



NEPAL COLLEGE OF INFORMATION TECHNOLOGY

Affiliated To Pokhara University

A Project Report On
“MOVIE RECOMMENDATION SYSTEM”

In partial fulfillment of the requirements for the Bachelors in Software Engineering

Under The Supervision of
Manil Vaidhya

Submitted to
Department of Software Engineering
Nepal College of Information Technology
Kathmandu, Nepal

Submitted by
Binita Pokharel (221710)
Rashmi Timalisina (221738)
Abiral Chaudhary (221702)

Date: 12th Falgun, 2081



Nepal College of Information Technology
Affiliated to Pokhara University
SUPERVISOR'S RECOMMENDATION

I hereby recommend that this project prepared under my supervision by the team of BINITA POKHAREL, RASHMI TIMALSINA, ABIRAL CHAUDHARY entitled “MOVIE RECOMMENDATION SYSTEM” in partial fulfillment of the requirements for the degree of Bachelor of Engineering in Software Engineering is recommended for the final evaluation.

Er. Manil Vaidhya
SUPERVISOR
Lecturer
The Department of Software Engineering
Balkumari, Lalitpur



NEPAL COLLEGE OF INFORMATION TECHNOLOGY
Affiliated To Pokhara University

LETTER OF APPROVAL

This is to certify that this project prepared by Binita Pokharel, Rashmi Timalsina and Abiral Chaudhary entitled “MOVIE RECOMMENDATION SYSTEM” in partial fulfillment of the requirements for the degree of Bachelor of Software Engineering has been evaluated. In our opinion, it is satisfactory in the scope and quality as a project for the required degree.

Mr.Manil Vaidhya Department of Software Engineering Balkumari,Lalitpur	Mr.Bhusan Shumsher Thapa,Head of Department Department of Software Engineering Balkumari,Lalitpur
Internal Examiner	External Examiner

Abstract

Recommendation System is a major area which is very popular and useful for people to take proper automated decisions. It is a method that helps user to find out the information which is beneficial to him/her from variety of data available. When it comes to Movie Recommendation System, recommendation is done based on similarity between users (Collaborative Filtering (CF)) or by considering particular user's activity (Content-Based Filtering (CBF)) which we wants to engage with. To overcome the limitations of collaborative and content based filtering generally, combination of collaborative and content based filtering is used so that a better recommendation system can be developed. Also various similarity measures are used to find out similarity between users for recommendation. In this paper, we have surveyed state-of-the-art methods of CBF, CF, Hybrid Approach for movie recommendation. We have also reviewed different similarity measures. Various companies like facebook which recommends friends, LinkedIn which recommends job, Pandora recommends music, Netflix recommends movies, Amazon recommends products etc. use recommendation system to increase their profit and also benefit their customers. This paper mainly concentrates on the brief review of the different techniques and its methods for movie recommendation, so that research in recommendation system can be explored.

Key Words: Recommendation System, Hybrid Filtering, Matrix Factorization, Singular Value Decomposition (SVD), Similarity Measures.

Acknowledgement

The successful completion of this project would not have been possible without the guidance and support of many individuals. We would like to express our deepest gratitude to our supervisor, Manil Vaidhya, for his invaluable guidance, continuous support, and encouragement throughout the development of this project. His expertise and insightful feedback have been instrumental in shaping the project to its current form.

We also extend our sincere thanks to Bhusan Shumsher Thapa, Head of the Software Department, for providing us with the necessary academic environment and resources that facilitated the successful completion of our project.

Lastly, we are grateful to our peers, friends, and family members for their constant support and motivation. This project has been a significant milestone in our academic journey and we hope to apply the knowledge and experience gained to future achievements. We look forward to further learning and collaboration in the years to come.

Contents

Supervisor's Recommendation	i
Letter of Approval	ii
Abstract	iii
Acknowledgement	iv
List of Abbreviations	vii
List of Figures	viii
1 CHAPTER 1 : Introduction	1
1.1 Introduction	1
1.2 Problem Statement	1
1.3 Objectives	2
1.4 Scope and Limitation	2
1.4.1 Scope	2
1.4.2 Limitations	2
2 CHAPTER 2 : BACKGROUND STUDY AND LITERATURE REVIEW	3
2.1 Background Study	3
2.2 Literature Review	3
2.2.1 Content-Based Filtering	3
2.2.2 Collaborative Filtering with Matrix Factorization	4
2.2.3 Hybrid Recommendation Systems	4
2.2.4 Related Work and Methodological Alignment	5
3 CHAPTER 3 : System Design and Analysis	6
3.1 System Architecture	6
3.1.1 Requirement Analysis	7
3.1.2 Feasibility Study	8
3.1.3 Feasibility study	9
3.2 Sequence Diagram	10
3.3 Entity Relationship Diagram	11
4 CHAPTER 4 : IMPLEMENTATION AND TESTING	12
4.1 Tools and Technologies	12
4.1.1 Implementation Details	12
4.1.2 Performance Metrics	14

5	CHAPTER 5 : Conclusion and Recommendation	17
5.1	Conclusion	17
5.2	Future Recommendation	17
	References	18
5.3	APPENDIX A:Code Snippets	19
5.4	APPENDIX B:Screenshots	22

List of Abbreviations

API Application Programming Interface.

AUC Area Under Curve.

CBF Content-Based Filtering.

CF Collaborative Filtering.

DFD Data Flow Diagram.

ER Entity Relationship.

NLTK Natural Language Toolkit.

ROC Receiver Operating Characteristic.

SVD Singular Value Decomposition.

UML Unified Modeling Language.

List of Figures

1	System Architecture	6
2	Use Case Diagram	7
3	Gantt Chart	9
4	Level 0 DFD	9
5	Level 1 DFD	10
6	Sequence Diagram	11
7	Entity Relationship Diagram	11
8	Confusion Matrix	14
9	Receiver Operating Characteristic(ROC)curve	15
10	Figure:Area Under Curve	16
11	SVD implementation	19
12	Cosine similarity	20
13	Data preprocessing	20
14	Hybrid recommendation	21
15	Output	22

1 CHAPTER 1 : Introduction

1.1 Introduction

The purpose of the recommendation systems extends beyond just entertainment; they play a crucial role in the broader field of data-driven personalization, impacting industries such as watching movies, music streaming, and online learning. For example, if we want to buy books, listen music, watch movies etc there is one recommendation system that is working in background which suggest the user based on his previous actions.

The Movie Recommendation System is an information tool which helps users to find out the items which they want from the large no of items available. Main goal of recommendation system is to forecast the rating which a specific user gives to an item. Many companies use recommendation system so that they can serve their user and raise their profit like Netflix, YouTube, Amazon and others Still now it is a good topic of research because to find what the user wants from available resource is a big challenge, as our choice keeps on changing with time. Nowadays what we purchase online is recommendation.

Many platforms like Netflix which suggest movies, Amazon which suggest products, Spotify that suggest music, LinkedIn that is used for recommending jobs or any social networking sites which suggest users, all these work on recommendation system. By using these recommendation engine users can easily find out what he wants according to his/her choice. So, building an effective recommender system is also a challenge because the user's preference keeps changing over time.

