

```
In [485]: 1 import pandas as pd
          2 from IPython.display import HTML
```

```
In [486]: 1 # Data Import and first Inspection
```

```
In [487]: 1 Data=pd.read_csv('movies_complete.csv')
```

### Some additional information on Features/Columns:

- **id:** The ID of the movie (clear/unique identifier).
- **title:** The Official Title of the movie.
- **tagline:** The tagline of the movie.
- **release\_date:** Theatrical Release Date of the movie.
- **genres:** Genres associated with the movie.
- **belongs\_to\_collection:** Gives information on the movie series/franchise the particular film belongs to.
- **original\_language:** The language in which the movie was originally shot in.
- **budget\_musd:** The budget of the movie in million dollars.
- **revenue\_musd:** The total revenue of the movie in million dollars.
- **production\_companies:** Production companies involved with the making of the movie.
- **production\_countries:** Countries where the movie was shot/produced in.
- **vote\_count:** The number of votes by users, as counted by TMDB.
- **vote\_average:** The average rating of the movie.
- **popularity:** The Popularity Score assigned by TMDB.
- **runtime:** The runtime of the movie in minutes.
- **overview:** A brief blurb of the movie.
- **spoken\_languages:** Spoken languages in the film.
- **poster\_path:** The URL of the poster image.
- **cast:** (Main) Actors appearing in the movie.
- **cast\_size:** number of Actors appearing in the movie.
- **director:** Director of the movie.
- **crew\_size:** Size of the film crew (incl. director, excl. actors).

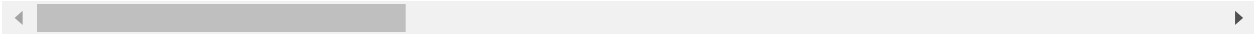
In [488]:

1Data.head(3)

Out[488]:

	id	title	tagline	release_date	genres	belongs_to_collection	original_title
0	862	Toy Story	NaN	1995-10-30	Animation Comedy Family	Toy Story Collection	
1	8844	Jumanji	Roll the dice and unleash the excitement!	1995-12-15	Adventure Fantasy Family		NaN
2	15602	Grumpier Old Men	Still Yelling. Still Fighting. Still Ready for...	1995-12-22	Romance Comedy	Grumpy Old Men Collection	

3 rows × 22 columns



In [489]:

1#Data=Data.set\_index('title')

In [490]: 1 Data.isnull().sum()

```
Out[490]: id                                0
          title                             0
          tagline                           24407
          release_date                       34
          genres                             2105
          belongs_to_collection              40228
          original_language                  10
          budget_musd                        35837
          revenue_musd                       37306
          production_companies               11335
          production_countries               5856
          vote_count                          0
          vote_average                       2614
          popularity                         0
          runtime                            1512
          overview                           951
          spoken_languages                   3597
          poster_path                        224
          cast                               2189
          cast_size                           0
          crew_size                           0
          director                           731
          dtype: int64
```

In [491]: 1 Data.shape

Out[491]: (44691, 22)

```
In [492]: 1 Data['profit']=Data['revenue_musd']-Data['budget_musd']
          2 Data['ROI']=(Data['revenue_musd']/Data['budget_musd'])*100
```





In [ ]: 1

```
In [493]: 1 #_function
          2
          3 def filter_data(Basis,Number_Of_Record,ascending,min_bud=0,minRating=0):
          4     B=Basis
          5     N=Number_Of_Record
          6     A=ascending
          7     subData=Data[(Data['budget_musd']>=min_bud)&(Data['vote_average']>=minRating)]
          8     Record=subData.sort_values(by=B,ascending=A)[0:N][['title','poster_path','E
          9     return HTML(Record.to_html(escape=False))
```

In [494]:

```
1 #__Movies Top 5 - Highest Revenue__
2 Top_5_movie_highest_Revenue=filter_data('revenue_musd',5,False)
3 Top_5_movie_highest_Revenue
4
```



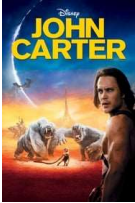
Out[494]:

	poster_path	revenue_musd
title		
Avatar		2787.965087
Star Wars: The Force Awakens		2068.223624
Titanic		1845.034188
The Avengers		1519.557910
Jurassic World		1513.528810

