

Project: TMDB Movie Dataset Analysis

Table of Contents

- [Introduction](#)
- [Data Wrangling](#)
- [Exploratory Data Analysis](#)
- [Conclusions](#)

Introduction

Dataset

I choose the TMDB movie dataset, the dataset contains information about 10,000 movies, including 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.

Questions:

- What are the 10 movies with the highest production budget?
- What are the 10 movies with the highest revenue?
- What are the top 10 most profitable movies in the market?
- What are the top ten movies with the highest popularity?
- What are the top ten movies with the highest vote average?
- What are the 10 movies with the highest runtime?
- What is the average movie runtime over time?
- What are the average movie profits over time?
- How many movies per year?
- In which year has the highest or lowest production of movies?
- What are the top 10 directors in 2014, according to the vote?
- How many movies per genre?
- How many movies per director?
- What are the top 10 most directors who have maximum movies?
- What are the most profitable months?

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

```
In [2]: from IPython.core.display import display
sns.set(rc={'figure.figsize':(14,6)})
```

Data Wrangling

```
In [3]: df = pd.read_csv('content/tmdb-movies.csv')
df.head()
```

	id	imdb_id	popularity	budget	revenue	original_title	cast	homepage	direct
0	135397	tt0369610	32.985763	15000000	1513528810	Jurassic World	Chris PrattBryce Dallas HowardIrfan Khan[...]	http://www.jurassicworld.com/	C Trevon
1	76341	tt1392190	28.419936	15000000	378436354	Mad Max: Fury Road	Tom HardyCharize TheronHugh Keays-Byrne[...]	http://www.madmaxmovie.com/	Geo M
2	262500	tt2098446	13.112507	11000000	295236201	Insurgent	Shailene WoodleyTheo JamesKate Winslet[...]	http://www.thedivergentseries.movie/insurgent	Rot Schwen
3	140067	tt2488496	11.173104	20000000	2068178225	Star Wars: The Force Awakens	Harrison FordMark HamillAdam Driver[...]	http://www.starwars.com/films/star-wars-episode...	Abra
4	168259	tt2820852	9.335014	19000000	1508249360	Furious 7	Vin DieselPaul WalkerJason Statham[...]	http://www.furious7.com/	Jan Y

```
In [4]: df.shape
Out[4]: (10866, 21)
```

```
In [5]: df.describe()
Out[5]:
```

	id	popularity	budget	revenue	runtime	vote_count	vote_average	release_year	budget_adj
count	10866.000000	10866.000000	1.086600e+04	1.086600e+04	10866.000000	10866.000000	10866.000000	10866.000000	1.086600e+04
mean	66064.177434	0.646441	1.462537e+07	3.986232e+08	102.070863	217.389748	5.974922	2001.322658	1.756106e+07
std	92130.136561	0.001085	3.091321e+07	1.170035e+08	31.381455	575.619058	0.935142	12.812941	3.430616e+07
min	5.000000	0.000000	0.000000e+00	0.000000e+00	0.000000	10.000000	1.500000	1960.000000	0.000000e+00
25%	10596.250000	0.207583	0.000000e+00	0.000000e+00	90.000000	17.000000	5.400000	1995.000000	0.000000e+00
50%	20969.000000	0.383585	0.000000e+00	0.000000e+00	99.000000	38.000000	6.000000	2006.000000	0.000000e+00
75%	75610.000000	0.713817	1.500000e+07	2.400000e+07	111.000000	145.750000	6.600000	2011.000000	2.085252e+07
max	417859.000000	32.985763	4.250000e+08	2.781506e+09	900.000000	9767.000000	9.200000	2015.000000	4.250000e+08

```
In [6]: df.dtypes
Out[6]:
```

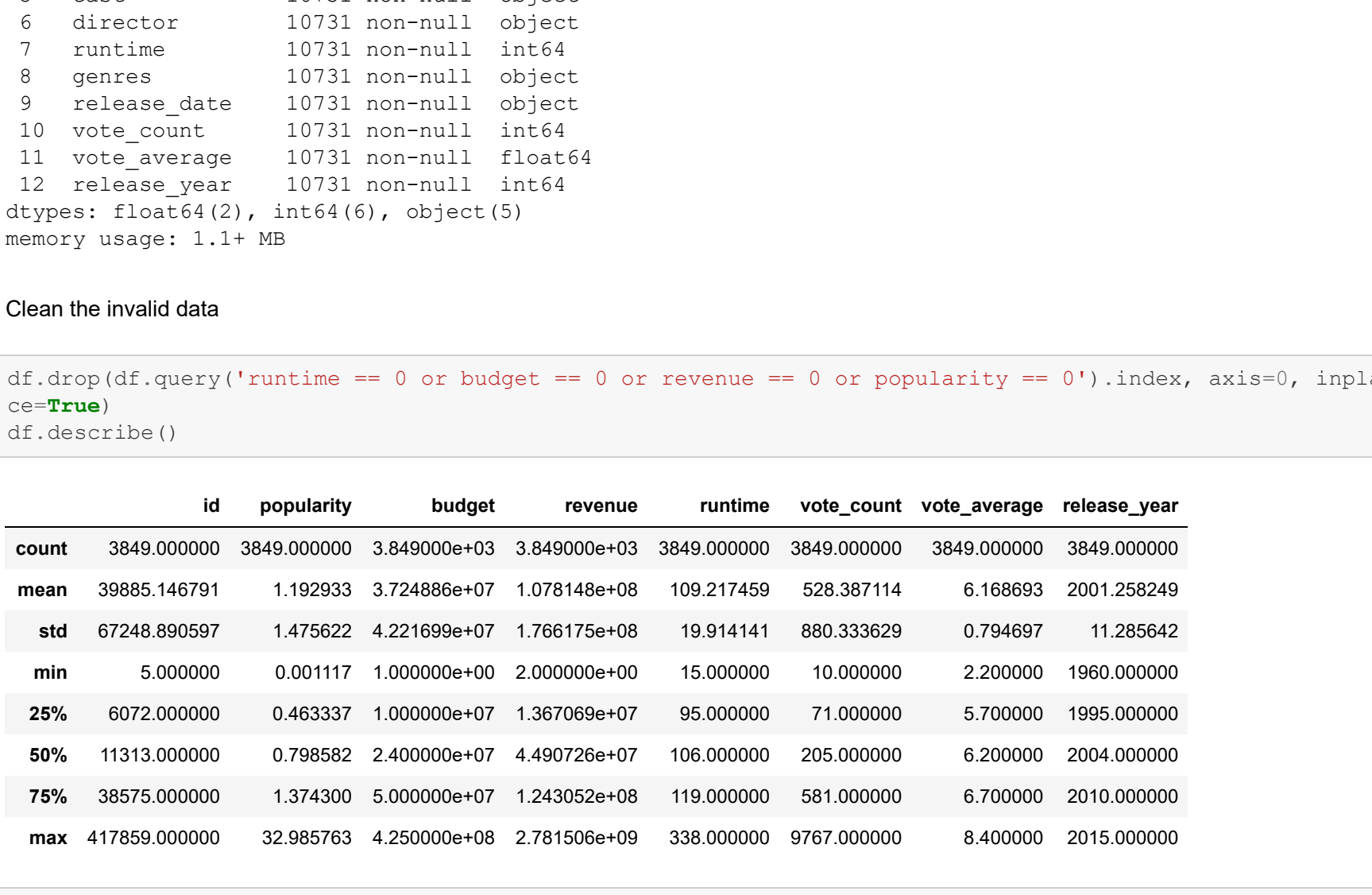
id	int64
imdb_id	object
popularity	float64
budget	int64
revenue	int64
runtime	int64
original_title	object
cast	object
homepage	object
director	object
tagline	object
keywords	object
overview	object
release_year	int64
genres	object
production_companies	object
release_date	object
vote_count	int64
vote_average	float64
release_year_adj	int64
budget_adj	float64
revenue_adj	float64
dtype:	object

```
In [7]: df.columns
Out[7]: 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 [8]: df.info()
Out[8]:
```

```
<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               10856 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               10842 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
dtype: float64(4), int64(6), object(11)
memory usage: 1.74 MB
```

```
In [9]: df.hist(figsize=(15,15))
Out[9]:
```



