

Project: Investigate a tmdb movies Dataset

Table of Contents

- Introduction
- Data Wrangling
- Exploratory Data Analysis
- Conclusions

Introduction

In this section of the report, We provide a brief introduction to the dataset we've selected for analysis. At the end of this section, we described the questions that we have plan on exploring over the course of the report.

```
In [110]: #Important libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

TMDB Movies Dataset

This data set contains information about 10,000 movies collected from The Movie Database (TMDB), including user ratings and revenue. we are here to provide helpful questions and its answers.

Questions

1. Which genres are most popular from year to year?
2. What kinds of properties are associated with movies that have high revenues?
3. What movies couldn't cover its budget(revenue = 0)?
4. What is the highest release genres?

Data Wrangling

In this section of the report, we will load in the data, and check for cleanliness, and then trim and clean our dataset for analysis. Then we will make sure that we document our steps carefully and justify our cleaning decisions.

General Properties

```
In [111]: #Explore and Load, inspect data from the dataset
df = pd.read_csv('tmdb-movies.csv')
df.head()
```

	id	imdb_id	popularity	budget	revenue	original_title	cast	homepage	director	tagline	overview	runtime	genres
0	135397	tt0395190	32.985763	150000000	1513528810	Jurassic World	Chris Pratt Bryce Dallas Howard Rianne Khan(\\...	http://www.jurassicworld.com/	Colin Trevorrow	The park is open...	Twenty-two years after the events of Jurassic...	124	Action Adventure Scienc...
1	76341	tt1382190	28.419936	150000000	378436354	Mad Max: Fury Road	Tom Hardy Charlize Theron Hugh Keays-Byrne \\...	http://www.madmaxmovie.com/	George Miller	What a Lovely Day...	An apocalyptic story set in the furthest reach...	120	Action Adventure Scienc...
2	262500	tt2098446	13.112507	110000000	295238021	Insurgent	Shailene Woodley Theo James Kate Winslet Ansel...	http://www.thedivergentseries.movie/Insurgent	Robert Schwentke	One Choice Can Destroy You	Beatrice Prior must confront her inner demons...	119	Adventure Scienc...
3	140607	tt448496	11.173104	200000000	2068178225	Star Wars: The Force Awakens	Harrison Ford Mark Hamill Carrie Fisher Adam D...	http://www.starwars.com/film/star-wars-awake...	J.J. Abrams	Every generation has a story.	Thirty years after defeating the Galactic Empi...	136	Action Adventure Scienc...
4	166259	tt2020882	9.335014	190000000	1506249360	Furious 7	Vin Diesel Paul Walker John Stamos Dwayne 'The Rock' Johnson	http://www.furious7.com/	James Wan	Vengeance Is My Name	Deckard Shaw seeks revenge against Dominic Tor...	137	Action Crime Thrill...

5 rows x 21 columns

```
In [112]: #We can see there is 10866 entries with 21 columns.
#cast, homepage, director, tagline, keywords, overview, genres and production_companies
#those are the columns with missing values
df.info()
```

```
<class 'pandas.core.data.frame'>
Int64Index: 10866 entries, 0 to 10865
Data columns (total 21 columns):
#  Column                Non-Null Count  Dtype
---  ---                ---
0  id                     10866 non-null    int64
1  imdb_id               10866 non-null    object
2  popularity            10866 non-null    float64
3  budget               10866 non-null    object
4  revenue              10866 non-null    int64
5  original_title        10866 non-null    object
6  cast                 10730 non-null    object
7  homepage             2936 non-null     object
8  director             10822 non-null    object
9  tagline              8942 non-null     object
10 keywords            9373 non-null     object
11 overview            10862 non-null    object
12 runtime             10866 non-null    int64
13 genres              10843 non-null    object
14 production_companies 9836 non-null     object
15 release_date         10866 non-null    object
16 vote_count          10866 non-null    int64
17 vote_average         10866 non-null    float64
18 release_year         10866 non-null    int64
19 budget_adj           10866 non-null    float64
20 revenue_adj          10866 non-null    float64
dtypes: float64(4), int64(6), object(11)
memory usage: 1.7+ MB
```

