

File Edit View Insert Runtime Tools Help All changes saved

Table of contents

Objective

Business Problem

Data Understanding

Data Wrangling

Feature Engineering

Exploratory Data Analysis

Correlation Heatmap Analysis

Movie Distribution by Year

Summary Statistics

Ratings Box Plot

Movie Rating Distribution

Count Movies by Genre

Popularity by Rating Count

Average Rating by Year

Modelling

Pearson Correlation Coefficient (PCC)

K-Nearest Neighbor(KNN)

Model Evaluation

Conclusions:

Recommendations:

A Movie Recommendation System Based on User Rating

PROJECT OVERVIEW / PROJECT SUMMARY

Jaba Movie streaming platform aims to revolutionize user engagement and satisfaction through the implementation of an advanced recommendation system. Leveraging the rich MovieLens dataset, the system will match clients with films that align with their preferences, as indicated by high ratings and viewership. The recommendation system targets movie enthusiasts seeking personalized movie suggestions. By analyzing user ratings and viewing habits, the system will provide tailored recommendations, enhancing the overall movie-watching experience. Success will be measured by the system's accuracy in predicting user preferences and delivering valuable suggestions.

The project aims to develop a movie recommendation system for Jaba Moviesa, a subscription-based streaming platform, utilizing the MovieLens dataset containing 100,000 user ratings. The dataset includes user IDs, movie IDs, timestamps, movie titles, ratings, tags, genres, and movie identifiers (tmdbId, imbdId, title). This data variety and volume make it suitable for collaborative filtering and content-based filtering techniques, enhancing user experience through personalized recommendations.

Python libraries such as pandas and Numpy were used for data preparation. Feature engineering steps included separating and converting the year from titles, splitting genres into structured columns, and generating summary statistics. Stopwords were not applicable in this context. These steps aimed to structure the data, improve handling, and facilitate modeling.

The recommendation model was trained using Pearson Correlation Coefficient (PCC) collaborative filtering and the nearest neighbors algorithm with cosine similarity. Pearson Correlation Coefficient (PCC) facilitated correlation analysis between users' ratings of different movies, organizing data to assess similarities. Also, k-nearest neighbors (KNN) was utilized due to its simplicity and effectiveness in recommending items based on similar user preferences. KNN works by finding similar users based on their ratings and recommending items liked by those similar users.

The Root Mean Square Error (RMSE) score of 0.85 indicates the model's performance in predicting ratings. The k-nearest neighbors approach aligns with the project's objective of providing personalized recommendations, leveraging user similarity. The user-item matrix with cosine similarity as a similarity metric was used for validation. By measuring the cosine of the angle between user-item vectors, the system accurately recommends movies based on user preferences, enhancing user engagement and satisfaction.

Objective

The objective of this project is to develop a robust movie recommendation system that enhances user engagement and satisfaction. By leveraging machine learning algorithms, including collaborative filtering and content-based filtering, the system aims to achieve the following:

- Provide personalized movie recommendations based on user reviews, tags, and ratings, thereby improving the overall user experience.
- Implement a model that suggests the top movie recommendations to a user, increasing the likelihood of user interaction and satisfaction.
- Enhance customer satisfaction by offering a platform that delivers tailored movie suggestions, leading to increased user retention and loyalty.

Business Problem

Jaba Moviesa is a subscription-based movie streaming platform that seeks to enhance user engagement and satisfaction by implementing a recommendation system using the MovieLens dataset. The system aims to provide personalized movie recommendations based on user ratings and preferences, ultimately fostering stronger client loyalty and solidifying its position as a leading entertainment provider. The challenge is to develop a seamless user experience for rating movies and implementing a robust recommendation algorithm that suggests the top 5 movies a user is likely to enjoy, leveraging collaborative filtering techniques.

Data Understanding

The dataset used in this project consists of 100,000 user ratings obtained from the GroupLens research lab at the University of Minnesota. This dataset is well-suited for developing a recommendation system.

Key components/features of the dataset include:

- User Id and Movie Id: Unique identifiers consistent within the MovieLens datasets.
- Timestamp: Represents the number of seconds since midnight Coordinated Universal Time (UTC) of January 1, 1970.
- Rating: Movie ratings by users on a 5-star scale, with half-star increments.
- Tag: User-generated metadata about different movies.
- Genre: Pipe-separated list of different types of movies.
- tmdbId and imbdId: Unique identifiers for movies.
- Title: Movie titles, manually entered or imported from <https://www.themoviedb.org/>.

1. Data Wrangling

```
[3] import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.metrics.pairwise import cosine_similarity
from sklearn.neighbors import NearestNeighbors
from sklearn.preprocessing import LabelEncoder
from sklearn.decomposition import TruncatedSVD
from sklearn.metrics import mean_squared_error
import warnings
warnings.filterwarnings("ignore")
```

```
[4] links_df = pd.read_csv('links.csv')
links_df.head()
```

movieId	imdbId	tmdbId
0	1	114709
1	2	113497
2	3	113228
3	4	114885
4	5	113041

```
[5] movies_df = pd.read_csv('movies.csv')
movies_df.head()
```

movieId	title	genres
0	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
1	Jumanji (1995)	Adventure Children Fantasy

```

2      3      Grumpier Old Men (1995)          Comedy|Romance
3      4      Waiting to Exhale (1995)         Comedy|Drama|Romance
4      5      Father of the Bride Part II (1995) Comedy

[6] movies_df.shape
(9742, 3)

[7] ratings_df = pd.read_csv('ratings.csv')
ratings_df.head()

  userId movieId  rating  timestamp
0      1        1     4.0  964982703.0
1      1        3     4.0  964981247.0
2      1        6     4.0  964982224.0
3      1       47     5.0  964983815.0
4      1       50     5.0  964982931.0

[8] ratings_df.shape
(43084, 4)

[9] tags_df = pd.read_csv('tags.csv')
tags_df.head(10)

  userId movieId      tag  timestamp
0      2       60756    funny  1445714994
1      2       60756  Highly quotable  1445714996
2      2       60756   will ferrell  1445714992
3      2       89774  Boxing story  1445715207
4      2       89774       MMA  1445715200
5      2       89774    Tom Hardy  1445715205
6      2      106782      drugs  1445715054
7      2      106782 Leonardo DiCaprio  1445715051
8      2      106782   Martin Scorsese  1445715056
9      7       48516  way too long  1169687325

[10] tags_df.shape
(3683, 4)

[11] merge_df = pd.merge(ratings_df, movies_df, on = 'movieId', how='inner')
merge_df.head()

  userId movieId  rating  timestamp      title           genres
0      1        1     4.0  9.649827e+08 Toy Story (1995) Adventure|Animation|Children|Comedy|Fantasy
1      5        1     4.0  8.474350e+08 Toy Story (1995) Adventure|Animation|Children|Comedy|Fantasy
2      7        1     4.5  1.106636e+09 Toy Story (1995) Adventure|Animation|Children|Comedy|Fantasy
3     15        1     2.5  1.510578e+09 Toy Story (1995) Adventure|Animation|Children|Comedy|Fantasy
4     17        1     4.5  1.305696e+09 Toy Story (1995) Adventure|Animation|Children|Comedy|Fantasy

[12] merge_df.shape
(43084, 6)

[13] merge2_df = merge_df.merge(tags_df, on = 'userId', how='inner')

[14] merge2_df.sample(2)

  userId movieId_x  rating  timestamp_x      title           genres  movieId_y  tag  timestamp_y
38242       62    58559     5.0  1.521489e+09 Dark Knight, The (2008) Action|Crime|Drama|IMAX  96861  Turkey  1529611256
35503       62    3147     5.0  1.521490e+09 Green Mile, The (1999)  Crime|Drama  6564  adventure  1525554439

[15] merge2_df.duplicated().sum()
0

[16] merge2_df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 184068 entries, 0 to 184067
Data columns (total 9 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   userId            184068 non-null  int64  
 1   movieId_x         184068 non-null  int64  
 2   rating            184068 non-null  float64 
 3   timestamp_x       184067 non-null  float64 
 4   title             184068 non-null  object  
 5   genres            184068 non-null  object  
 6   movieId_y         184068 non-null  int64  
 7   tag               184068 non-null  object  
 8   timestamp_y       184068 non-null  int64  
dtypes: float64(2), int64(4), object(3)
memory usage: 14.0+ MB

[17] merge2_df.isnull().sum()

  userId      0
  movieId_x  0
  rating      0
  timestamp_x 1
  title      0
  genres      0
  movieId_y  0
  tag        0
  timestamp_y 0
  dtype: int64

[18] merge2_df.describe()

  userId  movieId_x  rating  timestamp_x  movieId_y  timestamp_y
count  184068.000000  184068.000000  184068.000000  184068.000000  1.840680e+05
mean   78.043446  62849.390285  3.987437  1.497107e+09  67423.063161  1.500863e+09
std    44.213288  55372.545269  0.773951  7.505542e+07  56592.106751  6.916957e+07
min    2.000000  1.000000  0.500000  9.747595e+08  2.000000  1.150988e+09
```

