

Data mining - mini project 2 (Harsh Siroya)

The questions i am addressing are:

1. what is the most common genre of movies? and what is the total and mean box office revenue for each genre?
2. Is there a correlation between the actor age and the box office revenue? and what is the most common age of actors in the dataset? This will help us to know if the age of the actor has an impact on the box office revenue. Also the total box office revenue for each movie a specific actor has acted in. This will help us identify the most successful actors in the dataset. This we can also correlate with the number of movies an actor has acted in.
3. What is the total revenue per year? This will help us to know the trend of the box office revenue over the years, and if the revenue is increasing or decreasing.
4. Which movie is similar to the other movies in the dataset? This will help people watch movies that are similar to the movies they like. Just like the recommendation system in Netflix.

```
In [ ]:
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import ast

In [ ]:
character_metadata_cols = [
    "wikipedia Movie ID", "Freebase Movie ID", "Movie Release Date", "Character Name", "Actor DOB",
    "Actor Gender", "Actor Height", "Actor Ethnicity", "Actor Name",
    "Actor Age at Movie Release", "Freebase character sap id", "Freebase character id", "Freebase actor id"
]
character_metadata_df = pd.read_csv("MovieSummaries/character_metadata.csv", sep=";", header=None, names=character_metadata_cols)

movie_metadata_cols = [
    "wikipedia Movie ID", "Freebase Movie ID", "Movie Name", "Movie Release Date",
    "Movie Box Office Revenue", "Movie Runtime", "Movie Languages", "Movie Countries",
    "Movie Genres"
]
movie_metadata_df = pd.read_csv("MovieSummaries/movie_metadata.csv", sep=";", header=None, names=movie_metadata_cols)

plot_summaries_df = pd.read_csv("MovieSummaries/plot_summaries.txt", sep=";", header=None, names=["wikipedia Movie ID", "Plot Summary"])
names_clusters_df = pd.read_csv("MovieSummaries/names_clusters.txt", sep=";", header=None, names=["Name", "Cluster ID"])
vtropes_clusters_df = pd.read_csv("MovieSummaries/vtropes_clusters.txt", sep=";", header=None, names=["TV Trope", "Details"])

In [ ]:
#Exploratory Data Analysis
print(movie_metadata_df.head())
print(movie_metadata_df.info())
print("Summary Information of Movie Metadata:")
print(movie_metadata_df.isnull().sum())
print("Missing Values in Movie Metadata:")

Movie Metadata:
wikipedia Movie ID  Freebase Movie ID  \
0          975908          /w/8b3wyt  \
1          336793          /w/8b3wyt  \
2          2663795          /w/8b3wyt  \
3          9363483          /w/8b3wyt  \
4          262136          /w/8b3wyt  \

Movie Name Movie Release Date  \
0          Ghosts of Mars          2001-08-24  \
1          Getting Away with Murder: The JonBenet Ramsey ...          2000-02-16  \
2          White of the Eye          1987  \
3          A Woman in Flames          1983  \

