

IMDB Movie Analysis

```
In [1]: #Importing required packages
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import plotly.express as px

In [2]: #Importing csv to the dataframe df
df = pd.read_csv('Downloads/IMDB_Movies.csv')

In [3]: #checking the dataset
df

Out[3]:
```

	color	director_name	num_critic_for_reviews	gross	duration	director_facebook_likes	actor_3_facebook_likes	actor_2_name	actor_1_facebook_likes	gross
0	Color	James Cameron	723.0	178.0	0.0	855.0	Joni Dawe Moore	1000.0	76505647.0	
1	Color	Gore Verbinski	302.0	168.0	563.0	1000.0	Orlando Bloom	40000.0	309404152.0	
2	Color	Sam Mendes	602.0	148.0	0.0	161.0	Rory Kinnear	11000.0	309404152.0	
3	Color	Christopher Nolan	813.0	164.0	22000.0	23800.0	Christian Bale	27000.0	448110642.0	
4	NaN	Doug Walker	NaN	NaN	131.0	NaN	Rob Walker	131.0	NaN	
...	
5038	Color	Scott Smith	1.0	87.0	2.0	318.0	Dagmar Durnig	637.0	NaN	
5039	Color	NaN	43.0	43.0	NaN	319.0	Viviane Curry	841.0	NaN	
5040	Color	Benjamin Roberts	13.0	76.0	0.0	0.0	Naama Moody	0.0	NaN	
5041	Color	Daniel Hsia	14.0	100.0	0.0	489.0	Daniel Henney	940.0	10443.0	
5042	Color	Jon Gunn	43.0	90.0	16.0	16.0	Brian Huppert	86.0	8522.0	

5043 rows x 10 columns

A.Cleaning the data. This is one of the most important step to perform before moving forward with the analysis. Use your knowledge learned till now to do this. (Droping columns, removing null values, etc.

Hint: Clean the data

```
In [4]: # Checking the columns for analysis
of columns

Out[4]:
```

```
Index(['color', 'director_name', 'num_critic_for_reviews', 'duration', 'director_facebook_likes', 'actor_3_facebook_likes', 'actor_2_name', 'actor_1_facebook_likes', 'gross', 'genres', 'actor_1_name', 'movie_title', 'num_voted_users', 'num_user_for_reviews', 'language', 'country', 'content_rating', 'budget', 'title_year', 'actor_2_facebook_likes', 'imdb_score', 'aspect_ratio', 'movie_facebook_likes'],
      dtype='object')
```

```
In [5]: # dropping unnecessary columns
df = df.drop(['color', 'director_facebook_likes', 'actor_3_facebook_likes', 'actor_2_name', 'actor_1_facebook_likes', 'actor_2_facebook_likes', 'actor_3_facebook_likes', 'actor_2_name', 'cast_total_facebook_likes', 'actor_3_name', 'duration', 'faceumber_in_poster', 'content_ratio', 'country', 'movie_imdb_link', 'aspect_ratio', 'plot_keywords'], axis=1)
```

```
In [6]: #Finding the number on null values in decreasing order
df.isnull().sum(axis=0).sort_values(ascending=False)
```

```
Out[6]:
```

	gross	budget	title_year	director_name	num_critic_for_reviews	num_user_for_reviews	num_voted_users	movie_title	imdb_score	movie_facebook_likes	dtype
gross	884	492	188	184	158	29	12	7	0	0	int64
budget	492	188	184	158	29	12	7	0	0	0	int64
title_year	188	184	158	29	12	7	0	0	0	0	int64
director_name	184	158	29	12	7	0	0	0	0	0	int64
num_critic_for_reviews	158	29	12	7	0	0	0	0	0	0	int64
num_user_for_reviews	29	12	7	0	0	0	0	0	0	0	int64
num_voted_users	12	7	0	0	0	0	0	0	0	0	int64
movie_title	7	0	0	0	0	0	0	0	0	0	int64
imdb_score	0	0	0	0	0	0	0	0	0	0	int64
movie_facebook_likes	0	0	0	0	0	0	0	0	0	0	int64

```
In [7]: #since gross and budget have the most null values I have passed the data frame with not
# null constraint
df = df[df['gross'] > 0]
df = df[df['budget'] > 0]

In [8]: # Checking after dropping null values in gross and budget
df.isnull().sum(axis=0).sort_values(ascending=False)
```

```
Out[8]:
```

