

- Introduction
- Data Wrangling
- Exploratory Data Analysis
- Conclusions

## overview

Tip: this report is about data analysis of movies according to genres , revenue ,budget and rating and it help us to find the best of movies that people attracted to them and how much the production team pay for that we get data from this link <https://www.kaggle.com/tmdb/tmdb-movie-metadata>

## Question(s) for Analysis

- 1-what's the most produces year of movies?
- 2-what's the most profitable according to year?
- 3-what's the most profitable type of movies?
- 4-what's the most profitable type of movies ?
- 5-what's the most popular movies genre?

```
In [141]: import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import pandas as pd
```

## General properties

1-importing data and check it 2-look at the id the columns we need 3-delete from the table the columns that we wilnt use 4-check the column value and if it will affect on analysis we will remove this column 5-check data types and covert numerical value from float to integer and convert object to string

importing data and check it look at data and choose from columns we need.

```
In [141]: df = df.read_csv('tmdb-movies.csv')
df.head()
```

	id	imdb_id	popularity	popularity.1	revenue	original_title	cast	homepage	director	tagline	...	overview	runtime	genre
Out[141]:	0	135397	ht0369610	32.985763	33.985763	1513528810	Jurassic World	Chris Pratt Bryce Dallas Howard Julianne Moore	http://www.jurassicworld.com/	Colin Trevorrow	The park is open	Twenty-two years after the events of Jurassic Park	124	Action/Adventure/Science Fiction/Thriller
	1	76341	ht1392190	28.419936	29.419936	378436354	Mad Max: Fury Road	Mad Max: Fury Road	http://www.madmaxmovie.com/	George Miller	What a Lovely Day	apocalyptic story set in the furthest reaches of our planet	120	Action/Adventure/Science Fiction/Thriller
	2	262500	ht2698446	13.112507	14.112507	295238201	Insurgent	Shailene Woodley Theo James Kate Winslet	http://www.thedivergentseries.movie/insurgent	Robert Schwendke	One Choice Can Destroy You	Prior must confront her inner demons	119	Adventure/Science Fiction/Thriller
	3	140607	ht2489496	11.173104	12.173104	2068178225	Star Wars: The Force Awakens	Harrison Ford Mark Hamill Carrie Fisher Adam Driver	http://www.starwars.com/films/star-wars-episode-7	J.J. Abrams	Every generation has a story.	Thirty years after the events of the Galactic Empire...	136	Action/Adventure/Science Fiction/Fantasy
	4	168259	ht2620852	9.335014	10.335014	1506249360	Furiosa 7	Vin Diesel Paul Walker Aaron Taylor-Johnson	http://www.furiosa7.com/	James Wan	Vengeance Has Home	Decades later, Shaw seeks revenge against Dominic Toretto	137	Action/Crime/Thriller

5 rows × 15 columns

1-clearing data by convert nan value in numerical columns to zero 2-check null value 3-drop columns which has more null value and wilnt be benefit for analysis 4-check null value again and drop it 5-check duplicated rows and drop it 6-check data types and covert numerical value from float to integer and convert object to string

```
In [148]: movies_tmdb_rev['budget_adj'], 'revenue_adj', 'runtime']
df[movies_tmdb_rev!=df[movies_tmdb_rev].replace(0,np.nan)]
df.head()
```

