# GOVERNMENT COLLEGE OF ENGINEERING

# CHHATRAPATI SAMBHAJINAGAR

**(An Autonomous Institute of Government of Maharashtra)**

**The Project Report of the " Movie Recommendation System using Flask with Collaborative and Content-Based Filtering "**

*In partial fulfilment of the requirements for the degree of*

*Bachelor of Technology*

*In*

*Electronics and Telecommunication Engineering*

Submitted by

**Zaheer UL Islam Khan    [BE21F04F038]**

**Ashutosh munde            [BE21F04F048]**

Under the Esteemed Guidance of

**Dr. S. R. HIREKHAN**



**DEPARTMENT OF ELECTRONICS & TELECOMMUNICATION**

**ENGINEERING**

**GOVERNMENT COLLEGE OF ENGINEERING AURANGABAD (Chhatrapati Sambhajinagar) – 431001**

**| 2024-2025  |**

# CERTIFICATE

This is to certify that the report of the **" Movie Recommendation System using Flask with Collaborative and Content-Based Filtering "**, which is being submitted herewith for the award of the Bachelor's Degree of Electronics and Telecommunication Engineering at Government College of Engineering Aurangabad affiliated to Dr. Babasaheb Ambedkar Marathwada University, Chhatrapati Sambhajinagar, is the result of the study, work & contribution by **Zaheer UL Islam Khan [BE21F04F038]** and **Ashutosh Munde [BE21F04F048]** under my supervision and guidance.

**Place:**

**Date:**

 

**Dr. S.R.HIREKHAN**                    **Dr. S.R.HIREKHAN**

Guide                                                        Head Of Department
Department of Electronics and              Department of Electronics and
Telecommunication Engineering          Telecommunication Engineering

 

**Dr. S. S. Dambhare**

Principal
Government College of Engineering, Aurangabad (Chhatrapati Sambhajinagar)-431001

# THESIS APPROVAL SHEET

**Ashutosh Munde [BE21F04F048] and Zaheer UL Islam Khan [BE21F04F038]** have done the appropriate work related to **" Movie Recommendation System using Flask with Collaborative and Content-Based Filtering "**, for the award of Bachelor of Technology (Electronics & Telecommunication), is being submitted to Government College of Engineering, Chhatrapati Sambhajinagar.

**External Examiner:**

**Sign:**

**Guide: Prof.  Dr. S. R. HIREKHAN**

**Place:** Government College Of Engineering, Chhatrapati Sambhajinagar

**Date:** 03/05/2025

# DECLARATION

We hereby declare that the project work entitled **"Movie Recommendation System Using Collaborative and Content-Based Filtering",** submitted to the Government College of Engineering, Chhatrapati Sambhajinagar is a record of an original work done by us under the guidance of **Dr. S.R.Hirekhan**. This project work is submitted in the partial fulfilment of the requirements for the award of the degree of Bachelor of Technology in the Electronics and Telecommunication Engineering Department.

**Place:** Government College of Engineering, Chhatrapati Sambhajinagar

**Date:** 03/05/2025

# ACKNOWLEDGEMENT

We express our sincere gratitude to our guide, **Dr. S. R. Hirekhan**, for his invaluable guidance, support, and encouragement throughout the course of this project. Her expertise and insightful feedback helped shape this project to completion.

We also extend our heartfelt thanks to **Dr. S. S. Dambhare**, Principal of Government College of Engineering, Chhatrapati Sambhajinagar, for providing us with the opportunity and facilities to work on this project.

Special thanks to the faculty and staff of the Department of Electronics and Telecommunication Engineering for their constant support. We also acknowledge the help and motivation provided by our peers, friends, and family throughout this journey.

**Ashutosh Munde [BE21F04F048]**

**Zaheer UL Islam Khan [BE21F04F038]**

# ABSTRACT

In the era of digital entertainment, users are faced with an ever-increasing volume of movies available across numerous platforms. Navigating this vast landscape to find films aligned

with individual tastes presents a significant challenge. This project addresses this challenge by developing a content-based Movie Recommendation System using Python, Flask, and machine learning techniques. The core objective is to provide users with personalized movie suggestions by identifying films that share similar characteristics with a movie they already like.

The system employs a content-based filtering approach, which relies on analyzing the attributes of the movies themselves, such as genres, keywords, cast, and crew. Unlike collaborative filtering methods that depend on user interaction data (like ratings), the content-based approach can recommend new or niche movies without requiring extensive user feedback on those specific items.

The technical foundation of the project is built upon a robust stack including Python as the primary programming language, Flask as a lightweight web framework for creating the user interface, and essential libraries such as Pandas and NumPy for efficient data handling and manipulation. Machine learning functionalities are provided by Scikit-learn, specifically for tasks like feature extraction and calculating movie similarity. The Natural Language Toolkit (NLTK) is utilized for text processing of movie metadata. To ensure ease of deployment and portability, the entire application is containerized using Docker.

The system workflow involves several key stages: Initially, a comprehensive dataset containing movie titles, genres, keywords, and other relevant metadata is preprocessed. This involves cleaning the data, handling missing values, and combining relevant text features. Subsequently, feature engineering is performed using techniques like Term Frequency-

Inverse Document Frequency (TF-IDF) vectorization to convert the text data into numerical representations. A similarity matrix is then computed using cosine similarity, quantifying the likeness between every pair of movies based on their vectorized features. This matrix serves as the core of the recommendation engine. The preprocessed data and the similarity matrix are serialized using Python's pickle module for persistence and quick loading.