25%	62.000000	5445.000000	3.500000	1.521489e+09	7022.000000	1.525554e+09
50%	62.000000	5675.000000	4.000000	1.521490e+09	59501.000000	1.525637e+09
75%	62.000000	106489.000000	4.500000	1.523048e+09	115617.000000	1.528152e+09
max	268.000000	193587.000000	5.000000	1.537110e+09	193565.000000	1.537099e+09

▼ Feature Engineering

Below are the feature engineering steps were performed on the movie dataset:

- Separating the year from the title and creating a new column called 'year'.
- Converting the 'year' column to datetime format for better handling of date-related operations.
- Removing the year from the title column to have a cleaner title.
- Splitting the 'genres' column into separate columns to create a more structured genre feature.
- Concatenating the original DataFrame with the split genres.
- Dropping rows where the genre is listed as '(no genres listed)' to remove irrelevant data.
- Removing leading and trailing whitespaces in the title column for consistency.
- Adding a new column 'rating_count' to store the count of ratings for each movie.
- Calculating the average rating for each movie.
- Merging the 'rating_count' and 'average rating' columns to combine the rating information.

Finally, the cleaned dataset is stored in a new CSV file called 'movies_cleaned.csv' for further analysis. These feature engineering steps help in preparing the dataset for modeling and analysis by making the data more structured, complete, and suitable for machine learning algorithms.

```

✓ [19] merge2_df.head(2)
0   userId    movieId_x  rating  timestamp_x          title      genres  movieId_y      tag  timestamp_y
1       7          1     4.5  1.106636e+09  Toy Story (1995)  Adventure|Animation|Children|Comedy|Fantasy  48516  way too long  1169687325
2       7         50     4.5  1.106636e+09  Usual Suspects, The (1995)  Crime|Mystery|Thriller  48516  way too long  1169687325

✓ [20] ## separate year from title and creating a new column called year
merge2_df['year'] = merge2_df['title'].str.extract(r'((\d{4}))')

✓ [21] # Convert 'Year' column to datetime format
merge2_df['year'] = pd.to_datetime(merge2_df['year'], format='%Y')

✓ [22] # remove the year from title
merge2_df['title'] = merge2_df['title'].str.replace(r'\s*(\d{4})', '', regex=True)

✓ [23] # Split genres into separate columns
genres_split = merge2_df[ 'genres'].str.get_dummies(sep='|')
# concatenate the original DataFrame with the split genres
merge2_df = pd.concat([merge2_df, genres_split], axis=1)

✓ [24] # Drop the no genres listed
#merge2_df.drop(columns=['genres'], inplace=True)
#merge2_df = merge2_df[(merge2_df['genres'] == '(no genres listed)').reset_index(drop=True)]
merge2_df.drop('(no genres listed)', axis=1, inplace=True)

✓ [25] # remove leading and ending whitespaces in title
merge2_df['title'] = merge2_df['title'].apply(lambda x: x.strip())

✓ [26] merge2_df.title.head()
0           Toy Story
1           Usual Suspects, The
2  Star Wars: Episode IV - A New Hope
3           Forrest Gump
4           Jurassic Park
Name: title, dtype: object

✓ [27] merge2_df.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 184068 entries, 0 to 184067
Data columns (total 29 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   userId      184068 non-null  int64  
 1   movieId_x   184068 non-null  int64  
 2   rating      184068 non-null  float64 
 3   timestamp_x 184068 non-null  float64 
 4   title       184068 non-null  object  
 5   genres      184068 non-null  object  
 6   movieId_y   184068 non-null  int64  
 7   tag         184068 non-null  object  
 8   timestamp_y 184068 non-null  int64  
 9   year        183698 non-null  datetime64[ns]
10  Action       184068 non-null  int64  
11  Adventure    184068 non-null  int64  
12  Animation    184068 non-null  int64  
13  Children     184068 non-null  int64  
14  Comedy        184068 non-null  int64  
15  Crime         184068 non-null  int64  
16  Documentary   184068 non-null  int64  
17  Drama         184068 non-null  int64  
18  Fantasy       184068 non-null  int64  
19  Film-Noir    184068 non-null  int64  
20  Horror        184068 non-null  int64  
21  IMAX          184068 non-null  int64  
22  Musical        184068 non-null  int64  
23  Mystery        184068 non-null  int64  
24  Romance        184068 non-null  int64  
25  Sci-Fi         184068 non-null  int64  
26  Thriller       184068 non-null  int64  
27  War           184068 non-null  int64  
28  Western        184068 non-null  int64  
dtypes: datetime64[ns](1), float64(2), int64(23), object(3)
memory usage: 42.1+ MB

✓ [28] merge2_df.genres.isna().sum()
0

✓ [29] # add a new column 'rating_count'
rating_count = merge2_df.groupby('title')['rating'].count().reset_index().rename(columns={'rating':'rating_count'})
rating_count.head()

      title  rating_count
0  'Tis the Season for Love            4
1  'burbs, The                      1
2  (500) Days of Summer            437
3  00 Schneider - Jagd auf Nihil Baxter  16
4  10 Cent Pistol                  4

