

Movie Recommendation System

EDA and Prediction System (Notebook)

In [203]:

```
# Importing Libraries
import numpy as np
import pandas as pd
```

Load the datasets

In [204]:

```
# Reading Files
credits_df = pd.read_csv('tmdb_5000_credits.csv')
```

In [205]:

```
movies_df = pd.read_csv('tmdb_5000_movies.csv')
```

In [206]:

```
# Printing Top 5 Data
credits_df.head()
```

Out[206]:

	movie_id	title	cast	crew
0	19995	Avatar	[[{"cast_id": 242, "character": "Jake Sully", "...	[[{"credit_id": "52fe48009251416c750aca23", "de...
1	285	Pirates of the Caribbean: At World's End	[[{"cast_id": 4, "character": "Captain Jack Spa...	[[{"credit_id": "52fe4232c3a36847f800b579", "de...
2	206647	Spectre	[[{"cast_id": 1, "character": "James Bond", "cr...	[[{"credit_id": "54805967c3a36829b5002c41", "de...
3	49026	The Dark Knight Rises	[[{"cast_id": 2, "character": "Bruce Wayne / Ba...	[[{"credit_id": "52fe4781c3a36847f81398c3", "de...
4	49529	John Carter	[[{"cast_id": 5, "character": "John Carter", "c...	[[{"credit_id": "52fe479ac3a36847f813eaa3", "de...

In [207]:

```
# Printing Top 5 Data
movies_df.head()
```

Out[207]:

	budget	genres	homepage	id	keywords	original_language	original_title	overview	populari
0	237000000	[[{"id": 28, "name": "Action"}, {"id": 12, "nam...	http://www.avatarmovie.com/	19995	[[{"id": 1463, "name": "culture clash"}, {"id": ...	en	Avatar	In the 22nd century, a paraplegic Marine is di...	150.4375
1	300000000	[[{"id": 12, "name": "Adventure"}, {"id": 14, "...	http://disney.go.com/disneypictures/pirates/	285	[[{"id": 270, "name": "ocean"}, {"id": 726, "na...	en	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...	139.0826
2	245000000	[[{"id": 28, "name": "Action"}, {"id": 12, "nam...	http://www.sonypictures.com/movies/spectre/	206647	[[{"id": 470, "name": "spy"}, {"id": 818, "name...	en	Spectre	A cryptic message from Bond's past sends him o...	107.3767
3	250000000	[[{"id": 28, "name": "Action"}, {"id": 80, "nam...	http://www.thedarkknightrises.com/	49026	[[{"id": 849, "name": "dc comics"}, {"id": 853,...	en	The Dark Knight Rises	Following the death of District Attorney Harve...	112.3129
4	260000000	[[{"id": 28, "name": "Action"}, {"id": 12, "nam...	http://movies.disney.com/john-carter	49529	[[{"id": 818, "name": "based on novel"}, {"id": ...	en	John Carter	John Carter is a war-weary, former military ca...	43.9269

In [208]:

```
# Merging two different datasets into one to work as a one dataframe based on 'title' or 'id'
merge_df = pd.merge(movies_df, credits_df, on='title')
```

In [209]:

```
# Top 5 Merge Data
merge_df.head()
```

Out[209]:

	budget	genres	homepage	id	keywords	original_language	original_title	overview	populari
0	237000000	[[{"id": 28, "name": "Action"}, {"id": 12, "nam...	http://www.avatarmovie.com/	19995	[[{"id": 1463, "name": "culture clash"}, {"id": ...	en	Avatar	In the 22nd century, a paraplegic Marine is di...	150.4375
1	300000000	[[{"id": 12, "name": "Adventure"}, {"id": 14, "...	http://disney.go.com/disneypictures/pirates/	285	[[{"id": 270, "name": "ocean"}, {"id": 726, "na...	en	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...	139.0826
2	245000000	[[{"id": 28, "name": "Action"}, {"id": 12, "nam...	http://www.sonypictures.com/movies/spectre/	206647	[[{"id": 470, "name": "spy"}, {"id": 818, "name...	en	Spectre	A cryptic message from Bond's past sends him o...	107.3767
3	250000000	[[{"id": 28, "name": "Action"}, {"id": 80, "nam...	http://www.thedarkknightrises.com/	49026	[[{"id": 849, "name": "dc comics"}, {"id": 853,...	en	The Dark Knight Rises	Following the death of District Attorney Harve...	112.3129
4	260000000	[[{"id": 28, "name": "Action"}, {"id": 12, "nam...	http://movies.disney.com/john-carter	49529	[[{"id": 818, "name": "based on novel"}, {"id": ...	en	John Carter	John Carter is a war-weary, former military ca...	43.9269

5 rows × 23 columns

Overview

In [210]:

```
print(merge_df.info())

<class 'pandas.core.frame.DataFrame'>
Int64Index: 4809 entries, 0 to 4808
Data columns (total 23 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   budget          4809 non-null   int64
 1   genres          4809 non-null   object
 2   homepage        1713 non-null   object
```

Insight:

- The dataset contains information about movies, encompassing various aspects such as budget, genres, production details, release date, and popularity, cast, crew.
- There are numeric columns for quantitative analysis, categorical columns for categorization, and text data for potential natural language processing tasks.

(4809, 23)

Insight:

- There are in total 4809 movies with 23 features in it the dataframe

Selecting only require fields

```
['budget', 'genres', 'homepage', 'id', 'keywords', 'original_language', 'original_title', 'overview', 'popularity', 'production_companies', 'production_countries', 'release_date', 'revenue', 'runtime', 'spoken_languages', 'status', 'tagline', 'title', 'vote_average', 'vote_count', 'movie_id', 'cast', 'crew']
```

	id	title	genres	keywords	overview	cast	crew
0	19995	Avatar	["(id": 28, "name": "Action"), {"id": 12, "nam...	["(id": 1463, "name": "culture clash"), {"id": ...	In the 22nd century, a paraplegic Marine is di...	["(cast_id": 242, "character": "Jake Sully", "...	["(credit_id": "52fe48009251416c750aca23", "de...
1	285	Pirates of the Caribbean: At World's End	["(id": 12, "name": "Adventure"), {"id": 14, "...	["(id": 270, "name": "ocean"), {"id": 726, "na...	Captain Barbosa, long believed to be dead, ha...	["(cast_id": 4, "character": "Captain Jack Spa...	["(credit_id": "52fe4232c3a36847f800b579", "de...
2	206647	Spectre	["(id": 28, "name": "Action"), {"id": 12, "nam...	["(id": 470, "name": "spy"), {"id": 818, "name..."	A cryptic message from Bond's past sends him o...	["(cast_id": 1, "character": "James Bond", "cr...	["(credit_id": "54805967c3a36829b5002c41", "de...
3	49026	The Dark Knight Rises	["(id": 28, "name": "Action"), {"id": 80, "nam...	["(id": 849, "name": "dc comics"), {"id": 853, ...	Following the death of District Attorney Harvey...	["(cast_id": 2, "character": "Bruce Wayne / Ba...	["(credit_id": "52fe4781c336847f81398c3", "de...
4	49529	John Carter	["(id": 28, "name": "Action"), {"id": 12, "nam...	["(id": 818, "name": "based on novel"), {"id": "...	John Carter is a war-weary, former military ca...	["(cast_id": 5, "character": "John Carter", "c...	["(credit_id": "52fe479ac3a36847f813ea33", "de...