Data Cleaning

Drop unnecessary columns

```
In [10]: df.drop(['imdb_id', 'homepage', 'tagline', 'overview', 'budget_adj', 'revenue_adj', 'production_companies', 'keywords'], axis=1, inplace=True)
Out[11]:
```

```
df.duplicated().sum()
Out[12]: 1
```

```
In [12]: df.drop_duplicates(inplace=True)
df.duplicated().sum()
Out[12]: 0
```

```
In [13]: df.isnull().sum()
Out[13]:
```

id	0
popularity	0
budget	0
revenue	0
original_title	0
cast	76
director	44
runtime	0
genres	23
release_date	0
vote_count	0
vote_average	0
release_year	0
dtype:	int64

Drop null values

```
In [14]: df.dropna(subset=['cast', 'director', 'genres'], how='any', inplace=True)
In [15]: df.info()
```

id	0
popularity	0
budget	0
revenue	0
original_title	0
cast	76
director	44
runtime	0
genres	23
release_date	0
vote_count	0
vote_average	0
release_year	0
dtype:	float64(2), int64(6), object(5)
memory usage:	1.1+ MB

Clean the invalid data

```
In [16]: df.drop(df.query('runtime == 0 or budget == 0 or revenue == 0 or popularity == 0').index, axis=0, inplace=True)
df.describe()
Out[16]:
```

	id	popularity	budget	revenue	runtime	vote_count	vote_average	release_year
count	3849.000000	3849.000000	3.849000e+03	3.849000e+03	3849.000000	3849.000000	3849.000000	3849.000000
mean	39865.146791	1.192933	3.724886e+07	1.078148e+08	109.217459	528.387114	6.168663	2001.258249
std	67248.809597	1.475622	4.222169e+07	1.768175e+08	19.914141	880.333629	0.794697	11.286842
min	5.000000	0.001117	1.000000e+00	2.000000e+00	15.000000	10.000000	2.000000	1960.000000
25%	6072.000000	0.465337	1.000000e+00	1.367096e+07	95.000000	71.000000	5.700000	1995.000000
50%	11313.000000	0.795882	2.400000e+07	4.490702e+07	106.000000	205.000000	6.200000	2004.000000
75%	38575.000000	1.374300	5.000000e+07	1.243052e+08	119.000000	581.000000	6.700000	2010.000000
max	417859.000000	32.985763	4.250000e+08	2.781506e+09	338.000000	9767.000000	8.400000	2015.000000

```
In [17]: df.head(2)
Out[17]:
```

	id	popularity	budget	revenue	original_title	cast	director	runtime	genres	release_date	vote
0	135397	32.985763	15000000	1513528810	Jurassic World	Chris PrattBryce Dallas HowardIrfan Khan[...]	Colin Trevorrow	124	Action[Adventure Science Fiction Thriller]	6/9/15	
1	76341	28.419936	15000000	378436354	Mad Max: Fury Road	Tom HardyCharize TheronHugh Keays-Byrne[...]	George Miller	120	Action[Adventure Science Fiction Thriller]	5/13/15	

Changing format of release date into datetime format

```
In [18]: df['release_date'] = pd.to_datetime(df['release_date'])
df['release_date'].head()
Out[18]:
```

```
0    2015-06-09
1    2015-05-13
2    2015-03-18
3    2015-12-15
4    2015-04-01
Name: release_date, dtype: datetime64[ns]
```

