

INTRODUCTION

The dataset we would be exploring in this project is tmdb-movies which consists of 21 features and 10866 observations, during the course of the work we would be answering the following questions

- 1. Which year had the highest number of movies released?
- 1. Which month had the highest number of released movies?
- 1. Which month generated more revenue for the movie industry?
- 1. Does the number of movies released translate to more revenue?
- 1. Does vote_count affect revenue generated from a movie?
- 1. Which day accounts for more movies released?
- 1. Which movie title generated the highest revenue?
- 1. Which movie title has the longest runtime?
- 1. Does a longer runtime indicate more revenue?
- 1. Which movie has the highest runtime?

IMPORTING LIBRARIES

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import datetime
```

LOADING DATASET

```
In [2]: df = pd.read_csv('tmdb-movies.csv')
df.head()
```

	id	imdb_id	popularity	budget	revenue	original_title	cast	homepage	director
Out[2]:	0	135397	td369610	32.985763	150000000	1513528810	Jurassic World	Chris Pratt Byrne Dallas Howard Jeffrey Huang V.L. Kwan	Colin Trevorrow
	1	76341	tt1392190	28.419936	150000000	378436354	Mad Max: Fury Road	Tom Hardy Charlize Theron Hugh Keays-Byrne Nic...	George Miller
	2	262500	tt2908446	13.112507	110000000	295238201	Insurgent	Shailene Woodley Theo James Kate Winslet Jesse...	Robert Schwentke
	3	140607	tt2488496	11.173104	200000000	2068778225	Star Wars: The Force Awakens	Harrison Ford Mark Hamill Carrie Fisher Adam D...	J.J. Abrams
	4	168259	tt2820852	9.335014	190000000	1506249360	Furious 7	Vin Diesel Paul Walker Jason Statham Michelle...	James Wan

5 rows x 21 columns

DATA WRANGLING

```
In [3]: #Checking to see the number of features and observations
df.shape
```

Out[3]: (10866, 21)

```
In [4]: #Checking to see the names of all the columns
df.columns
```

```
Out[4]: Index(['id', 'imdb_id', 'popularity', 'budget', 'revenue', 'original_title',
          'cast', 'homepage', 'director', 'tagline', 'keywords', 'overview',
          'runtime', 'genres', 'production_companies', 'release_date',
          'vote_count', 'vote_average', 'release_year', 'budget_adj',
          'revenue_adj'],
          dtype='object')
```

```
In [5]: #Checking to see the data type of the dataset
df.dtypes
```

```
Out[5]: id                int64
imdb_id              object
popularity           float64
budget              int64
revenue             int64
original_title       object
cast                object
homepage             object
director             object
tagline             object
keywords            object
overview            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 [6]: #Checking to see a brief information of the dataset
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 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   int64
 4   revenue             10866 non-null   int64
 5   original_title       10866 non-null   object
 6   cast                10790 non-null   object
 7   homepage             2936 non-null   object
 8   director             10822 non-null   object
 9   tagline             8042 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
```

```
In [7]: #Checking for null values
df.isnull().sum()
```

```
Out[7]: id                0
imdb_id              10
popularity            0
budget               0
revenue              0
original_title        0
cast                 10
homepage             7930
director              44
tagline              2824
keywords             1493
overview              23
runtime              0
production_companies 1030
release_date         0
vote_count           0
vote_average         0
release_year         0
budget_adj           0
revenue_adj          0
dtype: int64
```

```
In [8]: #Checking to see a summary statistics of the dataset
df.describe()
```