	actor_1_name	language	num_critic_for_reviews	director_name	gross	genres	actor_1_name	movie_title	num_voted_users	num_user_for_reviews	dtype
actor_1_name	3	3	1	0	0	0	0	0	0	0	int64
language	3	3	1	0	0	0	0	0	0	0	int64
num_critic_for_reviews	1	0	0	0	0	0	0	0	0	0	int64
director_name	0	0	0	0	0	0	0	0	0	0	int64
gross	0	0	0	0	0	0	0	0	0	0	int64
genres	0	0	0	0	0	0	0	0	0	0	int64
movie_title	0	0	0	0	0	0	0	0	0	0	int64
num_voted_users	0	0	0	0	0	0	0	0	0	0	int64
num_user_for_reviews	0	0	0	0	0	0	0	0	0	0	int64
budget	0	0	0	0	0	0	0	0	0	0	int64
title_year	0	0	0	0	0	0	0	0	0	0	int64
imdb_score	0	0	0	0	0	0	0	0	0	0	int64
movie_facebook_likes	0	0	0	0	0	0	0	0	0	0	int64

```
In [9]: #dropping all duplicate records
df.drop_duplicates()
```

```
Out[9]:
```

	director_name	num_critic_for_reviews	gross	genres	actor_1_name	movie_title	num_voted_users	num_user_for_reviews	language
0	James Cameron	723.0	76505647.0	Action/Adventure/Fantasy/Sci-Fi	CCI Paounder	Avatar	885204	3054	Englis
1	Gore Verbinski	302.0	309404152.0	Action/Adventure/Fantasy	Johnny Depp	Pirates of the Caribbean: At World's End	471220	1238	Englis
2	Sam Mendes	602.0	200074175.0	Action/Adventure/Thriller	Christoph Waltz	Spectre	275869	994	Englis
3	Christopher Nolan	813.0	448110642.0	Action/Thriller	Tom Hardy	The Dark Knight Rises	1144337	2701	Englis
5	Andrew Stanton	462.0	73058679.0	Action/Adventure/Sci-Fi	Daryl Sabara	John Carter	212204	738	Englis
...
5038	Shane Carruth	143.0	424760.0	Drama/Sci-Fi/Thriller	Shane Carruth	Primer	72039	37	
5034	Nail Dala Liana	35.0	70701.0	Thriller	Ian Gamazon	Cavite	589	9	
5035	Robert Rodriguez	56.0	2040020.0	Action/Crime/Drama/Romance/Comedy	Cedric Belfrage	El Mariachi	52055	13	
5037	Edwin Burns	14.0	4584.0	Comedy/Drama	Kenny Bock	Nanowatts	1338	1	
5042	Jon Gunn	43.0	86322.0	Documentary	John August	My Date with the Doctor	4285	8	

3556 rows x 13 columns

B.Movies with highest profit. Create a new column called profit which contains the difference of the two columns: gross and budget. Sort the column using the profit column as reference. Plot profit (y-axis) vs budget (x-axis) and observe the outliers using the appropriate chart type