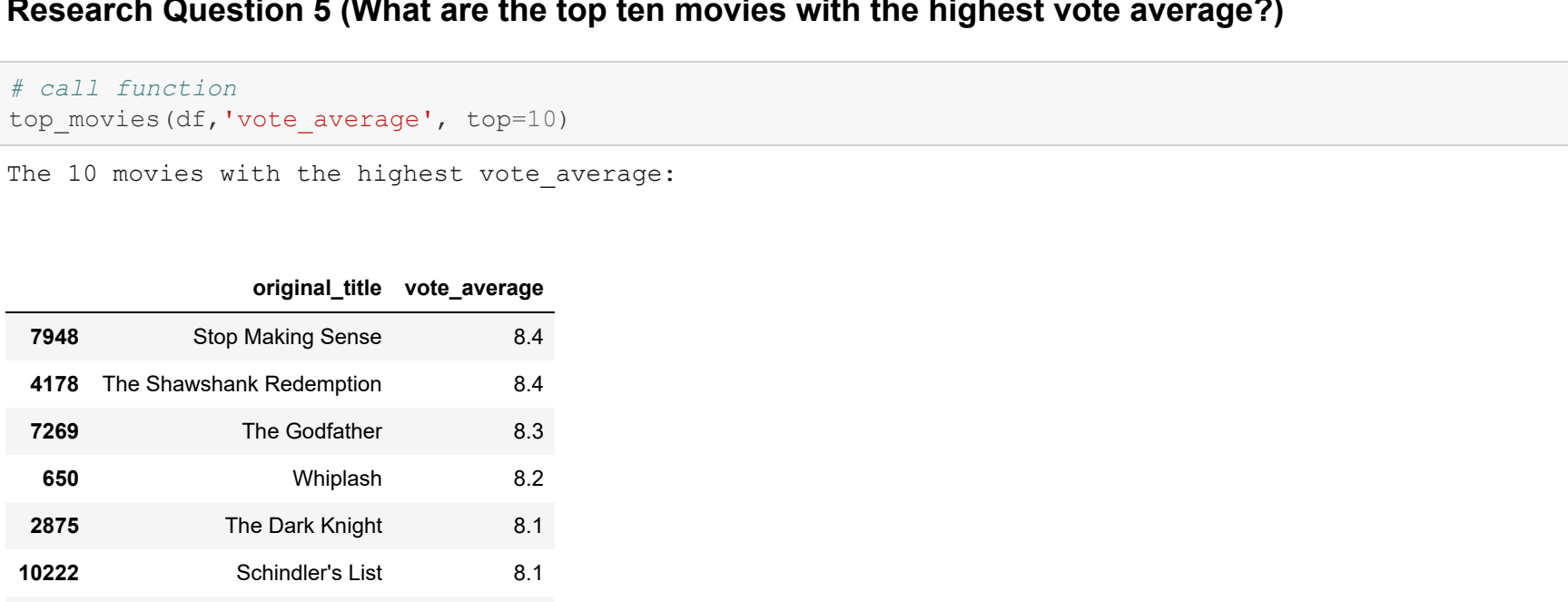
Exploratory Data Analysis

Research Question 1 (What are the 10 movies with the highest production budget?)

```
In [19]: def top_movies(df, column , top=10):
# Extract data
top = df.sort_values(column, ascending = False).loc[:, ['original_title', column]].head(10)
# Plot
ax = sns.pointplot(data = top, x = column, y = 'original_title')
# Add title and format it
ax.set_title("The 10 movies with the highest (1):\n".format(column))
# Add y label
ax.set_ylabel("\n".format(column))
# Display dataframe
print("The 10 movies with the highest (1):\n".format(column))
display(top)
# Call function
top_movies(df, 'budget', top=10)
```

The 10 movies with the highest budget:

	original_title	budget
2244	The Warrior's Way	42500000
3375	Pirates of the Caribbean: On Stranger Tides	38000000
7387	Pirates of the Caribbean: At World's End	30000000
14	Avengers: Age of Ultron	28000000
6570	Superman Returns	27700000
1929	Tangled	26000000
4411	John Carter	26000000
7394	Spider-Man 3	25800000
5408	The Lone Ranger	25500000
6503	X-Men: Days of Future Past	25000000

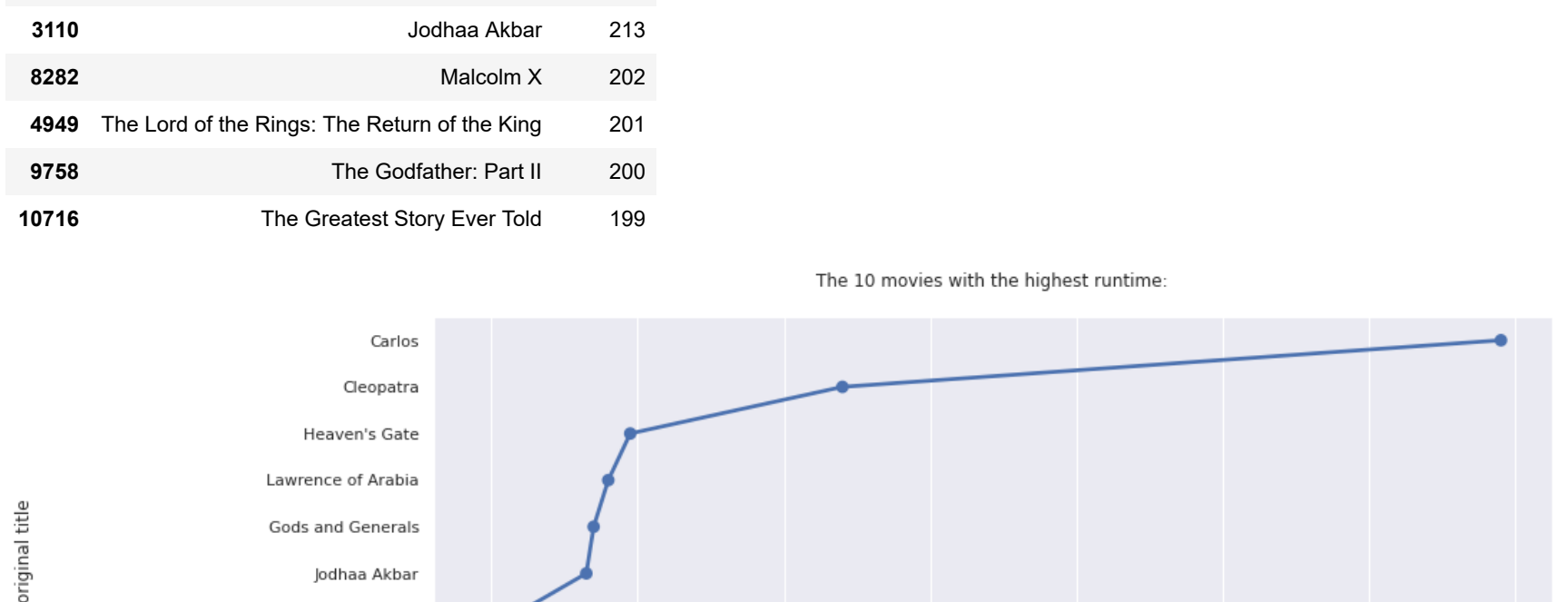


After seeing the plot and the output I can conclude that 'The Warrior's Way' is a movie that has the highest production budget.

Research Question 2 (What are the 10 movies with the highest revenue?)

```
In [20]: # Call function
top_movies(df, 'revenue', top=10)
The 10 movies with the highest revenue:
```

	original_title	revenue
1386	Avatar	2781505847
3	Star Wars: The Force Awakens	2068178225
5231	Titanic	1845034188
4361	The Avengers	1519557910
0	Jurassic World	1513528810
4	Furious 7	1508249360
14	Avengers: Age of Ultron	1405035767
3374	Harry Potter and the Deathly Hallows: Part 2	1327817822
5422	Frozen	1274219009
5425	Iron Man 3	1215439994



After seeing the plot and the output I can conclude that 'Avatar' is a movie that has the highest revenue.

Research Question 3 (What are the top 10 most profitable movies in the market?)

```
In [21]: # Add a new column (profits) into the dataframe
df['profits'] = df['revenue'] - df['budget']
# Call function
top_movies(df, 'profits', top=10)
The 10 movies with the highest profits:
```

	original_title	profits
1386	Avatar	2544050547
3	Star Wars: The Force Awakens	1868178225
5231	Titanic	1645034188
0	Jurassic World	1363528810
4	Furious 7	1316249360
4361	The Avengers	1299557910
3374	Harry Potter and the Deathly Hallows: Part 2	1208178222
14	Avengers: Age of Ultron	1125035767
5422	Frozen	1124219009
6094	The Net	1084279658



After seeing the plot and the output I can conclude that 'Avatar' is a movie that has the highest profits in the market.

