

## PREDICTING IMDb SCORES

## **PHASE 4 – DEVELOPMENT PART 2**



## Introduction

Predicting IMDb scores is a fascinating task that lies at the intersection of data science, machine learning, and the entertainment industry. IMDb, or the Internet Movie Database, is one of the most popular platforms for users to rate and review movies and television shows. Predicting IMDb scores involves developing models that can estimate the average user rating (on a scale of 1 to 10) a movie or TV show might receive based on various features and characteristics. In this introduction, we will provide some background information and a brief data description related to this task.

## Features Engineering :

Feature engineering is a crucial step in building predictive models, including predicting IMDb scores for movies. IMDb scores, or movie ratings, are influenced by various factors, and selecting or creating relevant features can significantly impact the quality of your predictions. Here are some feature engineering ideas for IMDb score prediction:

1. **Genre Information:** Movies often fall into different genres (e.g., action, drama, comedy). Create binary variables for each genre to indicate whether a movie belongs to that genre or not. The combination of genres can influence the IMDb score.
2. **Director's Reputation:** Consider the reputation of the movie's director. You can use data on the director's past movies, their IMDb ratings, and awards won. The director's track record can be an important feature.
3. **Actor's Popularity:** Similar to directors, the popularity and past performances of lead actors can be a good predictor. You can use actors' IMDb ratings, awards, or the number of movies they've starred in.
4. **Runtime:** The length of a movie can affect its rating. Some people prefer shorter movies, while others enjoy longer ones. You can use runtime as a feature.
5. **Release Date:** The release date can be significant. Seasonal effects or holiday releases can impact a movie's success. You can create features like the month or season of release.
6. **Budget:** The budget of a movie can influence its production quality and success. Consider including the movie's budget as a feature.
7. **Marketing and Promotion:** Features related to marketing efforts, such as the number of trailers, marketing budget, or social media engagement, can provide insights into a movie's expected IMDb score.
8. **Critical Reviews:** Use data on critic reviews, Metacritic scores, and Rotten Tomatoes ratings as features. These can be indicative of how well-received a movie is critically.
9. **User Reviews:** Incorporate user-generated content, such as IMDb user reviews, ratings, and the number of user reviews. This can reflect the audience's perception of the movie.
10. **Awards and Nominations:** Include features related to awards won or nominated for, such as Oscars, Golden Globes, or other prestigious awards.
11. **Language and Country:** Consider the language of the movie and the country of origin as features. Different countries have different cinematic traditions and styles.
12. **Word and Phrase Analysis:** Analyze movie descriptions, summaries, or reviews using natural language processing (NLP) techniques to extract keywords, sentiment, or themes.

13. **Box Office Performance:** The box office revenue of a movie can be a strong predictor of its popularity and success.
14. **Sequels and Franchises:** If the movie is part of a popular franchise or a sequel, consider creating a feature to indicate this.
15. **Cultural or Historical Context:** For period films or movies with historical or cultural significance, you can create features related to the time or place in which the story is set.

## Model Training :

Training a machine learning model for IMDb score prediction involves several key steps. Here's a general outline of the process:

### 1. Data Preprocessing:

Data sets = (<https://www.kaggle.com/datasets/luisortner/netflix-original-films-imdb-scores>)

