

Project: TMDB 5000 Movie Dataset

Table of Contents

- Introduction
- Data Wrangling
- Exploratory Data Analysis
- Conclusions

Introduction

Dataset Description

In this project, we will be analyzing data about TMDB movies. This dataset contains information about 5000 movies collected from The Movie Database (TMDB), including user ratings and revenue.

most columns in the dataset are self explanatory but there are certain columns, like 'genres' that contain multiple values separated by pipe (|) characters, but it's nothing to worry about we will clean it up. The original data can be found on kaggle <https://www.kaggle.com/datasets/tmdb/tmdb-movie-metadata>.

Question(s) for Analysis

In this project We are interested in exploring and answering the following questions using python libraries such as Pandas, numpy, and matplotlib etc:

Q1: Which movies have the highest budget spend on?

Q2: What are the top 5 highest profit making movies and the 5 lowest movies?

Q3: What is the relationships between, (Budget & Revenue), and (profit and revenue)?

Q4: Which year has the highest release of movies?

Q5: Which Month has the highest release movies?

Q6: Which months made the highest average profits?

Import Packages

```
## In this section we will import necessary packages we need to analyze the dataset
```

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
sns.set()
%matplotlib inline
```

Data Wrangling

In this section of the report, we will load, access, and explore the data to understand the characteristics of our data.

```
# load the data and read the first five rows.
df = pd.read_csv("tmdb_5000_movies.csv")
df.head(2)
```

	budget	genres	homepage	id	keywords	original_language	original_title	overview	pc
0	237000000	["id": 28, "name": "Action"], ["id": 12, "name": "Romance"]	http://www.avatarmovie.com/	19995	["id": 1463, "name": "culture clash"], ["id": 1718, "name": "culture clash"]	en	Avatar	In the 22nd century, a paraplegic Marine is di...	56
1	300000000	["id": 12, "name": "Adventure"], ["id": 14, "name": "Fantasy"], ["id": 14, "name": "Fantasy"]	http://disney.go.com/disneypictures/pirates/	285	["id": 270, "name": "culture clash"], ["id": 726, "name": "culture clash"]	en	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...	138

```
# check the data dimensions
df.shape
```

```
(4803, 20)
```

```
# let's inspect the data and get a good understanding of data
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4803 entries, 0 to 4802
Data columns (total 20 columns):
 # Column Non-Null Count Dtype
 0 budget 4803 non-null int64
 1 genres 4803 non-null object
 2 homepage 1712 non-null object
 3 id 4803 non-null int64
 4 keywords 4803 non-null object
 5 original_language 4803 non-null object
 6 original_title 4803 non-null object
 7 overview 4803 non-null object
 8 popularity 4803 non-null float64
 9 production_companies 4803 non-null object
10 production_countries 4803 non-null object
11 release_date 4802 non-null object
12 revenue 4803 non-null int64
13 runtime 4803 non-null float64
14 spoken_languages 4803 non-null object
15 status 4803 non-null object
16 tagline 3959 non-null object
17 title 4803 non-null object
18 title_original 4803 non-null object
19 vote_count 4803 non-null int64
20 vote_average float64(3), int64(4), object(13)
memory usage: 720.4+ KB
```

observing the data dimensions and the info() function result, we see that we have 4803 rows and 20 columns in our data. We also notice there are some columns that have some null values for instance, the homepage and tagline have a lot of missing value. The release date is missing 1 value and is an object(dtype) type, so in our next section on Data Cleaning we will work on this columns, drop the homepage, tagline and the Overview columns as they are not necessary for our analysis. We also change the date datatype and fill or drop the missing values for release_date and runtime as they only have few null values.

```
# Now let check our data statics using describe() and transpose() to make our data easy to read.
df.describe().transpose()
```

	count	mean	std	min	25%	50%	75%	max
budget	4803.0	2.904504e+07	4.072239e+07	0.0	7900000.0000	1.500000e+07	4.000000e+07	3.800000e+08
id	4803.0	5.716548e+04	8.869461e+04	5.0	90145000e+03	1.462900e+04	5.661050e+04	4.594880e+05
popularity	4803.0	2.149230e+01	3.181665e+01	0.0	4.66807	1.292159e+01	2.831350e+01	8.755813e+02
revenue	4803.0	8.228065e+07	1.628571e+08	0.0	0.00000	1.917000e+07	9.297159e+07	2.787965e+09
runtime	4803.0	1.068754e+02	1.834831e+01	0.0	94.00000	1.030000e+02	1.180000e+02	3.380000e+02
vote_average	4803.0	6.609217e+00	1.194613e+00	0.0	5.60000	6.200000e+00	6.800000e+00	1.000000e+01
vote_count	4803.0	6.902180e+02	1.234586e+03	0.0	54.00000	2.350000e+02	7.370000e+02	1.375200e+04

Observing the data summary we see we have some 0 values in budget, revenue and in runtime. we will treat this as missing data and change it Nans