Research Question 4 (What are the top ten movies with the highest popularity?)

```
In [22]: # Call function
top_movies(df, 'popularity', top=10)
The 10 movies with the highest popularity:
```

	original_title	popularity
0	Jurassic World	32.985763
1	Mad Max: Fury Road	28.419936
629	Interstellar	24.949134
630	Guardians of the Galaxy	14.311205
2	Insurgent	13.112507
631	Captain America: The Winter Soldier	12.971027
1329	Star Wars	12.037933
3	John Wick	11.422751
632	Star Wars: The Force Awakens	11.173104
633	The Hunger Games: Mockingjay - Part 1	10.739009



After seeing the plot and the output I can conclude that 'Jurassic World' is a movie that has the highest popularity.

Research Question 5 (What are the top ten movies with the highest vote average?)

```
In [23]: # Call function
top_movies(df, 'vote_average', top=10)
The 10 movies with the highest vote average:
```

	original_title	vote_average
7948	Stop Making Sense	8.4
4178	The Shawshank Redemption	8.4
7269	The Godfather	8.3
650	Whiplash	8.2
2875	The Dark Knight	8.1
10222	Schindler's List	8.1
9758	The Godfather: Part II	8.1
4177	Pulp Fiction	8.1
4179	Forrest Gump	8.1
2409	Fight Club	8.1

After seeing the plot and the output I can conclude that 'Stop Making Sense' is a movie that has the highest vote average.

Research Question 6 (What are the 10 movies with the highest runtime?)

```
In [24]: # Call function
top_movies(df, 'runtime', top=10)
The 10 movies with the highest runtime:
```

	original_title	runtime
2107	Carlos	338
10443	Cleopatra	248
7332	Heaven's Gate	219
9850	Lawrence of Arabia	216
5065	Gods and Generals	214
3110	Jodhaa Akbar	213
8283	Makoun X	202
4849	The Lord of the Rings: The Return of the King	201
9788	The Godfather: Part II	200
10716	The Greatest Story Ever Told	199

After seeing the plot and the output I can conclude that 'Carlos' is a movie that has the highest runtime.

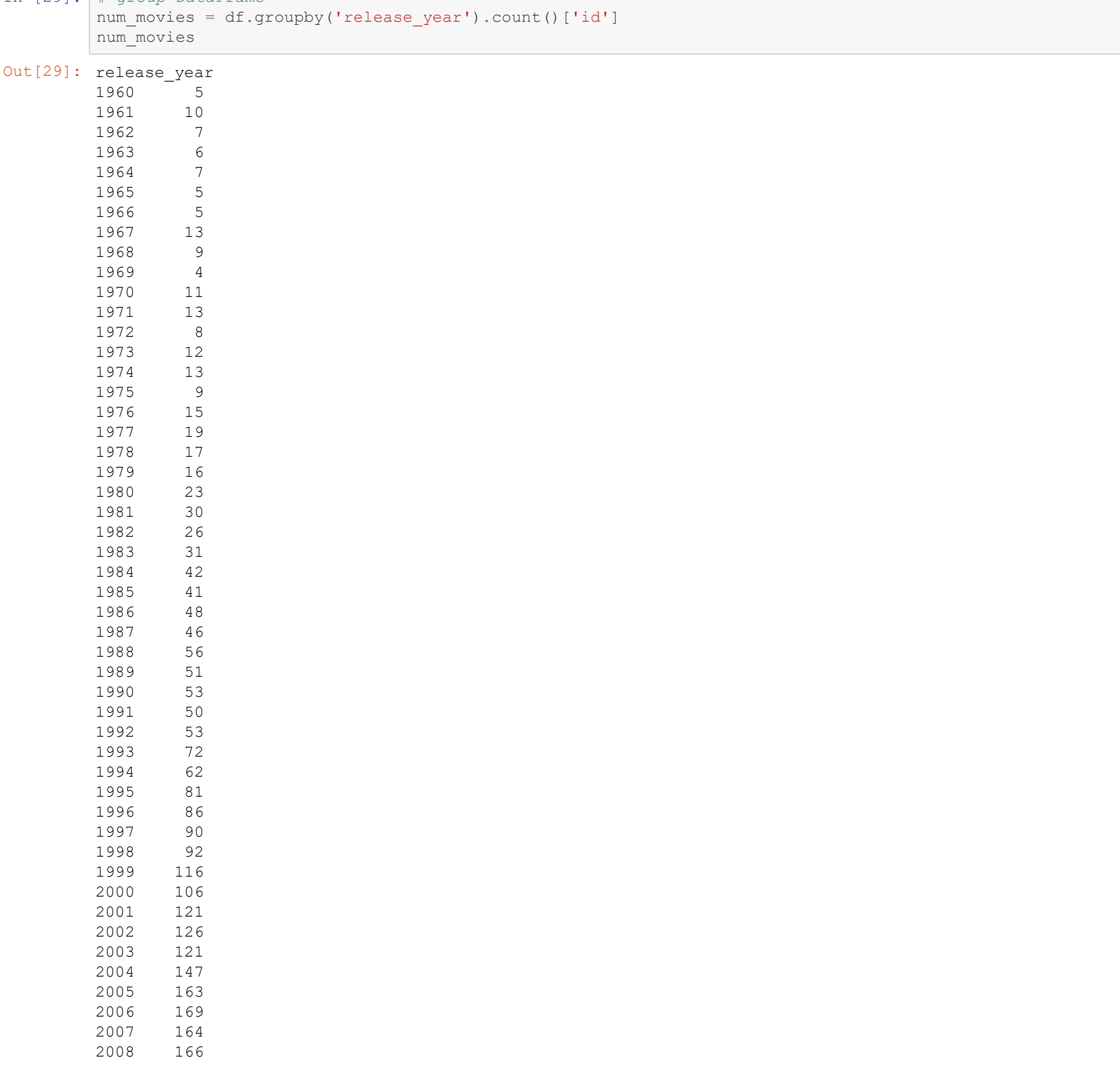
Research Question 7 (What is the average movie runtime over time?)

```
In [25]: # group DataFrame
average_runtime_over_time = df.groupby(['release_year'])['runtime'].mean()
average_release_year_over_time = df.groupby(['release_year'])['release_year'].mean()
Out[25]:
```