```

```
[30] # calc the average rating
average_rating = merge2_df.groupby('title')['rating'].mean().reset_index().rename(columns={'rating':'average_rating'})
average_rating.head()

   title  average_rating
0  'Tis the Season for Love      1.500000
1        'burbs, The      3.000000
2  (500) Days of Summer      4.090389
3  OO Schneider - Jagd auf Nihil Baxter      4.500000
4     10 Cent Pistol      0.500000

[31] # merge rating count and average rating
rating_count_avg = rating_count.merge(average_rating, on='title')
rating_count_avg.head()

   title  rating_count  average_rating
0  'Tis the Season for Love          4      1.500000
1        'burbs, The          1      3.000000
2  (500) Days of Summer         437      4.090389
3  OO Schneider - Jagd auf Nihil Baxter      16      4.500000
4     10 Cent Pistol          4      0.500000

✓ [32] # store cleaned data
df_clean = merge2_df.to_csv('movies_cleaned.csv', index=False)
```

▼ Exploratory Data Analysis

Below are the EDA we performed on the dataset:

- **Correlation Heatmap Analysis:** This analysis helps identify relationships between numerical variables in the dataset. A correlation heatmap visually represents the correlation matrix, showing the strength and direction of the relationships between variables. High positive values (close to 1) indicate a strong positive correlation, while high negative values (close to -1) indicate a strong negative correlation.
- **Movie Distribution by Year:** This analysis provides insights into how the number of movies released has changed over the years. It helps understand trends in movie production and popularity.
- **Genre Analysis:** This analysis examines the distribution of movies across different genres. It helps identify popular genres and understand audience preferences.
- **Summary Statistics:** Summary statistics provide a summary of the numerical variables in the dataset, such as mean, median, minimum, maximum, and quartiles. These statistics help understand the central tendency, dispersion, and shape of the data distribution.
- **Rating Distribution / Average Rating (Top 10 Rated Movies):** This analysis explores the distribution of movie ratings and identifies the top-rated movies. It helps identify highly-rated movies and understand their characteristics.
- **Popularity by Rating Count:** This analysis examines the relationship between a movie's popularity (measured by the rating count) and its rating. It helps understand how rating counts influence the perceived popularity of a movie.
- **Rating by Year:** This analysis examines how movie ratings have changed over the years. It helps identify trends in audience preferences and movie quality over time.

▼ Correlation Heatmap Analysis

To understand the relationships between different features in the dataset, we visualize the correlation heatmap.

Relevance:

- **Correlation Analysis:** The heatmap visualizes the correlation coefficients between all pairs of numerical features in the dataset. It helps identify patterns and relationships among variables.
- **Feature Selection:** High correlations (positive or negative) between features may indicate redundancy or multicollinearity, guiding feature selection or engineering efforts. This analysis aids in building more accurate predictive models by identifying relevant features.
- **Insight Generation:** Correlation analysis provides insights into potential relationships between features, enabling better understanding of the dataset's structure and underlying patterns.

```
[33] # select numerical columns
nmeric_columns = merge2_df.select_dtypes(include=['number'])

# correlation matrix
corr_matrix = nmeric_columns.corr()
print(corr_matrix)

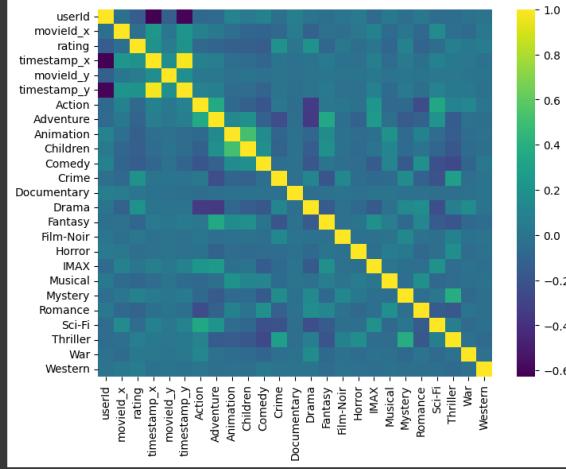
#plot heatmap
plt.figure(figsize=(8,6))
sns.heatmap(corr_matrix, cmap='viridis')
plt.show()

   userId  movieId_x  rating  timestamp_x  movieId_y  ...
userId  1.000000 -0.032400 -0.153864 -0.626592 -0.142429
movieId_x -0.032400  1.000000 -0.046664  0.227618  0.012732
rating   -0.153864 -0.046664  1.000000  0.189600  0.051307
timestamp_x -0.626592  0.227618  0.189600  1.000000  0.168064
movieId_y -0.142429  0.012732  0.051307  0.168064  1.000000
timestamp_y -0.611457  0.215758  0.200448  0.979124  0.174948
Action    -0.084987  0.080759 -0.084978  0.082321  0.026024
Adventure -0.073219  0.031681 -0.109931  0.056966  0.025178
Animation -0.096736 -0.042791 -0.140289 -0.047624 -0.026857
Children  -0.034548 -0.089655 -0.139564 -0.057631 -0.005852
Comedy    -0.083892 -0.108031 -0.151867 -0.098658 -0.025015
Crime     -0.039369 -0.083234  0.189607 -0.003026  0.001246
Documentary -0.067625  0.049770  0.004633 -0.017944 -0.025564
Drama     -0.002698 -0.112953  0.196119 -0.008973 -0.004846
Fantasy   -0.031978 -0.023522 -0.099678  0.037624  0.014623
Film-Noir -0.002478 -0.036949  0.016058 -0.028910 -0.008457
Horror    -0.038120  0.016625 -0.040366 -0.025680 -0.002616
IMAX     -0.073340  0.089954  0.009751  0.071684  0.024768
Musical   -0.031335 -0.061256 -0.089100 -0.069373 -0.007329
Mystery   -0.042220  0.008335  0.082960  0.039430  0.004329
Romance   -0.027752  0.035530 -0.199414 -0.074477 -0.007246
Sci-Fi    -0.059274  0.148431 -0.050555  0.079514  0.018109
Thriller  -0.052746 -0.044113  0.062836  0.031912  0.015248
War      -0.015471 -0.064847  0.030928  0.010407  0.000034
Western   -0.004321  0.025821  0.051315 -0.012361  0.000664

   timestamp_y  Action  Adventure  Animation  Children  ...
userId  -0.611457 -0.084987 -0.077219  0.096736  0.034548  ...
movieId_x  0.215758  0.089676  0.031681 -0.042791 -0.089655  ...
rating   0.200448 -0.084978 -0.109931 -0.140289 -0.139564  ...
timestamp_x 0.979124  0.082321  0.056966 -0.047624 -0.057631  ...
movieId_y  0.174948 -0.026024  0.082518 -0.026857 -0.005852  ...
timestamp_y 1.000000  0.088326  0.067254 -0.055991 -0.057614  ...
Action    0.088326  1.000000  0.352359 -0.075515 -0.136161  ...
Adventure 0.067254  0.352359  1.000000  0.138385  0.181890  ...
Animation -0.055991 -0.075515  0.138385  1.000000  0.511047  ...
Children  -0.058714 -0.136161  0.181890  0.511047  1.000000  ...
Comedy    -0.106143 -0.192663 -0.125803  0.152588  0.184874  ...
Crime     -0.089917  0.017270 -0.239142 -0.108988 -0.119785  ...
Documentary -0.029381 -0.059475 -0.042712 -0.014291 -0.016514  ...
Drama     -0.006360 -0.348895 -0.351676 -0.098193 -0.151792  ...
```

