

QS World University Rankings 2025: Machine Learning Analysis and Prediction

Presented by Youssef Masoud, Arnav Kucheriya, and Kushagra
Verma

Agenda

- Dataset Overview
- Data Cleaning and Preprocessing
- Exploratory Data Analysis
- Model Selection and Training
- Model Evaluation and Visualizations
- Feature Importance
- Conclusions

Project Overview

- Understand factors influencing university rankings
- Predict Academic Reputation Scores using machine learning
- Compare the performance of Linear Regression and Random Forest models
- Identify key predictors for academic reputation

Dataset Overview

- Source: QS World University Rankings 2025
- 27 Columns, 1,500+ Universities
- Key Features: Reputation scores, research impact, faculty ratios

Dataset Details

Feature Types:

- Categorical: Institution Name, Location
- Numerical: Various score features (e.g., citations, employment outcomes)

Special Notes:

- Missing values represented as '-'
- Some features highly correlated with each other

Data Cleaning and Preprocessing

- Convert score columns to numeric
- Handle missing values
- Label encode categorical features

Detailed Preprocessing

Missing Values:

- ~5% of records dropped due to incomplete data

Label Encoding:

- Institution Name and Location turned into integers

Rationale:

- Ensured clean numeric matrix input for modeling

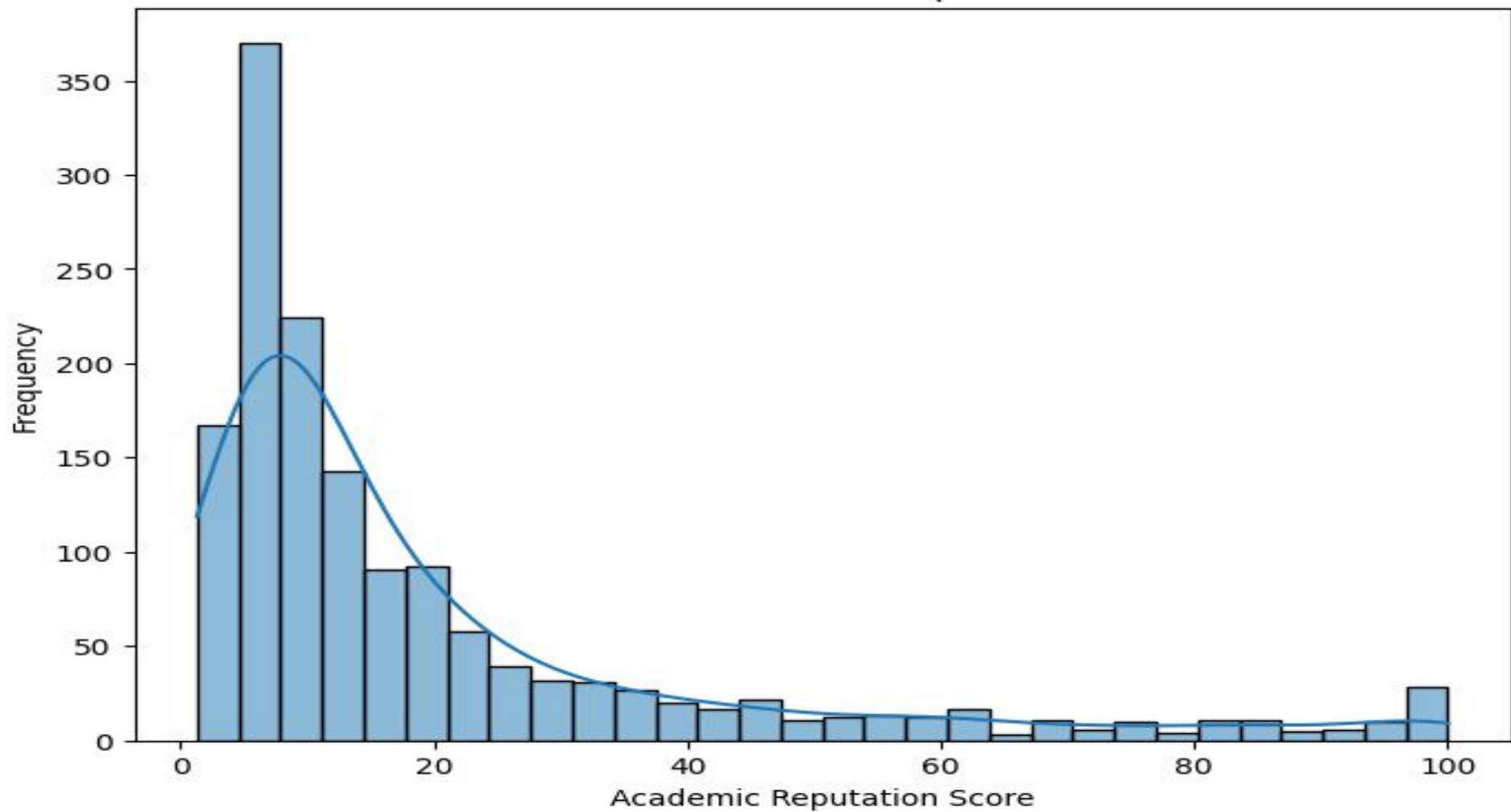
Why Data Preprocessing Matters

- Ensures model compatibility Improves accuracy
- Improves accuracy
- Reduces biases

Exploratory Data Analysis (EDA)

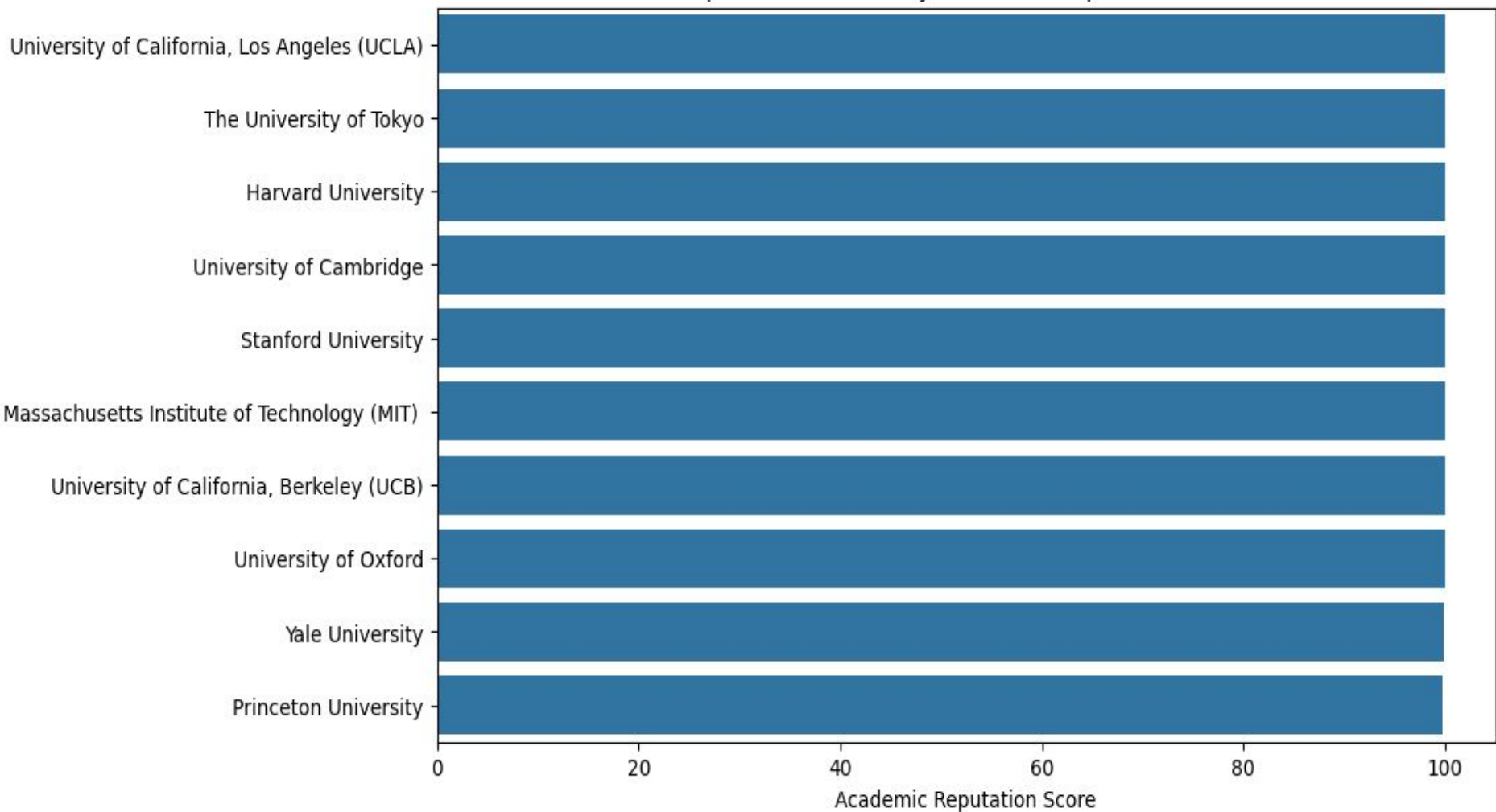
- Understand patterns and relationships
- Identify important variables
- Guide feature selection