The user interacts with the system through a simple and intuitive web interface developed using HTML and styled with Bootstrap. Users select a movie from a dropdown list populated with titles from the dataset. Upon submission, the Flask application processes the request, uses the pre-computed similarity matrix to find the top 20 most similar movies to the selected one, and retrieves additional information, including movie posters, by integrating with the TMDB API. The recommendations, along with their posters, are then displayed to the user on the web page.

For streamlined deployment, a Dockerfile is provided, enabling the application to be built into a portable container image. This ensures that the application runs consistently across different environments, eliminating dependency issues. The container can be easily built and run using standard Docker commands, making the system readily accessible via a web browser.

While the current implementation effectively demonstrates a content-based recommendation approach, future enhancements could include incorporating user-based collaborative filtering to provide more personalized recommendations, adding user authentication and the ability to rate movies, and exploring deployment on cloud platforms like Heroku or AWS to make the system publicly available.

This project successfully integrates data processing, machine learning, web development, and containerization to create a practical and functional movie recommendation system, serving as a valuable learning experience in applying these technologies to solve a real-world problem.

## TABLE OF CONTENTS

# 1. INTRODUCTION

## 1.1 Background of the Project:

In the digital age, the consumption of media, particularly movies, has surged

dramatically with the proliferation of streaming platforms and online libraries. This abundance, while offering unparalleled choice, also presents a significant challenge: how do users efficiently discover movies that align with their individual tastes and preferences? Recommendation systems have emerged as powerful tools to address this information overload by suggesting items that users are likely to enjoy. These systems analyze user data, item data, or a combination of both to generate personalized recommendations.

Movie recommendation systems, in particular, play a vital role in enhancing user experience on streaming services and movie databases. They help users navigate vast catalogs, discover hidden gems, and spend less time searching and more time watching. The effectiveness of a movie recommendation system directly impacts user engagement and satisfaction.

Recommendation systems can broadly be categorized into content-based filtering, collaborative filtering, and hybrid approaches. Content-based systems recommend items similar to those a user has liked in the past, based on item attributes. Collaborative filtering systems recommend items that users with similar tastes have liked. Hybrid systems combine aspects of both. This project focuses on developing a content-based movie recommendation system, leveraging the rich metadata associated with movies to provide relevant suggestions.

## 1.2     Motivation:

The motivation for undertaking this project is multifaceted. Firstly, it provides a practical opportunity to apply fundamental machine learning concepts and techniques, such as feature extraction, similarity calculation, and model serialization, to a tangible and widely applicable problem. Building a functional recommendation system from scratch offers valuable hands-on experience in the end-to-end process of developing a data-driven application.

Secondly, the project allows for exploration of web development using Flask, integrating a machine learning model into a web application, and interacting with external APIs (like TMDB) to enhance functionality. This integration of different technologies is representative of many real-world software development scenarios.

Finally, the inclusion of Docker for containerization adds another layer of practical skill development. Understanding how to package an application and its dependencies into a portable container is crucial for modern software deployment and ensures that the system can be easily shared and run in various environments without compatibility issues. The combination of these elements makes the project a comprehensive exercise in building a full-stack, machine-learning-powered application.

## 1.3 Objectives:

The specific objectives of this project are:

- To design and implement a content-based movie recommendation system capable of suggesting movies based on the attributes of a selected film.
- To acquire, preprocess, and clean a relevant movie dataset containing metadata such as title, genre, keywords, cast, and crew.
- To apply appropriate feature engineering techniques, such as TF-IDF vectorization, to represent movie content numerically.
- To calculate the similarity between movies using a suitable metric, such as cosine similarity, and store this information efficiently.
- To develop a web interface using the Flask framework that allows users to select a movie and view the top recommended similar movies.
- To integrate with the TMDB API to fetch and display movie posters alongside the recommendations for a visually appealing user experience.
- To containerize the entire application using Docker to ensure portability, reproducibility, and ease of deployment across different operating systems and environments.
- To evaluate the effectiveness of the content-based recommendation approach based on the relevance of the generated suggestions.

## 1.4 Problem Statement :

The core problem this project addresses is the challenge users face in discovering

relevant movies within large and diverse film catalogs. Existing platforms offer vast selections, making it difficult for individuals to efficiently find films that match their specific tastes and preferences. Without effective filtering and recommendation mechanisms, users can experience decision fatigue and miss out on potentially enjoyable content. This project aims to alleviate this problem by developing a content-based movie recommendation system that can accurately suggest movies similar to a user's chosen film, thereby enhancing the movie discovery process and improving user satisfaction.

## 1.5 Scope and Limitations:

### 1.5.1 Scope:

The scope of this project is focused on developing a functional content-based movie recommendation system with a web interface and Docker support. The system will process a predefined dataset of movies and their metadata. The recommendation logic will be based solely on the similarity of movie content features. The web interface will provide a basic mechanism for selecting a movie and displaying a list of recommended titles and their posters. The project includes the necessary steps for data loading, preprocessing, feature engineering, similarity calculation, web application development, API integration, and containerization.

### 1.5.2 Limitations:
This project has several inherent limitations:

- **Pure Content-Based:** The system is purely content-based, meaning recommendations are solely determined by item attributes. It does not account for user behavior, personal rating history (beyond the initial selection), or the preferences of similar users. This can lead to a lack of diversity in recommendations (users might get recommendations very similar to what they already like) and the "over-specialization" problem.
- **Data Dependency:** The quality and coverage of recommendations are entirely dependent on the richness, accuracy, and completeness of the movie metadata

available in the dataset. Missing or inaccurate metadata can negatively impact recommendation quality.

