# IMDB Movie Analysis

SOURAV PATTANAYAK

# Table of contents

# Project description

A designated data analyst has an important role in evaluating the processes and progresses of organizations, with the help of given datasets a data analyst performs few steps to reach the insight.

In this project I was provided with IMDb movie dataset containing various columns for my final project and by applying the 'Five whys' approach and data analysis skills the answers of asked queries were explained.

The entire dataset was cleaned before taking it any further to reach the required insights.

# Project objectives

- To inspect the **Root Cause Analysis**.

- To clean the data.

- To find the movies with highest profit.

- To find IMDb top 250 movies based on the number of voted users and extracting the non-english movies from the list.

- To find the best directors.

- To find popular movie genres.

- To find the critic favourite and audience favourite actors.

- To find the decade with highest votes.

# Approach

Basic and advanced data analysis methods were used to calculate the queries, inspect the queries and reach the desired insights, data was cleaned before the 'Five whys' approach was used extensively to understand the Root cause analysis.

Different tools were used to visualize and represent the results and later the information were gathered and loopholes were found to jot them down.
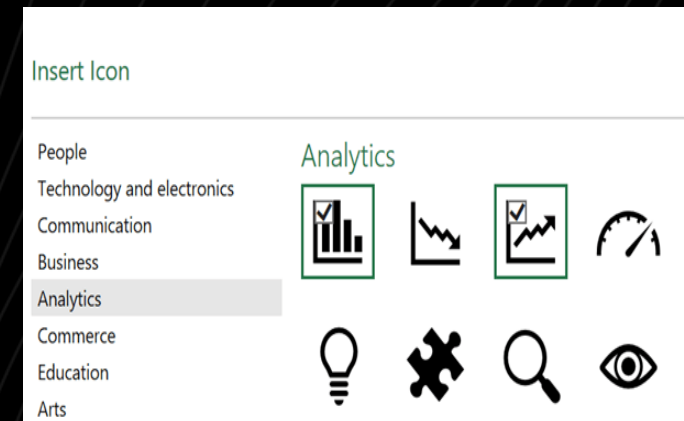
# Methodology and Tech-Stack used

- ✓ Microsoft excel tool was used to complete the entire project.

- ✓ To evaluate the queries Microsoft advanced excel formulas were used.

- ✓ Some of the queries were visualized using Microsoft excel charts for better explanation.

- ✓ Microsoft Powerpoint presentation was used to prepare the presentation.

# Insights

## 1.Data cleaning :

- Data cleaning is one of the most important step that needs to be taken before analyzing the dataset.

- In the given dataset the data was cleaned after revising the queries by deleting unnecessary columns, null values etc.

- Initially I had 5044 row and 28 columns and after cleaning I have 3785 rows with 14 columns.

- Basic excel methods(e.g:filter columns, sorting) were used to exclude  null values and sort the data based on requirements.

- The duplicate values were removed using 'Data' tab ⟶ Remove duplicate values/cells.

- *To go to the cleaned dataset after the operations were done,* click here.

# Insights

## Preview before cleaning:

# Insights

Preview after cleaning:

