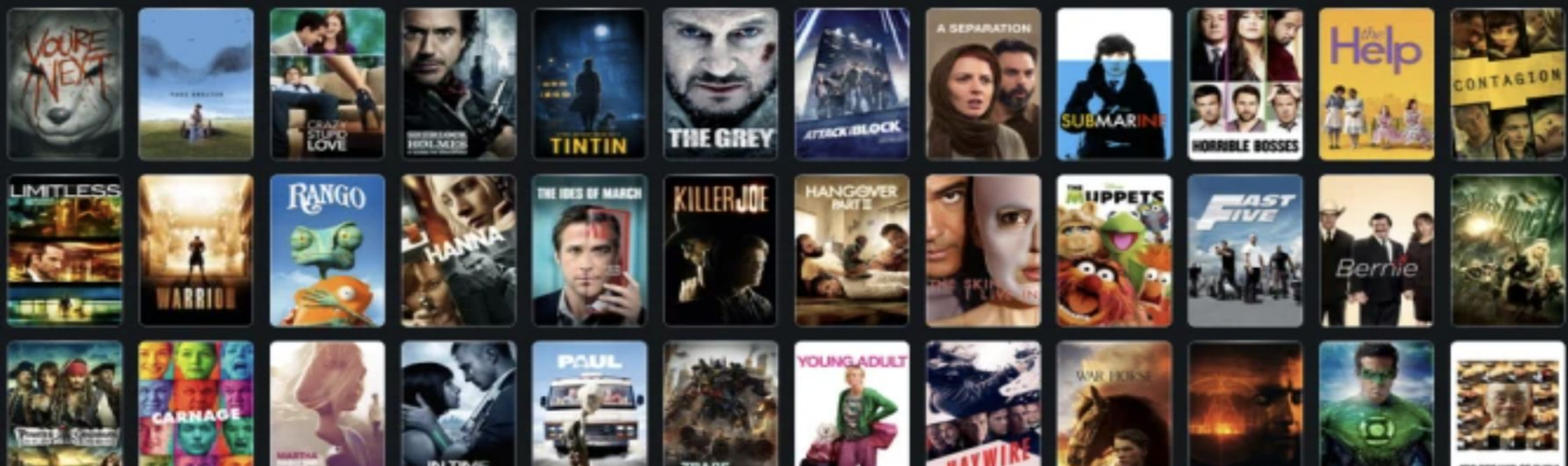# Predicting IMDb Movie Ratings

# Group Members

Adegbenga Ayoola
Grace Le
Michelle Raj
Syed Bari
Wardah Anis

# OBJECTIVE

The purpose of the project is to predict the IMDb Movie Rating from the features contained in the IMDb dataset, execute exploratory Data Analysis and Visualization for the IMDb Dataset
And finally utilise Machine Learning Models to achieve ideal model performance by comparing suitable Machine Learning Algorithms

# Why we chose this topic?

- IMDb, whose full English name is Internet Movie Database, is a web-based data set identified with films, TV shows, home recordings, games. Also, streaming internet-based substance, including full team: entertainers, creation group. What's more, IMDb gives a wide scope of film-related data, be it individual memoirs, plot synopses, tests, or client audits. From that point, you can likewise effectively see the sorts of rankings dependent on a wide range of rules, so new clients can undoubtedly see the issues and content that get the most client consideration.

- We chose this topic because it makes it easy for people to find good movies. Besides, with the current situation, the Covid-19 makes people still afraid to watch movies in cinemas, people often choose the movies to watch at home. Therefore, IMDb makes it easy for people to pick out which movies they love, with high ratings, or their favorite actors.

# ERD Diagram

# Percentage of voting

| Votes | Quantity | Total votes | Percentage |
|-------|----------|-------------|------------|
| 8 - 8.9 | 1686 | 85855 | 1.96% |
| 9 - 9.9 | 55 | 85855 | 0.06% |

