

Project: Investigate a Dataset (Tmdb Movies Data Set)

Table of Contents

Introduction

Data Cleaning

Data Wrangling

Data Exploratory

Introduction

The movie database (tmdb) is an records that contains about 10865 movies which contain information such as the id, imdb_id, popularity, budget, revenue, original title, cast, release_date, vote_count, vote_average, release_year, runtime, genres, director, The homepage, langughe, overview, production_companies are some of the additional information, some column contains multiple values separated with pipe(|), columns such as genres, cast, production_companies. At the last two columns a contain the budget_adj, revenue_adj which represent the budget and revenue of the associated movie in US Dollar base on inflation as at 2010

Motives

From the data set we can answer some question such as :

Top Ten Years which Movie are Produced.

Top Genres according to the Revenue Generated.

Average year which Movies are produced by year.

The relationship between popularity and vote_count.

The relationship with Budget and Revenue.

Top 10 Directors according to the Revenue Generated.

Top 10 genres according to the Revenue generated.

The range of movies released by Month.

Testing correlation between the numerical values.

```
In [15]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

movies = pd.read_csv('C:/Users/ADMIN/Desktop/My DATABASE/tmdb-movies.csv')
```

```
Out[16]:
```

	id	imdb_id	popularity	budget	revenue	original_title	cast	homepage	director	tagline	...	overview	runtime	genres	production_companies	release_date
0	139397	tt00961010	32.985763	150000000	1513528810	Jurassic World	Chris Pratt Dwayne Dallas Howard Tyronne Khuri ...	http://www.jurassicworld.com/	Colin Trevorrow	The park is open...	...	Twenty-two years after the events of Jurassic...	124	Action Adventure Science Fiction Thriller	Universal Studios Amblin Entertainment genda...	2015-06-09
1	70341	tt1382180	28.4481936	150000000	378436354	Mad Max: Fury Road	Tom Hardy Charlize Theron Hugh Keays-Byrne Nic...	http://www.madmaxmovie.com/	George Miller	What a Lovely Day...	...	An apocalyptic story set in the furthest reaches...	120	Action Adventure Science Fiction Thriller	Village Roadshow Pictures Kennedy Miller Produ...	2015-05-13
2	262500	tt2906446	13.112507	110000000	295238201	Insurgent	Shailene Woodley Theo James Kate Winslet Kane...	http://www.thevengentseries.moviesinsurgent...	Robert Schwendke	One Choice ... Can Destroy You	...	Revenge. Prior must confront her inner demons...	119	Adventure Science Fiction Thriller	Entertainment Maxwellville Film Media Rights ...	2015-03-18
3	145007	tt2488496	11.173104	200000000	2068178225	Star Wars: The Force Awakens	Harrison Ford Carrie Fisher Adam D...	http://www.starwars.com/finncaster-wars-ep8...	J.J. Abrams	Every generation has a story	...	Thirty years after defining the Galactic Empir...	136	Action Adventure Science Fiction Fantasy	Lucasfilm Twentieth Productions Bad Robot	2015-12-15
4	168059	tt2088621	9.335014	190000000	1506249960	Furious 7	Vin Diesel Paul Walker Jason Statham Michael...	http://www.furious7.com/	James Wan	Vengeance has no home	...	Deckard Shaw seeks revenge against Dominic To...	137	Action Crime Thriller	Universal Pictures Original Film Media Rights ...	2015-04-01

5 rows x 17 columns

```
In [17]: movies.isna().sum()
```

```
Out[17]:
```

id	8
imdb_id	18
popularity	8
budget	8
revenue	8
director	44
runtime	8
genres	23
production_companies	23
release_date	3038
vote_count	1838
vote_average	8
release_year	8
release_year	8
dtype: int64	

```
In [18]: movies.shape
```

```
Out[18]:
```