```
Out[8]:
```

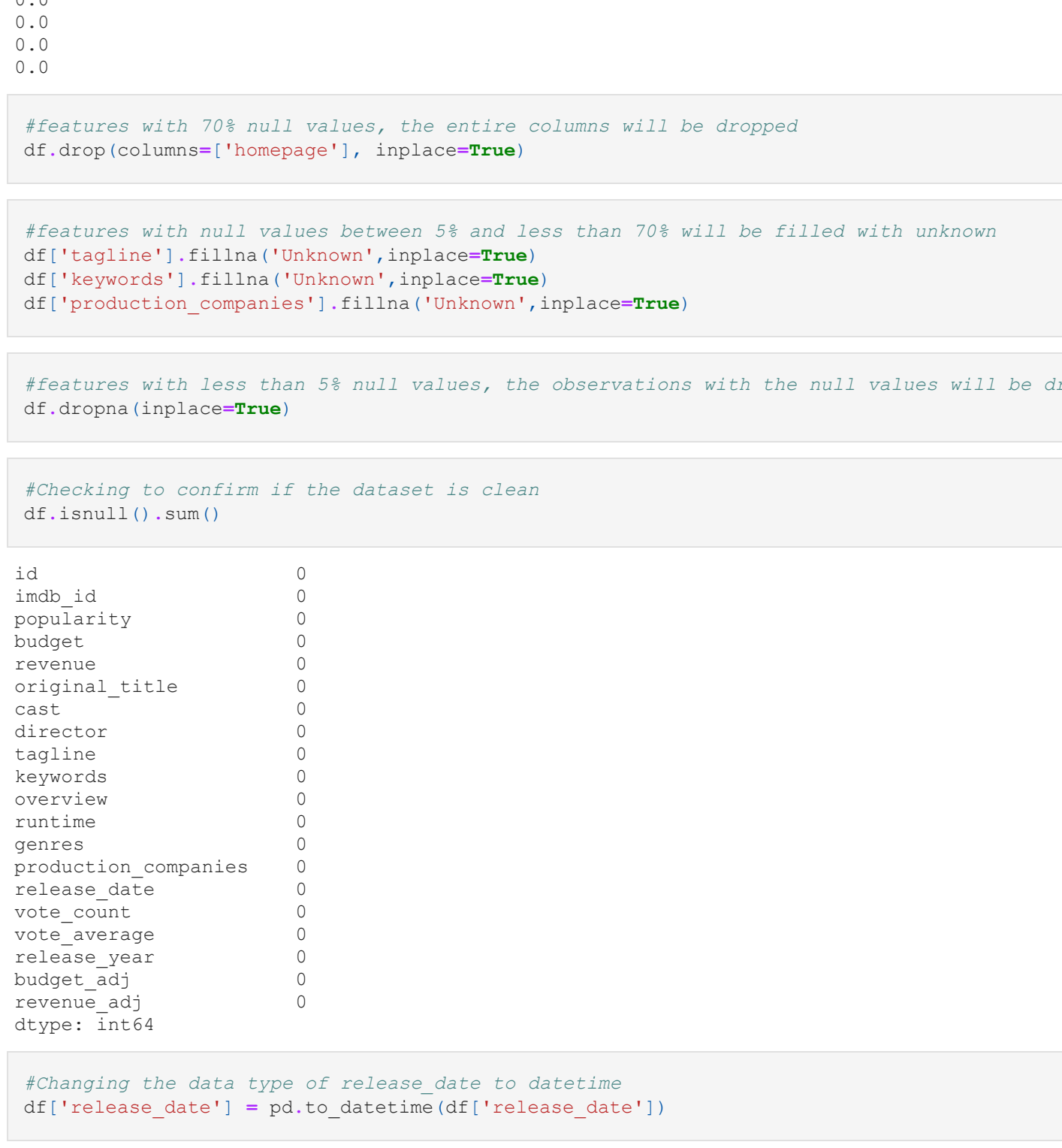
	id	popularity	budget	revenue	runtime	vote_count	vote_average	release_year	budget_adj	revenue_adj
count	10866.000000	10866.000000	1.086600e+04	1.086600e+04	108.66000000	10866.000000	10866.000000	10866.000000	1.086600e+04	1.086600e+04
mean	66064.177434	0.646441	1.462570e+07	1.388233e+08	102.070863	217.389748	5.974922	2001.322658	1.755104e+07	5.1364
std	92130.136561	0.000185	3.091321e+07	1.700035e+08	31.381405	575.619058	0.935142	12.812941	3.430616e+07	1.4463
min	5.000000	0.000000	0.000000e+00	0.000000e+00	0.000000	10.000000	0.935142	1960.000000	0.000000e+00	0.0000
25%	10596250000	0.207983	0.000000e+00	0.000000e+00	90.000000	17.000000	5.400000	1995.000000	0.000000e+00	0.0000
50%	20669.000000	0.383856	0.000000e+00	0.000000e+00	99.000000	38.000000	6.000000	2006.000000	0.000000e+00	0.0000
75%	175810.000000	0.713817	1.500000e+07	2.430000e+07	111.000000	145.750000	6.600000	2011.000000	2.085325e+07	3.8671
max	417859.000000	32.985763	4.250000e+08	2.781500e+09	900.000000	9767.000000	9.200000	2015.000000	4.250000e+08	2.8289

