# IMDb Movie Analysis - 10,000 Movies

Shoshana Farber

DATA 602 - Spring 2023

# Abstract

The world of movies is ever-evolving, with fresh releases hitting the scene week after week. Yet, not every film finds success, and some are deemed "flops" based on viewer reviews. This analysis aims to discover popular movies and to uncover patterns and insights related to movie popularity and ratings. The data, comprising 10,000 movies, includes features such as release year, rating, genre, certification, and critical production details.

Exploratory data analysis (EDA) is employed to understand the distributions of variables such as rating, metascore, and gross revenue for movies, and to determine how to best "wrangle" the data. Further analysis is done to determine factors relating to high movie ratings. Some notable findings include romance movies as the top rated genre, followed closely by mystery, drama, and crime. A model was created using release year, movie length, votes, genre, and certification to predict movie ratings. The model accounted for about 47% of the variability in movie ratings and could be further improved with additional movie information.

Additionally, the project introduces a basic recommender system based on IMDb's weighted rating formula and a content based recommender system using the movie description, genres, and other production information. Using these recommender systems, one could filter the data to come up with the top rated movies for their specifications or they can find similar movies to those they've enjoyed.

# Introduction

The dataset selected for this project encompasses IMDb movies and their associated attributes, including release year, rating, genre, certification, and key production details such as director, stars, runtime, and more. IMDb is an internet movie database featuring information on thousands of movies, TV shows, video games, and more worldwide. From timeless classics to the latest releases, IMDb spans the entire spectrum of entertainment. Registered users can rate movies, helping others discover popular and noteworthy films.

This project aims to see which movies/types of movies are most popular on IMDb and to identify features which may influence film popularity. I will also be attempting to come up with a few basic recommender systems based on the rating and content of each movie.

# Required Libraries

**General:**

```python
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
import re
```

**Modeling:**

```python
import sklearn.linear_model as lm
from sklearn.model_selection import train_test_split as tts
from sklearn.metrics import mean_squared_error
```

**Recommender:**

```python
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import linear_kernel
```

# Dataset and Variables

**Movie Name** - name of the movie

**Year of Release** - year movie was released

**Run Time in minutes** - length of the movie

**Movie Rating** - rating of movie

**Votes** - number of users who rated the movie

**MetaScore** - metascore of the movie (metric of rating based on movie critics)

**Gross** - how much the movie made

**Genre** - genre(s) of movie

**Certification** - content rating of movie

**Director** - director(s) of movie

**Stars** - main actor(s) in movie

**Description** - movie bio

| | Movie Name | Year of Release | Run Time in minutes | Movie Rating | Votes | MetaScore | Gross | Genre | Certification | Director | Stars | Description |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | The Shawshank Redemption | 1994 | 142 | 9.3 | 2804443 | 82.0 | 28340000.0 | ['Drama'] | R | ['Frank Darabont'] | ['Tim Robbins', 'Morgan Freeman', 'Bob Gunton'... | ['Over', 'the', 'course', 'of', 'several', 'ye... |
| 1 | The Godfather | 1972 | 175 | 9.2 | 1954174 | 100.0 | 134970000.0 | ['Crime', ' Drama'] | R | ['Francis Ford Coppola'] | ['Marlon Brando', 'Al Pacino', 'James Caan', '... | ['Don', 'Vito', 'Corleone,', 'head', 'of', 'a'... |
| 2 | Ramayana: The Legend of Prince Rama | 1993 | 135 | 9.2 | 12995 | NaN | NaN | ['Animation', ' Action', ' Adventure'] | PG | ['Ram Mohan', 'Yûgô Sakô', 'Koichi Saski'] | ['Arun Govil', 'Nikhil Kapoor', 'Edie Mirman',... | ['An', 'anime', 'adaptation', 'of', 'the', 'Hi... |
| 3 | The Chaos Class | 1975 | 87 | 9.2 | 42231 | NaN | NaN | ['Comedy', ' Drama'] | NaN | ['Ertem Egilmez'] | ['Kemal Sunal', 'Münir Özkul', 'Halit Akçatepe... | ['Lazy,', 'uneducated', 'students', 'share', '... |
| 4 | The Dark Knight | 2008 | 152 | 9.0 | 2786129 | 84.0 | 534860000.0 | ['Action', ' Crime', ' Drama'] | PG-13 | ['Christopher Nolan'] | ['Christian Bale', 'Heath Ledger', 'Aaron Eckh... | ['When', 'the', 'menace', 'known', 'as', 'the'... |