| | director_name | num_critic_for_reviews | duration | gross | genres | actor_1_name | movie_title | num_voted_users | num_user_for_reviews | language | country | budget | title_year | imdb_score | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | director_name | num_critic_for_reviews | duration | gross | genres | actor_1_name | movie_title | num_voted_users | num_user_for_reviews | language | country | budget | title_year | imdb_score | |
| 2 | James Cameron | 723 | 178 | 7.6E+08 | Action\|Adventure | CCH Pounder | Avatar | 886204 | 3054 | English | USA | 2.4E+08 | 2009 | 7.9 | |
| 3 | Gore Verbinski | 302 | 169 | 3.1E+08 | Action\|Adventure | Johnny Depp | Pirates of the Caribbean: At World's End | 471220 | 1238 | English | USA | 3E+08 | 2007 | 7.1 | |
| 4 | Sam Mendes | 602 | 148 | 2E+08 | Action\|Adventure | Christoph Waltz | Spectre | 275868 | 994 | English | UK | 2.5E+08 | 2015 | 6.8 | |
| 5 | Christopher Nolan | 813 | 164 | 4.5E+08 | Action\|Thriller | Tom Hardy | The Dark Knight Rises | 1144337 | 2701 | English | USA | 2.5E+08 | 2012 | 8.5 | |
| 7 | Andrew Stanton | 462 | 132 | 7.3E+07 | Action\|Adventure | Daryl Sabara | John Carter | 212204 | 738 | English | USA | 2.6E+08 | 2012 | 6.6 | |
| 8 | Sam Raimi | 392 | 156 | 3.4E+08 | Action\|Adventure | J.K. Simmons | Spider-Man 3 | 383056 | 1902 | English | USA | 2.6E+08 | 2007 | 6.2 | |
| 9 | Nathan Greno | 324 | 100 | 2E+08 | Adventure\|Animation | Brad Garrett | Tangled | 294810 | 387 | English | USA | 2.6E+08 | 2010 | 7.8 | |
| 10 | Joss Whedon | 635 | 141 | 4.6E+08 | Action\|Adventure | Chris Hemsworth | Avengers: Age of Ultron | 462669 | 1117 | English | USA | 2.5E+08 | 2015 | 7.5 | |
| 11 | David Yates | 375 | 153 | 3E+08 | Adventure\|Family | Alan Rickman | Harry Potter and the Half-Blood Prince | 321795 | 973 | English | UK | 2.5E+08 | 2009 | 7.5 | |
| 12 | Zack Snyder | 673 | 183 | 3.3E+08 | Action\|Adventure | Henry Cavill | Batman v Superman: Dawn of Justice | 371639 | 3018 | English | USA | 2.5E+08 | 2016 | 6.9 | |
| 13 | Bryan Singer | 434 | 169 | 2E+08 | Action\|Adventure | Kevin Spacey | Superman Returns | 240396 | 2367 | English | USA | 2.1E+08 | 2006 | 6.1 | |
| 14 | Marc Forster | 403 | 106 | 1.7E+08 | Action\|Adventure | Giancarlo Giannini | Quantum of Solace | 330784 | 1243 | English | UK | 2E+08 | 2008 | 6.7 | |
| 15 | Gore Verbinski | 313 | 151 | 4.2E+08 | Action\|Adventure | Johnny Depp | Pirates of the Caribbean: Dead Man's Chest | 522040 | 1832 | English | USA | 2.3E+08 | 2006 | 7.3 | |
| 16 | Gore Verbinski | 450 | 150 | 8.9E+07 | Action\|Adventure | Johnny Depp | The Lone Ranger | 181792 | 711 | English | USA | 2.2E+08 | 2013 | 6.5 | |
| 17 | Zack Snyder | 733 | 143 | 2.9E+08 | Action\|Adventure | Henry Cavill | Man of Steel | 548573 | 2536 | English | USA | 2.3E+08 | 2013 | 7.2 | |
| 18 | Andrew Adamson | 258 | 150 | 1.4E+08 | Action\|Adventure | Peter Dinklage | The Chronicles of Narnia: Prince Caspian | 149922 | 438 | English | USA | 2.3E+08 | 2008 | 6.6 | |
| 19 | Joss Whedon | 703 | 173 | 6.2E+08 | Action\|Adventure | Chris Hemsworth | The Avengers | 995415 | 1722 | English | USA | 2.2E+08 | 2012 | 8.1 | |
| 20 | Rob Marshall | 448 | 136 | 2.4E+08 | Action\|Adventure | Johnny Depp | Pirates of the Caribbean: On Stranger Tides | 370704 | 484 | English | USA | 2.5E+08 | 2011 | 6.7 | |
| 21 | Barry Sonnenfeld | 451 | 106 | 1.8E+08 | Action\|Adventure | Will Smith | Men in Black 3 | 268154 | 341 | English | USA | 2.3E+08 | 2012 | 6.8 | |
| 22 | Peter Jackson | 422 | 164 | 2.6E+08 | Adventure\|Fantasy | Aidan Turner | The Hobbit: The Battle of the Five Armies | 354228 | 802 | English | New Zealand | 2.5E+08 | 2014 | 7.5 | |
| 23 | Marc Webb | 599 | 153 | 2.6E+08 | Action\|Adventure | Emma Stone | The Amazing Spider-Man | 451803 | 1225 | English | USA | 2.3E+08 | 2012 | 7 | |
| 24 | Ridley Scott | 343 | 156 | 1.1E+08 | Action\|Adventure | Mark Addy | Robin Hood | 211765 | 546 | English | USA | 2E+08 | 2010 | 6.7 | |
| 25 | Peter Jackson | 509 | 186 | 2.6E+08 | Adventure\|Fantasy | Aidan Turner | The Hobbit: The Desolation of Smaug | 483540 | 951 | English | USA | 2.3E+08 | 2013 | 7.9 | |
| 26 | Chris Weitz | 251 | 113 | 7E+07 | Adventure\|Family | Christopher Lee | The Golden Compass | 149019 | 666 | English | USA | 1.8E+08 | 2007 | 6.1 | |
| 27 | Peter Jackson | 446 | 201 | 2.2E+08 | Action\|Adventure | Naomi Watts | King Kong | 316018 | 2618 | English | New Zealand | 2.1E+08 | 2005 | 7.2 | |
| 28 | James Cameron | 315 | 194 | 6.6E+08 | Drama\|Romance | Leonardo DiCaprio | Titanic | 793059 | 2528 | English | USA | 2E+08 | 1997 | 7.7 | |