 $(4809, 7)$ **Insight:**

- There are 4807 movies with selecting only 7 features

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 4809 entries, 0 to 4808
Data columns (total 7 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   id               4809 non-null   int64
1   title            4809 non-null   object
2   genres           4809 non-null   object
3   keywords         4809 non-null   object
4   overview         4806 non-null   object
5   cast             4809 non-null   object
6   crew             4809 non-null   object
dtypes: int64(1), object(6)
memory usage: 300.6+ KB
```

```
id          0
title       0
genres      0
keywords    0
overview    3
cast        0
crew        0
dtype: int64
```

Insight

- There are three missing values in the overview column
- Different Imputation techniques could be performed to fill those values.
- For smaller amount of overview, we could remove it

```
/var/folders/yb/v5f667y93j393zcsz2t_26qw0000gn/T/ipykernel_13438/3891106807.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
movies.dropna(inplace=True)
```

In [217]:

```
# Rechecking for missing values
movies.isnull().sum()
```

Out[217]:

```
id          0
title       0
genres      0
keywords    0
overview    0
cast        0
crew        0
dtype: int64
```

In [218]:

```
# Checking for duplicates
movies.duplicated().sum()
```

Out[218]: 0

In [219]:

```
movies.iloc[0].genres
```

Out[219]: '[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]'

Insight:

- Since it is in form of list of dictionaries, we only need the name key's value for drawing the genres name associated with the movie.
- `ast.literal_eval` is used to safely evaluate the expression represented as a string containing a Python literal or container display

In [220]:

```
import ast
# Function to get the name key's value in a List
def convert(obj):
    L = []
    for i in ast.literal_eval(obj):
        L.append(i['name'])
    return L
```

In [221]:

```
# Applying the function and updating the columns
movies['genres'] = movies['genres'].apply(convert)
movies['keywords'] = movies['keywords'].apply(convert)
```

/var/folders/yb/v5f667y93j393zcsz2t_26qw0000gn/T/ipykernel_13438/2687346810.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
movies['genres'] = movies['genres'].apply(convert)
/var/folders/yb/v5f667y93j393zcsz2t_26qw0000gn/T/ipykernel_13438/2687346810.py:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
movies['keywords'] = movies['keywords'].apply(convert)
```

In [222]:

```
# Checking for the changes
movies.head()
```

Out[222]:

	id	title	genres	keywords	overview	cast	crew
0	19995	Avatar	[Action, Adventure, Fantasy, Science Fiction]	[culture clash, future, space war, space colon...	In the 22nd century, a paraplegic Marine is di...	[{"cast_id": 242, "character": "Jake Sully", "...	[{"credit_id": "52fe48009251416c750aca23", "de...
1	285	Pirates of the Caribbean: At World's End	[Adventure, Fantasy, Action]	[ocean, drug abuse, exotic island, east india ...	Captain Barbossa, long believed to be dead, ha...	[{"cast_id": 4, "character": "Captain Jack Spa...	[{"credit_id": "52fe4232c3a36847f800b579", "de...
2	206647	Spectre	[Action, Adventure, Crime]	[spy, based on novel, secret agent, sequel, mi...	A cryptic message from Bond's past sends him o...	[{"cast_id": 1, "character": "James Bond", "cr...	[{"credit_id": "54805967c3a36829b5002c41", "de...
3	49026	The Dark Knight Rises	[Action, Crime, Drama, Thriller]	[dc comics, crime fighter, terrorist, secret i...	Following the death of District Attorney Harve...	[{"cast_id": 2, "character": "Bruce Wayne / Ba...	[{"credit_id": "52fe4781c3a36847f81398c3", "de...
4	49529	John Carter	[Action, Adventure, Science Fiction]	[based on novel, mars, medallion, space travel...	John Carter is a war-weary, former military ca...	[{"cast_id": 5, "character": "John Carter", "c...	[{"credit_id": "52fe479ac3a36847f813eaa3", "de...