# DATA WRANGLING

# Column Manipulation

- Change **Run Time in minutes** to **Runtime Mins**

- Clean up and turn **Genres, Directors, Stars, Description** columns into lists

- Join words in **Description** columns list

- Add **Year of Release** into **Movie Name**

- Fill any NaN **Certification** with "Not Rated"

# Column Creation

- **Main Genre** - first genre in list

- **Num Directors** - number of directors

- **Num Stars** - number of main stars

- Dummy variables for **Genre**

- Dummy variables for **Certification**

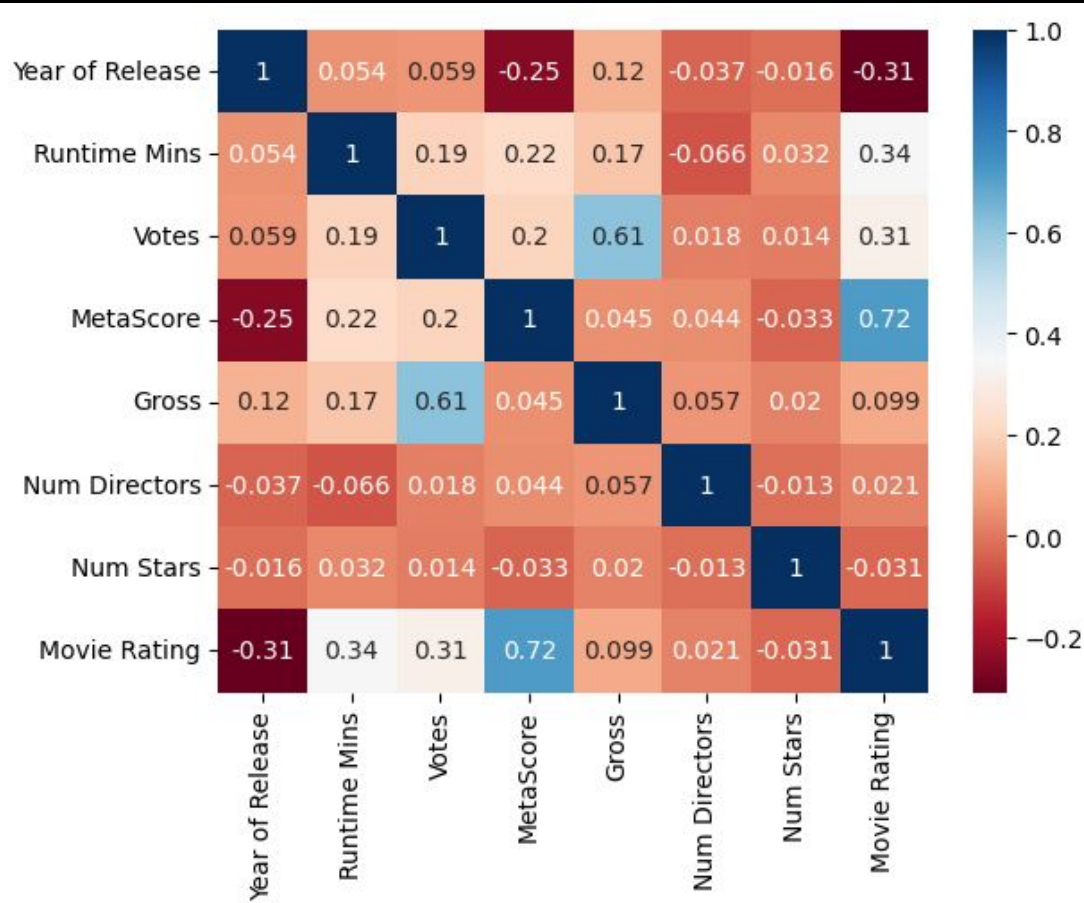| | Movie Name | Year of Release | Runtime Mins | Movie Rating | Votes | MetaScore | Gross | Genre | Certification | Director | Stars | Description | Main Genre | Num Directors | Num Stars |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | The Shawshank Redemption (1994) | 1994 | 142 | 9.3 | 2804443 | 82.0 | 28340000.0 | [DRAMA] | R | [Frank Darabont] | [Tim Robbins, Morgan Freeman, Bob Gunton, W... | over the course of several years two c... | DRAMA | 1 | 4 |
| 1 | The Godfather (1972) | 1972 | 175 | 9.2 | 1954174 | 100.0 | 134970000.0 | [CRIME, DRAMA] | R | [Francis Ford Coppola] | [Marlon Brando, Al Pacino, James Caan, Dian... | don vito corleone head of a mafia fami... | CRIME | 1 | 4 |
| 2 | Ramayana: The Legend of Prince Rama (1993) | 1993 | 135 | 9.2 | 12995 | NaN | NaN | [ANIMATION, ACTION, ADVENTURE] | PG | [Ram Mohan, Yûgô Sakô, Koichi Saski] | [Arun Govil, Nikhil Kapoor, Edie Mirman, Ra... | an anime adaptation of the hindu epic t... | ANIMATION | 3 | 4 |
| 3 | The Chaos Class (1975) | 1975 | 87 | 9.2 | 42231 | NaN | NaN | [COMEDY, DRAMA] | Not Rated | [Ertem Egilmez] | [Kemal Sunal, Münir Özkul, Halit Akçatepe, ... | lazy uneducated students share a very c... | COMEDY | 1 | 4 |
| 4 | The Dark Knight (2008) | 2008 | 152 | 9.0 | 2786129 | 84.0 | 534860000.0 | [ACTION, CRIME, DRAMA] | PG-13 | [Christopher Nolan] | [Christian Bale, Heath Ledger, Aaron Eckhart... | when the menace known as the joker wrea... | ACTION | 1 | 4 |

