



Audience Rating **PREDICTION**

By Elamathi. R

TABLE OF CONTENTS

01	PROBLEM STATEMENT
02	OBJECTIVE
03	MACHINE LARNING
04	CODE EXPLANATION
05	EDA ANALYSIS RESULT
06	MODEL BUILDING
07	CONCLUSION & FUTURE WORK

Problem STATEMENT

- With the given dataset, build a model to predict '**audience_rating**'. Demonstrate the working of the pipeline with a notebook, also validate the model for its accuracy.

project **OBJECTIVE**

- The objective of this project is to predict audience ratings for movies using machine learning models, based on various features like cast, genre, and directors.
- The goal is to identify the most accurate model to understand how these factors influence audience sentiment and improve decision-making in the entertainment industry.



Random Forest Regressor



Gradient Boosting Regressor



AdaBoost Regressor



Linear Regression



Ridge Regression



Decision Tree Regressor



Lasso Regression



Support Vector Regressor (SVR)

Machine Learning Models for Audience Rating Prediction



Machine learning, a subset of artificial intelligence, enables systems to learn and improve from data without explicit programming.

In this project, machine learning models are used to predict audience ratings by analyzing historical data and evaluating performance using metrics like R2 score, MSE, and MAE.



CODE EXPLANATION

Data Preprocessing: The dataset is cleaned by handling missing values and transforming categorical features using one-hot encoding and scaling numerical features.

Exploratory Data Analysis (EDA): Basic statistics, correlation analysis, and feature distributions are examined to identify patterns and relationships in the data.

Model Selection: Various machine learning models, including Random Forest, Gradient Boosting, and Linear Regression, are trained to predict audience ratings.

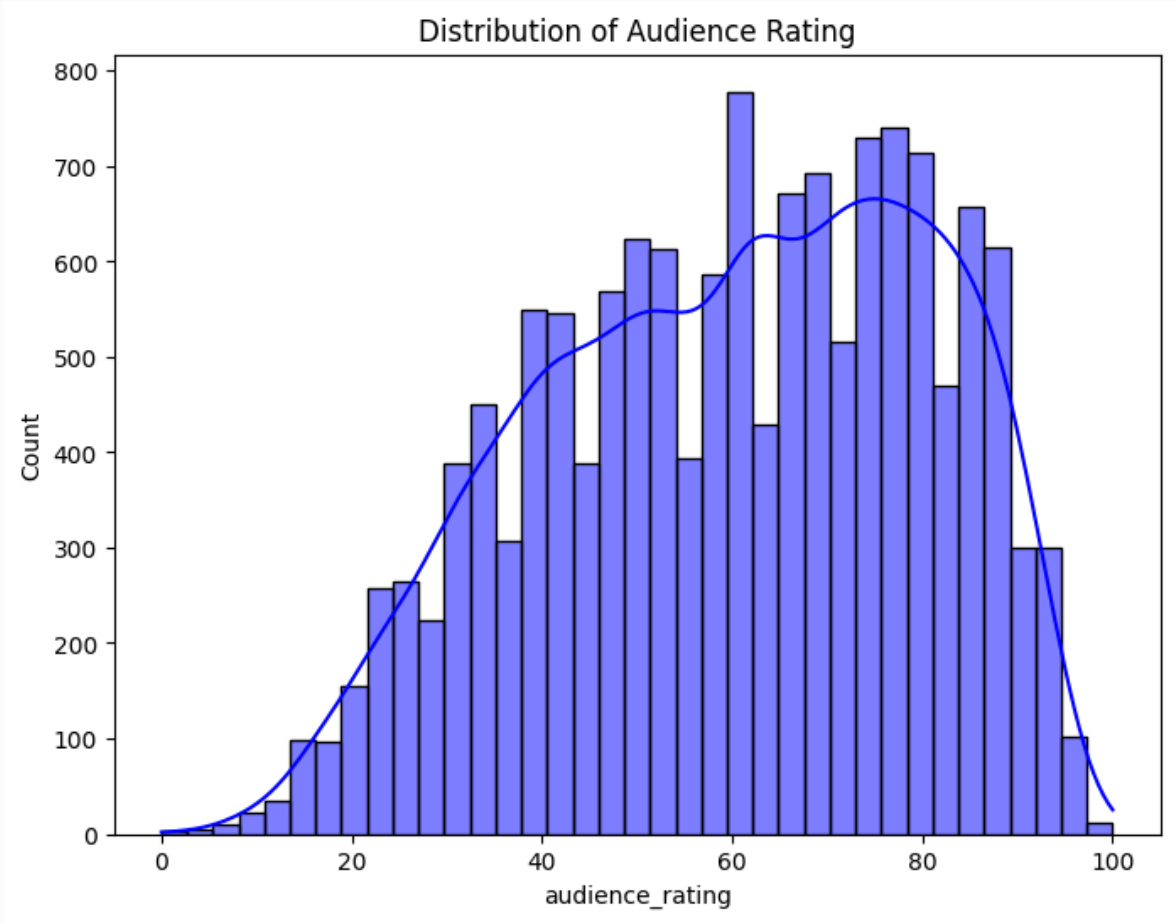
Model Evaluation: Performance metrics like MSE, MAE, and R2 are used to evaluate and compare model effectiveness.

