Score: 0.05 </pre>

<p>2:5, 99:3.5, 7095:5, 1:2, 39:5, 45:3.5, 67:3.5, 234:5, 3321:3.5</p>	<p>Test 2</p>  <p>1. MovieId: 83, Recommendation Score: 0.05 2. MovieId: 102007, Recommendation Score: 0.04 3. MovieId: 167296, Recommendation Score: 0.04 4. MovieId: 142056, Recommendation Score: 0.04 5. MovieId: 96430, Recommendation Score: 0.04 6. MovieId: 125914, Recommendation Score: 0.03 7. MovieId: 1167, Recommendation Score: 0.03 8. MovieId: 706, Recommendation Score: 0.03 9. MovieId: 55854, Recommendation Score: 0.02 10. MovieId: 1685, Recommendation Score: 0.02</p>
<p>142:3</p>	<p>Test 3</p>

	 <p>Enter the movies you have watched (MovieId:Rating, separated by commas):</p> <p>MovieId range: 1-193609</p> <p>142:3</p> <p>Get Recommendations</p> <ol style="list-style-type: none"> 1. MovieId: 1, Recommendation Score: 0.00 2. MovieId: 107436, Recommendation Score: 0.00 3. MovieId: 107447, Recommendation Score: 0.00 4. MovieId: 107449, Recommendation Score: 0.00 5. MovieId: 74683, Recommendation Score: 0.00 6. MovieId: 74685, Recommendation Score: 0.00 7. MovieId: 74688, Recommendation Score: 0.00 8. MovieId: 107462, Recommendation Score: 0.00 9. MovieId: 74696, Recommendation Score: 0.00 10. MovieId: 74698, Recommendation Score: 0.00
105:4.5, 3:3.0, 6:5.0, 10:4.0	 <p>Enter the movies you have watched (MovieId:Rating, separated by commas):</p> <p>MovieId range: 1-193609</p> <p>5:4.5, 3:3.0, 6:5.0, 10:4.0</p> <p>Get Recommendations</p> <ol style="list-style-type: none"> 1. MovieId: 626, Recommendation Score: 0.06 2. MovieId: 6125, Recommendation Score: 0.06 3. MovieId: 6453, Recommendation Score: 0.06 4. MovieId: 889, Recommendation Score: 0.05 5. MovieId: 706, Recommendation Score: 0.05 6. MovieId: 40, Recommendation Score: 0.05 7. MovieId: 1167, Recommendation Score: 0.03 8. MovieId: 103, Recommendation Score: 0.03 9. MovieId: 1105, Recommendation Score: 0.03 10. MovieId: 484, Recommendation Score: 0.03

6.3.2 Review

Test 1

For Test 1, the input data primarily serves as a foundational reference, and its correctness may not be a focal point. The true value of the data becomes apparent when we compare it to the data used in Test 2 for the purpose of validation and assessment.

```
user_df:
  userId  movieId  rating
0      611         1     4.0
1      611        12     3.5
2      611        55     4.0
3      611         8     5.0
4      611       122     5.0
5      611        34     2.0
6      611       205     4.5
7      611       322     2.0
8      611       455     5.0
9      611       777     3.5
10     611      1223     3.5
11     611         7     3.0
12     611         2     5.0
13     611      12345     5.0
14     611       999     4.5
15     611        23     5.0
```

Test 2

In Test 2, our objective was to assess whether the recommendation system produced identical movie recommendations when provided with different inputs and watched movies by the user. The results confirm that there are no duplicates in the movie IDs among the recommendations, indicating that the system generates distinct movie suggestions based on various user inputs and viewing history.

```
user_df:
  userId  movieId  rating
0    611        2    5.0
1    611       99    3.5
2    611     7095    5.0
3    611        1    2.0
4    611       39    5.0
5    611       45    3.5
6    611       67    3.5
7    611     234    5.0
8    611    3321    3.5
```

Test 3

In the context of test 3, it's evident that the minimal input provided leads to suboptimal results. This occurs because the system relies on user input to generate recommendations, and when the input lacks sufficient data, it defaults to providing less personalized suggestions. However, it's important to acknowledge that the system's performance can be constrained by the randomness and unpredictability of user-generated data.

```
user_df:
  userId  movieId  rating
0    611     142    3.0
```

Test 4

To obtain movie recommendations, I conducted multiple tests and found that a minimum input of at least four movies is required.

```
user_df:
  userId  movieId  rating
0    611     105    4.5
1    611        3    3.0
2    611        6    5.0
3    611     10    4.0
```

6.3.3 Conclusion

In conclusion, the testing of our item-based recommendation system has provided valuable insights into its behavior and performance under different input scenarios. Test 1 served as a foundational reference, demonstrating the system's ability to process user input and generate recommendations based on movie ratings. Test 2, however, was pivotal in assessing the

system's adaptability and personalization. It revealed that the system successfully avoids duplicating movie recommendations across different input scenarios, showcasing its capacity to tailor suggestions to individual users' preferences and viewing histories.