# Highest genre of votes from 8 to 10

| Genre | Max Quantiry |
|-------|--------------|
| Drama | 325 |

# Comparing the duration from 8 to 10

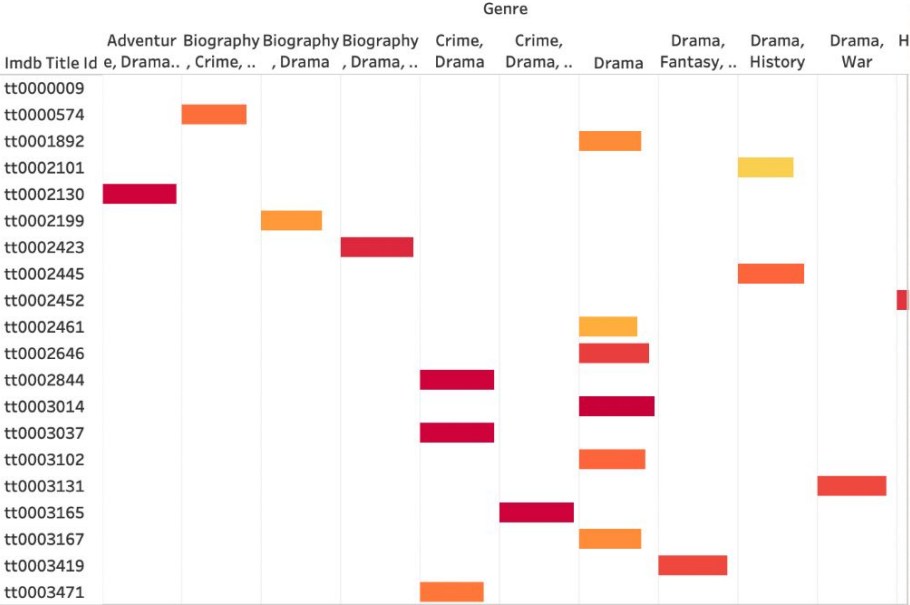| Duration | Total |
|---|---|
| Less than 60 minutes | 19 |
| Greater than 100 minutes | 1281 |

# DataBase AWS S3

S3 is Amazon's cloud file storage service that uses key-value pairs. Files are stored on multiple servers and have a high rate of availability of more than 99.9%. S3 is also very scalable—you are not limited to the memory of one computer. As data flows in, more and more can be stored, as opposed to a local computer that is limited by available memory. Additionally, it offers availability—several team members can access massive amounts of data from one central location
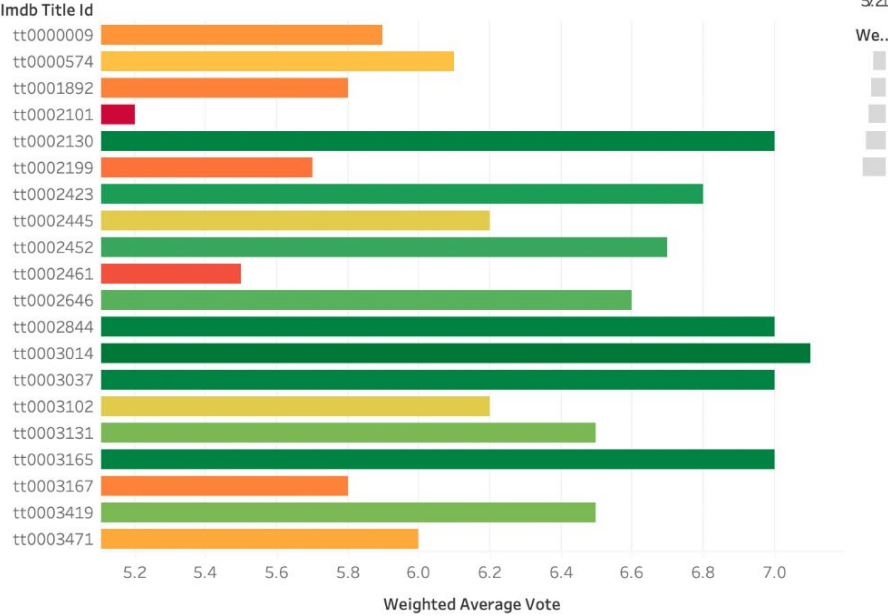
AWS link for the IMDb datasets