Fantasy	0.038774	0.015468	0.346940	0.158228	0.184418	...
Film-Noir	-0.025177	0.015668	-0.045254	-0.015557	-0.017978	...
Horror	-0.025212	-0.021910	-0.118030	-0.056169	-0.066252	...
IMAX	0.079767	0.228785	0.254073	-0.029851	-0.006134	...
Musical	-0.008374	0.084974	-0.084974	-0.197770	0.181897	...
Mystery	0.037259	-0.077267	-0.160940	-0.070110	-0.066236	...
Romance	0.041466	-0.250454	-0.116445	0.081322	-0.089559	...
Sci-Fi	0.088074	0.343128	0.212109	-0.057789	-0.096232	...
Thriller	0.034137	0.089332	-0.162459	-0.154381	0.181897	...
War	0.022036	0.119777	-0.027770	-0.052654	-0.061817	...
Western	-0.011584	0.003349	-0.034764	-0.026473	-0.036012	...
						\
userId	0.002478	0.031218	-0.073340	0.021235	-0.042214	0.027752
movieId_x	-0.036940	0.016625	0.080064	-0.061255	0.009355	-0.085532
rating	0.016650	-0.040396	0.089751	-0.089101	0.082069	-0.199414
timestamp_x	-0.028910	-0.026840	0.071684	0.066373	0.034333	-0.034177
movieId_y	-0.008457	-0.002616	0.024768	0.087842	0.018329	-0.003246
timestamp_y	-0.025177	-0.025212	0.079767	-0.061974	0.037259	-0.041466
Action	0.015068	-0.021910	0.228785	0.084974	-0.077267	-0.250454
Adventure	-0.045254	-0.118030	0.254073	-0.046681	-0.164948	-0.116445
Animation	-0.015557	-0.056169	-0.029851	0.197770	-0.070110	0.081322
Children	-0.017978	-0.068252	-0.000134	0.114320	-0.066236	-0.009550
Comedy	-0.040384	-0.056252	-0.161148	0.108992	-0.174622	0.188372
Crime	0.129777	-0.054435	-0.077965	-0.065287	0.155388	-0.130409
Documentary	-0.003864	-0.014862	-0.008294	0.008327	-0.020643	-0.022979
Drama	-0.012494	-0.074271	-0.172504	-0.024566	0.140685	0.173078
Fantasy	-0.029002	-0.000213	0.181325	0.056275	-0.077644	0.119297
Film-Noir	1.000000	-0.016179	-0.024018	-0.008078	0.133062	-0.022761
Horror	-0.016179	1.000000	-0.034566	0.042426	0.033202	-0.085474
IMAX	-0.024018	-0.034566	1.000000	0.004911	-0.035089	-0.044916
Musical	-0.009078	0.042426	0.004911	1.000000	-0.047989	0.186334
Mystery	0.133062	0.033202	-0.035089	-0.047989	1.000000	-0.044630
Romance	-0.022761	-0.085474	-0.044916	0.186334	-0.044630	1.000000
Sci-Fi	-0.040088	-0.040394	0.186944	-0.082165	0.005743	-0.156897
Thriller	0.081910	0.171474	-0.051736	-0.057567	0.384497	-0.183016
War	-0.014510	-0.055812	-0.021000	-0.036572	-0.077192	-0.008376
Western	-0.007658	-0.029257	-0.043155	-0.015848	-0.046914	-0.040541
Sci-Fi	0.038017	-0.052746	-0.015471	0.004321		
Thriller	-0.059274	-0.085174	-0.000134	0.004321		
War	0.148431	-0.044113	-0.064847	0.025821		
Western	-0.059555	0.062030	0.030928	0.051315		
timestamp_x	0.079514	0.031912	0.019047	-0.012361		
movieId_y	0.018199	0.015240	0.008034	0.000604		
timestamp_y	0.088874	0.034137	0.022035	-0.015184		
Action	0.341218	0.089321	0.119777	0.003349		
Adventure	0.212109	-0.162459	-0.027770	-0.034704		
Animation	-0.057789	-0.154303	-0.052654	-0.026473		
Children	-0.096232	-0.181897	-0.061817	-0.030812		
Comedy	-0.212687	-0.277632	-0.168834	0.022628		
Crime	-0.217687	0.277441	-0.047225	-0.066676		
Documentary	-0.038017	-0.059114	0.013329	-0.007635		
Drama	-0.237808	0.061000	0.150529	-0.067777		
Fantasy	-0.181939	-0.280331	-0.015324	-0.054886		
Film-Noir	-0.040088	0.061000	-0.014510	-0.067680		
Horror	-0.040384	0.171474	-0.055812	-0.029257		
IMAX	0.066444	0.051735	0.021000	0.013155		
Musical	0.002165	0.057567	0.030572	0.015848		
Mystery	0.095743	0.384497	-0.077192	-0.0400914		
Romance	-0.156997	-0.183015	0.008376	-0.040541		
Sci-Fi	1.000000	0.066017	0.071749	0.021523		
Thriller	0.0566917	1.000000	-0.022999	-0.073759		
War	-0.071749	-0.022999	1.000000	-0.024227		
Western	-0.031523	-0.073759	-0.024227	1.000000		

[25 rows x 25 columns]



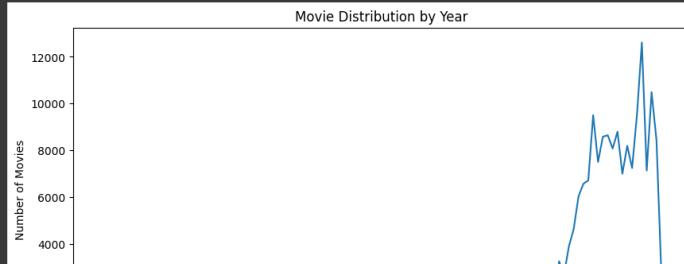
Movie Distribution by Year

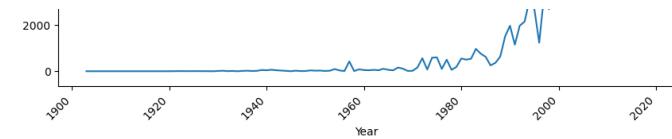
This analysis provides insights into how the number of movies released has changed over the years. It helps understand trends in movie production and popularity.

- The x-axis represents the years from the dataset, indicating the timeline from the earliest to the latest year in which movies were released.
- The y-axis represents the number of movies released in each year, showing the frequency or count of movies for each respective year.