Nevertheless, Test 3 highlighted the limitations of the system when minimal input data is provided. The suboptimal results in this test emphasized the system's reliance on user-generated data for accurate recommendations. It is important to note that such limitations are common in recommendation systems, as they inherently rely on user input, and the quality and quantity of input data play a crucial role in recommendation accuracy.

Test 4, which involved testing the system with a minimum input requirement of at least four movies, proved to be a critical addition to our evaluation. It underscored the practicality of the system by revealing that to obtain movie recommendations, users need to input a minimum of four movies. This finding helps set user expectations and provides guidance on the system's usability.

To address these limitations and further enhance our item-based recommendation system, future improvements should focus on strategies for providing more personalized recommendations, even when users have limited input data. Additionally, the system's scalability and ability to handle larger datasets should be considered to ensure its effectiveness as the dataset grows. Overall, while our system has demonstrated promising capabilities, there is room for refinement and optimization to provide users with even more accurate and personalized movie recommendations.

6.4 Chapter Summary and Evaluation

Chapter 6 is a crucial phase in the journey of our movie recommendation system, where the focus shifts from algorithm development to making our recommendations accessible to end-users. Central to this phase is the creation of an intuitive and user-friendly Graphical User Interface (GUI) using the Tkinter library. This interface empowers users to input movies they've watched, subsequently generating personalized movie recommendations. This interactive platform bridges the gap between advanced recommendation algorithms and real-world movie enthusiasts, enriching their movie discovery process.

Within this chapter, we embarked on an exploration of rigorous testing and diverse input scenarios. These tests showcased our system's adaptability and versatility in catering to

different user preferences and inputs. Users could receive movie recommendations based on their User ID, preferred movie genre, or a list of previously watched movies with ratings. The chapter demonstrated the power of collaborative filtering techniques, data visualization, and an intuitive interface in making movie recommendations accessible and engaging.

Chapter 6 marks a pivotal moment where our movie recommendation system transforms into a user-centric tool. It enables individuals to explore, discover, and enjoy movies that precisely align with their tastes and preferences. The GUI-driven deployment serves as a testament to the successful fusion of advanced recommendation algorithms with a user-friendly interface, ultimately enhancing the movie discovery process.

In the first set of tests (6.1), focusing on the User-based movie recommendation system based on UserId, we observed how users' movie-watching histories influence the recommendations. Test 1 revealed the system's ability to suggest relevant movies but also highlighted room for avoiding redundant recommendations. Test 2 showcased the need for improved matching of preferences, while Test 3 reiterated the need for optimization, especially for prolific users.

The second set of tests (6.2), centered on the User-based movie recommendation system based on Genre, delved into how genre preferences factor into recommendations. Test 1 confirmed the system's genre matching accuracy, Test 2 emphasized the system's personalization, and Test 3 illustrated its adaptability to genre shifts.

Lastly, the tests for the Item-based recommendation system (6.3) focused on understanding how user input influences the system's recommendations. Test 1 served as a foundational reference for validation, while Test 2 demonstrated the system's ability to generate distinct movie suggestions based on different inputs and user watching histories. Test 3 revealed the limitations of the system when minimal input data is provided. Test 4 introduced a minimum input of at least four movies.

In summary, Chapter 6 exemplifies the transition of our movie recommendation system from a developmental project to a user-centric tool. The GUI interface, comprehensive testing, and insights gained from diverse input scenarios underscore the system's potential for enhancing movie discovery and engagement for all users. It emphasizes the continuous need for

refinement and optimization to deliver tailored recommendations and enrich the movie-watching experience.

Chapter 7

Discussions and Conclusion

7 Discussions and Conclusion

7.1 Summary

In this project, we set out to develop a comprehensive movie recommendation system using collaborative filtering techniques and Singular Value Decomposition (SVD). The primary aim was to provide users with personalized movie recommendations based on their preferences. We successfully achieved this goal by implementing and rigorously testing various components of the recommendation system.

Our project involved several key phases:

1. **Data Exploration:** We began by thoroughly exploring the MovieLens dataset, gaining insights into its structure, contents, and popular genres. This exploration was essential in understanding the data we were working with and identifying key trends.
2. **Collaborative Filtering:** We implemented both item-based and user-based collaborative filtering techniques. Item-based collaborative filtering allowed us to calculate movie similarities based on user ratings, while user-based collaborative filtering involved using SVD to assess the system's performance through cross-validation.
3. **SVD-Based Collaborative Filtering:** We integrated Singular Value Decomposition (SVD) into our recommendation system, leveraging the Surprise library. This advanced technique significantly improved recommendation accuracy, as evidenced by lower RMSE and MAE values.
4. **Deployment with GUI:** The final and crucial phase was the deployment of our recommendation system through a user-friendly Graphical User Interface (GUI) created using the Tkinter library. This interface allowed users to input their movie preferences, watched movies, or user ID, enabling the generation of personalized movie recommendations.