```
In [9]: #Checking to see correlation between features
df.corr()
```

```
Out[9]:
```

	id	popularity	budget	revenue	runtime	vote_count	vote_average	release_year	budget_adj	revenue_adj
id	1.000000	-0.014350	-0.141351	-0.099227	-0.088360	-0.035551	-0.058363	0.511364	-0.189015	-0.138477
popularity	-0.014350	1.000000	0.545472	0.663358	0.139033	0.808028	0.209511	0.089801	0.513550	0.609083
budget	-0.141351	0.545472	1.000000	0.734901	0.190333	0.632702	0.081014	0.115931	0.968963	0.622505
revenue	-0.099227	0.663358	0.734901	1.000000	0.162838	0.791175	0.172564	0.057048	0.706427	0.919110
runtime	-0.088360	0.139033	0.191283	0.162838	1.000000	0.162278	0.156835	-0.117204	0.221114	0.175766
vote_count	-0.035551	0.808028	0.632702	0.791175	0.162278	1.000000	0.253823	0.107948	0.587051	0.707942
vote_average	-0.058363	0.209511	0.081014	0.172564	0.156835	0.253823	1.000000	-0.117632	0.091939	0.193085
release_year	0.511364	0.089801	0.115931	0.057048	-0.117632	-0.117632	-0.117632	1.000000	0.016793	-0.066256
budget_adj	-0.189015	0.513550	0.968963	0.706427	0.221114	0.587051	0.093039	0.016793	1.000000	0.646607
revenue_adj	-0.138477	0.609083	0.622505	0.919110	0.175766	0.707942	0.193085	-0.066256	0.646607	1.000000

```
In [10]: #Visualizing correlation between features
plt.figure(figsize=(12,8))
sns.heatmap(df.corr(), cmap = 'coolwarm', annot = True)
plt.show()
```



```
In [11]: #Checking the number of unique features
df.nunique()
```

```
Out[11]: id                10866
imdb_id              10865
popularity           10814
budget              10571
revenue             10719
original_title       10571
cast                10719
homepage             2896
director             5067
tagline             2997
keywords            1884
overview            10847
runtime             247
genres              2039
production_companies 1445
release_date         5909
vote_count           1289
vote_average         72
release_year         56
budget_adj           2614
revenue_adj          4840
dtype: int64
```

DATA CLEANING

```
In [12]: null_val = df.isnull().sum()
```

```
In [13]: #Determining percentage null values in dataset
for item in null_val:
    print((item/10866)*100)
```

```
0.0
0.09203018590097553
0.0
0.0
0.0
0.0
0.6994234128474139
0.279819474947359
0.40493281796429226
25.989924498435313
13.740180755013645
0.03881207436039021
0.2116694275724366
9.479109147800479
0.0
0.0
0.0
0.0
0.0
```

```
In [14]: #Features with 70% null values, the entire column will be dropped
df.drop(columns=['homepage'], inplace=True)
```

```
In [15]: #Features with null values between 5% and less than 70% will be filled with unknown
df['tagline'].fillna('Unknown',inplace=True)
df['keywords'].fillna('Unknown',inplace=True)
df['production_companies'].fillna('Unknown',inplace=True)
```

```
In [16]: #Features with less than 5% null values, the observations with the null values will be dropped since the percent is less than 5%
df.dropna(inplace=True)
```

```
In [17]: #Checking to confirm if the dataset is clean
df.isnull().sum()
```

```
Out[17]: id                0
imdb_id              10
popularity            0
budget               0
revenue              0
original_title        0
cast                 10
director              44
tagline              2824
keywords             1493
overview              23
runtime              0
production_companies 1030
release_date         0
vote_count           0
vote_average         0
release_year         0
budget_adj           0
revenue_adj          0
dtype: int64
```

