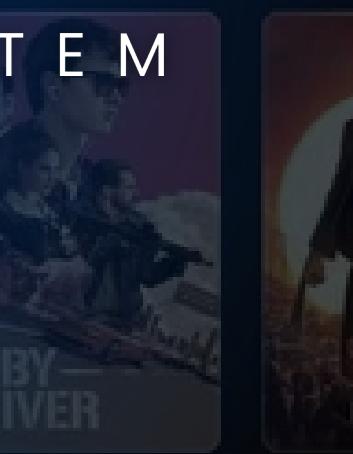
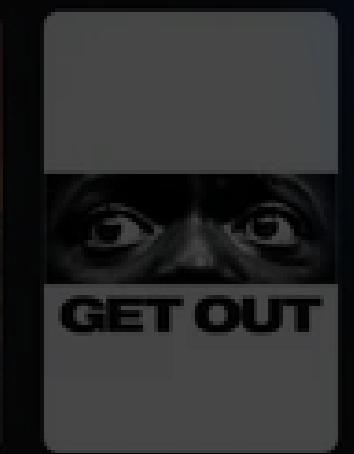
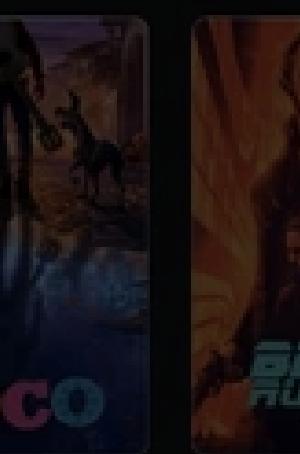
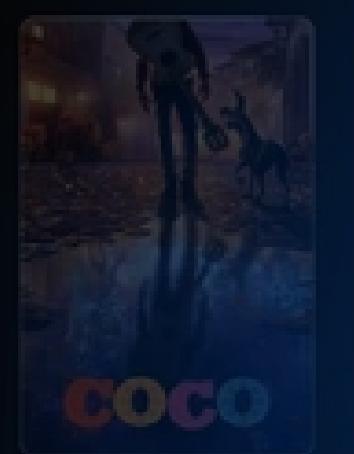
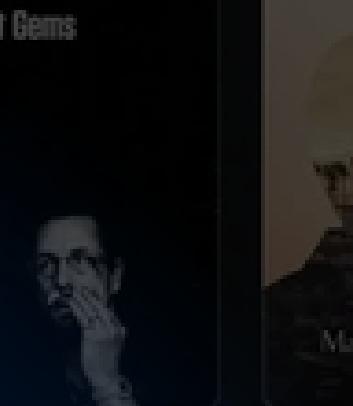
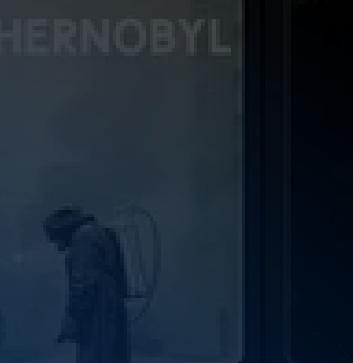


Uncut Gems



CHALUNGSANG Tenzin - KHIN Vannkiriya - LIU Kelly - PHORIMAVONG Gwendoline - TRINH Laurie

# Introduction

The main goal is to suggest similar movies to users by analyzing movie metadata, such as genre, director, actors, and description. This study is based on the Kaggle notebook "Getting Started with a Movie Recommendation System."

To achieve this, we use a dataset movies released up to July 2017, containing detailed information such as cast, production team, plot keywords, budget, revenue, posters, and release dates.

Through this study, we aim to understand how recommendation systems use this information to generate relevant suggestions and apply commonly used techniques in this field.



## Séries

Genres ▾

Top 10 des séries TV en France aujourd'hui



Tendances actuelles



Séries policières dramatiques



Top 10 Movies in the U.S. Today



Netflix Originals



# The Original Notebook



The datasets used in the Kaggle notebook are files containing information about movies, users, and movie ratings. They help train the recommendation model by providing a database from which it can learn to make predictions.

# Datasets Used:

## TMDB 5000 Movies:

TMDB 5000 Movies: Contains movie details (title, genre, duration, popularity).

This helps the model understand what types of movies users prefer.

## TMDB 5000 Credits:

TMDB 5000 Credits: Includes actor and director information.

This helps connect people involved in the movie to its success and can influence recommendations.

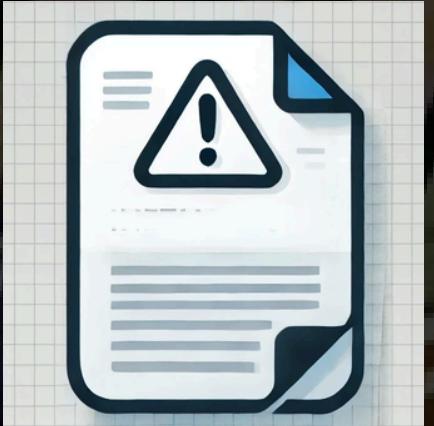
## MovieLens Ratings (Small):

MovieLens Ratings (Small): Stores user ratings for movies.

This information is used to understand user preferences and improve movie recommendations.

These datasets help train the model by combining movie details, cast/crew data, and audience ratings to improve recommendations.

# Problems Encountered on Google Colab



## Initial Issue :

On Kaggle, everything is ready to use.

On Colab, datasets are not recognized automatically, requiring us to manually download and prepare the data from Kaggle.



## Access to Files:

On Kaggle datasets are preloaded and easy to use.

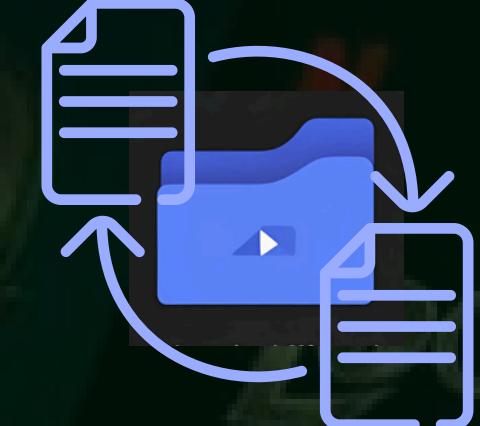
On Colab we used Kaggle API + kaggle.json to download files, then unzipped them.