Results: The models' performance is visualized, and the best model is selected based on the highest R2 score.

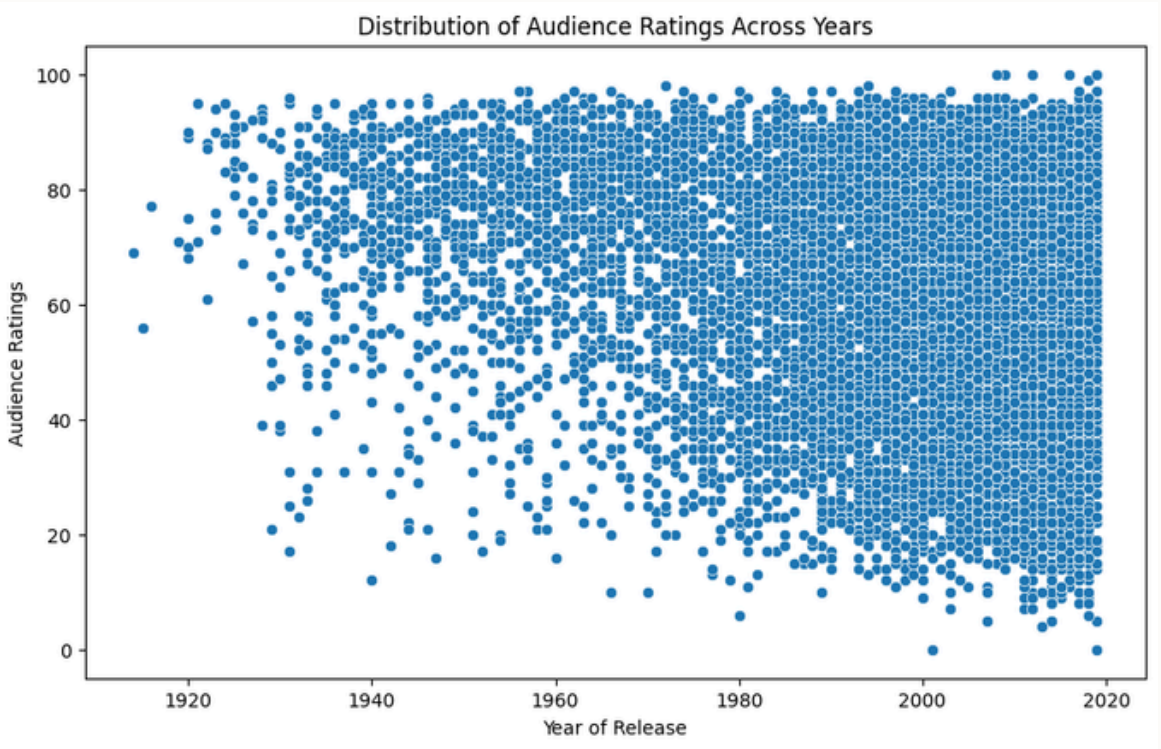
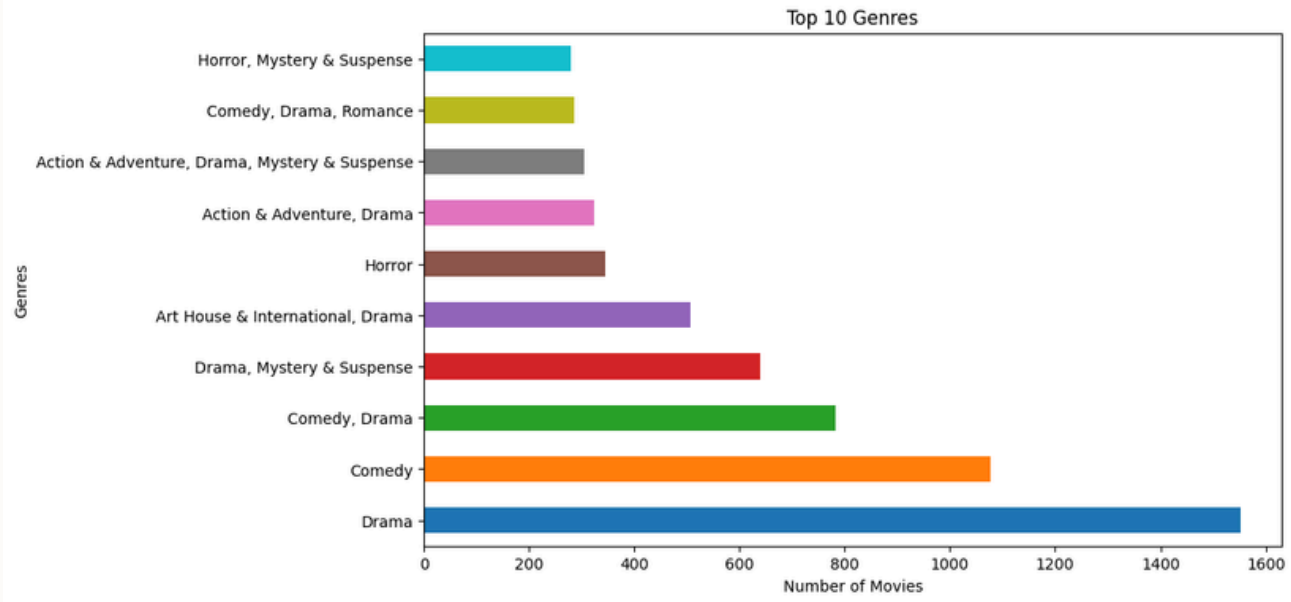
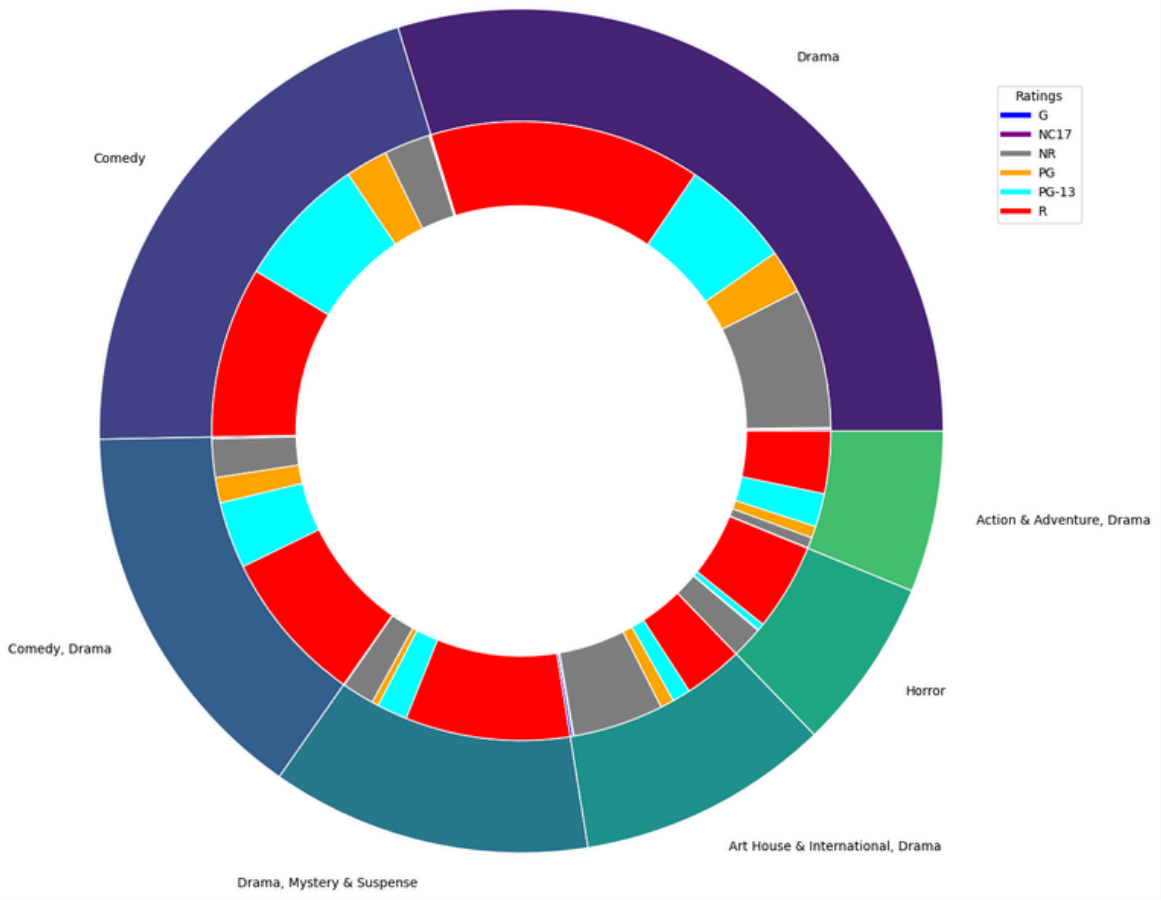
```
1 # Calculate median for numerical columns
2 median_runtime = audience_data['runtime_in_minutes'].median()
3 median_rating = audience_data['audience_rating'].median()
```

```
1 # Convert the 'in_theaters_date' column to datetime format and extract the release year
2 # This helps in aggregating data by release year
3 audience_data['release_year'] = pd.to_datetime(audience_data['in_theaters_date'], format='%d-%m-%Y').dt.year