In [495]:






```
1 #__Movies Top 5 - Highest Budget__
2 Top_5_movie_highest_Budget=filter_data('budget_musd',5,False)
3 Top_5_movie_highest_Budget
```

Out[495]:

title		poster_path	budget_musd
Pirates of the Caribbean: On Stranger Tides			380.0
Pirates of the Caribbean: At World's End			300.0
Avengers: Age of Ultron			280.0
Superman Returns			270.0
John Carter			260.0

```
In [496]: 1 #__Movies Top 5 - Highest Profit__
          2 Top_5_movie_highest_Profit=filter_data('profit',5,False)
          3 Top_5_movie_highest_Profit
```

Out[496]:

	poster_path	profit
title		
Avatar		2550.965087
Star Wars: The Force Awakens		1823.223624
Titanic		1645.034188
Jurassic World		1363.528810
Furious 7		1316.249360

In [497]:

```
1 #__Movies Top 5 - Lowest Profit__
2 Top_5_movie_Lowest_Profit=filter_data('profit',5,True)
3 Top_5_movie_Lowest_Profit
```

Out[497]:

	poster_path	profit
title		
The Lone Ranger		-165.710090
The Alamo		-119.180039
Mars Needs Moms		-111.007242
Valerian and the City of a Thousand Planets		-107.447384
The 13th Warrior		-98.301101

In [498]:

1

`#__Movies Top 5 - Highest ROI__`

2

`#Highest Return on Investment (=Revenue / Budget) (only movies with Budget >`




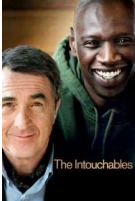

3

`Top_5_movie_Highest_ROI=filter_data('ROI',5,False,10)`

4

`Top_5_movie_Highest_ROI`




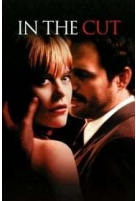

Out[498]:

	poster_path	ROI
title		
E.T. the Extra-Terrestrial		7552.050724
Star Wars		7049.072791
Pretty Woman		3307.142857
The Intouchables		3280.622085
The Empire Strikes Back		2991.111111



```
In [499]: 1 #Lowest Return on Investment (=Revenue / Budget) (only movies with Budget >=
2 #__Movies Top 5 - Lowest ROI__
3 Top_5_movie_Lowest_ROI=filter_data('ROI',5,True,10)
4 Top_5_movie_Lowest_ROI
```






Out[499]:

	poster_path	ROI
title		
Chasing Liberty		0.000052
The Cookout		0.000075
Deadfall		0.000180
In the Cut		0.000192
The Samaritan		0.021008

In [500]:

```
1 #__Movies Top 5 - Most Votes__
2 #Lowest Rating (only movies with 10 or more Ratings)
3 Top_5_movie_Most_Votes=filter_data('vote_count',5,False)
4 Top_5_movie_Most_Votes
```

Out[500]:

	poster_path	vote_count
title		
Inception		14075.0
The Dark Knight		12269.0
Avatar		12114.0
The Avengers		12000.0
Deadpool		11444.0

In [501]:

1

`#__Movies Top 5 - Highest Rating__`

2

`#Highest Rating (only movies with 10 or more Ratings)`

3





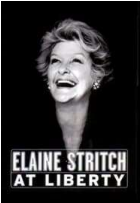
`Top_5_movie_Highest_Rating=filter_data('vote_average',5,False,minRating=0)`

4

`Top_5_movie_Highest_Rating`






5

Out[501]:

	poster_path	vote_average
	title	
	<div><div></div><div>Time Pass</div><div></div></div>	10.0
	<div><div></div><div>Shuttlecock Boys</div><div></div></div>	10.0
	<div><div></div><div>Forever</div><div></div></div>	10.0
	<div><div></div><div>Souls of Zen: Ancestors and Agency in Contemporary Japanese Temple Buddhism</div><div></div></div>	10.0
	<div><div></div><div>Elaine Stritch: At Liberty</div><div></div></div>	10.0






```
In [502]: 1 #Lowest Rating (only movies with 10 or more Ratings)
2 Top_5_movie_Lowest_Rating=filter_data('vote_average',5,True,minRating=0)
3 Top_5_movie_Lowest_Rating
4
```