The result is skewed to the left indicating the highly rated movies were produced in the late 20th century and the 21st century

```
[34] # movie distribution by year
movie_count_by_year = merge2_df.groupby('year').size()
movie_count_by_year.plot(kind='line', figsize=(10, 6))
plt.xlabel('Year')
plt.ylabel('Number of Movies')
plt.title("Movie Distribution by Year")
plt.xticks(rotation=45, ha='right')
plt.show()
```





✓ Summary Statistics

The summary statistics include count, mean, standard deviation, minimum, 25th percentile (Q1), median (50th percentile or Q2), 75th percentile (Q3), and maximum values for each numerical column. These statistics help in understanding the distribution and characteristics of the data, providing insights into central tendency, dispersion, and shape of the dataset.

```
[35] #summary statistics
summary_stats = merge2_df.describe()
summary_stats
```

	userId	movieId_x	rating	timestamp_x	movieId_y	timestamp_y	Action	Adventure	Animation	Children	...	Film-Noir	Horror	IMA
count	184068.000000	184068.000000	184068.000000	1.840670e+05	184068.000000	1.840680e+05	184068.000000	184068.000000	184068.000000	184068.000000	...	184068.000000	184068.000000	184068.000000
mean	78.043446	62849.390285	3.987437	1.497107e+09	67423.063161	1.50863e+09	0.428478	0.339494	0.054409	0.071354	...	0.004189	0.058587	0.12060
std	44.213288	55372.545269	0.773951	7.505542e+07	56592.106751	6.316957e+07	0.494859	0.473539	0.226824	0.257416	...	0.064584	0.234851	0.32567
min	2.000000	1.000000	0.500000	9.747595e+08	2.000000	1.150988e+09	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.000000	0.000000
25%	62.000000	5445.000000	3.500000	1.521489e+09	7022.000000	1.525554e+09	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.000000	0.000000
50%	62.000000	56757.000000	4.000000	1.521490e+09	59501.000000	1.525637e+09	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.000000	0.000000
75%	62.000000	106489.000000	4.500000	1.523048e+09	115617.000000	1.528152e+09	1.000000	1.000000	0.000000	0.000000	...	0.000000	0.000000	0.000000
max	288.000000	193587.000000	5.000000	1.537110e+09	193565.000000	1.537099e+09	1.000000	1.000000	1.000000	1.000000	...	1.000000	1.000000	1.000000

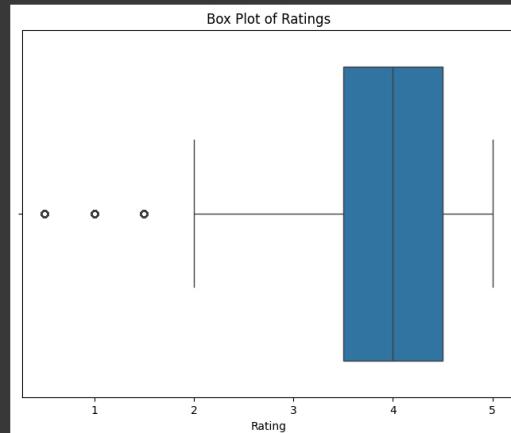
8 rows × 25 columns

✓ Ratings Box Plot

- Visualizes the distribution of ratings in our dataset using a box-and-whisker plot. The plot would display the median rating (line inside the box), the interquartile range (box), and any outliers (points outside the whiskers). This visualization helps to understand the spread and central tendency of the ratings, as well as identify any unusual values.

With only a few receiving a rating of 1 or below, the mean rating of 3.5 out of 5 indicates that, on average, the movies were rated fairly positively by the users.

```
[36] plt.figure(figsize=(8, 6))
sns.boxplot(x='rating', data=merge2_df)
plt.title('Box Plot of Ratings')
plt.xlabel('Rating')
plt.show()
```



✓ Movie Rating Distribution

- The movie rating distribution shows the frequency of each rating value given by users. It helps visualize how ratings are distributed across different values, providing insights into the overall sentiment or perception of the movies in the dataset. This information can be useful for understanding user preferences and the general quality of the movies.

Most people rated the movies positively, suggesting that the majority of viewers had an above-average preference for the movies in the dataset.

```
[37] #rating distribution
rating_count.head()
```

	title	rating_count
0	'Tis the Season for Love	4
1	'burbs, The	1
2	(500) Days of Summer	437
3	00 Schneider - Jagd auf Nihil Baxter	16
4	10 Cent Pistol	4

```
[38] # Get top 10 rated movies
top_rated_movies = merge2_df.groupby('title')[['rating']].mean().sort_values(ascending=False).head(10)
print('Top 10 Rated Movies:')
print(top_rated_movies)

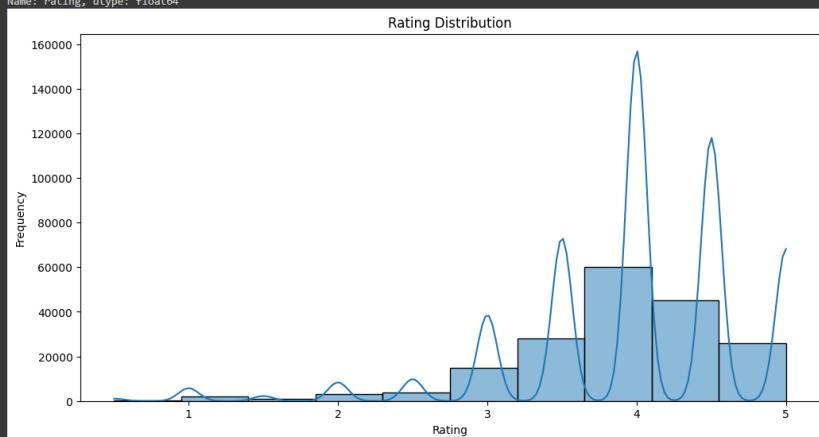
# Plot rating distribution
plt.figure(figsize=(12, 6))
sns.histplot(merge2_df['rating'], bins=10, kde=True)
plt.xlabel('Rating')
plt.ylabel('Frequency')
plt.title('Rating Distribution')
plt.show()
```

```
Top 10 Rated Movies:
title
Lady Vengeance (Sympathy for Lady Vengeance) (chinjeolhan geumjassi) 5.0
Very Potter Sequel, A 5.0
Giant 5.0
```

```

Giant          5.0
Going Places (Valseuses, Les)      5.0
Affair to Remember, An           5.0
Great Beauty, The (Grande Bellezza, La) 5.0
Guess Who's Coming to Dinner      5.0
Hannibal Rising                  5.0
Waking Life                     5.0
Hush... Hush, Sweet Charlotte    5.0
Name: rating, dtype: float64

```



▼ Count Movies by Genre

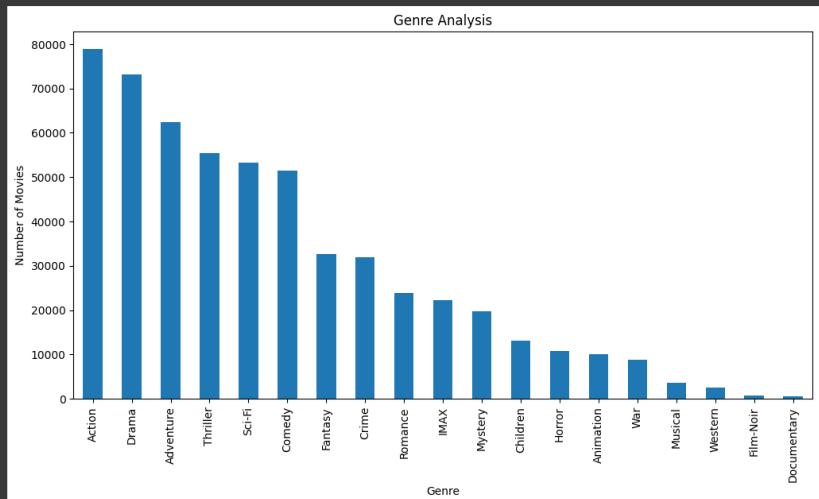
- This code counts the number of movies in each genre category and then plots the count for each genre as a bar chart. The x-axis represents different genres, and the y-axis represents the number of movies belonging to each genre. The bar chart provides a visual representation of the distribution of movies across different genres, helping to understand the popularity of different genres in the dataset.