(16866, 17)

```
In [19]: movies.drop(columns = ['imdb_id','original_title','cast', 'homepage','tagline','keywords','overview','budget_adj','revenue_adj'], inplace = True)
```

```
In [20]: movies.isna().sum()
```

```
Out[20]:
```

id	8
popularity	8
budget	8
revenue	8
director	44
runtime	8
genres	23
production_companies	23
release_date	3038
vote_count	1838
vote_average	8
release_year	8
dtype: int64	

```
In [21]: movies.nunique()
```

```
Out[21]:
```

id	10865
popularity	10814
budget	857
revenue	4752
director	567
runtime	247
genres	2309
production_companies	7445
release_date	5989
vote_count	129
vote_average	72
release_year	56
dtype: int64	

```
In [22]: movies.columns
```

```
Out[22]:
```

Index(['id', 'popularity', 'budget', 'revenue', 'director', 'runtime', 'genres', 'production_companies', 'release_date', 'vote_count', 'vote_average', 'release_year'], dtype='object')

```
In [23]: movies.describe()
```

```
Out[23]:
```

	id	popularity	budget	revenue	runtime	vote_count	vote_average	release_year
count	10865.000000	10865.000000	1.086600e+04	1.086600e+04	10865.000000	10865.000000	10865.000000	10865.000000
mean	66064.177434	1.165453	0.146207e+07	3.982332e+07	102.070863	217.389748	5.974922	2001.322658
std	82130.130561	1.000185	3.091321e+07	1.170035e+08	31.381405	575.619058	0.935142	12.812941
min	5.000000	0.000005	0.000000e+00	0.000000e+00	0.000000	10.000000	1.500000	1960.000000
25%	10396.250000	0.207383	0.000000e+00	0.000000e+00	90.000000	17.000000	5.400000	1995.000000
50%	29999.000000	0.380895	0.000000e+00	0.000000e+00	90.000000	20.000000	5.900000	2005.000000
75%	75650.000000	0.713817	1.500000e+07	2.400000e+07	111.000000	145.750000	6.900000	2011.000000
max	417889.000000	32.985763	4.250000e+08	2.781506e+09	900.000000	9787.000000	9.200000	2015.000000

```
In [24]: movies[movies['revenue']==0].head(5)
```

```
Out[24]:
```

	id	popularity	budget	revenue	director	runtime	genres	production_companies	release_date	vote_count	vote_average	release_year
47	265208	2.832340	300000000	0	Simon West	92	Thriller Crime Drama	Current Entertainment Longjett Serra / Affin...	2015-01-14	481	5.3	2015
67	334074	2.336136	200000000	0	James McTeigue	96	Crime Thriller Action	No Image Film Winkler Films Millennium Film...	2015-05-21	280	5.4	2015
74	347096	2.165453	0	0	Anne K. Black	108	Action Adventure Fantasy	Arrowstorm Entertainment	2015-06-24	27	5.1	2015
76	290959	2.141506	0	0	Alfonso Gomez-Rejon	105	Comedy Drama	Indica Panambush	2015-06-12	569	7.7	2015
92	370667	1.876037	0	0	A. Todd Smith	0	Fantasy Action Adventure	Arrowstorm Entertainment Camera 40 Productions...	2015-12-19	11	6.4	2015

```
In [25]: movies[movies['budget']==0].head(5)
```

```
Out[25]:
```

	id	popularity	budget	revenue	director	runtime	genres	production_companies	release_date	vote_count	vote_average	release_year
30	290996	3.927333	0	23955203	Bill Condon	103	Mystery Drama	BBC Film See-Saw Films Fimination Entertainme...	2015-06-19	425	6.4	2015
36	339527	3.358321	0	22354572	Alonso Poyat	101	Crime Drama Mystery	Eden Rock Media Fimination Entertainment Flyn...	2015-09-03	474	6.2	2015
72	284289	2.272044	0	45895	Jean-Baptiste L��onetti	95	Thriller	Further Films	2015-04-17	81	5.5	2015
74	347096	2.165453	0	0	Anne K. Black	108	Action Adventure Fantasy	Arrowstorm Entertainment	2015-06-24	27	5.1	2015
75	308369	2.141506	0	0	Alfonso Gomez-Rejon	105	Comedy Drama	Indica Panambush	2015-06-12	569	7.7	2015

```
In [26]: movies[movies['budget']==0].head(5)
```

```
Out[26]:
```