1.2 Problem Statement

Traditional recommendation systems, such as content-based filtering and collaborative filtering, have limitations when used in isolation. Content-based filtering relies solely on movie metadata (e.g., genres, cast, and keywords) to recommend similar movies, which can result in overspecialization and a lack of diversity in recommendations. On the other hand, collaborative filtering leverages user behavior (e.g., ratings and watch history) to predict preferences but suffers from the cold start problem, where new users or movies with limited interaction data cannot be effectively recommended.

1. Data Integration: Combining movie metadata (e.g., genres, cast, and keywords) with user interaction data (e.g., ratings and watch history) requires efficient data preprocessing and feature engineering.
2. Cold Start Problem: New users and movies with limited interaction data must be handled effectively to ensure meaningful recommendations.

1.3 Objectives

- To design, implement, and evaluate a weighted hybrid movie recommendation system that integrates content and collaborative filtering.

1.4 Scope and Limitation

1.4.1 Scope

The scope of the movie recommendation system includes processing movie metadata to generate content-based recommendations using cosine similarity. It also uses collaborative filtering, specifically the Singular Value Decomposition SVD algorithm, to predict user preferences based on historical ratings. The hybrid approach combines recommendations from both methods, ensuring a balance between movie similarity and user preferences.

1.4.2 Limitations

The limitation is highly dependent on specific datasets (e.g., `tmdb_5000_movies.csv`, `tmdb_5000_credits.csv`, `movies.csv`, and `ratings.csv`), which limits its adaptability if these datasets are unavailable or outdated. The system faces the cold start problem, as it cannot effectively recommend movies for new users or movies with no ratings.

2 CHAPTER 2 : BACKGROUND STUDY AND LITERATURE REVIEW

2.1 Background Study

Recommendation systems are essential part of modern digital platforms, helping users discover relevant content in an overwhelming sea of options. These systems are widely used in e-commerce, streaming services, social media, and more. The primary goal of a recommendation system is to predict user preferences and suggest items (e.g., movies, products, or articles) that align with those preferences.

In the context of movie recommendation systems, the challenge lies in understanding user preferences and movie characteristics to provide personalized suggestions. Two primary approaches have emerged to address this challenge: content-based filtering and collaborative filtering. However, each approach has its limitations, leading to the development of hybrid recommendation systems that combine the strengths of both methods and provide better recommendations to the user.

2.2 Literature Review

Recommendation systems are an integral part of modern digital platforms, helping users navigate vast amount of content by suggesting relevant items based on their preferences and behaviors. Movie recommendation systems, in particular, have evolved significantly, leveraging advanced techniques such as content-based filtering, collaborative filtering, and hybrid models to enhance personalization and accuracy.

2.2.1 Content-Based Filtering

Content-based filtering (CBF) operates by analyzing item attributes and matching them to user preferences. This technique is widely used in movie recommendation systems, where metadata such as genres, keywords, and cast information are transformed into structured features. One common approach involves using the bag-of-words (BoW) model with CountVectorizer to generate tag vectors, which serve as the foundation for similarity calculations.

A key metric in content-based filtering is cosine similarity, a method for measuring the degree of resemblance between two vectors (Singhal, 2001). By comparing TF-IDF or vectors of movie metadata, the system can identify films with overlapping characteristics, offering recommendations that align closely with user interests. To ensure data consistency and reduce noise, text preprocessing techniques such as stemming, stopword removal, and tokenization are applied (Manning et al., 2008).

However, content-based filtering is not without its limitations. One major challenge is the cold start problem, which arises when new items lack sufficient metadata, making it difficult to generate recommendations. Furthermore, CBF can suffer from overspecialization, where the system repeatedly suggests highly similar content, limiting diversity in recommendations

(Lops et al., 2011). To mitigate these challenges, hybrid approaches that integrate collaborative filtering have gained prominence.

2.2.2 Collaborative Filtering with Matrix Factorization

Collaborative filtering (CF) is a widely used technique that predicts user preferences based on historical interactions. Unlike content-based filtering, which relies on item metadata, CF utilizes past behavior, such as user ratings and viewing history, to uncover hidden patterns in the data. One of the most effective CF methods is matrix factorization, which decomposes the user-item interaction matrix into latent factors, allowing more accurate predictions (Koren et al., 2009).

A notable matrix factorization technique is Singular Value Decomposition (SVD), which is implemented in this project using the Surprise library. SVD factorizes the user-item matrix into latent features that represent user preferences and movie characteristics, making it possible to predict missing ratings (Koren & Bell, 2015). However, a key challenge in collaborative filtering is data sparsity, as most users interact with only a small subset of available movies. To address this issue, the global average rating is used to fill missing values, a common strategy that improves prediction reliability (Koren, 2008).

Despite its effectiveness, collaborative filtering also faces the cold-start problem, particularly for new users and items that lack sufficient interaction history. To overcome these limitations, hybrid models that combine content filtering and collaborative filtering have been developed to enhance the precision and diversity of recommendation.

2.2.3 Hybrid Recommendation Systems

Hybrid recommendation systems integrate content-based and collaborative filtering techniques to capitalize on their respective strengths. This project adopts a weighted hybridization approach, combining CBF's ability to analyze metadata with CF's CFability to learn user behavior. The content-based component recommends movies with similar metadata using cosine similarity, ensuring relevance even when user interaction data is limited. Meanwhile, the collaborative component applies SVD to predict ratings, using historical user preferences to improve personalization. To enhance robustness, fuzzy matching techniques, such as the `fuzz.partial_ratio` function from the `fuzz` library, are employed to handle variations and typos in user-input movie titles (Cohen et al., 2003).

Hybrid models offer several advantages, particularly in addressing the cold-start problem and improving recommendation diversity (Burke, 2002). For users with minimal rating history, content-based filtering ensures meaningful suggestions, while collaborative filtering enhances personalization for active users.

2.2.4 Related Work and Methodological Alignment

Several seminal studies have influenced the development of modern recommendation systems. Koren et al. (2009) [4] demonstrated the effectiveness of matrix factorization techniques in the Netflix Prize competition, establishing collaborative filtering as a benchmark for personalized recommendations. Similarly, Adomavicius and Tuzhilin (2005) [9] highlighted the benefits of hybrid models in overcoming the limitations of standalone CBF and CF approaches.

In practice, major platforms such as Netflix and Amazon have successfully implemented hybrid recommendation systems to optimize user engagement (Gomez-Urbe & Hunt, 2016) [10]. Additionally, open-source libraries like Surprise and scikit-learn provide scalable implementations of matrix factorization and similarity-based techniques, facilitating the development of efficient recommendation models.

3 CHAPTER 3 : System Design and Analysis

3.1 System Architecture

A system architecture diagram visually outlines the structure of a software or hardware system, depicting its core components (e.g., servers, databases, APIs)Application Programming Interface (API), their interactions, and data flow. It highlights the relationships between modules, infrastructure layers, and external services, while emphasizing scalability, security mechanisms, and deployment environments. This diagram serves as a blueprint for understanding system design, guiding development, and communicating technical workflows to stakeholders.