# Insights

## 2A.Top 10 movies with highest profit:

- We have gross and budget columns in our dataset, the profit was calculated by subtracting budget from the gross amount and was shown as a new column.

- The profit column was sorted with top 10 after selecting the column and this has given us the top 10 movies with highest profit.

- A graph of Profit vs Top 10 movies was plotted for better understanding.



**Top 10 most profitable movies**

# Insights

## 2B. Observing the outliers:

✓ Major outliers are on points:

  2127.52,2400,2500,1100 and 4200

  having the R2 Score 0.0813.

✓ The profit and budget were divided by 100000

  For easier calculation.



Profit VS Budget

# Insights

## 3. IMDb Top 250 movies:

- ✓ Using sort and filter method,the movies >25000 num_voted_users were filtered.

- ✓ The IMDb scores were arranged in a decreasing order and only top 250 rows were chosen for evaluation.

- ✓ A new Rank column was added to rank the movies from 1 to 250(Using cell+1,drag and fill method)

- ✓ By unselecting the 'English' language from language column we get the movies of foreign_language under IMDbs top 250 list.

# Top 250 of IMDb preview:

| Rank | IMDb_Top_250 | imdb_score | num_voted_users | Top_Foreign_Lang_Film |
|---|---|---|---|---|
| 1 | The Shawshank Redemption | 9.3 | 1689764 | The Good, the Bad and the Ugly |
| 2 | The Godfather | 9.2 | 1155770 | City of God |
| 3 | The Dark Knight | 9 | 1676169 | Seven Samurai |
| 4 | The Godfather: Part II | 9 | 790926 | Spirited Away |
| 5 | The Lord of the Rings: The Return of the King | 8.9 | 1215718 | Samsara |
| 6 | Schindler's List | 8.9 | 865020 | |
| 7 | Pulp Fiction | 8.9 | 1324680 | |
| 8 | The Good, the Bad and the Ugly | 8.9 | 503509 | |
| 9 | Inception | 8.8 | 1468200 | |
| 10 | The Lord of the Rings: The Fellowship of the Ring | 8.8 | 1238746 | |
| 11 | Fight Club | 8.8 | 1347461 | |
| 12 | Forrest Gump | 8.8 | 1251222 | |
| 13 | Star Wars: Episode V - The Empire Strikes Back | 8.8 | 837759 | |
| 14 | The Lord of the Rings: The Two Towers | 8.7 | 1100446 | |
| 15 | The Matrix | 8.7 | 1217752 | |
| 16 | Goodfellas | 8.7 | 728685 | |
| 17 | Star Wars: Episode IV - A New Hope | 8.7 | 911097 | |
| 18 | One Flew Over the Cuckoo's Nest | 8.7 | 680041 | |
| 19 | City of God | 8.7 | 533200 | |
| 20 | Seven Samurai | 8.7 | 229012 | |

