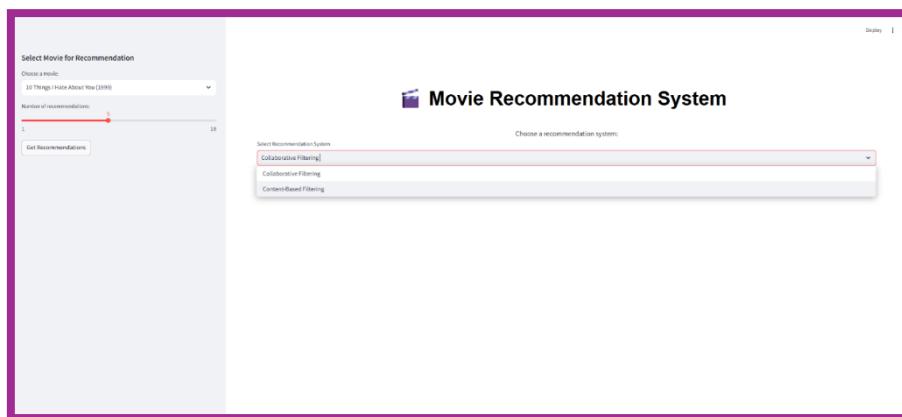


Movie Recommendation System

Abstract

This project presents an innovative approach to movie recommendation using a combination of Machine Learning (ML) techniques, specifically Collaborative Filtering and Content-Based Filtering, with a focus on enhancing user experience through personalized recommendations. The project leverages datasets obtained from Kaggle and incorporates cutting-edge algorithms such as K-Nearest Neighbors (KNN) and TF-IDF Vectorization to suggest movies based on user preferences and content similarity. Additionally, the project employs Streamlit for the development of a user-friendly web application, enabling users to interactively receive movie suggestions. Power BI is utilized for insightful data visualization, providing an in-depth understanding of movie trends, user preferences, and other significant dataset metrics. The goal of this project is to deliver an AI-driven recommendation system that can assist users in discovering relevant movies while gaining valuable insights from the movie dataset.



By: NUMAIRA ZAIB

1. Introduction

This project aims to build a comprehensive movie recommendation system using a combination of collaborative and content-based filtering approaches. The system uses Python and Streamlit for the web app interface, integrating various datasets for both recommendation methods. Additionally, Power BI has been used to visualize dataset insights in detail.

2. Objective

The objective of this project is to:

1. Build a recommendation system to suggest movies based on user ratings and content similarity.
2. Develop a user-friendly application using Streamlit for displaying movie recommendations.
3. Visualize and analyze movie data using Power BI to gain insights into trends, genres, and user preferences.

3. Datasets Used

Collected data using www.kaggle.com

Collaborative Filtering Datasets

1. Ratings Dataset (ratings.csv):

- o Columns: **userId, movieId, rating, timestamp**
- o Contains user ratings for various movies, enabling Collaborative Filtering by identifying similar users or movies based on ratings.

2. Movies Dataset (movies.csv):

- o Columns: **movieId, title, genres**
- o Provides details about each movie's title and genre.

Content-Based Filtering Datasets

1. TMDB Movies Dataset (tmdb_5000_movies.csv):

- o Columns: **budget, genres, keywords, overview, popularity, vote_average**, etc.
- o Includes essential movie features for content-based analysis, such as genre, popularity, budget, and description.

4. Models and Algorithms Used

4.1 Collaborative Filtering

Collaborative Filtering provides recommendations based on user interaction data. In this project, we implemented K-Nearest Neighbors (KNN) with cosine similarity for rating-based recommendations.

Algorithm: K-Nearest Neighbors (KNN)

- Purpose: KNN identifies similar movies or users based on past rating data.
- Functionality:
 - A pivot table of users vs. movies is created, where each cell represents a user's rating for a movie.
 - Cosine similarity is used to calculate the similarity between movies based on ratings.
 - For a given movie, KNN finds the top-K movies with similar ratings, and these similar movies are recommended to the user.

Steps in Code for Collaborative Filtering:

1. Data Loading and Preparation:

- The ratings and movie data are merged based on movieId.
- Movies with fewer ratings are filtered out to avoid sparsity issues.

2. Pivot Table Creation:

- A pivot table with userId and title is created.
- This matrix is fed into the KNN model to find movies with similar user interaction patterns.

3. Model Training and Prediction:

- Using KNN with cosine similarity, similar movies are identified for each movie in the dataset.
- This allows recommending top-K movies based on user similarity for Collaborative Filtering.

4.2 Content-Based Filtering

Content-based filtering recommends movies by identifying items with similar characteristics. This approach utilizes metadata such as genre, keywords, cast, and crew.