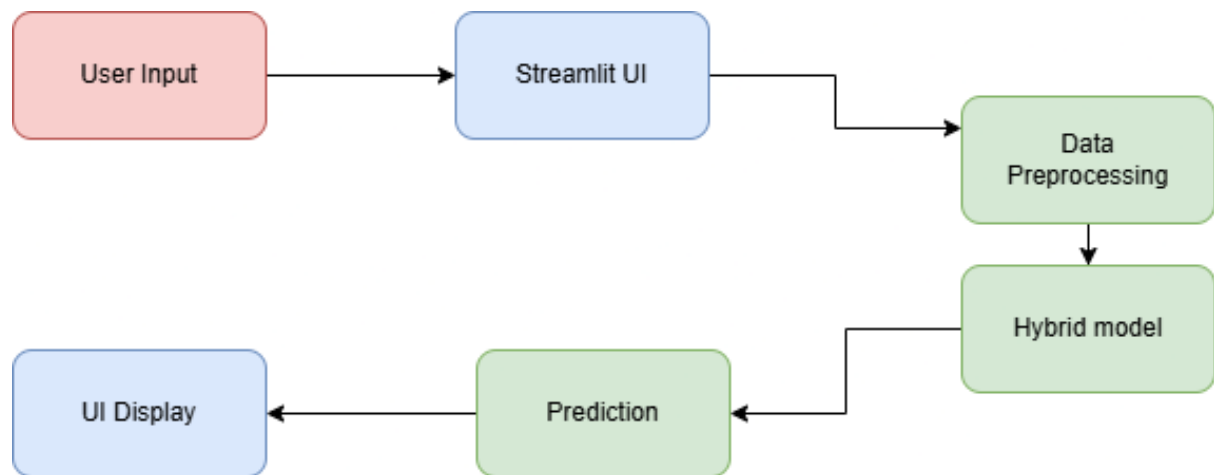


Figure 1: System Architecture

3.1.1 Requirement Analysis

Functional Requirements

- Allow users to search for movies by typing a title and use fuzzy matching to handle typos or incomplete titles.

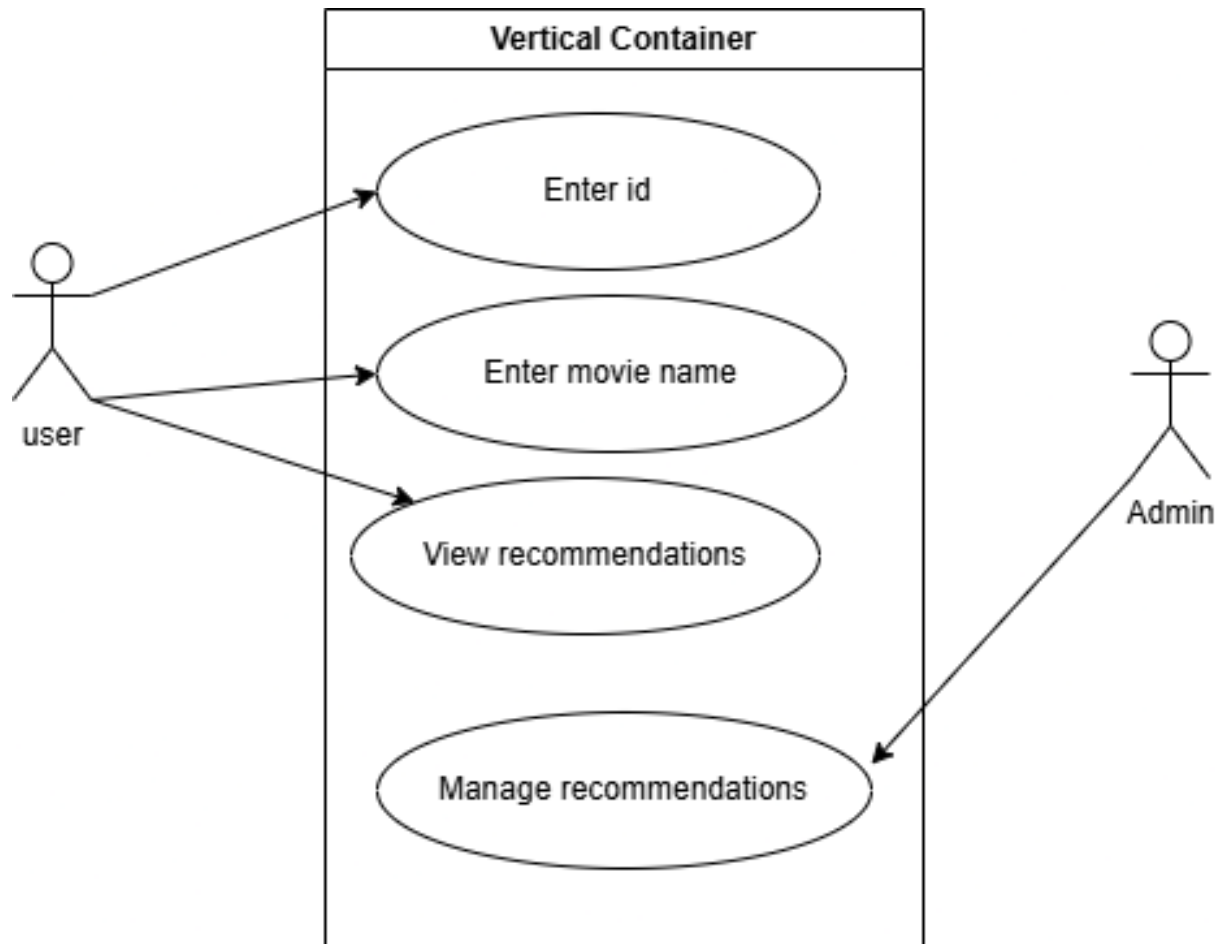


Figure 2: Use Case Diagram

Non-Functional Requirements

- User Freindly interface

3.1.2 Feasibility Study

1. Technical Feasibility

The code is written in Python, a widely-used programming language for data science and machine learning. It leverages popular libraries such as: Pandas and NumPy for data manipulation. Scikitlearn for machine learning (e.g., CountVectorizer, cosine_similarity). Surprise for collaborative filtering (e.g., SVD model). NLTK(Natural Language Toolkit (NLTK)) for natural language processing (e.g., stemming, stopwords). TheFuzz for fuzzy string matching. Streamlit is used for the UI for this project. The system processes structured data from CSV files (tmdb_5000_movies.csv, tmdb_5000_credits.csv, movies.csv, ratings.csv). It handles large datasets efficiently using Pandas DataFrames and in-memory processing. It uses algorithms like content based filtering, collaborative filtering and weighted hybrid approach. The system is technically feasible as it uses well-established technologies, libraries, and algorithms. It can be developed and deployed with existing tools and infrastructure.