```
# replace 0 in budget, revenue, runtime
df['budget'] = df['budget'].replace(0, np.NaN)
df['revenue'] = df['revenue'].replace(0, np.NaN)
df['runtime'] = df['runtime'].replace(0, np.NaN)
```

```
df.describe().transpose()
```

	count	mean	std	min	25%	50%	75%	max
budget	3766.0	3.704284e+07	4.264651e+07	0.0	7900000.0000	2.300000e+07	5.661050e+07	3.800000e+08
id	4803.0	5.716548e+04	8.869461e+04	5.0	90145000e+03	1.462900e+04	5.661050e+04	4.594880e+05
popularity	4803.0	2.149230e+01	3.181665e+01	0.0	4.66807	1.292159e+01	2.831350e+01	8.755813e+02
revenue	3766.0	1.170314e+08	1.834831e+08	0.0	1.97367788	1.401651e+08	2.787965e+09	
runtime	4766.0	1.076670e+02	2.074942e+01	14.0	94.00000e+00	1.040000e+02	1.180000e+02	3.380000e+02
vote_average	4803.0	6.609217e+00	1.194613e+00	0.0	5.600000e+00	6.200000e+00	6.800000e+00	1.000000e+01
vote_count	4803.0	6.902180e+02	1.234586e+03	0.0	54.00000e+01	2.350000e+02	7.370000e+02	1.375200e+04

```
# check for duplicates
df.duplicated().sum()
```

```
0
```

Data Cleaning

In this section we will clean our data, remove unused columns, and change data types where applicable.

```
# drop unnecessary Columns
df.drop(['homepage', 'tagline', 'overview'],axis=1, inplace=True)
```

```
# confirm the columns are removed
for c in df.columns:
    print(c)
```

```
budget
genres
id
keywords
original_language
original_title
popularity
production_companies
production_countries
release_date
revenue
runtime
spoken_languages
status
title
vote_average
vote_count
dtype: object
```

```
# Convert the release date to datetime format
df['release_date'] = pd.to_datetime(df['release_date'])
```

```
# confirm
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4803 entries, 0 to 4802
Data columns (total 17 columns):
 # Column Non-Null Count Dtype
 ---  ---
 0 budget 3766 non-null float64
 1 genres 4803 non-null object
 2 id 4803 non-null int64
 3 keywords 4803 non-null object
 4 original_language 4803 non-null object
 5 original_title 4803 non-null object
 6 popularity 4803 non-null float64
 7 production_companies 4803 non-null object
 8 production_countries 4803 non-null object
 9 release_date 4802 non-null datetime64[ns]
10 revenue 3376 non-null float64
11 runtime 4766 non-null float64
12 spoken_languages 4803 non-null object
13 status 4803 non-null object
14 title 4803 non-null object
15 vote_average 4803 non-null float64
16 vote_count 4803 non-null int64
dtypes: datetime64[ns](1), float64(5), int64(2), object(9)
memory usage: 638.0+ KB
```

```
release_date_mode = df['release_date'].mode()
```

```
release_date_mode
```

```
0 2006-01-01
Name: release_date, dtype: datetime64[ns]
```

```
# lets fill the missing release date with mode date
df.fillna({'release_date': '2006-01-01', inplace=True)
```

```
df.isnull().sum()
```

```
budget      1037
genres      0
id           0
keywords     0
original_language  0
original_title  0
popularity   0
production_companies  0
production_countries  0
release_date  0
revenue     1427
runtime      37
spoken_languages  0
status       0
title        0
vote_average  0
vote_count   0
dtype: object
```

```
# similarly lets fill the runtime missing values with the run time mean value
df['runtime'] = df['runtime'].mean()
```

```
runtime
```

```
Out[179]: 107.66072177926983
```

```
# Get the runtime mean to replace the runtime missing value as well.
df.fillna({'runtime': runtime, inplace=True)
```

```
# Check if we have any null values in our data
df.isnull().sum()
```

```
budget      1037
genres      0
id           0
keywords     0
original_language  0
original_title  0
popularity   0
production_companies  0
production_countries  0
release_date  0
revenue     1427
runtime      0
spoken_languages  0
status       0
title        0
vote_average  0
vote_count   0
dtype: object
```

Now that we have cleaned data, we can move the next step EDA and answer our posted questions.