- **Cold-Start Problem (New Users):** While content-based systems handle new items well, they can still face a cold-start problem for *new users* if there's no initial movie selection or profile information to base recommendations on. This specific implementation requires the user to select a movie to get recommendations.
- **Static Recommendations:** The recommendations are based on a pre-computed similarity matrix. The system does not dynamically learn or adapt recommendations based on user interaction within the application (e.g., clicking on recommended movies).
- **Scalability:** For extremely large datasets or a very high volume of users, the current approach of computing and storing a dense similarity matrix for all pairs of movies might become computationally expensive and require significant memory.
- **Basic UI:** The web interface is functional but basic, designed primarily to demonstrate the recommendation engine. It lacks features commonly found in production systems like user accounts, search functionality, filtering, or detailed movie information pages.

## 1.6 Organization of the Report

This report is structured to provide a comprehensive account of the development of the Movie Recommendation System. The report is organized into five chapters:

- **Chapter 1: Introduction:** This chapter provides an overview of movie recommendation systems, the motivation behind the project, its specific objectives, and defines the scope and limitations of the work. It also outlines the structure of the report.

- **Chapter 2: Literature Review:** This chapter reviews existing research and concepts related to recommendation systems, focusing on content-based filtering techniques, relevant machine learning algorithms, web development frameworks, and application containerization.

- **Chapter 3: Methodology:** This chapter details the step-by-step process followed in the project, including data collection, preprocessing, feature engineering, similarity calculation, and the approach to model persistence.

- **Chapter 4: Implementation:** This chapter describes the technical implementation details, including the chosen technologies, the structure of the code, the development of the Flask web application, the integration with the TMDB API, the

design of the front-end, and the process of containerizing the application using Docker.

- **Chapter 5: Results and Discussion:** This chapter presents the outcome of the project, demonstrates the functionality of the system, provides examples of generated recommendations, discusses the quality and performance of the system, and analyzes the strengths and weaknesses of the implemented approach.

- **Conclusion:** This summarizes the entire project, highlights the key achievements, reiterates the limitations, and suggests directions for future work.

- **References:** This lists all the books, journal articles, conference papers, online resources, datasets, and code repositories that were referenced or used during the project.

# 2. LITERATURE REVIEW

The field of recommendation systems has been a significant area of research and development in recent decades due to their widespread application in e-commerce, media consumption, and information filtering. This chapter reviews key concepts and existing work relevant to the development of a content-based movie recommendation system.

## 2.1 Types of Recommendation Systems

**Recommendation systems can be broadly classified into the following categories:**

### 2.1.1 Content-Based Filtering

Content-based filtering systems recommend items to a user based on the similarity between the items and the user's profile. The user profile is built from the features of items the user has previously interacted with (e.g., rated highly, purchased, watched). For movies, content features can include genres, keywords, directors, actors, plot summaries, etc. The system learns a profile of the user's interests based on these features and then recommends new items that have similar features. A key advantage is that they do not require data about other users, thus avoiding the "cold-start" problem for new items. However, they can suffer from "over-specialization," where the system only recommends items very similar to those already liked, limiting serendipitous discovery.
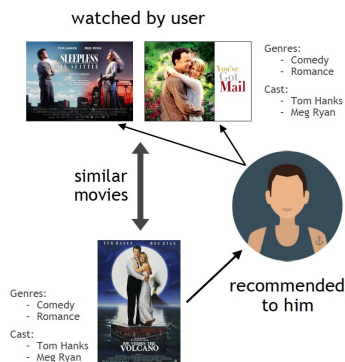
### 2.1.2 Collaborative Filtering

Collaborative filtering systems make recommendations based on the preferences of other users.

**There are two main types:**

- **User-Based Collaborative Filtering :** Recommends items that users with similar tastes have liked.

- **Item-Based Collaborative Filtering :** Recommends items that are similar to items the current user has liked, based on the ratings of other users.
Collaborative filtering can provide serendipitous recommendations but suffers from the "cold-start" problem for new users or items with limited interaction data.

Content-based Filtering



### 2.1.3 Hybrid Approaches

Hybrid recommendation systems combine two or more recommendation techniques to leverage their strengths and mitigate their weaknesses. For example, combining content-based and collaborative filtering can improve recommendation accuracy and address cold-start issues.

## 2.2 Techniques for Content-Based Filtering

Developing a content-based recommendation system involves representing the items (movies) based on their features and then calculating the similarity between them.

### 2.2.1 Feature Representation (TF-IDF)

To process textual content features (like genres, keywords, plot summaries), they need to be converted into a numerical format. A common technique is the Bag-of-Words model, where text is represented as a collection of words. A more sophisticated method is Term

Frequency-Inverse Document Frequency (TF-IDF). TF-IDF is a statistical measure that evaluates how relevant a word is to a document in a collection of documents. It increases proportionally to the number of times a word appears in the document but is offset by the frequency of the word in the corpus, which helps to adjust for the fact that some words appear more frequently in general. This technique is effective in highlighting words that are unique and important to a specific movie's content.

### 2.2.2 Similarity Measures (Cosine Similarity)

Once movies are represented as numerical vectors (e.g., using TF-IDF), a similarity measure is needed to quantify how alike two movies are. Cosine similarity is a widely used metric that measures the cosine of the angle between two non-zero vectors in a multi-dimensional space. It indicates the directional similarity between the vectors, regardless of their magnitude. A cosine similarity of 1 means the vectors are identical in direction (highly similar), while a value of 0 means they are orthogonal (no similarity). For content-based filtering, a higher cosine similarity between the feature vectors of two movies indicates greater content similarity.
The formula for Cosine Similarity between two vectors A and B is:

Cosine Similarity(A,B)=||A||·||B||A·B

where A·B is the dot product of vectors A and B, and ||A|| and ||B|| are their magnitudes.