# Insights

## 4.Finding the best directors:

To find the best directors based on the average imdb scores,

- ✓ A pivot table of Director_name and imdb_score was created.

- ✓ In the values section of the pivot table the mean(average) of imdb_score was calculated.

- ✓ The values were sorted in descending order.

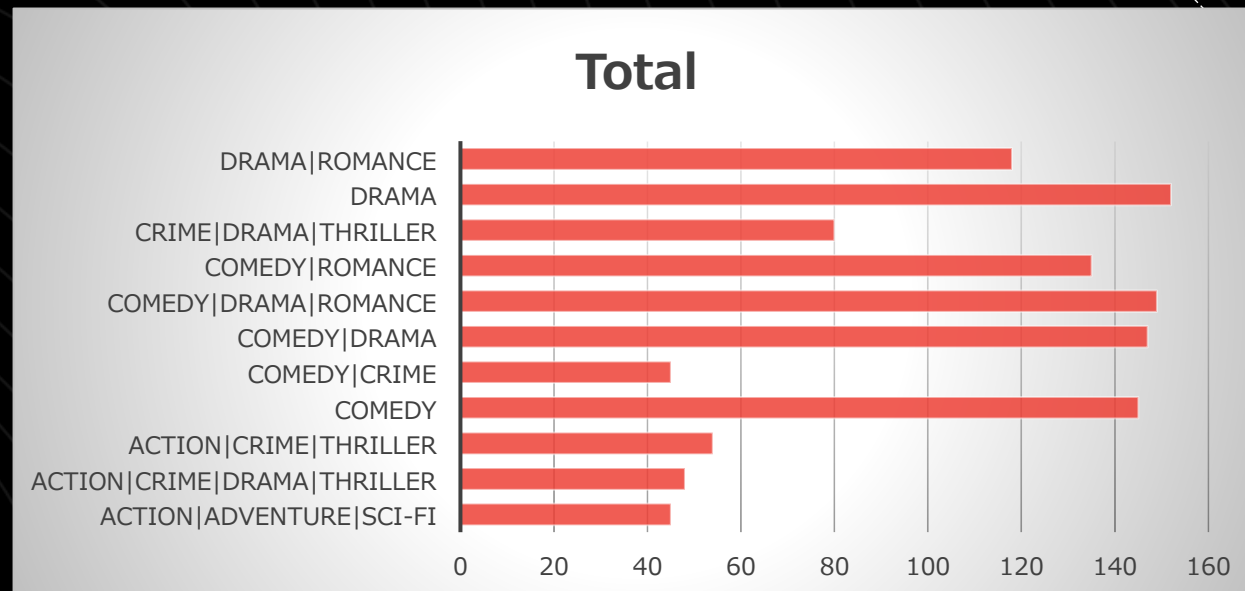| 3 | Director_name | Average of imdb_score |
|---|---------------|----------------------|
| 4 | Tony Kaye | 8.60 |
| 5 | Charles Chaplin | 8.60 |
| 6 | Alfred Hitchcock | 8.50 |
| 7 | Ron Fricke | 8.50 |
| 8 | Damien Chazelle | 8.50 |
| 9 | Majid Majidi | 8.50 |
| 10 | Sergio Leone | 8.43 |
| 11 | Christopher Nolan | 8.43 |
| 12 | S.S. Rajamouli | 8.40 |
| 13 | Richard Marquand | 8.40 |
| 14 | Marius A. Markevicius | 8.40 |
| 15 | Asghar Farhadi | 8.40 |

# Insights

## 5.Finding popular genres:

To find the popular genre based on the average imdb scores,

✓ A pivot table of genres and imdb_score was created.