```
# After discussing the structure of the data and any problems that need to be cleaned, perform those cleaning steps in the second part of this section.
```

Exploratory Data Analysis

Question 1. Which top 5 movies have the highest budget?

```
df.columns
```

```
Index(['budget', 'genres', 'id', 'keywords', 'original_language', 'original_title', 'popularity', 'production_companies', 'production_countries', 'release_date', 'revenue', 'runtime', 'spoken_languages', 'status', 'title', 'vote_average', 'vote_count'],
      dtype='object')
```

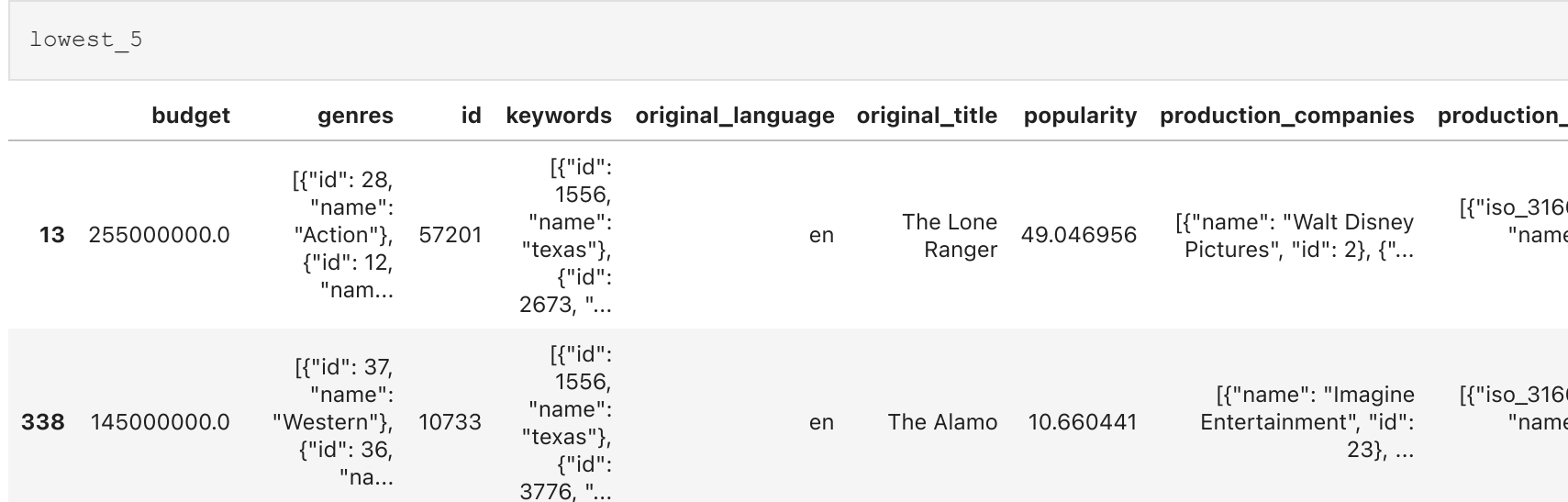
```
# In order to answer this question we will need to find out the top 5 movies with the highest budget
We probably don't need the whole columns to answer this question let take a subset of our data frame
# by creating a new data frame
df1 = df[['release_date', 'title', 'genres', 'runtime', 'popularity', 'budget', 'revenue']]
df1
```

	release_date	title	genres	runtime	popularity	budget	revenue
0	2009-12-10	Avatar	["id": 28, "name": "Action"], ["id": 12, "name": "Romance"]	162.0	150.437577	237000000.0	2.787965e+09
1	2007-05-19	Pirates of the Caribbean: At World's End	["id": 12, "name": "Adventure"], ["id": 14, "name": "Fantasy"], ["id": 14, "name": "Fantasy"]	169.0	139.0782615	300000000.0	9.610000e+08
2	2015-10-26	Spectre	["id": 28, "name": "Action"], ["id": 12, "name": "Romance"]	148.0	107.376788	245000000.0	8.806746e+08
3	2012-07-16	The Dark Knight Rises	["id": 28, "name": "Action"], ["id": 12, "name": "Romance"]	165.0	112.312950	250000000.0	1.084939e+09
4	2012-03-07	John Carter	["id": 28, "name": "Action"], ["id": 12, "name": "Romance"]	132.0	43.926995	260000000.0	2.841391e+08

	release_date	title	genres	runtime	popularity	budget	revenue
4798	1992-09-04	El Mariachi	["id": 28, "name": "Action"], ["id": 80, "name": "Comedy"]	81.0	14.269792	220000.0	2.040920e+06
4799	2011-12-26	Newlyweds	["id": 35, "name": "Comedy"], ["id": 18, "name": "Comedy"]	85.0	0.642552	9000.0	NaN
4800	2013-10-13	Signed, Sealed, Delivered	["id": 35, "name": "Comedy"], ["id": 18, "name": "Comedy"]	120.0	1.444476	NaN	NaN
4801	2012-05-03	Shanghai Calling	[]	98.0	0.857008	NaN	NaN
4802	2005-08-05	My Date with Drew	["id": 99, "name": "Documentary"]	90.0	1.929883	NaN	NaN

4803 rows x 7 columns