4
5 # Creating a histogram to visualize the number of movies released each year
6 plt.figure(figsize=(10, 6))
7
8 # Counting movies for each release year and sorting by year
9 release_year_counts = audience_data['release_year'].value_counts().sort_index()
10 release_year_counts.plot(kind='bar', color='skyblue')
11
12 # Adding plot title and axis labels
13 plt.title("Number of Movies by Year of Release")
14 plt.xlabel("Year")
15 plt.ylabel("Number of Movies")
16
17 # Customizing x-axis labels to show every 5th year for better readability
18 years = release_year_counts.index # Get unique years from the data
19 plt.xticks(
20     ticks=range(0, len(years), 5), # Adjust tick spacing to every 5th year
21     labels=[str(year) for year in years[::5]] # Display labels for every 5th year
22 )
23
24 # Display the plot
25 plt.show()
```

```
# Dictionary to store different regression models
models = {
    'RandomForest': RandomForestRegressor(random_state=42),
    'GradientBoosting': GradientBoostingRegressor(random_state=42),
    'AdaBoost': AdaBoostRegressor(random_state=42),
    'LinearRegression': LinearRegression(),
    'Ridge': Ridge(),
    'Lasso': Lasso(),
    'DecisionTree': DecisionTreeRegressor(random_state=42),
    'SVR': SVR()
}
```