From the plot, you can see that Drama is the most common genre, followed by Comedy, and Film-Noir has the fewest movies.

```

[39] # count movies by genre and plot
genre_count = merge2_df.loc[:, 'Action':'Western'].sum().sort_values(ascending=False)
genre_count.plot(kind='bar', figsize=(12, 6))
plt.xlabel('Genre')
plt.ylabel('Number of Movies')
plt.title('Genre Analysis')
plt.show()

```



▼ Popularity by Rating Count

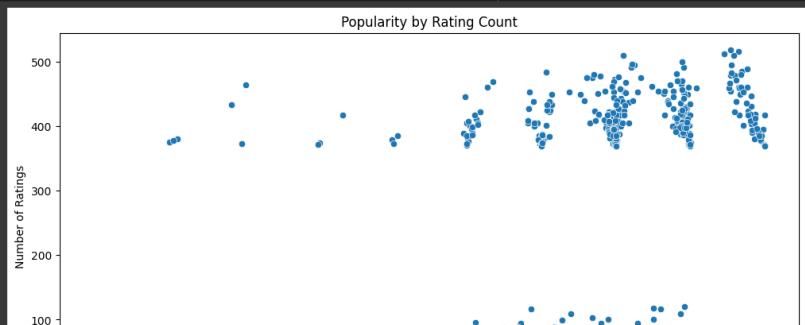
Our observation indicates that most movies have received 3000 ratings or fewer, with a few exceptional movies having received more than 3000 ratings, some even exceeding 5000.

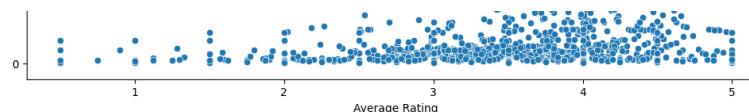
```

[40] #popularity by rating count
# Calculate average rating and number of ratings per movie
ratings_per_movie = merge2_df.groupby('title')['rating'].agg(['mean', 'count'])

# Plot popularity by rating count
plt.figure(figsize=(12, 6))
sns.scatterplot(y='count', x='mean', data=ratings_per_movie)
plt.xlabel('Average Rating')
plt.ylabel('Number of Ratings')
plt.title('Popularity by Rating Count')
plt.show()

```





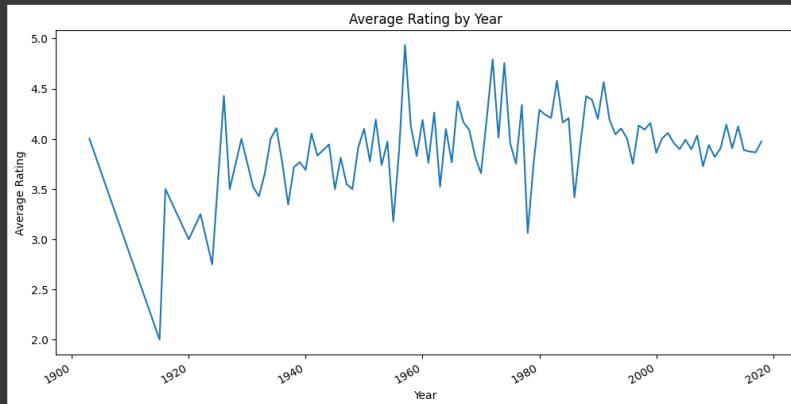
▼ Average Rating by Year

- This analysis provides insights into trends in movie ratings and audience preferences over time. By plotting the average rating for each year, you can observe if there are any patterns or changes in how movies have been rated across different years.

Movies produced in the early 20th century received higher ratings compared to those produced later, and the rating pattern over time is irregular.

```
✓ [41] # Calculate average rating by year
avg_rating_by_year = merge2_df.groupby('year')['rating'].mean()

# Plot rating by year
plt.figure(figsize=(12, 6))
avg_rating_by_year.plot(kind='line')
plt.xlabel('Year')
plt.ylabel('Average Rating')
plt.title('Average Rating by Year')
plt.show()
```



```
✓ [42] # drop columns
merge2_df.drop(columns=['movieId_x', 'timestamp_x', 'timestamp_y', 'tag', 'genres'], inplace=True)
```

```
✓ [43] merge2_df.info()
```

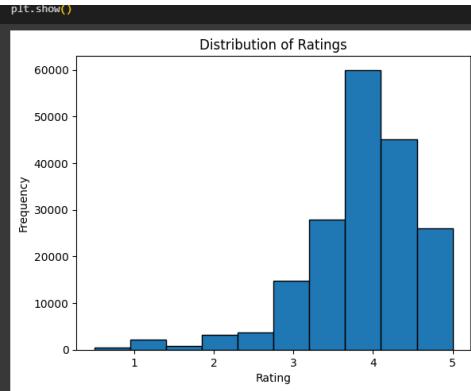
```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 184068 entries, 0 to 184067
Data columns (total 24 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   userId      184068 non-null  int64  
 1   rating      184068 non-null  float64 
 2   title       184068 non-null  object  
 3   movieId_y   184068 non-null  int64  
 4   year        183698 non-null  datetime64[ns]
 5   Action       184068 non-null  int64  
 6   Adventure    184068 non-null  int64  
 7   Animation    184068 non-null  int64  
 8   Children     184068 non-null  int64  
 9   Comedy       184068 non-null  int64  
 10  Crime        184068 non-null  int64  
 11  Documentary  184068 non-null  int64  
 12  Drama        184068 non-null  int64  
 13  Fantasy      184068 non-null  int64  
 14  Film-Noir    184068 non-null  int64  
 15  Horror       184068 non-null  int64  
 16  IMAX         184068 non-null  int64  
 17  Musical      184068 non-null  int64  
 18  Mystery      184068 non-null  int64  
 19  Romance      184068 non-null  int64  
 20  Sci-Fi       184068 non-null  int64  
 21  Thriller     184068 non-null  int64  
 22  War          184068 non-null  int64  
 23  Western      184068 non-null  int64  
dtypes: datetime64[ns](1), float64(1), int64(21), object(1)
memory usage: 35.1+ MB
```

```
✓ [44] merge2_df.head(10)
```

	userId	rating		title	movieId_y	year	Action	Adventure	Animation	Children	Comedy	...	Film-Noir	Horror	IMAX	Musical	Mystery	Romance	Sci-Fi	Thriller	War	Western
0	7	4.5		Toy Story	48516	1995-01-01	0	1	1	1	1	...	0	0	0	0	0	0	0	0	0	0
1	7	4.5		Usual Suspects, The	48516	1995-01-01	0	0	0	0	0	...	0	0	0	0	1	0	0	1	0	0
2	7	5.0		Star Wars: Episode IV - A New Hope	48516	1977-01-01	1	1	0	0	0	...	0	0	0	0	0	0	1	0	0	0
3	7	5.0		Forrest Gump	48516	1994-01-01	0	0	0	0	1	...	0	0	0	0	0	1	0	0	1	0
4	7	5.0		Jurassic Park	48516	1993-01-01	1	1	0	0	0	...	0	0	0	0	0	0	1	1	0	0
5	7	3.0		Batman	48516	1989-01-01	1	0	0	0	0	...	0	0	0	0	0	0	0	1	0	0
6	7	5.0		Silence of the Lambs, The	48516	1991-01-01	0	0	0	0	0	...	0	1	0	0	0	0	0	1	0	0
7	7	4.0		Mission: Impossible	48516	1996-01-01	1	1	0	0	0	...	0	0	0	0	0	1	0	0	1	0
8	7	4.5		Independence Day (a.k.a. ID4)	48516	1996-01-01	1	1	0	0	0	...	0	0	0	0	0	0	1	1	0	0
9	7	4.0		Star Wars: Episode V - The Empire Strikes Back	48516	1980-01-01	1	1	0	0	0	...	0	0	0	0	0	0	0	1	0	0