```
# Let's get the top 10 movies with the highest budget spent on, sort descending and plot
highest_budget10 = df1.nlargest(n=5, columns='budget')
highest_budget10
```



Observations

"Pirates of the Caribbean: On Stranger Tides", "At World's End", "Avengers: Age of Ultron", "Superman Returns", "John Carter" are top 5 highest based the budget.

Question 2 What are the top 5 highest profit making movies and the 5 lowest movies

```
# calculate Profit for each of the movie
# since line of code is going to create a new column in our data called 'Profit'
df['profit'] = df['revenue'] - df['budget']
```

```
df1.columns
```

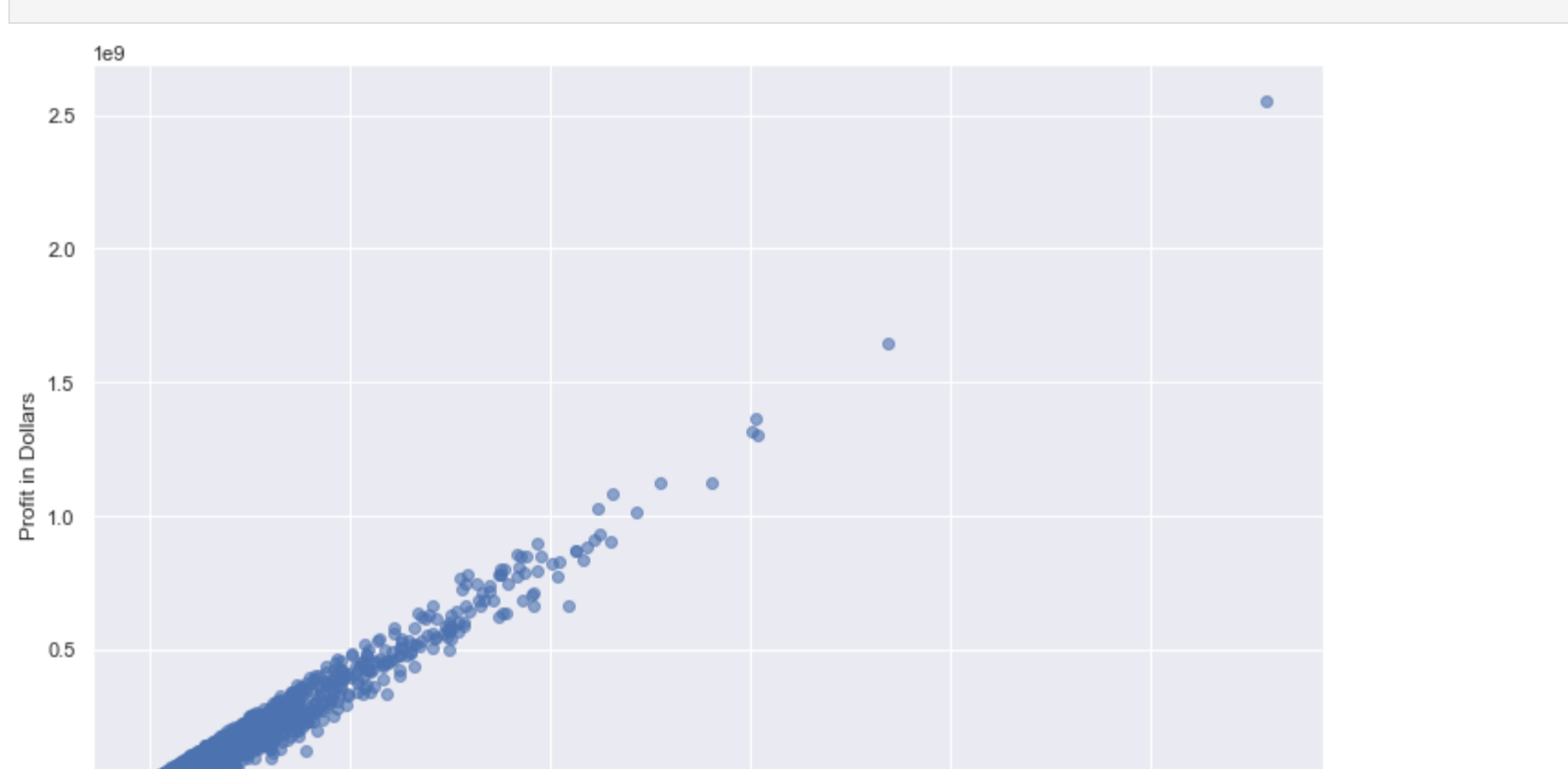
```
Index(['release_date', 'title', 'genres', 'runtime', 'popularity', 'budget', 'revenue', 'profit', 'year', 'released_month'],
      dtype='object')
```

