

# **PREDICTING IMDB SCORES**

## **TEAM MEMBERS**

**1.HARISH G**

**2.HARI PRASATH S**

**3.SANGAMESHWAR**

**4.SAKTHEESWARAN C**

**5.ROHITH VARSHAN**

## **PHASE 3 SUBMISSION DOCUMENT**

## **PHASE-3 DEVELOPMENT PART 1**



# INTRODUCTION

## IMDb Score Prediction Project:

An IMDb score prediction project is a data-driven initiative that focuses on forecasting the rating a movie or TV show is likely to achieve on IMDb, a prominent platform for film and television ratings and reviews.

## Data Analysis and Machine Learning:

This project involves extensive data analysis, encompassing factors like the cast, crew, genre, budget, and other relevant information from a comprehensive dataset of films and TV shows. Machine learning techniques, such as regression models or deep learning algorithms, are employed to make predictions based on these key features.

## Valuable Insights for the Entertainment Industry:

The primary objective of this project is to develop a reliable model capable of accurately estimating IMDb scores. These predictions are invaluable to filmmakers, production studios, and audiences seeking insights into the anticipated quality of a particular production before its release, aiding in informed decision-making and viewing choices.

Data set link: ( <https://www.kaggle.com/datasets/luisortier/netflix-original-films-imdb-scores>)

**Step 1: Data Preprocessing:** Data preprocessing is the initial step in any data science project. In this context, it involves loading the dataset and performing data cleaning. This includes addressing missing data, dealing with inconsistencies, and ensuring that the dataset is ready for analysis and modelling.

**Step 2: Exploratory Data Analysis (EDA):** Exploratory Data Analysis is a pivotal phase for understanding the dataset. Through data visualization and summary statistics, you gain insights into IMDb score distributions, relationships between features like genre, runtime, and IMDb scores, and potential outliers that may require special attention during modelling.

**Step 3: Feature Selection and Engineering:** Determining which features to include in your model is a crucial decision. In this step, you'll choose the most relevant attributes for IMDb score prediction. Additionally, you may opt to engineer new

features that can enhance prediction accuracy, such as deriving release month or conducting textual analysis on titles and descriptions.

**Step 4: Data Splitting:** To train and evaluate your model effectively, you'll divide the dataset into training and testing subsets. The training set will be used to teach your IMDb score prediction model, while the testing set helps evaluate its performance and generalizability.

**Step 5: Model Selection:** Selecting an appropriate machine learning algorithm is pivotal to your project's success. Choices range from linear regression to more complex methods like decision trees, random forests, gradient boosting, or neural networks. The selection should align with your specific project goals and dataset characteristics.

**Step 6: Model Training:** In this phase, your chosen model is trained on the training dataset. The process involves fine-tuning hyperparameters to optimize performance. Repeated iterations may be required to achieve the best results.

**Step 7: Model Evaluation:** The performance of your IMDb score prediction model is assessed using the testing dataset. Evaluation metrics like mean squared error (MSE) and R-squared (R<sup>2</sup>) are common for regression tasks. Interpret the results to understand which features significantly impact the prediction and unveil any patterns discovered by the model.

**Step 8: Model Interpretation:** Model interpretation is key to understanding the inner workings of your IMDb score prediction model. Gain insights into which features carry the most weight in predicting IMDb scores and discover patterns that the model has identified. This is essential for making informed decisions and gaining a deeper understanding of the data.

**Step 9: Deployment (Optional):** If desired, you can take your model to the next level by deploying it, creating a tool for predicting IMDb scores for new Netflix content. This step requires expertise in deployment platforms and technologies.

**Step 10: Iterate and Improve:** Machine learning models are rarely perfect from the start. Continuously refine your model based on feedback, new data, and advances in machine learning techniques. Experiment with different algorithms and features to improve prediction accuracy over time.

**Step 11: Documentation and Reporting:** The final step is to create a comprehensive report that summarizes your findings, methodology, and the model's performance. Sharing insights and results with stakeholders or the wider community is crucial for knowledge dissemination and project impact.

## NECESSARY STEPS TO FOLLOW:

To start data analysis for IMDb score prediction, you'll need to follow several necessary steps. Below, I outline the general steps, along with sample Python code for each step. We'll use a hypothetical dataset for demonstration purposes. In practice, you'd have your dataset, which you can load using appropriate libraries or from a file.