```
In [18]: #Changing the data type of release_date to datetime
df['release_date'] = df['release_date'].astype('datetime64[ns]')
```

```
In [19]: #Confirming the change in release_date data type
df.dtypes
```

```
Out[19]: id                int64
imdb_id              object
popularity           float64
budget              int64
revenue             int64
original_title       object
cast                object
director             object
tagline             object
keywords            object
overview            object
runtime             int64
genres              object
production_companies object
release_date         datetime64[ns]
vote_count           int64
vote_average         float64
release_year         int64
budget_adj           float64
revenue_adj          float64
dtype: object
```

```
In [20]: #Creating new feature release_day from release_date
df['release_day'] = df['release_date'].dt.weekday
```

```
In [21]: #Converting the numeric variables in release_day to string
mapper = {'0':'Monday','1':'Tuesday','2':'Wednesday','3':'Thursday','4':'Friday','5':'Saturday'}
df['release_day'] = df['release_day'].map(mapper)
```

```
In [22]: #Confirming the change
df.head()
```

```
#Removing duplicate values
df.drop_duplicates(inplace=True)
sum(df.duplicated())
```

0

EXPLORATORY DATA ANALYSIS

Which year had the most number of movies released?

```
pd.DataFrame(df[["release_year"]].value_counts(ascending=False)[0:16])
```

	release_year
2014	682
2013	648
2015	617
2012	574
2011	530
2009	523
2008	486
2010	475
2007	430
2006	404
2005	361
2004	307
2003	281
2002	264
2001	241
1999	224

5 rows x 21 columns

```
In [23]: #Creating new feature release_month from release_date
df['release_month'] = df['release_date'].dt.month
```

```
In [24]: #Converting the numeric variables in release_month to string
mapper = {'1':'January','2':'February','3':'March','4':'April','5':'May','6':'June','7':'July','8':'August','9':'September','10':'October','11':'November'}
df['release_month'] = df['release_month'].map(mapper)
```

```
In [25]: #Confirming the change in release_month
df.head()
```

Which month the highest number of released movies?

```
pd.DataFrame(df['release_month'].value_counts())
```

	release_month
September	1318
October	1137
December	974
August	906
January	888
June	822
March	813

5 rows x 22 columns

```
In [26]: #Checking for duplicates in the dataset
sum(df.duplicated())
```

Out[26]: 1

```
In [27]: #Removing duplicate values
sum(df.duplicated(inplace=True))
```

Out[27]: 0

EXPLORATORY DATA ANALYSIS

Which year had the most number of movies released?

```
In [28]: pd.DataFrame(df['release_year'].value_counts(ascending=False))
```

```
Out[28]:
```

```
release_year
2014    682
2013    648
2015    617
2012    574
2011    530
2009    523
2008    486
2010    475
2007    430
2006    404
2005    361
2004    307
2003    281
2002    264
2001    241
1999    224
```

```
In [29]: #Visualizing the result
plt.figure(figsize=(12,10))
df['release_year'].value_counts().plot(kind='bar')
plt.title('TOP 30 MOVIE YEARS')
```