```
df.describe()
```

	id	popularity	budget	revenue	runtime	vote_count	vote_average	release_year	budget_adj	revenue_adj
count	10866.000000	10866.000000	1.086600e+04	1.086600e+04	10866.000000	10866.000000	10866.000000	10866.000000	1.086600e+04	1.086600e+04
mean	60604.717344	0.646641	1.462570e+07	3.962332e+07	102.070893	217.389748	5.974822	2001.322550	1.755104e+07	5.136439e+07
std	92130.136951	1.000185	3.091321e+07	1.170039e+08	31.381405	575.613698	0.930142	12.811294	3.430016e+07	1.446325e+08
min	0	0.000000	0.000000e+00	0.000000e+00	0.000000	0.000000	1.000000	1960.000000	0.000000e+00	0.000000e+00
25%	10957	5.000000	0.207603	0.000000e+00	0.000000e+00	50.000000	0.000000	1995.000000	0.000000e+00	0.000000e+00
50%	20669.000000	0.363866	0.000000e+00	0.000000e+00	99.000000	36.000000	6.000000	2006.000000	0.000000e+00	0.000000e+00
75%	75410.000000	0.713817	1.000000e+07	2.400000e+07	111.000000	145.750000	6.000000	2011.000000	2.083250e+07	3.969710e+07
max	471869.000000	32.985763	4.250000e+08	2.783090e+09	900.000000	9767.000000	9.200000	2015.000000	4.200000e+08	2.827124e+09

```
In [114]: #TMDB id and release_date are string we can convert them or we may drop them since we don't need those columns.
df.dtypes
```

```
Out[114]:
id                int64
imdb_id          object
popularity        float64
budget           int64
revenue          int64
original_title    object
cast             object
homepage         object
director         object
tagline          object
keywords         object
runtime          int64
genres           object
production_companies  object
release_date     object
vote_count       int64
vote_average     float64
release_year     int64
budget_adj       float64
revenue_adj      float64
dtype: object
```

```
In [115]: #Number of unique values in every column
df.nunique()
```

	id	popularity	budget	revenue	runtime	vote_count	vote_average	release_year	budget_adj	revenue_adj
count	10866.000000	10866.000000	1.086600e+04	1.086600e+04	10866.000000	10866.000000	10866.000000	10866.000000	1.086600e+04	1.086600e+04
mean	60604.717344	0.646641	1.462570e+07	3.962332e+07	102.070893	217.389748	5.974822	2001.322550	1.755104e+07	5.136439e+07
std	92130.136951	1.000185	3.091321e+07	1.170039e+08	31.381405	575.613698	0.930142	12.811294	3.430016e+07	1.446325e+08
min	0	0.000000	0.000000e+00	0.000000e+00	0.000000	0.000000	1.000000	1960.000000	0.000000e+00	0.000000e+00
25%	10957	5.000000	0.207603	0.000000e+00	0.000000e+00	50.000000	0.000000	1995.000000	0.000000e+00	0.000000e+00
50%	20669.000000	0.363866	0.000000e+00	0.000000e+00	99.000000	36.000000	6.000000	2006.000000	0.000000e+00	0.000000e+00
75%	75410.000000	0.713817	1.000000e+07	2.400000e+07	111.000000	145.750000	6.000000	2011.000000	2.083250e+07	3.969710e+07
max	471869.000000	32.985763	4.250000e+08	2.783090e+09	900.000000	9767.000000	9.200000	2015.000000	4.200000e+08	2.827124e+09