## Access to Surprise Library:

Scikit-surprise not pre-installed on Colab, so we had to install it manually.

File paths differ between platforms, so we adjusted code to handle both environments.



## Changes in Surprise Library:

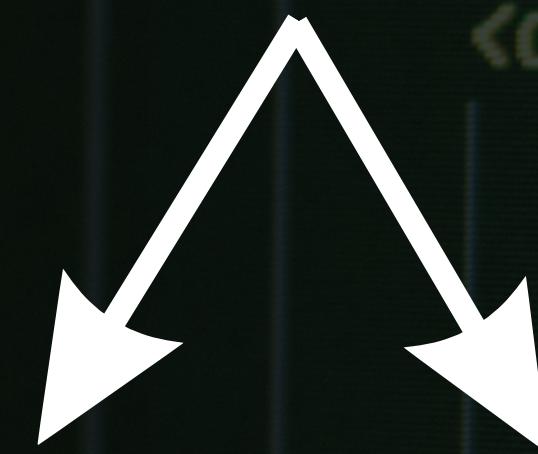
The method **evaluate()** was replaced by **cross\_validate()** for model validation.

The function **split(n\_folds=5)** was removed and replaced by **cross\_validate()**.

We adapted the code for Colab by managing file access, installing dependencies, and updating outdated functions.

# Models

The notebook contains two machine learning models



Content-Based Filtering Model

Collaborative Filtering Model

# Content-Based Filtering Model

**Recommends movies similar to those a user has already watched**

The content-based filtering model works by recommending movies similar to those a user has already watched, based only on their descriptions. In our model, we analyze elements such as the title, synopsis, and genres of movies.

**TF-IDF Vectorizer:**

Converts movie descriptions into numerical vectors

Reduces the importance of frequent words, highlights specific words

**Cosine Similarity:**

Measures similarity between movies (smaller angle = more similar)

Uses linear\_kernel for greater efficiency

# Collaborative Filtering

**Recommends movies based on user ratings**

**Uses a User-Movie Matrix: each user rates a set of movies**

**Identifies similar users to recommend new movies**

**Problem: Many missing ratings (users only rate a few movies)**

## SVD

SVD (Singular Value Decomposition) decomposes the matrix into smaller matrices

- Captures latent factors (e.g., genre preferences)
- Predicts ratings for unseen movies based on similar users' tastes

## Model Validation

Uses cross-validation to test the model's accuracy

RMSE and MAE metrics evaluate predictions

Lower RMSE & MAE = better prediction accuracy

# IV. A look at the notebook

## Data loading and preparation

### To use data from Kaggle:

1. Download it and open it in a format that we can work with using Pandas.
2. Install the Kaggle library, which helps us get files automatically. 3. Add a special file called **kaggle.json**, which acts as a key to connect to Kaggle. This file is stored in a hidden **.kaggle** folder, and we change its settings to keep it secure.
4. Download the dataset as a **ZIP file**. Since ZIP files are compressed, we need to
5. Extract them using **zipfile.ZipFile** before we can read the data.
6. Open the CSV files: **tmdb\_5000\_credits.csv** and **tmdb\_5000\_movies.csv**
7. Store these files in two Pandas DataFrames:
  - df1 for the credits file (actors and directors)
  - df2 for the movies file (titles, ratings, descriptions)

Finally, we check that everything is loaded correctly by displaying the first few rows of the data.

```
# 1 Installation de la bibliothèque Kaggle (si non installée)
!pip install kaggle

# 2 Téléchargement de la clé API de Kaggle
from google.colab import files
files.upload() # L'utilisateur doit importer son fichier kaggle.json

# 3 Configuration de l'authentification
import os
os.makedirs('/root/.kaggle', exist_ok=True) # Création du dossier .kaggle s'il n'existe pas
!mv kaggle.json /root/.kaggle/ # Déplacement du fichier au bon emplacement
os.chmod('/root/.kaggle/kaggle.json', 600) # Sécurisation du fichier

# 4 Téléchargement des datasets depuis Kaggle
!kaggle datasets download -d tmdb/tmdb-movie-metadata # Remplace le chemin de Kaggle

# 5 Décompression du fichier ZIP
import zipfile
with zipfile.ZipFile('tmdb-movie-metadata.zip', 'r') as zip_ref:
    zip_ref.extractall('.')

# 6 Chargement des fichiers dans Pandas
import pandas as pd
import numpy as np

df1 = pd.read_csv('tmdb_5000_credits.csv') # Chargement des crédits
df2 = pd.read_csv('tmdb_5000_movies.csv') # Chargement des films

# Vérification du chargement
print(df1.head())
print(df2.head())
```

```

Downloading tmdb-movie-metadata.zip to /content
 0% 0.00/8.89M [00:00<?, ?B/s]
100% 8.89M/8.89M [00:00<00:00, 108MB/s]
movie_id          title \
0      19995        Avatar
1       285  Pirates of the Caribbean: At World's End
2     206647        Spectre
3      49026  The Dark Knight Rises
4      49529        John Carter

cast \
0 [{"cast_id": 242, "character": "Jake Sully", "de...
1 [{"cast_id": 4, "character": "Captain Jack Spa...
2 [{"cast_id": 1, "character": "James Bond", "cr...
3 [{"cast_id": 2, "character": "Bruce Wayne / Ba...
4 [{"cast_id": 5, "character": "John Carter", "c...

crew \
0 [{"credit_id": "52fe48009251416c750aca23", "de...
1 [{"credit_id": "52fe4232c3a36847f800b579", "de...
2 [{"credit_id": "54805967c3a36829b5002c41", "de...
3 [{"credit_id": "52fe4781c3a36847f81398c3", "de...
4 [{"credit_id": "52fe479ac3a36847f813eaa3", "de...

budget \
0 237000000 [{"id": 28, "name": "Action"}, {"id": 12, "nam...
1 300000000 [{"id": 12, "name": "Adventure"}, {"id": 14, "...
2 245000000 [{"id": 28, "name": "Action"}, {"id": 12, "nam...
3 250000000 [{"id": 28, "name": "Action"}, {"id": 80, "nam...
4 260000000 [{"id": 28, "name": "Action"}, {"id": 12, "nam...

genres \
0 237000000 [{"id": 28, "name": "Action"}, {"id": 12, "nam...
1 300000000 [{"id": 12, "name": "Adventure"}, {"id": 14, "...
2 245000000 [{"id": 28, "name": "Action"}, {"id": 12, "nam...
3 250000000 [{"id": 28, "name": "Action"}, {"id": 80, "nam...
4 260000000 [{"id": 28, "name": "Action"}, {"id": 12, "nam...

homepage    id \
0 http://www.avatarmovie.com/ 19995
1 http://disney.go.com/disneypictures/pirates/ 285
2 http://www.sonypictures.com/movies/spectre/ 206647
3 http://www.thedarkknightrises.com/ 49026
4 http://movies.disney.com/john-carter 49529