```
Out[29]: Text(0.5, 1.0, 'TOP 30 MOVIE YEARS')
```



We can see that more movies were released in the year 2014 with a count of 682 movies released

Which month had the highest number of released movies?

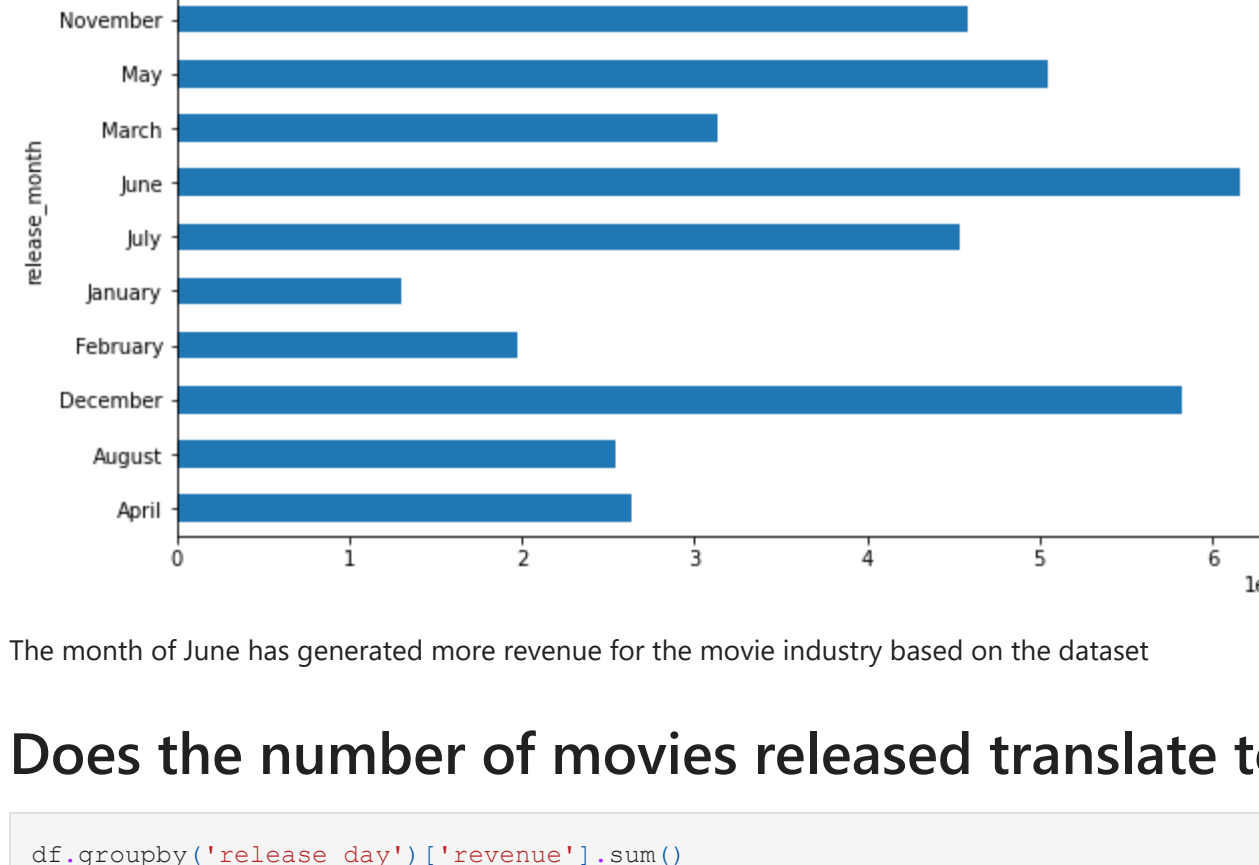
```
In [30]: pd.DataFrame(df['release_month'].value_counts())
```

```
Out[30]:
```

```
release_month
September    1318
October      1137
December     974
August       906
January      888
June         823
March        813
May          807
November     798
April        791
July         790
February     680
```

```
In [31]: df['release_month'].value_counts(ascending=True).plot(kind='bar', figsize=(10,6))
```

```
Out[31]: <AxesSubplot>
```



The month of September had the highest number of movies released

Which day accounts for more movie release?

```
In [32]: df['release_day'].value_counts()
```

```
Out[32]: Friday      4225
Thursday      1674
Wednesday     1513
Tuesday       1189
Saturday       782
Monday         653
Name: release_day, dtype: int64
```