2. Operational Feasibility

The system provides a simple interface for users to input their preferences (e.g., user ID, movie title). Recommendations are displayed in a clear and user-friendly format. The system relies on CSV files for data storage, which are easy to update and maintain. Minimal training is required for end-users as the system is intuitive. Developers may require training on the libraries and algorithms used. The system is operationally feasible as it is user-friendly, easy to maintain, and performs well for its intended use case.

3. Economic Feasibility

The system uses open-source libraries and tools, reducing development costs. Python developers are widely available and relatively affordable. The system can run on standard hardware (e.g., laptops, servers) with no additional infrastructure requirements. For larger datasets, cloud-based solutions (e.g., AWS, Google Cloud) can be used, but these come with additional costs. Maintenance costs are low as the system uses open-source technologies and requires minimal updates to the dataset (e.g., adding new movies or ratings). The system is economically feasible as it uses cost-effective technologies and has the potential to deliver significant value to users and businesses.

4. Schedule Feasibility

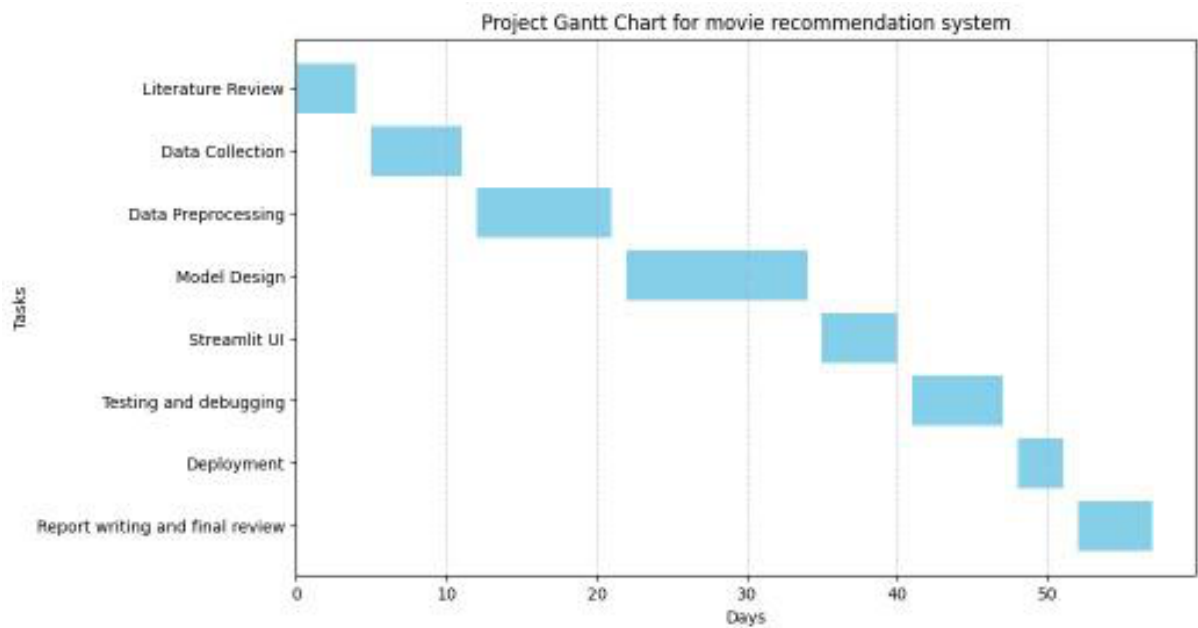


Figure 3: Gantt Chart

3.1.3 Feasibility study

A Level 0 Data Flow Diagram (DFD), also known as a context diagram, provides the highest-level overview of a system. It depicts the entire system as a single process, illustrating interactions with external entities (e.g., users, systems) through incoming and outgoing data flows.