Algorithm: TF-IDF Vectorization and K-Nearest Neighbors (KNN)

- TF-IDF Vectorization: The overview, keywords, genres, cast, and crew columns are vectorized to create a content feature set.
 - TF-IDF (Term Frequency-Inverse Document Frequency) assigns weights to words based on their frequency, making less common but relevant words more impactful.
- K-Nearest Neighbors (KNN) with Cosine Similarity:
 - KNN identifies movies with the closest vectorized content, suggesting those with similar characteristics.

Steps in Code for Content-Based Filtering:

1. Data Preprocessing:

- The genres, keywords, cast, and crew columns are converted to text.
- A combined column metadata is created by concatenating these columns to represent the movie content.

2. TF-IDF Vectorization:

- TF-IDF is applied to metadata, creating a sparse matrix of feature weights.
- This matrix serves as the input for KNN to calculate similarity scores.

3. Similarity Calculation and Recommendation:

- The KNN model finds the top-K most similar movies based on cosine similarity.
- Given a movie, the model recommends other movies with similar content profiles.

5. Streamlit Application

The Streamlit app provides an interactive interface for users to obtain recommendations based on collaborative or content-based filtering. Key components of the Streamlit app include:

5.1 Features and User Interface

1. Collaborative Filtering Section:

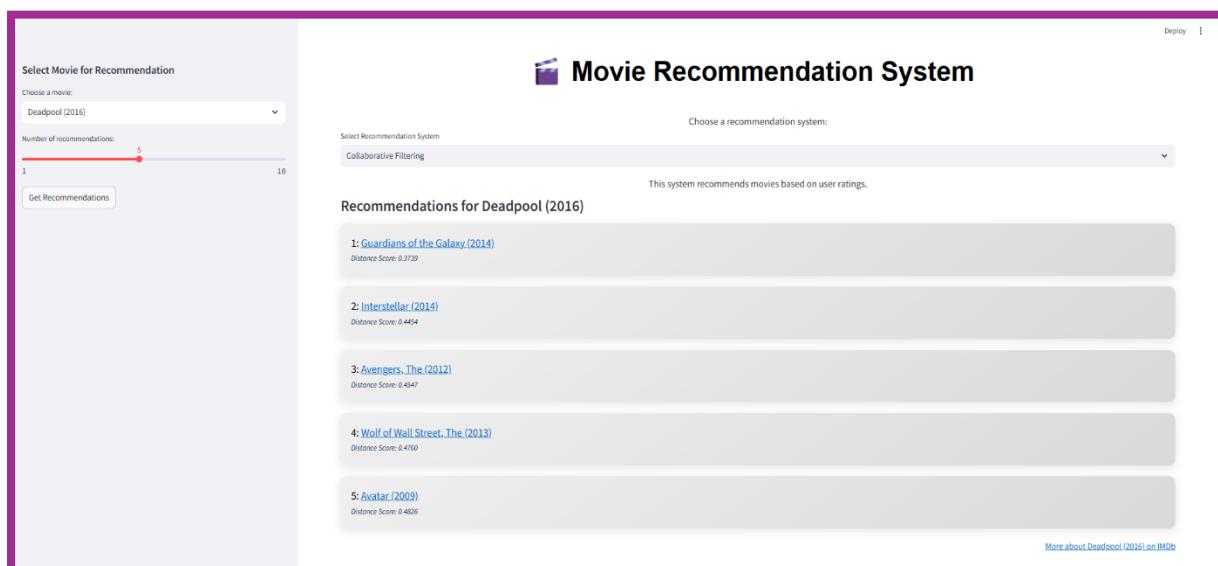
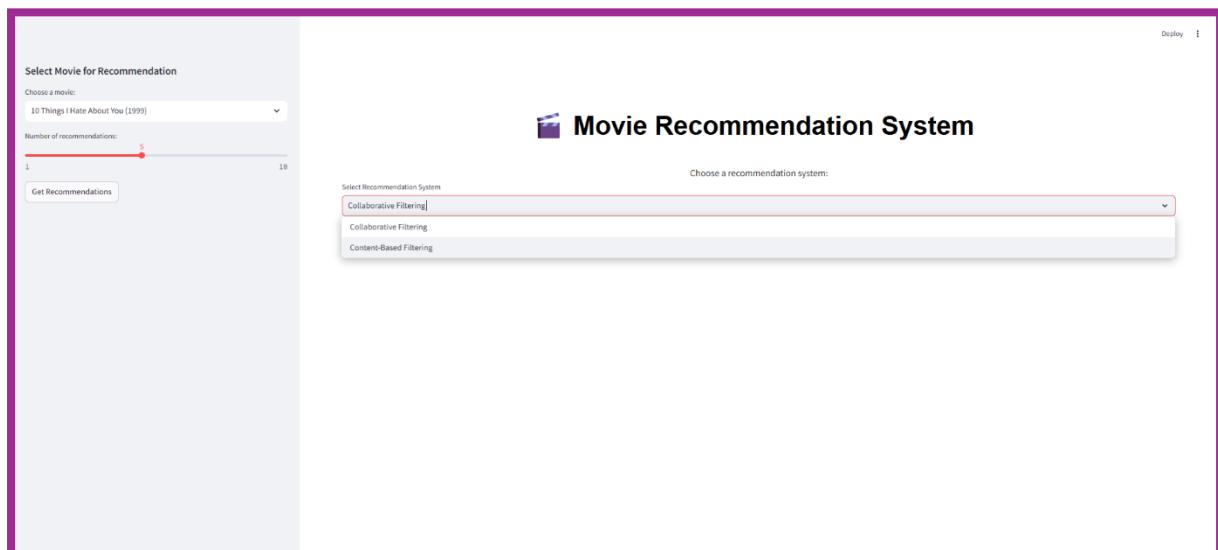
- Allows users to select a movie and receive recommendations based on user rating patterns.
- An input slider lets users specify the number of recommendations.