```
In [116]: #TMDB id and release_date are string we can convert them or we may drop them since we don't need those columns.
df.dtypes
```

	id	popularity	budget	revenue	original_title	director	runtime	genres	production_companies	vote_count	release_year
0	135397	32.985763	150000000	1513528810	Jurassic World	Colin Trevorrow	124	Action Adventure Science Fiction Thriller	Universal Studios Amblin Entertainment Legenda...	5562	2015
1	76341	28.419936	150000000	378436354	Mad Max: Fury Road	George Miller	120	Action Adventure Science Fiction Thriller	Village Roadshow Pictures Kennedy Miller Produ...	6185	2015
2	262500	13.112507	110000000	295238021	Insurgent	Robert Schwentke	119	Adventure Science Fiction Thriller	Summit Entertainment Mandeville Films Red Wag...	2480	2015
3	140607	11.173104	200000000	2068178225	Star Wars: The Force Awakens	J.J. Abrams	136	Action Adventure Science Fiction Fantasy	Lucasfilm Twentieth Productions Bad Robot	5292	2015
4	166259	9.335014	190000000	1506249360	Furious 7	James Wan	137	Action Crime Thriller	Universal Pictures Original Film Media Rights ...	2947	2015

```
In [117]: #No. of duplicates in dataset
sum(df.duplicated())
```

```
Out[117]:
1
```

```
In [118]: #Drop duplicates
df.drop_duplicates(inplace=True)
```

```
#Check duplicates
sum(df.duplicated())
```

```
Out[118]:
0
```

```
In [119]: #Drop null values since they are not numbers I don't suggest filling them with the mean
df.dropna(inplace=True)
```

```
df.shape
```

```
Out[119]:
(9896, 11)
```

```
In [120]: df.isnull().sum()
```

```
Out[120]:
id                0
popularity        0
budget            0
revenue           0
original_title    0
cast             0
homepage         0
director         0
tagline          0
keywords         0
runtime          0
genres           0
production_companies 0
vote_count       0
vote_average     0
release_year     0
dtype: int64
```

Exploratory Data Analysis

Research Question 1 Which genres are most popular from year to year?

```
In [121]: #those are the genres count in every year before splitting
gn_df = df.groupby(df['release_year']).genres.value_counts()
gn_df
```

release_year	genres	
1960	Horror	3
	Comedy	2
	Comedy Drama Romance	2
	Comedy Romance	2
	Drama	2
2015	War Action	1
	War Adventure Science Fiction	1
	War Drama	1
	Western Drama	1
	Western Drama Adventure Thriller	1
Name: genres, Length: 5697, dtype: int64		

```
In [122]: # We can split values in genres column to see the most popular genre
df_copy = df.copy()
ser = df_copy['genres'].str.split('|')
ser.index = ser.index.droplevel(-1)
ser.name = 'genres'
del df_copy['genres']
df_copy = df_copy.join(ser)
```

```
# Check on split
df_copy.head()
```

	id	popularity	budget	revenue	original_title	director	runtime	production_companies	vote_count	release_year	genres
0	135397	32.985763	150000000	1513528810	Jurassic World	Colin Trevorrow	124	Universal Studios Amblin Entertainment Legenda...	5562	2015	Action
1	135397	32.985763	150000000	1513528810	Jurassic World	Colin Trevorrow	124	Universal Studios Amblin Entertainment Legenda...	5562	2015	Adventure
0	135397	32.985763	150000000	1513528810	Jurassic World	Colin Trevorrow	124	Universal Studios Amblin Entertainment Legenda...	5562	2015	Science Fiction
0	135397	32.985763	150000000	1513528810	Jurassic World	Colin Trevorrow	124	Universal Studios Amblin Entertainment Legenda...	5562	2015	Thriller
1	76341	28.419936	150000000	378436354	Mad Max: Fury Road	George Miller	120	Village Roadshow Pictures Kennedy Miller Produ...	6185	2015	Action