✓ In the values section of the pivot table the mean(average) of imdb_score was calculated.

✓ The values were sorted in descending order.

✓ For better understanding, genre-wise count and average profit was also calculated using pivot table.



Total

| | A | B | C | D |
|---|---|---|---|---|
| 1 | Genre | Average imdb_score | genres count | Average Profit |
| 2 | Action\|Adventure\|Sci-Fi | 6.67 | 45 | 41416969.69 |
| 3 | Action\|Crime\|Drama\|Thriller | 6.52 | 48 | -1516065.438 |
| 4 | Action\|Crime\|Thriller | 6.40 | 54 | 8055116.611 |
| 5 | Comedy | 5.84 | 145 | 21034780.79 |
| 6 | Comedy\|Crime | 6.04 | 45 | 16373313.49 |
| 7 | Comedy\|Drama | 6.58 | 147 | 8344097.81 |
| 8 | Comedy\|Drama\|Romance | 6.50 | 149 | 10352425.22 |
| 9 | Comedy\|Romance | 5.90 | 135 | 20002148.13 |
| 10 | Crime\|Drama\|Thriller | 6.87 | 80 | 7115124.925 |
| 11 | Drama | 7.04 | 152 | 1314534.941 |
| 12 | Drama\|Romance | 6.95 | 118 | 9366491.754 |
| 13 | | | | |

# Insights

## 6A. Inspecting the critic and audience favourite actors:

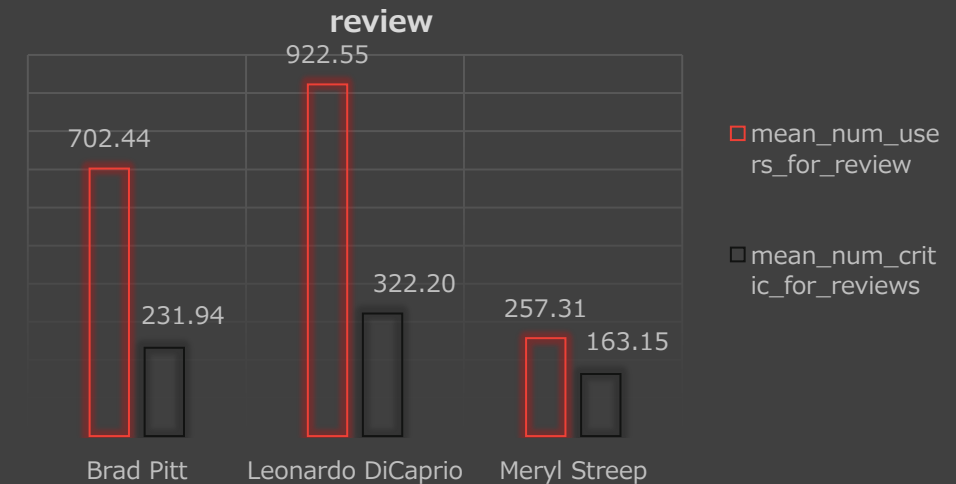Since the inspection was asked on specific three actors,

✓ Three separate columns Meryl_Streep, Leo_Caprio, and Brad_Pitt were created.

✓ These three columns contain movies where the said actors were lead actors , I have used pivot table for this action.

✓ The actor_1_name column containing the lead actor names under respective movies was used for extracton.