| Genre | Movie Name | ACTION | ADVENTURE | ANIMATION | BIOGRAPHY | COMEDY | CRIME | DRAMA | FAMILY | FANTASY | ... | HORROR | MUSIC | MUSICAL | MYSTERY | ROMANCE | SCI-FI | SPORT | THRILLER | WAR |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | #Alive (2020) | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | ... | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 1 | #Home (2021) | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 2 | '71 (2014) | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 1.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 3 | '83 (2021) | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 1.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 4 | (T)Raumschiff Surprise - Periode 1 (2004) | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 |

| Certification | Movie Name | 13+ | 16+ | 18+ | Approved | G | GP | M | M/PG | MA-17 | ... | R | TV-13 | TV-14 | TV-G | TV-MA | TV-PG | TV-Y7 | TV-Y7-FV | Unrated | X |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | #Alive (2020) | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 1 | #Home (2021) | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 2 | '71 (2014) | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 3 | '83 (2021) | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 4 | (T)Raumschiff Surprise - Periode 1 (2004) | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |

# Final Data Frame

- Join main data frame with
dataframes of dummy variables

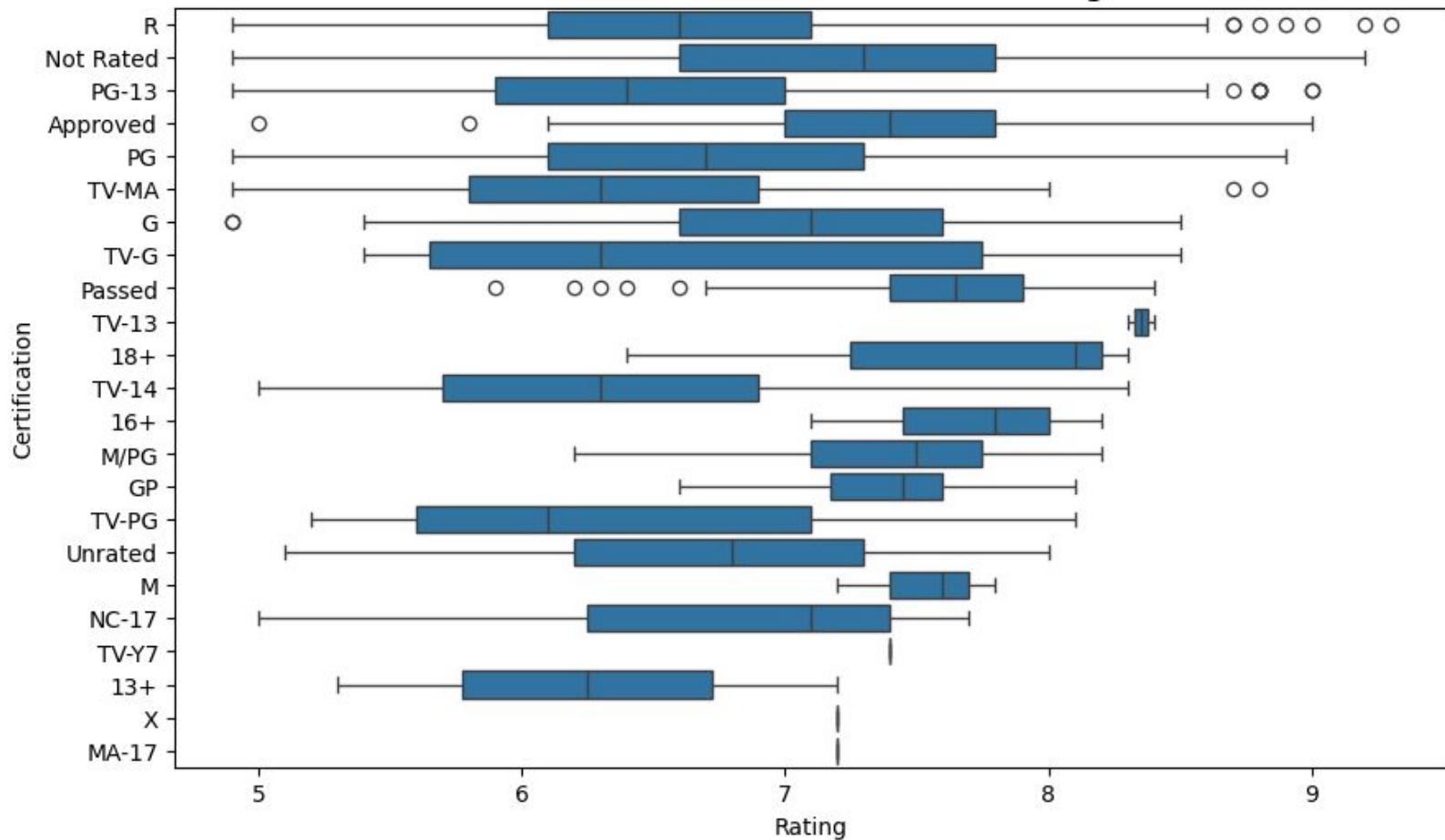| | Movie Name | Year of Release | Runtime Mins | Movie Rating | Votes | MetaScore | Gross | Genre | Certification | Director | ... | R | TV-13 | TV-14 | TV-G | TV-MA | TV-PG | TV-Y7 | TV-Y7-FV | Unrated | X |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | The Shawshank Redemption (1994) | 1994 | 142 | 9.3 | 2804443 | 82.0 | 28340000.0 | [DRAMA] | R | [Frank Darabont] | ... | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 1 | The Godfather (1972) | 1972 | 175 | 9.2 | 1954174 | 