```

To analyze movies along with their actors and directors, we need to combine the two datasets. We do this by merging the two Pandas DataFrames using the id column, which is present in both files and helps match the correct movie information.

The merge() function links the movie details from df2 (titles, ratings, descriptions) with the credits from df1 (actors, directors). After merging, each row in the new DataFrame contains all the information about a movie in one place, making it easier to analyze.

Finally, we display a preview of the merged DataFrame to make sure everything is correct before moving on to the next steps. This helps confirm that the data is well-organized and ready for analysis or creating a model. By merging and preparing the data, we now have a complete dataset with all the important details, like movie titles, ratings, actors, and directors, which we can use in the following stages of the project.

	budget	genres	homepage	id	keywords	original_language	original_title	overview	popularity	production_companies	...	runtime	spoken_languages	status	tagline	title	vote_average	vote_count	title	cast	crew
0	237000000	[{"id": 28, "name": "Action"}, {"id": 12, "name": "Advent...]	http://www.avatarmovie.com/	19995	[{"id": 1463, "name": "culture clash"}, {"id": ...]	en	Avatar	In the 22nd century, a paraplegic Marine is di...	150.437577	[{"name": "Ingenious Film Partners", "id": 299...]	...	162.0	[{"iso_639_1": "en", "name": "English"}, {"iso_639_1": "zh", "name": "中文"}]	Released	Enter the World of Pandora.	Avatar	7.2	11000	Avatar	[{"cast_id": 242, "character": "Jake Sully", "credit_id": "52fe48009251416c750aca23", "de...]	[{"credit_id": "52fe4232c3a36847f800b579", "de...]
1	300000000	[{"id": 12, "name": "Adventure"}, {"id": 14, "name": "Sci-Fi"}, {"id": 80, "name": "Thriller"]]	http://disney.go.com/disneypictures/pirates/	285	[{"id": 270, "name": "ocean"}, {"id": 726, "name": "sea"}, {"id": 100, "name": "pirate"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Sci-Fi"}, {"id": 80, "name": "Thriller"}]	en	Pirates of the Caribbean: At World's End	Captain Barbosa, long believed to be dead, ha...	139.032615	[{"name": "Walt Disney Pictures", "id": 2}, {"name": "Buena Vista Motion Pictures", "id": 100}, {"name": "Columbia Pictures", "id": 5}, {"name": "The Walt Disney Company", "id": 1463}, {"name": "Disney Channel", "id": 12}, {"name": "Disney XD", "id": 14}, {"name": "Disney Parks, Experiences and Products", "id": 1000}, {"name": "Walt Disney Animation Studios", "id": 1001}, {"name": "Walt Disney Imagineering", "id": 1002}, {"name": "Walt Disney Studios Motion Pictures", "id": 1003}, {"name": "Walt Disney Television", "id": 1004}, {"name": "Walt Disney Television Worldwide Media Networks", "id": 1005}, {"name": "Walt Disney Television Worldwide Entertainment", "id": 1006}, {"name": "Walt Disney Television Worldwide Studios", "id": 1007}, {"name": "Walt Disney Television Worldwide Home Entertainment", "id": 1008}, {"name": "Walt Disney Television Worldwide Media", "id": 1009}, {"name": "Walt Disney Television Worldwide", "id": 1010}, {"name": "Walt Disney Imagineering", "id": 1011}, {"name": "Walt Disney Imagineering", "id": 1012}, {"name": "Walt Disney Imagineering", "id": 1013}, {"name": "Walt Disney Imagineering", "id": 1014}, {"name": "Walt Disney Imagineering", "id": 1015}, {"name": "Walt Disney Imagineering", "id": 1016}, {"name": "Walt Disney Imagineering", "id": 1017}, {"name": "Walt Disney Imagineering", "id": 1018}, {"name": "Walt Disney Imagineering", "id": 1019}, {"name": "Walt Disney Imagineering", "id": 1020}, {"name": "Walt Disney Imagineering", "id": 1021}, {"name": "Walt Disney Imagineering", "id": 1022}, {"name": "Walt Disney Imagineering", "id": 1023}, {"name": "Walt Disney Imagineering", "id": 1024}, {"name": "Walt Disney Imagineering", "id": 1025}, {"name": "Walt Disney Imagineering", "id": 1026}, {"name": "Walt Disney Imagineering", "id": 1027}, {"name": "Walt Disney Imagineering", "id": 1028}, {"name": "Walt Disney Imagineering", "id": 1029}, {"name": "Walt Disney Imagineering", "id": 1030}, {"name": "Walt Disney Imagineering", "id": 1031}, {"name": "Walt Disney Imagineering", "id": 1032}, {"name": "Walt Disney Imagineering", "id": 1033}, {"name": "Walt Disney Imagineering", "id": 1034}, {"name": "Walt Disney Imagineering", "id": 1035}, {"name": "Walt Disney Imagineering", "id": 1036}, {"name": "Walt Disney Imagineering", "id": 1037}, {"name": "Walt Disney Imagineering", "id": 1038}, {"name": "Walt Disney Imagineering", "id": 1039}, {"name": "Walt Disney Imagineering", "id": 1040}, {"name": "Walt Disney Imagineering", "id": 1041}, {"name": "Walt Disney Imagineering", "id": 1042}, {"name": "Walt Disney Imagineering", "id": 1043}, {"name": "Walt Disney Imagineering", "id": 1044}, {"name": "Walt Disney Imagineering", "id": 1045}, {"name": "Walt Disney Imagineering", "id": 1046}, {"name": "Walt Disney Imagineering", "id": 1047}, {"name": "Walt Disney Imagineering", "id": 1048}, {"name": "Walt Disney Imagineering", "id": 1049}, {"name": "Walt Disney Imagineering", "id": 1050}, {"name": "Walt Disney Imagineering", "id": 1051}, {"name": "Walt Disney Imagineering", "id": 1052}, {"name": "Walt Disney Imagineering", "id": 1053}, {"name": "Walt Disney Imagineering", "id": 1054}, {"name": "Walt Disney Imagineering", "id": 1055}, {"name": "Walt Disney Imagineering", "id": 1056}, {"name": "Walt Disney Imagineering", "id": 1057}, {"name": "Walt Disney Imagineering", "id": 1058}, {"name": "Walt Disney Imagineering", "id": 1059}, {"name": "Walt Disney Imagineering", "id": 