2. Content-Based Filtering Section:

- Allows users to select a movie and receive recommendations based on similar metadata features.
- Displays detailed information about each recommended movie.

3. Custom Styling:

- Custom CSS styling is used to enhance the visual layout and create hover effects for interactivity.



The top screenshot shows a dropdown menu titled "Select Movie for Recommendation" with the following options:

- Choose a movie:
- Avatar
- Pirates of the Caribbean: At World's End
- Spectre
- The Dark Knight Rises
- John Carter
- Spider-Man 3
- Tangled
- Avengers: Age of Ultron

The bottom screenshot shows the same interface after a recommendation has been made for "Pirates of the Caribbean: At World's End". The dropdown menu now shows "Pirates of the Caribbean: At World's End". The "Number of recommendations" slider is set to 5. The "Get Recommendations" button is visible. The "Movie Recommendation System" title and logo are at the top right. The "Select Recommendation System" dropdown is set to "Content-Based Filtering". A descriptive text below states: "This system recommends movies based on content (genres, keywords, overview, cast, crew)." The "Recommendations for Pirates of the Caribbean: At World's End" section lists the following:

- 1: [Pirates of the Caribbean: Dead Man's Chest](#)
Distance Score: 0.2128
- 2: [Pirates of the Caribbean: The Curse of the Black Pearl](#)
Distance Score: 0.3147
- 3: [Contact](#)
Distance Score: 0.3489
- 4: [The Fifth Element](#)
Distance Score: 0.3572
- 5: [American Pie 2](#)
Distance Score: 0.3952

A link "More about Pirates of the Caribbean: At World's End on IMDb" is located at the bottom right of the recommendations section.

6. Power BI Dashboard for Data Visualization

Power BI is used to create data visualizations that provide insights into the movie dataset. Key metrics and charts include:

6.1 Measures and Calculations

- AvgRating = AVERAGE(ratings[rating])
Calculates the average rating given by users across all movies
- MaxRating = MAX(Ratings[rating])
Retrieves the highest rating value from the ratings dataset, highlighting the maximum score given by any user
- RatingCount = COUNT(Ratings[rating])
Counts the total number of individual ratings submitted by users for movies

Numaira Zaib

- St.DevRating = STDEV.P(Ratings[rating])
Calculates the population standard deviation of all movie ratings, providing insight into rating variability
- AvgVote = AVERAGE(tmdb_5000_movies[vote_average])
Calculates the average number of votes received across all movies in the tmdb_5000_movies dataset
- PopularityScore= AVERAGE(tmdb_5000_movies[popularity])
Calculates the average popularity score for movies in the dataset, representing user engagement

6.2 Visualizations

1. Treemap of Rating Count by Rating:

- Shows the distribution of rating counts, indicating popular rating ranges.

2. Donut Chart (Budget vs Revenue):

- Compares the average budget against revenue, offering insights into profitability across movies.

3. Clustered Bar Chart (Revenue by Genre):

- Displays the total revenue generated by different genres, highlighting successful genres.

4. Line and Stacked Column Chart (Vote Count by Vote Average):

- Shows user engagement through vote count and average rating.

5. Genre Slicer:

- Enables filtering by genre, allowing users to analyze genre-specific trends.