Hint: Find the movies with the highest profit?

```
In [10]: #creating a new column profit which is difference between gross and budget
df['profit'] = df['gross'] - df['budget']
df.head()
```

```
Out[10]:
```

	director_name	num_critic_for_reviews	gross	genres	actor_1_name	movie_title	num_voted_users	num_user_for_reviews	language
0	James Cameron	723.0	76505647.0	Action/Adventure/Fantasy/Sci-Fi	CCI Paounder	Avatar	885204	3054	Englis
1	Gore Verbinski	302.0	309404152.0	Action/Adventure/Fantasy	Johnny Depp	Pirates of the Caribbean: At World's End	471220	1238	Englis
2	Sam Mendes	602.0	200074175.0	Action/Adventure/Thriller	Christoph Waltz	Spectre	275869	994	Englis
3	Christopher Nolan	813.0	448110642.0	Action/Thriller	Tom Hardy	The Dark Knight Rises	1144337	2701	Englis
5	Andrew Stanton	462.0	73058679.0	Action/Adventure/Sci-Fi	Daryl Sabara	John Carter	212204	738	Englis

```
In [11]: #Sorting on basis of profit (Decreasing Order)
df.sort_values(by=['profit'],ascending=False)
```

```
In [12]: #Plotting Profit vs Budget of the films having profit greater than 200M
df = df[df['profit']>200000000]
fig = px.scatter(df, x='budget', y='profit', title='Profit vs Budget', markers=True)
fig.show()
```

```
In [13]: top100 = df[['movie_title', 'budget', 'profit']]
top100.head(10)
```

```
Out[13]:
```

	movie_title	budget	profit
8	Avatar	237000000.0	52055647.0
29	Jurassic World	150000000.0	920177271.0
26	Warrior	200000000.0	459873932.0
3034	Star Wars: Episode IV - A New Hope	11000000.0	49995669.0
3080	E.T. the Extra-Terrestrial	10500000.0	42449459.0
794	The Avengers	220000000.0	40237947.0
67	The Avengers	220000000.0	40237947.0
590	The Lone King	45000000.0	37718377.0
240	Star Wars: Episode I - The Phantom Menace	115000000.0	35954677.0
66	The Dark Knight	185000000.0	348316061.0

C.Top 250. Create a new column IMDB_Top_250 and store the top 250 movies with the highest IMDB Rating (corresponding to the column: imdb_score). Also make sure that for all of these movies, the num_voted_users is greater than 25,000. Also add a Rank column containing the values 1 to 250 indicating the ranks of the corresponding films.

Extract all the movies in the IMDB_Top_250 column which are not in the English language and store them in a new column named Top_Foreign_Lang_Film. You can use your own imagination about

Hint: Find IMDB_Top_250

```
In [14]: # Creating a dataframe IMDB_Top_250 having voted users greater than 25000
IMDB_Top_250 = df[df['num_voted_users'] > 25000].sort_values(by='imdb_score',ascending=False).head(250)
IMDB_Top_250
```

```
Out[14]:
```

	director_name	num_critic_for_reviews	gross	genres	actor_1_name	movie_title	num_voted_users	num_user_for_reviews
1327	Frank Darabont	199.0	20341489.0	Crime/Drama	Morgan Freeman	The Shawshank Redemption	1689764	414
1466	Francis Ford Coppola	208.0	134021952.0	Crime/Drama	Al Pacino	The Godfather	1155770	221
66	Christopher Nolan	645.0	53331093.0	Action/Crime/Drama/Thriller	Christian Bale	The Dark Knight	1670169	491
2837	Francis Ford Coppola	149.0	57300000.0	Crime/Drama	Robert De Niro	The Godfather: Part II	790926	61
1874	Steven Spielberg	174.0	96067179.0	Biography/Drama/History	Liam Neeson	Schindler's List	865020	127
...
1171	Yimou Zhang	283.0	84961.0	Action/Adventure/History	Jie Li	Hero	148414	84
1748	F Gary Gray	349.0	16102670.0	Biography/Crime/Drama/History/Music	Aids Hodge	Straight Outta Compton	119928	31
788	Cameron Crowe	149.0	35252352.0	Adventure/Comedy/Drama/Music	Philip Seymour Hoffman	Almost Famous	207287	81
639	Michael Mann	209.0	29665197.0	Biography/Drama/Thriller	Al Pacino	The Insider	133526	51
525	Alfonso Cuarón	372.0	35286428.0	Drama/Sci-Fi/Thriller	Charlie Hunnam	Children of Men	361787	12