In [223]:

```
# 0th index values
movies['cast'][0]
```

Out[223]: '[{"cast_id": 242, "character": "Jake Sully", "credit_id": "5602a8a7c3a3685532001c9a", "gender": 2, "id": 65731, "name": "Sam Worthington", "order": 0}, {"cast_id": 3, "character": "Neytiri", "credit_id": "52fe48009251416c750ac9cb", "gender": 1, "id": 8691, "name": "Zoe Saldana", "order": 1}, {"cast_id": 25, "character": "Dr. Grace Augustine", "credit_id": "52fe48009251416c750aca39", "gender": 1, "id": 10205, "name": "Sigourney Weaver", "order": 2}, {"cast_id": 4, "character": "Col. Quaritch", "credit_id": "52fe48009251416c750ac9cf", "gender": 2, "id": 32747, "name": "Stephen Lang", "order": 3}, {"cast_id": 5, "character": "Trudy Chacon", "credit_id": "52fe48009251416c750ac9d3", "gender": 1, "id": 17647, "name": "Michelle Rodriguez", "order": 4}, {"cast_id": 8, "character": "Selfridge", "credit_id": "52fe48009251416c750ac9e1", "gender": 2, "id": 1771, "name": "Giovanni Ribisi", "order": 5}, {"cast_id": 7, "character": "Norm Spellman", "credit_id": "52fe48009251416c750ac9dd", "gender": 2, "id": 59231, "name": "Joel David Moore", "order": 6}, {"cast_id": 9, "character": "Moat", "credit_id": "52fe48009251416c750ac9e5", "gender": 1, "id": 30485, "name": "CCH Pounder", "order": 7}, {"cast_id": 11, "character": "Eytukan", "credit_id": "52fe48009251416c750ac9ed", "gender": 2, "id": 15853, "name": "Wes Studi", "order": 8}, {"cast_id": 10, "character": "Tsu\Tey", "credit_id": "52fe48009251416c750ac9e9", "gender": 2, "id": 10964, "name": "Laz Alonso", "order": 9}, {"cast_id": 12, "character": "Dr. Max Patel", "credit_id": "52fe48009251416c750ac9f1", "gender": 2, "id": 95697, "name": "Dileep Rao", "order": 10}, {"cast_id": 13, "character": "Lyle Wainfleet", "credit_id": "52fe48009251416c750ac9f5", "gender": 2, "id": 98215, "name": "Matt Gerald", "order": 11}, {"cast_id": 32, "character": "Private Fike", "credit_id": "52fe48009251416c750aca5b", "gender": 2, "id": 154153, "name": "Sean Anthon ...", "order": 12}, {"cast_id": 33, "character": "Miles", "credit_id": "52fe48009251416c750aca5c", "gender": 2, "id": 154154, "name": "Sean Anthon ...", "order": 13}, {"cast_id": 34, "character": "Miles", "credit_id": "52fe48009251416c750aca5d", "gender": 2, "id": 154155, "name": "Sean Anthon ...", "order": 14}, {"cast_id": 35, "character": "Miles", "credit_id": "52fe48009251416c750aca5e", "gender": 2, "id": 154156, "name": "Sean Anthon ...", "order": 15}, {"cast_id": 36, "character": "Miles", "credit_id": "52fe48009251416c750aca5f", "gender": 2, "id": 154157, "name": "Sean Anthon ...", "order": 16}, {"cast_id": 37, "character": "Miles", "credit_id": "52fe48009251416c750aca60", "gender": 2, "id": 154158, "name": "Sean Anthon ...", "order": 17}, {"cast_id": 38, "character": "Miles", "credit_id": "52fe48009251416c750aca61", "gender": 2, "id": 154159, "name": "Sean Anthon ...", "order": 18}, {"cast_id": 39, "character": "Miles", "credit_id": "52fe48009251416c750aca62", "gender": 2, "id": 154160, "name": "Sean Anthon ...", "order": 19}, {"cast_id": 40, "character": "Miles", "credit_id": "52fe48009251416c750aca63", "gender": 2, "id": 154161, "name": "Sean Anthon ...", "order": 20}, {"cast_id": 41, "character": "Miles", "credit_id": "52fe48009251416c750aca64", "gender": 2, "id": 154162, "name": "Sean Anthon ...", "order": 21}, {"cast_id": 42, "character": "Miles", "credit_id": "52fe48009251416c750aca65", "gender": 2, "id": 154163, "name": "Sean Anthon ...", "order": 22}, {"cast_id": 43, "character": "Miles", "credit_id": "52fe48009251416c750aca66", "gender": 2, "id": 154164, "name": "Sean Anthon ...", "order": 23}, {"cast_id": 44, "character": "Miles", "credit_id": "52fe48009251416c750aca67", "gender": 2, "id": 154165, "name": "Sean Anthon ...", "order": 24}, {"cast_id": 45, "character": "Miles", "credit_id": "52fe48009251416c750aca68", "gender": 2, "id": 154166, "name": "Sean Anthon ...", "order": 25}, {"cast_id": 46, "character": "Miles", "credit_id": "52fe48009251416c750aca69", "gender": 2, "id": 154167, "name": "Sean Anthon ...", "order": 26}, {"cast_id": 47, "character": "Miles", "credit_id": "52fe48009251416c750aca6a", "gender": 2, "id": 154168, "name": "Sean Anthon ...", "order": 27}, {"cast_id": 48, "character": "Miles", "credit_id": "52fe48009251416c750aca6b", "gender": 2, "id": 154169, "name": "Sean Anthon ...", "order": 28}, {"cast_id": 49, "character": "Miles", "credit_id": "52fe48009251416c750aca6c", "gender": 2, "id": 154170, "name": "Sean Anthon ...", "order": 29}, {"cast_id": 50, "character": "Miles", "credit_id": "52fe48009251416c750aca6d", "gender": 2, "id": 154171, "name": "Sean Anthon ...", "order": 30}, {"cast_id": 51, "character": "Miles", "credit_id": "52fe48009251416c750aca6e", "gender": 2, "id": 154172, "name": "Sean Anthon ...", "order": 31}, {"cast_id": 52, "character": "Miles", "credit_id": "52fe48009251416c750aca6f", "gender": 2, "id": 154173, "name": "Sean Anthon ...", "order": 32}, {"cast_id": 53, "character": "Miles", "credit_id": "52fe48009251416c750aca70", "gender": 2, "id": 154174, "name": "Sean Anthon ...", "order": 33}, {"cast_id": 54, "character": "Miles", "credit_id": "52fe48009251416c750aca71", "gender": 2, "id": 154175, "name": "Sean Anthon ...", "order": 34}, {"cast_id": 55, "character": "Miles", "credit_id": "52fe48009251416c750aca72", "gender": 2, "id": 154176, "name": "Sean Anthon ...", "order": 35}, {"cast_id": 56, "character": "Miles", "credit_id": "52fe48009251416c750aca73", "gender": 2, "id": 154177, "name": "Sean Anthon ...", "order": 36}, {"cast_id": 57, "character": "Miles", "credit_id": "52fe48009251416c750aca74", "gender": 2, "id": 154178, "name": "Sean Anthon ...", "order": 37}, {"cast_id": 58, "character": "Miles", "credit_id": "52fe48009251416c750aca75", "gender": 2, "id": 154179, "name": "Sean Anthon ...", "order": 38}, {"cast_id": 59, "character": "Miles", "credit_id": "52fe48009251416c750aca76", "gender": 2, "id": 154180, "name": "Sean Anthon ...", "order": 39}, {"cast_id": 60, "character": "Miles", "credit_id": "52fe48009251416c750aca77", "gender": 2, "id": 154181, "name": "Sean Anthon ...", "order": 40}, {"cast_id": 61, "character": "Miles", "credit_id": "52fe48009251416c750aca78", "gender": 2, "id": 154182, "name": "Sean Anthon ...", "order": 41}, {"cast_id": 62, "character": "Miles", "credit_id": "52fe48009251416c750aca79", "gender": 2, "id": 154183, "name": "Sean Anthon ...", "order": 42}, {"cast_id": 63, "character": "Miles", "credit_id": "52fe48009251416c750aca7a", "gender": 2, "id": 154184, "name": "Sean Anthon ...", "order": 43}, {"cast_id": 