# IMDB MOVIE ANALYSIS

*Project description*: This project is based upon the analysis of imdb movies,

Where a dataset has been provided which contains 29 columns and 5044 rows.Well it's a big dataset so,in the dataset we may get noisey data or unclean data or missing value or any errors present so first step is to we should clean the data then we should find patterns and analyse them.

Cleaning data: In the data I got so many missing values and noisy data,where I got rid of them by using some cleaning tasks. After then using statistics I analysed which column is useful which column is relatable by using correlations and I removed some of the columns which are not useful for analysing and which are not interrelated to each other. After cleaning the data my data got reduced to 15 columns and 3837 rows which made me felt more easy to analyse the data. Now we can say my data is clean data or good data.
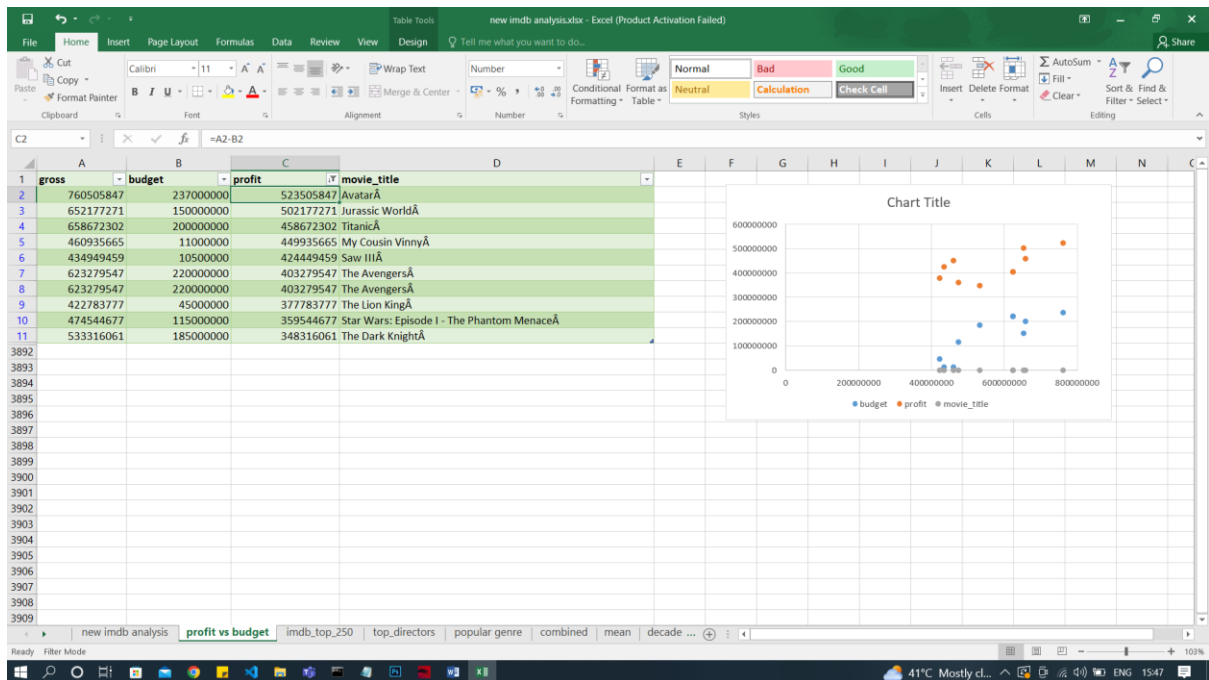
*Tech-stack used*: Ms-Execl

*Approach*:

After completion of cleaning I started performing the tasks which are:

A. Cleaning the data

B. Movies with the highest profit

C. Top 250

D. Best Directors

E. Popular Genres

F. Charts

As above I said I have completed the cleaning task. Then the next step is to find the movies with the highest profit.

**High profit**: first I created a new column called profit. Which was difference between the two columns of gross and budget. Then after I plotted a scatter plot chart of profit vs budget and observed which movie got the highest profit. The movie with the highest profit was Avatar followed by Jurassic park and titanic.

**Top 250**: In this task I have created a table consisting of movie name, language, imdb_score, num_voted_users, and created a column called rank where I ranked all the movies from 1 to 250. And I made sure that the num_voted_users are greater than 25000. After extracting all the 250 movies I also extracted the movies from top 250 movies which are not English language and stored them in a table called top_foreign_language.



To calculate rank I have used the formula:
=RANK.EQ($C2,$C$2:$C$253)+COUNTIFS($C$2:$C$253,$C2,$D$2:$D$253,">"&

$D2). The top 1 movie is The Shawshank Redemption with the imdb score of 9.3.

**Best Directors**: In this task I have retrieved the director name column,their movies and imdb_score and after I created a pivot table and then created a pivot table and then after I calculated mean of the imdb_score and then sorted them in alphabetical order.



The top director is Akira Kurosawa of 8.7 score of imdb.