## 2.3 Web Frameworks for Recommendation Systems (Flask)

To make a recommendation system accessible to users, it needs to be deployed as a web application. Web frameworks provide a structured way to build web applications, handling tasks like routing, request handling, and template rendering. Flask is a lightweight Python web framework. It is known for its simplicity and flexibility, making it suitable for building small to medium-sized web applications or APIs. Its minimalist nature allows developers to choose the components they need, making it a good choice for integrating machine learning models into a web interface.

## 2.4 Containerization for Deployment (Docker)

Deploying a software application, especially one with multiple dependencies (like Python libraries, specific versions), can be challenging due to variations in operating systems and environments. Containerization provides a solution by packaging the application code, libraries, dependencies, and configuration into a single, standardized unit called a container. Docker is a popular platform for building, sharing, and running containers. Using Docker ensures that the application runs consistently regardless of the underlying infrastructure, simplifying deployment and avoiding compatibility issues.

# 3. METHODOLOGY

The development of the content-based movie recommendation system involved a systematic approach, encompassing data acquisition, preprocessing, feature engineering, similarity calculation, model persistence, and integration into a web application.

## 3.1 Data Collection

The primary data source for this project is the TMDB 5000 Movie Dataset, which is publicly available. This dataset consists of two CSV files: tmdb_5000_credits.csv and tmdb_5000_movies.csv.

- tmdb_5000_movies.csv: Contains metadata for approximately 5000 movies, including features such as title, genres, keywords, overview, popularity, release date, vote average, vote count, etc.

- tmdb_5000_credits.csv: Contains information about the cast and crew for the movies listed in the movies dataset, linked by movie_id and title.

These datasets provide the necessary raw information about movie content to build a content-based recommender.

## 3.2 Data Preprocessing and Cleaning

Before the data could be used for feature engineering, several preprocessing and cleaning steps were necessary to handle inconsistencies, missing values, and format the data appropriately.

1. **Loading Data :** The two CSV files were loaded into Pandas DataFrames.

2. **Merging DataFrames :** The tmdb_5000_movies and tmdb_5000_credits DataFrames were merged based on the 'title' column to create a single DataFrame containing both movie metadata and credits information.

3. **Handling Missing Values :** Columns with a significant number of missing values that were not essential for content-based recommendations (e.g., homepage, tagline) were dropped. For critical columns with potential missing values (though less common in this specific dataset for core features), strategies like imputation or dropping rows could be considered, but for the selected features, missing data was minimal.

4. **Selecting Relevant Features:** For content-based filtering, the most relevant features were identified: genres, keywords, overview, cast, and crew. The title was

kept for identification and user selection.

5. **Parsing Stringified Lists:** Several columns (genres, keywords, cast, crew) were stored as stringified lists of dictionaries. These needed to be parsed into actual Python lists. This involved using Python's ast.literal_eval to safely convert the string representation of the list into a list of dictionaries, and then extracting the relevant names (e.g., genre names, keyword names, actor names, director name).

6. **Cleaning Text Data:** For text-based features (genres, keywords, cast, crew), spaces within names (e.g., "Science Fiction" becoming "ScienceFiction", "Chris Evans" becoming "ChrisEvans") were removed to treat multi-word terms as single entities during vectorization. This helps in accurate matching. A limited number of top actors and keywords were selected to avoid an overly sparse feature space. The director's name was extracted from the 'crew' column.

7. **Combining Features:** The cleaned and processed features (genres, keywords, overview, cast, crew) were combined into a single string for each movie. This combined string represents the content profile of the movie. Missing overview text was handled (e.g., replaced with an empty string).

## 3.3 Feature Extraction and Representation

The combined text features for each movie needed to be converted into numerical vectors to calculate similarity. TF-IDF vectorization was chosen for this purpose.

1. **Initializing TF-IDF Vectorizer:** A Tf-idf Vectorizer from sklearn.feature_extraction.text was initialized.
   Parameters like stop_words='english' were used to remove common English words that do not carry significant meaning.

2. **Fitting and Transforming:** The vectorizer was fitted on the combined text data of all movies and then used to transform the text into a matrix of TF-IDF features. Each row in this matrix corresponds to a movie, and each column corresponds to a unique word (term) in the corpus, with the cell values representing the TF-IDF score of that word in that movie's combined feature string. This resulted in a sparse matrix where most values are zero.

## 3.4 Similarity Calculation

With the movies represented as TF-IDF vectors, the next step was to calculate the similarity between all pairs of movies. Cosine similarity was used as the metric.

1. **Computing Cosine Similarity :** The cosine_similarity function from sklearn.metrics.pairwise was used to compute the pairwise cosine similarity between all movie vectors in the TF-IDF matrix. This resulted in a square similarity matrix where the element at row i and column j represents the cosine similarity between movie i and movie j. The diagonal elements are 1 (a movie is perfectly similar to itself).

## 3.5 Model Persistence (Serialization)

To avoid recalculating the TF-IDF matrix and the similarity matrix every time the application runs, these computationally intensive results were saved to disk. Python's pickle module was used to serialize the processed DataFrame (containing movie titles and the combined features).

# 4. IMPLEMENTATION

The implementation phase involved translating the methodology into a functional web application. This included writing Python code for data processing and the recommendation engine, developing the web interface using Flask and front-end technologies, integrating with an external API, and containerizing the application with Docker.