1. IMDb_ratings.csv [Link](Link)

2. IMDb_movies.csv [Link](Link)
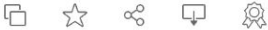
3. Encoded_data.csv [Link](Link)

4. IMDb_data.csv [Link](Link)

# Dashboard

# language vs average rating by Wardah Anis

## Sheet 1

Language 1



| Language 1 | Av Rating |
|---|---|
| Aboriginal | |
| Afrikaans | |
| Akan | |
| Albanian | |
| Algonquin | |
| American Sign Language | |
| Amharic | |
| Arabic | |
| Aramaic | |
| Armenian | |
| Aromanian | |
| Assamese | |
| Aymara | |
| Azerbaijani | |
| Bambara | |
| Basque | |
| Belarusian | |
| Bengali | |
| Berber languages | |
| Bosnian | |
| Brazilian Sign Language | |
| Bulgarian | |
| Burmese | |
| Cantonese | |
| Catalan | |

Av Rating

Country vs Average Rating

# Average Vote by Country



| Count of IM.. | Action | Action, Adventu.. | Action, Adventu.. | Action, Adventu.. | Action, Adventu.. | Action, Adventu.. | Action, Adventu.. | Action, Adventu.. | Action, Adventu.. | Action, Adventu.. | Action, Adventu.. |

**Director**

**Country**

**Count of IMDb movies.csv**

**Country**
- Null
- Afghanistan, France
- Afghanistan, Franc...
- Afghanistan, Iran
- Afghanistan, Irela...
- Albania
- Albania, Austria, F...
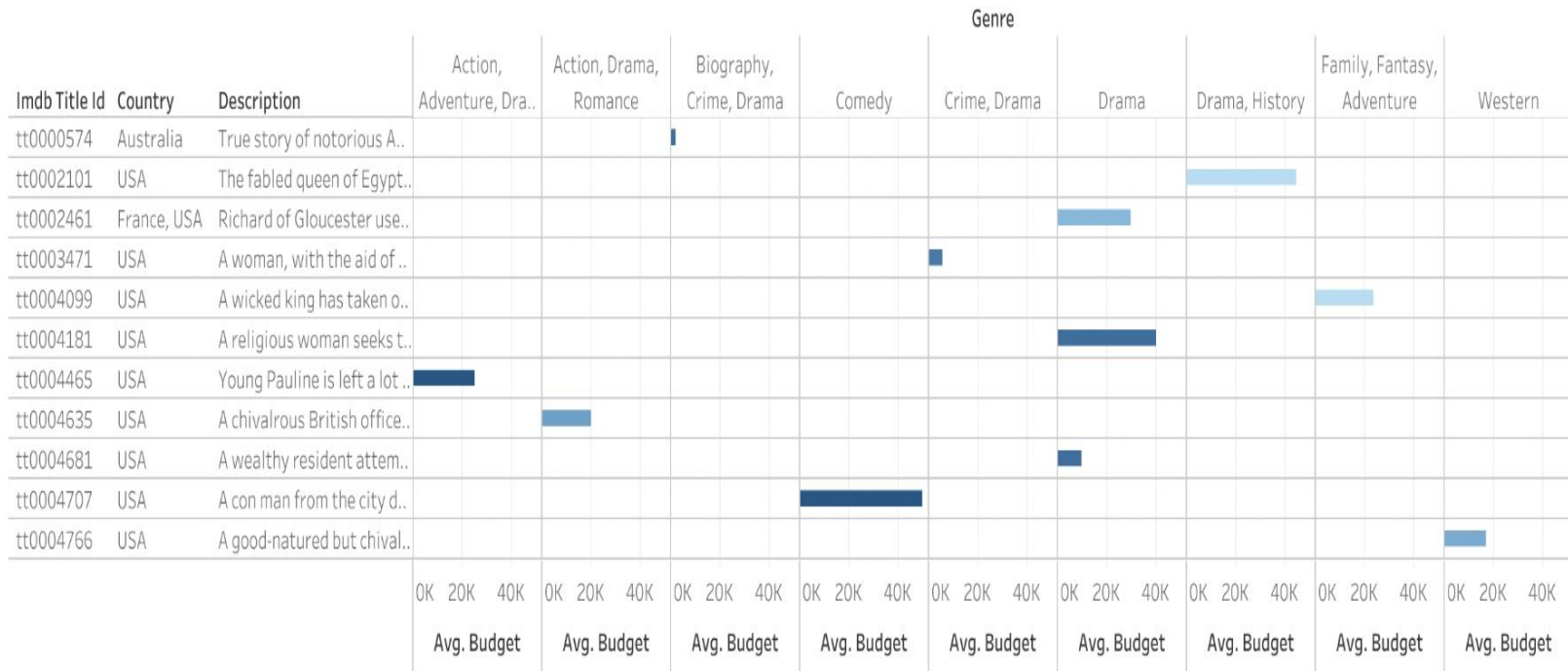- Albania, Czech Rep...
- Albania, Greece, Fr...
- Albania, Italy, Koso...
- Albania, Romania, ...
- Albania, Soviet Uni...
- Albania, UK
- Albania, USA
- Algeria
- Algeria, Belgium, F...
- Algeria, France
- Algeria, France, Ge...
- Algeria, France, Ge...
- Algeria, France, M...
- Algeria, France, Qa...