Movie Box Office Revenue Movie Runtime  \
0          14018832.0          98.0  \
1          NaN          95.0  \
2          NaN          81.0  \
3          NaN          118.0  \
4          NaN          106.0  \

Movie Languages  \
0          ("en/8b3wyt": "English Language")  \
1          ("en/8b3wyt": "English Language")  \
2          ("en/8b3wyt": "Norwegian Language")  \
3          ("en/8b3wyt": "English Language")  \
4          ("en/8b3wyt": "German Language")  \

Movie Countries  \
0          ("en/8b3wyt": "United States of America")  \
1          ("en/8b3wyt": "United States of America")  \
2          ("en/8b3wyt": "Norway")  \
3          ("en/8b3wyt": "Germany")  \
4          ("en/8b3wyt": "Germany")  \

Movie Genres  \
0          ("en/8b3wyt": "Thriller", "/w/8b3wyt": "Science...")  \
1          ("en/8b3wyt": "Mystery", "/w/8b3wyt": "Mystery...")  \
2          ("en/8b3wyt": "Crime Fiction", "/w/8b3wyt": "C...")  \
3          ("en/8b3wyt": "Thriller", "/w/8b3wyt": "Thriller...")  \
4          ("en/8b3wyt": "Drama")  \

Summary Information of Movie Metadata:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 81741 entries, 0 to 81740
Data columns (total 9 columns):
#   Column              Non-Null Count  Dtype  \
--  --              \
0   wikipedia Movie ID  81741 non-null  int64  \
1   Freebase Movie ID  81741 non-null  object \
2   Movie Name         78879 non-null  object \
3   Movie Release Date  78879 non-null  object \
4   Movie Box Office Revenue  8893 non-null  float64 \
5   Movie Runtime       81741 non-null  float64 \
6   Movie Languages     81741 non-null  object \
7   Movie Countries     81741 non-null  object \
8   Movie Genres        81741 non-null  object \
dtypes: float64(1), int64(1), object(6)
memory usage: 3.6+ MB
None

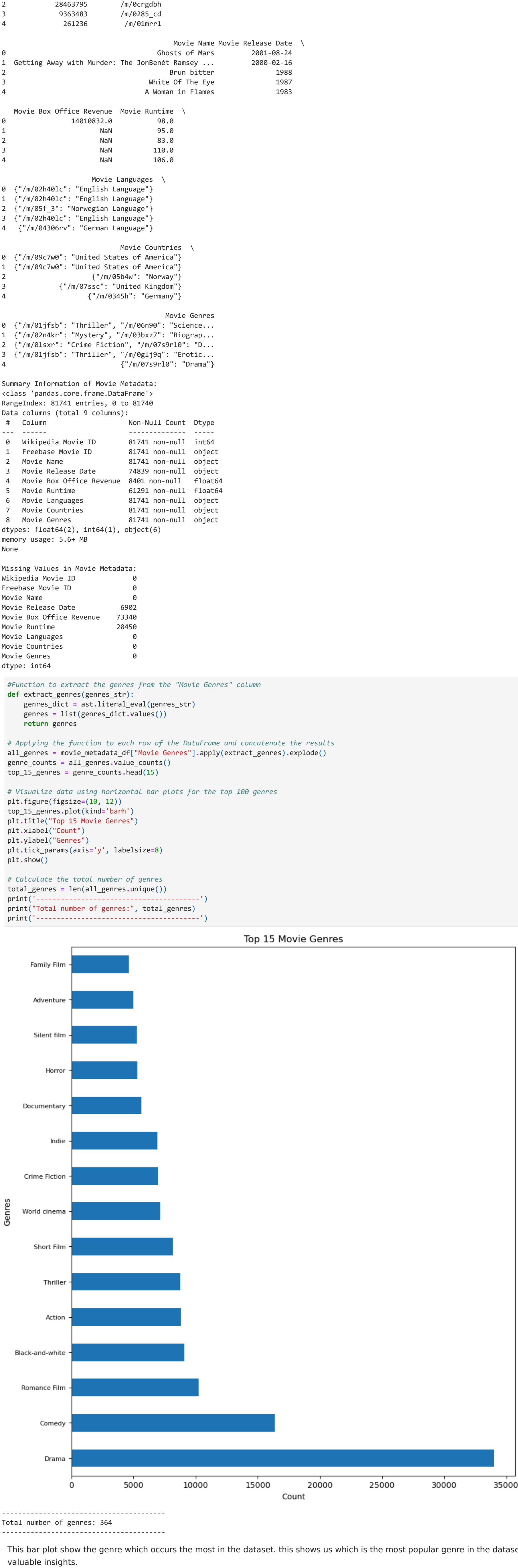
Missing Values in Movie Metadata:
wikipedia Movie ID  0
Freebase Movie ID  0
Movie Name         2348
Movie Release Date  6902
Movie Box Office Revenue  72488
Movie Runtime       28458
Movie Languages     0
Movie Countries     0
Movie Genres        0
dtype: int64

In [ ]:
#Function to extract the genres from the "Movie Genres" column
def extract_genres(genres_str):
    genres_list = ast.literal_eval(genres_str)
    genres = list(genres_list.values())
    return genres

# Applying the function to each row of the dataframe and concatenate the results
all_genres = movie_metadata_df["Movie Genres"].apply(extract_genres).explode()
genre_counts = all_genres.value_counts()
top_15_genres = genre_counts.head(15)

# Visualize data using horizontal bar plots for the top 100 genres
plt.figure(figsize=(10, 12))
top_15_genres.plot(kind='bar')
plt.title("Top 15 Movie Genres")
plt.xlabel("Genre")
plt.ylabel("Count")
plt.show()

# Create the total number of genres
total_genres = len(all_genres.unique())
print("Total number of genres: ", total_genres)
print("Total number of genres: ", total_genres)
print("Total number of genres: ", total_genres)
```