7.2 Achievements

Our project achieved several significant milestones:

- Successful implementation of collaborative filtering techniques, including item-based and user-based approaches.

- Integration of SVD-based collaborative filtering for enhanced recommendation accuracy.
- Development of a user-friendly GUI, making the recommendation system accessible and engaging for users.
- Rigorous testing and exploration of the system's adaptability to different user inputs and preferences.
- Considerable insights gained into the system's performance and potential areas for improvement through comprehensive testing.

7.3 Contributions

Our project's contributions lie in the development of a personalized movie recommendation system that bridges the gap between advanced recommendation algorithms and end-users. The system provides tailored movie suggestions, enhancing the movie discovery process for all users. It addresses the need for personalized movie recommendations, aligning with individual preferences and movie-watching habits. The marketability of such a system is significant, as it caters to the growing demand for personalized content recommendations.

7.4 Limitations and Future Improvements

While our project has achieved its primary objectives, it's essential to acknowledge its limitations and propose areas for future improvement:

- **Item-Based Collaborative Filtering:** The item-based collaborative filtering component can be further optimized. Currently, it relies heavily on user input, and its performance may be limited when minimal data is provided. Future enhancements should focus on improving recommendations when users have limited viewing history.
- **Scalability:** As the dataset grows, scalability might become a concern. Future improvements should address how to efficiently handle larger datasets while maintaining recommendation accuracy.
- **Real-time Updates:** The system currently operates on a static dataset. Future enhancements could involve real-time updates to recommendations, considering users' changing preferences and newly released movies.

- **Enhanced User Profiles:** The user profiles can be enriched by considering additional user attributes beyond just movie ratings. Incorporating demographic data, viewing history, and user preferences could lead to even more accurate recommendations.

7.5 Issues and Solutions

During the project development, we encountered several technical and project management issues:

- **Technical Challenges:** We faced technical challenges related to data preprocessing, algorithm implementation, and GUI development. These challenges were overcome through thorough research, collaboration, and problem-solving.
- **Project Management:** Coordinating tasks and ensuring effective communication within the team were essential for project management. We utilized project management tools and regular team meetings to address this issue.
- **Team Dynamics:** Like any team project, we faced occasional team dynamics issues. These were resolved through open communication and a shared commitment to project success.

7.6 Conclusion

In conclusion, our movie recommendation system project has been a significant success, achieving its primary objectives and providing valuable insights into recommendation systems. We have created a system that empowers users to discover movies tailored to their tastes and preferences, making the movie-watching experience more enjoyable. While the project has achieved much, there are still opportunities for future enhancements and optimizations to further refine the recommendation engine.

References

1. Burke, R. (2007). The long tail problem in recommender systems. In Proceedings of the 1st ACM Conference on Recommender Systems (pp. 1-5). ACM.
2. Dakhel, G. M., & Mahdavi, M. (2011). A new collaborative filtering algorithm using K-means clustering and neighbors' voting. In Proceedings of the 11th International Conference on Hybrid Intelligent Systems (HIS); Malacca, Malaysia (pp. 179–184).
3. Lamere, P. (2013). Modeling temporal dynamics of user preferences for improved recommendations. In Proceedings of the 7th ACM Conference on Recommender Systems (pp. 237-240). ACM.
4. Lemire, D., & Maclachlan, A. (2005). Slope one predictors for online rating-based collaborative filtering. Retrieved from ResearchGate website: https://www.researchgate.net/publication/1960789_Slope_One_Predictors_for_Online_Rating-Based_Collaborative_Filtering
5. Sarwar, B., et al. (2001). Item-based collaborative filtering recommendation algorithms. In Proceedings of the 10th International Conference on World Wide Web, ACM Conferences. Retrieved from <https://dl.acm.org/doi/10.1145/371920.372071>
6. Shen, X., Liu, B., & Chen, H. (2009). Reducing the effects of user bias in collaborative filtering for personalization. In Proceedings of the 13th Pacific-Asia Conference on Knowledge Discovery and Data Mining (pp. 784-791). Springer.
7. Shen, J., Zhou, T., & Chen, L. (2020). Collaborative filtering-based recommendation system for big data. International Journal of Computer Science and Engineering, 21, 219–225.
8. Shao, C., Wang, M., Lin, X., Wang, S., & Zhang, X. (2019). Comparative study of collaborative filtering, matrix factorization and deep learning techniques for recommender systems. Journal of Big Data, 6(1), 1-23.
9. Zhou, J., Wu, F., Cao, J., & Zhuang, Y. (2010). A novel approach for movie recommendation with item-based collaborative filtering. In Proceedings of the International Conference on Multimedia and Expo (ICME) (pp. 1764-1769). IEEE.

This page is intentionally left blank to indicate the back cover. Ensure that the back cover is black in color.