# Machine Learning Model Proposal



The purpose of the project is to find the best Machine learning Algorithms that can predict the imdb rating based on the features given in the four tables taken from our IMDB Dataset. To achieve this from Various Models the model with the lowest root mean squared error, best accuracy, and best confusion matrix is selected.

# IMDB Movie Dataset Tables
## Datasets selected from Kaggle

### IMDB Movies.csv

| A imdb_title_id | # weighted_averag... | # total_votes | # mean_vote | # median_vote | # votes_10 |
|---|---|---|---|---|---|
| title ID on IMDb | total weighted average rating | total votes received | total mean vote | total median vote | number of votes with rating equal to 10 |
| **85855** unique values | 1 — 9.9 | 99 — 2.28m | 1 — 9.8 | 1 — 10 | 0 — 1.26m |

### IMDB Ratings.csv

| A imdb_title_id | A title | A original_title | # year | A date_published | a genre |
|---|---|---|---|---|---|
| title ID on IMDb | title name | original title name | year of release | date of release | movie genre |
| **85855** unique values | **82094** unique values | **80852** unique values | 1894 — 2020 | **22012** unique values | Drama 15% / Comedy 9% / Other (65619) 76% |

# Machine Learning Models

## Linear Regression, Logistic Regression, Random Forest,  SVM, Deep Learning, Gradient Boosting

**Target:** IMDb Total Average Weighted Rating

**Feature**:  genre, duration, country, language, weighted_average_rating, tot_voters_below_30, tot_voters_below_18, tot_voters_above_45, tot_voters_below_45, tot_male_voters, tot_female_voters
Categorical Feature: title, year, genre, country, language, direction, reviews
Quantitative Feature: Date, Duration, vote, budget, gross income, total votes, us_voters rating,

**Results:** Accuracy, Confusion Matrix, Mean Squared Error and Mean Absolute Error

# Datatype of Inputs and Outputs

## Inputs

```
title                      object
original_title             object
year                       object
date_published             object
genre                      object
duration                    int64
country                    object
language                   object
director                   object
writer                     object
production_company         object
actors                     object
description                object
avg_vote                  float64
votes                       int64
budget                     object
usa_gross_income           object
worlwide_gross_income      object
metascore                 float64
reviews_from_users        float64
reviews_from_critics      float64
dtype: object
```

## Outputs

```
weighted_average_vote     float64
total_votes                 int64
```

# Preliminary Data Preprocessing

- The very first steps include importing the libraries and importing the datasets which are IMDB_movies.csv a total of 22 columns and IMDB_ratings.csv a total of 49 columns

- The IMDB_movies.csv contained the features titel, year, genre, duration, country, language, director, actor, writer, description, average vote, budget, review from users, review from critics