	id	imdb_id	popularity	popularity.1	revenue	original_title	cast	homepage	director	tagline	...	overview	runtime	genre
Out[148]:	0	135397	ht0369610	32.985763	33.985763	1513528810	Jurassic World	Chris Pratt Bryce Dallas Howard Julianne Moore	http://www.jurassicworld.com/	Colin Trevorrow	The park is open	Twenty-two years after the events of Jurassic Park	124.0	Action/Adventure/Science Fiction/Thriller
	1	76341	ht1392190	28.419936	29.419936	378436354	Mad Max: Fury Road	Mad Max: Fury Road	http://www.madmaxmovie.com/	George Miller	What a Lovely Day	An apocalyptic story set in the furthest reaches of our planet	120.0	Action/Adventure/Science Fiction/Thriller
	2	262500	ht2698446	13.112507	14.112507	295238201	Insurgent	Shailene Woodley Theo James Kate Winslet	http://www.thedivergentseries.movie/insurgent	Robert Schwendke	One Choice Can Destroy You	Beatrice Prior must confront her inner demons	119.0	Adventure/Science Fiction/Thriller
	3	140607	ht2489496	11.173104	12.173104	2068178225	Star Wars: The Force Awakens	Harrison Ford Mark Hamill Carrie Fisher Adam Driver	http://www.starwars.com/films/star-wars-episode-7	J.J. Abrams	Every generation has a story.	Thirty years after the events of the Galactic Empire...	136.0	Action/Adventure/Science Fiction/Fantasy
	4	168259	ht2620852	9.335014	10.335014	1506249360	Furiosa 7	Vin Diesel Paul Walker Aaron Taylor-Johnson	http://www.furiosa7.com/	James Wan	Vengeance Has Home	Decades later, Shaw seeks revenge against Dominic Toretto	137.0	Action/Crime/Thriller

5 rows × 15 columns

delete from the table the columns that we wilnt need

```
In [150]: df.isnull().sum()
```

id	0
imdb_id	0
popularity	0
popularity.1	0
revenue	0
original_title	0
cast	76
homepage	7939
director	44
tagline	2824
keywords	1493
overview	31
genres	233
production_companies	1030
release_date	0
vote_count	0
release_year	0
budget_adj	5696
revenue_adj	6616
dtype:	int64

```
In [162]: df.movies.drop(['imdb_id','cast','homepage','tagline','overview','release_date','director','keywords','production_companies'],axis=1)
movies_tmdb.head()
```

	id	popularity	popularity.1	revenue	original_title	runtime	genres	vote_count	vote_average	release_year	budget_adj	revenue_adj
Out[162]:	0	135397	32.985763	33.985763	1513528810	124.0	Action/Adventure/Science Fiction/Thriller	5562	6.5	2015	137999939.3	1392445893
	1	76341	28.419936	29.419936	378436354	120.0	Action/Adventure/Science Fiction/Thriller	6185	7.1	2015	137999939.3	3481613e+08
	2	262500	13.112507	14.112507	295238201	119.0	Adventure/Science Fiction/Thriller	2480	6.3	2015	101199955.5	2.716130e+08
	3	140607	11.173104	12.173104	2068178225	136.0	Action/Adventure/Science Fiction/Fantasy	5292	7.5	2015	183999919.0	1.902723e+09
	4	168259	9.335014	10.335014	1506249360	137.0	Action/Crime/Thriller	2947	7.3	2015	174799923.1	1.385749e+09

```
In [163]: df.movies.isnull().sum()
```

id	0
popularity	0
popularity.1	0
revenue	0
original_title	0
runtime	31
genres	23
vote_count	0
vote_average	0
release_year	0
budget_adj	5696
revenue_adj	6616
dtype:	int64

```
In [164]: movies_tmdb.drop(movies_tmdb.isnull().sum())
movies_tmdb.head()
```

	id	popularity	popularity.1	revenue	original_title	runtime	genres	vote_count	vote_average	release_year	budget_adj	revenue_adj
Out[164]:	0	135397	32.985763	33.985763	1513528810	124.0	Action/Adventure/Science Fiction/Thriller	5562	6.5	2015	137999939.3	1.392445e+09
	1	76341	28.419936	29.419936	378436354	120.0	Action/Adventure/Science Fiction/Thriller	6185	7.1	2015	137999939.3	3.481613e+08
	2	262500	13.112507	14.112507	295238201	119.0	Adventure/Science Fiction/Thriller	2480	6.3	2015	101199955.5	2.716130e+08
	3	140607	11.173104	12.173104	2068178225	136.0	Action/Adventure/Science Fiction/Fantasy	5292	7.5	2015	183999919.0	1.902723e+09
	4	168259	9.335014	10.335014	1506249360	137.0	Action/Crime/Thriller	2947	7.3	2015	174799923.1	1.385749e+09

```
In [165]: print(movies_tmdb.shape)
(3855, 12)
```

```
In [166]: movies_tmdb.isnull().sum()
```

id	0
popularity	0
popularity.1	0
revenue	0
original_title	0
runtime	0
genres	0
vote_count	0
vote_average	0
release_year	0
budget_adj	0
revenue_adj	0
dtype:	int64

```
In [167]: movies_tmdb.duplicated().sum()
```