1060}, {"name": "Walt Disney Imagineering", "id": 1061}, {"name": "Walt Disney Imagineering", "id": 1062}, {"name": "Walt Disney Imagineering", "id": 1063}, {"name": "Walt Disney Imagineering", "id": 1064}, {"name": "Walt Disney Imagineering", "id": 1065}, {"name": "Walt Disney Imagineering", "id": 1066}, {"name": "Walt Disney Imagineering", "id": 1067}, {"name": "Walt Disney Imagineering", "id": 1068}, {"name": "Walt Disney Imagineering", "id": 1069}, {"name": "Walt Disney Imagineering", "id": 1070}, {"name": "Walt Disney Imagineering", "id": 1071}, {"name": "Walt Disney Imagineering", "id": 1072}, {"name": "Walt Disney Imagineering", "id": 1073}, {"name": "Walt Disney Imagineering", "id": 1074}, {"name": "Walt Disney Imagineering", "id": 1075}, {"name": "Walt Disney Imagineering", "id": 1076}, {"name": "Walt Disney Imagineering", "id": 1077}, {"name": "Walt Disney Imagineering", "id": 1078}, {"name": "Walt Disney Imagineering", "id": 1079}, {"name": "Walt Disney Imagineering", "id": 1080}, {"name": "Walt Disney Imagineering", "id": 1081}, {"name": "Walt Disney Imagineering", "id": 1082}, {"name": "Walt Disney Imagineering", "id": 1083}, {"name": "Walt Disney Imagineering", "id": 1084}, {"name": "Walt Disney Imagineering", "id": 1085}, {"name": "Walt Disney Imagineering", "id": 1086}, {"name": "Walt Disney Imagineering", "id": 1087}, {"name": "Walt Disney Imagineering", "id": 1088}, {"name": "Walt Disney Imagineering", "id": 1089}, {"name": "Walt Disney Imagineering", "id": 1090}, {"name": "Walt Disney Imagineering", "id": 1091}, {"name": "Walt Disney Imagineering", "id": 1092}, {"name": "Walt Disney Imagineering", "id": 1093}, {"name": "Walt Disney Imagineering", "id": 1094}, {"name": "Walt Disney Imagineering", "id": 1095}, {"name": "Walt Disney Imagineering", "id": 1096}, {"name": "Walt Disney Imagineering", "id": 1097}, {"name": "Walt Disney Imagineering", "id": 1098}, {"name": "Walt Disney Imagineering", "id": 1099}, {"name": "Walt Disney Imagineering", "id": 1100}, {"name": "Walt Disney Imagineering", "id": 1101}, {"name": "Walt Disney Imagineering", "id": 1102}, {"name": "Walt Disney Imagineering", "id": 1103}, {"name": "Walt Disney Imagineering", "id": 1104}, {"name": "Walt Disney Imagineering", "id": 1105}, {"name": "Walt Disney Imagineering", "id": 1106}, {"name": "Walt Disney Imagineering", "id": 1107}, {"name": "Walt Disney Imagineering", "id": 1108}, {"name": "Walt Disney Imagineering", "id": 1109}, {"name": "Walt Disney Imagineering", "id": 1110}, {"name": "Walt Disney Imagineering", "id": 1111}, {"name": "Walt Disney Imagineering", "id": 1112}, {"name": "Walt Disney Imagineering", "id": 1113}, {"name": "Walt Disney Imagineering", "id": 1114}, {"name": "Walt Disney Imagineering", "id": 1115}, {"name": "Walt Disney Imagineering", "id": 1116}, {"name": "Walt Disney Imagineering", "id": 1117}, {"name": "Walt Disney Imagineering", "id": 1118}, {"name": "Walt Disney Imagineering", "id": 1119}, {"name": "Walt Disney Imagineering", "id": 1120}, {"name": "Walt Disney Imagineering", "id": 1121}, {"name": "Walt Disney Imagineering", "id": 1122}, {"name": "Walt Disney Imagineering", "id": 1123}, {"name": "Walt Disney Imagineering", "id": 1124}, {"name": "Walt Disney Imagineering", "id": 1125}, {"name": "Walt Disney Imagineering", "id": 1126}, {"name": "Walt Disney Imagineering", "id": 1127}, {"name": "Walt Disney Imagineering", "id": 1128}, {"name": "Walt Disney Imagineering", "id": 1129}, {"name": "Walt Disney Imagineering", "id": 1130}, {"name": "Walt Disney Imagineering", "id": 1131}, {"name": "Walt Disney Imagineering", "id": 1132}, {"name": "Walt Disney Imagineering", "id": 1133}, {"name": "Walt Disney Imagineering", "id": 1134}, {"name": "Walt Disney Imagineering", "id": 1135}, {"name": "Walt Disney Imagineering", "id": 1136}, {"name": "Walt Disney Imagineering", "id": 1137}, {"name": "Walt Disney Imagineering", "id": 1138}, {"name": "Walt Disney Imagineering", "id": 1139}, {"name": "Walt Disney Imagineering", "id": 1140}, {"name": "Walt Disney Imagineering", "id": 1141}, {"name": "Walt Disney Imagineering", "id": 1142}, {"name": "Walt Disney Imagineering", "id": 1143}, {"name": "Walt Disney Imagineering", "id": 1144}, {"name": "Walt Disney Imagineering", "id": 1145}, {"name": "Walt Disney Imagineering", "id": 1146}, {"name": "Walt Disney Imagineering", "id": 1147}, {"name": "Walt Disney Imagineering", "id": 1148}, {"name": "Walt Disney Imagineering", "id": 1149}, {"name": "Walt Disney Imagineering", "id": 1150}, {"name": "Walt Disney Imagineering", "id": 1151}, {"name": "Walt Disney Imagineering", "id": 1152}, {"name": "Walt Disney Imagineering", "id": 1153}, {"name": "Walt Disney Imagineering", "id": 1154}, {"name": "Walt Disney Imagineering", "id": 1155}, {"name": "Walt Disney Imagineering", "id": 1156}, {"name": "Walt Disney Imagineering", "id": 1157}, {"name": "Walt Disney Imagineering", "id": 1158}, {"name": "Walt Disney Imagineering", "id": 1159}, {"name": "Walt Disney Imagineering", "id": 1160}, {"name": "Walt Disney Imagineering", "id": 1161}, {"name": "Walt Disney Imagineering", "id": 1162}, {"name": "Walt Disney Imagineering", "id": 1163}, {"name": "Walt Disney Imagineering", "id": 1164}, {"name": "Walt Disney Imagineering", "id": 1165}, {"name": "Walt Disney Imagineering", "id": 1166}, {"name": "Walt Disney Imagineering", "id": 1167}, {"name": "Walt Disney Imagineering", "id": 1168}, {"name": "Walt Disney Imagineering", "id": 1169}, {"name": "Walt Disney Imagineering", "id": 1170}, {"name": "Walt Disney Imagineering", "id": 1171}, {"name": "Walt