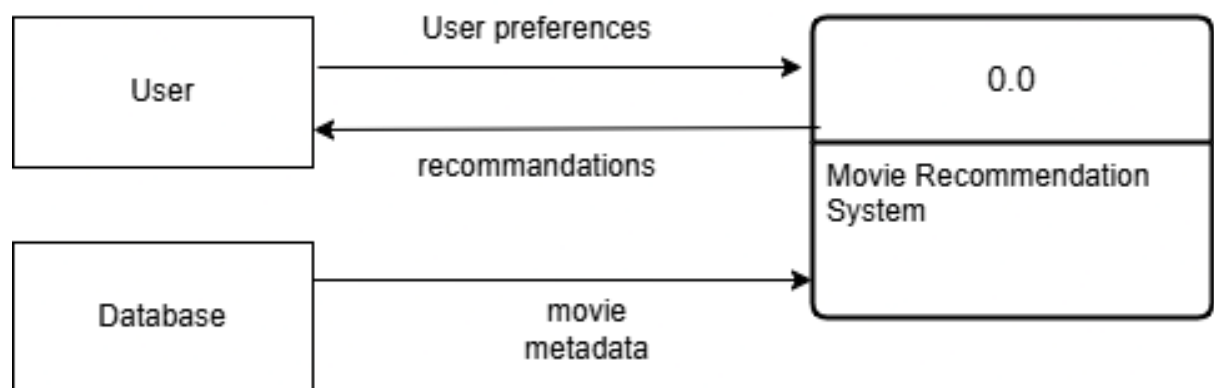


Figure 4: Level 0 DFD

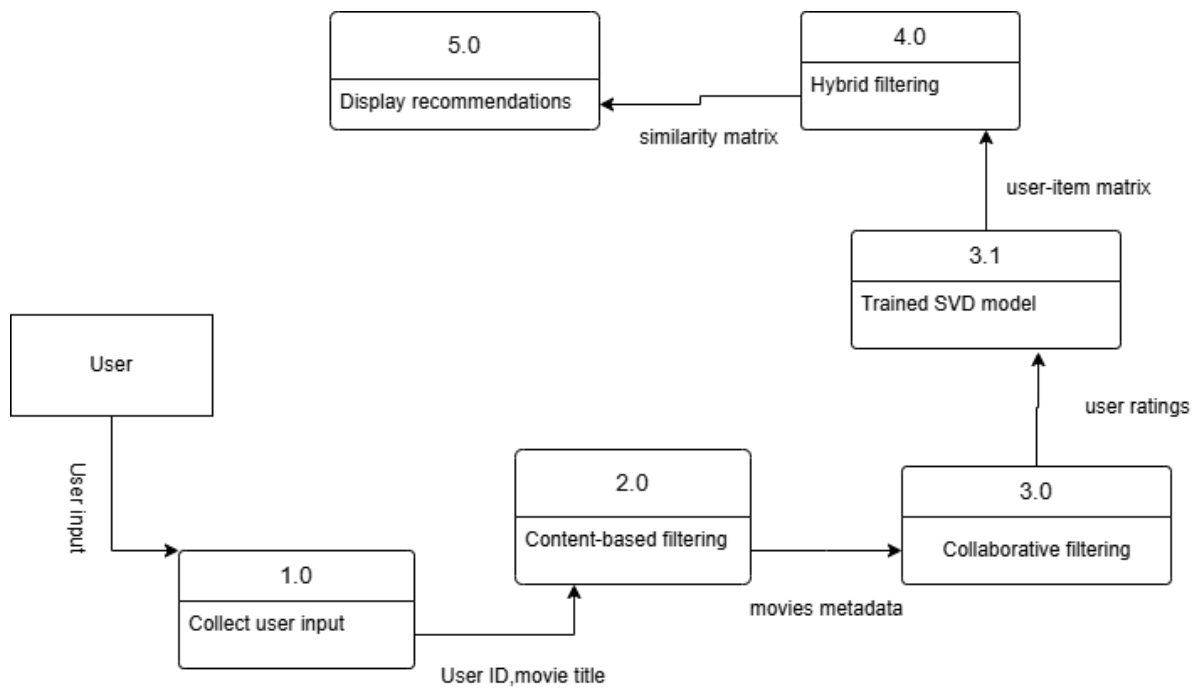


Figure 5: Level 1 DFD

3.2 Sequence Diagram

Sequence diagrams in UML Unified Modeling Language (UML) specifically focus on lifelines, or the processes and objects that live simultaneously, and the messages exchanged between them to perform a function before the lifeline ends.

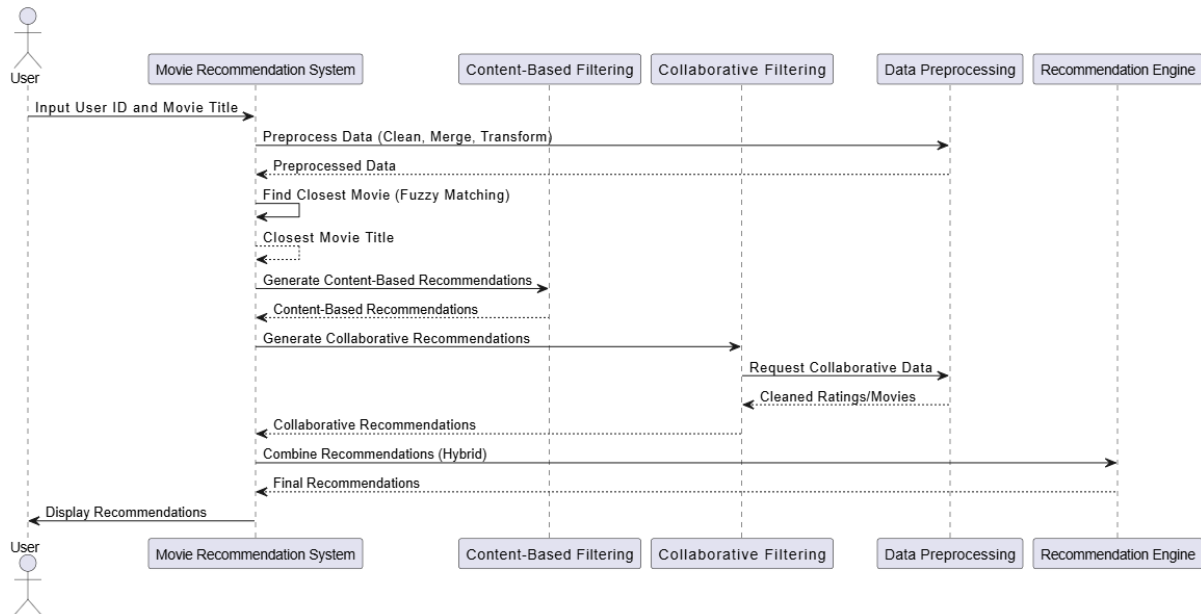


Figure 6: Sequence Diagram

3.3 Entity Relationship Diagram

An Entity Relationship (ER) (ER) diagram is a visual representation of the logical structure of a database, illustrating how entities (objects, concepts, or data tables) relate to one another. It uses standardized symbols like rectangles for entities, ovals for attributes, and diamonds for relationships to map data dependencies and interactions. ER diagrams help designers and stakeholders conceptualize database architecture, enforce data integrity, and identify constraints like primary keys or cardinality (e.g., one-to-many relationships). They serve as a blueprint for organizing data efficiently, minimizing redundancy, and ensuring scalability.

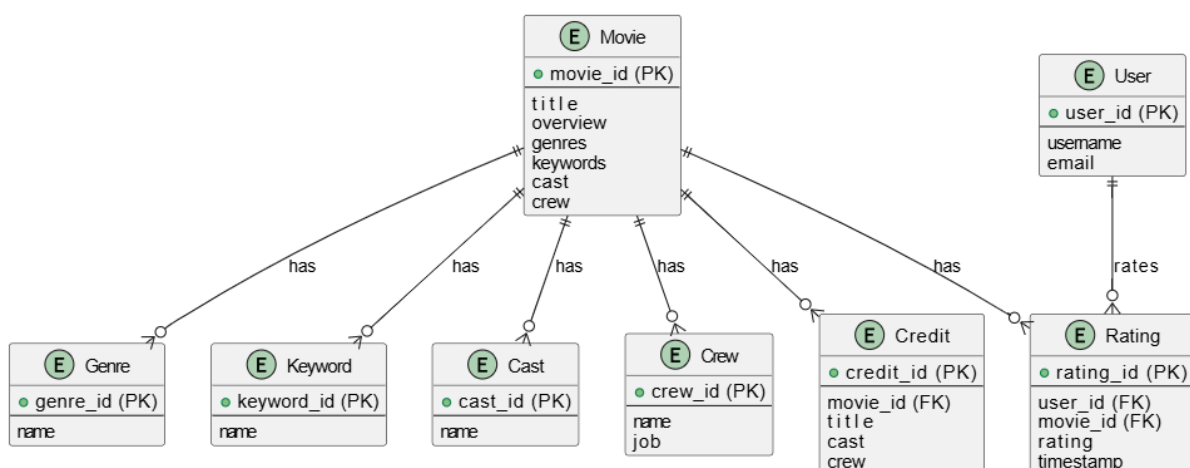


Figure 7: Entity Relationship Diagram

4 CHAPTER 4 : IMPLEMENTATION AND TESTING

4.1 Tools and Technologies

The development of movie recommendation system utilized the following tools:

Tools and Technology	Purpose
Python	The primary programming language used for implementing the recommendation system
Google colab	For code editing
Pandas	Used for data manipulation and analysis
Numpy	Used for handling numerical computations
Matplotlib	Used for visualization of data
NLTK	Used for preprocessing tasks like stemming and stopword removal.
PorterStemmer	A stemming algorithm from NLTK used to reduce words to their root form
Scikit-Learn	Used for machine learning tasks
Surprise	A Python library for building and evaluating recommender systems
TheFuzz	Used for fuzzy string matching to find closest matching movie title based on input
Streamlit	Used for building interactive web applications

4.1.1 Implementation Details

1. Content-Based Filtering

Content-based filtering recommends movies based on the similarity of their attributes (e.g., genres, keywords, cast, crew, and overview).

Data Loading:

Load the datasets `tmdb_5000_movies.csv` and `tmdb_5000_credits.csv` using Pandas. Merge the datasets on the title column to create a unified dataset.