64, "character": "Miles", "credit_id": "52fe48009251416c750aca7b", "gender": 2, "id": 154185, "name": "Sean Anthon ...", "order": 44}, {"cast_id": 65, "character": "Miles", "credit_id": "52fe48009251416c750aca7c", "gender": 2, "id": 154186, "name": "Sean Anthon ...", "order": 45}, {"cast_id": 66, "character": "Miles", "credit_id": "52fe48009251416c750aca7d", "gender": 2, "id": 154187, "name": "Sean Anthon ...", "order": 46}, {"cast_id": 67, "character": "Miles", "credit_id": "52fe48009251416c750aca7e", "gender": 2, "id": 154188, "name": "Sean Anthon ...", "order": 47}, {"cast_id": 68, "character": "Miles", "credit_id": "52fe48009251416c750aca7f", "gender": 2, "id": 154189, "name": "Sean Anthon ...", "order": 48}, {"cast_id": 69, "character": "Miles", "credit_id": "52fe48009251416c750aca80", "gender": 2, "id": 154190, "name": "Sean Anthon ...", "order": 49}, {"cast_id": 70, "character": "Miles", "credit_id": "52fe48009251416c750aca81", "gender": 2, "id": 154191, "name": "Sean Anthon ...", "order": 50}, {"cast_id": 71, "character": "Miles", "credit_id": "52fe48009251416c750aca82", "gender": 2, "id": 154192, "name": "Sean Anthon ...", "order": 51}, {"cast_id": 72, "character": "Miles", "credit_id": "52fe48009251416c750aca83", "gender": 2, "id": 154193, "name": "Sean Anthon ...", "order": 52}, {"cast_id": 73, "character": "Miles", "credit_id": "52fe48009251416c750aca84", "gender": 2, "id": 154194, "name": "Sean Anthon ...", "order": 53}, {"cast_id": 74, "character": "Miles", "credit_id": "52fe48009251416c750aca85", "gender": 2, "id": 154195, "name": "Sean Anthon ...", "order": 54}, {"cast_id": 75, "character": "Miles", "credit_id": "52fe48009251416c750aca86", "gender": 2, "id": 154196, "name": "Sean Anthon ...", "order": 55}, {"cast_id": 76, "character": "Miles", "credit_id": "52fe48009251416c750aca87", "gender": 2, "id": 154197, "name": "Sean Anthon ...", "order": 56}, {"cast_id": 77, "character": "Miles", "credit_id": "52fe48009251416c750aca88", "gender": 2, "id": 154198, "name": "Sean Anthon ...", "order": 57}, {"cast_id": 78, "character": "Miles", "credit_id": "52fe48009251416c750aca89", "gender": 2, "id": 154199, "name": "Sean Anthon ...", "order": 58}, {"cast_id": 79, "character": "Miles", "credit_id": "52fe48009251416c750aca8a", "gender": 2, "id": 154200, "name": "Sean Anthon ...", "order": 59}, {"cast_id": 80, "character": "Miles", "credit_id": "52fe48009251416c750aca8b", "gender": 2, "id": 154201, "name": "Sean Anthon ...", "order": 60}, {"cast_id": 81, "character": "Miles", "credit_id": "52fe48009251416c750aca8c", "gender": 2, "id": 154202, "name": "Sean Anthon ...", "order": 61}, {"cast_id": 82, "character": "Miles", "credit_id": "52fe48009251416c750aca8d", "gender": 2, "id": 154203, "name": "Sean Anthon ...", "order": 62}, {"cast_id": 83, "character": "Miles", "credit_id": "52fe48009251416c750aca8e", "gender": 2, "id": 154204, "name": "Sean Anthon ...", "order": 63}, {"cast_id": 84, "character": "Miles", "credit_id": "52fe48009251416c750aca8f", "gender": 2, "id": 154205, "name": "Sean Anthon ...", "order": 64}, {"cast_id": 85, "character": "Miles", "credit_id": "52fe48009251416c750aca90", "gender": 2, "id": 154206, "name": "Sean Anthon ...", "order": 65}, {"cast_id": 86, "character": "Miles", "credit_id": "52fe48009251416c750aca91", "gender": 2, "id": 154207, "name": "Sean Anthon ...", "order": 66}, {"cast_id": 87, "character": "Miles", "credit_id": "52fe48009251416c750aca92", "gender": 2, "id": 154208, "name": "Sean Anthon ...", "order": 67}, {"cast_id": 88, "character": "Miles", "credit_id": "52fe48009251416c750aca93", "gender": 2, "id": 154209, "name": "Sean Anthon ...", "order": 68}, {"cast_id": 89, "character": "Miles", "credit_id": "52fe48009251416c750aca94", "gender": 2, "id": 154210, "name": "Sean Anthon ...", "order": 69}, {"cast_id": 90, "character": "Miles", "credit_id": "52fe48009251416c750aca95", "gender": 2, "id": 154211, "name": "Sean Anthon ...", "order": 70}, {"cast_id": 91, "character": "Miles", "credit_id": "52fe48009251416c750aca96", "gender": 2, "id": 154212, "name": "Sean Anthon ...", "order": 71}, {"cast_id": 92, "character": "Miles", "credit_id": "52fe48009251416c750aca97", "gender": 2, "id": 154213, "name": "Sean Anthon ...", "order": 72}, {"cast_id": 93, "character": "Miles", "credit_id": "52fe48009251416c750aca98", "gender": 2, "id": 154214, "name": "Sean Anthon ...", "order": 73}, {"cast_id": 94, "character": "Miles", "credit_id": "52fe48009251416c750aca99", "gender": 2, "id": 154215, "name": "Sean Anthon ...", "order": 74}, {"cast_id": 95, "character": "Miles", "credit_id": "52fe48009251416c750aca9a", "gender": 2, "id": 154216, "name": "Sean Anthon ...", "order": 75}, {"cast_id": 96, "character": "Miles", "credit_id": "52fe48009251416c750aca9b", "gender": 2, "id": 154217, "name": "Sean Anthon ...", "order": 76}, {"cast_id": 97, "character": "Miles", "credit_id": "52fe48009251416c750aca9c", "gender": 2, "id": 154218, "name": "Sean Anthon ...", "order": 77}, {"cast_id": 98, "character": "Miles", "credit_id": "52fe48009251416c750aca9d", "gender": 2, "id": 154219, "name": "Sean Anthon ...", "order": 78}, {"cast_id": 99, "character": "Miles", "credit_id": "52fe48009251416c750aca9e", "gender": 2, "id": 154220, "name": "Sean Anthon ...", "order": 79}, {"cast_id": 100, "character": "Miles", "credit_id": "52fe48009251416c750aca9f", "gender": 2, "id": 154221, "name": "Sean Anthon ...", "order": 80}, {"cast_id": 101, "character": "Miles", "credit_id": "52fe48009251416c750acaa0", "gender": 2, "id": 154222, "name": "Sean Anthon ...", "order": 81}, {"cast_id": 102, "character": "Miles", "credit_id": "52fe48009251416c750acaa1", "gender": 2, "id": 154223, "name": "Sean Anthon ...", "order": 82}, {"cast_id": 103, "character": "Miles", "credit_id": "52fe48009251416c750acaa2", "gender": 2, "id": 154224, "name": "Sean Anthon ...", "order": 83}, {"cast_id": 104, "character": "Miles", "credit_id": "52fe48009251416c750acaa3", "gender": 2, "id": 154225, "name": "Sean Anthon ...", "order": 84}, {"cast_id": 105, "character": "Miles", "credit_id": "52fe48009251416c750acaa4", "gender": 2, "id": 154226, "name": "Sean Anthon ...", "order": 85}, {"cast_id": 106, "character": "Miles", "credit_id": "52fe48009251416c750acaa5", "gender": 2, "id": 154227, "name": "Sean Anthon ...", "order": 86}, {"cast_id": 107, "character": "Miles", "credit_id": "52fe48009251416c750acaa6", "gender": 2, "id": 154228, "name": "Sean Anthon ...", "order": 87}, {"cast_id": 108, "character": "Miles", "credit_id": "52fe48009251416c750acaa7", "gender": 2, "id": 