-----  
Total number of genres: 34  
-----  
This bar plot shows the genre which occurs the most in the dataset. This shows us which is the most popular genre in the dataset. Directors can use this information to make movies that are more likely to be successful, along with more valuable insights.

```
In [ ]:
# To understand the character metadata
print(character_metadata_df.head())
print(character_metadata_df.info())
print(character_metadata_df.isnull().sum())

Movie Metadata:
wikipedia Movie ID  Freebase Movie ID  Movie Release Date  \
0          975908          /w/8b3wyt          2001-08-24  \
1          336793          /w/8b3wyt          2001-08-24  \
2          2663795          /w/8b3wyt          2001-08-24  \
3          9363483          /w/8b3wyt          2001-08-24  \
4          262136          /w/8b3wyt          2001-08-24  \

Character Name  Actor DOB Actor Gender  Actor Height  \
0          Al Pacino          1940-04-25          M          1.83  \
1          Lieutenant Melvin Belland          1974-08-15          F          1.78  \
2          Donald Williams          1960-08-15          M          1.77  \
3          Sgt Jericho Butler          1967-09-12          M          1.75  \
4          Bashira Kincaid          1977-09-25          F          1.69  \

Actor Ethnicity  Actor Name Actor Age at Movie Release  \
0          NaN          Nana Du Jusu          42.0  \
1          /w/8b3wyt          Nana Du Jusu          42.0  \
2          /w/8b3wyt          Ice Cube          32.0  \
3          NaN          Jason Statham          33.0  \
4          NaN          Eric Dujail          22.0  \

Freebase character sap id Freebase character id Freebase actor id  \
0          /w/8b3wyt          /w/8b3wyt          /w/8b3wyt  \
1          /w/8b3wyt          /w/8b3wyt          /w/8b3wyt  \
2          /w/8b3wyt          /w/8b3wyt          /w/8b3wyt  \
3          /w/8b3wyt          /w/8b3wyt          /w/8b3wyt  \
4          /w/8b3wyt          /w/8b3wyt          /w/8b3wyt  \

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 81609 entries, 0 to 81608
Data columns (total 13 columns):