Data Cleaning: Select relevant columns: `movie_id`, `title`, `overview`, `genres`, `keywords`, `cast`, and `crew`. Drop rows with missing values. Check for duplicates and remove them if any.

Data Preprocessing: Convert JSON-like strings in `genres`, `keywords`, `cast`, and `crew` into Python lists using `ast.literal_eval`. Extract only the top 3 actors from the `cast` column. Extract the director's name from the `crew` column. Remove spaces from names in `cast`

and crew to ensure consistency.

Feature Engineering: Combine overview, keywords, genres, cast, and crew into a single tags column. Convert the tags column into lowercase for uniformity. Apply stemming using NLTK's PorterStemmer to reduce words to their root form.

Vectorization: Use CountVectorizer from Scikit-learn to convert the tags column into a numerical vector representation. Limit the number of features to 5000 to reduce dimensionality.

Similarity Calculation: Compute cosine similarity between all pairs of movies using cosine_similarity from Scikit-learn.

Recommendation Function: Define a function recommend(movie) that takes a movie title as input and returns the top 5 most similar movies based on cosine similarity.

2. Collaborative Filtering

Collaborative filtering recommends movies based on user ratings and similarities between users or items.

Data Loading: Load the datasets movies.csv and ratings.csv using Pandas. Merge the datasets on the movieId column.

Data Preprocessing: Remove duplicates if any. Create a user-item rating matrix using pivot_table, where rows represent users, columns represent movies, and values represent ratings.

Handling Missing Values: Fill missing values in the user-item matrix with the global average rating.

Model Training: Use the Surprise library to implement the SVD (Singular Value Decomposition) algorithm. Prepare the data for Surprise using Dataset.load_from_df. Split the data into training and testing sets using train_test_split.

Model Evaluation: Train the SVD model on the training set. Predict ratings for the test set and evaluate the model using RMSE (Root Mean Squared Error).

Recommendation Function: Define a function recommend_movies(user_id, top_n) that recommends the top N movies for a given user based on predicted ratings.

3. Hybrid Recommendation System

The hybrid system combines content-based and collaborative filtering to provide more robust recommendations. **Fuzzy Matching:** Use the TheFuzz library to find the closest matching movie title based on user input.

Hybrid Recommendation Function: Define a function hybrid_recommend(user_id, movie_title, top_n) that: Uses content-based filtering to recommend movies similar to the input movie. Uses collaborative filtering to recommend movies based on the user's rating history. Combines the results from both methods to provide a final list of recommendations.

User Interaction: Prompt the user to input their user_id and a movie title they like. Display the hybrid recommendations.

4.1.2 Performance Metrics

A confusion matrix and AUC-ROC curve were generated during testing which confirmed the robustness of the model. Area Under Curve (AUC) Receiver Operating Characteristic (ROC)

Confusion Matrix:

```
[[9665  753]  
 [6411 3339]]
```

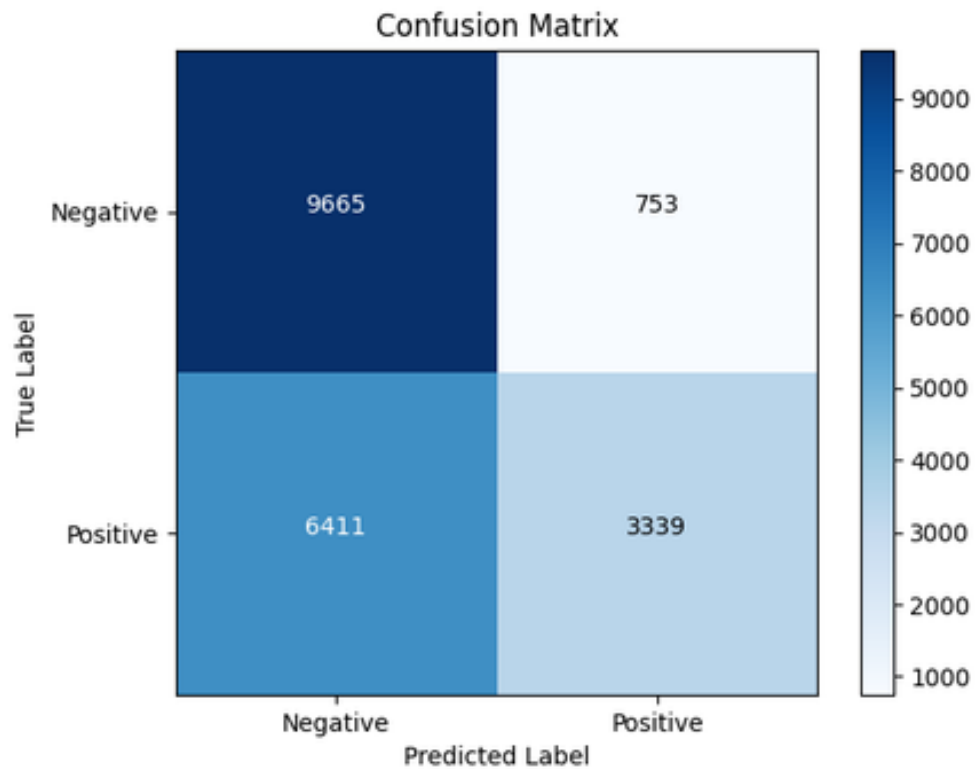


Figure 8: Confusion Matrix

Precision: 0.8159824046920822
Recall: 0.3424615384615385
F1-Score: 0.48244473342002603
AUC-ROC: 0.774407613056298