✓ In the values section of pivot table, mean of num_critic_for_reviews and mean of num_user_for_reviews were put against the actor names.

| | | |
|---|---|---|
| 2 | | |
| 3 | Row Labels | mean_num_users_for_review | mean_num_critic_for_reviews |
| 4 | Brad Pitt | 702.44 | 231.94 |
| 5 | Leonardo DiCaprio | 922.55 | 322.20 |
| 6 | Meryl Streep | 257.31 | 163.15 |
| 7 | | |
| 8 | | |

**Comparison between mean users review and mean critics review**

922.55

702.44

231.94

322.20

257.31

163.15

☐ mean_num_users_for_review

☐ mean_num_critic_for_reviews

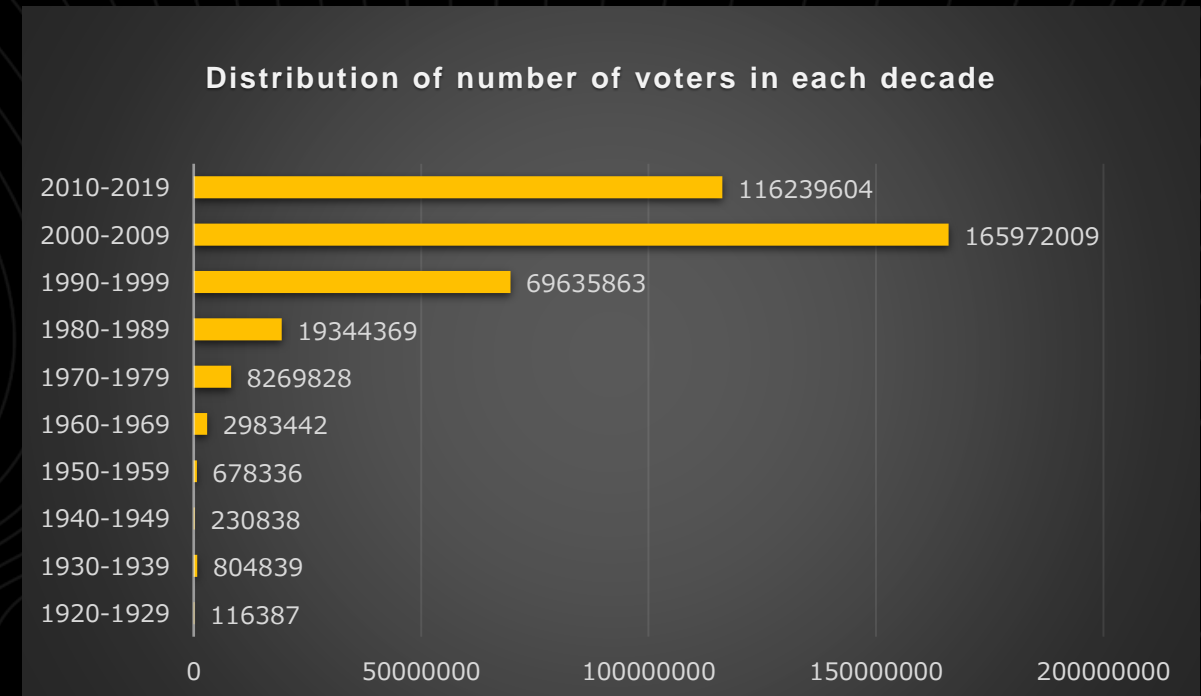Brad Pitt     Leonardo DiCaprio     Meryl Streep