#   Column              Non-Null Count  Dtype  \
--  --              \
0   wikipedia Movie ID  496809 non-null  int64  \
1   Freebase Movie ID  496809 non-null  object \
2   Movie Release Date  446074 non-null  object \
3   Character Name       153754 non-null  object \
4   Actor DOB           344524 non-null  object \
5   Actor Gender         405960 non-null  object \
6   Actor Height         154824 non-null  float64 \
7   Actor Ethnicity      348441 non-null  object \
8   Actor Name           449441 non-null  object \
9   Actor Age at Movie Release  292556 non-null  float64 \
10  Freebase character sap id  496809 non-null  object \
11  Freebase character id  152884 non-null  object \
12  Freebase actor id     440454 non-null  object \
dtypes: float64(2), int64(1), object(10)
memory usage: 41.7+ MB
None

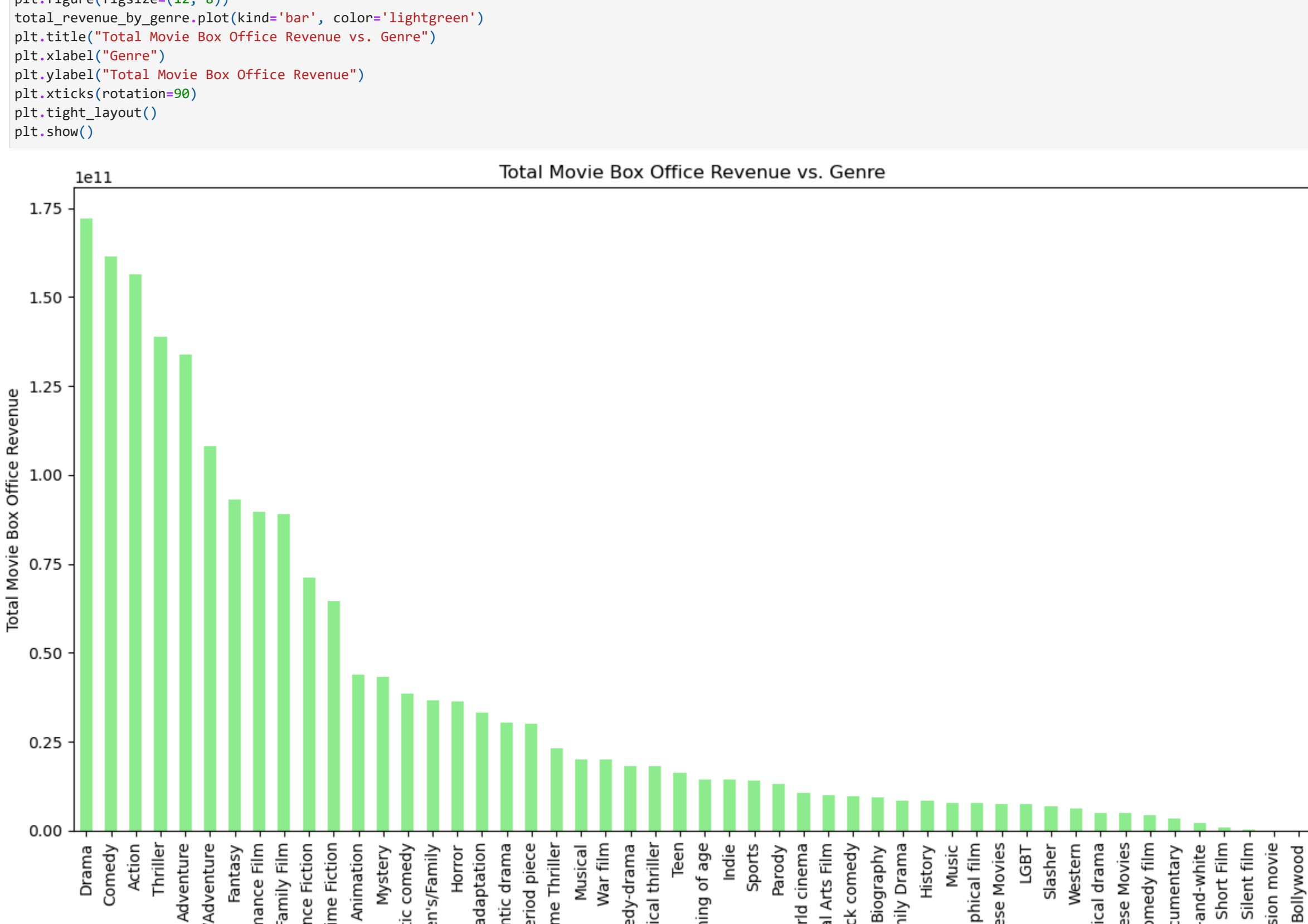
Movie Metadata:
wikipedia Movie ID  0
Freebase Movie ID  0
Movie Release Date  9995
Character Name       27875
Actor DOB           28165
Actor Gender         45069
Actor Height         295845
Actor Ethnicity      34441
Actor Name           1228
Actor Age at Movie Release  15813
Features: character sap id  0
Freebase character id  27585
Freebase actor id    85
dtype: int64

In [ ]:
movie_metadata_df["genres"] = movie_metadata_df["Movie Genres"].apply(extract_genres)
all_genres = movie_metadata_df.explode("genres")
movie_genres_df = movie_genres_df.dropna(subset=["Movie Box Office Revenue"])

top_50_genres = genre_counts.head(50)
movie_genres_df = movie_genres_df[movie_genres_df["genres"].isin(top_50_genres)]
movie_genres_df["Movie Box Office Revenue"] = pd.to_numeric(movie_genres_df["Movie Box Office Revenue"])

# Calculate the total revenue for each genre
total_revenue_by_genre = movie_genres_df.groupby("genres")["Movie Box Office Revenue"].sum().sort_values(ascending=False)

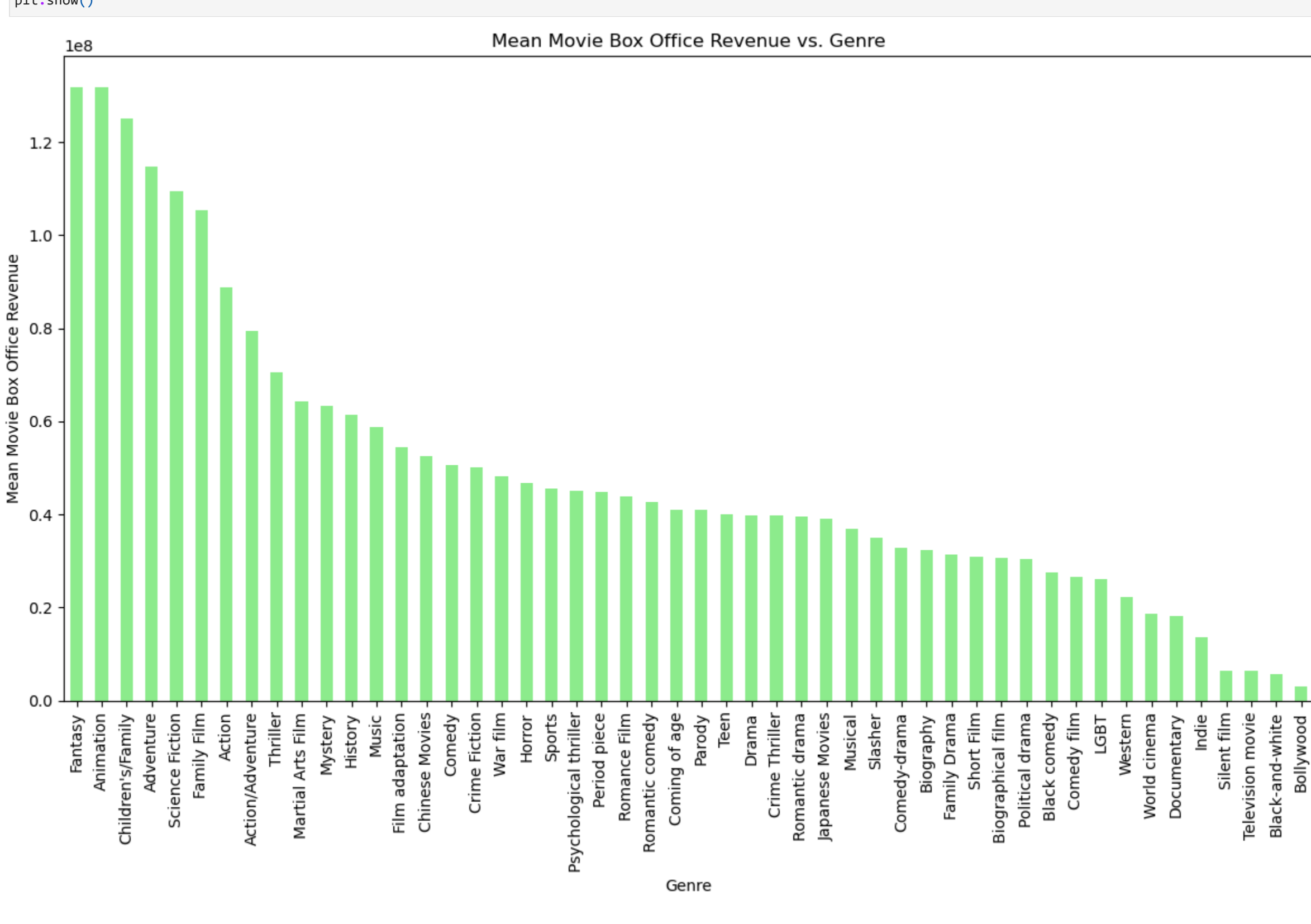
# Plot the bar plot
plt.figure(figsize=(12, 8))
total_revenue_by_genre.plot(kind="bar", color="lightgreen")
plt.title("Total Movie Box Office Revenue vs. Genre")
plt.xlabel("Genre")
plt.ylabel("Total Movie Box Office Revenue")
plt.xticks(rotation=90)
plt.tight_layout()
plt.show()
```