```
In [33]: pd.DataFrame(df['release_day'].value_counts(ascending=False, normalize=True)*100)
```

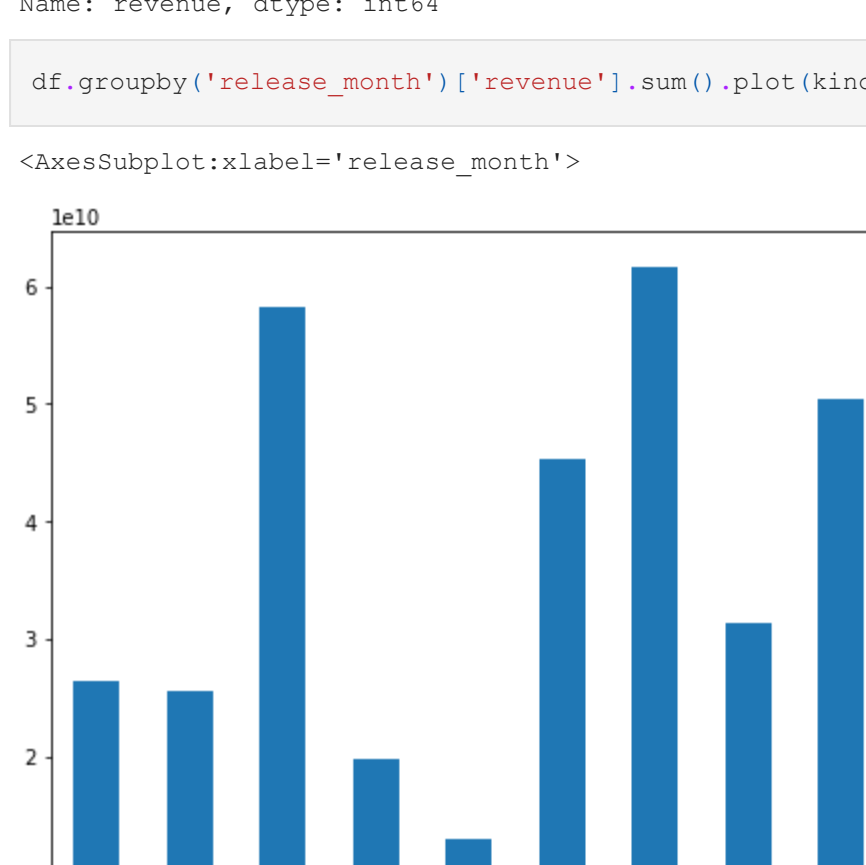
```
Out[33]:
```

```
release_day
Friday      42.098446
Thursday    16.679952
Wednesday   15.075727
Tuesday     11.847350
Saturday     7.791949
Monday      6.506576
```

```
In [34]: #Visualizing the result
release_day = pd.DataFrame(df.groupby('release_day')['release_day'].size().sort_values(ascending=False))
release_day.columns = ['Count']
```

```
#Define Seaborn color palette to use
colors = sns.color_palette('pastel')[0:6]
```

```
#Create pie chart
plt.figure(figsize=(8,8))
plt.pie(release_day['Count'], labels = labels, colors = colors, autopct='%0.0f%%')
plt.show()
```



More movies were released on Friday which account for about 42%

Which month has generated more revenue in the movie industry?

```
In [35]: df.groupby('release_month')['revenue'].sum()
```

```
Out[35]: release_month
April          26392671428
August         25477769000
December       59211863204
February       19793785507
January        12968479332
July           45337825481
June           61660585217
March          31394443375
November       50454865815
May            45896250786
October        29353709677
September      25731469855
Name: revenue, dtype: int64
```

```
In [36]: df.groupby('release_month')['revenue'].sum().plot(kind='bar', figsize=(10,6))
```

```
Out[36]: <AxesSubplot>
```


The month of June has generated more revenue for the movie industry based on the dataset

Does the number of movies released translate to more revenue?

```
In [37]: df.groupby('release_day')['revenue'].sum()
```

```
Out[37]: release_day
Friday      128852378314
Monday      24995065456
Saturday    13060362270
Thursday    9662729340
Tuesday     49507764714
Wednesday   105661470111
Name: revenue, dtype: int64
```

```
In [38]: df.groupby('release_day')['revenue'].sum().plot(kind='bar', figsize=(8,6))
```

```
Out[38]: <AxesSubplot>
```



```
In [39]: df.groupby('release_month')['revenue'].sum()
```

```
Out[39]: release_month
April          26392671428
August         25477769000
December       59211863204
February       19793785507
January        12968479332
July           45337825481
June           61660585217
March          31394443375
November       50454865815
May            45896250786
October        29353709677
September      25731469855
Name: revenue, dtype: int64</
```



```
In [43]: df['vote_count'].max()
Out[43]: 9767
In [44]: df['vote_count'].min()
Out[44]: 10
```