release_year	1960.000000
1961	132.500000
1962	141.285714
1963	153.500000
1964	122.428571
1965	167.600000
1966	132.200000
1967	109.394615
1968	130.777778
1969	127.500000
1970	114.722222
1971	112.615385
1972	113.000000
1973	111.916667
1974	109.165811
1975	119.333333
1976	118.333333
1977	106.789474
1978	122.882353
1979	115.437500
1980	117.701220
1981	105.133333
1982	112.653846
1983	107.677419
1984	108.610454
1985	104.658537
1986	105.125000
1987	109.304343
1988	104.628571
1989	106.705882
1990	108.207547
1991	


```
In [27]: # group DataFrame
average_profits_over_time = df.groupby('release_year')['profits'].mean()
average_profits_over_time

Out [27]:
release_year
1960    2.163961e+07
1961    2.908950e+07
1962    2.383998e+07
1963    1.923531e+07
1964    4.209051e+07
1965    7.985637e+07
1966    1.049374e+07
1967    4.894826e+07
1968    1.863666e+07
1969    5.171575e+07
1970    4.127909e+07
1971    1.586805e+07
1972    5.715449e+07
1973    8.814051e+07
1974    5.004477e+07
1975    1.880081e+07
1976    4.196177e+07
1977    9.175308e+07
1978    4.912359e+07
1979    5.939563e+07
1980    4.356807e+07
1981    1.408938e+07
1982    6.968294e+07
1983    4.855077e+07
1984    4.363529e+07
1985    1.063607e+07
1986    3.469489e+07
1987    4.670754e+07
1988    4.133906e+07
1989    7.066551e+07
1990    6.518689e+07
1991    6.056232e+07
1992    7.307634e+07
1993    6.355317e+07
1994    7.587717e+07
1995    6.923578e+07
1996    5.507705e+07
1997    7.089157e+07
1998    5.732767e+07
1999    5.545870e+07
2000    5.764882e+07
2001    6.629789e+07
2002    7.148870e+07
2003    7.605087e+07
2004    6.663779e+07
2005    6.606646e+07
2006    5.163491e+07
2007    7.125673e+07
2008    7.130493e+07
2009    8.085908e+07
2010    7.495069e+07
2011    7.520922e+07
2012    1.053156e+08
2013    8.768191e+07
2014    1.010679e+08
2015    1.185059e+08
Name: profits, dtype: float64
```

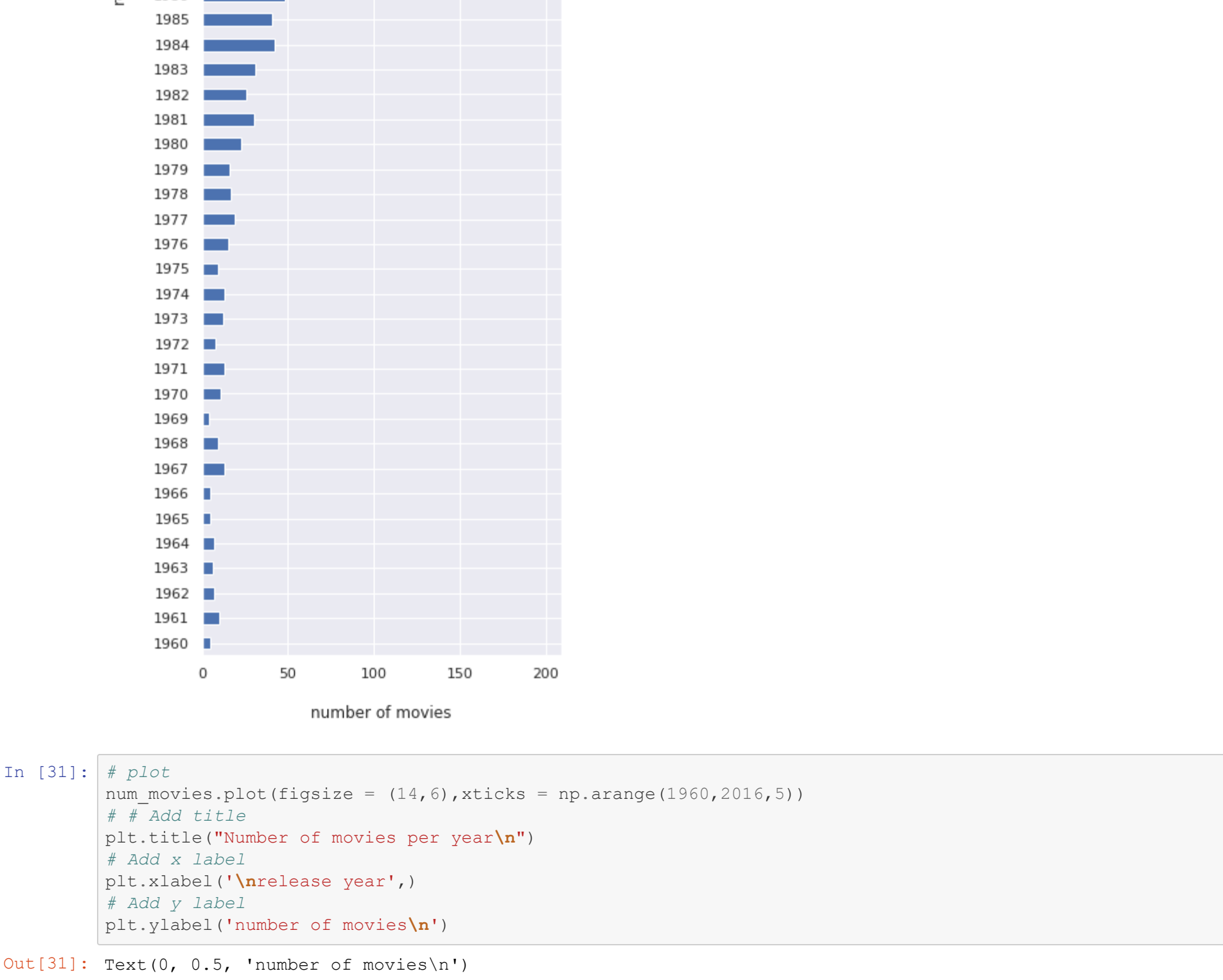


After seeing the plot and the output I can conclude that the average movie profits over time vary, but most of them increase.

Research Question 9 (How many movies per year?)

```
In [29]: # group DataFrame
num_movies = df.groupby('release_year').count()['id']
num_movies

Out [29]:
release_year
1960     5
1961    10
1962     7
1963     6
1964     7
1965     5
1966     5
1967    13
1968     9
1969     4
1970    11
1971    13
1972     8
1973    12
1974    13
1975     9
1976    15
1977    19
1978    17
1979    16
1980    15
1981    30
1982    26
1983    31
1984    42
1985    41
1986    48
1987    46
1988    56
1989    51
1990    53
1991    50
1992    53
1993    72
1994    62
1995    81
1996    86
1997    90
1998    92
1999   116
2000   106
2001   121
2002   126
2003   121
2004   147
2005   163
2006   169
2007   178
2008   199
2009   174
2010   179
2011   199
2012   157
2013   180
2014   165
2015   160
Name: id, dtype: int64
```



After seeing the plot and the output I can conclude that the number of movies is increasing every year.