	id	popularity	budget	revenue	director	runtime	genres	production_companies	release_date	vote_count	vote_average	release_year
30	290996	3.927333	0	23955203	Bill Condon	103	Mystery Drama	BBC Film See-Saw Films Fimination Entertainme...	2015-06-19	425	6.4	2015
36	339527	3.358321	0	22354572	Alonso Poyat	101	Crime Drama Mystery	Eden Rock Media Fimination Entertainment Flyn...	2015-09-03	474	6.2	2015
72	284289	2.272044	0	45895	Jean-Baptiste L��onetti	95	Thriller	Further Films	2015-04-17	81	5.5	2015
74	347096	2.165453	0	0	Anne K. Black	108	Action Adventure Fantasy	Arrowstorm Entertainment	2015-06-24	27	5.1	2015
75	308369	2.141506	0	0	Alfonso Gomez-Rejon	105	Comedy Drama	Indica Panambush	2015-06-12	569	7.7	2015

Errors Detected in the Dataset

The data contains missing values

Some of the columns are not needed for analysis

The data set contain one duplicate row

Two or more genre in genre column separated with pipe(|)

Zero value at runtime, budget, revenue columns this could probably be a data entry errors

Data Cleaning

```
In [27]: movies.columns
```

```
Out[27]:
```

Index(['id', 'popularity', 'budget', 'revenue', 'director', 'runtime', 'genres', 'production_companies', 'release_date', 'vote_count', 'vote_average', 'release_year'], dtype='object')

```
In [28]: movie_copy = movies.copy()
```

```
Out[28]:
```

Replacing Zeros in Runtime,Revenue,Budget with the Mean Values Respectively

```
In [29]: def rep_zero(col_name:list, data: pd.DataFrame):
    for name in col_name:
        for i in range(data.shape[0]):
            if(data[name].iloc[i]==0):
                data[name].iloc[i]=data[name].mean()
rep_zero(['revenue','budget','runtime'], movie_copy)
```

```
C:/Users/ADMIN/AppData/Local/Temp/ipykernel_6132/1811297984.py:5: 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

```
data[name].iloc[i]=data[name].mean()
```

Dropping the Duplicated Rows

```
In [30]: movies.drop_duplicates(inplace = True)
```

Dropping the Rows with null values

```
In [31]: movies.dropna(inplace = True)
```

Remove Pipe the Genres Columns

```
In [32]: df_genres = movies.copy()[['genres']]
movies['genres'] = df_genres.apply(lambda x : x.split('|')[0])
```

```
In [33]: movies['genres'].value_counts()
```

```
Out[33]:
```

Drama	2225
Comedy	2084
Action	1479
Horror	853
Adventure	551
Thriller	6184731096
Animation	362
Crime	1829055609
Documentary	278
Fantasy	254
Science Fiction	198
Romance	161
Family	128
Mystery	118
Music	79
TV Movie	99
War	57
Western	42
History	39
Foreign	8

Name: genres, dtype: int64