1
---

```
In [168]: movies_tmdb=movies_tmdb.drop_duplicates()
movies_tmdb.head()
```

	id	popularity	popularity.1	revenue	original_title	runtime	genres	vote_count	vote_average	release_year	budget_adj	revenue_adj
Out[168]:	0	135397	32.985763	33.985763	1513528810	124.0	Action/Adventure/Science Fiction/Thriller	5562	6.5	2015	137999939.3	1.392445e+09
	1	76341	28.419936	29.419936	378436354	120.0	Action/Adventure/Science Fiction/Thriller	6185	7.1	2015	137999939.3	3.481613e+08
	2	262500	13.112507	14.112507	295238201	119.0	Adventure/Science Fiction/Thriller	2480	6.3	2015	101199955.5	2.716130e+08
	3	140607	11.173104	12.173104	2068178225	136.0	Action/Adventure/Science Fiction/Fantasy	5292	7.5	2015	183999919.0	1.902723e+09
	4	168259	9.335014	10.335014	1506249360	137.0	Action/Crime/Thriller	2947	7.3	2015	174799923.1	1.385749e+09

```
In [169]: movies_tmdb.dtypes
```

id	int64
popularity	float64
popularity.1	float64
revenue	int64
original_title	object
runtime	float64
genres	object
vote_count	int64
vote_average	float64
release_year	int64
budget_adj	float64
revenue_adj	float64
dtype:	object

```
In [170]: integers=['revenue_adj','budget_adj']
movies_tmdb[integer]=movies_tmdb[integer].applymap(np.int64)
movies_tmdb.head()
```

	id	popularity	popularity.1	revenue	original_title	runtime	genres	vote_count	vote_average	release_year	budget_adj	revenue_adj
Out[170]:	0	135397	32.985763	33.985763	1513528810	124.0	Action/Adventure/Science Fiction/Thriller	5562	6.5	2015	137999939.3	1.392445e+09
	1	76341	28.419936	29.419936	378436354	120.0	Action/Adventure/Science Fiction/Thriller	6185	7.1	2015	137999939.3	3.481613e+08
	2	262500	13.112507	14.112507	295238201	119.0	Adventure/Science Fiction/Thriller	2480	6.3	2015	101199955.5	2.716130e+08
	3	140607	11.173104	12.173104	2068178225	136.0	Action/Adventure/Science Fiction/Fantasy	5292	7.5	2015	183999919.0	1.902723e+09
	4	168259	9.335014	10.335014	1506249360	137.0	Action/Crime/Thriller	2947	7.3	2015	174799923.1	1.385749e+09

```
In [171]: movies_tmdb['profit'] = movies_tmdb['revenue']-movies_tmdb['budget_adj']
movies_tmdb['profit'] = movies_tmdb['profit'].apply(np.int64)
movies_tmdb.head()
```

	id	popularity	popularity.1	revenue	original_title	runtime	genres	vote_count	vote_average	release_year	budget_adj	revenue_adj	profit
Out[171]:	0	135397	32.985763	33.985763	1513528810	124.0	Action/Adventure/Science Fiction/Thriller	5562	6.5	2015	137999939.3	1392445893	1379528871
	1	76341	28.419936	29.419936	378436354	120.0	Action/Adventure/Science Fiction/Thriller	6185	7.1	2015	137999939.3	348161292	240436415
	2	262500	13.112507	14.112507	295238201	119.0	Adventure/Science Fiction/Thriller	2480	6.3	2015	101199955.5	271619025	194038246
	3	140607	11.173104	12.173104	2068178225	136.0	Action/Adventure/Science Fiction/Fantasy	5292	7.5	2015	183999919.0	190272330	1884178306
	4	168259	9.335014	10.335014	1506249360	137.0	Action/Crime/Thriller	2947	7.3	2015	174799923.1	1385748801	1331449437

```
In [172]: movies_tmdb['genres']=movies_tmdb['genres'].apply(np.str_)
movies_tmdb.dtypes
```

id	int64
popularity	float64
popularity.1	float64
revenue	int64
original_title	object
runtime	float64
genres	object
vote_count	int64
vote_average	float64
release_year	int64
budget_adj	int64
revenue_adj	int64
profit	object
dtype:	object

```
In [173]: print(movies_tmdb.shape)
(3854, 12)
```