Research Question 10 (In which year has the highest or lowest production of movies?)

```
In [32]: num_movies = num_movies.sort_values()
print("A year has the lowest production:\n",num_movies.head(1))
print("A year has the highest production:\n",num_movies.tail(1))

A year has the lowest production:
release_year
1969     4
Name: id, dtype: int64

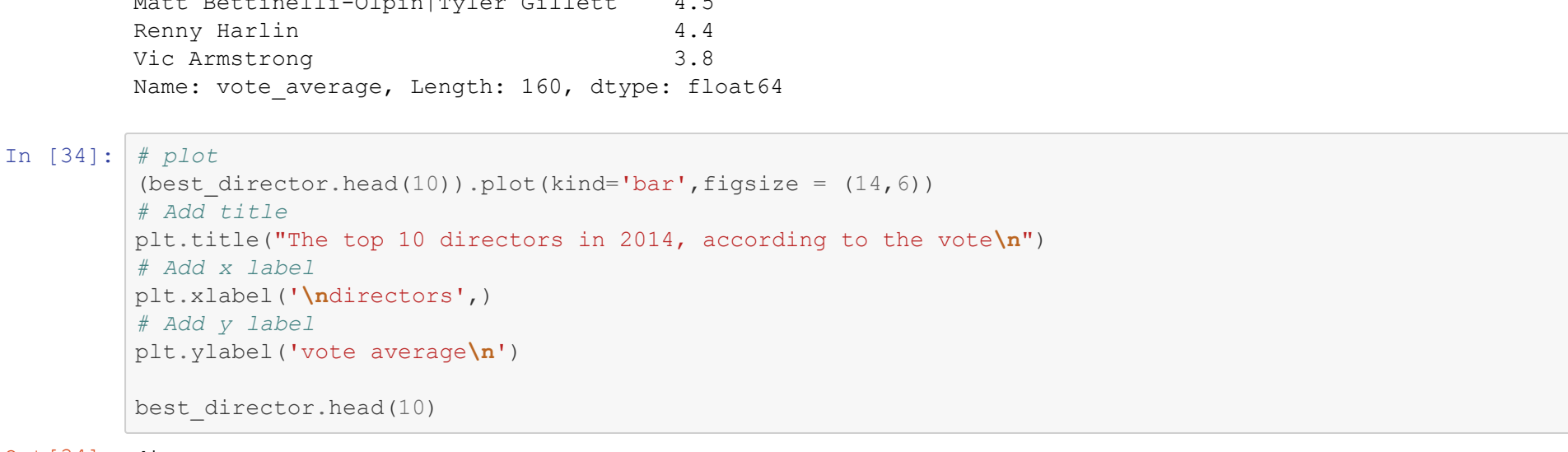
A year has the highest production:
release_year
2014   199
2011   199
Name: id, dtype: int64
```

After seeing the output I can conclude that 2014 is the year that has the highest production of movies and 1969 is the year that has the lowest production of movies.

Research Question 11 (What are the top 10 directors in 2014, according to the vote?)

```
In [33]: # group DataFrame
best_director = df.loc[df['release_year']==2014].groupby('director').max().vote_average
# sort values
best_director = best_director.sort_values(ascending=False)
best_director

Out [33]:
director
Damen Chazelle      8.2
Xavier Dolan        8.0
Christopher Nolan    8.0
Morten Tyldum        7.9
David Fincher        7.9
Grégory Levasseur    ...
Brian A Miller        4.6
Matt Bettinelli-Olpin/Tyler Gillett  4.5
Renny Harlin          4.4
Vic Armstrong         3.8
Name: vote_average, Length: 160, dtype: float64
```