Title	Genre	Premiere	Runtime	IMDb Score	Language
Enter The Anime	Documentary	8/5/19	58	2.5	English/Japanese
Dark Forces	Thriller	8/21/20	81	2.6	Spanish
The App	Science Fiction/Drama	12/26/19	79	2.6	Italian
The Open House	Horror Thriller	1/19/18	94	3.2	English
Kaali Khuhi	Mystery	10/30/20	90	3.4	Hindi
Drive	Action	11/1/19	147	3.5	Hindi
Leyla Everlasting	Comedy	12/4/20	112	3.7	Turkish
The Last Days Of American Crime	Heist Film/Thriller	6/5/20	149	3.7	English
Paradox	Musical/Western/Fantasy	3/23/18	73	3.9	English
Sardar Ka Grandson	Comedy	5/18/21	139	4.1	Hindi
Searching For Sheela	Documentary	4/22/21	58	4.1	English
The Call	Drama	11/27/20	112	4.1	Korean
Whipped	Romantic Comedy	9/18/20	97	4.1	Indonesian
All Because Of You	Action Comedy	10/1/20	101	4.2	Malay
Mercy	Thriller	11/22/16	90	4.2	English
After The Raid	Documentary	12/19/19	25	4.3	Spanish
Ghost Stories	Horror Anthology	January 1, 2020	144	4.3	Hindi
The Last Thing He Wanted	Political Thriller	2/21/20	115	4.3	English
What Happened To Mr. Cha?	Comedy	January 1, 2021	102	4.3	Korean
Death Note	Horror Thriller	8/25/17	100	4.4	English
Hello Privilege. It's Me, Chelsea	Documentary	9/13/19	64	4.4	English
Secret Obsession	Thriller	7/18/19	97	4.4	English
Sextuplets	Comedy	8/16/19	99	4.4	English

## 2. Model Selection:

- Choose an appropriate regression algorithm for IMDb score prediction. Common choices include linear regression, decision trees, random forests, gradient boosting, support vector regression, and neural networks.

## 3. Model Training:

- Use the training data to train your selected machine learning model. During training, the model learns the relationship between the input features (movie attributes) and the target variable (IMDb scores).

## 4. Hyperparameter Tuning:

- Fine-tune the hyperparameters of your model to optimize its performance. This may involve adjusting parameters like learning rate, max depth (for decision trees), number of estimators (for random forests), and so on.

## 5. Model Evaluation:

- Use the testing dataset to assess the performance of your trained model. Common evaluation metrics for regression tasks include:
  - I. **Mean Absolute Error (MAE):** Measures the average absolute difference between predicted and actual IMDb scores.
  - II. **Root Mean Square Error (RMSE):** Provides a measure of the standard deviation of the prediction errors.
  - III. **R-squared ( $R^2$ ):** Indicates the proportion of the variance in IMDb scores that is explained by the model. A higher  $R^2$  value indicates a better fit.

## 6. Cross-Validation:

- Perform k-fold cross-validation to ensure the model's generalization performance. This helps to assess the model's stability and its performance on unseen data.

## 7. Model Interpretability:

- Analyze feature importance to understand which attributes have the most significant impact on IMDb scores. Techniques like SHAP values and feature importance plots can be useful.

## 8. Iteration:

- If the initial model's performance is not satisfactory, iterate through the process by experimenting with different algorithms, feature engineering approaches, and hyperparameter tuning.

#### 9. Model Deployment (Optional):

- If you intend to use the model for IMDb score predictions in practice, create a deployment strategy. This might involve building a user-friendly interface or API for inputting movie attributes and obtaining IMDb score predictions.

#### 10. Monitoring and Maintenance:

- Once deployed, monitor the model's performance in real-world scenarios. Make updates as necessary, especially if IMDb's rating system or user preferences change over time.

It's important to note that IMDb score prediction is a complex task because it depends on a multitude of factors, including subjective user ratings. The quality of your dataset, feature engineering, and model selection play crucial roles in the model's success. Additionally, continuous evaluation and potential retraining are essential to maintain the model's accuracy as new data becomes available.