```
# the top 5 highest profit making movies
top_5 = df1.nlargest(5, 'profit')
```

	budget	genres	id	keywords	original_language	original_title	popularity	production_companies	production_co
0	237000000.0	["id": 28, "name": "Action"], ["id": 12, "name": "Romance"]	19995	["id": 1463, "name": "culture clash"], ["id": 1718, "name": "culture clash"]	en	Avatar	150.437577	["name": "Ingenious Film Partners", "id": 289...]	["iso_3166_1", "name": "Sta...
25	200000000.0	["id": 18, "name": "Drama"], ["id": 12, "name": "Adventure"], ["id": 10749, "name": "Fantasy"]	597	["id": 2580, "name": "shipwreck"], ["id": 288, "name": "shipwreck"]	en	Titanic	100.025899	["name": "Paramount Pictures", "id": 4], ["na...	["iso_3166_1", "name": "Sta...
28	150000000.0	["id": 28, "name": "Action"], ["id": 12, "name": "Romance"]	135397	["id": 1299, "name": "monster"], ["id": 1718, "name": "monster"]	en	Jurassic World	418.708552	["name": "Universal Studios", "id": 13], ["na...	["iso_3166_1", "name": "Sta...
44	190000000.0	["id": 28, "name": "Action"], ["id": 12, "name": "Romance"]	168259	["id": 830, "name": "car race"], ["id": 3428, "name": "car race"]	en	Furious 7	102.322217	["name": "Universal Pictures", "id": 33], ["na...	["iso_3166_1", "name": "Sta...
16	220000000.0	["id": 878, "name": "Science Fiction"], ["id": 12, "name": "Adventure"]	24428	["id": 242, "name": "new york"], ["id": 5539, "name": "new york"]	en	The Avengers	144.448633	["name": "Paramount Pictures", "id": 4], ["na...	["iso_3166_1", "name": "Sta...