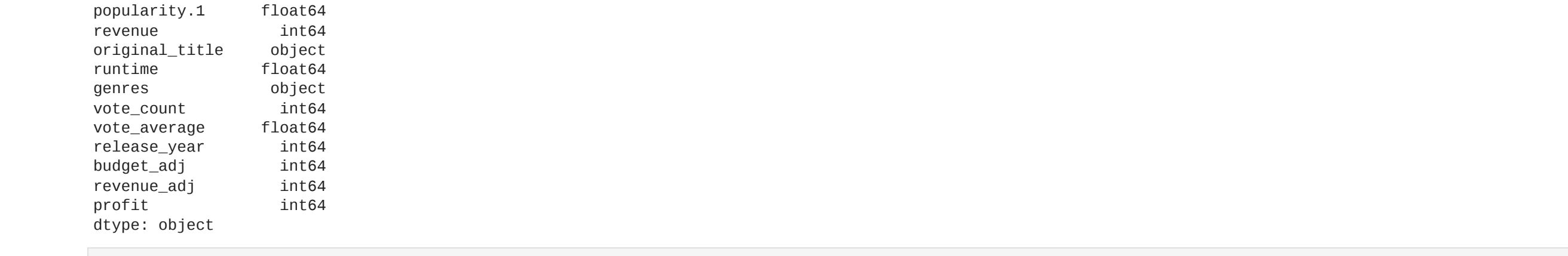
```
In [175]: movies_tmdb.duplicated().sum()
```

0
---

```
In [176]: movies_tmdb.hist(figsize=(18,8))
```

what's the most year of production of movies?

what's the most produces year of movies ? it appears for that number of movies increases gradually by year



tmdb\_genres=tmdb

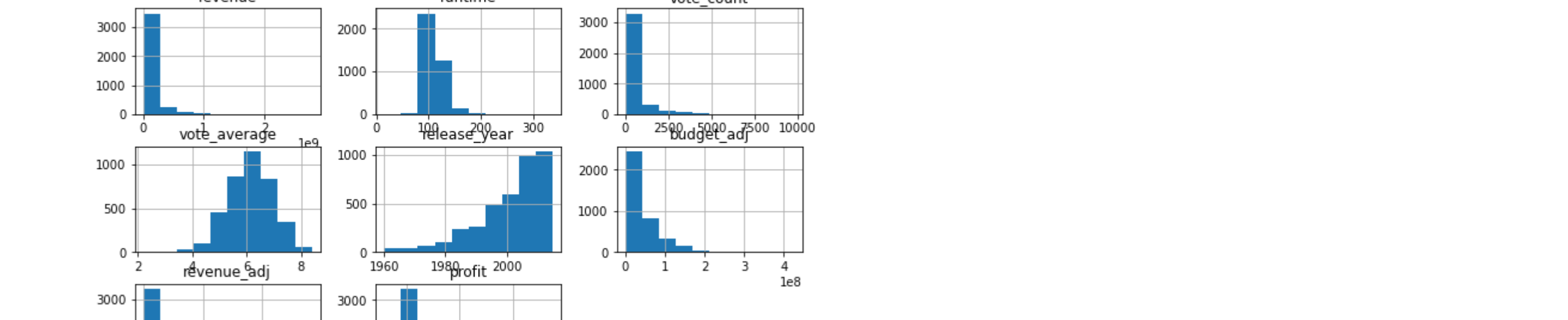
what's the most profitable according to year? we found the profit of movies increases gradually by year