## Evaluation :

In IMDb score prediction, the evaluation of your predictive model is crucial to determine how well it performs in estimating IMDb scores for movies or TV shows. You can use various evaluation metrics to assess the model's accuracy and reliability. Common evaluation metrics for regression tasks like IMDb score prediction include:

#### 1. Mean Absolute Error (MAE) :

- MAE measures the average absolute difference between the predicted IMDb scores and the actual IMDb scores. It gives you a sense of how far, on average, your predictions are from the true values. A lower MAE indicates better performance.
- Formula:  $MAE = (1/n) * \sum |Actual\ IMDb\ Score - Predicted\ IMDb\ Score|$

#### 2. Root Mean Square Error (RMSE):

- RMSE is similar to MAE but places more weight on large errors. It penalizes larger prediction errors more than smaller ones. Like MAE, a lower RMSE is desirable.
- Formula:  $RMSE = \sqrt{(1/n) * \sum (Actual\ IMDb\ Score - Predicted\ IMDb\ Score)^2}$

#### 3. R-squared (R<sup>2</sup>):

- R-squared is a statistical measure that indicates the proportion of the variance in IMDb scores explained by the model. An R<sup>2</sup> value closer to 1 signifies that the model fits the data well, while an R<sup>2</sup> value close to 0 suggests a poor fit.
- Formula:  $R^2 = 1 - (\text{Sum of Squares Residuals} / \text{Total Sum of Squares})$

- I. Sum of Squares Residuals: The sum of the squared differences between the actual IMDb scores and the predicted IMDb scores.
- II. Total Sum of Squares: The total variance in the actual IMDb scores.

#### 4. Mean Absolute Percentage Error (MAPE):

- MAPE measures the percentage difference between the predicted IMDb scores and the actual IMDb scores. It is useful when you want to understand the relative error.
- Formula:  $MAPE = (1/n) * \sum (|Actual\ IMDb\ Score - Predicted\ IMDb\ Score| / |Actual\ IMDb\ Score|) * 100$

#### 5. Coefficient of Determination (CD):

- CD provides an alternative way to assess the model's goodness of fit. A CD value closer to 1 indicates a better model fit, while a value close to 0 suggests a poor fit.
- Formula:  $CD = 1 - (1 - R^2) * ((n - 1) / (n - p - 1))$ 
  - I. n: Number of observations.
  - II. p: Number of predictor variables.

#### 6. Adjusted R-squared (Adj. R<sup>2</sup>):

- Adjusted R-squared adjusts the R-squared value based on the number of predictor variables in the model. It helps account for the complexity of the model and penalizes overfitting.
- Formula:  $Adj.\ R^2 = 1 - (1 - R^2) * ((n - 1) / (n - p - 1))$ 
  - I. n: Number of observations.
  - II. p: Number of predictor variables.

### Importants Of Feature Engineering, Model Training, Evaluation In IMDb Score Prediction :

In the context of IMDb score prediction, feature engineering, model training, and evaluation are crucial steps to build an accurate and effective movie rating prediction system. Here's why each of these components is important:

#### 1. Feature Engineering: Feature engineering involves selecting, transforming, and creating input variables (features) for your predictive model. In the context of IMDb score prediction, feature engineering can include:

- Movie-specific features: These may include information about the cast, crew, genre, release date, budget, and more.
- User-specific features: Information about the user's historical ratings, preferences, and viewing history.
- Text-based features: Sentiment analysis of user reviews or movie descriptions.

The importance of feature engineering lies in extracting meaningful information from the available data to improve the model's ability to make accurate predictions. Effective feature

engineering can lead to better model performance and a more nuanced understanding of the factors that influence IMDb scores.

**2. Model Training:** Model training refers to the process of training a machine learning or statistical model on a labeled dataset. In IMDb score prediction, the model is trained to learn the relationships between the selected features and IMDb scores. Key points to consider during model training include:

- **Model selection:** Choosing an appropriate machine learning algorithm or model architecture for the task.
- **Hyperparameter tuning:** Adjusting hyperparameters to optimize model performance.
- **Cross-validation:** Splitting the dataset into training and validation sets to assess the model's generalization ability.

The model's ability to capture patterns and relationships in the data is essential for making accurate predictions of IMDb scores. Effective training ensures that the model can generalize well to unseen data.