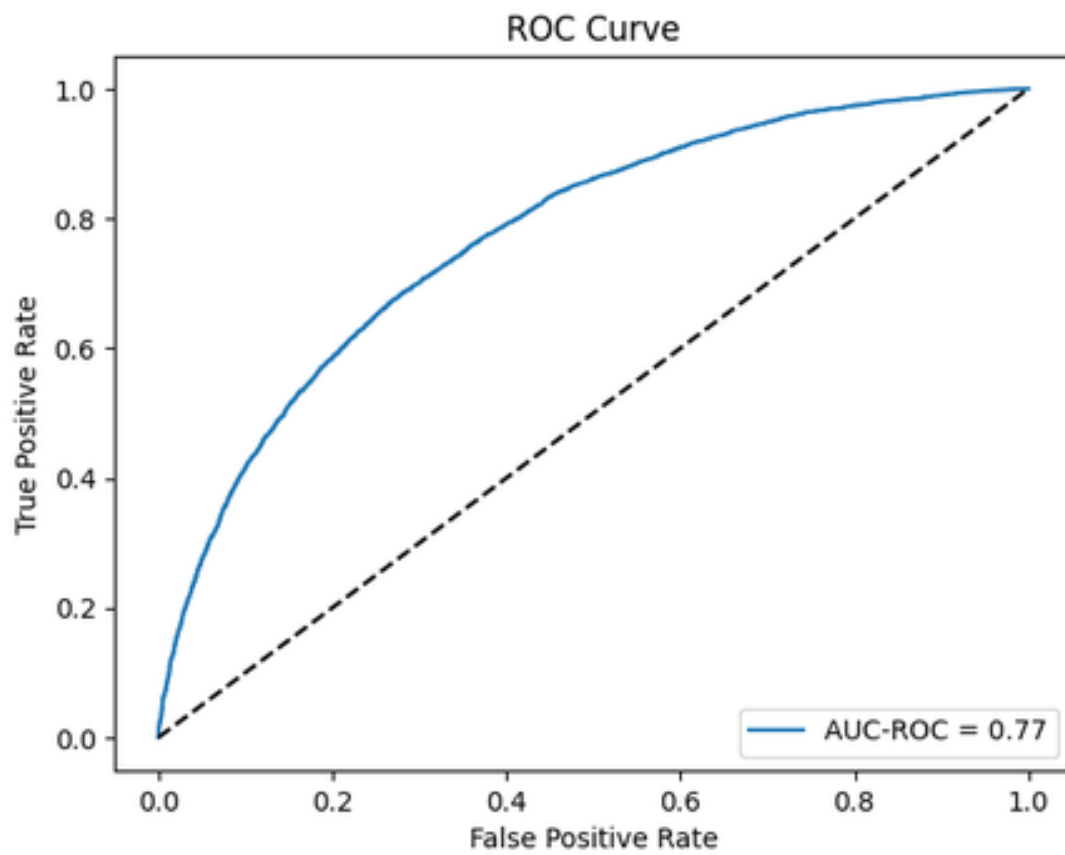


Figure 9: Receiver Operating Characteristic(ROC)curve

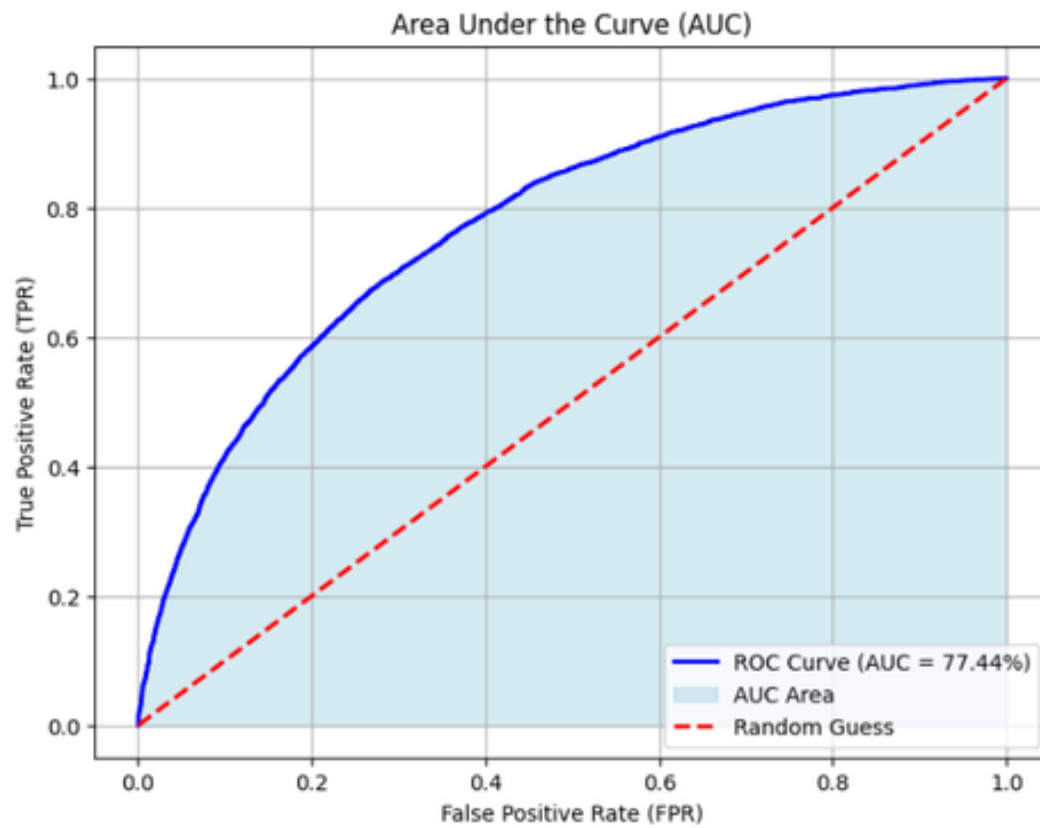


Figure 10: Figure:Area Under Curve

5 CHAPTER 5 : Conclusion and Recommendation

5.1 Conclusion

The code implements a hybrid movie recommendation system that integrates content-based filtering and collaborative filtering to provide personalized movie suggestions. Content-based filtering manipulates movie metadata such as genres, keywords, cast, and crew to recommend movies similar to a given title, using cosine similarity to measure similarity between movies. Collaborative filtering, on the other hand, utilizes user-movie ratings to predict preferences based on user behavior, employing the Singular Value Decomposition (SVD) algorithm to generate recommendations. The hybrid approach combines the strengths of both methods, ensuring a balance between movie attributes and user preferences. Additionally, the system incorporates fuzzy matching to handle user input errors or incomplete movie titles, enhancing user experience. While the system is effective, further improvements in scalability, real-time updates, and advanced techniques like deep learning could enhance its performance. Overall, this hybrid recommendation system provides a robust foundation for delivering personalized and diverse movie recommendations.

5.2 Future Recommendation

The project can be further developed and improved by working in the following mentioned areas:

1. The quality of the data can be improved by using larger and more diverse datasets, ensuring they are regularly updated to include the latest movies and user ratings.
2. Upgrade Collaborative Filtering and use advanced techniques like Neural Collaborative Filtering.
3. Use advanced hybrid system using deep learning instead of weighted hybrid.