154229, "name": "Sean Anthon ...", "order": 88}, {"cast_id": 109, "character": "Miles", "credit_id": "52fe48009251416c750acaa8", "gender": 2, "id": 154230, "name": "Sean Anthon ...", "order": 89}, {"cast_id": 110, "character": "Miles", "credit_id": "52fe48009251416c750acaa9", "gender": 2, "id": 154231, "name": "Sean Anthon ...", "order": 90}, {"cast_id": 111, "character": "Miles", "credit_id": "52fe48009251416c750acaaa", "gender": 2, "id": 154232, "name": "Sean Anthon ...", "order": 91}, {"cast_id": 112, "character": "Miles", "credit_id": "52fe48009251416c750acaaa1", "gender": 2, "id": 154233, "name": "Sean Anthon ...", "order": 92}, {"cast_id": 113, "character": "Miles", "credit_id": "52fe48009251416c750acaaa2", "gender": 2, "id": 154234, "name": "Sean Anthon ...", "order": 93}, {"cast_id": 114, "character": "Miles", "credit_id": "52fe48009251416c750acaaa3", "gender": 2, "id": 154235, "name": "Sean Anthon ...", "order": 94}, {"cast_id": 115, "character": "Miles", "credit_id": "52fe48009251416c750acaaa4", "gender": 2, "id": 154236, "name": "Sean Anthon ...", "order": 95}, {"cast_id": 116, "character": "Miles", "credit_id": "52fe48009251416c750acaaa5", "gender": 2, "id": 154237, "name": "Sean Anthon ...", "order": 96}, {"cast_id": 117, "character": "Miles", "credit_id": "52fe48009251416c750acaaa6", "gender": 2, "id": 154238, "name": "Sean Anthon ...", "order": 97}, {"cast_id": 118, "character": "Miles", "credit_id": "52fe48009251416c750acaaa7", "gender": 2, "id": 154239, "name": "Sean Anthon ...", "order": 98}, {"cast_id": 119, "character": "Miles", "credit_id": "52fe48009251416c750acaaa8", "gender": 2, "id": 154240, "name": "Sean Anthon ...", "order": 99}, {"cast_id": 120, "character": "Miles", "credit_id": "52fe48009251416c750acaaa9", "gender": 2, "id": 154241, "name": "Sean Anthon ...", "order": 100}, {"cast_id": 121, "character": "Miles", "credit_id": "52fe48009251416c750acaab", "gender": 2, "id": 154242, "name": "Sean Anthon ...", "order": 101}, {"cast_id": 122, "character": "Miles", "credit_id": "52fe48009251416c750acaab1", "gender": 2, "id": 154243, "name": "Sean Anthon ...", "order": 102}, {"cast_id": 123, "character": "Miles", "credit_id": "52fe48009251416c750acaab2", "gender": 2, "id": 154244, "name": "Sean Anthon ...", "order": 103}, {"cast_id": 124, "character": "Miles", "credit_id": "52fe48009251416c750acaab3", "gender": 2, "id": 154245, "name": "Sean Anthon ...", "order": 104}, {"cast_id": 125, "character": "Miles", "credit_id": "52fe48009251416c750acaab4", "gender": 2, "id": 154246, "name": "Sean Anthon ...", "order": 105}, {"cast_id": 126, "character": "Miles", "credit_id": "52fe48009251416c750acaab5", "gender": 2, "id": 154247, "name": "Sean Anthon ...", "order": 106}, {"cast_id": 127, "character": "Miles", "credit_id": "52fe48009251416c750acaab6", "gender": 2, "id": 154248, "name": "Sean Anthon ...", "order": 107}, {"cast_id": 128, "character": "Miles", "credit_id": "52fe48009251416c750acaab7", "gender": 2, "id": 154249, "name": "Sean Anthon ...", "order": 108}, {"cast_id": 129, "character": "Miles", "credit_id": "52fe48009251416c750acaab8", "gender": 2, "id": 154250, "name": "Sean Anthon ...", "order": 109}, {"cast_id": 130, "character": "Miles", "credit_id": "52fe48009251416c750acaab9", "gender": 2, "id": 154251, "name": "Sean Anthon ...", "order": 110}, {"cast_id": 131, "character": "Miles", "credit_id": "52fe48009251416c750acaaba", "gender": 2, "id": 154252, "name": "Sean Anthon ...", "order": 111}, {"cast_id": 132, "character": "Miles", "credit_id": "52fe48009251416c750acaaba1", "gender": 2, "id": 154253, "name": "Sean Anthon ...", "order": 112}, {"cast_id": 133, "character": "Miles", "credit_id": "52fe48009251416c750acaaba2", "gender": 2, "id": 154254, "name": "Sean Anthon ...", "order": 113}, {"cast_id": 134, "character": "Miles", "credit_id": "52fe48009251416c750acaaba3", "gender": 2, "id": 154255, "name": "Sean Anthon ...", "order": 114}, {"cast_id": 135, "character": "Miles", "credit_id": "52fe48009251416c750acaaba4", "gender": 2, "id": 154256, "name": "Sean Anthon ...", "order": 115}, {"cast_id": 136, "character": "Miles", "credit_id": "52fe48009251416c750acaaba5", "gender": 2, "id": 154257, "name": "Sean Anthon ...", "order": 116}, {"cast_id": 137, "character": "Miles", "credit_id": "52fe48009251416c750acaaba6", "gender": 2, "id": 154258, "name": "Sean Anthon ...", "order": 117}, {"cast_id": 138, "character": "Miles", "credit_id": "52fe48009251416c750acaaba7", "gender": 2, "id": 154259, "name": "Sean Anthon ...", "order": 118}, {"cast_id": 139, "character": "Miles", "credit_id": "52fe48009251416c750acaaba8", "gender": 2, "id": 154260, "name": "Sean Anthon ...", "order": 119}, {"cast_id": 140, "character": "Miles", "credit_id": "52fe48009251416c750acaaba9", "gender": 2, "id": 154261, "name": "Sean Anthon ...", "order": 120}, {"cast_id": 141, "character": "Miles", "credit_id": "52fe48009251416c750acaaba", "gender": 2, "id": 154262, "name": "Sean Anthon ...", "order": 121}, {"cast_id": 142, "character": "Miles", "credit_id": "52fe48009251416c750acaaba1", "gender": 2, "id": 154263, "name": "Sean Anthon ...", "order": 122}, {"cast_id": 143, "character": "Miles", "credit_id": "52fe48009251416c750acaaba2", "gender": 2, "id": 154264, "name": "Sean Anthon ...", "order": 123}, {"cast_id": 144, "character": "Miles", "credit_id": "52fe48009251416c750acaaba3", "gender": 2, "id": 154265, "name": "Sean Anthon ...", "order": 124}, {"cast_id": 145, "character": "Miles", "credit_id": "52fe48009251416c750acaaba4", "gender": 2, "id": 154266, "name": "Sean Anthon ...", "order": 125}, {"cast_id": 146, "character": "Miles", "credit_id": "52fe48009251416c750acaaba5", "gender": 2, "id": 154267, "name": "Sean Anthon ...", "order": 126}, {"cast_id": 147, "character": "Miles", "credit_id": "52fe48009251416c750acaaba6", "gender": 2, "id": 154268, "name": "Sean Anthon ...", "order": 127}, {"cast_id": 148, "character": "Miles", "credit_id": "52fe48009251416c750acaaba7", "gender": 2, "id": 154269, "name": "Sean Anthon ...", "order": 128}, {"cast_id": 149, "character": "Miles", "credit_id": "52fe48009251416c750acaaba8", "gender": 2, "id": 154270, "name": "Sean Anthon ...", "order": 129}, {"cast_id": 150, "character": "Miles", "credit_id": "52fe48009251416c750acaaba9", "gender": 2, "id": 154271, "name": "Sean Anthon ...", "order": 130}, {"cast_id": 151, "character": "Miles", "credit_id": "52fe48009251416c750acaaba", "gender": 2, "id": 154272, "name": "Sean Anthon ...", "order": 131}, {"cast_id": 152, "character": "Miles", "credit_id": "52fe4800925141