# Description of preliminary feature engineering and selection

- Identify the dependent and the independent variable. After looking through the dataset the features removed were USA gross income, World gross income, metascore due to the null data. Other features were removed 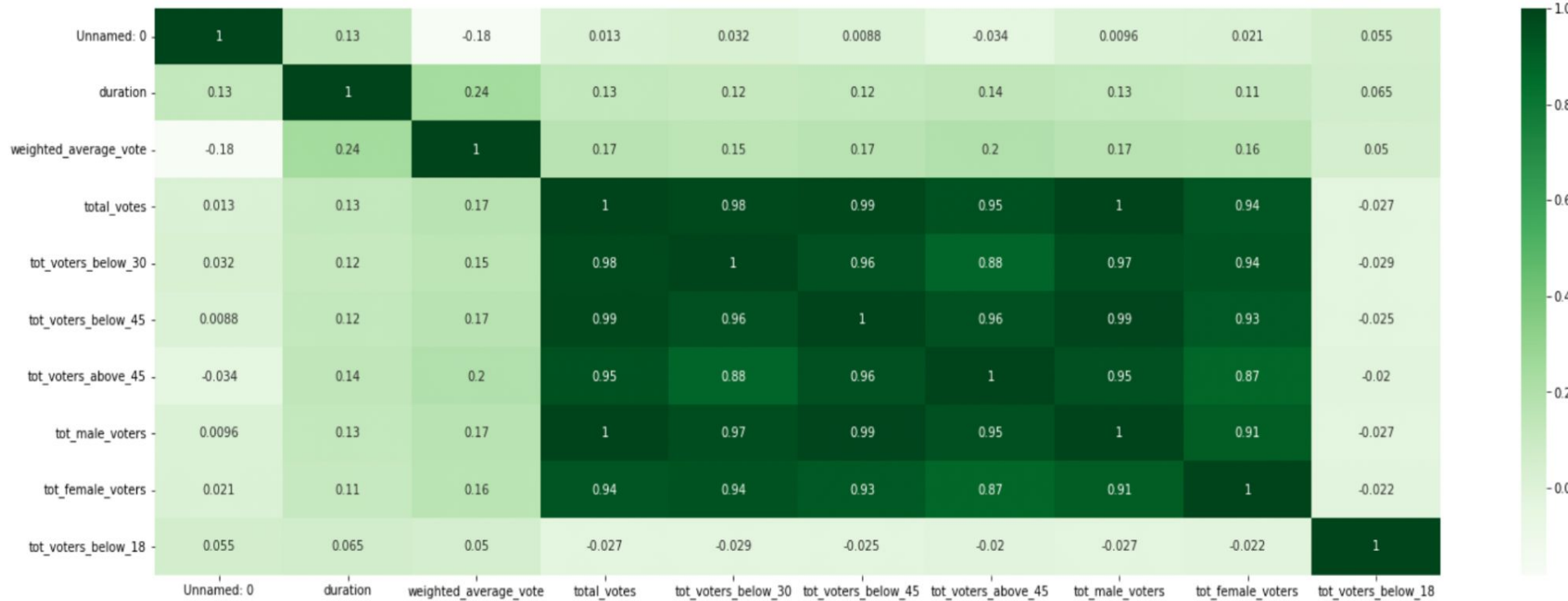due to low correlation to the dependent variable "weighed average rating". The independent variable are as follows [genre, duration, country, language, weighted_average_rating, tot_voters_below_30, tot_voters_below_18, tot_voters_above_45, tot_voters_below_45, tot_male_voters, tot_female_voters]

- The dataset where the cleaned by addressing the null values and dropping hte null rows

- After the cleaning process the categorical data like genre, country and language carrying multiple values were addressed using binary encoding with help of excel sheets

# Correlation Matrix for Cleaned Data

# Correlation matrix of Encoded Data

# 1. Multiple Linear Regression

Linear Regression is a very simple algorithm that can be implemented very easily to give satisfactory results.Furthermore, these models can be trained easily and efficiently even on systems with relatively low computational power when compared to other complex algorithms.Linear regression has a considerably lower time complexity when compared to some of the other machine learning algorithms.The mathematical equations of Linear regression are also fairly easy to understand and interpret

```
r2 socre is  0.24489492727110707
mean_sqrd_error is== 1.2087212257382316
root_mean_squared error of is== 1.0994185853159986
```

The reason for the low score would be due to insufficent data relating to the average rating variable and relevance to the problem although the score has increased from previous linear regression model by 20% more.

**Benefits**

- Models can be trained easily and efficiently even on systems with relatively low computational power when compared to other complex algorithms
- Linear regression fits linearly separable datasets almost perfectly and is often used to find the nature of the relationship between variables
- Overfitting is a situation that arises when a machine learning model fits a dataset very closely and hence captures the noisy data as well

**Limitation**

- A situation that arises when a machine learning model fails to capture the data properly.This typically occurs when the hypothesis function cannot fit the data well.
- **Outliers** of a data set are anomalies or extreme values that deviate from the other data points of the distribution.Data outliers can damage the performance
- Outliers can have a very big impact on linear regression's performance and hence they must be dealt with appropriately before linear regression is applied on the dataset.