The bar plot displays the total movie box office revenue in USD for each genre, sorted in descending order of total revenue.

```
In [ ]:
# Calculate the total revenue for each genre
mean_revenue_by_genre = movie_genres_df.groupby("genres")["Movie Box Office Revenue"].mean().sort_values(ascending=False)

# Plot the bar plot
plt.figure(figsize=(12, 8))
mean_revenue_by_genre.plot(kind="bar", color="lightgreen")
plt.title("Mean Movie Box Office Revenue vs. Genre")
plt.xlabel("Genre")
plt.ylabel("Mean Movie Box Office Revenue")
plt.xticks(rotation=90)
plt.tight_layout()
plt.show()
```



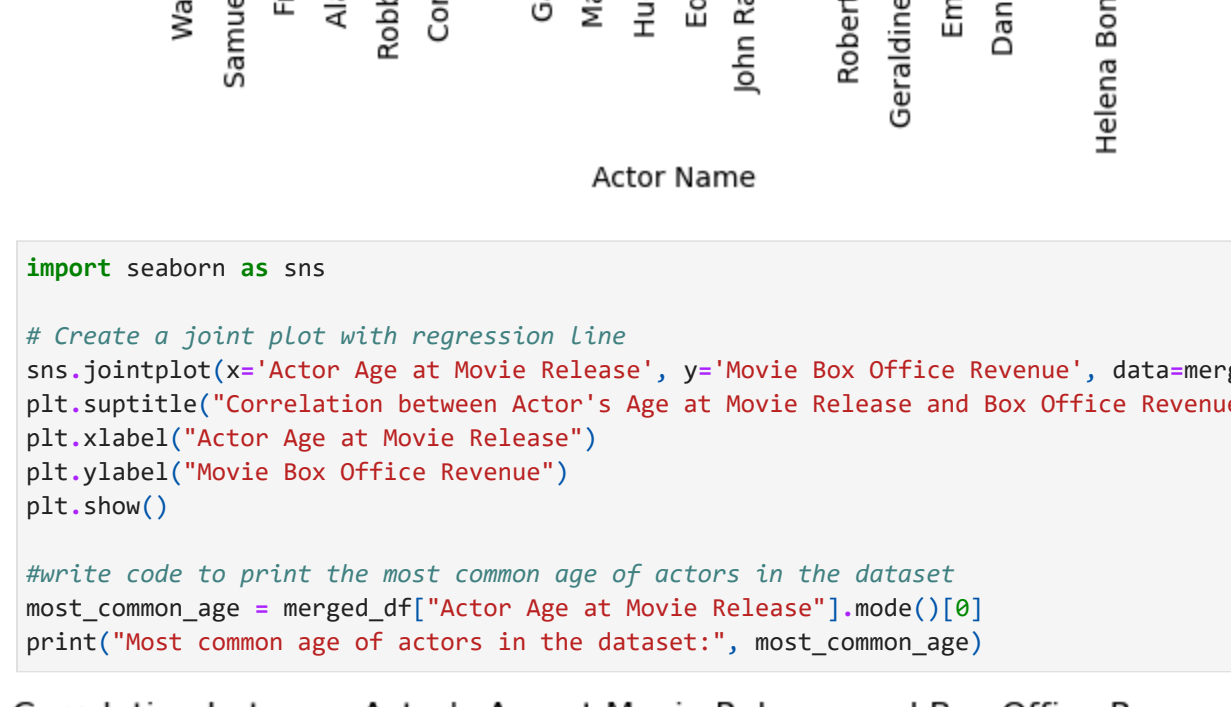
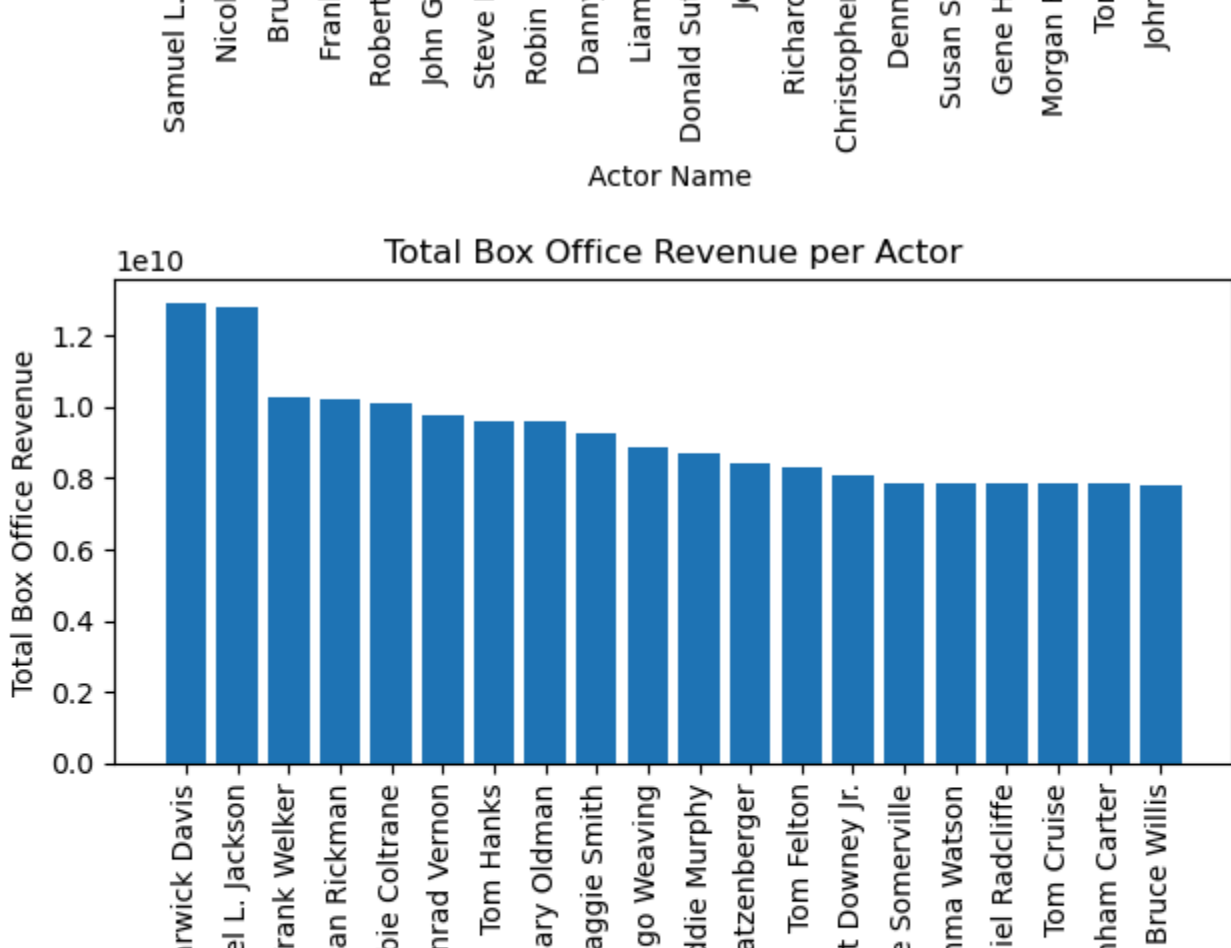
The bar plot displays the mean movie box office revenue in USD for each genre, sorted in descending order of mean revenue.

```
In [ ]:
merged_df = pd.merge(character_metadata_df, movie_metadata_df, on="wikipedia Movie ID", how="inner")
merged_df.dropna(subset=["Actor Name", "Movie Release Date", "Movie Box Office Revenue"], inplace=True)

# Exploratory Data Analysis
print(movie_metadata_df.head())
print(movie_metadata_df.info())
print(movie_metadata_df.isnull().sum())

# Analyze the distribution of the number of movies each actor has appeared in
movie_per_actor = merged_df.groupby("Actor Name")["Movie Box Office Revenue"].count().sort_values(ascending=False)
plt.bar(movie_per_actor.index[:20], movie_per_actor.values[:20])
plt.title("Number of Movies per Actor")
plt.xlabel("Actor Name")
plt.ylabel("Number of Movies")
plt.xticks(rotation=90)
plt.tight_layout()
plt.show()
```

```
# Analyze box office success of actors
box_office_revenue_per_actor = merged_df.groupby("Actor Name")["Movie Box Office Revenue"].sum().sort_values(ascending=False)
plt.bar(box_office_revenue_per_actor.index[:20], box_office_revenue_per_actor.values[:20])
plt.title("Total Box Office Revenue per Actor")
plt.xlabel("Actor Name")
plt.ylabel("Total Box Office Revenue")
plt.xticks(rotation=90)
plt.tight_layout()
plt.show()
```