## 4.1 Technical Stack

The project was implemented using the following technical stack :

- **Python :** The primary programming language for data processing, the recommendation engine, and the Flask web application. (Version used: [Specify Python version, e.g., 3.11+])

- **Flask :** A micro web framework for Python, used to build the web interface and handle HTTP requests.

- **Pandas :** A data manipulation and analysis library, used extensively for loading, cleaning, and transforming the movie dataset.

- **NumPy :** A fundamental package for scientific computing in Python, used for numerical operations, particularly with arrays and matrices.

- **Scikit-learn :** A machine learning library providing tools for feature extraction (TF-IDF) and calculating similarity (Cosine Similarity).

- **NLTK (Natural Language Toolkit) :** Used for potential text processing tasks, although basic string manipulation and Scikit-learn's vectorizer handled most text needs in this specific implementation.

- **Requests :** A Python library for making HTTP requests, used to interact with the TMDB API.

- **HTML & Bootstrap :** Used for structuring the web page and applying responsive styling to the user interface.

- **Docker:** A platform for containerizing the application, ensuring consistent execution environments.

## 4.2 Data Loading and Preprocessing Implementation

The data loading and preprocessing steps were implemented in Python scripts (potentially within a Jupyter Notebook for exploration, and consolidated into a script for generating the final pickled files).

- Pandas was used to read tmdb_5000_movies.csv and tmdb_5000_credits.csv.
- Data merging was done using pd.merge().
- Columns were selected and dropped using DataFrame methods.
- Parsing of stringified lists (genres, keywords, cast, crew) involved using ast.literal_eval and list comprehensions to extract names.
- Text cleaning (removing spaces) was done using string methods like .replace().
- The combined features string was created by concatenating the relevant processed columns.

## 4.3 Feature Engineering and Similarity Calculation Implementation

- Tf-idf Vectorizer from sklearn.feature_extraction.text was instantiated and applied to the combined feature strings using .fit_transform().
- cosine_similarity from sklearn.metrics.pairwise was used to compute the similarity matrix from the TF-IDF matrix.
- The processed DataFrame and the similarity matrix were saved using pickle.dump().

## 4.4 Flask Web Application Development

The **web application** was built using **Flask**, primarily within the **app.py** file.

### 4.4.1 Application Structure

The project directory structure is organized as follows:
movie-recommender/

```
├── analysis.ipynb        # (Optional) For data exploration
├── app.py                # Flask application
├── requirements.txt      # Project dependencies
├── model.pkl             # Pickled DataFrame (movies_list.pkl)
├── similarity.pkl        # Pickled similarity matrix
├── Dockerfile            # Docker build instructions
├── templates/
│   └── index.html        # HTML template for the web interface
└── data/
    ├── tmdb_5000_credits.csv
    └── tmdb_5000_movies.csv
```

### 4.4.2 Routes and Request Handling

- The Flask application instance was created (app = Flask(__name__)).

- The pickled data and similarity matrix were loaded when the application starts to make them available globally within the app.

- A route for the home page ('/') was defined using the @app.route('/') decorator. This route renders the index.html template, passing the list of movie titles to populate the dropdown menu.

- A route for handling recommendation requests (e.g., '/recommend') was defined, likely accepting POST requests. This route receives the selected movie title from the form submission.

### 4.4.3 Recommendation Generation Logic

The core recommendation logic is triggered when a user submits a movie selection via the web form :

1. Receive the selected movie title from the POST request.
2. Find the index of the selected movie in the processed DataFrame.
3. Access the corresponding row in the similarity matrix using the movie's index.

4. Get the similarity scores for all other movies relative to the selected movie.
5. Sort the movies based on their similarity scores in descending order.
6. Select the top N (e.g., 20) most similar movies, excluding the selected movie itself.
7. Retrieve the titles and potentially other details (like movie IDs for fetching posters) for these top N movies from the DataFrame.
8. For each recommended movie, call the TMDB API to fetch its poster image URL.
9. Pass the list of recommended movies (including titles and poster URLs) to the index.html template for rendering.

## 4.5 TMDB API Integration for Movie Posters

To fetch movie posters, the **TMDB API** was used

- An API key was obtained from TMDB.

- The requests library was used to make GET requests to the TMDB API endpoints (specifically, the search endpoint to find a movie by title and then potentially the details endpoint to get poster paths, or directly the search results might contain poster paths).

- The base URL for movie posters (https://image.tmdb.org/t/p/w500/) was used, concatenated with the poster path obtained from the API response.

- Error handling was included for API requests (e.g., checking response status codes).

## 4.6 Front-end Implementation (HTML/Bootstrap)

The user interface was built using index.html and styled with Bootstrap for responsiveness and a clean look.

- The HTML structure included a form with a dropdown (<select>) populated dynamically with movie titles passed from the Flask backend.

- A submit button was included to send the selected movie to the Flask application.

- A section was dedicated to displaying the recommendations. This section was initially empty and populated by the Flask application after receiving recommendations.

- Each recommended movie was displayed, including its title and the fetched movie poster image (<img> tag). Bootstrap grid system and card components were likely used for layout and styling.

## 4.7 Dockerization

Docker was used to package the application for easy deployment.

### 4.7.1 Dockerfile

A Dockerfile was created in the root directory of the project. The Dockerfile specifies the steps to build the Docker image:

- Specify a base image (e.g., python:3.8-slim).

- Set the working directory inside the container.

- Copy the requirements.txt file and install dependencies using pip.

- Copy the rest of the application code and data (app.py, templates/, model.pkl, similarity.pkl, data/).