```
In [34]: TopMoviesGenres = movies['genres'].value_counts()[1:10]
TopMoviesGenres_order=TopMoviesGenres.index
TopMoviesGenres.index
```

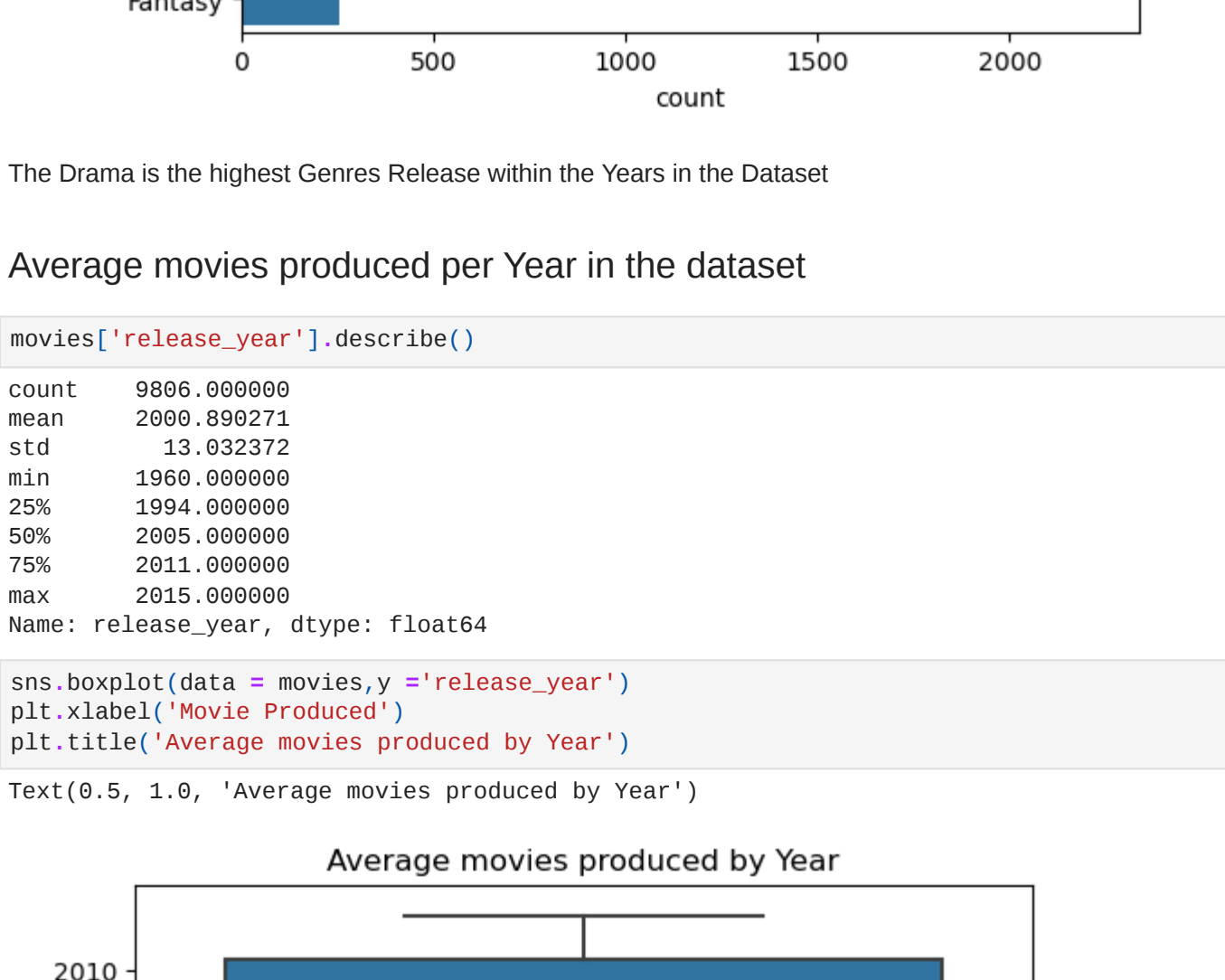
```
Out[34]:
```

Index(['Drama', 'Comedy', 'Action', 'Horror', 'Adventure', 'Thriller', 'Animation', 'Crime', 'Documentary', 'Fantasy'], dtype='object')

```
In [35]: base_color = sns.color_palette()[0]
sns.countplot(data = movies, x = 'genres', color=base_color, order=TopMoviesGenres_order)
```