# 2. Support Vector Machine

In this project we used Support Vector Machine algorithm as a binary classifier. We divided our target into two groups:

I. Ratings above 7

II. Ratings below 7

We took a random sample of 1500 rows from our dataset to run this model on, since it was taking immense amount of time to run on the complete dataset.

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.90 | 0.93 | 0.92 | 296 |
| 1 | 0.71 | 0.61 | 0.65 | 79 |
| | | | | |
| accuracy | | | 0.86 | 375 |
| macro avg | 0.80 | 0.77 | 0.78 | 375 |
| weighted avg | 0.86 | 0.86 | 0.86 | 375 |

# Benefits

1. SVM works relatively well when there is a clear margin of separation between classes.

2. SVM is more effective in high dimensional spaces.

3. SVM is effective in cases where the number of dimensions is greater than the number of samples.

4. SVM is relatively memory efficient

# Limitation

1. SVM algorithm is not suitable for large data sets.

2. SVM does not perform very well when the data set has more noise i.e. target classes are overlapping.

3. In cases where the number of features for each data point exceeds the number of training data samples, the SVM will underperform.

4. As the support vector classifier works by putting data points, above and below the classifying hyperplane there is no probabilistic explanation for the classification

# 3. Random Forest Regression

Random Forest Regression is a supervised learning algorithm that uses ensemble learning method for regression. Ensemble learning method is a technique that combines predictions from multiple machine learning algorithms to make a more accurate prediction than a single model. A random forest regressor works with data having a numeric or continuous output and they cannot be defined by classes.

Random Forest algorithm

 - Pick at random k data points from the training set.

- Build a decision tree associated to these k data points.

- Choose the number N of trees you want to build and repeat steps 1 and 2.

- For a new data point, make each one of your N-tree trees predict the value of y for the data point in question and assign the new data point to the average across all of the predicted y values.

# Description Data splitting and Modelling

- Dependent and Independent variables were assigned to X(features except weighted average rating) and Y (weighted average rating) variables respectively

- Split the dataset into the Training set and Test set. The training set contains known output from which the model learns, test set then tests the model's predictions based on what it learned from the training set, with the random state=78.

- Fitting the Standard Scaler with the training data. Scaling the data.

- Create a random forest regressor, did model fitting with n_estimator from 200 to 800 with the random state being 1

# Results

## Metrics for Regression Model used

1. Mean Squared Error
2. Root Mean Squared Error
3. Mean Absolute Error

**Output for Estimators** (200, 800, 100)
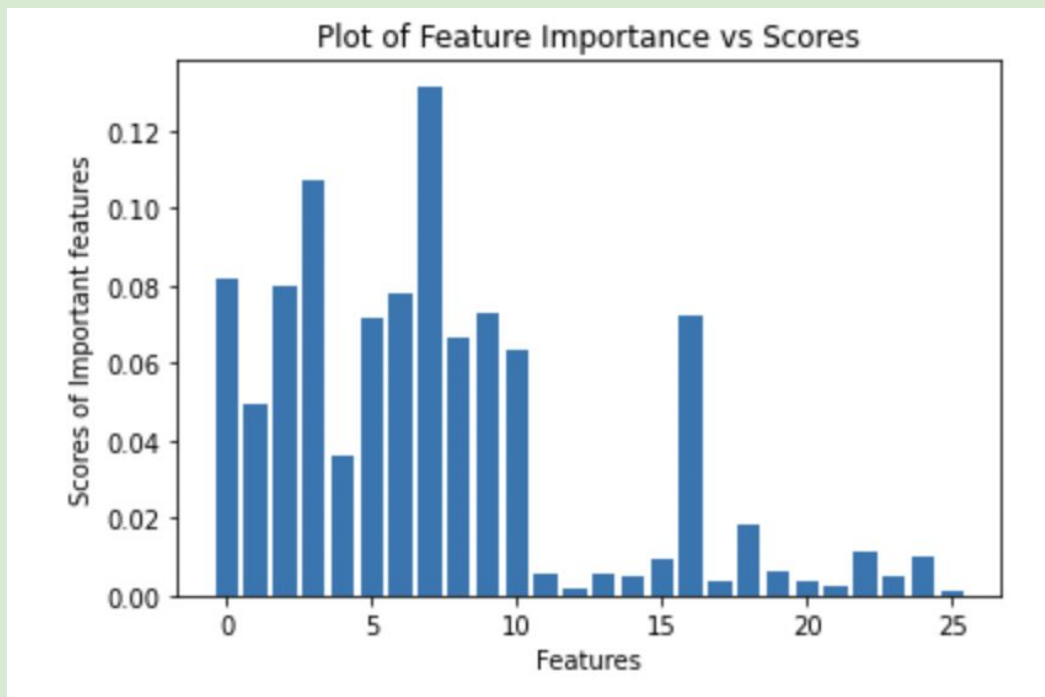
[0.49415757514316994,0.4947008957127814, 0.4958920277984976,