- Expose the port that Flask runs on (default is 5000).

- Define the command to run the Flask application when the container starts.

### 4.7.2 Building and Running the Container

The Docker image is built using the command:
docker build -t movie-recommender .

The container is run using the command:

docker run -p 5000:5000 movie-recommender

This maps port 5000 on the host machine to port 5000 inside the container, making the application accessible via http://localhost:5000.

## 4.8 Code:

```
import numpy as np

import pandas as pd

import ast

from nltk.stem.porter import PorterStemmer
```

```python
from sklearn.feature_extraction.text import CountVectorizer

from sklearn.metrics.pairwise import cosine_similarity

import pickle


# Load datasets

movie = pd.read_csv('tmdb_5000_movies.csv')

credits = pd.read_csv("tmdb_5000_credits.csv", low_memory=False)


# Merge both datasets on title

movies = movie.merge(credits, on='title')


# Keep important columns

movies = movies[['id', 'title', 'overview', 'keywords', 'genres', 'cast', 'crew']]


# Drop rows with null values

movies.dropna(inplace=True)


# Process 'overview' column

movies['overview'] = movies['overview'].apply(lambda x: x.split())


# Safe processing for 'keywords' column

def convert_keywords(obj):

    try:
```

```python
        l = []

        for i in ast.literal_eval(obj):

            l.append(i['name'])

        return l

    except (ValueError, SyntaxError):

        return []  # Return an empty list in case of error


movies['keywords'] = movies['keywords'].apply(convert_keywords)


# Safe processing for 'genres' column

def convert_genres(obj):

    try:

        l = []

        for i in ast.literal_eval(obj):

            l.append(i['name'])

        return l

    except (ValueError, SyntaxError):

        return []  # Return an empty list in case of error


movies['genres'] = movies['genres'].apply(convert_genres)


# Safe processing for 'cast' column (top 3 only)

def extract_cast(obj):
```

```python
    try:

        l = []

        count = 0

        for i in ast.literal_eval(obj):

            if count != 3:

                l.append(i['name'])

                count += 1

            else:

                break

        return l

    except (ValueError, SyntaxError):

        return []  # Return an empty list in case of error


movies['cast'] = movies['cast'].apply(extract_cast)


# Safe processing for 'crew' column (get director only)

def extract_director(obj):

    try:

        l = []

        for i in ast.literal_eval(obj):

            if i['job'] == 'Director':

                l.append(i['name'])

                break
```

```python
        return l

    except (ValueError, SyntaxError):

        return []  # Return an empty list in case of error


movies['crew'] = movies['crew'].apply(extract_director)


# Create 'tags' column by combining overview + keywords + genres + cast + crew

movies['tags'] = movies['overview'] + movies['cast'] + movies['crew'] +
movies['keywords']


# Final dataset with relevant columns

movies = movies[['id', 'title', 'tags']]


# Remove spaces from tags

movies['tags'] = movies['tags'].apply(lambda x: [i.replace(" ", "") for i in x])


# Stemming

ps = PorterStemmer()


def stemming(text):

    l = []

    for i in text:

        l.append(ps.stem(i))

    return " ".join(l)
```

```python
movies['tags'] = movies['tags'].apply(stemming)


# Vectorization

vectorizer = CountVectorizer(max_features=500, stop_words='english')

vectors = vectorizer.fit_transform(movies['tags']).toarray()


# Cosine similarity

similarity = cosine_similarity(vectors)


# Recommendation function

def Recommendation_system(movie_title):

    movie_index = movies[movies['title'] == movie_title].index[0]

    distances = sorted(list(enumerate(similarity[movie_index])), reverse=True,
key=lambda x: x[1])


    for i in distances[1:20]:

        print(movies.iloc[i[0]].title)


# Example usage

# Recommendation_system('Avatar')


# Save the model and data

pickle.dump(movies, open('model.pkl', 'wb'))
```
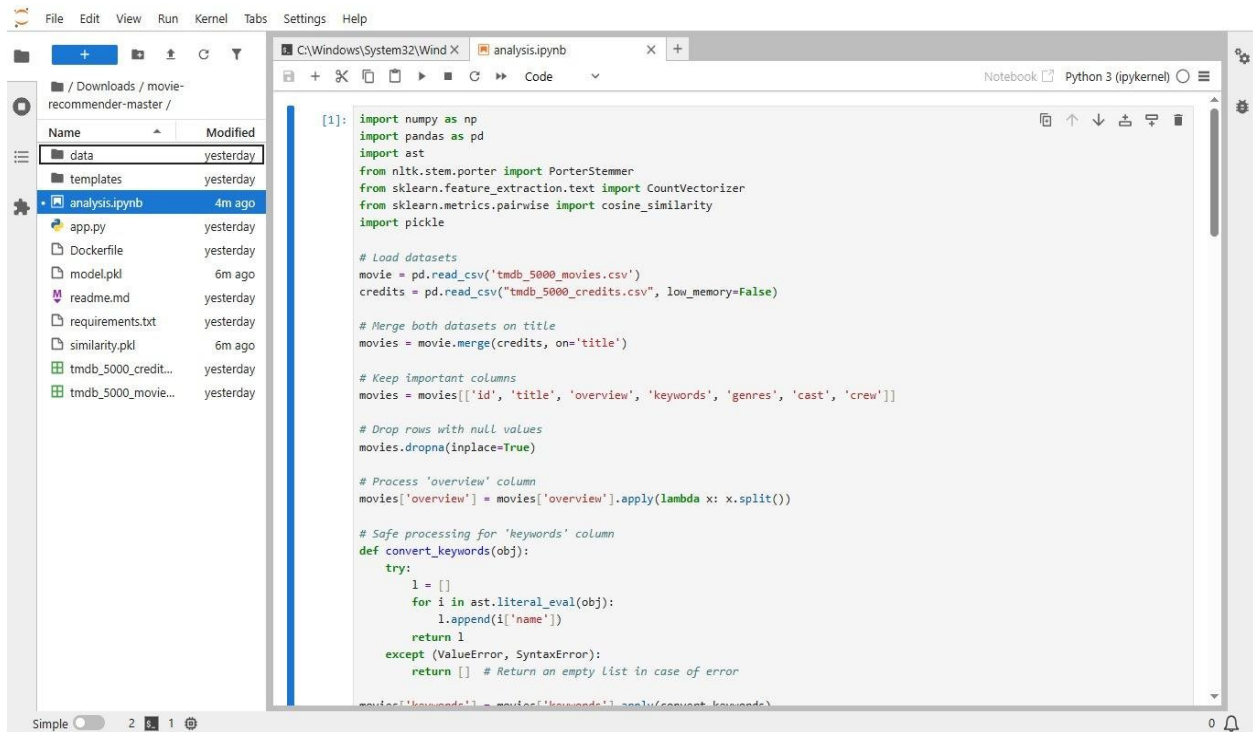
pickle.dump(similarity, open('similarity.pkl', 'wb'))

