## PREDICTING IMDB SCORES

## PHASE 2: INNOVATION

Predicting IMDb scores and ratings is a common problem in data science, and there have been several innovations and techniques developed to address this challenge. Here are some key innovations and methods that we have implemented to predict IMDb scores and ratings in our model are :

    i.     Graph Data
    ii.     Natural Language Processing (NPL)
    iii.     Temporal Analysis

### 1. GRAPH DATA

Graph-based approaches for predicting IMDb scores involve using techniques that leverage the relationships between various entities in a movie dataset. Here's a high-level overview of how you can use graph data theory for this purpose:

### i.     Data Representation:

Represent your movie dataset as a graph where nodes represent entities (e.g., movies, actors, directors) and edges represent relationships between them (e.g., actor acted in a movie, director directed a movie).

Assign attributes to nodes and edges, such as movie features (budget, genre), actor information (age, awards), and director information (experience, previous movie ratings).
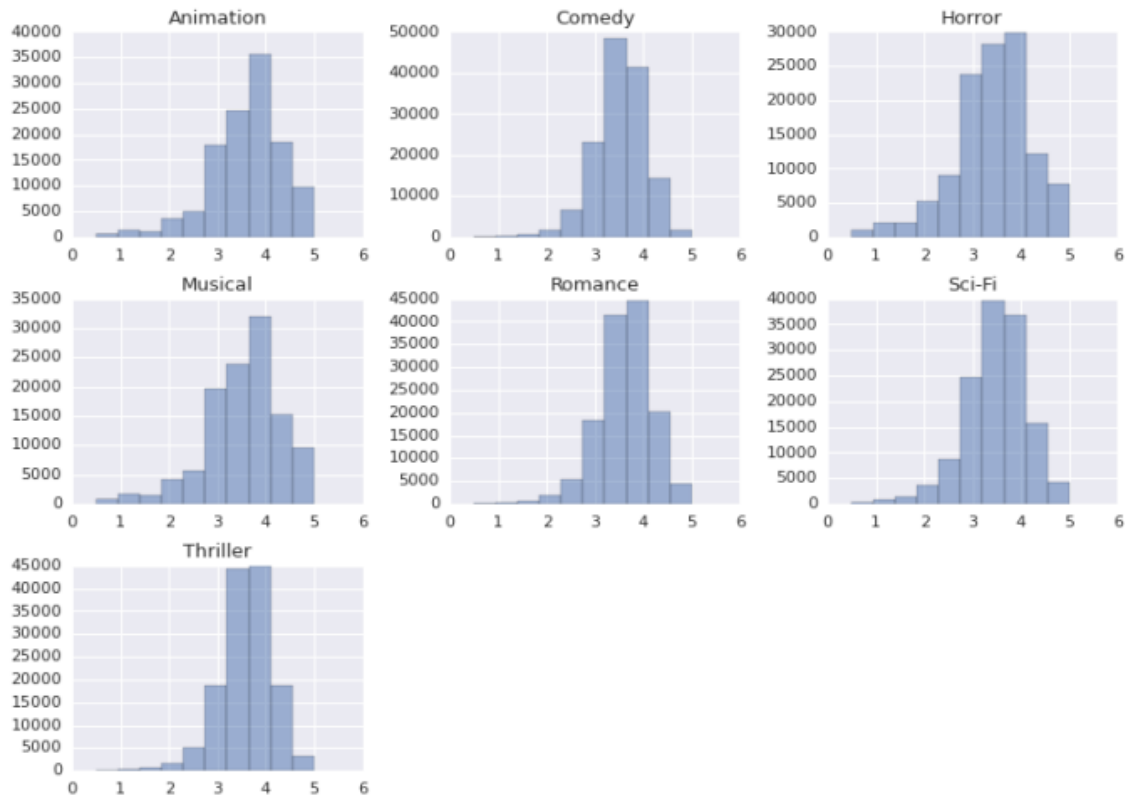
### ii.     Graph Algorithms:

Apply graph algorithms to uncover hidden patterns and relationships. For example, use community detection algorithms to identify groups of movies with similar characteristics.

Utilize link prediction algorithms to predict potential collaborations between directors and actors that might result in high IMDb scores.

### iii.     Visualization:

Visualize the graph structure and important nodes using tools like network visualization libraries (e.g., NetworkX, Gephi) to gain insights into the relationships that affect IMDb scores.

SAMPLE DATA: Based on Movie Genres



## 2. NATURAL LANGUAGE PROCESSING

Natural Language Processing (NLP) can be used in predicting IMDb scores in several ways:

**i. Sentiment Analysis:**

NLP can be used to analyze the sentiment of movie reviews. Positive sentiment in reviews may correlate with higher IMDb scores, while negative sentiment may correlate with lower scores. Machine learning models can be trained to predict IMDb scores based on the sentiment of the reviews.

**ii. Text Features**:

NLP can extract features from movie descriptions, reviews, and other textual data. These features can include word frequencies, TF-IDF scores, or word embeddings. Machine learning models can then use these features to predict IMDb scores.

### iii. Topic Modeling:

 NLP techniques like topic modeling can identify the main topics or themes in movie reviews. If certain topics are associated with higher or lower IMDb scores, this information can be used in prediction models.