```
In [45]: df[df['vote_count'] == 9767]
Out[45]:
```

	id	imdb_id	popularity	budget	revenue	original_title	cast	director	tagline	keywords	...
1919	27205	tt1375666	9.363643	16000000	82550000	Inception	Leonardo DiCaprioJoseph Gordon-LevittEllen P...	Christopher Nolan	Your mind is the scene of the crime.	love dream sleep subconscious	lo...

1 rows × 22 columns

```
In [46]: df[df['revenue'] == 2781505847]
Out[46]:
```

	id	imdb_id	popularity	budget	revenue	original_title	cast	director	tagline	keywords	...
1386	19995	tt0499549	9.432768	237000000	2781505847	Avatar	WorthingtonZoe SaldanaSigourney WeaverS...	James Cameron	Enter the World of Pandora.	culture clash future space war space colony jo...	Action Ad...

1 rows × 22 columns

```
In [47]: df[df['vote_count'] == 10]
Out[47]:
```

	id	imdb_id	popularity	budget	revenue	original_title	cast	director	tagline	keywords	...
240	363689	tt4229298	0.532700	0	0	The Unspoken	Jodelle FerlandSammy SijmNeil McDonoughMa...	Sheldon Wilson	Unknown		Unknown
259	360603	tt5133572	0.476341	0	0	Crown for Christmas	Danica McKellarRupert Penry-JonesEllie Bote...	Alex Zamm	Unknown		Unknown
298	354220	tt3826866	0.370258	0	0	The Girl in the Photographs	Kai PemaClaudia LeeKenney WormaelDoby Herring...	Nick Simon	Unknown	serial killer fied feet tied up while barefoot	
345	361263	tt4505830	0.342530	0	0	Jesse Stone: Lost in Paradise	Tom SelleckWilliam DevaneGloria RubenJill Henning...	Robert Harmon	Unknown		Unknown
386	316885	tt4180576	0.291244	0	0	Bridgeend	Harshad MurraySteven WaddingtonAdrian Rawlins...	Jeppe R��nde	Unknown	suicide adolescence based on true story magic ...	
10802	54000	tt0077539	0.227338	6727000	0	Fedora	William HoldenMartha KellnerKilledgeard KnellJo...	Billy Wilder	been a habit of hers for so long tha...	hollywood actor down on his	
10808	40060	tt0077904	0.129123	0	0	The Manito	Tony CurtisMichael AnsaraJoan StrasserKevin StrasserBer...	William Girdler	Evil does not die!�� it waits to be re-born!	female nudity based on novel nudity barot card...	
10809	31428	tt0078295	0.128231	0	0	Someone's Watching Me!	Lauren HuttonDavid AnspaughKevin BarbeauJill...	John Carpenter	Unknown	stalker suspense stalking made for television	
10812	31948	tt0077629	0.318883	0	5438927	Gray Lady Down	Charlton HestonDavid CarradineTacy KeachJill...	David Greene	Trapped underwater... with time running out.	submarine drowning based on novel rescue disaster	
10857	5921	tt0060748	0.131378	0	0	Nevada Smith	Steve McQueenKarl MaldenBrian KeithArthur K...	Henry Hathaway	Some called him savage- and some called him sa...	revenge revenge native american wild west ha...	

489 rows × 22 columns