0.4960607807872609,

0.4961116558077189,

0.4958916103852188]

The reason for the low score would be due to insuffient data relating to the average rating variable and relevance to the problem although the score has increased from previous linear regression model by 20% more.



Plot of Feature Importance vs Scores

**Benefits**

- Reduces overfitting in decision trees and improved performance

- Random Forests are not influenced by outliers to a fair degree. It does this by binning the variables

- Automates missing values present in the data

- Normalising of data is not required as it uses a rule-based approach.

**Limitations**

- Random forest is like a black box algorithm, you have very little control over what the model does

- Requires much time for training as it combines a lot of decision trees to determine the class

- Due to ensemble of decision trees, it also suffers interpretability and fails to determine the significance of each variable.

# 4. Gradient Boosting

Gradient boosting is a greedy algorithm and can overfit a training dataset quickly.

It can benefit from regularization methods that penalize various parts of the algorithm and generally improve the performance of the algorithm by reducing overfitting. Boosting is a method of converting weak learners into strong learners. In boosting, each new tree is a fit on a modified version of the original data set.

# Benefits

- Often provides predictive accuracy that cannot be beat.
- Lots of flexibility - can optimize on different loss functions and provides several hyperparameter tuning options that make the function fit very flexible.
- No data pre-processing required - often works great with categorical and numerical values as is.
- Handles missing data - imputation not required.

## Limitation

- GBMs will continue improving to minimize all errors. This can overemphasize outliers and cause overfitting. Must use cross-validation to neutralize.
- Computationally expensive - GBMs often require many trees (>1000) which can be time and memory exhaustive.
- The high flexibility results in many parameters that interact and influence heavily the behavior of the approach

# 5. Deep Learning

Deep learning is a subset of ML, which is essentially a neural network with three or more layers. These neural networks attempt to simulate the behavior of the human brain—albeit far from matching its ability—allowing it to "learn" from large amounts of data. While a neural network with a single layer can still make approximate predictions, additional hidden layers can help to optimize and refine for accuracy.

# Benefits

- Features are automatically deduced and optimally tuned for desired outcome.
- Features are not required to be extracted ahead of time. This avoids time consuming machine learning techniques.
- Robustness to natural variations in the data is automatically learned.
- The same neural network based approach can be applied to many different applications and data types.

# Limitation

- It requires very large amount of data in order to perform better than other techniques.

- It is extremely expensive to train due to complex data models.

- There is no standard theory to guide you in selecting right deep learning tools as it requires knowledge of topology, training method and other parameters

# Summary of Results for Machine Learning Models

1.    **Multiple linear Regression** Best Score : 0.244

2.    **Support Vector Machine** Accuracy : 0.86

3.    **Random Forest Regression**   Best Score : 0.4961116558077189

4.    **Gradient Boosting**  Accuracy : 0.40275730924649394

5.    **Deep Learning** Accuracy : **0.98**

# Summary of Significant Steps

1. Join the four tables to create a single table, analyze the datasets
2. Data Cleaning, and dropping unwanted rows and columns
3. Find the relationship between the feature and the target, find the importance of each feature, drop the columns according to the importance ranking
4. Data splitting into training and testing sets
5. Train and Fit the Machine Learning Model using the processed and cleaned data
6. Calculated the Mean Score Error balanced accuracy score along with the confusion matrix
7. Compare accuracy in different models ; SVM, Random Forest and Neural Networks
8. Add results to database AWS
9. Reports outcomes in Tableau for Visualization.
10. Final Summary

# Thank You