### **1. Import Libraries:**

Start by importing the necessary libraries for data analysis and machine learning.

```
```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
```
```

### **2. Load the Dataset:**

Load your IMDb dataset. You can use the `pandas` library for this task.

```
```python
```

```
data = pd.read_csv('imdb_data.csv') # Replace 'imdb_data.csv' with your dataset file
```

```
...
```

### 3. Explore the Data:

It's essential to understand your data. Explore it to get a sense of its structure and the types of features you have.

```
```python
```

```
print(data.head()) # Display the first few rows of the dataset
```

```
print(data.info()) # Get information about the dataset (data types, missing values, etc.)
```

```
print(data.describe()) # Summary statistics
```

```
...
```

### 4. Data Preprocessing:

Prepare your data for analysis. This includes handling missing values, encoding categorical features, and scaling numeric features if needed.

```
```python
```

```
# Handle missing values
```

```
data.dropna(inplace=True)
```

```
# Encode categorical variables if any
```

```
data = pd.get_dummies(data, columns=['categorical_column'])
```

```
# Scaling numeric features
```

```
scaler = StandardScaler()
data['numeric_column'] = scaler.fit_transform(data[['numeric_column']])
...
```

## 5. Feature Selection:

Select the relevant features for your analysis. Choose features that are likely to influence IMDb scores.

```
```python
# Select relevant features
X = data[['feature1', 'feature2', 'feature3']]
y = data['IMDb_score']
...
```

## 6. Train-Test Split:

Split the dataset into training and testing sets. The typical split is 80% for training and 20% for testing.

```
```python
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
...
```

## 7. Model Building:

Build a predictive model. In this example, we'll use a simple linear regression model.

```
```python
model = LinearRegression()
model.fit(X_train, y_train)
...

```

## 8. Model Evaluation:

Evaluate the model's performance on the test dataset.

```
```python
y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print(f"Mean Squared Error: {mse}")
print(f"R-squared: {r2}")
...

```

## 9. Visualization:

Visualize the model's performance or explore the data using libraries like `matplotlib` and `seaborn`.

```
```python
plt.scatter(y_test, y_pred)
plt.xlabel('Actual IMDb Scores')
plt.ylabel('Predicted IMDb Scores')
plt.title('Actual vs. Predicted IMDb Scores')
plt.show()
...

```

## 10. Predict IMDb Scores:

Once you are satisfied with your model's performance, you can use it to predict IMDb scores for new data.

```
```python
new_data = pd.DataFrame({'feature1': [value1], 'feature2': [value2], 'feature3':
[value3]})
predicted_score = model.predict(new_data)
```
```

Certainly, here's a sample conclusion for your IMDb score prediction data analysis:

## Conclusion

In this data analysis guide, we have walked through the essential steps for predicting IMDb scores using a hypothetical dataset. The journey can be summarized as follows:

1. Importing Libraries: We started by importing the necessary Python libraries, including pandas for data manipulation, scikit-learn for machine learning, and matplotlib and seaborn for data visualization.
2. Loading the Dataset: We loaded the IMDb dataset (replace 'imdb\_data.csv' with your dataset file) to begin our analysis.



3. Exploring the Data: Understanding our dataset is crucial. We inspected the first few rows, gathered information about the data types and missing values, and computed summary statistics.

4. Data Preprocessing: We handled missing values, encoded categorical variables if any, and scaled numeric features to prepare the data for modeling.

5. Feature Selection: Selecting relevant features that could impact IMDb scores is key to the prediction process.

6. Train-Test Split: We divided the dataset into training and testing sets to evaluate our model's performance.

7. Model Building: We constructed a linear regression model, which serves as a foundation for IMDb score prediction. In practice, more advanced models can be explored based on the dataset characteristics and goals.

8. Model Evaluation: We assessed the model's performance using metrics like Mean Squared Error (MSE) and R-squared ( $R^2$ ) on the test dataset.

9. Data Visualization: Effective data visualization plays a vital role in understanding the relationships within the data and interpreting model results. We demonstrated a simple scatter plot of actual vs. predicted IMDb scores as an example.

Throughout this guide, it's important to remember that IMDb score prediction is a simplified example, and real-world scenarios may involve more intricate data preprocessing, feature engineering, and model selection. This guide provides a starting point for your data analysis journey.

For a comprehensive IMDb score prediction system, further exploration, tuning, and experimentation with various models and data enhancements may be required. Successful IMDb score prediction can provide valuable insights for content creators, filmmakers, and the entertainment industry to produce content that resonates with audiences.

---