**Popular Genres**: In this task I created a new table which consists director name,genres,gross,num_voted_users and then created a pivot table in the existing worksheet of Genres and sum_of_voted_users and sum_of_gross. According to the data we can see that the most watched genres is action
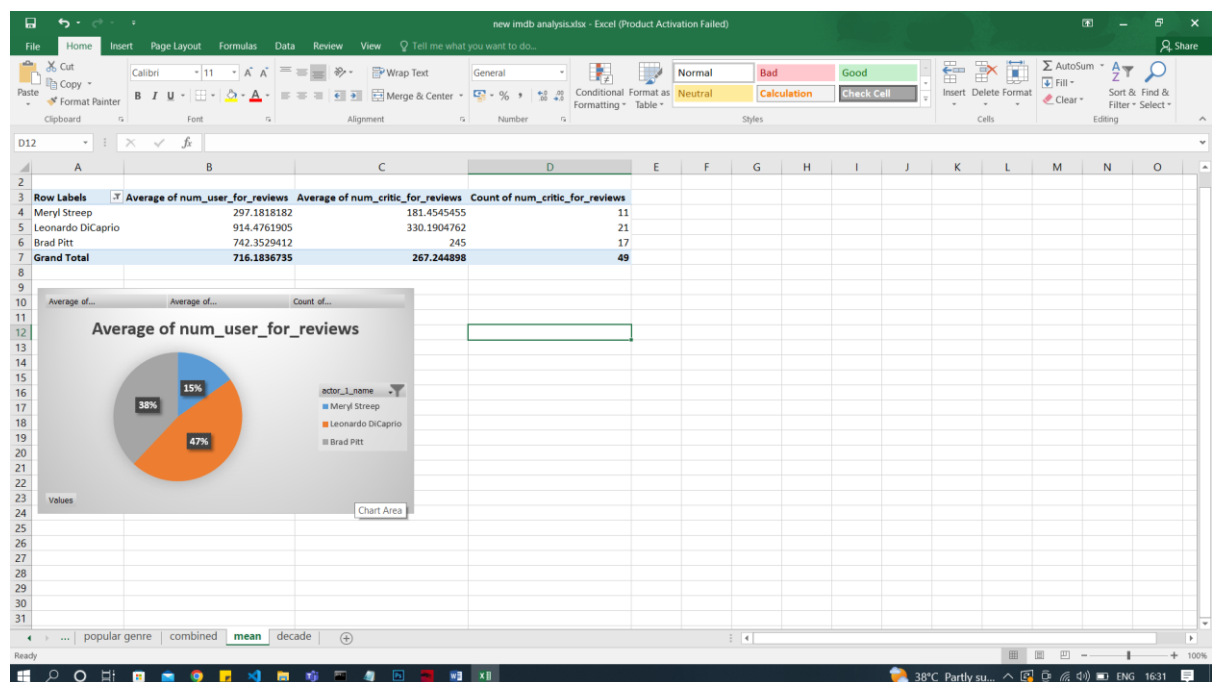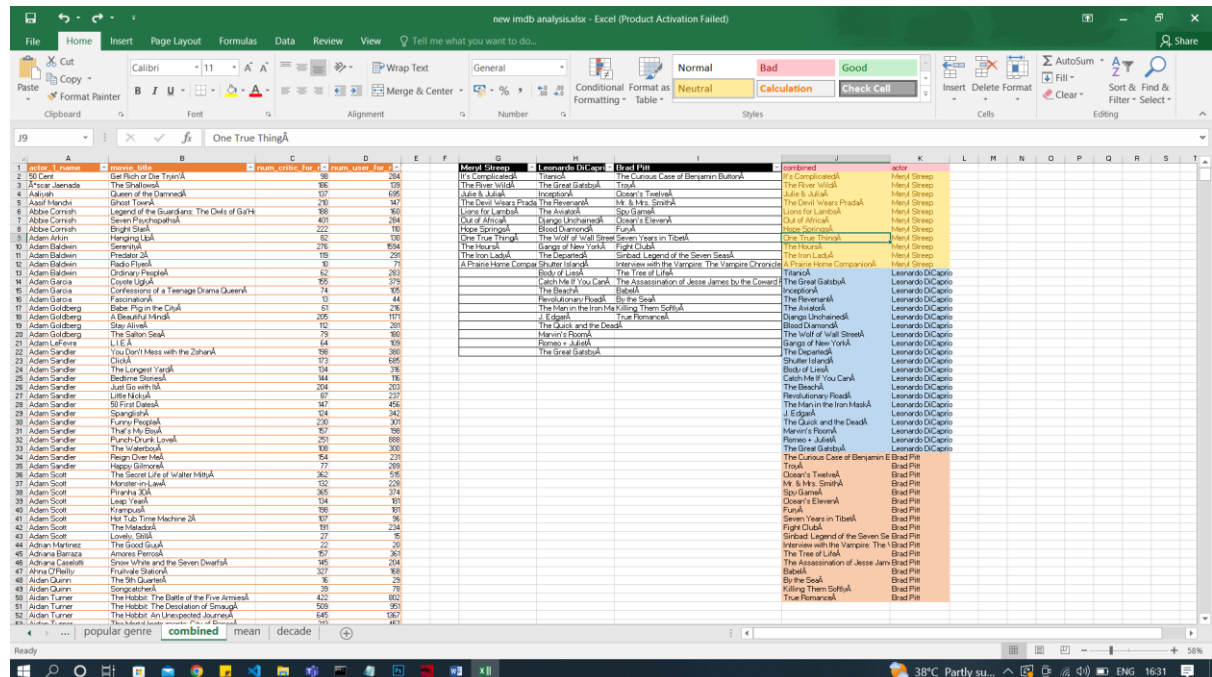
adventure and scifi.



**Charts**: In this task I have created three new columns called Meryl_streep,Leonardo_Di_Caprio and brad_pitt by retrieving from the table I created with the actor name,movie title ,num_critic_for_review and num_user_from_review and then I sorted the particular movies to particular actors and after that I created another column called combined where I combined the actors and their movies by group by and then calculated mean of num_user_for_reviews and num_critic_for_review by using pivot table and then analysed who is the critic favourite and audience favourite and created a pie chart. The favourite actor is Leonardo di caprio with 47% and also has the

highest mean .





After that I have created pivot table for years and sum of num_voted_users

And then I created a timeline  For example, the title_year year 1923, 1925 should be stored as 1920s. which I named the column called decade and sorted it and calculated sum of users voted in each decade and created a bar chart.

From 2001-2010 has the highest total no.of voted users of 179965615.