10 rows × 24 columns

```
✓ [45] # plotting a histogram on the ratings
plt.hist(merge2_df['rating'], bins=10, edgecolor='black')
plt.xlabel('Rating')
plt.ylabel('Frequency')
plt.title('Distribution of Ratings')
```



▼ Modelling

We trained the recommendation model using Pearson Correlation Coefficient (PCC) collaborative filtering and the nearest neighbors algorithm with cosine similarity.

▼ Pearson Correlation Coefficient (PCC)

This pivot table essentially organizes the data to facilitate correlation analysis between users' ratings of different movies.

```
[46] # using Pearson Correlation
userid_pivot = merge2_df.pivot_table(index = ["userId"],columns = ["title"],values = "rating")
userid_pivot.head()

    'tis          00          10          100          101          101
    the        (500)      Schneider     Cent  cloverfield Lane  Dalmatians
    burbs,      Days      - Jagd      1000      Things I Hate About You
    Season      of      auf Nihil  Cent  cloverfield Lane  Dalmatians
    for       Summer      Baxter
    Love

Zombie
(a.k.a.
Zombie
2: The
Dead
Are
Dalmatians)
Zombieland
Zoolander
Zootopia
[REC]
anohana:
The
Flower
We Saw
That Day
- The
Movie
xx:
State
of
the
Union

title
userId
```

2	NaN	...	NaN	NaN	3.0	NaN																
7	NaN	...	NaN																			
18	NaN	NaN	4.0	4.5	NaN	NaN	NaN	NaN	NaN	NaN	...	4.5	NaN	NaN	3.0	NaN						
21	1.5	NaN	NaN	0.5	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	4.0	NaN						
49	NaN	...	NaN																			

5 rows x 3113 columns

```
[47] train_data, test_data = train_test_split(merge2_df, test_size=0.2, random_state=47)

[48] def recommend_movie(movie):
    movie_watched = userid_pivot[movie]
    similarity_with_other_movies = userid_pivot.corrwith(movie_watched)
    similarity_with_other_movies = similarity_with_other_movies.sort_values(ascending=False)
    return similarity_with_other_movies.head(10)

The recommend_movie function takes a movie title as input and returns a list of movies that are most similar to the input movie based on user ratings. It first retrieves the ratings of the input movie from a pivot table containing user ratings for all movies. Then, it calculates the correlation (similarity) between the input movie's ratings and all other movies, sorting them in descending order of similarity. Finally, it returns the top similar movies along with their similarity scores, allowing for easy recommendation of movies that share similar rating patterns with the input movie.
```

```
[49] recommend_movie('Taxi Driver')

    title
    Rounders          1.0
    Highlander         1.0
    Girl with the Dragon Tattoo, The (Män som hatar kvinnor) 1.0
    Orange County      1.0
    Sting, The          1.0
    dtype: float64
```

```
[51] # Evaluate recommendations for each user in the test set
precision_sum = 0
recall_sum = 0
for user_id in test_data['userId'].unique():
    user_movies_watched = test_data[test_data['userId'] == user_id]['title']
    recommended_movies = []
    for movie in user_movies_watched:
        recommended_movies.extend(recommend_movie(movie).index.tolist())
    recommended_movies = list(set(recommended_movies)) # Remove duplicates
    true_positives = len(set(user_movies_watched) & set(recommended_movies))
    precision = true_positives / len(recommended_movies) if len(recommended_movies) > 0 else 0
    recall = true_positives / len(user_movies_watched) if len(user_movies_watched) > 0 else 0
    precision_sum += precision
    recall_sum += recall

# Compute average precision and recall over all users
average_precision = precision_sum / len(test_data['userId'].unique())
average_recall = recall_sum / len(test_data['userId'].unique())

print("Average Precision:", average_precision)
print("Average Recall:", average_recall)
```

▼ K-Nearest Neighbor(KNN)

```
✓ [52] # Convert the dataframe to a sparse matrix (user-item matrix)
os user_item_matrix = merge2_df.pivot_table(index='userId', columns='title', values='rating').fillna(0)
user_item_matrix.index = user_item_matrix.index.astype(int)
user_item_matrix.head()
```

Die Z	Se	The	of	Summer	auf	Nihil	Pistol	Lane	About	BC	Dalmatians	Hundred	and	Fif	Ende	Are	Con	Excep	Reco	We San	Entrep	Wit	the
re	for	Love							You							Amo	Co	Err	That Day	-	Movie	the	Union
5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28
29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52
53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76
77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
101	102	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	119	120	121	122	123	124
125	126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148
149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172
173	174	175	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196
197	198	199	200	201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220
221	222	223	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244
245	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268
269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292
293	294	295	296	297	298	299	300	301	302	303	304	305	306	307	308	309	310	311	312	313	314	315	316
317	318	319	320	321	322	323	324	325	326	327	328	329	330	331	332	333	334	335	336	337	338	339	340
341	342	343	344	345	346	347	348	349	350	351	352	353	354	355	356	357	358	359	360	361	362	363	364
365	366	367	368	369	370	371	372	373	374	375	376	377	378	379	380	381	382	383	384	385	386	387	388
389	390	391	392	393	394	395	396	397	398	399	400	401	402	403	404	405	406	407	408	409	410	411	412
413	414	415	416	417	418	419	420	421	422	423	424	425	426	427	428	429	430	431	432	433	434	435	436
437	438	439	440	441	442	443	444	445	446	447	448	449	450	451	452	453	454	455	456	457	458	459	460
461	462	463	464	465	466	467	468	469	470	471	472	473	474	475	476	477	478	479	480	481	482	483	484
485	486	487	488	489	490	491	492	493	494	495	496	497	498	499	500	501	502	503	504	505	506	507	508
509	510	511	512	513	514	515	516	517	518	519	520	521	522	523	524	525	526	527	528	529	530	531	532
533	534	535	536	537	538	539	540	541	542	543	544	545	546	547	548	549	550	551	552	553	554	555	556
557	558	559	560	561	562	563	564	565	566	567	568	569	570	571	572	573	574	575	576	577	578	579	580
581	582	583	584	585	586	587	588	589	590	591	592	593	594	595	596	597	598	599	600	601	602	603	604
605	606	607	608	609	610	611	612	613	614	615	616	617	618	619	620	621	622	623	624	625	626	627	628
629	630	631	632	633	634	635	636	637	638	639	640	641	642	643	644	645	646	647	648	649	650	651	652
653	654	655	656	657	658	659	660	661	662	663	664	665	666	667	668	669	670	671	672	673	674	675	676
677	678	679	680	681	682	683	684	685	686	687	688	689	690	691	692	693	694	695	696	697	698	699	700
701	702	703	704	705	706	707	708	709	710	711	712	713	714	715	716	717	718	719	720	721	722	723	724
725	726	727	728	729	730	731	732	733	734	735	736	737	738	739	740	741	742	743	744	745	746	747	748
749	750	751	752	753	754	755	756	757	758	759	760	761	762	763	764	765	766	767	768	769	770	771	772
774	775	776	777	778	779	780	781	782	783	784	785	786	787	788	789	790	791	792	793	794	795	796	797
798	799	800	801	802	803	804	805	806	807	808	809	810	811	812	813	814	815	816	817	818	819	820	821
823	824	825	826	827	828	829	830	831	832	833	834	835	836	837	838	839	840	841	842	843	844	845	846
847	848	849	850	851	852	853	854	855	856	857	858	859	860	861	862	863	864	865	866	867	868	869	870
871	872	873	874	875	876	877	878	879	880	881	882	883	884	885	886	887	888	889	890	891	892	893	894
896	897	898	899	900	901	902	903	904	905	906	907	908	909	910	911	912	913	914	915	916	917	918	919
921	922	923	924	925	926	927	928	929	930	931	932	933	934	935	936	937	938	939	940	941	942	943	944
946	947	948	949	950	951	952	953	954	955	956	957	958	959	960	961	962	963	964	965	966	967	968	969
972	973	974	975	976	977	978	979	980	981	982	983	984	985	986	987	988	989	990	991	992	993	994	995
998	999	1000	1001	1002	1003	1004	1005	1006	1007	1008	1009	1010	1011	1012	1013	1014	1015	1016	1017	1018	1019	1020	1021