100.0 | 134970000.0 | [CRIME, DRAMA] | R | [Francis Ford Coppola] | ... | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 2 | Ramayana: The Legend of Prince Rama (1993) | 1993 | 135 | 9.2 | 12995 | NaN | NaN | [ANIMATION, ACTION, ADVENTURE] | PG | [Ram Mohan, Yûgô Sakô, Koichi Saski] | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 3 | The Chaos Class (1975) | 1975 | 87 | 9.2 | 42231 | NaN | NaN | [COMEDY, DRAMA] | Not Rated | [Ertem Egilmez] | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 4 | The Dark Knight (2008) | 2008 | 152 | 9.0 | 2786129 | 84.0 | 534860000.0 | [ACTION, CRIME, DRAMA] | PG-13 | [Christopher Nolan] | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |

ANALYSIS

# Numerical Correlations

- Large positive correlation between **MetaScore** and **Movie Rating**

- Slight positive correlation between **Runtime Mins** and **Movie Rating**, **Votes** and **Movie Rating**

- Very small positive correlation between **Gross** and **Movie Rating**

- There is a slight negative correlation between **Year of Release** and **Movie Rating**

**Distribution of Genre vs. Rating for Top 10 Genres**

Distribution of Certification vs. Rating

# Modeling

```python
# variables of interest
variables = full[['Movie Rating', 'Year of Release', 'Runtime Mins', 'Votes', 'ACTION', 'ADVENTURE',
'ANIMATION', 'BIOGRAPHY', 'COMEDY', 'CRIME', 'DRAMA', 'FAMILY', 'FANTASY', 'FILM-NOIR', 'HISTORY',
'HORROR', 'MUSIC', 'MUSICAL', 'MYSTERY', 'ROMANCE', 'SCI-FI', 'SPORT', 'THRILLER', 'WAR', 'WESTERN',
'13+', '16+', '18+', 'Approved', 'G', 'GP', 'M', 'M/PG', 'MA-17', 'NC-17', 'Not Rated', 'PG', 'PG-13',
'Passed', 'R', 'TV-13', 'TV-14', 'TV-G', 'TV-MA', 'TV-PG', 'TV-Y7', 'TV-Y7-FV', 'Unrated', 'X']]


# define X, y and split into train:test data
X = variables.drop(columns=['Movie Rating']).values
y = variables['Movie Rating'].values
X_train, X_test, y_train, y_test = tts(X, y, test_size=0.3, random_state=23)


# initiate model and fit data
model = lm.LinearRegression()
model.fit(X_train, y_train)
```