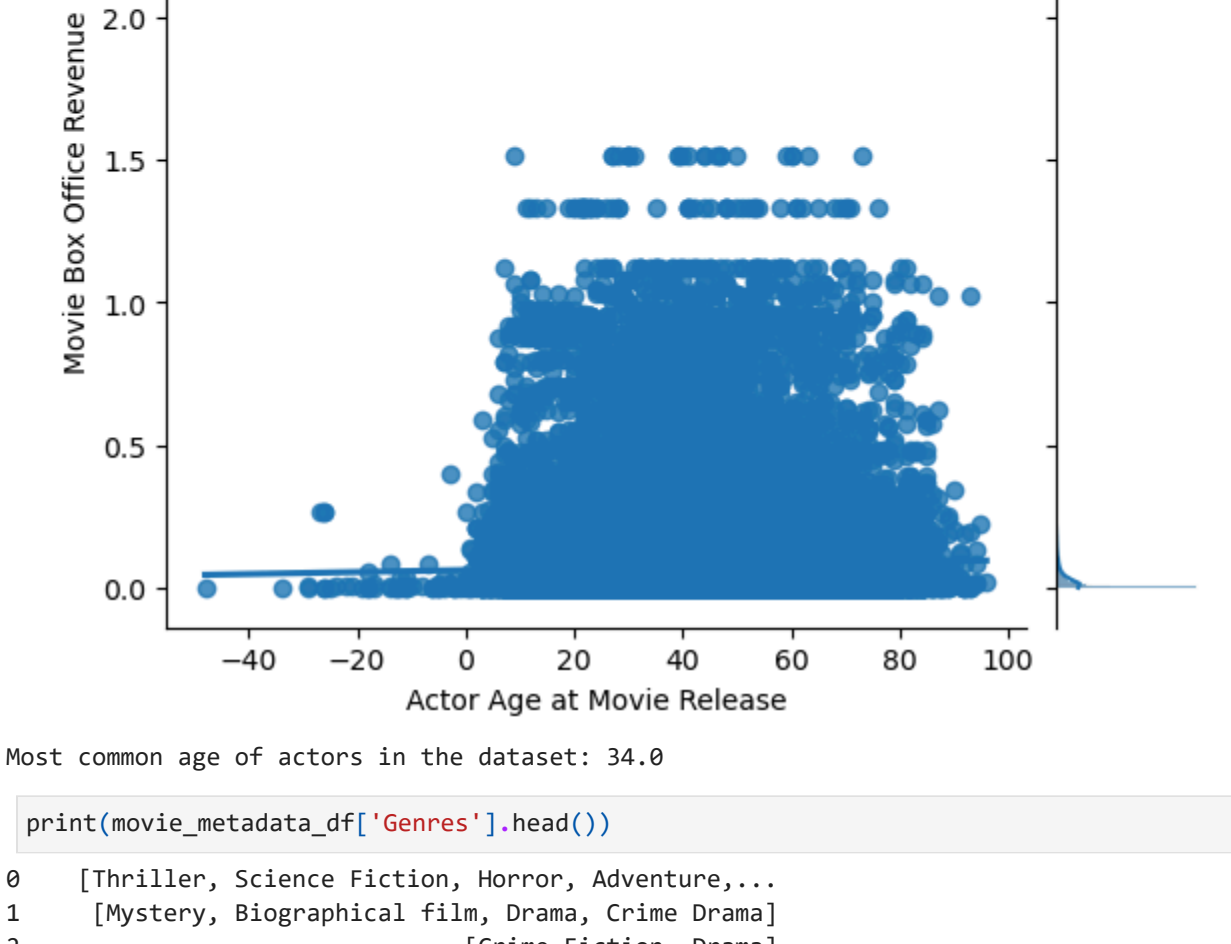
```
In [ ]:
import seaborn as sns

# Create a joint plot with regression line
sns.jointplot(x="Actor Age at Movie Release", y="Movie Box Office Revenue", data=merged_df, kind="reg")
plt.suptitle("Correlation between Actor's Age at Movie Release and Box Office Revenue", x=0.5, y=1.02)
plt.xlabel("Actor Age at Movie Release")
plt.ylabel("Movie Box Office Revenue")
plt.show()
```