```
In [48]: df[df['revenue'] == 0]
Out[48]:
```

	id	imdb_id	popularity	budget	revenue	original_title	cast	director	tagline	keywords	...
48	265208	tt231253	2.932340	30000000	0	Wild Card	Jason StathamMichael Jai WhiteToby Kebbell...	Simon West	Never bet against a killer hand.	gambling bo	
67	334074	tt3247714	2.331636	20000000	0	Survivor	Pierce BrosnanMilla JovovichDylan McDermott...	James McTeigue	His Next Target is Now Hunting Him	evil re sh downto t	
74	347096	tt3478232	2.165433	0	0	Mythica: The Darkspore	Melanie StoneKevin SorboAdam JohnsonJake St...	Anne K. Black	Unknown	sword magic sorcery recrea...	
75	308369	tt2562496	2.141506	0	0	Me and Earl and the Dying Girl	Thomas MannRJ CylerOlivia CookeConor Brit...	Alfonso Gomez-Rejon	A Little Friendship Never Killed Anyone.	school leukemi	
92	370687	tt3608646	1.876037	0	0	Mythica: The Necromancer	Melanie StoneAdam JohnsonKevin SorboNicola ...	A. Todd Smith	Unknown	sword magic sorcery recrea...	
10861	21	tt0060371	0.080598	0	0	The Endless Summer	Michael HymonRobert August and "Toby Ho B...	Bruce Brown	Unknown	surfer s	
10862	20379	tt0060472	0.065543	0	0	Grand Prix	James GarnerEva Marie SaintTyesa MontanaTol...	John Frankenheimer	Cinemas sweeps YOU into a drama of speed and ...	car race	
10863	39768	tt0060161	0.065141	0	0	Bergris Avtomobilya	Innokenty SmolturnovskiyOleg EfremovGeorgi Z...	Eldar Ryazanov	Unknown	car tr	
10864	21449	tt0061177	0.064317	0	0	What's Up, Tiger Lily?	Tat��y�� MihashiAkiko WakabayashiMie HamaJo��...	Woody Allen	WOODY ALLEN STRIKES BACK!		
10865	22293	tt0060666	0.035919	19000	0	Manos: The Hands of Fate	Harold P. WarrenTom NeymanJohn ReynoldsDian...	Harold P. Warren	It's Shocking! It's Beyond Your Imagination!	fire gun drives s	

5881 rows × 22 columns

The vote_count does not affect the revenue as we can see that the minimum revenue is 0 which has quite a large number of vote_count while the minimum value for vote_count is 10 some having revenue of 0 and some having a large amount as revenue

Which movie generated the highest amount of revenue?

```
In [49]: df['revenue'].max()
Out[49]: 2781505847
In [50]: df[df['revenue'] == 2781505847]
Out[50]:
```

	id	imdb_id	popularity	budget	revenue	original_title	cast	director	tagline	keywords	...
1386	19995	tt0499549	9.432768	237000000	2781505847	Avatar	WorthingtonZoe SaldanaSigourney WeaverS...	James Cameron	Enter the World of Pandora.	culture clash future space war space colony jo...	Action Ad...

1 rows × 22 columns

The movie Avatar generated the highest amount of revenue with the amount 2781505847

Does longer runtime indicate more revenue?

```
In [51]: plt.figure(figsize=(10,8))
sns.set(style='whitegrid')
ax = sns.lmplot(x='runtime', y='revenue', data=df)
<Figure size 720x576 with 0 Axes>
```

The longer the runtime doesn't indicate more revenue

What is the longest runtime of a movie?

```
In [52]: df['runtime'].max()
Out[52]: 900
The longest runtime for a movie is 900 minutes
```

Which movie has the longest runtime?

```
In [53]: df[df['runtime'] == 900]
Out[53]:
```

	id	imdb_id	popularity	budget	revenue	original_title	cast	director	tagline	keywords	...	genres	produ...
3894	125336	tt2044056	0.006925	0	0	The Story of Film: An Odyssey	Mark CousinsJean-Michel FrodonCarla Beauchamp...	Mark Cousins	Unknown	cinema nouvelle vague hindi cinema cinema nov...	Documentary		

1 rows × 22 columns

The Story of Film:An Odyssey has the longest runtime

SUMMARY OF FINDINGS AND CONCLUSION

- More movies were released in the year 2014.
- The month of September recorded the highest number of movies released.
- 42% of the movies were released on Friday.
- The month of June generated more revenue for the movie industry based on the dataset.
- From the release_day, more movies were released on Friday which translated to more revenue, but also more movies were released in the month of September but this didn't translate to more revenue hence we can't say from the analysis if more movies released can translate to more revenue.
- The vote_count does not affect the revenue.
- The longer the runtime doesn't indicate more revenue.
- The longest runtime for a movie is 900 minutes.
- The Story of Film:An Odyssey had the longest runtime.
- The movie Avatar generated the highest amount of revenue.

Conclusion: Basic questions were just answered in this project based on the scope of the work, machine learning model can be built to better predict how more features can affect the target variable.