Out[502]:

	poster_path	vote_average
title		
Extinction: Nature Has Evolved		0.0
		
Roukli		0.0
		0.0
		0.0
		
Unrated II: Scary as Hell		0.5

```
In [503]: 1 #__Movies Top 5 - Most Popular__
          2 Top_5_movie_Most_Popular=filter_data('popularity',5,False)
          3 Top_5_movie_Most_Popular
```

Out[503]:





	poster_path	popularity
title		
		
Minions		547.488298
		
Wonder Woman		294.337037
		
Beauty and the Beast		287.253654
		
Baby Driver		228.032744
		
Big Hero 6		213.849907

## Find your next Movie

```
In [504]: 1 #__Search 1: Science Fiction Action Movie with Bruce Willis (sorted from hig
```

```
In [505]: 1 Gn=Data[Data['genres'].str.contains("Action") & Data['genres'].str.contains(
2 Ac=Data['cast'].str.contains('Bruce Willis')
3 sub=Data.loc[(Ac & str(Gn)), ["title", "vote_average", 'poster_path']].sort_v
4 HTML(sub.to_html(escape=False))
```


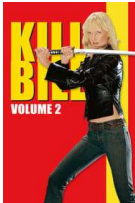

Out[505]:

	title	vote_average	poster_path
291	Pulp Fiction	8.3	
2620	The Sixth Sense	7.7	
18846	Moonrise Kingdom	7.6	
35664	B.B. King: The Life of Riley	7.5	
998	Die Hard	7.5	

```
In [506]: 1 #__Search 2: Movies with Uma Thurman and directed by Quentin Tarantino (Low
```

```
In [507]: 1 sub=Data[(Data['cast'].str.contains('Uma Thurman')) & (Data['director']=='Qu  
2 sub=sub[['title','release_date','runtime','poster_path']]  
3 HTML(sub.to_html(escape=False))
```

Out[507]:






	title	release_date	runtime	poster_path
6667	Kill Bill: Vol. 1	2003-10-10	111.0	
7208	Kill Bill: Vol. 2	2004-04-16	136.0	
291	Pulp Fiction	1994-09-10	154.0	

```
In [508]: 1 #__Search 3: Most Successful Pixar Studio Movies between 2010 and 2015 (high
```

```
In [509]: [1'production_companies']==( 'Pixar Animation Studios')].sort_values(by='revenue_mu
```

```
In [510]: 1 sub=sub[sub['release_date'].between("2010-01-01", "2015-12-31")][['title','p
2 HTML(sub.to_html(escape=False))
```

Out[510]:






	title	production_companies	poster_path
21694	The Blue Umbrella	Pixar Animation Studios	
22489	La luna	Pixar Animation Studios	
24252	Hawaiian Vacation	Pixar Animation Studios	
24254	Small Fry	Pixar Animation Studios	
25516	The Legend of Mor'du	Pixar Animation Studios	

```
In [511]: 1 #__Search 4: Action or
2 #Thriller Movie with original Language English and minimum Rating of 7.5 (mo
```



```
In [512]: 1 sub=Data[Data['genres'].str.contains('Action') & Data['genres'].str.contains
2 sub=sub[(sub['spoken_languages']=='English')& (sub['vote_average']>=7.5)][
3 HTML(sub.to_html(escape=False))
```

Out[512]:

		genres	poster_path	vote_average	release_date
title					
The Dark Knight Rises		Action Crime Drama Thriller		7.6	2012-07-16
Ghost Recon: Alpha		Action Science Fiction Thriller War		7.5	2012-05-03
Oxy-Morons		Action Thriller		8.0	2011-10-02
Inception		Action Thriller Science Fiction Mystery Adventure		8.1	2010-07-14
Prison Break: The Final Break		Action Drama Thriller		7.5	2009-05-26

```
1 Data.head(2)
```

```
In [ ]: 1
```

```
In [513]: 1 #4. __Analyze__ the Dataset and __find out whether Franchises
2 #(Movies that belong to a # collection)are more successful
3 #than stand-alone movies__ in terms of:
```