```
In [123]: #This shows the most popular genre in each year
s_df = df_copy[['popularity']].str.split('|')
pop_df = pd.DataFrame({'popularity': s_df.groupby(['s_df['release_year']]).popularity.max(), reset_index()})
pop_df.index = pop_df.pop_df['popularity']
s_df = s_df[s_df['popularity'].isin(pop_df.index)]
# Most popular genres in each year
s_df
```

	release_year	genres	popularity
0	2015	Action	32.985763
0	2015	Adventure	32.985763
0	2015	Science Fiction	32.985763
0	2015	Thriller	32.985763
629	2014	Adventure	24.949134
...
10724	1969	Thriller	1.778745
10755	1978	Music	1.697618
10820	1966	Animation	1.227582
10820	1966	Family	1.227582
10820	1966	Comedy	1.227582

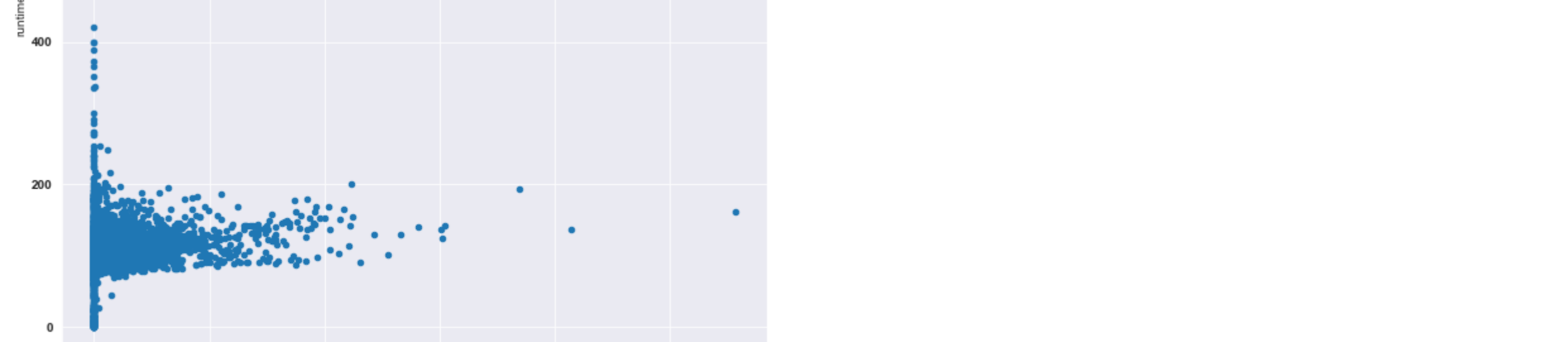
166 rows x 3 columns

```
In [124]: # Most popular genres from year to year based on count and popularity
s_df['genres'].value_counts()
```

```
Out[124]:
Adventure    31
Action       28
Science Fiction  28
Thriller     29
Drama        14
Family       11
Fantasy      10
Animation    8
Crime        7
Comedy       6
Romance      3
Mystery      3
Horror       3
Music        3
Name: genres, dtype: int64
```

```
In [125]: # Plot the bar graph.
s_df['genres'].value_counts().plot(kind='bar', fontsize=11, figsize=(8,6))
```