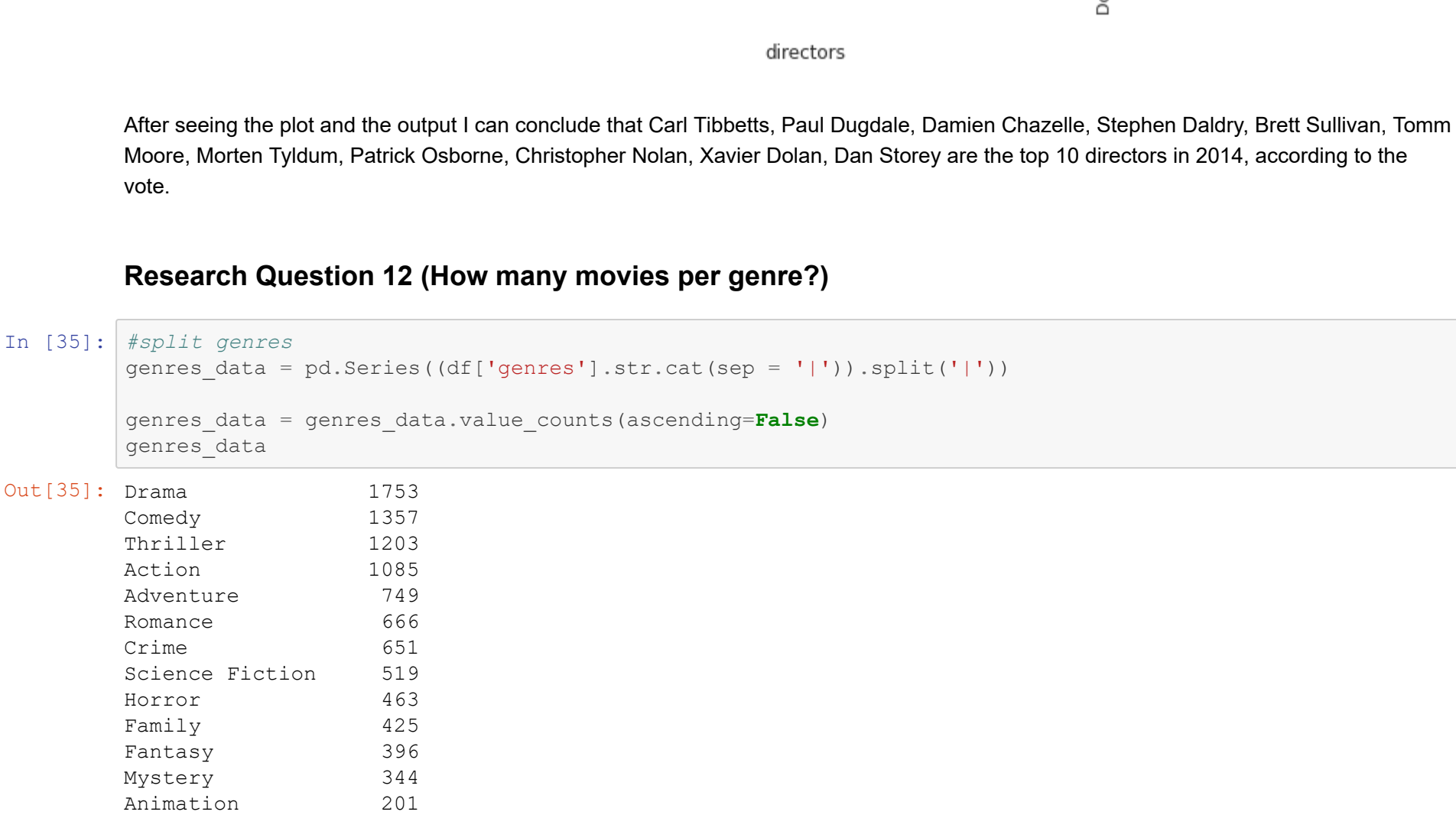


After seeing the plot and the output I can conclude that Carl Tibbetts, Paul Dugdale, Damien Chazelle, Stephen Daldry, Brett Sullivan, Tomm Moore, Morten Tyldum, Patrick Osborne, Christopher Nolan, Xavier Dolan, Dan Storey are the top 10 directors in 2014, according to the vote.

Research Question 12 (How many movies per genre?)

```
In [35]: #split genres
genres_data = pd.Series((df['genres'].str.cat(sep = '|')).split('|'))
genres_data = genres_data.value_counts(ascending=False)
genres_data

Out [35]:
genres_data
Drama      1703
Comedy     1357
Thriller   1203
Action     1085
Adventure   749
Romance    666
Crime       651
Science Fiction  519
Horror      463
Family      425
Fantasy     396
Mystery     344
History     329
War         119
Western     52
Documentary  31
Foreign     12
TV Movie    1
dtype: int64
```



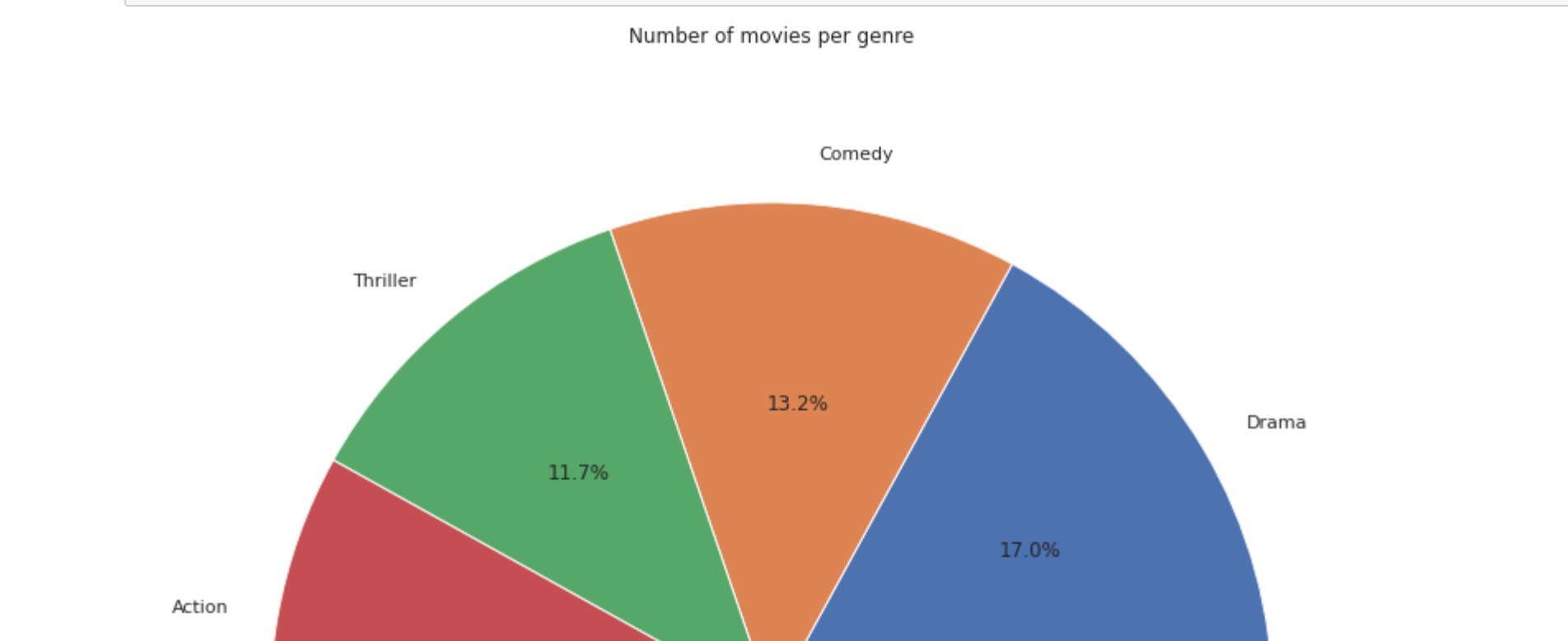
After seeing the bar plot, pie plot and output I can conclude that Drama is the most popular genre, following by Comedy, Thriller, Action, Romance, Horror, etc.

Research Question 13 (How many movies per director?)

```
In [38]: # group DataFrame
num_movies_per_director = df.groupby('director').count()['id']
# sort values
num_movies_per_director = num_movies_per_director.sort_values(ascending=False)
num_movies_per_director

Out [38]:
director
Steven Spielberg    27
Clint Eastwood      24
Ridley Scott        21
Woody Allen         18
Steven Soderbergh   17
Marc Evans          1
Marcel Langenecker  1
Marco Brambilla     1
Marco Schnabel      1
Frédéric Jardin     1
Name: id, Length: 1710, dtype: int64
```

Research Question 14 (What are the top 10 most directors who have maximum movies?)