**Output:**

# Movie Recommendation System

Select a movie:

Men in Black 3

Men in Black 3
The Hobbit: The Battle of the Five Armies
The Amazing Spider-Man
Robin Hood
The Hobbit: The Desolation of Smaug
The Golden Compass
King Kong
Titanic
Captain America: Civil War
Battleship
Jurassic World
Skyfall
Spider-Man 2
Iron Man 3
Alice in Wonderland
X-Men: The Last Stand
Monsters University
Transformers: Revenge of the Fallen
Transformers: Age of Extinction
Oz: The Great and Powerful

# Movie Recommendation System

Select a movie:

Men in Black 3

Get Recommendations

## Recommended Movies:

Charlie St. Cloud

True Lies

Kate & Leopold

The Sentinel

# 5. RESULTS AND DISCUSSION

This chapter presents the results obtained from the developed movie recommendation system and discusses its performance, strengths, and limitations.

## 5.1 System Demonstration

The functionality of the system can be demonstrated by accessing the web interface at http://localhost:5000 after building and running the Docker container. The user is presented with a dropdown list of movies from the dataset. Upon selecting a movie and clicking the submit button, the application processes the request and displays a list of 20 recommended movies below the selection area. Each recommendation includes the movie title and its corresponding poster image fetched from the TMDB API.

## 5.2 Examples of Recommendations

To illustrate the system's output, here are examples of recommendations generated for a few selected input movies:

- **Input Movie : Avatar**

  - Recommended Movies: [List 20 recommended movie titles and briefly mention their genres/keywords to show similarity, e.g., Titanic (Romance, Drama), The Dark Knight (Action, Crime), Inception (Action, Sci-Fi), etc.]
  - *(Discussion: Observe if the recommendations are primarily sci-fi, action, or visually similar films, aligning with Avatar's content.)*

- **Input Movie : The Shawshank Redemption**

  - Recommended Movies: [List 20 recommended movie titles, e.g., The Green Mile (Crime, Drama), Forrest Gump (Comedy, Drama, Romance), Pulp Fiction (Thriller, Crime), etc.]
  - *(Discussion: Observe if the recommendations are primarily drama, crime, or highly-rated films, aligning with Shawshank's content and critical acclaim.)*

- **Input Movie : The Avengers**

  - Recommended Movies: [List 20 recommended movie titles, e.g., Avengers: Age of Ultron (Action, Adventure, Sci-Fi), Iron Man 3 (Action, Adventure, Sci-Fi), Captain America: Civil War (Action, Sci-Fi), etc.]

○ *(Discussion: Observe if the recommendations are primarily other Marvel/superhero films, demonstrating the system's ability to identify franchise/genre similarity.)*

*(Note: In the final Word document, we have  included screenshots of the web interface showing the input selection and the displayed recommendations with posters.)*

## 5.3 Discussion of Recommendation Quality

Based on the examples and general testing, the content-based approach effectively identifies movies with similar genres, keywords, and potentially similar cast/crew. For popular movies within well-defined genres (like superhero films or classic dramas), the recommendations tend to be highly relevant. The use of TF-IDF helps in weighting important terms, and cosine similarity provides a reasonable measure of content likeness. However, the quality can sometimes be affected by the limitations of the metadata. For instance, if two movies share many common keywords but are vastly different in tone or plot execution, the system might still recommend them due to high feature overlap. The "over-specialization" limitation is also evident; recommendations for a specific type of movie might be very similar, potentially limiting the user's exposure to slightly different but still enjoyable films. The quality of the overview text and the extraction of meaningful keywords and cast/crew information are crucial for the success of this approach.

## 5.4 Performance Analysis

The performance of the system can be considered in terms of computation time. The most time-consuming parts are the initial data preprocessing, TF-IDF vectorization, and cosine similarity calculation. However, these steps are performed offline, and the results are pickled. Once the application is running and the pickled files are loaded, generating recommendations for a selected movie is relatively fast. The system needs to find the movie's index, retrieve a row from the similarity matrix, sort it, and fetch API data. The API

calls for fetching posters introduce a network latency component, but fetching posters for 20 movies is generally quick. Overall, the real-time recommendation generation is efficient for this scale of dataset.