```
lowest_5 = df1.nsmallest(5, 'profit')
```

```
# plot the top 5 and the 5 lowest profit making movies.
```



```
top_5
```

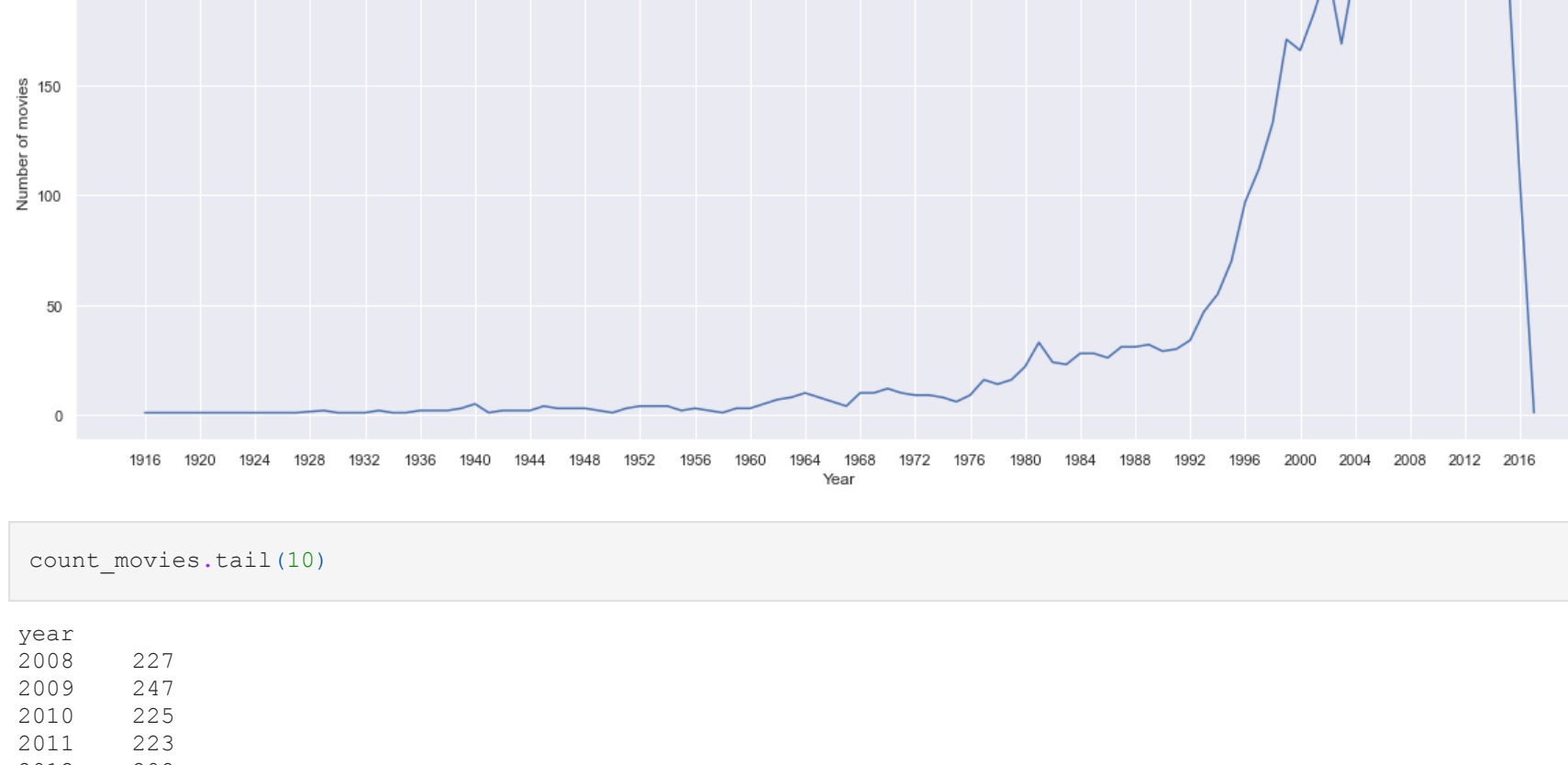
	budget	genres	id	keywords	original_language	original_title	popularity	production_companies	production_cou
0	237000000.0	["id": 28, "name": "Action"], ["id": 12, "name": "Romance"]	19995	["id": 1463, "name": "culture clash"], ["id": 1718, "name": "culture clash"]	en	Avatar	150.437577	["name": "Ingenious Film Partners", "id": 289...]	["iso_3166_1", "name": "Sta...
25	200000000.0	["id": 18, "name": "Drama"], ["id": 12, "name": "Adventure"], ["id": 10749, "name": "Fantasy"]	597	["id": 2580, "name": "shipwreck"], ["id": 288, "name": "shipwreck"]	en	Titanic	100.025899	["name": "Paramount Pictures", "id": 4], ["na...	["iso_3166_1", "name": "Sta...
28	150000000.0	["id": 28, "name": "Action"], ["id": 12, "name": "Romance"]	135397	["id": 1299, "name": "monster"], ["id": 1718, "name": "monster"]	en	Jurassic World	418.708552	["name": "Universal Studios", "id": 13], ["na...	["iso_3166_1", "name": "Sta...
44	190000000.0	["id": 28, "name": "Action"], ["id": 12, "name": "Romance"]	168259	["id": 830, "name": "car race"], ["id": 3428, "name": "car race"]	en	Furious 7	102.322217	["name": "Universal Pictures", "id": 33], ["na...	["iso_3166_1", "name": "Sta...
16	220000000.0	["id": 878, "name": "Science Fiction"], ["id": 12, "name": "Adventure"]	24428	["id": 242, "name": "new york"], ["id": 5539, "name": "new york"]	en	The Avengers	144.448633	["name": "Paramount Pictures", "id": 4], ["na...	["iso_3166_1", "name": "Sta...

Which movies made the hight profit?

Observations

As we see on the Avatar has made the highest profit with profit of 2550985087 USD, followed by the Titanic, Jurassic, Furious, and The Avengers.

```
# Let us also plot the movies making the lowest profit.
# plot the top 5 and the 5 lowest profit making movies.
```