# Modeling

```python
# predict
y_pred = model.predict(X_test)

# check accuracy
r2 = model.score(X_test, y_test)
mse = mean_squared_error(y_test, y_pred)

print(f'R-squared: {r2}')
print(f'Mean Squared Error: {mse}')

>>>
R-squared: 0.4671394895179628
Mean Square Error: 0.3639993450869121
>>>
```

**RECOMMENDER SYSTEM**

# Basic Recommender - Weighted Rating

$$WR = \frac{v}{v+m} * R + \frac{m}{v+m} * C$$

- R = average rating for movie
- v = number of votes for movie
- m = minimum votes to be listed in Top 250 (25,000 votes)
- C = mean vote across dataset

```python
# function for weighted rating
def weighted_rating(df, m=25000):
    R = df['Movie Rating']
    v = df['Votes']
    C = df['Movie Rating'].mean()

    return (v/(v+m) * R) + (m/(v+m) * C)
```

DRAMA:
Twilight (2008)
Showgirls (1995)
Rocky V (1990)
The Circle (2017)
Flatliners (2017)

COMEDY:
Year One (2009)
The Flintstones (1994)
Scary Movie 4 (2006)
Sex Tape (2014)
High School Musical 3: Senior Year (2008)

ACTION:
Wild Wild West (1999)
The Twilight Saga: Eclipse (2010)
Anaconda (1997)
Charlie's Angels: Full Throttle (2003)
Transformers: The Last Knight (2017)

# Content Based Recommender

**Step 1:** Combine information into one column (**Description, Directors, Actors**)

**Step 2:** Vectorize words in new column

```python
tfidf = TfidfVectorizer(stop_words='english')

# make TF-IDF matrix
tfidf_matrix = tfidf.fit_transform(movies_df['Content'])
```

**Step 3:** Use cosine similarity scores to compare similar information

```python
# cosine similarity matrix
cosine_sim = linear_kernel(tfidf_matrix, tfidf_matrix)

# reverse map of indices and movie titles
indices = pd.Series(movies_df.index, index=movies_df['Movie Name'])
```

# Content Based Recommender

**Step 4:** Function to get top 10 recommendations based on similarity scores

```python
# function that takes in movie title as input and outputs most similar movies
def get_recommendations(title, cosine_sim=cosine_sim):
    index = indices[title]  # get index of movie
    sim_scores = list(enumerate(cosine_sim[index]))  # get similarity scores
    sim_scores = sorted(sim_scores, key=lambda x: x[1], reverse=True)   # sort based on similarity scores
    sim_scores = sim_scores[1:11]   # scores of 10 most similar movies
    movie_indices = [i[0] for i in sim_scores]  # movie indices for similar movies
    return movies_df['Movie Name'].iloc[movie_indices]  # return top 10 similar movies
```

**Step 5:** Implement function

```python
get_recommendations('X-Men: Days of Future Past (2014)')
```

```
2460                                    X-Men (2000)
2143                                       X2 (2003)
5008                            The Wolverine (2013)
350                                     Logan (2017)
4180                         X-Men: Apocalypse (2016)
5448                     X-Men: The Last Stand (2006)
6         The Lord of the Rings: The Return of the King ...
5145                                Apt Pupil (1998)
5884                  X-Men Origins: Wolverine (2009)
8649                             Dark Phoenix (2019)
```

# Conclusions

This analysis of IMDb movie data provides insights into some of the factors determining movie popularity and ratings. Through exploratory data analysis, we identified the top genres, with romance movies emerging as the highest rated, closely followed by mystery, drama, and crime. The distribution of ratings for different genres was visualized, revealing interesting patterns and variations. The correlation analysis highlighted factors such as movie length and the number of votes, indicating their influence on higher ratings. A model was created to predict movie rating based on release year, movie length, votes, genre, and certification. The model accounted for 47% of the variation in the data and had a relatively high mean square error but could be improved using other demographic or user information which was not sourced for this project.

A basic recommender system was implemented to suggest movies based on highest weighted rating and a content based recommender system was implemented to suggest movies based on previously watched movies.

In conclusion, this project provides a comprehensive exploration of IMDb movie data, offering valuable insights for movie enthusiasts, analysts, and industry professionals. The findings contribute to a deeper understanding of the dynamics that drive movie ratings and popularity on the IMDb platform.

# Sources

1. Data: https://www.kaggle.com/datasets/dk123891/10000-movies-data?select=data.csv

2. IMDb Weighted Ratings: https://www.quora.com/How-does-IMDbs-rating-system-work

3. Recommender Systems Tutorial: https://www.datacamp.com/tutorial/recommender-systems-python