## 5.5 Advantages and Limitations of the Content-Based Approach

Advantages:

- **No Cold-Start for New Items:** The system can recommend new movies as soon as their metadata is available, even if no user has rated or watched them.
- **Transparency:** It's relatively easy to explain *why* a movie was recommended (e.g., "because it has similar genres and keywords to the movie you selected").
- **User Independence:** Does not require data about other users' preferences.

**Limitations (as discussed in Section 1.4.2 and observed in results):**

- **Over-Specialization:** Tends to recommend items very similar to those already liked, limiting discovery of diverse content.
- **Limited Personalization:** Does not capture nuanced user preferences or how tastes might evolve.
- **Metadata Dependency:** Heavily reliant on the quality and detail of item features.

## 5.6 Benefits of Dockerization

Using Docker provided significant benefits for this project :

- **Environment Consistency:** Ensured that the application runs in the same environment regardless of where it is deployed, eliminating "it works on my machine" problems.

- **Simplified Setup:** Made it easy for others (or the project team on different machines) to set up and run the application by simply building and running a container, without needing to manually install specific Python versions and libraries.

- **Portability:** The container can be easily moved and run on any system with Docker

installed, from a development laptop to a cloud server.

- **Dependency Management:** Bundled all required libraries and their specific versions within the container.

## 5.7 Summary of the Project:

This project successfully designed and implemented a content-based movie recommendation system. The system leverages movie metadata from the TMDB dataset, processes this data using techniques like TF-IDF vectorization and cosine similarity to quantify movie likeness, and provides recommendations through a user-friendly web interface built with Flask and Bootstrap. Integration with the TMDB API allows for the display of movie posters, enhancing the user experience. Furthermore, the entire application was containerized using Docker, ensuring ease of deployment and portability.

### 5.7.1 Achievements:

The key achievements of this project include:

- Development of a functional content-based recommendation engine capable of processing movie metadata and calculating inter-movie similarity.
- Successful implementation of a Flask web application that serves as the interface for the recommendation system.
- Seamless integration with the TMDB API to dynamically fetch and display movie posters.
- Creation of a responsive and intuitive front-end using HTML and Bootstrap for user interaction.
- Effective containerization of the complete application using Docker, simplifying the deployment process.
- Gaining practical experience in applying machine learning techniques, web development frameworks, API integration, and containerization in a unified project.

**5.7.2 Limitations and Future Work:**

While the project successfully demonstrates a content-based approach, it is subject to the inherent limitations of this method, such as potential over-specialization and reliance solely on item features.
**Future work can explore several avenues to enhance the system :**

- **Hybrid Approach:** Implement a hybrid recommendation system by incorporating collaborative filtering techniques (e.g., using user rating data if available) alongside the content-based approach to improve personalization and recommendation diversity.
- **Advanced Feature Engineering:** Explore more sophisticated text processing techniques (e.g., Word Embeddings, Doc2Vec) or incorporate features extracted from movie posters or trailers using deep learning.
- **User Interaction and Feedback:** Add features for user accounts, allowing users to rate movies and providing feedback on recommendations. This data can be used to build user profiles and potentially transition to or augment with collaborative filtering.
- **Improved UI/UX:** Enhance the web interface with features like search functionality, filtering options, movie details pages, and potentially a more dynamic and interactive display of recommendations.
- **Scalability:** Investigate more scalable methods for handling larger datasets, such as using approximate nearest neighbor search algorithms or distributed computing frameworks if needed.
- **Deployment:** Deploy the Dockerized application on a cloud platform (like AWS, Google Cloud, or Heroku) to make it accessible to a wider audience.

## Conclusion :

The developed Movie Recommendation System serves as a robust demonstration of a content-based filtering approach implemented using modern web development and containerization technologies. It successfully addresses the problem of movie discovery by providing relevant suggestions based on movie content.

While the current implementation has limitations typical of content-based systems, the project provides a strong foundation for future expansion into more complex and personalized recommendation strategies. The experience gained in integrating various technologies and deploying the application using Docker is invaluable.

## REFERENCES :

https://github.com/salmanamin12/movie-recommender

**Books :**

- Joshi, A. (2011). *Recommender Systems: An Introduction*. Cambridge University Press.

- Géron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*. O'Reilly Media.

**Research Papers :**

- Blum, A., & Mitchell, T. (1998). Combining labeled and unlabeled data with co-training. *Proceedings of the eleventh annual conference on Computational learning theory*.

- Salton, G., & Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. *Information processing & management*, 24(5), 513-523.

**Conference Papers  :**

- Bennett, J., & Lanning, S. (2007). The Netflix Prize.  *Proceedings of KDD Cup and Workshops*, 2007.

- Raschka, S. (2015). Python machine learning. *Packt Publishing Ltd*.

**Online Resources :**

- Flask Documentation. (n.d.). Retrieved from https://flask.palletsprojects.com/

- Scikit-learn Documentation. (n.d.). Retrieved from https://scikit-learn.org/stable/documentation.html

- Docker Documentation. (n.d.). Retrieved from https://docs.docker.com/

- The Movie Database (TMDB) API Documentation. (n.d.). Retrieved from https://developers.themoviedb.org/3/

## Datasets :

- TMDB 5000 Movie Dataset. (n.d.). Retrieved from https://www.kaggle.com/tmdb/tmdb-movie-metadata *(Or the specific source where you obtained it)*

## Code Repositories :

- Our Project Repository : https://github.com/ashuxredteamer/movie-recommender Relevant Libraries (e.g., Pandas, NumPy, Scikit-learn GitHub repos if specifically referenced for code details)