```
lowest_5
```

	budget	genres	id	keywords	original_language	original_title	popularity	production_companies	production_cou
13	255000000.0	["id": 37, "name": "Action"], ["id": 12, "name": "Adventure"], ["id": 36, "name": "Adventure"]	27201	["id": 1556, "name": "teen"], ["id": 2873, "name": "teen"]	en	The Lone Ranger	49.046956	["name": "Walt Disney Pictures", "id": 2], ["na...	["iso_3166_1", "name": "i Stat...
338	145000000.0	["id": 37, "name": "Western"], ["id": 36, "name": "Western"], ["id": 16, "name": "Western"]	10793	["id": 1556, "name": "teen"], ["id": 1718, "name": "teen"]	en	The Alamo	10.660441	["name": "Imagine Entertainment", "id": 23], ["na...	["iso_3166_1", "name": "i Stat...
141	150000000.0	["id": 12, "name": "Adventure"], ["id": 16, "name": "Adventure"]	60321	["id": 5202, "name": "boy"], ["id": 5961, "name": "boy"]	en	Mars Needs Moms	12.362599	["name": "Walt Disney Animation Studios", "id": ...]	["iso_3166_1", "name": "i Stat...
208	160000000.0	["id": 12, "name": "Adventure"], ["id": 16, "name": "Adventure"]	1911	["id": 616, "name": "witch"], ["id": 1964, "name": "witch"]	en	The 13th Warrior	27.220157	["name": "Touchstone Pictures", "id": 9195]	["iso_3166_1", "name": "i Stat...
311	100000000.0	["id": 28, "name": "Action"], ["id": 35, "name": "Action"], ["id": 16, "name": "Action"]	11692	["id": 305, "name": "moon"], ["id": 585, "name": "moon"]	en	The Adventures of Pluto Nash	12.092241	["name": "Village Roadshow Pictures", "id": 7...]	["iso_3166_1", "name": "Austri...

Observations

Observing the graph we see that The Lone Ranger is the biggest loser, making the lowest profit of -16571009 USD; followed by: The Alamo, Mars Needs Moms, the 13th warrior, the adventures of Pluto Nash.

```
# Let us check out if there is a relationship between budget and Revenue
plt.figure(figsize=(12, 8))
plt.scatter(df['budget'], df['revenue'], alpha=0.7,
```



Observing the data points there is no relationship between profit and profits



As we can observe the plot, september has the highest number of release movies.

Which Month Made The Highest average profit?

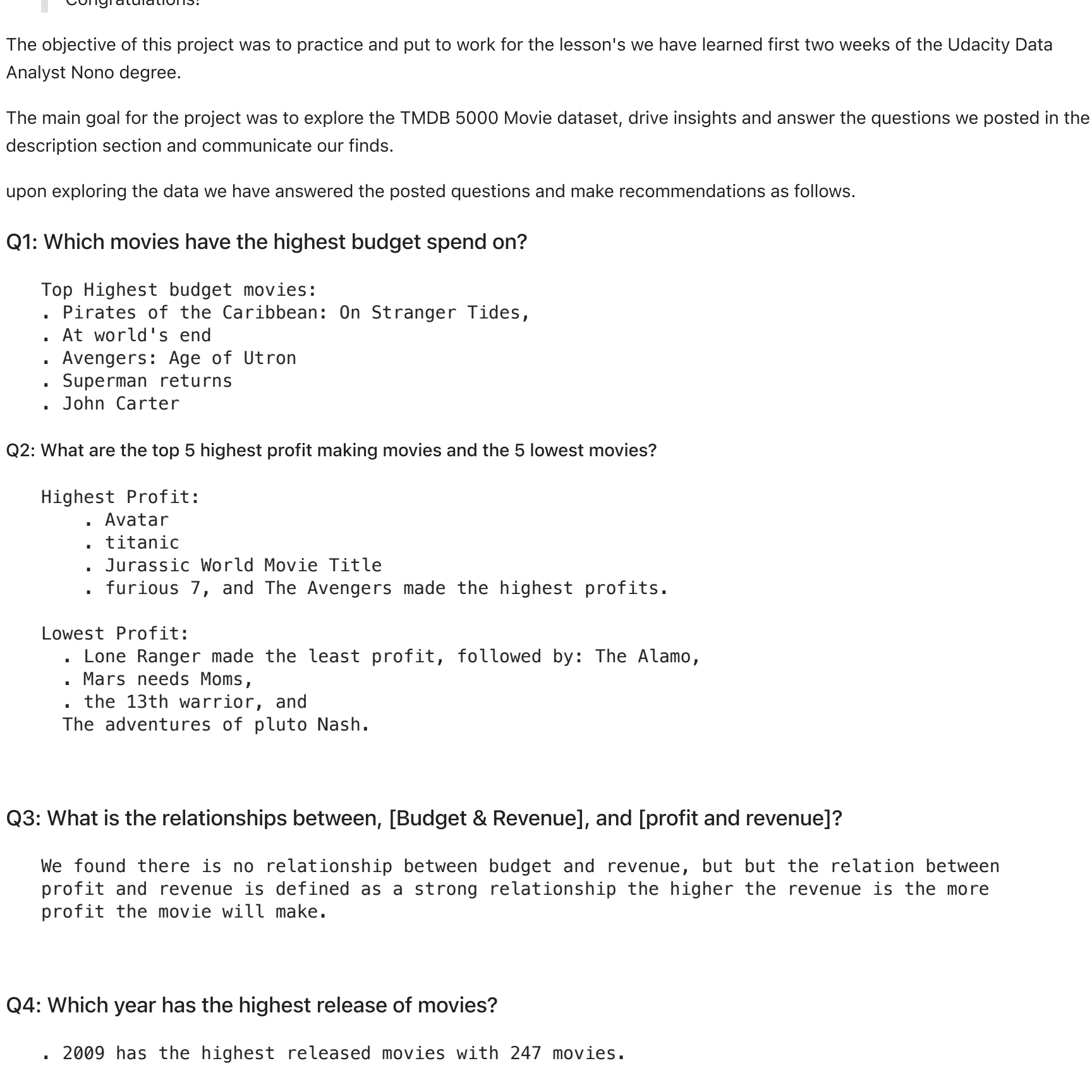
```
df2 = df[["title", "released_month", "revenue", "profit"]]
df2
```

	title	released_month	revenue	profit
0	Avatar	December	2.787965e+09	2.550965e+09
1	Pirates of the Caribbean: At World's End	May	9.610000e+08	6.610000e+08
2	Spectre	October	8.806746e+08	6.356746e+08
3	The Dark Knight Rises	July	1.084939e+09	8.349391e+08
4	John Carter	March	2.841391e+08	2.413910e+07
...
4798	El Mariachi	September	2.040920e+06	1.820920e+06
4799	Newlyweds	December	NaN	NaN
4800	Signed, Sealed, Delivered	October	NaN	NaN
4801	Shanghai Calling	May	NaN	NaN
4802	My Date with Drew	August	NaN	NaN

4803 rows x 4 columns