- mean revenue
- median Return on Investment
- mean budget raised
- mean popularity
- mean rating

```
In [514]: 1 Data['Franchise']=Data['belongs_to_collection'].notna()
```

```
In [515]: 1 #__Franchise vs. Stand-alone: Average Revenue__
```

```
In [516]: 1 Data.groupby(Data['Franchise'])['revenue_musd'].mean().sort_values(ascending
```

```
Out[516]: Franchise
True      165.708193
False     44.742814
Name: revenue_musd, dtype: float64
```

```
In [517]: 1 #__Franchise vs. Stand-alone: Return on Investment / Profitability__
```

```
In [518]: 1 Data.groupby(Data['Franchise'])['ROI'].median().sort_values(ascending=False)
```

```
Out[518]: Franchise
True      370.919508
False     161.969933
Name: ROI, dtype: float64
```

```
In [519]: 1 #__Franchise vs. Stand-alone: Average Budget__
2 Data.groupby(Data['Franchise'])['budget_musd'].mean().sort_values(ascending=
```

```
Out[519]: Franchise
True      38.319847
False     18.047741
Name: budget_musd, dtype: float64
```

```
In [520]: 1 #__Franchise vs. Stand-alone: Average Popularity__
2 Data.groupby(Data['Franchise'])['popularity'].mean().sort_values(ascending=F
```

```
Out[520]: Franchise
True      6.245051
False     2.592726
Name: popularity, dtype: float64
```

```
In [521]: 1 #_Franchise vs. Stand-alone: Average Rating__
          2 Data.groupby(Data['Franchise'])['vote_average'].mean().sort_values(ascending
```

```
Out[521]: Franchise
False      6.008787
True       5.956806
Name: vote_average, dtype: float64
```

## Most Successful Franchises

5. Find the most successful Franchises in terms of

- total number of movies
- total & mean budget
- total & mean revenue
- mean rating

```
In [522]: 1 Fr=Data.groupby(Data['belongs_to_collection']).agg({'title':"count",'budget_
```

```
In [523]: 1 #total number of movies
          2 Fr.sort_values(by=('title', 'count'),ascending=False).head(1)
```

```
Out[523]:
```

	title	budget_musd		revenue_musd		vote_average	
	count	sum	mean	sum	mean	mean	
<b>belongs_to_collection</b>							
	The Bowery Boys	29	0.0	NaN	0.0	NaN	6.675

```
In [524]: 1 #total & mean budget
          2 Fr.sort_values(by=('budget_musd', 'mean'),ascending=False).head(1)
```

```
Out[524]:
```

	title	budget_musd		revenue_musd		vote_average	
	count	sum	mean	sum	mean	mean	
<b>belongs_to_collection</b>							
	Tangled Collection	2	260.0	260.0	591.794936	591.794936	7.25

```
In [525]: 1 #mean rating
          2 Fr.sort_values(by=('vote_average', 'mean'), ascending=False).head(1)
```

Out[525]:

	title	budget_musd		revenue_musd		vote_average
	count	sum	mean	sum	mean	mean
belongs_to_collection						
Argo Collection	1	0.0	NaN	0.0	NaN	9.3

## Most Successful Directors

6. Find the most successful Directors in terms of

- **total number of movies**
- **total revenue**
- **mean rating**

```
In [526]: 1 Dire=Data.groupby(Data['director']).agg({'title':'count','revenue_musd':'sum'
```

```
In [527]: 1 #- __total number of movies__
          2 Dire.sort_values(by='title',ascending=False).head(1)
```

Out[527]:

	title	revenue_musd	vote_average
director			
John Ford	66	85.170757	6.381818

```
In [528]: 1 #- __total revenue__
          2 Dire.sort_values(by='revenue_musd',ascending=False).head(1)
```

Out[528]:

director	title	revenue_musd	vote_average
Steven Spielberg	33	9256.621422	6.893939

```
In [529]: 1 #- __mean rating__
          2 Dire.sort_values(by='vote_average',ascending=False).head(1)
          3
```

Out[529]:

	title	revenue_musd	vote_average
director			
Antonis Sotiropoulos	1	0.0	10.0

END

In [ ]:

1	
---	--