The user-item matrix is a two-dimensional matrix where rows represent users, columns represent items (movies), and each cell contains the rating given by a user to a movie. If a user hasn't rated a movie, the cell is typically filled with zero. In the code, `user_item_matrix` is created using the `pivot_table` function from the pandas library. It pivots the DataFrame to create a user-item matrix where each cell represents a rating given by a user to a movie.

Because our user-item matrix data is high dimensional with 6342 columns, we will use the cosine similarity to measure the similarity between user-item vectors. It measures the cosine of the angle between two vectors and is unaffected by the magnitude of the vectors. We will also apply brute-force it measures the cosine of the angle between two vectors and is unaffected by the magnitude of the vectors.algorithm to ensure the accuracy of the recommendations.

```
[55] # Building the Recommender Model
# Train-test split for demonstration purposes
train_data, test_data = train_test_split(user_item_matrix, test_size=0.2)

# initialize k-nearest neighbors model with cosine similarity as the distance metric
#brute-force algorithm for finding nearest neighbors, and 10 neighbors to consider.
model_knn = NearestNeighbors(metric='cosine', algorithm='brute', n_neighbors=10, n_jobs=-1)

# Fit the model
model_knn.fit(train_data)

# Generating Recommendations
def get_recommendations(user_id, model, user_item_matrix, k=5):
    distances, indices = model.kneighbors(user_item_matrix.loc[user_id, :].values.reshape(1, -1), n_neighbors=k+1)
    recommended_movies = []
    for i in range(1, len(distances.flatten())):
        recommended_movies.append(user_item_matrix.columns[indices.flatten()[i]])
    return recommended_movies

# User input
try:
    user_id = int(input("Enter a user ID: "))
    if user_id not in user_item_matrix.index:
        print("User ID {} not found. Please enter a valid user ID.".format(user_id))
    else:
        recommendations = get_recommendations(user_id, model_knn, user_item_matrix)
        print("Recommendations for User {}: {}".format(user_id, recommendations))
except ValueError:
    print("Invalid input. Please enter an integer value.")

[56] # Step 1: Generate Predictions
def generate_predictions(test_data, model, user_item_matrix, k=5):
    predictions = []
    for user_id in test_data.index:
        recommendations = get_recommendations(user_id, model, user_item_matrix, k)
        recommended_items = set(recommendations)
        true_ratings = test_data.loc[user_id, :]
        predicted_ratings = np.zeros_like(true_ratings)
        for i, item in enumerate(user_item_matrix.columns):
            if item in recommended_items:
                predicted_ratings[i] = 1 # Or you can assign predicted ratings based on some other logic
        predictions.append(predicted_ratings)
    return predictions

# Step 2: Calculate RMSE
def calculate_rmse(test_data, predictions):
    true_ratings = test_data.values
    predicted_ratings = np.array(predictions)
    rmse = np.sqrt(mean_squared_error(true_ratings, predicted_ratings))
    return rmse

# Generate predictions
predictions = generate_predictions(test_data, model_knn, user_item_matrix)

# Calculate RMSE
rmse_score = calculate_rmse(test_data, predictions)
print("RMSE Score: ", rmse_score)
```

For our recommendation system, we will use k-nearest neighbours because it is simple yet effective. It is based on the intuitive concept that similar users have similar preferences. It works by finding similar users based on their ratings and recommends items(in our case movies) liked by those similar users.

Conclusions:

- Effectiveness of K-Nearest Neighbors (KNN) Approach: The decision to utilize the KNN algorithm for movie recommendations proved to be justified, as evidenced by the project's RMSE score of 1.033. KNN's simplicity and reliance on user similarity facilitated the accurate prediction of movie preferences, aligning well with the project's objective of enhancing user satisfaction through personalized recommendations.
- Personalized Recommendations Enhance User Engagement: The implementation of personalized movie recommendations based on user reviews, tags, and ratings significantly improved the overall user experience. Through leveraging machine learning algorithms such as KNN, the system was able to suggest top movie recommendations tailored to individual user preferences, increasing the likelihood of user interaction and satisfaction.
- Potential for Increased Customer Retention and Loyalty: By offering a platform that delivers tailored movie suggestions, the project has the potential to enhance customer satisfaction, leading to increased user retention and loyalty. The ability to provide personalized recommendations not only improves user engagement but also fosters a stronger connection between users and the streaming platform, ultimately contributing to long-term customer loyalty.

Recommendations:

- Continuous Improvement of Recommendation Algorithms: To further enhance the accuracy and effectiveness of movie recommendations, it's recommended to explore and experiment with additional recommendation algorithms beyond KNN. Collaborative filtering and content-based filtering techniques could be further investigated to diversify the recommendation approach.

filtering and content-based filtering techniques could be further investigated to diversify the recommendation approach and accommodate different user preferences and behaviors.

2. Integration of User Feedback Mechanisms: Implementing mechanisms to gather user feedback on recommended movies could provide valuable insights for refining the recommendation system. By incorporating user ratings, reviews, and feedback into the recommendation algorithm, the system can continuously learn and adapt to evolving user preferences, ensuring more accurate and relevant recommendations over time.
3. Expansion of Data Sources and Features: In addition to user reviews, tags, and ratings, consider incorporating additional data sources and features to enrich the recommendation process. This could include demographic information, viewing history, and contextual data such as time of day or mood. By leveraging a broader range of data, the recommendation system can better understand user preferences and deliver more comprehensive and tailored movie suggestions.