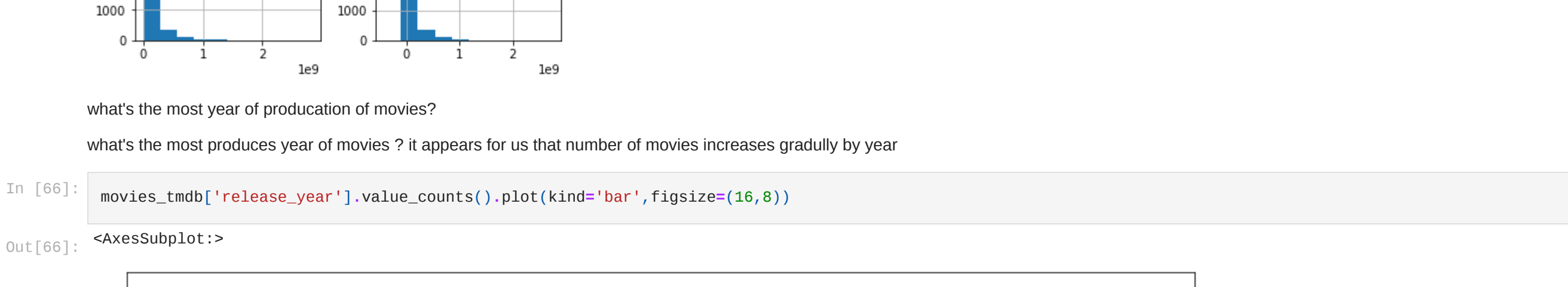
we can estimate the profit mean as below

```
In [188]: movies_tmdb['profit'].describe()
```

count	3.854906e+03
mean	6.344662e+07
std	1.508996e+08
min	-4.129124e+07
25%	-5.129906e+06
50%	1.443327e+07
75%	7.451931e+07
max	2.540619e+09
Name:	profit, dtype: float64



we can check the relation between profit and budget for this plot



we use split to split column of genres to make it easy for analysing data and add the column after splitting to table to make more analysis

```
In [198]: genres_df = movies_tmdb['genres'].str.split(" ", expand=True)
genres_movies = genres_df.stack()
genres_movies = pd.DataFrame(genres_df)
genres_movies.head()
```

	0	1	2	3	4
Out[198]:	Action	Adventure	Science Fiction	Thriller	None
	1	Action	Adventure	Science Fiction	Thriller
	2	Adventure	Science Fiction	Thriller	None
	3	Action	Adventure	Science Fiction	Fantasy
	4	Action	Crime	Thriller	None

```
In [200]: genres_movies_count = genres_df.stack()
genres_movies_count = pd.DataFrame(genres_movies_count)
genres_movies_count.head()
```

	0
Out[200]:	0
	1
	2
	3
	4

what's the most productive type of movies? from the table we find drama movies is the most productive type of movies

Out[199]:	movies_title.head()												
	id	popularity	popularity1	revenue	original_title	runtime	genres	vote_count	vote_average	release_year	budget_adj	revenue_adj	profit
	8	135997	32.99573	33.99573	1513628910		Jurassic World	124.0					
	1	76341	28.41936	29.41936	379436354		Mad Max: Fury Road	120.0					
	2	263500	13.112507	14.112507	295238201		Infernal Desire	119.0					
	3	140607	11.173104	12.173104	2068178225		Star Wars: The Force Awakens	136.0					
	4	166259	9.133014	10.133014	1500249400		Furiosa	137.0					
	5	146299	7.139233	8.139233	1500249400		Action/Adventure/Science Fiction/Thriller	5562	6.5	2015	173999939	1392446450	1379520871
	6	263500	13.112507	14.112507	295238201		Adventure/Science Fiction/Thriller	341816	7.1	2015	137999939	1392446450	260436345
	7	140607	11.173104	12.173104	2068178225		Action/Adventure/Science Fiction/Thriller	2480	6.3	2015	101199955	271619052	194038246
	8	140607	11.173104	12.173104	2068178225		Action/Adventure/Science Fiction/Thriller	5292	7.5	2015	183899919	1390723130	1884178906
	9	146299	7.139233	8.139233	1500249400		Action/Crime/Thriller	2947	7.3	2014	17499923	180077380	133144934