Disney Imagineering", "id": 1172}, {"name": "Walt Disney Imagineering", "id": 1173}, {"name": "Walt Disney Imagineering", "id": 1174}, {"name": "Walt Disney Imagineering", "id": 1175}, {"name": "Walt Disney Imagineering", "id": 1176}, {"name": "Walt Disney Imagineering", "id": 1177}, {"name": "Walt Disney Imagineering", "id": 1178}, {"name": "Walt Disney Imagineering", "id": 1179}, {"name": "Walt Disney Imagineering", "id": 1180}, {"name": "Walt Disney Imagineering", "id": 1181}, {"name": "Walt Disney Imagineering", "id": 1182}, {"name": "Walt Disney Imagineering", "id": 1183}, {"name": "Walt Disney Imagineering", "id": 1184}, {"name": "Walt Disney Imagineering", "id": 1185}, {"name": "Walt Disney Imagineering", "id": 1186}, {"name": "Walt Disney Imagineering", "id": 1187}, {"name": "Walt Disney Imagineering", "id": 1188}, {"name": "Walt Disney Imagineering", "id": 1189}, {"name": "Walt Disney Imagineering", "id": 1190}, {"name": "Walt Disney Imagineering", "id": 1191}, {"name": "Walt Disney Imagineering", "id": 1192}, {"name": "Walt Disney Imagineering", "id": 1193}, {"name": "Walt Disney Imagineering", "id": 1194}, {"name": "Walt Disney Imagineering", "id": 1195}, {"name": "Walt Disney Imagineering", "id": 1196}, {"name": "Walt Disney Imagineering", "id": 1197}, {"name": "Walt Disney Imagineering", "id": 1198}, {"name": "Walt Disney Imagineering", "id": 1199}, {"name": "Walt Disney Imagineering", "id": 1200}, {"name": "Walt Disney Imagineering", "id": 1201}, {"name": "Walt Disney Imagineering", "id": 1202}, {"name": "Walt Disney Imagineering", "id": 1203}, {"name": "Walt Disney Imagineering", "id": 1204}, {"name": "Walt Disney Imagineering", "id": 1205}, {"name": "Walt Disney Imagineering", "id": 1206}, {"name": "Walt Disney Imagineering", "id": 1207}, {"name": "Walt Disney Imagineering", "id": 1208}, {"name": "Walt Disney Imagineering", "id": 1209}, {"name": "Walt Disney Imagineering", "id": 1210}, {"name": "Walt Disney Imagineering", "id": 1211}, {"name": "Walt Disney Imagineering", "id": 1212}, {"name": "Walt Disney Imagineering", "id": 1213}, {"name": "Walt Disney Imagineering", "id": 1214}, {"name": "Walt Disney Imagineering", "id": 1215}, {"name": "Walt Disney Imagineering", "id": 1216}, {"name": "Walt Disney Imagineering", "id": 1217}, {"name": "Walt Disney Imagineering", "id": 1218}, {"name": "Walt Disney Imagineering", "id": 1219}, {"name": "Walt Disney Imagineering", "id": 1220}, {"name": "Walt Disney Imagineering", "id": 1221}, {"name": "Walt Disney Imagineering", "id": 1222}, {"name": "Walt Disney Imagineering", "id": 1223}, {"name": "Walt Disney Imagineering", "id": 1224}, {"name": "Walt Disney Imagineering", "id": 1225}, {"name": "Walt Disney Imagineering", "id": 1226}, {"name": "Walt Disney Imagineering", "id": 1227}, {"name": "Walt Disney Imagineering", "id": 1228}, {"name": "Walt Disney Imagineering", "id": 1229}, {"name": "Walt Disney Imagineering", "id": 1230}, {"name": "Walt Disney Imagineering", "id": 1231}, {"name": "Walt Disney Imagineering", "id": 1232}, {"name": "Walt Disney Imagineering", "id": 1233}, {"name": "Walt Disney Imagineering", "id": 1234}, {"name": "Walt Disney Imagineering", "id": 1235}, {"name": "Walt Disney Imagineering", "id": 1236}, {"name": "Walt Disney Imagineering", "id": 1237}, {"name": "Walt Disney Imagineering", "id": 1238}, {"name": "Walt Disney Imagineering", "id": 1239}, {"name": "Walt Disney Imagineering", "id": 1240}, {"name": "Walt Disney Imagineering", "id": 1241}, {"name": "Walt Disney Imagineering", "id": 1242}, {"name": "Walt Disney Imagineering", "id": 1243}, {"name": "Walt Disney Imagineering", "id": 1244}, {"name": "Walt Disney Imagineering", "id": 1245}, {"name": "Walt Disney Imagineering", "id": 1246}, {"name": "Walt Disney Imagineering", "id": 1247}, {"name": "Walt Disney Imagineering", "id": 1248}, {"name": "Walt Disney Imagineering", "id": 1249}, {"name": "Walt Disney Imagineering", "id": 1250}, {"name": "Walt Disney Imagineering", "id": 1251}, {"name": "Walt Disney Imagineering", "id": 1252}, {"name": "Walt Disney Imagineering", "id": 1253}, {"name": "Walt Disney Imagineering", "id": 1254}, {"name": "Walt Disney Imagineering", "id": 1255}, {"name": "Walt Disney Imagineering", "id": 1256}, {"name": "Walt Disney Imagineering", "id": 1257}, {"name": "Walt Disney Imagineering", "id": 1258}, {"name": "Walt Disney Imagineering", "id": 1259}, {"name": "Walt Disney Imagineering", "id": 1260}, {"name": "Walt Disney Imagineering", "id": 1261}, {"name": "Walt Disney Imagineering", "id": 1262}, {"name": "Walt Disney Imagineering", "id": 1263}, {"name": "Walt Disney Imagineering", "id": 1264}, {"name": "Walt Disney Imagineering", "id": 1265}, {"name": "Walt Disney Imagineering", "id": 1266}, {"name": "Walt Disney Imagineering", "id": 1267}, {"name": "Walt Disney Imagineering", "id": 1268}, {"name": "Walt Disney Imagineering", "id": 1269}, {"name": "Walt Disney Imagineering", "id": 1270}, {"name": "Walt Disney Imagineering", "id": 1271}, {"name": "Walt Disney Imagineering", "id": 1272}, {"name": "Walt Disney Imagineering", "id": 1273}, {"name": "Walt Disney Imagineering", "id": 1274}, {"name": "Walt Disney Imagineering", "id": 1275}, {"name": "Walt Disney Imagineering", "id": 1276}, {"name": "Walt Disney Imagineering", "id": 1277}, {"name": "Walt Disney Imagineering", "id": 1278}, {"name": "Walt Disney Imagineering", "id": 1279}, {"name": "Walt Disney Imagineering", "id": 1280}, {"name": "Walt Disney Imagineering", "id": 1281}, {"name": "Walt Disney Imagineering", "id": 1282}, {"name": "Walt Disney Imagineering", "											