Distribution of Academic Reputation Scores

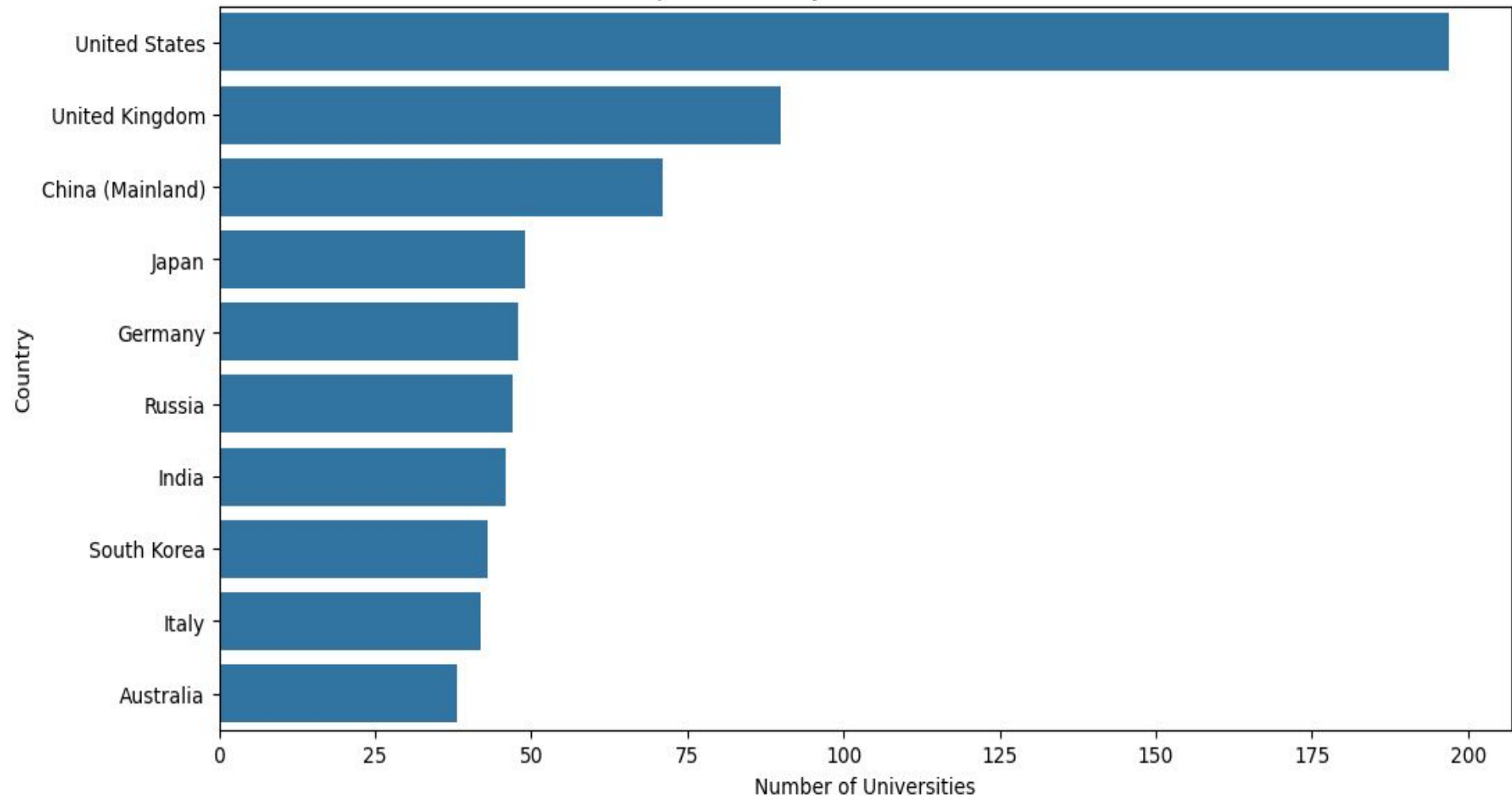


Top 10 Universities by Academic Reputation Score

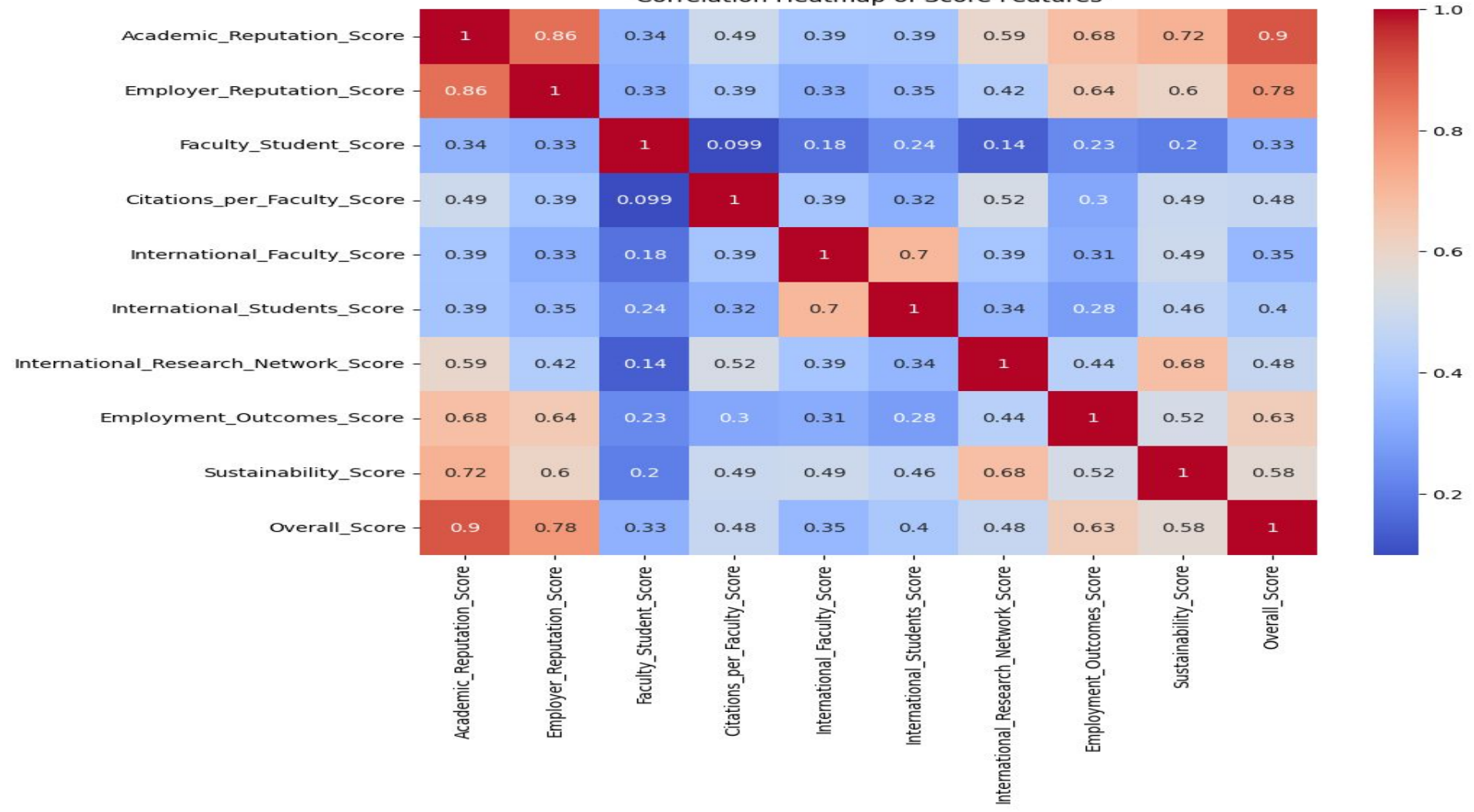
University



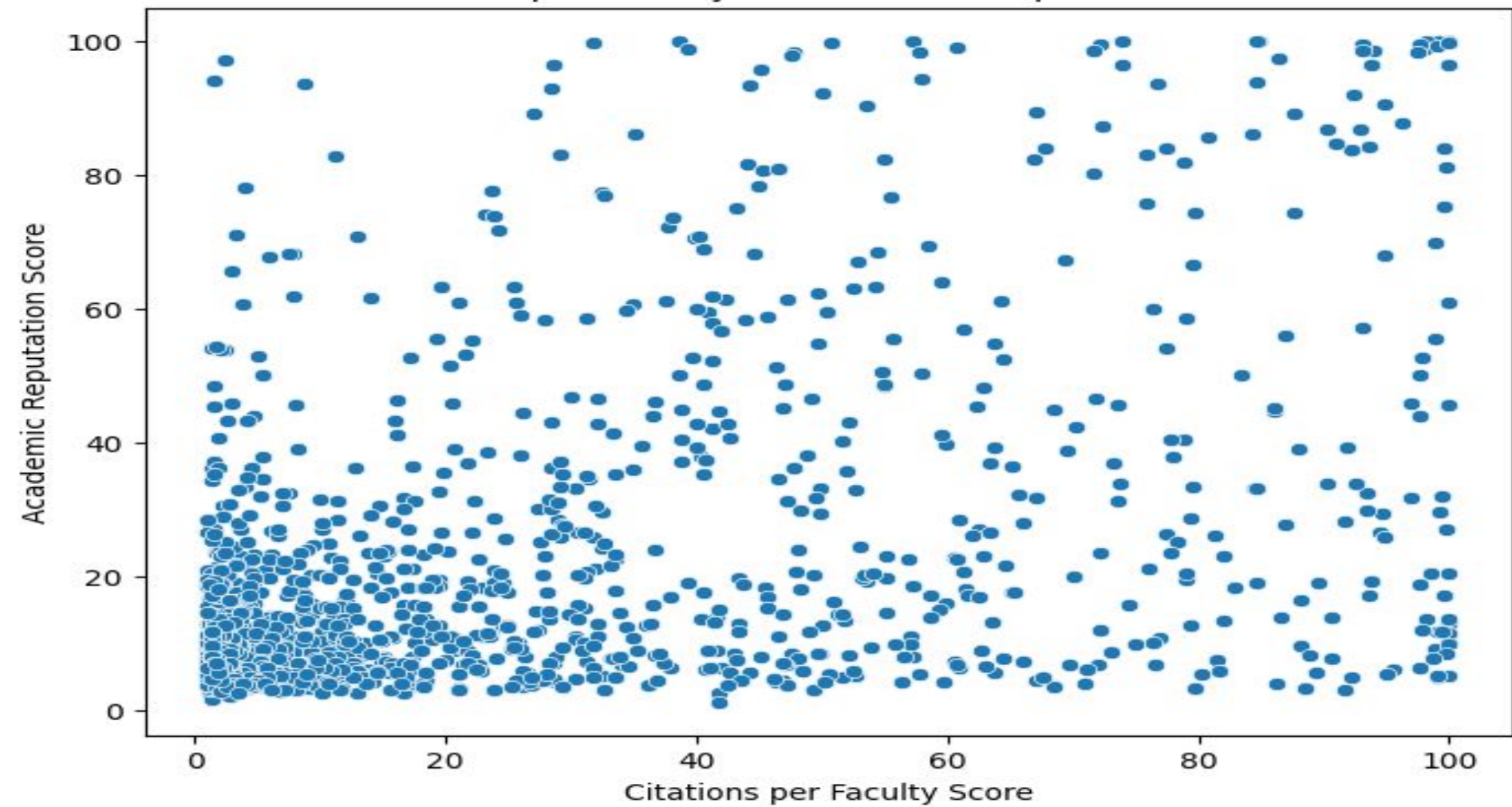
Top Countries by Number of Universities



Correlation Heatmap of Score Features



Citations per Faculty vs Academic Reputation Score



Model Training Summary

- Linear Regression: Quick training, simple assumptions
- Random Forest: Handles complexity, avoids overfitting
- Purpose: Predict Academic Reputation Scores accurately

Training Code: Linear Regression

```
from sklearn.linear_model import LinearRegression  
lr_model = LinearRegression()  
lr_model.fit(X_train, y_train)
```


Training Code: Random Forest

```
from sklearn.ensemble import RandomForestRegressor  
rf_model = RandomForestRegressor(n_estimators=100, random_state=  
rf_model.fit(X_train, y_train)
```

Model Evaluation Code

```
from sklearn.metrics import mean_squared_error, r2_score

# Predictions
y_pred_lr = lr_model.predict(X_test)
y_pred_rf = rf_model.predict(X_test)

# Evaluation Metrics
rmse_lr = mean_squared_error(y_test, y_pred_lr, squared=False)
r2_lr = r2_score(y_test, y_pred_lr)

rmse_rf = mean_squared_error(y_test, y_pred_rf, squared=False)
r2_rf = r2_score(y_test, y_pred_rf)
```