**3. Evaluation:** Model evaluation is critical for assessing the performance and reliability of your IMDb score prediction model. Common evaluation metrics for regression tasks like IMDb score prediction include Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and R-squared (R<sup>2</sup>). The importance of model evaluation lies in:

- **Identifying model performance:** Evaluation metrics quantify how well the model is predicting IMDb scores.
- **Comparing models:** You can use evaluation results to compare different models and choose the best one.
- **Identifying overfitting:** Evaluation helps in determining whether the model is overfitting the training data.

An accurate evaluation process ensures that your IMDb score prediction system provides reliable and consistent results, which are crucial for user satisfaction and decision-making.

## **Challenges Involved In Feature Engineering, Model Training, Evaluation in IMDb Score Prediction :**

Feature engineering, model training, and evaluation in IMDb score prediction present several challenges that need to be addressed for a successful movie rating prediction system. Here are some of the challenges associated with each of these components:

### **Feature Engineering Challenges:**

- a. Data quality:** The quality and completeness of the data can be a significant challenge. Incomplete or noisy data can lead to inaccurate feature engineering.
- b. Data volume:** The dataset may be limited in terms of the number of movies and user ratings, which can make it challenging to extract meaningful features.



c. Feature selection: Determining which features to include and which to exclude can be challenging. Including too many irrelevant features may lead to overfitting, while excluding important features can reduce prediction accuracy.

d. Handling categorical data: Encoding and dealing with categorical features such as movie genres, directors, or actors can be complex, and the choice of encoding method can impact model performance.

e. Temporal data: Handling time-related features, such as release date or user rating history, presents challenges in capturing temporal patterns and trends effectively.

### **Model Training Challenges :**

a. Model selection: Choosing the most appropriate machine learning algorithm or deep learning architecture can be challenging. The optimal model can vary depending on the dataset and the complexity of the problem.

b. Hyperparameter tuning: Finding the right hyperparameters for the chosen model is a trial-and-error process that can be time-consuming and computationally expensive.

c. Imbalanced data: In IMDb score prediction, the distribution of movie ratings may be skewed, leading to imbalanced datasets. Dealing with this imbalance is important to avoid model bias.

d. Overfitting: Preventing overfitting is crucial. Ensuring that the model generalizes well to unseen data can be a challenge, particularly when dealing with limited data.

### **Evaluation Challenges :**

a. Evaluation metric choice: Selecting the most appropriate evaluation metric for IMDb score prediction can be challenging. Different metrics may emphasize different aspects of prediction quality.

b. User-specific evaluation: Personalization is a key aspect of IMDb score prediction, and evaluating how well the model caters to individual user preferences can be complex.

c. Model interpretability: Understanding why a model makes certain predictions is important, especially for user trust. Many machine learning models, such as deep neural networks, lack transparency.

d. Handling outliers: Extreme outlier ratings may exist in the dataset, and these can disproportionately affect evaluation metrics. Deciding how to handle outliers can be challenging.



## Conclusion :

Certainly, let's provide a concise conclusion for each stage of IMDb score prediction:

### 1. Feature Engineering:

- Feature engineering is a crucial step in IMDb score prediction, involving the selection and transformation of relevant features.
- Key features may include movie attributes like genre, director, actors, release date, and textual data like movie descriptions and reviews.
- Domain knowledge and creativity in feature selection can significantly impact model performance.

### 2. Model Training:

- Model training is the process of selecting and training a suitable machine learning or deep learning model.
- Common choices include regression models, decision trees, random forests, gradient boosting, and neural networks.
- Proper hyperparameter tuning and cross-validation are essential for optimizing model performance.

### 3. Evaluation:

- Evaluation assesses the model's performance and generalization capabilities.
- Common regression metrics like MAE, MSE, and  $R^2$  are used to measure accuracy.
- Testing on a separate dataset and cross-validation help determine how well the model generalizes.
- Model interpretation and error analysis can provide insights for model improvement.