```
Out[36]:
```

<Axes: xlabel='count', ylabel='genres'>



The Drama is the Highest Genres Release within the Dataset

Average movies produced per Year in the dataset

```
In [36]: movies['release_year'].describe()
```

```
Out[36]:
```

count 8886.000000
mean 2000.890271
std 11.632372
min 1960.000000
25% 1994.000000
50% 2005.000000
75% 2011.000000
max 2015.000000
Name: release_year, dtype: float64

```
In [37]: sns.boxplot(data = movies, y = 'release_year')
plt.xlabel('Movie Produced')
plt.title('Average movies produced by Year')
```

```
Out[37]:
```

Text(0.5, 1.0, 'Average movies produced by Year')



The average year of movies produced is around 2008, and maximum year of movies produce is around 2015

Top 10 Movies which Movies are produced By Year

```
In [38]: movies['release_year'].value_counts()[1:10]
```

```
Out[38]:
```

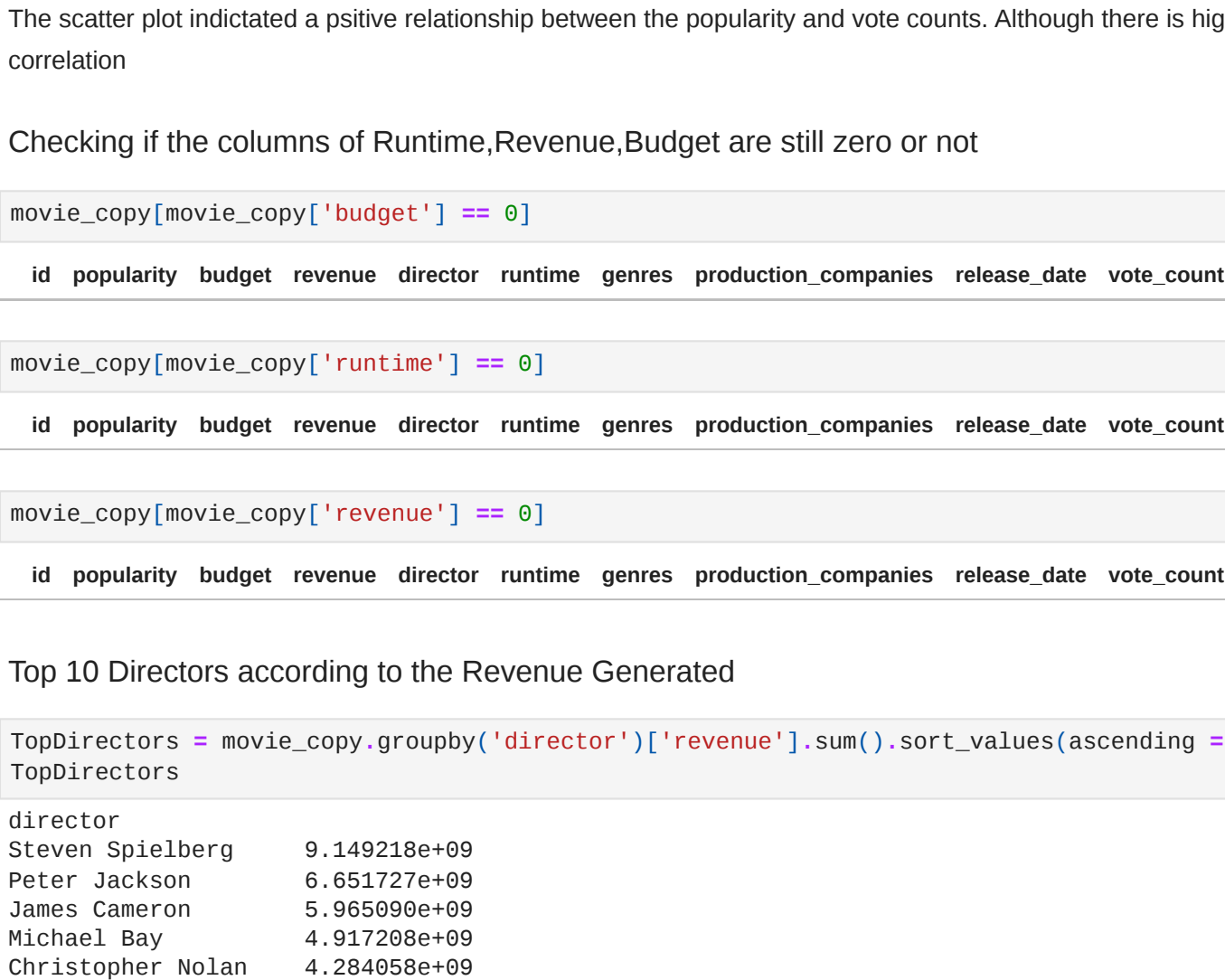
2014 638
2013 548
2015 561
2012 598
2009 474
2011 451
2008 443
2010 392
2006 251
2007 392
Name: release_year, dtype: int64