250 rows x 14 columns

```
In [15]: # Adding a new column Rank in IMDB_Top_250 as per IMDB Score
IMDB_Top_250['Rank'] = IMDB_Top_250['IMDB_Score'].rank(method='first',ascending=False)
IMDB_Top_250
```

```
Out[15]:
```

	director_name	num_critic_for_reviews	gross	genres	actor_1_name	movie_title	num_voted_users	num_user_for_reviews
1327	Frank Darabont	199.0	20341489.0	Crime/Drama	Morgan Freeman	The Shawshank Redemption	1689764	414
1466	Francis Ford Coppola	208.0	134021952.0	Crime/Drama	Al Pacino	The Godfather	1155770	221
66	Christopher Nolan	645.0	53331093.0	Action/Crime/Drama/Thriller	Christian Bale	The Dark Knight	1670169	491
2837	Francis Ford Coppola	149.0	57300000.0	Crime/Drama	Robert De Niro	The Godfather: Part II	790926	61
1874	Steven Spielberg	174.0	96067179.0	Biography/Drama/History	Liam Neeson	Schindler's List	865020	127
...
1171	Yimou Zhang	283.0	84961.0	Action/Adventure/History	Jie Li	Hero	148414	84
1748	F Gary Gray	349.0	16102670.0	Biography/Crime/Drama/History/Music	Aids Hodge	Straight Outta Compton	119928	31
788	Cameron Crowe	149.0	35252352.0	Adventure/Comedy/Drama/Music	Philip Seymour Hoffman	Almost Famous	207287	81
639	Michael Mann	209.0	29665197.0	Biography/Drama/Thriller	Al Pacino	The Insider	133526	51
525	Alfonso Cuarón	372.0	35286428.0	Drama/Sci-Fi/Thriller	Charlie Hunnam	Children of Men	361787	12

250 rows x 15 columns

```
In [16]: # Creating a dataframe named Top_Foreign_Lang_Film containing all non english films in top 250
Top_Foreign_Lang_Film = IMDB_Top_250[IMDB_Top_250['language']!='English']
Top_Foreign_Lang_Film
```

```
Out[16]:
```