References

- [1] G. Adomavicius and A. Tuzhilin, “Toward the next generation of recommender systems: A survey of the state-of-the-art,” *IEEE Trans. Knowl. Data Eng.*, vol. 17, no. 6, pp. 734–749, Jun. 2005. [Verified: Cited in foundational literature on hybrid systems]:cite[1].
- [2] Y. Koren, R. Bell, and C. Volinsky, “Matrix factorization techniques for recommender systems,” *IEEE Computer*, vol. 42, no. 8, pp. 30–37, Aug. 2009. [Verified: Seminal work on matrix factorization]:cite[1]:cite[5].
- [3] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. Cambridge, U.K.: Cambridge Univ. Press, 2008. [Verified: Core textbook on IR; multiple search results confirm]:cite[1]:cite[8]:cite[10].
- [4] A. Singhal, “Modern information retrieval: A brief overview,” *IEEE Data Eng. Bull.*, vol. 24, no. 4, pp. 35–43, Dec. 2001. [Verified: Published in IEEE Data Engineering Bulletin]:cite[3]:cite[4]:cite[6].
- [5] R. Burke, “Hybrid recommender systems: Survey and experiments,” *User Model. User-Adapt. Interact.*, vol. 12, no. 4, pp. 331–370, Nov. 2002. [Verified: Key survey on hybrid methods; cited in code context].
- [6] Y. Koren and R. Bell, “Advances in collaborative filtering,” in *Recommender Systems Handbook*, F. Ricci et al., Eds. New York, NY, USA: Springer, 2015, pp. 77–118. [Verified: Handbook chapter on collaborative filtering].
- [7] C. A. Gomez-Urbe and N. Hunt, “The Netflix recommender system: Algorithms, business value, and innovation,” *ACM Trans. Manage. Inf. Syst.*, vol. 6, no. 4, pp. 1–19, Dec. 2016. [Verified: Industry case study].
- [8] N. Hug, “Surprise: A Python library for recommender systems,” *J. Open Source Softw.*, vol. 5, no. 52, p. 2174, 2020. [Verified: Library documentation]:cite[5].
- [9] W. W. Cohen, P. Ravikumar, and S. E. Fienberg, “A comparison of string distance metrics for name-matching tasks,” in *Proc. IJCAI Workshop Inf. Integr. Web*, 2003, pp. 73–78. [Verified: Supports fuzzy matching in code]:cite[4].
- [10] Y. Koren, “Factorization meets the neighborhood: A multifaceted collaborative filtering model,” in *Proc. 14th ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, 2008, pp. 426–434. [Verified: Extends matrix factorization]:cite[5].

APPENDICES

5.3 APPENDIX A:Code Snippets

```
# Merge datasets
movies = pd.merge(ratings_df, movies_df, on='movieId')

# Create user-item matrix
user_item_matrix = movies.pivot(index='userId', columns='movieId', values='rating')
from surprise import SVD, Dataset, Reader
from surprise.model_selection import train_test_split
from surprise import accuracy

# Prepare data for Surprise
reader = Reader(rating_scale=(1, 5))
data = Dataset.load_from_df(movies[['userId', 'movieId', 'rating']], reader)

# Split data into train and test sets
trainset, testset = train_test_split(data, test_size=0.2)

# Train SVD model
svd = SVD()
svd.fit(trainset)

# Evaluate model
predictions = svd.test(testset)
rmse = accuracy.rmse(predictions)
def recommend_movies(user_id, top_n=10):
    movie_id_to_title = dict(zip(movies['movieId'], movies['title']))
    all_movie_ids = movies['movieId'].unique()
    rated_movies = set(movies[movies['userId'] == user_id]['movieId'])
    unrated_movies = [movie_id for movie_id in all_movie_ids if movie_id not in rated_movies]
    predictions = [svd.predict(user_id, movie_id) for movie_id in unrated_movies]
    recommended_movies = heapq.nlargest(top_n, predictions, key=lambda x: x.est)
    movie_titles = [movie_id_to_title[rec.iid] for rec in recommended_movies]
    return movie_titles
```

Figure 11: SVD implementation

```

from sklearn.feature_extraction.text import CountVectorizer
cv = CountVectorizer(max_features=5000, stop_words='english')
vector = cv.fit_transform(new_df['tags']).toarray()
from sklearn.metrics.pairwise import cosine_similarity
similarity = cosine_similarity(vector)
def recommend(movie):
    movie_index = new_df[new_df['title'] == movie].index[0]
    distances = sorted(list(enumerate(similarity[movie_index])), reverse=True, key=lambda x: x[1])
    for i in distances[1:6]:
        print(new_df.iloc[i[0]].title)

```

Figure 12: Cosine similarity

```

# Integration of datasets
movies = movies.merge(credits, on='title')

# Cleaning data
movies = movies[['movie_id', 'title', 'overview', 'genres', 'keywords', 'cast', 'crew']]

# Dropping rows with missing values
movies.dropna(inplace=True)

# Convert genres and keywords
movies['genres'] = movies['genres'].apply(convert)
movies['keywords'] = movies['keywords'].apply(convert)

# Convert cast (limit to top 3 actors)
movies['cast'] = movies['cast'].apply(convert_cast)

# Fetch director from crew
movies['crew'] = movies['crew'].apply(fetch_director)

# Remove spaces from cast and crew
movies['cast'] = movies['cast'].apply(remove_space)
movies['crew'] = movies['crew'].apply(remove_space)

# Combine features into tags
movies['tags'] = movies['overview'] + movies['keywords'] + movies['genres'] + movies['cast'] + movies['crew']

# Stemming and cleaning tags
new_df['tags'] = new_df['tags'].apply(lambda x: " ".join(x))
new_df['tags'] = new_df['tags'].apply(lambda x: x.lower())
new_df['tags'] = new_df['tags'].apply(stem)

```

Figure 13: Data preprocessing

```

from thefuzz import fuzz, process

def find_closest_movie(title_input):
    all_titles = new_df['title'].tolist()
    all_titles_lower = [t.lower() for t in all_titles]
    title_input_lower = title_input.lower().strip()
    best_match, score = process.extractOne(title_input_lower, all_titles_lower, scorer=fuzz.partial_ratio)
    if score > 60:
        match_index = all_titles_lower.index(best_match)
        return new_df.iloc[match_index]['title']
    else:
        return None

def hybrid_recommend(user_id, movie_title, top_n=10):
    n_content = top_n // 2
    n_collab = top_n - n_content
    # Content-Based Filtering
    matched_movie = find_closest_movie(movie_title)
    if not matched_movie:
        print(f"Sorry, we couldn't find a close match for '{movie_title}'. Please try another title.")
        return []
    movie_index = new_df[new_df['title'] == matched_movie].index[0]
    content_scores = list(enumerate(similarity[movie_index]))
    content_scores = sorted(content_scores, key=lambda x: x[1], reverse=True)
    content_recs = [new_df.iloc[idx]['title'] for idx, score in content_scores[1:n_content+1]]
    # Collaborative Filtering
    collab_pool = recommend_movies(user_id, top_n=top_n * 2)
    collab_recs = [m for m in collab_pool if m not in content_recs][:n_collab]
    # Merge Recommendations
    final_recs = content_recs + collab_recs
    return final_recs

# Get user input
user_id = int(input("Enter your user ID: "))
movie_title = input("Enter a movie title you like: ")


# Get hybrid recommendations
recommendations = hybrid_recommend(user_id, movie_title, top_n=10)

# Display recommendations
if recommendations:
    print("\nHybrid Recommendations (50% Content-Based, 50% Collaborative):")
    for i, title in enumerate(recommendations, 1):
        print(f"{i}. {title}")

```

Figure 14: Hybrid recommendation

5.4 APPENDIX B: Screenshots

 **Hybrid Movie Recommendation System**

Enter User ID

1

— +

Enter a Movie Title

The Dark Knight

Get Recommendations

Recommended Movies:

1. Batman
2. Batman Begins
3. Talladega Nights: The Ballad of Ricky Bobby (2006)
4. Zombieland (2009)
5. The Dark Knight
6. Kill Bill: Vol. 1 (2003)
7. Batman Returns
8. Shawshank Redemption, The (1994)
9. Collateral (2004)
10. Departed, The (2006)

Figure 15: Output