```
In [39]: x_order = [2014,2013,2015,2012,2011,2008,2010,2007,2006]
```

```
In [40]: base_color = sns.color_palette()[0]
sns.countplot(data = movies, x = 'release_year', color = base_color, order=x_order)
plt.title('Distribution of the movies produced in the ten 10years')
```

```
Out[40]:
```

Text(0.5, 1.0, 'Distribution of the movies produced in the ten 10years')



The Relationship between popularity and vote_counts

```
In [41]: sns.regplot(data = movies, x = 'popularity', y = 'vote_count')
```

```
Out[41]:
```

<Axes: xlabel='popularity', ylabel='vote_count'>

```
In [42]: movie_copy[movie_copy['budget']==0]
```

```
Out[42]:
```

	id	popularity	budget	revenue	director	runtime	genres	production_companies	release_date	vote_count	vote_average	release_year
--	----	------------	--------	---------	----------	---------	--------	----------------------	--------------	------------	--------------	--------------

```
In [43]: movie_copy[movie_copy['runtime']==0]
```

```
Out[43]:
```

	id	popularity	budget	revenue	director	runtime	genres	production_companies	release_date	vote_count	vote_average	release_year
--	----	------------	--------	---------	----------	---------	--------	----------------------	--------------	------------	--------------	--------------

```
In [44]: movie_copy[movie_copy['revenue']==0]
```

```
Out[44]:
```

	id	popularity	budget	revenue	director	runtime	genres	production_companies	release_date	vote_count	vote_average	release_year
--	----	------------	--------	---------	----------	---------	--------	----------------------	--------------	------------	--------------	--------------

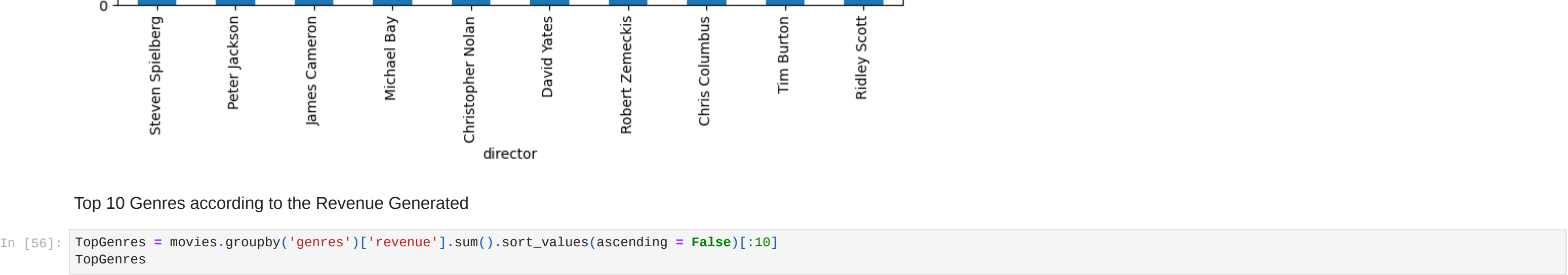
```
In [45]: Top10Directors = movie_copy.groupby(['director'])['revenue'].sum().sort_values(ascending = False)[1:10]
Top10Directors
```