# IV. A look at the notebook

## Processing Popularity Data

**Weighted\_rating** calculates an adjusted score for each movie.

It takes in a movie (x) and two numbers, **m** and **C**, which represent the minimum number of votes a movie must have to be considered, and the average rating of all movies, respectively.

Inside the function

**v**: number of votes for a movie.

**R**: movie's average rating.

The formula is based on **IMDb's** ranking system, where the movie's rating **R** is combined with the overall average **C**, but the number of votes **v** is also considered. If a movie has a lot of votes, its rating matters more. If it has fewer votes, the overall average rating **C** plays a bigger role in adjusting the score.

Once the function is defined, we apply it to the dataset using:

**q\_movies['score'] = q\_movies.apply(weighted\_rating, axis=1)**

This creates a new column called **score**, which contains the weighted rating of each movie.

Next, we sort the movies by this score:

**q\_movies = q\_movies.sort\_values('score', ascending=False)**

This provides a ranked list of movies from the most to the least popular.

Finally, we display the top 10 movies based on popularity:

**q\_movies[['title', 'vote\_count', 'vote\_average', 'score']].head(10).**

This allows us to see the highest-rated movies based on their weighted popularity score.

This method helps us recommend movies based on their actual popularity, considering both ratings and the number of votes.

```
[ ] def weighted_rating(x, m=m, C=C):
    v = x['vote_count']
    R = x['vote_average']
    # Calculation based on the IMDB formula
    return (v/(v+m) * R) + (m/(m+v) * C)

[ ] # Define a new feature 'score' and calculate its value with `weighted_rating()`
q_movies['score'] = q_movies.apply(weighted_rating, axis=1)
```

▶ #Sort movies based on score calculated above

q\_movies = q\_movies.sort\_values('score', ascending=False)

#Print the top 15 movies

q\_movies[['title', 'vote\_count', 'vote\_average', 'score']].head(10)

		title	vote_count	vote_average	score
1881		The Shawshank Redemption	8205	8.5	8.059258
662		Fight Club	9413	8.3	7.939256
65		The Dark Knight	12002	8.2	7.920020
3232		Pulp Fiction	8428	8.3	7.904645
96		Inception	13752	8.1	7.863239
3337		The Godfather	5893	8.4	7.851236
95		Interstellar	10867	8.1	7.809479
809		Forrest Gump	7927	8.2	7.803188
329	The Lord of the Rings: The Return of the King		8064	8.1	7.727243
1990		The Empire Strikes Back	5879	8.2	7.697884

# Content-Based Recommendation (Similarity)

**Recommends movies similar to those a user has already watched**

**Based on movie descriptions:**

Title, synopsis, genres

**TF-IDF Vectorizer:**

Converts movie descriptions into numerical vectors

Reduces the importance of frequent words,  
highlights specific words

**Cosine Similarity:**

Measures similarity between movies (smaller angle = more similar)

Uses linear\_kernel for greater efficiency

```
[ ] #Import TfIdfVectorizer from scikit-learn
from sklearn.feature_extraction.text import TfidfVectorizer

#Define a TF-IDF Vectorizer Object. Remove all english stop words such as 'the', 'a'
tfidf = TfidfVectorizer(stop_words='english')

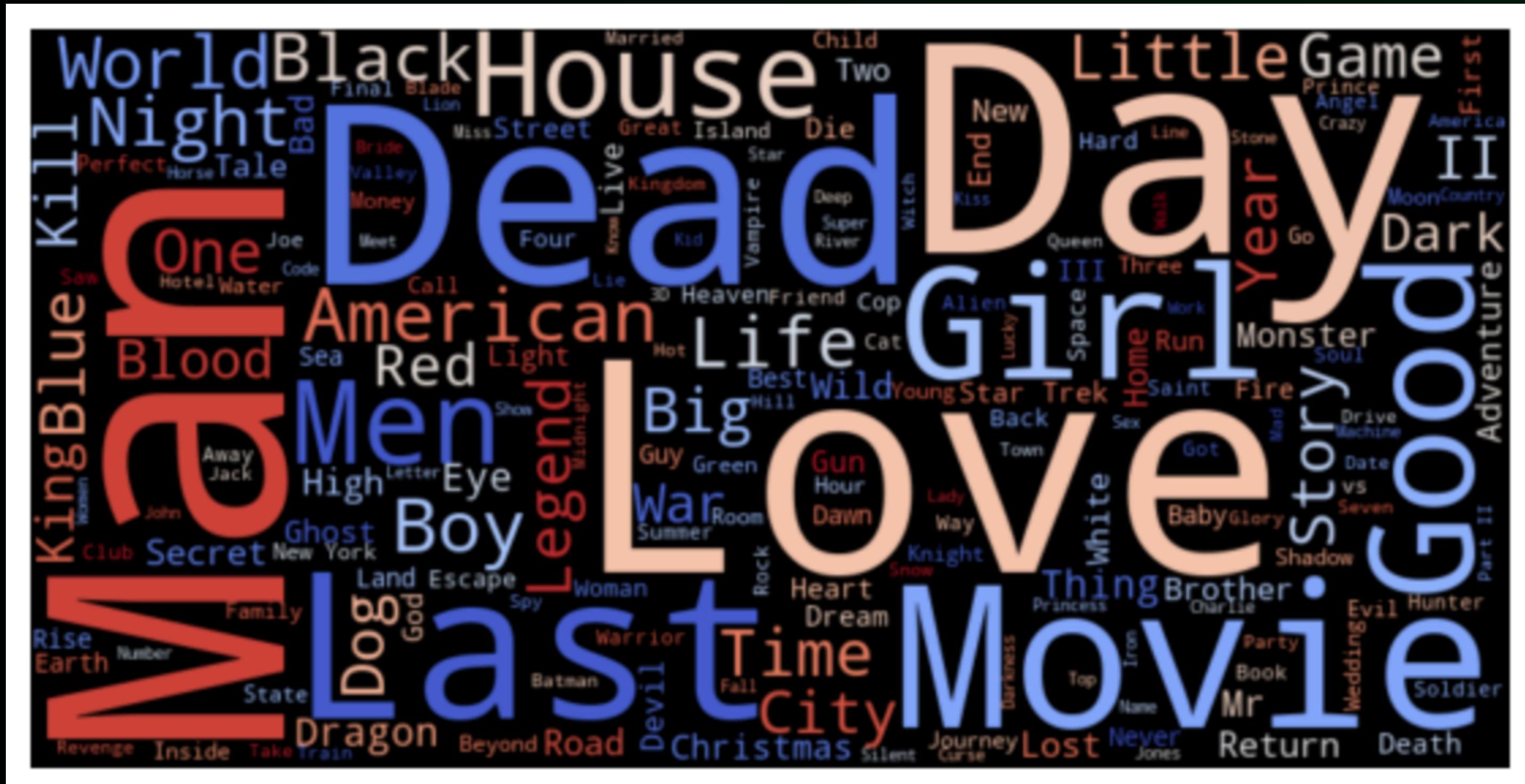
#Replace NaN with an empty string
df2['overview'] = df2['overview'].fillna('')

#Construct the required TF-IDF matrix by fitting and transforming the data
tfidf_matrix = tfidf.fit_transform(df2['overview'])

#Output the shape of tfidf_matrix
tfidf_matrix.shape
```

→ (4803, 20978)

# WordCloud



# Models

## Similarity matrix for recommendations

```
▶ # Import linear_kernel
from sklearn.metrics.pairwise import linear_kernel

# Compute the cosine similarity matrix
cosine_sim = linear_kernel(tfidf_matrix, tfidf_matrix)

[ ] #Construct a reverse map of indices and movie titles
indices = pd.Series(df2.index, index=df2['title']).drop_duplicates()

[ ] # Function that takes in movie title as input and outputs most similar movies
def get_recommendations(title, cosine_sim=cosine_sim):
    # Get the index of the movie that matches the title
    idx = indices.get_loc(title)

    # Get the pairwise similarity scores of all movies with that movie
    sim_scores = list(enumerate(cosine_sim[idx]))

    # Sort the movies based on the similarity scores
    sim_scores = sorted(sim_scores, key=lambda x: x[1], reverse=True)

    # Get the scores of the 10 most similar movies
    sim_scores = sim_scores[1:11]

    # Get the movie indices
    movie_indices = [i[0] for i in sim_scores]

    # Return the top 10 most similar movies
    return df2['title'].iloc[movie_indices]
```

Once the summaries are converted into vectors, and we've done the wordcloud, we calculate the cosine similarity between them.

This measure evaluates how close two movies are by comparing their descriptions as numerical vectors.

The higher the value, the more similar the movies.

```
[ ] get_recommendations('The Dark Knight Rises')
```

	title
65	The Dark Knight
299	Batman Forever
428	Batman Returns
1359	Batman
3854	Batman: The Dark Knight Returns, Part 2
119	Batman Begins
2507	Slow Burn
9	Batman v Superman: Dawn of Justice
1181	JFK
210	Batman & Robin

dtype: object

Recommend movies similar to a given movie.

Retrieve the most similar movies based on their similarity score and display the closest matches to the selected film.

Thus, if a user likes a movie, this model will suggest other movies with similar themes, stories, or genres.

# Enhanced Recommendations with Actors, Genres, and Directors

## DATA CLEANING

Standardizes all names by converting them to lowercase and removing unnecessary spaces

Columns: actors, keywords, directors, and genres

All this information is combined into a single column called "soup."

CountVectorizer to convert this text into a numerical matrix

From this matrix, we compute the cosine similarity between movies, allowing us to measure their similarity based on combined characteristics (actors, director, genre, keywords). The higher the similarity, the more alike the movies are.

```
# Function to convert all strings to lower case and strip names of spaces
def clean_data(x):
    if isinstance(x, list):
        return [str.lower(i.replace(" ", "")) for i in x]
    else:
        #Check if director exists. If not, return empty string
        if isinstance(x, str):
            return str.lower(x.replace(" ", ""))
        else:
            return ''
```

```
# Apply clean_data function to your features.
features = ['cast', 'keywords', 'director', 'genres']

for feature in features:
    df2[feature] = df2[feature].apply(clean_data)

[ ] def create_soup(x):
    return ' '.join(x['keywords']) + ' ' + ' '.join(x['cast']) + ' ' + x['director'] + ' ' + ' '.join(x['genres'])

df2['soup'] = df2.apply(create_soup, axis=1)
```

```
[ ] # Import CountVectorizer and create the count matrix
from sklearn.feature_extraction.text import CountVectorizer

count = CountVectorizer(stop_words='english')
count_matrix = count.fit_transform(df2['soup'])

[ ] # Compute the Cosine Similarity matrix based on the count_matrix
from sklearn.metrics.pairwise import cosine_similarity

cosine_sim2 = cosine_similarity(count_matrix, count_matrix)

[ ] # Reset index of our main DataFrame and construct reverse mapping as before
df2 = df2.reset_index()
indices = pd.Series(df2.index, index=df2['title'])
```

Finally, we use this similarity matrix to recommend movies most similar to a user's preferred film—for instance, suggesting films with shared actors, genres, or directors—resulting in more personalized recommendations that consider a broader range of factors.

title	
65	The Dark Knight
119	Batman Begins
4638	Amidst the Devil's Wings
1196	The Prestige
3073	Romeo Is Bleeding
3326	Black November
1503	Takers
1986	Faster
303	Catwoman
747	Gangster Squad

dtype: object

# Collaborative Filtering Model

- Uses user ratings to predict preferences
- User-movie matrix:
  - Rows = users, Columns = movies, Cells = ratings
- Challenges: Sparse matrix (most ratings are missing)

The ratings\_small.csv file contains user ratings for various movies. We apply Singular Value Decomposition (SVD) to train a model that can predict missing ratings.

The process begins by loading the data using pandas and creating a Dataset with Surprise, a specialized library for recommendation systems. We then use an SVD object to build the model, and the cross\_validate method to evaluate its performance using RMSE (Root Mean Squared Error) and MAE (Mean Absolute Error). These metrics indicate the model's prediction accuracy on test data.

## CROSS-VALIDATION

- Validates the model's accuracy using cross\_validate()
- Splits the data into subsets to test the model
- Metrics: RMSE (Root Mean Squared Error) and MAE (Mean Absolute Error)
- Goal: Check the model's ability to predict ratings for unseen movies

```
▶ from surprise import SVD
from surprise import Dataset
from surprise import Reader
from surprise.model_selection import cross_validate
from surprise import accuracy # Import the accuracy module

# Load the ratings data
import pandas as pd

reader = Reader()
ratings = pd.read_csv('ratings_small.csv') # Updated path to the extracted file
# Create the dataset
data = Dataset.load_from_df(ratings[['userId', 'movieId', 'rating']], reader)
svd = SVD()

# Evaluate the model using cross-validation and print RMSE and MAE
# Perform cross-validation (no need to call data.split)
# For example, to evaluate an SVD algorithm:
results = cross_validate(svd, data, measures=['RMSE', 'MAE'], cv=5, verbose=True)

# Print the average RMSE and MAE across folds
print(f"Average RMSE: {results['test_rmse'].mean()}")
print(f"Average MAE: {results['test_mae'].mean()}")
```

◀ Evaluating RMSE, MAE of algorithm SVD on 5 split(s).

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	Std
RMSE (testset)	0.8903	0.9007	0.8981	0.8964	0.8979	0.8967	0.0035
MAE (testset)	0.6846	0.6948	0.6932	0.6888	0.6928	0.6908	0.0037
Fit time	1.56	1.17	0.97	0.95	0.98	1.13	0.23
Test time	0.34	0.07	0.20	0.07	0.07	0.15	0.10
Average RMSE:	0.896672844867779						
Average MAE:	0.6908373208183228						

```
▶ trainset = data.build_full_trainset()
svd.fit(trainset)

◀ <surprise.prediction_algorithms.matrix_factorization.SVD at 0x7f5d7f6c5710>
```

ratings[ratings['userId'] == 1]				
	userId	movieId	rating	timestamp
0	1	31	2.5	1260759144
1	1	1029	3.0	1260759179
2	1	1061	3.0	1260759182
3	1	1129	2.0	1260759185
4	1	1172	4.0	1260759205
5	1	1263	2.0	1260759151
6	1	1287	2.0	1260759187
7	1	1293	2.0	1260759148
8	1	1339	3.5	1260759125
9	1	1343	2.0	1260759131
10	1	1371	2.5	1260759135
11	1	1405	1.0	1260759203
12	1	1953	4.0	1260759191

After performing cross-validation, the model is trained on the entire dataset using fit(), allowing us to make predictions.

User 1 would give to movie the rate of 302.

This process enables us to predict movies a user might enjoy based on the past ratings of other users. Using matrix decomposition via SVD, the model identifies movies with similar rating patterns.

# CONCLUSION

Recommendation methods			
	<b>Popularity-based recommendations</b>	<b>Content-based recommendations</b>	<b>Collaborative filtering (SVD)</b>
Explanation	Rely on global audience votes.	Identify movies with similar descriptions	Personalize recommendations based on user preferences.
Advantages	Simple and fast but lacks personalization	Useful for users looking for similar movies	The most advanced, as it learns individual user preferences

By combining these approaches, we can create a hybrid recommendation system similar to those used by Netflix or Amazon Prime, improving recommendation accuracy and user experience.

And finally, this project helps us to understand each specific term and how to run a program.

# THANK YOU

FOR YOUR ATTENTION