EDA ANALYSIS OUTPUT



Model BUILDING

Data Split: The dataset is split into training and testing sets using an 80-20 ratio, ensuring that the model is trained on a large portion of the data and tested on unseen data for validation.

Training Pipeline: A pipeline is used to preprocess the data (handling missing values, scaling numerical features, one-hot encoding categorical features) and then train the model.

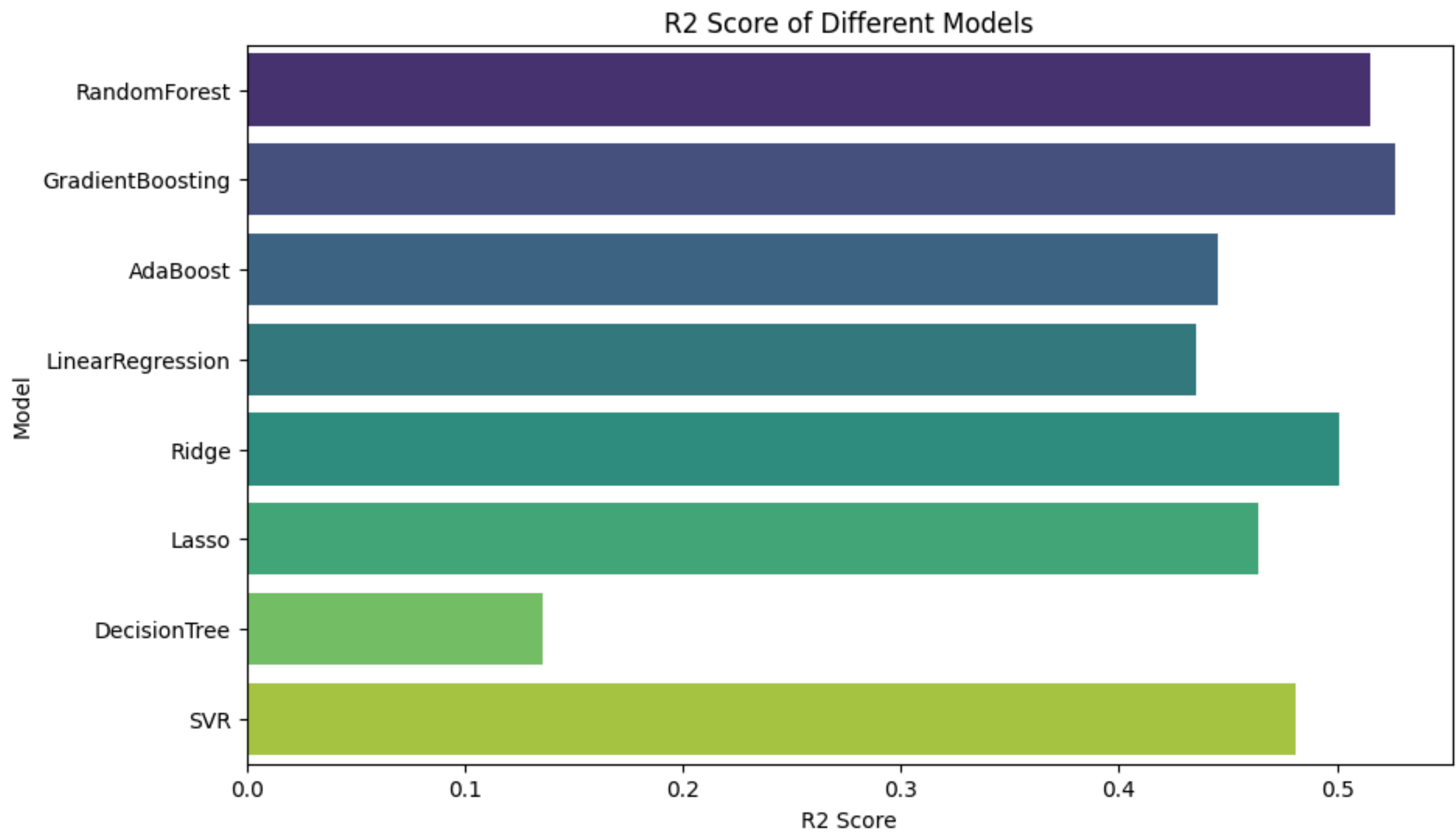
Performance Evaluation: Each model is evaluated using performance metrics such as:

- Mean Squared Error (MSE)
- Mean Absolute Error (MAE)
- R-squared (R²) score

Cross-Validation: Cross-validation is performed for the best model to assess its generalizability.

Best Model: The model with the highest R² score (Gradient Boosting) was selected for further analysis and deployment.

OVERALL MODEL PERFORMANCE



CONCLUSION

Best Model

GradientBoostingRegressor
achieved the highest R2 score
of 0.53.

Weakest Model

DecisionTree
performed poorly with
an R2 of 0.14.

FUTURE WORK

Model Tuning

Fine-tune hyperparameters of the GradientBoostingRegressor and RandomForest models for improved performance.

Feature Engineering

Explore additional features and advanced techniques to enhance predictive accuracy.

Ensemble Methods

Combine multiple models using ensemble techniques to boost model robustness and performance.