# Analysis Report on Movies' Overviews Dataset

## Introduction:

The success of a movie plays a crucial role in movie industry because of the huge investments and the fame. In recent days, movies are produced and released in a high rate. Even though the movies are costing over million dollars to produce, most of them are still unsuccessful or worse, flop. Movie success depends upon the various attributes. So, to analyse this success, different predictions are made, and new patterns are discovered in data. These results can be used by the moviemakers before planning a movie.

In this project, the three analysis methods are performed on the dataset: Correlation analysis, Association analysis, and K-means clustering analysis.

From the correlation analysis, we want to know how much the quantitative attributes: the popularity, vote_count and vote_average, are correlated to each other. This correlation will be helped in determining the success of the movie.

From the Association analysis, we want to find the data items that have affinity for each other, using the association rules. Depending upon the rules, the data items that have high probability are taken into consideration while planning out a movie.

From the k-means clustering analysis, we need to know the categorized groups of data items from the overview of the movie story. This clustering helps in which kind of movie stories that audience are more interested in and what kind of stories are popular.

## Dataset:

The dataset used in this Term project, is TMDB movies' overview dataset. This dataset is taken from https://www.kaggle.com/writuparnabanerjee/hollywood-movies and was generated from the TMDB Movie Database API. This dataset contains the information about popular movies along with their overviews and the popularity, explained in 7 columns and 10,001 rows.

The columns of the dataset include:

Title: The title of the movie

Overview: The overview of the movie

Original language: The language in which the movie was originally released

Release date: The date on which the movie was released

Popularity: The factor shows how popular the movie is.
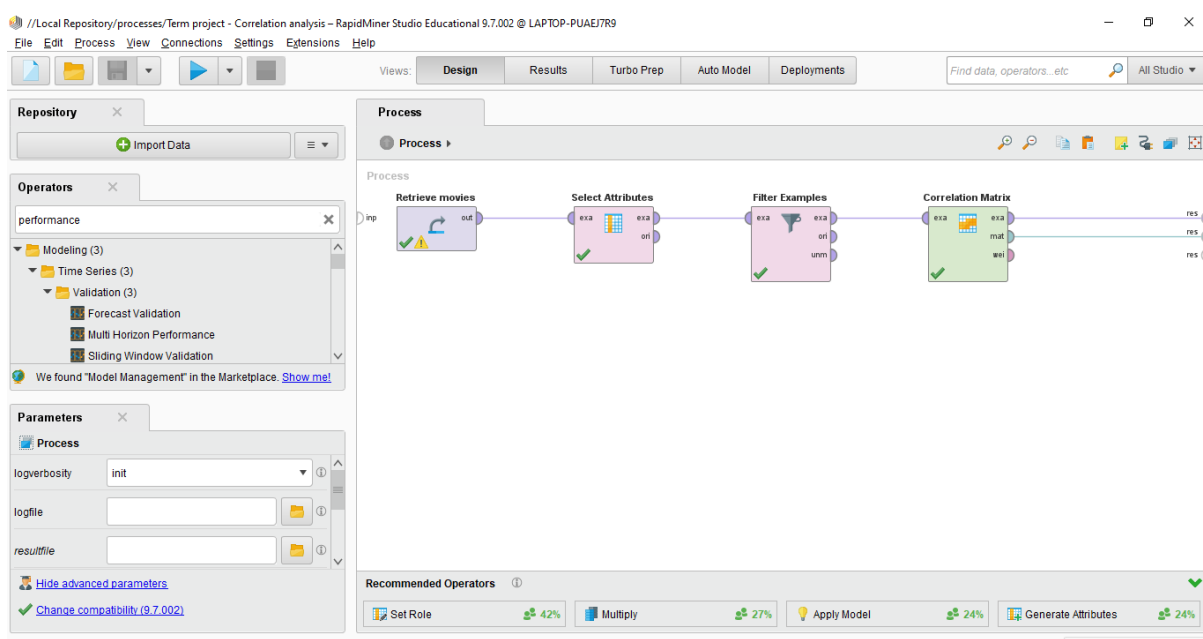
Vote_count: The number of votes for the movie

Vote_average: The average vote

## Data Cleaning:

The dataset used in this project, is cleaned by removing the first column since that column does not have any useful information. This dataset is already free from missing values and special characters so there is no requirement to replace them.

## Datamining techniques:

## <u>Correlation analysis:</u>

The variable Popularity is positively correlated with the variable vote_count with coefficient: 0.418 and is also correlated with the variable vote_average with coefficient: 0.080. The variable vote_count is positively correlated with the variable vote_average like coefficient: 0.228. These correlations are explained through matrix visualization also. It is seen that there is least correlation between popularity and vote_average.

## **Association analysis:**

From the association rules, it is known that only data items 'woman' and 'young' have affinity for each other, with less confidence of 0.321. These rules depend upon the criteria I have given earlier while performing the analysis.

# K-means clustering analysis:

Centroid table gives the results showing the attributes and its cluster values, taken from the output of the performance operator.

From the performance vector, the average within the centroid distance is -0.0987.

Since the 'k' value is given as 5, there are four clusters: cluster_0, cluster_1, cluster_2, cluster_3 and cluster_4.

In Cluster_0, there are 3738 items and the series of first five words are young, life, woman, love, and school. This cluster may be about the movies with teenage love stories.

In Cluster_1, there are 1568 items and the series of first five words are police, agent, murder, detective, and killer. This cluster may be about the Crime thriller movies.

In Cluster_2, there are 1193 items and the series of first five words are film, story, movie, based, and documentary. This cluster may be about the biopic movies of film actors.

In Cluster_3, there are 2356 items and the series of first words are world, earth, group, team, evil, battle, and mission. This cluster may be about fictional character movies.

In Cluster_4, there are 1109 items and the series of first words are family, town, life, young, father, and mother. This cluster may be about the family drama movies.

## Conclusion:

From the TMDB dataset, the popular movies are analysed based on various attributes. With the help of those attributes, the required criteria are calculated.

In this project, three different analysis methods are performed on the dataset to discover various patterns in the data. The predictions are made to know what attributes are correlated and associated. The overview of the movie story will be clustered to identify the discrete groups. However, success cannot be predicted based on a particular attribute. So, these models are based on interesting relations between attributes.

The movie industry persons like the directors, the script writers, the producers so on, can use these results to modify the movie criteria and plan a movie in advance to attain the successful blockbusters.

Not only the movie makers, even the audience can benefit from these models. Before buying a ticket and investing time, the movie watchers can know about the film according to their preferences and determine whether it is good movie, or a flop one.