### iv. User Reviews:

 Analysing user reviews for specific movies and aggregating user sentiment or opinion can be a valuable source of information for predicting IMDb scores.
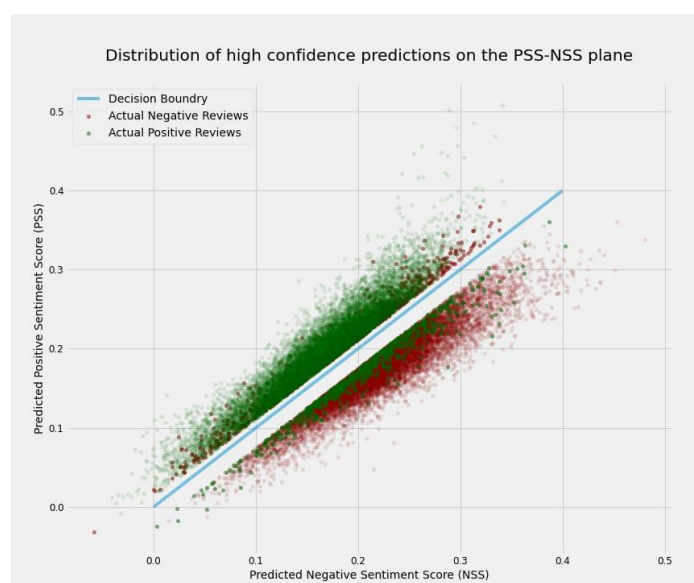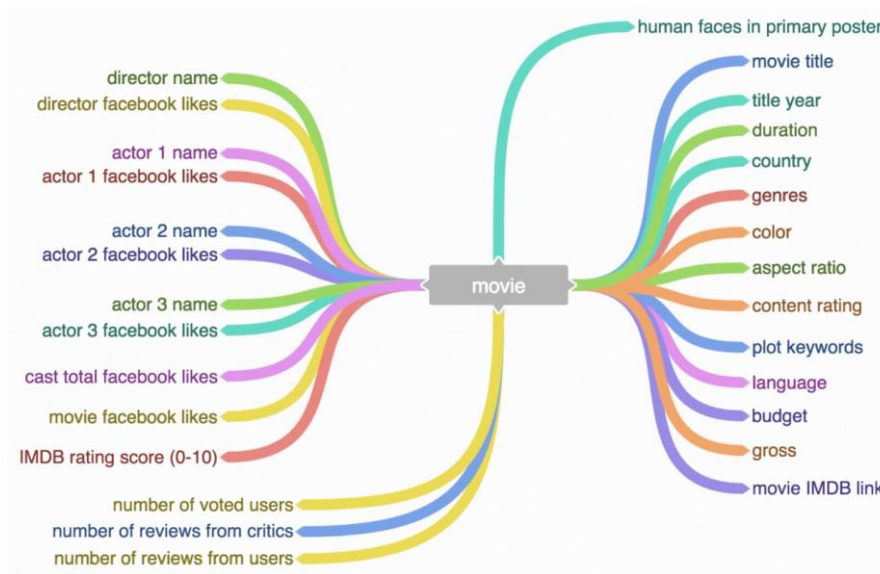
### v. Text Regression Models:

Regression models, such as linear regression or gradient boosting, can be used to predict IMDb scores based on textual features and sentiment scores extracted from reviews and descriptions.

### vi. Natural Language Understanding:

Advanced NLP models like BERT or GPT can be used to understand the context and nuances of movie reviews more comprehensively, potentially improving prediction accuracy.

It's important to note that predicting IMDb scores based solely on NLP features may not be highly accurate, as IMDb scores are influenced by various factors beyond text data, such as director, actors, and marketing. Combining NLP with other data sources and features can improve the accuracy of such predictions. Additionally, it's crucial to have a quality dataset and appropriate evaluation metrics to assess the performance of your prediction model.

### 3. TEMPORAL ANALYSIS

Predicting IMDb scores through temporal analysis involves examining how a movie's attributes and factors change over time and how they correlate with IMDb ratings. Here's a simplified approach:

    **i.**   **Data Collection:**

        Gather a dataset containing information about movies, including release date, genre, director, cast, budget, runtime, and initial IMDb scores.

    **ii.**   **Data Preprocessing:**

        Clean the dataset by handling missing values and outliers.
        Convert categorical variables into numerical format using techniques like one-hot encoding.
        Calculate additional features like the age of the movie at the time of release.

    **iii.**   **Temporal Analysis:**

        Group movies by release year, month, or season to capture temporal trends.
        Compute statistics and trends over time for each group, such as average IMDb scores, budget distribution, genre popularity, etc.

    **iv.**   **Feature Engineering:**

        Create new features that capture temporal dynamics, like the number of movies released in a particular month/year, average IMDb scores

of movies in the same season, or the cumulative budget spent over time.

v.     **Machine Learning Model:**
Select a suitable regression algorithm (e.g., linear regression, random forest, or neural networks).Train the model on the historical data, using features from steps 3 and 4 to predict IMDb scores.

vi.     **Model Evaluation:**
Use metrics like Mean Absolute Error (MAE) or Root Mean Square Error (RMSE) to evaluate the model's performance on a test dataset. Employ time-based cross-validation to assess model stability over different time periods.

vii.     **Feature Importance:**
Analyze feature importance to understand which temporal factors have the most significant impact on IMDb scores.

viii.     **Predictions**:
Apply the trained model to new movies to predict their IMDb scores based on their attributes and the current temporal context.

ix.     **Continuous Monitoring and Updating**:
Periodically retrain the model with new data to account for changing trends and preferences in movie ratings.