In [225]:

```
# Updating the cast column
movies['cast'] = movies['cast'].apply(convert4)

/var/folders/yb/v5f667y93j393zcsz2t_26qw000gn/T/ipykernel_13438/2674671810.py:2: SettingWithCopyWarnin
g:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
movies['cast'] = movies['cast'].apply(convert4)
```

In [226]:

```
# Checking for the first value
movies['cast'][0]
```

Out[226]:

```
['Sam Worthington', 'Zoe Saldana', 'Sigourney Weaver', 'Stephen Lang']
```

- In the crew section we only need to derive the Director name for "job" value

In [227]:

```
movies['crew'][0]
```

Out[227]:

```
{'credit_id': '52fe48009251416c750aca23', 'department': 'Editing', 'gender': 0, 'id': 1721, 'job': 'Editor', 'name': 'Stephen E. Rivkin'}, {'credit_id': '539c47ecc3a36810e3001f87', 'department': 'Art', 'gender': 2, 'id': 496, 'job': 'Production Design', 'name': 'Rick Carter'}, {'credit_id': '54491c89c3a3680fb4001cf7', 'department': 'Sound', 'gender': 0, 'id': 900, 'job': 'Sound Designer', 'name': 'Christopher Boyes'}, {'credit_id': '54491cb70e0a267480001bd0', 'department': 'Sound', 'gender': 0, 'id': 900, 'job': 'Supervising Sound Editor', 'name': 'Christopher Boyes'}, {'credit_id': '539c4a4cc3a36810c9002101', 'department': 'Production', 'gender': 1, 'id': 1262, 'job': 'Casting', 'name': 'Mali Finn'}, {'credit_id': '5544ee3b925141499f0008fc', 'department': 'Sound', 'gender': 2, 'id': 1729, 'job': 'Original Music Composer', 'name': 'James Horner'}, {'credit_id': '52fe48009251416c750ac9c3', 'department': 'Directing', 'gender': 2, 'id': 2710, 'job': 'Director', 'name': 'James Cameron'}, {'credit_id': '52fe48009251416c750ac9d9', 'department': 'Writing', 'gender': 2, 'id': 2710, 'job': 'Writer', 'name': 'James Cameron'}, {'credit_id': '52fe48009251416c750aca17', 'department': 'Editing', 'gender': 2, 'id': 2710, 'job': 'Editor', 'name': 'James Cameron'}, {'credit_id': '52fe48009251416c750aca29', 'department': 'Production', 'gender': 2, 'id': 2710, 'job': 'Producer', 'name': 'James Cameron'}, {'credit_id': '52fe48009251416c750aca3f', 'department': 'Writing', 'gender': 2, 'id': 2710, 'job': 'Screenplay', 'name': 'James Cameron'}, {'credit_id': '539c4987c3a36810ba0021a4', 'department': 'Art', 'gender': 2, 'id': 7236, 'job': 'Art Direction', 'name': 'Andrew Menzies'}, {'credit_id': '549598c3c3a3686ae9004383', 'department': 'Visual Effects', 'gender': 0, 'id': 6690, 'job': 'Visual Effects Producer', 'name': 'Jill Brooks'}, {'credit_id': '52fe48009251416c750aca4b', 'department': 'Production', 'gender': 2, 'id': 6217, 'job': 'Casting', 'name': 'Morgan Siskind'}, {'credit_id': '570b6f6100f1417d'}
```

In [228]:

```
# Only need the director name
def get_director(obj):
    L = []
    for i in ast.literal_eval(obj):
        if i['job'] == 'Director':
            L.append(i['name'])
            break
    return L
```

In [229]:

```
# Updating Crew
movies['crew'] = movies['crew'].apply(get_director)

/var/folders/yb/v5f667y93j393zcsz2t_26qw000gn/T/ipykernel_13438/4113385147.py:2: SettingWithCopyWarnin
g:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
movies['crew'] = movies['crew'].apply(get_director)
```

In [230]:

```
# Checking for the crew
movies.head()
```

Out[230]:

	id	title	genres	keywords	overview	cast	crew
0	19995	Avatar	[Action, Adventure, Fantasy, Science Fiction]	[culture clash, future, space war, space colon...	In the 22nd century, a paraplegic Marine is di...	[Sam Worthington, Zoe Saldana, Sigourney Weave...	[James Cameron]
1	285	Pirates of the Caribbean: At World's End	[Adventure, Fantasy, Action]	[ocean, drug abuse, exotic island, east india ...	Captain Barbossa, long believed to be dead, ha...	[Johnny Depp, Orlando Bloom, Keira Knightley, ...	[Gore Verbinski]
2	206647	Spectre	[Action, Adventure, Crime]	[spy, based on novel, secret agent, sequel, mi...	A cryptic message from Bond's past sends him o...	[Daniel Craig, Christoph Waltz, Léa Seydoux, R...	[Sam Mendes]
3	49026	The Dark Knight Rises	[Action, Crime, Drama, Thriller]	[dc comics, crime fighter, terrorist, secret i...	Following the death of District Attorney Harve...	[Christian Bale, Michael Caine, Gary Oldman, A...	[Christopher Nolan]
4	49529	John Carter	[Action, Adventure, Science Fiction]	[based on novel, mars, medallion, space travel...	John Carter is a war-weary, former military ca...	[Taylor Kitsch, Lynn Collins, Samantha Morton,...	[Andrew Stanton]

In [231]:

```
# Converting text to the list
movies['overview'] = movies['overview'].apply(lambda x:x.split())

/var/folders/yb/v5f667y93j393zcsz2t_26qw000gn/T/ipykernel_13438/1971172765.py:2: SettingWithCopyWarnin
g:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
movies['overview'] = movies['overview'].apply(lambda x:x.split())
```

In [232]:

```
movies.head()
```

Out[232]:

	id	title	genres	keywords	overview	cast	crew
0	19995	Avatar	[Action, Adventure, Fantasy, Science Fiction]	[culture clash, future, space war, space colon...	[In, the, 22nd, century,, a, paraplegic, Marin...	[Sam Worthington, Zoe Saldana, Sigourney Weave...	[James Cameron]
1	285	Pirates of the Caribbean: At World's End	[Adventure, Fantasy, Action]	[ocean, drug abuse, exotic island, east india ...	[Captain, Barbossa,, long, believed, to, be, d...	[Johnny Depp, Orlando Bloom, Keira Knightley, ...	[Gore Verbinski]
2	206647	Spectre	[Action, Adventure, Crime]	[spy, based on novel, secret agent, sequel, mi...	[A, cryptic, message, from, Bond's, past, send...	[Daniel Craig, Christoph Waltz, Léa Seydoux, R...	[Sam Mendes]
3	49026	The Dark Knight Rises	[Action, Crime, Drama, Thriller]	[dc comics, crime fighter, terrorist, secret i...	[Following, the, death, of, District, Attorney...	[Christian Bale, Michael Caine, Gary Oldman, A...	[Christopher Nolan]
4	49529	John Carter	[Action, Adventure, Science Fiction]	[based on novel, mars, medallion, space travel...	[John, Carter, is, a, war-weary,, former, mili...	[Taylor Kitsch, Lynn Collins, Samantha Morton,...	[Andrew Stanton]