```
# Set the labels and titles
plt.title('Genres vs Number Of Genres Releases Each Year Based On Popularity', fontsize=15)
plt.xlabel('Genres', fontsize=13)
plt.ylabel('Number of genres releases', fontsize=13)
plt.legend()
plt.show()
sns.set_style('darkgrid')
```



Research Question 2 What kinds of properties are associated with movies that have high revenues?

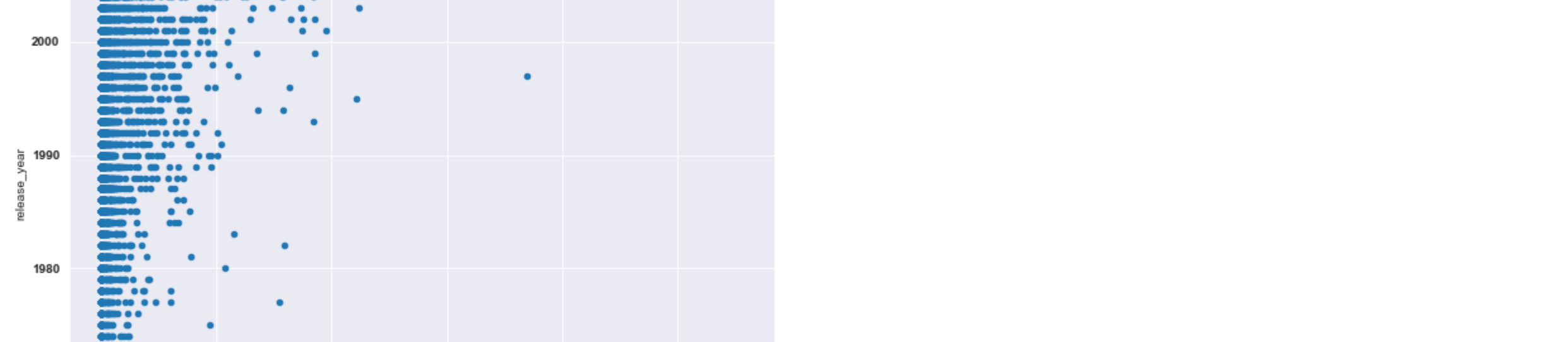
```
In [126]: # Sort df based on 'revenue'
df.sort_values('revenue', ascending=False, inplace=True)
df.head(10)
```

	id	popularity	budget	revenue	original_title	director	runtime	genres	production_companies	vote_count	release_year
1286	19955	9.432708	237000000	2718359847	Avatar	James Cameron	162	Action Adventure Fantasy Science Fiction	Ingenious Film Partners Twentieth Century Fox ...	8458	2009
3	140607	11.173104	200000000	2068178225	Star Wars: The Force Awakens	J.J. Abrams	136	Action Adventure Science Fiction Fantasy	Lucasfilm Twentieth Productions Bad Robot	5292	2015
5231	597	4.358219	200000000	24949134	Titanic	James Cameron	194	Drama Romance Thriller	Paramount Pictures Twentieth Century Fox Film ...	4654	1997
4361	24428	7.637767	220000000	151957910	The Avengers	Joss Whedon	143	Science Fiction Action Adventure	Marvel Studios	8903	2012
108	1286	9.432708	150000000	1513528810	Jurassic World	Colin Trevorrow	124	Action Adventure Science Fiction Thriller	Universal Studios Amblin Entertainment Legenda...	5562	2015
4	166259	9.335014	190000000	1506249360	Furious 7	James Wan	137	Action Crime Thriller	Universal Pictures Original Film Media Rights ...	2947	2015
14	9961	5.944827	280000000	1405035767	Avengers: Age of Ultron	Joss Whedon	141	Action Adventure Science Fiction	Marvel Studios Pine Focus Revolution Sun Studios	4304	2015
3374	12445	5.711315	150000000	1327817822	Harry Potter and the Deathly Hallows: Part 2	Chris Yates	130	Adventure Family Fantasy	Warner Bros. Heyday Films Kobal Picture Comp...	3750	2011
5422	109445	6.112766	150000000	1274219009	Frozen	Chris Buck Jennifer Lee	102	Animation Adventure Family	Walt Disney Pictures Walt Disney Animation Stu...	3769	2013
5425	68721	4.946136	200000000	1215439994	Iron Man 3	Shane Black	130	Action Adventure Science Fiction	Marvel Studios	6882	2013

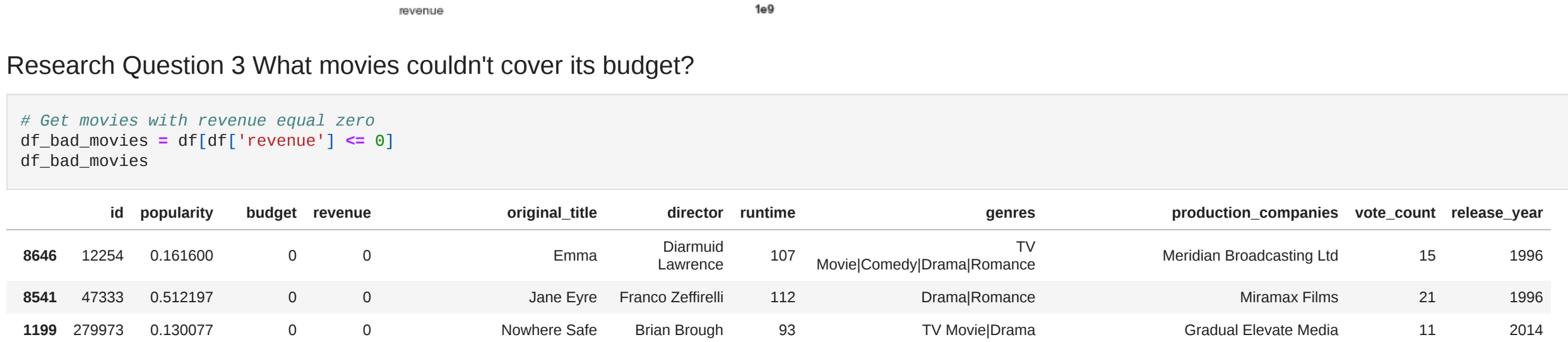
```
In [127]: # Lets see budget and revenue correlation
def scatter_plot(x,y):
    # Make a scatter plot.
    df.plot(x=x, y=y, kind='scatter', figsize=(10, 10));
```