Another code to print the most common age of actors in the dataset

```
most_common_age = merged_df["Actor Age at Movie Release"].mode()[0]
print("Most common age of actors in the dataset: ", most_common_age)
```

Correlation between Actor's Age at Movie Release and Box Office Revenue



Most common age of actors in the dataset: 34.0

```
In [ ]:
print(movie_metadata_df["genres"].head())

0   (Thriller, Science Fiction, Horror, Adventure,...)
1   (Mystery, Biographical Film, Drama, Crime Dram...)
2   [Crime Fiction, Drama]
3   (Thriller, Erotic thriller, Psychological Thri...)
4   Genres, dtype: object
```

```
In [ ]:
# Selecting relevant features
features = ["Movie Name", "Movie Release Date", "Movie Languages", "Movie Countries", "Movie Box Office Revenue"]

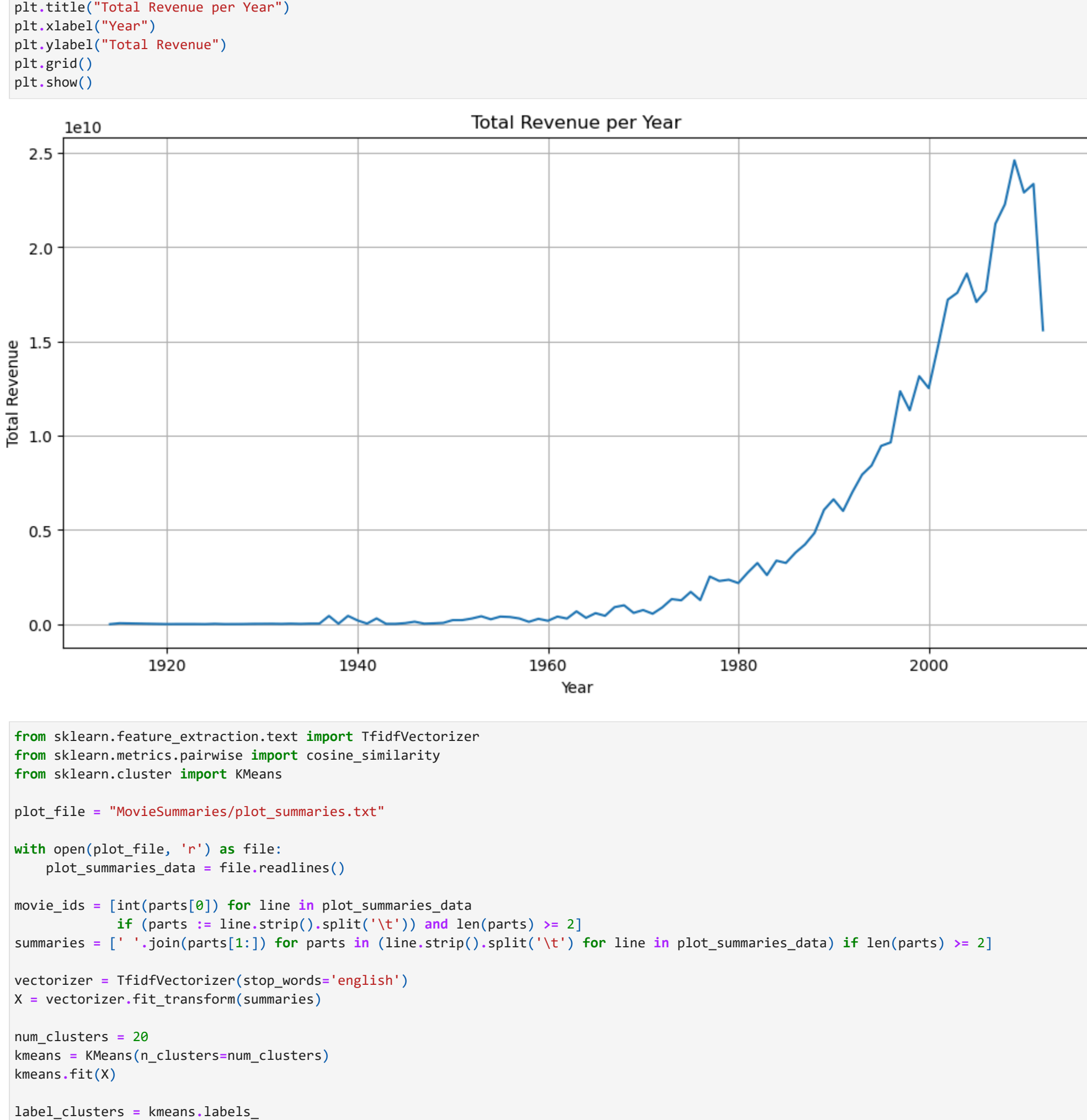
# Extracting only the relevant columns
data = movie_metadata_df[features].dropna()

# Preprocessing the data for each column
data["Movie Release Date"] = data["Movie Release Date"].split("(")[0]
data["Movie Release Date"] = pd.to_numeric(data["Movie Release Date"])

data["Movie Release Date"] = data["Movie Release Date"].drop(data[data["Movie Release Date"] < 1980].index)
```

```
# For each release year, calculate the total revenue in that year
revenue_per_year = data.groupby("Movie Release Date")["Movie Box Office Revenue"].sum()

# Plot the total revenue per year
plt.figure(figsize=(12, 6))
revenue_per_year.plot()
plt.title("Total Revenue per Year")
plt.xlabel("Year")
plt.ylabel("Total Revenue")
plt.grid()
plt.show()
```



```
In [ ]:
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity
from sklearn.cluster import KMeans

plot_file = "MovieSummaries/plot_summaries.txt"

with open(plot_file, "r") as file:
    plot_summaries_data = file.readlines()

movie_ids = [int(parts[0]) for line in plot_summaries_data
              if len(parts) > 1]
summaries = [" ".join(parts[1:]) for parts in (line.strip().split("\t") for line in plot_summaries_data) if len(parts) > 2]

vectorizer = TfidfVectorizer(stop_words='english')
X = vectorizer.fit_transform(summaries)

num_clusters = 20
kmeans = KMeans(n_clusters=num_clusters)
kmeans.fit(X)

label_clusters = kmeans.labels_
movie_names = pd.DataFrame({"wikipedia Movie ID": movie_ids})
movie_names["movie_names"] = movie_names["wikipedia Movie ID"].map(lambda x: " ".join(movie_metadata_df[movie_metadata_df["wikipedia Movie ID"] == x].Movie Name))

movies_clusters = {}
for i, cluster_label in enumerate(label_clusters):
    if cluster_label not in movies_clusters:
        movies_clusters[cluster_label] = []
        movies_clusters[cluster_label].append(movie_names.iloc[i]["Movie Name"])

for cluster_label, movies in movies_clusters.items():
    print(f"Cluster {cluster_label} = {movies}")
    movie_names = movies[:10]
    print(", ".join(movie_names))
    print()
```

~/opt/miniconda3/lib/python3.12/site-packages/sklearn/cluster/\_kmeans.py:1442: FutureWarning: The default value of 'n\_init' will change from 10 to 'auto' in 1.4. Set the value of 'n\_init' explicitly to suppress the warning

super().check\_params\_vs\_init(X, default\_n\_init=10)

Cluster 0: The Biggest Fan, A la salida nos vemos, Girl, Positive, The Incredibly True Adventure of Two Girls in Love, Classmates, Friendship, The Brotherhood III: Young Demons, Fish Hook, American Pie Presents: The Book of Love

Cluster 1: Meet John Doe, Darkness, A Modern Hero, Saturday the 14th, The Tie That Binds, All Quiet on the Western Front, Let Sleeping Corpses Lie, Heat Lightning, Food of Love, Princess of Thieves

Cluster 2: Ghost In The Moonday Sun, Class of '61, Against Her Will: An Incident in Baltimore, La Cité de la peur, A Slumdog Millionaire Goes Dancing, Twist, Mysterious Mose, In the Land of the Deaf, Changes, The Emperor Jones

Cluster 3: Eros, Green Dragon, The Rats of Tobruk, A Merry Mixup, The Good, the Bad, the Weird, Rendezvous, The Bloody Fists, The Assassin, Jacob, the Liar, Ghetto Stories: The Movie

Cluster 4: The Lemon Drop Kid, A Cry in the Dark, End Game, House Party 2, Pieces, Killjoy, Jaws: The Revenge, Hush, Hush, Sweet Charlotte, Treed Murray, Expired

Cluster 5: The

Traceback (most recent call last)

Cell In[38], line 35

88 print(cluster\_label, cluster\_label + 1:"]