In [233]:

```
# Taking care of the blank spaces in the columns
movies['genres'] = movies['genres'].apply(lambda x:[i.replace(" ", "")for i in x])
movies['keywords'] = movies['keywords'].apply(lambda x:[i.replace(" ", "")for i in x])
movies['cast'] = movies['cast'].apply(lambda x:[i.replace(" ", "")for i in x])
movies['crew'] = movies['crew'].apply(lambda x:[i.replace(" ", "")for i in x])

/var/folders/yb/v5f667y93j393zcsz2t_26qw0000gn/T/ipykernel_13438/2267989278.py:2: SettingWithCopyWarnin
g:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.
html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.h
tml#returning-a-view-versus-a-copy)
movies['genres'] = movies['genres'].apply(lambda x:[i.replace(" ", "")for i in x])
/var/folders/yb/v5f667y93j393zcsz2t_26qw0000gn/T/ipykernel_13438/2267989278.py:3: SettingWithCopyWarnin
g:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.
html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.h
tml#returning-a-view-versus-a-copy)
movies['keywords'] = movies['keywords'].apply(lambda x:[i.replace(" ", "")for i in x])
/var/folders/yb/v5f667y93j393zcsz2t_26qw0000gn/T/ipykernel_13438/2267989278.py:4: SettingWithCopyWarnin
g:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.
html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.h
tml#returning-a-view-versus-a-copy)
movies['cast'] = movies['cast'].apply(lambda x:[i.replace(" ", "")for i in x])
/var/folders/yb/v5f667y93j393zcsz2t_26qw0000gn/T/ipykernel_13438/2267989278.py:5: SettingWithCopyWarnin
g:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.
html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.h
tml#returning-a-view-versus-a-copy)
movies['crew'] = movies['crew'].apply(lambda x:[i.replace(" ", "")for i in x])
```

In [234]:

movies.head()

Out[234]:

	id	title	genres	keywords	overview	cast	crew
0	19995	Avatar	[Action, Adventure, Fantasy, ScienceFiction]	[cultureclash, future, spacewar, spacecolony, ...	[In, the, 22nd, century,, a, paraplegic, Marin...	[SamWorthington, ZoeSaldana, SigourneyWeaver, ...	[JamesCameron]
1	285	Pirates of the Caribbean: At World's End	[Adventure, Fantasy, Action]	[ocean, drugabuse, exoticisland, eastindiatrad...	[Captain, Barbossa,, long, believed, to, be, d...	[JohnnyDepp, OrlandoBloom, KeiraKnightley, Ste...	[GoreVerbinski]
2	206647	Spectre	[Action, Adventure, Crime]	[spy, basedonnovel, secretagent, sequel, mi6, ...	[A, cryptic, message, from, Bond's, past, send...	[DanielCraig, ChristophWaltz, LéaSeydoux, Ralp...	[SamMendes]
3	49026	The Dark Knight Rises	[Action, Crime, Drama, Thriller]	[dccomics, crimefighter, terrorist, secretiden...	[Following, the, death, of, District, Attorney...	[ChristianBale, MichaelCaine, GaryOldman, Anne...	[ChristopherNolan]
4	49529	John Carter	[Action, Adventure, ScienceFiction]	[basedonnovel, mars, medallion, spacetravel, p...	[John, Carter, is, a, war-weary,, former, mili...	[TaylorKitsch, LynnCollins, SamanthaMorton, Wi...	[AndrewStanton]

In [235]:

```
# Combining all the columns into the one columns
movies['tags'] = movies['overview'] + movies['genres'] + movies['keywords'] + movies['cast'] + movies['crew']

/var/folders/yb/v5f667y93j393zcsz2t_26qw0000gn/T/ipykernel_13438/2167919987.py:2: SettingWithCopyWarnin
g:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.
html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.h
tml#returning-a-view-versus-a-copy)
movies['tags'] = movies['overview'] + movies['genres'] + movies['keywords'] + movies['cast'] + movie
s['crew']
```

In [236]:

movies.head()

Out[236]:

	id	title	genres	keywords	overview	cast	crew	tags
0	19995	Avatar	[Action, Adventure, Fantasy, ScienceFiction]	[cultureclash, future, spacewar, spacecolony, ...	[In, the, 22nd, century,, a, paraplegic, Marin...	[SamWorthington, ZoeSaldana, SigourneyWeaver, ...	[JamesCameron]	[In, the, 22nd, century,, a, paraplegic, Marin...
1	285	Pirates of the Caribbean: At World's End	[Adventure, Fantasy, Action]	[ocean, drugabuse, exoticisland, eastindiatrad...	[Captain, Barbossa,, long, believed, to, be, d...	[JohnnyDepp, OrlandoBloom, KeiraKnightley, Ste...	[GoreVerbinski]	[Captain, Barbossa,, long, believed, to, be, d...
2	206647	Spectre	[Action, Adventure, Crime]	[spy, basedonnovel, secretagent, sequel, mi6, ...	[A, cryptic, message, from, Bond's, past, send...	[DanielCraig, ChristophWaltz, LéaSeydoux, Ralp...	[SamMendes]	[A, cryptic, message, from, Bond's, past, send...
3	49026	The Dark Knight Rises	[Action, Crime, Drama, Thriller]	[dccomics, crimefighter, terrorist, secretiden...	[Following, the, death, of, District, Attorney...	[ChristianBale, MichaelCaine, GaryOldman, Anne...	[ChristopherNolan]	[Following, the, death, of, District, Attorney...
4	49529	John Carter	[Action, Adventure, ScienceFiction]	[basedonnovel, mars, medallion, ...	[John, Carter, is, a, war-weary,, former, ...	[TaylorKitsch, LynnCollins, ...	[AndrewStanton]	[John, Carter, is, a, war-weary...

Science Fiction
Space
Travel
Future
Space Colony
Society
Space Travel
Futuristic
Romance
Space Alien
Tribe
Alien Planet
CGI
Marine
Soldier
Battle
Love Affair
Antiwar
Power Relations
Mind and Soul
3D
Sam Worthington
Zoe Saldana
Sigourney Weaver
Stephen Lang
James Cameron

Avatar
Pirates of the Caribbean: At World's End
Spectre
The Dark Knight Rises
John Carter

mili...
Samantha Morton, Willem Dafoe, Michael Fassbender, Matt Smith, Benedict Cumberbatch, ...
former, mili...

In [237]:

Creating final Dataset
final_df = movies[['id','title','tags']]

In [238]:

final_df.head()

Out[238]:

	id	title	tags
0	19995	Avatar	[In, the, 22nd, century,, a, paraplegic, Marin...
1	285	Pirates of the Caribbean: At World's End	[Captain, Barbossa,, long, believed, to, be, d...
2	206647	Spectre	[A, cryptic, message, from, Bond's, past, send...
3	49026	The Dark Knight Rises	[Following, the, death, of, District, Attorney...
4	49529	John Carter	[John, Carter, is, a, war-weary,, former, mili...