```
# Set the title and the labels of the scatter plot.
plt.title('Highest Release Genres', fontsize=13)
plt.xlabel(x)
plt.ylabel(y)
plt.show()
```

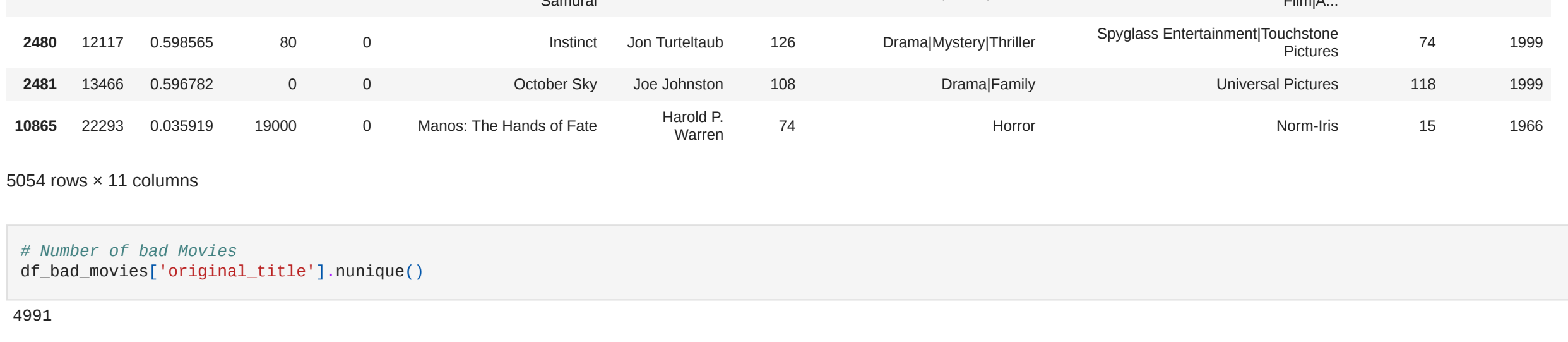
```
scatter_plot('revenue', 'budget')
```



```
In [128]: # Make a scatter plot.
scatter_plot('revenue', 'runtime')
```



```
In [129]: # revenue vs vote_count
scatter_plot('revenue', 'vote_count')
```



```
In [130]: # revenue vs release_year
scatter_plot('revenue', 'release_year')
```



Research Question 3 What movies couldn't cover its budget?

```
In [131]: # Get movies with revenue equal zero
df_bad_movies = df[df['revenue'] == 0]
df_bad_movies
```

Genre	Percentage
Foreign TV Movie	1%
Action	1%
Drama	1%
Comedy	1%
History	1%
Documentary	1%
Music	1%
Animation	1%
Mystery	1%
Fantasy	1%
Science Fiction	1%
Crime	1%
Adventure	1%
Horror	1%
Romance	1%