	director_name	num_critic_for_reviews	gross	genres	actor_1_name	movie_title	num_voted_users	num_user_for_reviews
4488	Sergio Leone	181.0	6100000.0	Western	Clint Eastwood	The Good, the Bad and the Ugly	503509	119
4029	Fernando Meriles	214.0	756397.0	Crime/Drama	Alma Roca	City of God	533200	
4747	Alisa Kutsawa	153.0	209051.0	Action/Adventure/Drama	Takashi Shimura	Seven Samurai	229012	
2373	Hayao Miyazaki	246.0	1040886.0	Adventure/Animation/Family/Fantasy	Bunta Sugawara	Spaced Out	419791	
4259	Florian Danneberg	215.0	1120457.0	Drama/Thriller	Sebastian Koch	The Lives of Others	259379	
4921	Majid Majidi	46.0	925402.0	Drama/Family	Bahare Sedigh	Children of Heaven	27882	
4155	Chen-wok Pak	305.0	2181290.0	Drama/Mystery/Thriller	Mink Choi	Oldboy	356131	
1296	Jean-Pierre Jeunet	242.0	3320161.0	Comedy/Romance	Mathieu Kassovitz	Amélie	534262	
4970	Walter Parkes	96.0	11423134.0	Adventure/Drama/Thriller/War	Gregory Peck	Dad Boat	165203	
4568	Farhad Fattahi	364.0	7094892.0	Drama/Mystery	Shahab Hosseini	A Separation	151812	
1329	S.S. Rajamouli	44.0	6488000.0	Action/Adventure/Drama/Fantasy/War	Tanishk Bhatia	Baahubali: The Beginning	62756	
2233	Hayao Miyazaki	174.0	2298191.0	Adventure/Animation/Fantasy	Minnie Driver	Picasso	221562	
4023	Thomas Vinterberg	349.0	610968.0	Drama	Thomas Bo Larsen	The Hunt	170355	
2734	Fritz Lang	260.0	26435.0	Drama/Sci-Fi	Ernst Han	Metropolis	113141	
2823	Oliver Hirschbiegel	181.0	5501940.0	Biography/Drama/History/War	Thomas Kretschmer	Downtfall	248564	
2921	Gulistan del Toro	406.0	67223143.0	Drama/Fantasy/War	Yanna Barreto	Para Luytén	467234	
3560	Denis Villeneuve	226.0	6867096.0	Drama/Mystery/War	Lubna Azabal	Inconceivable	80429	
4000	Juan José Campanella	262.0	20167424.0	Drama/Mystery/War	Ricardo Darín	The Secret in Their Eyes	131611	
2047	Hayao Miyazaki	212.0	4710405.0	Adventure/Animation/Family/Fantasy	Christian Bale	How's My Trip	214091	
2914	Je-kyu Kang	86.0	110186.0	Action/Drama/War	Mink Choi	The Outlaw Jose Yuma	31543	
3553	José Padilha	142.0	8060.0	Action/Crime/Drama/Thriller	Wagner Moura	Elite Squad	81644	
2830	Alejandro Amenábar	157.0	2083435.0	Biography/Drama/Romance	Belen Rueda	The Sea Inside	64566	
4267	Alexandre G. Ruffin	157.0	5381854.0	Drama/Thriller	Adriana Barrios	Amores Perros	173561	
4461	Thomas Vinterberg	98.0	1647780.0	Drama	Ulrich Thomsen	The Copenhagen	65961	
3423	Kabushiro Oosuo	150.0	439162.0	Action/Adventure/Sci-Fi/Thriller	Mitsuo Ikeda	Alita	126160	
4987	Sergio Leone	122.0	3500000.0	Action/Drama/Western	Clint Eastwood	A Fistful of Dollars	147566	
3344	Karan Johar	210.0	4016995.0	Adventure/Drama/Thriller	Shah Rukh Khan	My Name is Khan	69759	
4144	Walter Salles	71.0	5595429.0	Drama	Fernanda Montenegro	Central Station	28961	
1456	Vincent Paronnaud	242.0	444503.0	Animation/Biography/Drama/War	Catherine Deneuve	Persepolis	70394	
4284	Ali Fidan	231.0	2281276.0	Animation/Biography/Documentary/Drama/History/War	Ali Fidan	Walt with a Hat	46107	
3264	Michael Haneke	447.0	225377.0	Drama/Romance	Isabelle Huppert	Amour	70382	
4425	Fabrizio Borsari	94.0	1221261.0	Crime/Drama/Thriller	Ricardo Darín	Nine Queens	38215	
3677	Christopher YOUNG	112.0	3629758.0	Drama/Music	José Ignacio Marín	The Chorus	44151	
4640	Cristian Murgu	233.0	1183763.0	Drama	Aranzazu Alcaraz	4 Months, 3 Weeks and 2 Days	44703	
3510	Yosh Chikara	29.0	2921738.0	Drama/Musical/Romance	Shah Rukh Khan	Veer-Zaara	34449	
2963	Clint Eastwood	251.0	1375931.0	Drama/History/War	Yuki Matsukita	Letters from Iwo Jima	132149	
2493	Yimou Zhang	283.0	84961.0	Action/Adventure/History	Jie Li	Hero	148414	
1171	Yimou Zhang	283.0	84961.0	Action/Adventure/History	Jie Li	Hero	148414	

D.Best Directors: Group the column