```
Out[45]:
```

director	Steven Spielberg	9.149215e+89
	Peter Jackson	6.651272e+89
	James Cameron	5.960880e+89
	Michael Bay	4.917288e+89
	Christopher Nolan	4.248055e+89
	David Yates	4.238152e+89
	Robert Zemeckis	3.999550e+89
	Chris Columbus	3.957944e+89
	Tim Burton	3.851186e+89
	Ridley Scott	3.693227e+89

Name: revenue, dtype: float64

Steven Spielberg generated the highest Revenue.



Top 10 Genres according to the Revenue Generated

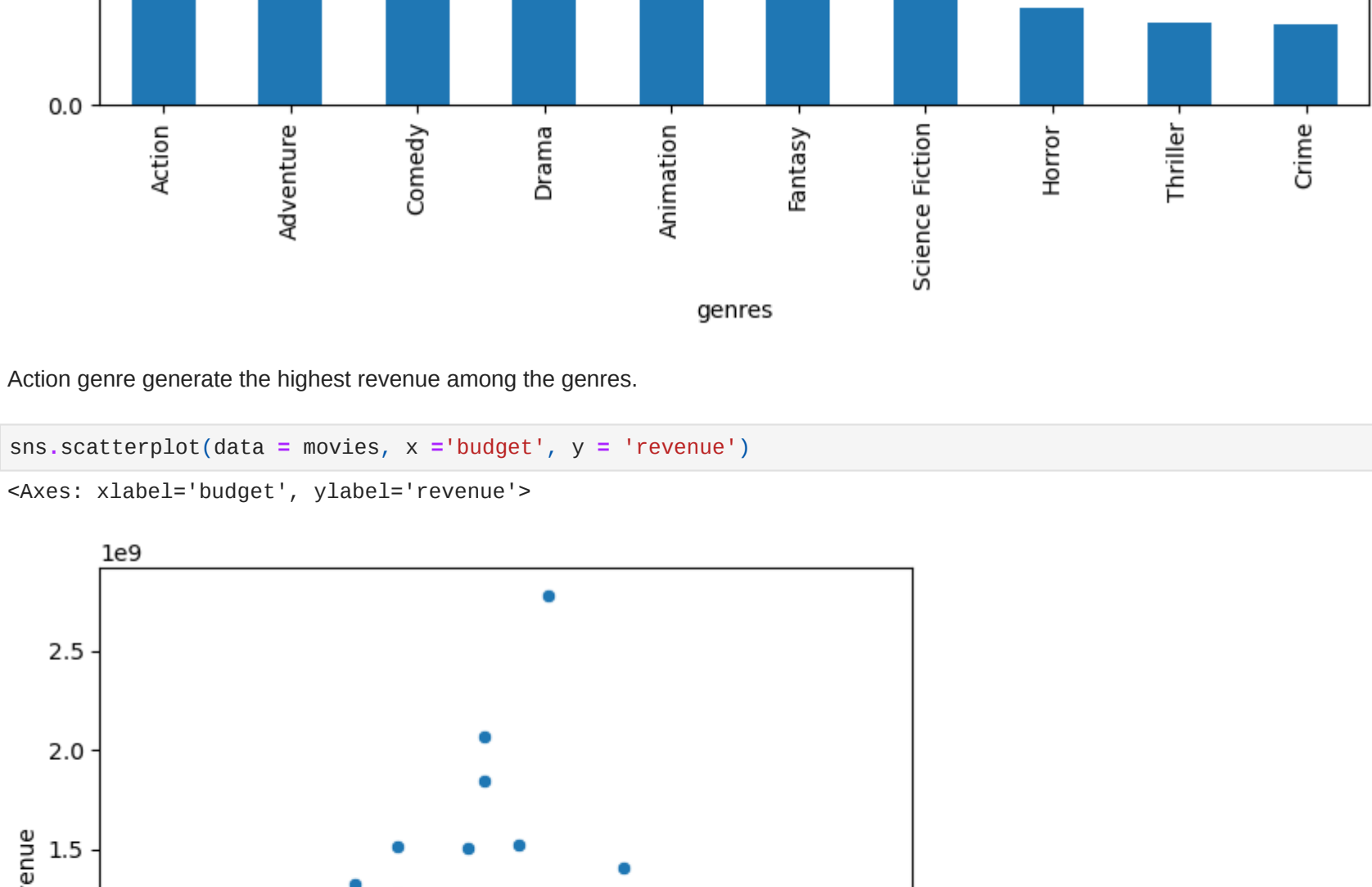
```
In [46]: TopGenres = movies.groupby(['genres'])['revenue'].sum().sort_values(ascending = False)[1:10]
TopGenres
```