# Insights

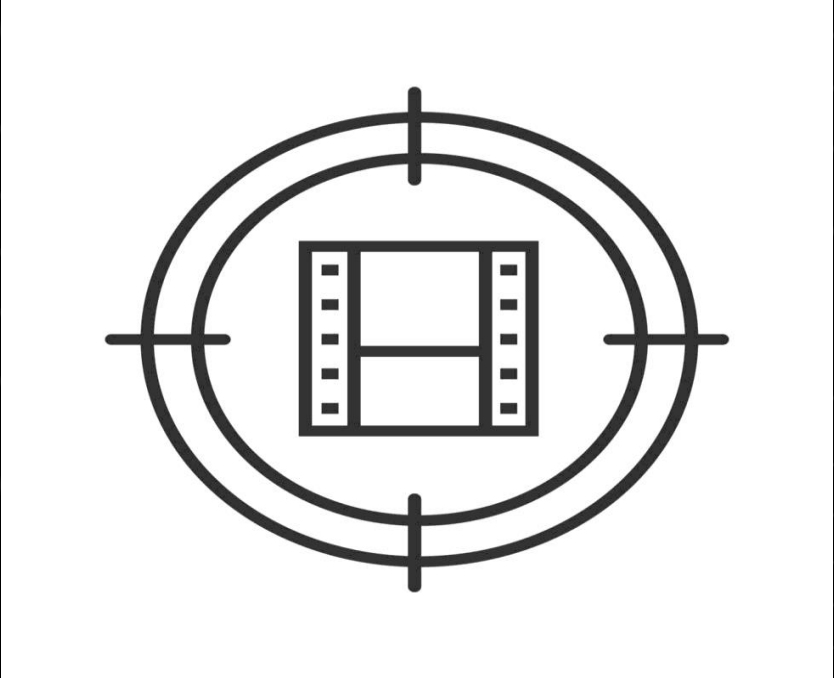## 6B. Calculating decade-wise voted users:

✓ For this calculation I have created one new column names 'Decade' that divides the title_year into decades of 10 years(e.g: 1920-1929).

**Used code for creating the decade column:**

```
=CONCATENATE(LEFT(A3, 3), "9-", LEFT(A3, 3), "0")
```

✓ Dragged and filled the rows using this formula/code.

✓ Created a pivot table, dragged the Decade column and sorted it , in the values section calculated the sum of num_voted_users to get Decade wise number of voted users.



**Distribution of number of voters in each decade**

| Decade | Voters |
|--------|--------|
| 2010-2019 | 116239604 |
| 2000-2009 | 165972009 |
| 1990-1999 | 69635863 |
| 1980-1989 | 19344369 |
| 1970-1979 | 8269828 |
| 1960-1969 | 2983442 |
| 1950-1959 | 678336 |
| 1940-1949 | 230838 |
| 1930-1939 | 804839 |
| 1920-1929 | 116387 |

# Key findings

➤ The movie with highest profit is **Avatar(**523505847) followed by **Jurassic world(502177271) ,Titanic(458672302)** and seven other movies.

➤ **The Shawshank Redemption** is the highest rated movie and also number one movie in IMDb and IMDb top 250 list with highest number of voted users.

➤ **The foreign language movies in IMDbs top 250 list are:** 1.The Good the Bad and the Ugly 2. Seven Samurai 3.Samsara 4.City of God 5.Spirited Away

➤ **Tony Kaye and  Charles Chaplin** are the top two Best directors amount the list of ten holding the same average IMDb score(8.60)

➤ **Drama** is the most popular movie genre holding 152 counts.

➤ **Leonardo DiCaprio i**s the highest user and critic rated actor , the most popular actor.

➤ The highest number of movie voters were seen in the **2000-2009 decade** holding 165972009 votes.

# Achievements

- **Data handling**: The entire project has helped me in learning methods to handle large dataset and tract records of each column.

- **Data cleaning**: This project has helped in understanding the data cleaning process after inspecting required queries.

- **Uses of Microsoft excel as business tool**: learned how Microsoft excel is being used with it's entire diversity in the process of statistical calculations.

- **Excel filters , sorting and formulas:** Learned different methods of filtering and sorting data based on requirements and implementation of excel formulas when needed.

- **The power of Powerpoint**: Learned how Powerpoint can strongly help explaining complex business queries in simple ways and slides.

# THANK YOU