```
profit_per_month = df2.groupby("released_month")["profit"].mean()
```

```
# plot the result
plt.figure(figsize=(20, 8))
profit_per_month.plot(kind="bar")
plt.xlabel("Movie released Month")
plt.ylabel("Profits made per month")
plt.show()
```



Insights

Movies released on July made the highest profits, followed by Jun, May, thus we can say that those months are the best time to release movies.

Conclusions

Submitting your Project

Tip: Before you submit your project, you need to create a .html or .pdf version of this notebook in the workspace here. To do that, run the code cell below. If it worked correctly, you should get a return code of 0, and you should see the generated .html file in the workspace directory (click on the orange Jupyter icon in the upper left).

Tip: Alternatively, you can download this report as .html via the **File > Download as** submenu, and then manually upload it into the workspace directory by clicking on the orange Jupyter icon in the upper left, then using the Upload button.

Tip: Once you've done this, you can submit your project by clicking on the "Submit Project" button in the lower right here. This will create and submit a zip file with this .ipynb doc and the .html or .pdf version you created. Congratulations!

The objective of this project was to practice and put to work for the lesson's we have learned first two weeks of the Udacity Data Analyst Nano degree.

The main goal for the project was to explore the TMDb 5000 Movie dataset, drive insights and answer the questions we posted in the description section and communicate our finds.

upon exploring the data we have answered the posted questions and make recommendations as follows.

Q1: Which movies have the highest budget spend on?

- Top Highest budget movies:
- Pirates of the Caribbean: On Stranger Tides,
 - At world's end
 - Avengers: Age of Ultron
 - Superman returns
 - John Carter

Q2: What are the top 5 highest profit making movies and the 5 lowest movies?

- Highest Profit:
- Avatar
 - Titanic
 - Jurassic World Movie Title
 - Furious 7, and The Avengers made the highest profits.

- Lowest Profit:
- Lone Ranger made the least profit, followed by: The Alamo,
 - Mars Needs Moms,
 - the 13th warrior, and
 - The adventures of Pluto Nash.

Q3: What is the relationships between, [Budget & Revenue], and [profit and revenue]?

We found there is no relationship between budget and revenue, but but the relation between profit and revenue is defined as a strong relationship the higher the revenue is the more profit the movie will make.

Q4: Which year has the highest release of movies?

- 2009 has the highest released movies with 247 movies.

Q5: Which Month has the highest release movies?

- September

Q6: Which months made the highest average profits?

- Movies released on July made the highest profits, followed by Jun, May. So practically movies released during summer makes the most profits, so our recommendation to release movies to during this months to make the most profits.

```
In [247]: from subprocess import call
          call(['python', '-m', 'nbconvert', 'Investigate_a_Dataset.ipynb'])
```

Out[247]: 255

In []: