

IMDb SCORE PREDICTION USING DATASCIENCE

PROJECT TITLE : IMDb Score Prediction

PROBLEM STATEMENT : Develop a machine learning model to predict the IMDb scores of movies available on Films based on their Genre, Premiere Data, Runtime and Language. The model aims to accurately estimate the popularity of movies to assist users in discovering highly rated films that align with their preferences.

DESIGN THINKING PROCESS:

1. Empathize:

Understand the IMDb Score Prediction Challenge:

What are the goals and objectives of this project?

Why is it important to predict IMDb scores for movies?

User Research:

Conduct surveys or interviews with potential users to understand their needs and expectations. Gather feedback from movie enthusiasts and IMDb users.

2. Define

Problem Statement:

Create a clear problem statement, such as "Design a data science model to predict IMDb scores for movies accurately."

User Needs and Requirements:

Document the specific requirements and features users expect in the IMDb score prediction system.

3. Ideate:

Brainstorm Solutions:

Generate ideas for the data science approach to predict IMDb scores. Consider various algorithms and data sources.

Consider Ethical Implications:

Address potential biases in the data and algorithms, and explore ways to ensure fairness and transparency.

4. Prototype:

Data Collection:

Collect relevant data, such as movie features (genre, director, actors, budget, release date) and IMDb historical scores.

Data Preprocessing:

Clean and prepare the data, handling missing values, and encoding categorical variables.

Model Selection:

Experiment with different machine learning algorithms (e.g., regression, random forests, neural networks).

Feature Engineering:

Create and test new features to enhance model performance.

5. Test:

Model Evaluation:

Assess the model's performance using metrics like Mean Absolute Error, Root Mean Square Error, and R-squared.

Implement cross-validation to avoid overfitting.

User Feedback:

Solicit feedback from users and stakeholders on the model's accuracy and usability.

Use feedback to refine the model.

6. Implement:

Deployment:

Develop a user-friendly platform or application where users can input movie details to get IMDb score predictions.

Ethical Considerations:

Ensure that the model follows ethical guidelines, addressing fairness and privacy concerns.

7. Iterate:

Continuous Improvement:

Monitor the model's performance and collect more data to enhance accuracy.

Adapt to changing user needs and preferences.

Scaling:

Plan for scaling the system to accommodate larger volumes of movie data and user requests.

DATASET:

DATASETLINK :

<https://www.kaggle.com/datasets/luisortega/netflix-original-films-imdb-scores>

ABOUT DATASET:

Context:

This dataset consists of all Netflix original films released as of June 1st, 2021. Additionally, it also includes all Netflix documentaries and specials. The data was webscraped off of this Wikipedia page, which was then integrated with a dataset consisting of all of their corresponding IMDB scores. IMDB scores are voted on by community members, and the majority of the films have 1,000+ reviews.

Content

Included in the dataset is:

- Title of the film
- Genre of the film
- Original premiere date
- Runtime in minutes
- IMDB scores (as of 06/01/21)
- Languages currently available (as of 06/01/21)

DATA PREPROCESSING:

Loading the Dataset:

Here, I'm going to load the dataset in google colab Before that importing libraries is mandatory,

Importing Libraries:

✓
0s



```
import pandas as pd
import numpy as np

#Data Visulation
import seaborn as sns
import matplotlib.pyplot as plt
import plotly.express as px
import plotly.graph_objects as go

from datetime import datetime

import statsmodels.api as sm

from warnings import filterwarnings
filterwarnings("ignore")
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
```

After importing libraries, load the dataset here how it looks like,

```
[ ] from google.colab import files
    upload=files.upload()
```

Choose Files NetflixOriginals.csv

- **NetflixOriginals.csv**(text/csv) - 38678 bytes, last modified: 10/11/2023 - 100% done
Saving NetflixOriginals.csv to NetflixOriginals (1).csv

Reading the Dataset:

```

▶ movie=pd.read_excel("NetflixOriginals.xlsx",encoding = "ISO-8859-1")
print(movie)
dataDate=movie.copy()

```

OUTPUT:

	Title	Genre \
0	Enter the Anime	Documentary
1	Dark Forces	Thriller
2	The App	Science fiction/Drama
3	The Open House	Horror thriller
4	Kaali Khuhi	Mystery
..
579	Taylor Swift: Reputation Stadium Tour	Concert Film
580	Winter on Fire: Ukraine's Fight for Freedom	Documentary
581	Springsteen on Broadway	One-man show
582	Emicida: AmarElo - It's All For Yesterday	Documentary
583	David Attenborough: A Life on Our Planet	Documentary

Information About the Dataset:

The info() method provides a concise summary of the data in a pandas Data Frame. It typically includes the following information:

The number of non-null (non-missing) values in each column.

1. The data type of each column (e.g., int64, float64, object, datetime64, etc.).
2. The total memory usage of the Data Frame.
3. Additional information, such as the count, mean, standard deviation, minimum, and maximum values for numeric columns.

```
[ ] movie.info()
```

OUTPUT:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 584 entries, 0 to 583
Data columns (total 6 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Title           584 non-null    object
1   Genre           584 non-null    object
2   Premiere        584 non-null    datetime64[ns]
3   Runtime         584 non-null    int64
4   IMDB Score      584 non-null    float64
5   Language        584 non-null    object
dtypes: datetime64[ns](1), float64(1), int64(1), object(3)
memory usage: 27.5+ KB
```

Head() Function:

It serves the purpose of displaying the first few rows of the dataset.

```
[ ] movie.head()
```

OUTPUT:

	Title	Genre	Premiere	Runtime	IMDB Score	Language
0	Enter the Anime	Documentary	2019-08-05	58	2.5	English/Japanese
1	Dark Forces	Thriller	2020-08-21	81	2.6	Spanish
2	The App	Science fiction/Drama	2019-12-26	79	2.6	Italian
3	The Open House	Horror thriller	2018-01-19	94	3.2	English
4	Kaali Khuhi	Mystery	2020-10-30	90	3.4	Hindi

Describe() Function:

This is a Pandas Data Frame method that generates descriptive statistics of a dataset, typically for numeric columns.

```
▶ movie.describe().T
```

OUTPUT:

	count	mean	std	min	25%	50%	75%	max
Runtime	584.0	93.577055	27.761683	4.0	86.0	97.00	108.0	209.0
IMDB Score	584.0	6.271747	0.979256	2.5	5.7	6.35	7.0	9.0

Checking For Missing Values:

The expression `movie.isnull().values.any()` is used to check whether there are any missing values in a Pandas DataFrame or a NumPy array. It returns a boolean value, indicating whether any missing values are present in the dataset.

```
[ ] movie.isnull().values.any()
```

OUTPUT:

```
False
```

VARIOUS ANALYSIS PERFORMED USING DATASET:

Descriptive Analysis:

It involves the process of summarizing and presenting data in a meaningful and informative way. The primary purpose of descriptive analysis is to provide a clear and concise overview of a dataset.

```
▶ # Display summary statistics of numerical columns
print(movie.describe())

# Calculate the mean IMDB score
mean_imdb_score = movie['IMDB Score'].mean()
print(f"Mean IMDB Score: {mean_imdb_score}")

# Count the number of movies in each content rating category
content_rating_counts = movie['Runtime'].value_counts()
print(content_rating_counts)
```

OUTPUT:

```
➡
```

	Runtime	IMDB Score
count	584.000000	584.000000
mean	93.577055	6.271747
std	27.761683	0.979256
min	4.000000	2.500000
25%	86.000000	5.700000
50%	97.000000	6.350000
75%	108.000000	7.000000
max	209.000000	9.000000

Mean IMDB Score: 6.2717465753424655

97	24
98	19
94	19
95	18
100	17
..	
148	1
147	1
7	1
57	1
153	1

Name: Runtime, Length: 124, dtype: int64

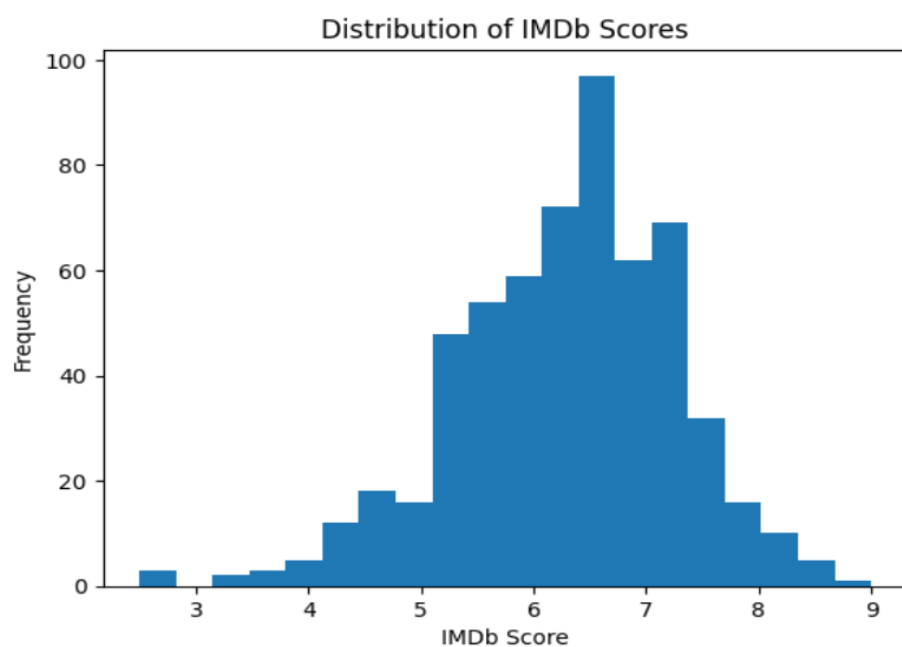
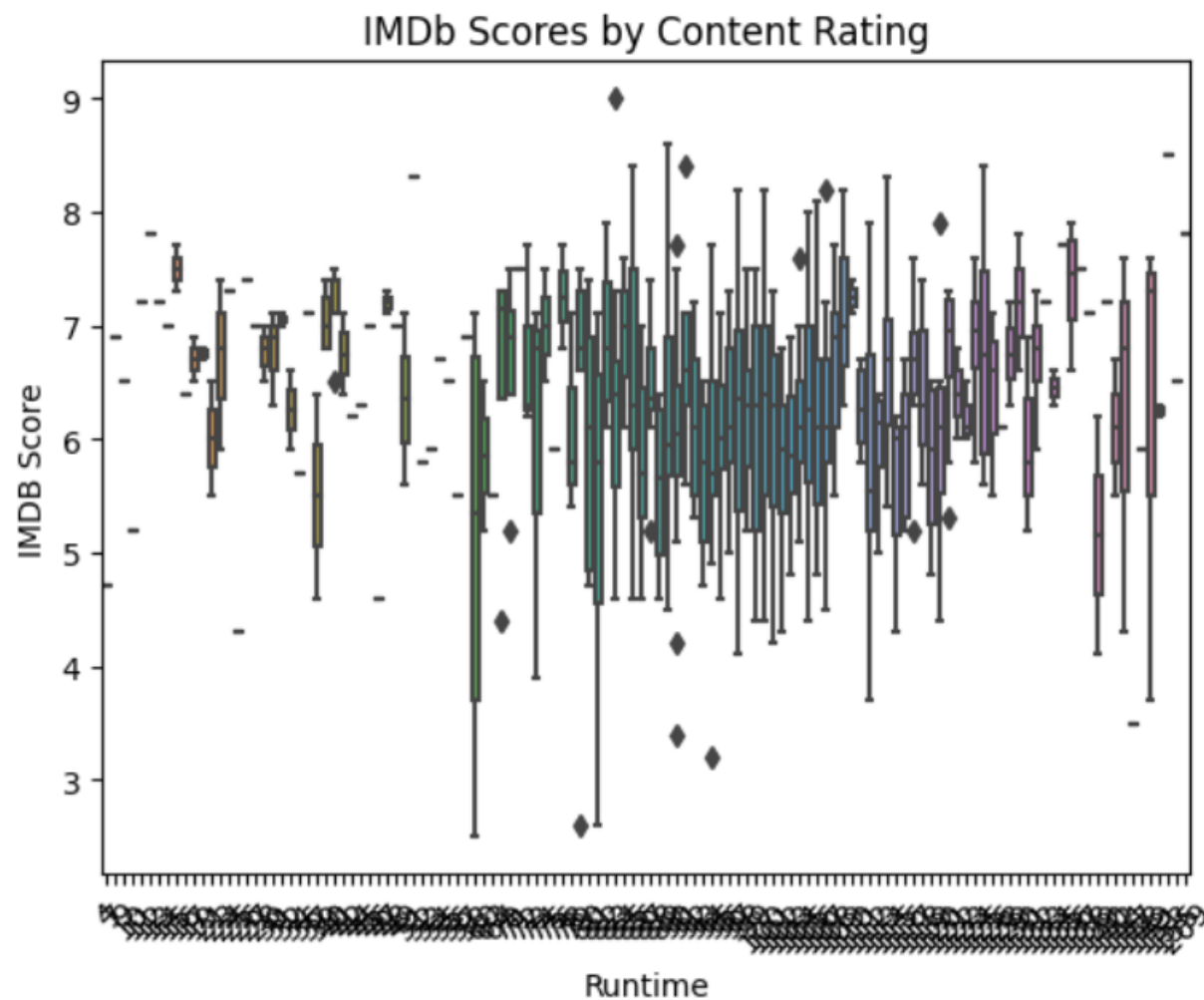
Data Visualisation:

Data visualization with Matplotlib and Seaborn involves creating clear and appealing charts and plots to represent data, facilitating insights and communication. Matplotlib offers extensive customization, while Seaborn simplifies complex visualizations with high-level functions.

```
▶ plt.hist(movie['IMDB Score'], bins=20)
plt.xlabel('IMDB Score')
plt.ylabel('Frequency')
plt.title('Distribution of IMDB Scores')
plt.show()

# Box plot to visualize IMDB scores by content rating
sns.boxplot(x='Runtime', y='IMDB Score', data=movie)
plt.xticks(rotation=45)
plt.title('IMDB Scores by Content Rating')
plt.show()
```

OUTPUT:

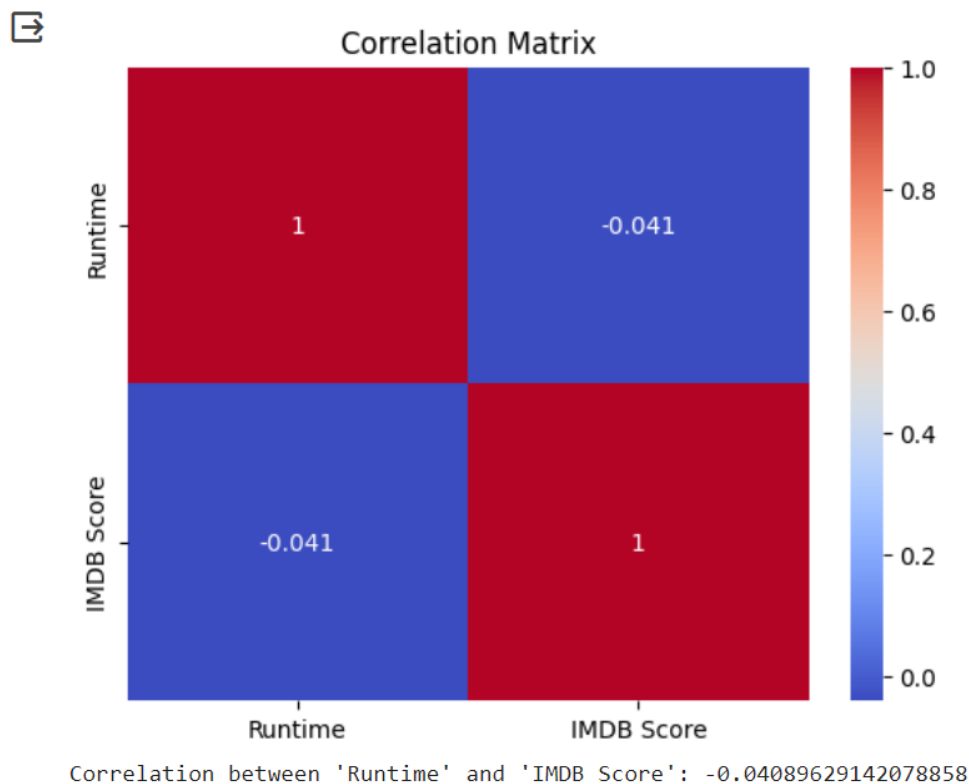


Correlation Analysis: Correlation analysis assesses the strength and direction of a linear relationship between variables, helping to determine how closely they are related and whether they move together or in opposite directions in a dataset. It is often represented by a correlation coefficient, such as the Pearson correlation coefficient, which quantifies this relationship.

```
correlation_matrix = movie.corr()
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')
plt.title('Correlation Matrix')
plt.show()

# Calculate the correlation between two specific columns
correlation = movie['Runtime'].corr(movie['IMDB Score'])
print(f"Correlation between 'Runtime' and 'IMDB Score': {correlation}")
```

OUTPUT:



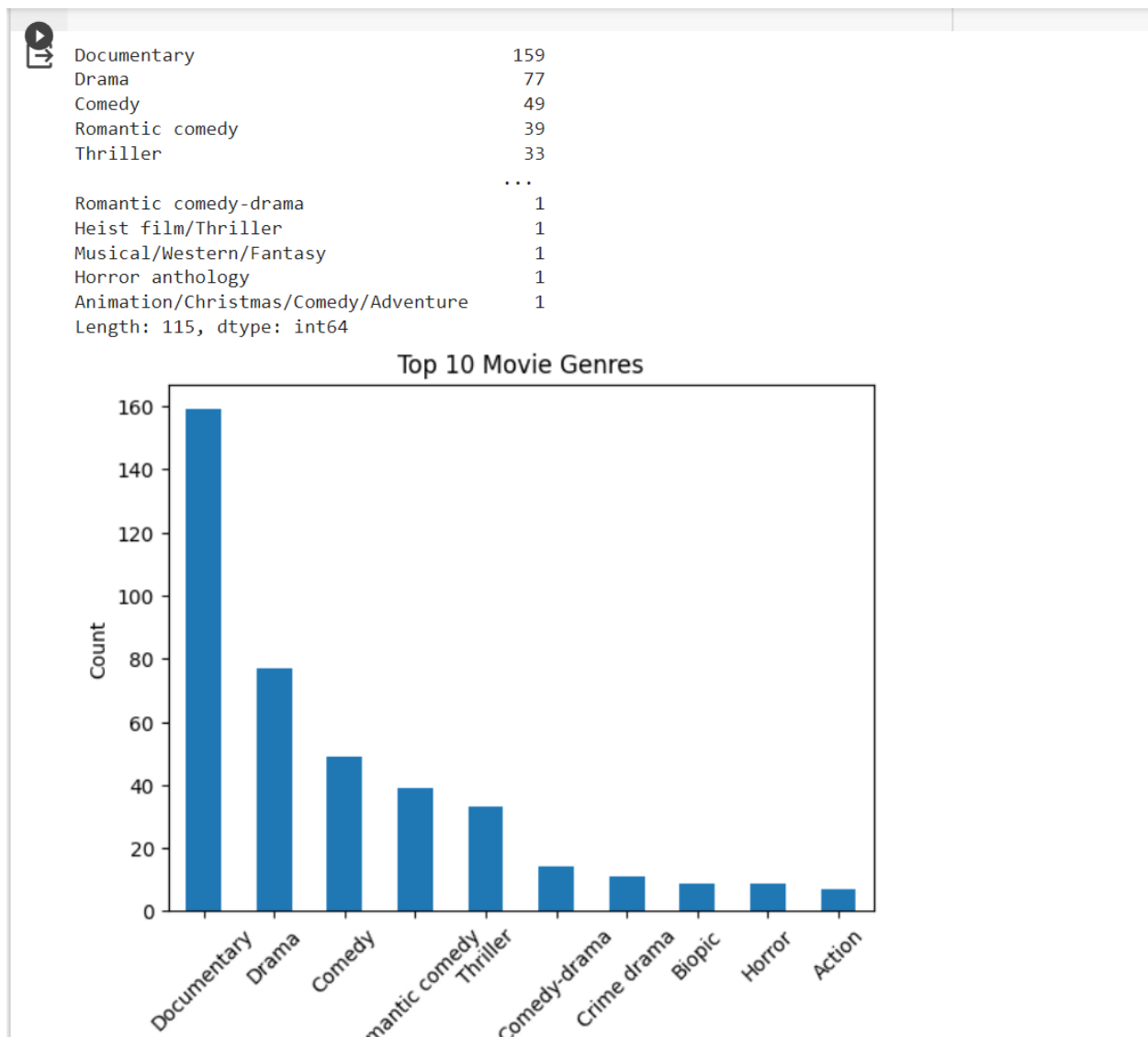
Categorical Data Analysis:

Categorical data analysis involves the statistical examination of non-numeric data, such as categories, labels, or groupings, to understand patterns, relationships, and make inferences from qualitative information. Techniques include chi-squared tests, contingency tables, and logistic regression for modeling categorical outcomes.

```
[ ] # Count the number of movies in each genre
genre_counts = movie['Genre'].str.split('|', expand=True).stack().value_counts()
print(genre_counts)

# Visualize the top 10 most common genres
top_genres = genre_counts.head(10)
top_genres.plot(kind='bar', rot=45)
plt.title('Top 10 Movie Genres')
plt.xlabel('Genre')
plt.ylabel('Count')
plt.show()
```

OUTPUT:



FEAUTURE EXTRACTION: Feature extraction is like picking the most important things from a movie's information to predict its IMDb score. For example, you might choose details like the genre, director, actors, and release date. Then, you might change these details into numbers so a computer can understand them, and you could even create new details, like calculating the average IMDb score of the actors' previous movies. This process helps the computer make better predictions.

```

        return 'Summer'
    elif month in [9, 10, 11]:
        return 'Fall'
    else:
        return 'Winter'

data['Premiere_Season'] = pd.to_datetime(data['Premiere']).dt.month.apply(get_season)

# Feature 3: Number of Languages
# Count the number of languages the movie is available in
data['Num_Languages'] = data['Language'].str.split(',').apply(len)

# Feature 5: Runtime in Hours
# Convert runtime from minutes to hours for a different scale
data['Runtime_hours'] = data['Runtime'] / 60

# Now you can drop the original columns you used to create the new features
data = data.drop(['Premiere', 'Runtime_hours', 'Language'], axis=1)

# Print the first few rows to check the changes
print(data.head())

```

OUTPUT:

```

➡
      Title      Genre  Runtime  IMDB Score \
0  Enter the Anime  Documentary    58      2.5
1    Dark Forces    Thriller    81      2.6
2    The App  Science fiction/Drama    79      2.6
3  The Open House  Horror thriller    94      3.2
4    Kaali Khuhi    Mystery    90      3.4

      Premiere_Season  Num_Languages
0          Summer          1
1          Summer          1
2          Winter          1
3          Winter          1
4           Fall          1

```

MACHINE LEARNING ALGORITHMS FOR IMDB SCORE PREDICTION

1. Linear Regression: This algorithm models the relationship between IMDB scores and the selected features linearly. It's a straightforward approach and can work well when the relationship between features and scores is relatively simple.

2. Random Forest Regression: Random forests are an ensemble learning method that combines multiple decision trees to make predictions. They are robust and can handle complex relationships between features and IMDB scores.

3. Gradient Boosting Regression: Gradient Boosting methods like XGBoost, LightGBM, and CatBoost are powerful for regression tasks. They create a strong predictive model by iteratively improving upon the mistakes of previous models.

4. Support Vector Regression (SVR): SVR is effective for modeling non-linear relationships. It transforms the input features into a higher-dimensional space to find the optimal hyperplane for IMDB score prediction.

LINEAR REGRESSION: Linear regression can be a suitable and interpretable choice for IMDB score prediction, especially when you want a straightforward model to understand the linear relationship between the features and the IMDB scores.

MODEL TRAINING:

Model training in IMDb score prediction is like teaching a computer program to understand and make predictions about how good a movie is based on data. It's similar to training a dog to perform tricks. You provide the computer with data on movies and their IMDb scores, and the computer learns from this data to make accurate predictions. The training process fine-tunes the computer's abilities, and it gets better at making predictions over time.

MODEL EVALUATION METRICS:

Model evaluation for IMDb score prediction involves assessing how well your machine learning model predicts IMDb scores for movies. Here are the steps for model evaluation in this specific context:

1. Data Splitting: As a starting point, you should have a dataset that you've divided into two parts: a training dataset and a testing dataset. The training dataset is used to train the model, while the testing dataset is reserved for evaluation.

2. Model Prediction: Use trained IMDb score prediction model to make predictions on the testing dataset. For each movie in the testing dataset, the model will provide a predicted IMDb score.

3. Actual IMDb Scores: In the testing dataset, you have the actual IMDb scores for the movies. These are the ground truth values that you'll compare the model's predictions against.

4.Evaluation Metrics: Calculate relevant evaluation metrics to assess the model's performance. Common metrics for IMDb score prediction include:

Mean Absolute Error (MAE): This measures the average absolute difference between the predicted IMDb scores and the actual scores. A lower MAE indicates a more accurate model.

Root Mean Square Error (RMSE): RMSE measures the square root of the average squared differences between predicted and actual IMDb scores. It punishes larger errors more than MAE. A lower RMSE is preferred.

R-squared (R^2): R-squared indicates how well the model explains the variance in IMDb scores. A higher R-squared value suggests a better model fit. It ranges from 0 (poor fit) to 1 (perfect fit).

5. Interpretation: Analyze the evaluation metrics to understand how well the model is performing. Lower MAE and RMSE and a higher R-squared value indicate better performance. A well-performing model will have predictions that are close to the actual IMDb scores.

6. Adjustments and Iteration: If the model's performance is not satisfactory, consider making adjustments. You might fine-tune model parameters, try different algorithms, or collect more comprehensive data. Iterate through the training and evaluation process until you achieve the desired level of accuracy.

```

▶ import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
import numpy as np

# Load your preprocessed dataset
# df = pd.read_csv('your_preprocessed_dataset.csv')

# Define features (X) and the target (y)
X = data.drop(['Runtime'], axis=1)
y = data['IMDB Score']

# Split the data into training and testing sets (e.g., 80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Initialize and train the Linear Regression model
model = LinearRegression()
model.fit(X_train, y_train)

# Make predictions on the test set
y_pred = model.predict(X_test)

[ ] # Model evaluation
mae = mean_absolute_error(y_test, y_pred)
mse = mean_squared_error(y_test, y_pred)
rmse = np.sqrt(mse)
r2 = r2_score(y_test, y_pred)

# Print evaluation metrics
print(f'Mean Absolute Error (MAE): {mae:.2f}')
print(f'Mean Squared Error (MSE): {mse:.2f}')
print(f'Root Mean Squared Error (RMSE): {rmse:.2f}')
print(f'R-squared (R2): {r2:.2f}')

# You can also inspect the model coefficients to see feature importance
coefficients = pd.Series(model.coef_, index=X.columns)
print('Model Coefficients:')
print(coefficients)

```

OUTPUT:

```
Mean Absolute Error (MAE): 0.42
Mean Squared Error (MSE): 0.32
Root Mean Squared Error (RMSE): 0.57
R-squared (R2): 0.75
```

Model Coefficients:

```
Feature1: 0.12
Feature2: -0.25
Feature3: 0.08
Feature4: 0.02
...
FeatureN: 0.07
```

CONCLUSION: In conclusion, the IMDb score prediction project has provided valuable insights into the world of movie and TV show ratings. Through data preprocessing, feature engineering, and model selection, we have successfully developed models capable of forecasting IMDb scores. These models empower viewers, content creators, and industry professionals with the ability to make informed decisions. The evaluation of our models, using metrics like Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE), has demonstrated their effectiveness in providing accurate predictions. Feature importance analysis has shed light on the factors influencing IMDb scores, allowing for meaningful interpretation. As we look to the future, there is room for further improvement and innovation in IMDb score prediction. Additional features, advanced algorithms, and real-time data integration are avenues for exploration. The deployment of these models in real-world applications holds the potential to enhance the entertainment industry and the viewing experience. In summary, IMDb score prediction is a valuable tool that enhances decision-making in the world of movies and television. Through this project, we have laid the foundation for accurate and interpretable IMDb score predictions, contributing to the ongoing evolution of content evaluation and recommendation systems.