```
Out[46]:
```

genres	revenue
Action	9648724872
Adventure	72848959448
Comedy	6794064785
Drama	6158173106
Animation	29748093451
Fantasy	1829055609
Science Fiction	1784641613
Horror	15920512240
Thriller	13083268725
Crime	1279554239

Name: revenue, dtype: int64

```
In [47]: TopGenres.plot(kind='bar', title='Genres with higher revenue',ylabel='Revenue (Naira)',figsize=(16,5));
```



Action genre generate the highest revenue among the genres.

```
In [48]: sns.scatterplot(data = movies, x = 'budget', y = 'revenue')
```

```
Out[48]:
```

<Axes: xlabel='budget', ylabel='revenue'>

```
In [49]: movie_copy[movie_copy['new_year']<10]
```

```
Out[49]:
```

	id	popularity	budget	revenue	runtime	vote_count	vote_average	release_year
popularity	-0.014280	1.000000	-0.121446	-0.082060	-0.078415	-0.025551	-0.058363	0.511364
budget	-0.014280	1.000000	1.000000	0.466611	0.519159	0.135045	0.800828	0.209511
revenue	-0.026260	0.016158	0.590375	1.000000	0.137418	0.759561	0.547550	0.029297
runtime	-0.078415	0.135045	0.166857	0.137418	1.000000	0.162584	0.154669	-0.112178
vote_count	-0.025551	0.800828	0.595913	0.759561	0.162584	1.000000	0.253823	0.107948
vote_average	0.058363	0.209511	0.074102	0.134750	0.154669	0.253823	1.000000	0.117632
release_year	0.511364	0.029297	0.086602	0.029297	-0.112178	0.107948	-0.117632	1.000000

```
In [50]: movie_copy['release_date'] = pd.to_datetime(movie_copy['release_date'])
movie_copy['MonthOfRelease'] = movie_copy['release_date'].dt.month_name()
movie_copy['MonthOfRelease'].value_counts()
```

```
Out[50]:
```

September	1331
October	1153
December	985
January	919
August	928
June	827
March	809
November	814
May	809
July	799
April	797
February	693

Name: MonthOfRelease, dtype: int64

The month of September is the Highest Month of Release while February has the lowest amount of movie released

```
In [51]: MonthOfRevenue = movie_copy.groupby(['MonthOfRelease'])['revenue'].sum().sort_values(ascending = False)[1:10]
MonthOfRevenue
```

```
Out[51]:
```

MonthOfRelease	revenue
December	8.562889e+10
June	6.402448e+10
May	7.393869e+10
November	7.008423e+10
July	6.743746e+10
October	6.636874e+10
September	6.376454e+10
March	5.539392e+10
April	5.121777e+10
August	5.028738e+10

Name: revenue, dtype: float64

The month of December has the highest generated Revenue. Despite September being the Highest month of Release