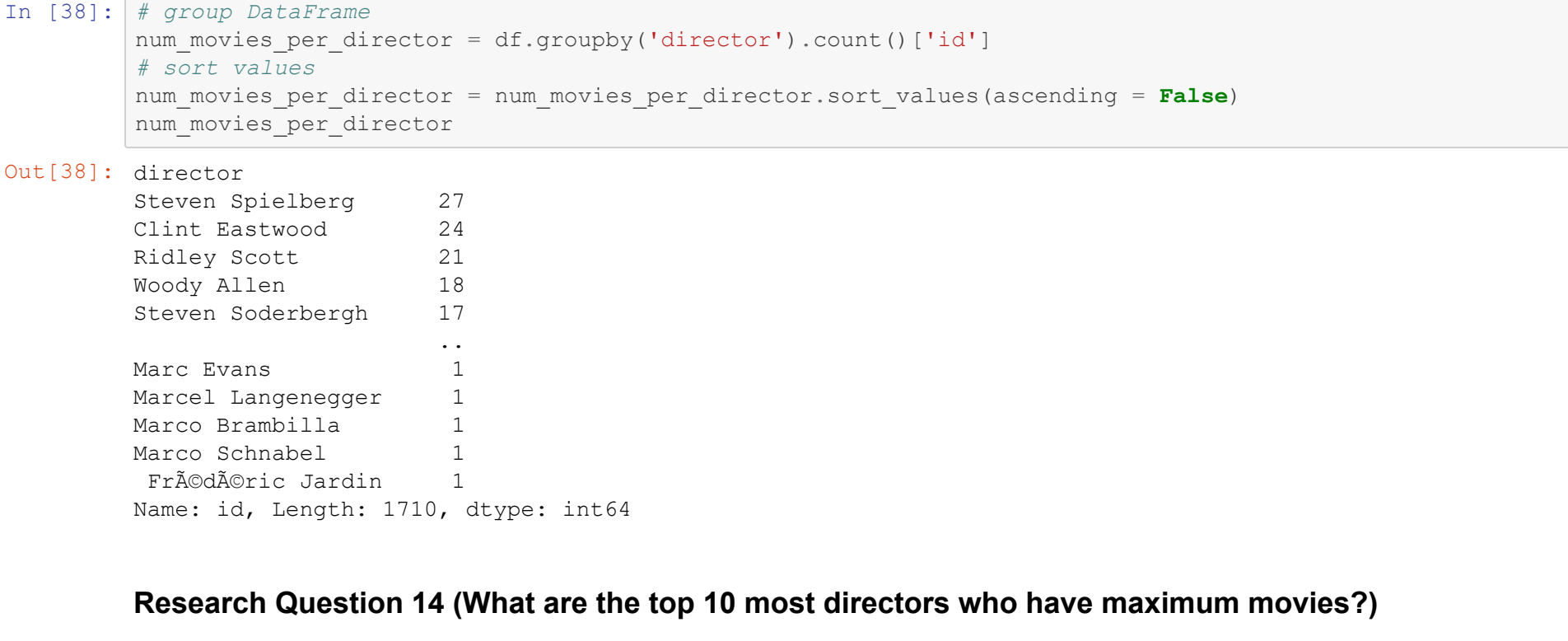


After seeing the plot and the output I can conclude that Woody Allen, Clint Eastwood, Steven Spielberg, Martin Scorsese, Ridley Scott, Ron Howard, Steven Soderbergh, Joel Schumacher, Brian De Palma, Wes Craven are the top 10 most directors who have maximum movies.

Research Question 15 (What are the most profitable months?)

```
In [40]: # add a new column (month) in the dataframe
df['month'] = df['release_date'].dt.month
# group DataFrame
profits_in_months = df.groupby('month')['profits'].mean()
profits_in_months

Out [40]:
month
3    3.051471e+07
2    4.056733e+07
3    6.559951e+07
4    5.997196e+07
5    1.187381e+08
6    1.266208e+08
7    9.404078e+07
8    4.074106e+07
9    2.941643e+07
10   4.730233e+07
11   1.016390e+08
12   9.198651e+07
Name: profits, dtype: float64
```



After seeing the plot and the output I can conclude that Month 6 is the most profitable Month, following by Month 5, Month 12, Month 11, Month 7, Month 3, Month 4, Month 8, Month 2, Month 10, Month 9, Month 1.

Conclusions

- The Warrior's Way, Pirates of the Caribbean: On Stranger Tides, Pirates of the Caribbean: At World's End, Avengers: Age of Ultron, Superman Returns, John Carter, Tangled, Spider-Man 3, The Lone Ranger, Harry Potter and the Half-Blood Prince are movies with the highest production budget.
- Avatar, Star Wars: The Force Awakens, Titanic, The Avengers, Jurassic World, Furious 7, Avengers: Age of Ultron, Harry Potter and the Deathly Hallows: Part 2, Frozen, Iron Man 3 are movies with the highest revenue.
- Avatar, Star Wars: The Force Awakens, Titanic, The Avengers, Jurassic World, Furious 7, Avengers: Age of Ultron, Harry Potter and the Deathly Hallows: Part 2, Frozen, The Net are the top 10 most profitable movies in the market.
- Jurassic World, Mad Max: Fury Road, Interstellar, Guardians of the Galaxy, Insurgent, Captain America: The Winter Soldier, Star Wars, John Wick, Star Wars: The Force Awakens, The Hunger Games: Mockingjay - Part 1 are the top ten movies with the highest popularity.
- The Story of Film: An Odyssey, Black Mirror: White Christmas, Pink Floyd: Pulse, The Art of Flight, A Personal Journey with Martin Scorsese, Dave Chappelle: Killin' Them Softly, Queen - Rock Montreal, The Shawshank Redemption, Rush: Beyond the Lighted Stage, The Jinx: The Life and Deaths of Robert Durst are the top ten movies with the highest vote average.
- The Story of Film: An Odyssey, Taken, Band of Brothers, Shoaib, Planet Earth, The Pacific, John Adams, Life, Generation Kill, The Pillars of the Earth are the 10 movies with the highest runtime.
- I've noticed that the average movie runtime over time is a decreasing.
- I've noticed that the average movie profits over time vary, but most of them increase.
- I've noticed that the number of movies is increasing every year.
- 2014 is a year that has the highest production of movies.
- 1969 is a year that has the lowest production of movies.
- Carl Tibbetts, Paul Dugdale, Damien Chazelle, Stephen Daldry, Brett Sullivan, Tomm Moore, Morten Tyldum, Patrick Osborne, Christopher Nolan, Xavier Dolan, Dan Storey are the top 10 directors in 2014, according to the vote.
- Drama is the most popular genre, following by Comedy, Thriller, Action, Romance, Horror, etc.
- Woody Allen, Clint Eastwood, Steven Spielberg, Martin Scorsese, Ridley Scott, Ron Howard, Steven Soderbergh, Joel Schumacher, Brian De Palma, Wes Craven are the top 10 most directors who have maximum movies.
- Month 6 is the most profitable Month, following by Month 5, Month 12, Month 11, Month 7, Month 3, Month 4, Month 8, Month 2, Month 10, Month 9, Month 1.

Limitations

- Missing values in the dataset affect the results.
- Duplicates values in the dataset affect the results.
- Some budget and revenue values in the dataset have zeros values.
- Incorrect data type.

Resources

- <https://pandas.pydata.org/docs/>
- <https://seaborn.pydata.org/>
- <https://github.com/usacity/data-analysis>