In [239]:

final_df['tags'].apply(lambda x:" ".join(x))

Out[239]:

0 In the 22nd century, a paraplegic Marine is di...
1 Captain Barbossa, long believed to be dead, ha...
2 A cryptic message from Bond's past sends him o...
3 Following the death of District Attorney Harve...
4 John Carter is a war-weary, former military ca...
...
4804 El Mariachi just wants to play his guitar and ...
4805 A newlywed couple's honeymoon is upended by th...
4806 "Signed, Sealed, Delivered" introduces a dedic...
4807 When ambitious New York attorney Sam is sent t...
4808 Ever since the second grade when he first saw ...
Name: tags, Length: 4806, dtype: object

In [240]:

final_df['tags'] = final_df['tags'].apply(lambda x:" ".join(x))

/var/folders/yb/v5f667y93j393zcsz2t_26qw0000gn/T/ipykernel_13438/2951742378.py:1: SettingWithCopyWarnin
g:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy ([https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.h
tml#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy))
final_df['tags'] = final_df['tags'].apply(lambda x:" ".join(x))

In [241]:

final_df.head()

Out[241]:

	id	title	tags
0	19995	Avatar	In the 22nd century, a paraplegic Marine is di...
1	285	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...
2	206647	Spectre	A cryptic message from Bond's past sends him o...
3	49026	The Dark Knight Rises	Following the death of District Attorney Harve...
4	49529	John Carter	John Carter is a war-weary, former military ca...

In [242]:

Checking for the 0th index
final_df['tags'][0]

Out[242]:

'In the 22nd century, a paraplegic Marine is dispatched to the moon Pandora on a unique mission, but be
comes torn between following orders and protecting an alien civilization. Action Adventure Fantasy Scie
nce Fiction culture clash future space war space colony society space travel futuristic romance space alien
tribe alien planet cgi marine soldier battle love affair antiwar power relations mind and soul 3d Sam Worthing
ton Zoe Saldana Sigourney Weaver Stephen Lang James Cameron'

Converting into Lower Case

In [243]:

final_df['tags'] = final_df['tags'].apply(lambda x:x.lower())

/var/folders/yb/v5f667y93j393zcsz2t_26qw0000gn/T/ipykernel_13438/2213019244.py:1: SettingWithCopyWarnin
g:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy ([https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.h
tml#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy))
final_df['tags'] = final_df['tags'].apply(lambda x:x.lower())

In [244]:

final_df.head()

Out[244]:

	id	title	tags
0	19995	Avatar	in the 22nd century, a paraplegic marine is di...
1	285	Pirates of the Caribbean: At World's End	captain barbossa, long believed to be dead, ha...
2	206647	Spectre	a cryptic message from bond's past sends him o...
3	49026	The Dark Knight Rises	following the death of district attorney harve...
4	49529	John Carter	john carter is a war-weary, former military ca...

In [245]:

!pip install nltk

Requirement already satisfied: nltk in /Users/varundesai/anaconda3/lib/python3.10/site-packages (3.7)
Requirement already satisfied: click in /Users/varundesai/anaconda3/lib/python3.10/site-packages (from
nltk) (8.0.4)
Requirement already satisfied: joblib in /Users/varundesai/anaconda3/lib/python3.10/site-packages (from
nltk) (1.1.1)
Requirement already satisfied: regex<=2021.8.3 in /Users/varundesai/anaconda3/lib/python3.10/site-packa
ges (from nltk) (2022.7.9)
Requirement already satisfied: tqdm in /Users/varundesai/anaconda3/lib/python3.10/site-packages (from n
ltk) (4.64.1)

In [246]:

Importing Libraries
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem import PorterStemmer
import string

```
In [247]: # Removing Stopwords, punctuation marks, applying stemming operation
def pre_process(text):
    tokens = word_tokenize(text)
    tokens = [word for word in tokens if word.isalnum()]

    stop_words = set(stopwords.words('english'))
    tokens = [word for word in tokens if word.lower() not in stop_words]

    stemmer = PorterStemmer()
    tokens = [stemmer.stem(word) for word in tokens]

    processed_text = ' '.join(tokens)

    return processed_text
```

```
In [248]: final_df['tags'] = final_df['tags'].apply(pre_process)

/var/folders/yb/v5f667y93j393zcsz2t_26qw0000gn/T/ipykernel_13438/3906571443.py:1: SettingWithCopyWarnin
g:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
final_df['tags'] = final_df['tags'].apply(pre_process)
```

```
In [249]: final_df.head()
```

Out[249]:

	id		title	tags
0	19995		Avatar	22nd centuri parapleg marin dispatch moon pand...
1	285	Pirates of the Caribbean: At World's End		captain barbossa long believ dead come back li...
2	206647		Spectre	cryptic messag bond past send trail uncov sini...
3	49026	The Dark Knight Rises		follow death district attorney harvey dent bat...
4	49529		John Carter	john carter former militari captain inexplic t...

```
In [250]: final_df['tags'][0]
```

Out[250]: '22nd centuri parapleg marin dispatch moon pandora uniqu mission becom torn follow order protect alien civil action adventur fantasi sciencefict cultureclash futur spacewar spacecoloni societi spacetravel futurist romanc space alien tribe alienplanet cgi marin soldier battl loveaffair antiwar powerrel mindan dsoul 3d samworthington zoesaldana sigourneyweav stephenlang jamescameron'

Training

- Different techniques such as Bag of Words, IF-IDF, Word2Vec could be used to perform word embedding.

```
In [251]: from sklearn.feature_extraction.text import TfidfVectorizer
# Max_features are 5000
tfidf_vectorizer = TfidfVectorizer(max_features=5000)

# Creating Matrix for the vector values
tfidf_matrix = tfidf_vectorizer.fit_transform(final_df['tags'])

# Converting into the numpy array
tfidf_df = pd.DataFrame(tfidf_matrix.toarray(), columns=tfidf_vectorizer.get_feature_names_out())
```

```
In [252]: tfidf_df.head()
```

Out[252]:

	007	10	100	11	12	13	15	16	17	17th	...	zhangyim	zhangziyi	zion	zoe	zoesaldana	zombi	zombieapocalyps	zone	zoo	zoc
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	0.207132	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	0.0	0.0	0.0	0.000000	0.0		0.0	0.0	0.0
2	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	0.000000	0.0		0.0	0.0	0.0
3	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	0.000000	0.0		0.0	0.0	0.0
4	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	0.000000	0.0		0.0	0.0	0.0

5 rows × 5000 columns

- Since 5000 features have been selected now, we can calculate the Cosine Similarity, in order to find the similar vectors values

```
In [253]: from sklearn.metrics.pairwise import cosine_similarity
cosine_similarity(tfidf_df).shape
```

Out[253]: (4806, 4806)

```
In [254]: similarity = cosine_similarity(tfidf_df)
```

```
In [255]: similarity[0]
```

Out[255]: array([1. , 0.02218662, 0.02952472, ..., 0.04533449, 0.00551997,
0.])

Insight:

- 0th index value is same as itself therefore values is 1
- 0th index is similar to the 1st index item with value of 0.02218..
- Higher the number, more similar it is to the item

```
In [256]: # Creating Recommendation Function
def recommend(movie):
    movie_index = final_df[final_df['title'] == movie].index[0]
    distances = similarity[movie_index]
    movies_list = sorted(list(enumerate(distances)),reverse=True,key = lambda x:x[1])[1:6]

    for i in movies_list:
        print(i[0])
```

```
In [257]: recommend('Batman')
```

1363
210

428
3
1364

```
In [258]: final_df.iloc[428].title
```

Out[258]: 'Batman Returns'

Insight:

- 428 indexed value has the title "Batman Returns" which is similar to the input "Batman"

```
In [259]: # Storing the model in .pkl format
import pickle
pickle.dump(final_df,open('movies.pkl','wb'))
```

```
In [ ]:
```

Type *Markdown* and LaTeX: α^2

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```