Model Performance Comparison

Linear Regression:

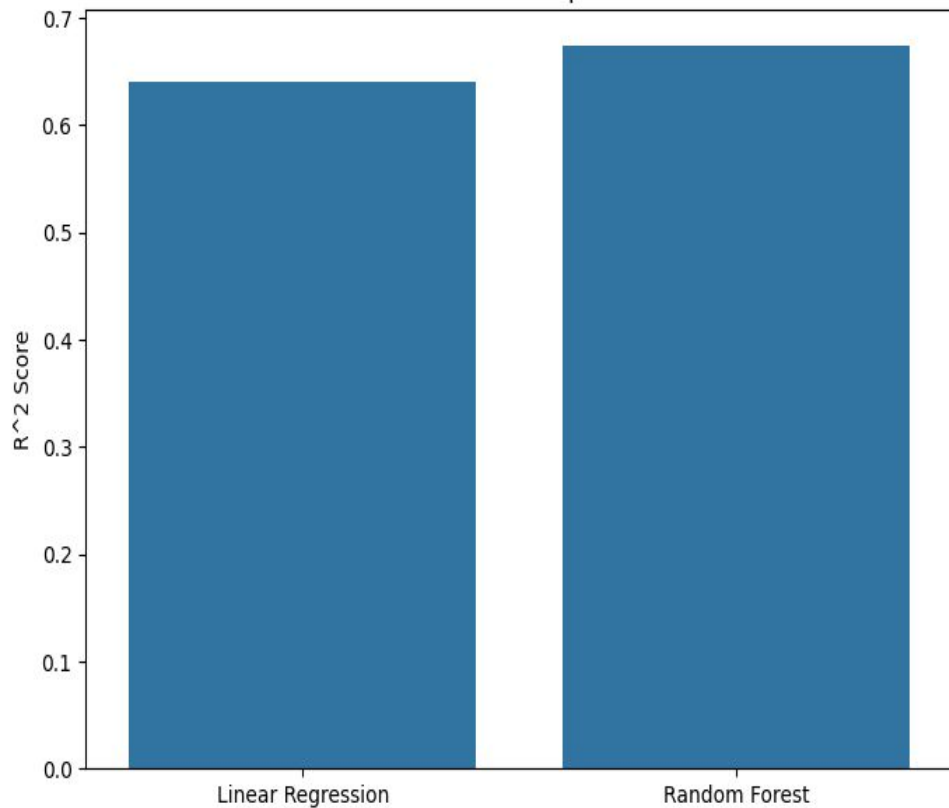
- RMSE: 15.53
- R^2 Score: 0.64

Random Forest:

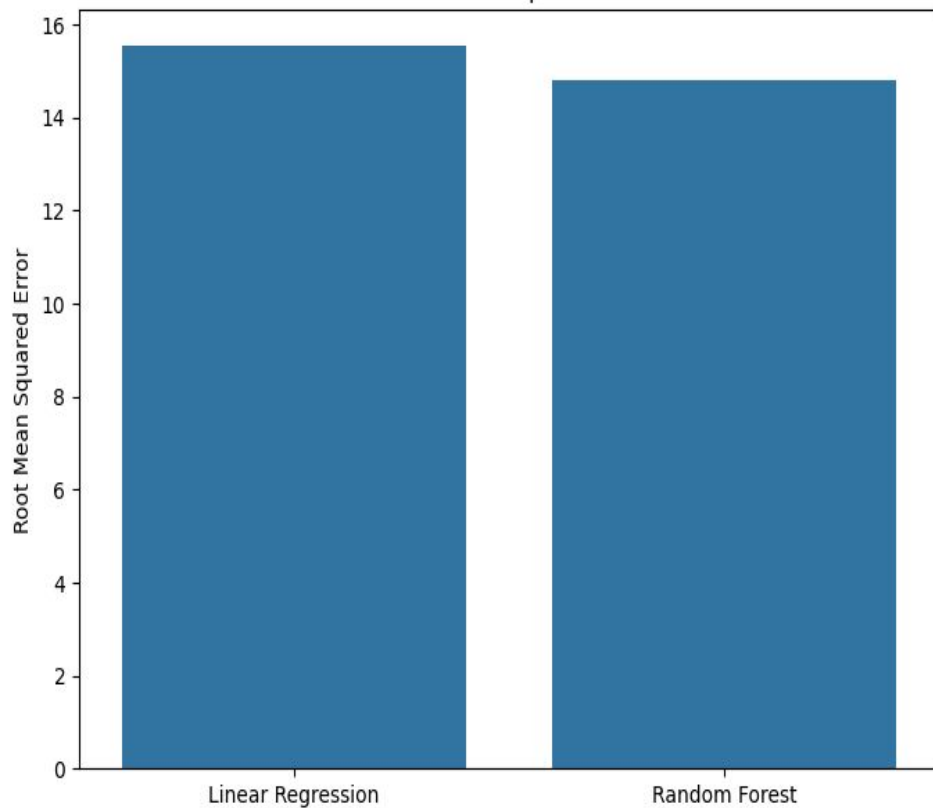
- RMSE: 14.80
- R^2 Score: 0.67

RMSE and R² Score Comparison

R² Score Comparison



RMSE Comparison



Why These Models?

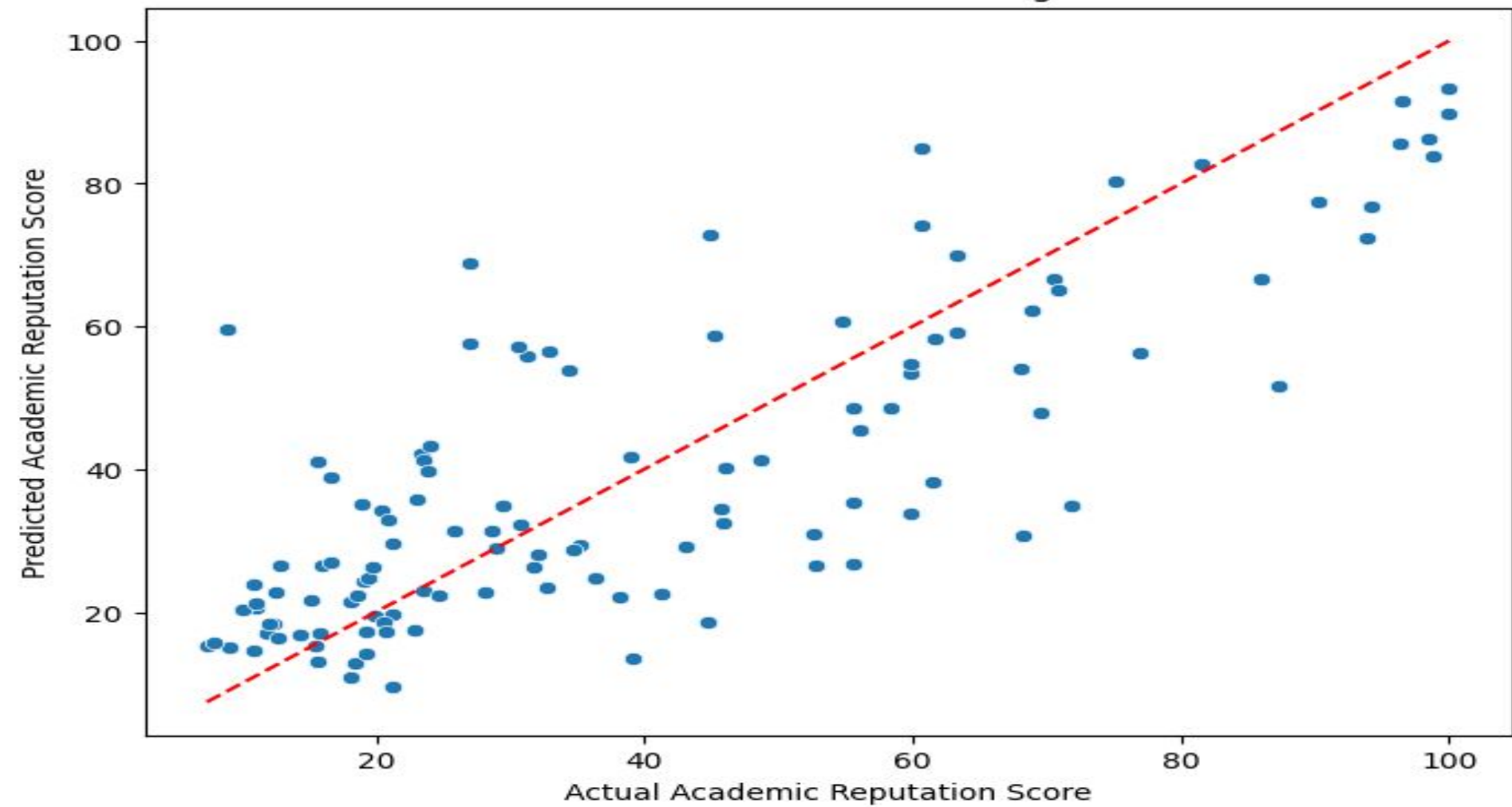
Linear Regression:

- Easy to interpret, good baseline

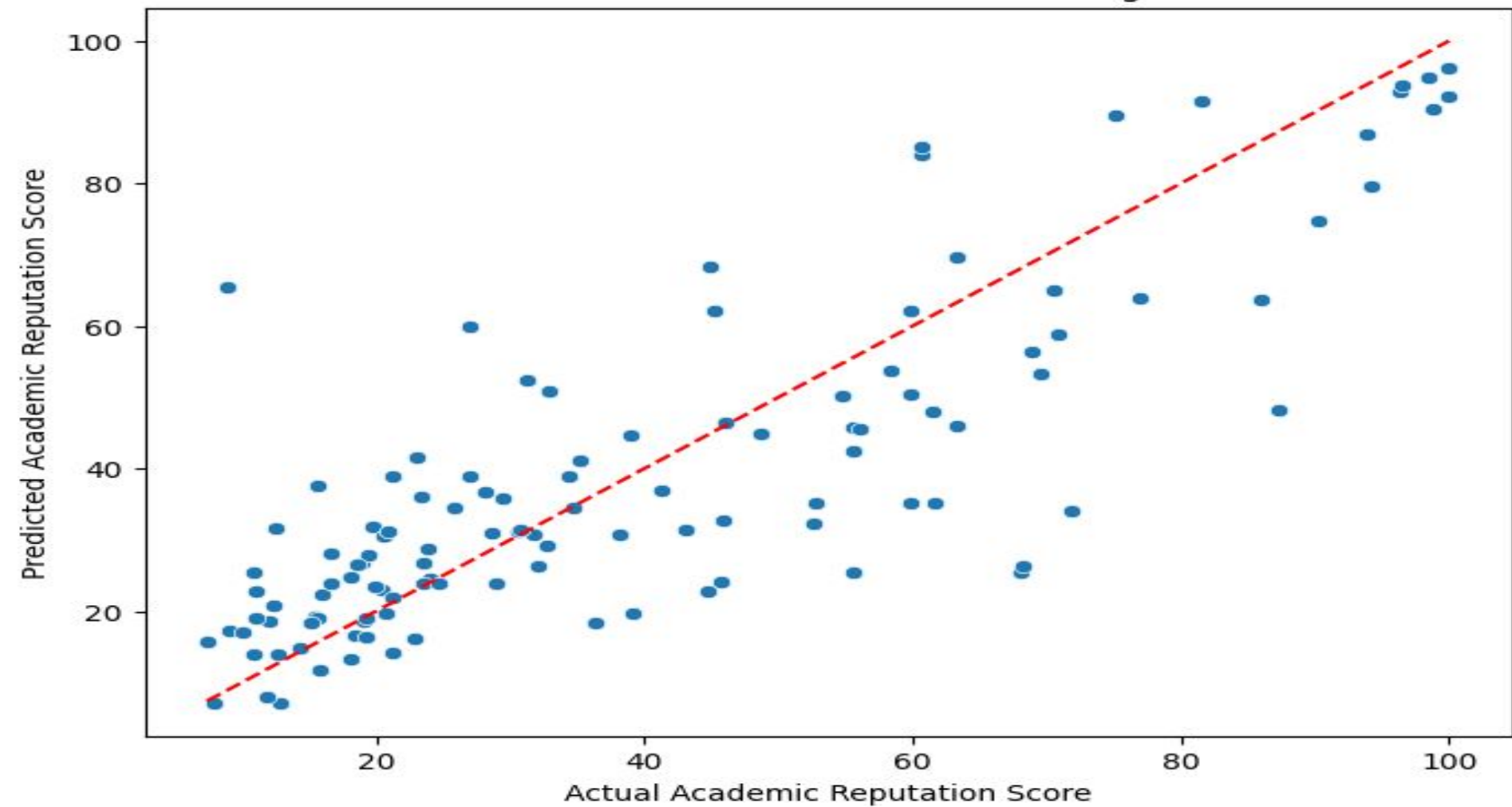
Random Forest:

- Captures non-linearities, robust to overfitting
- Comparison provides balanced insight

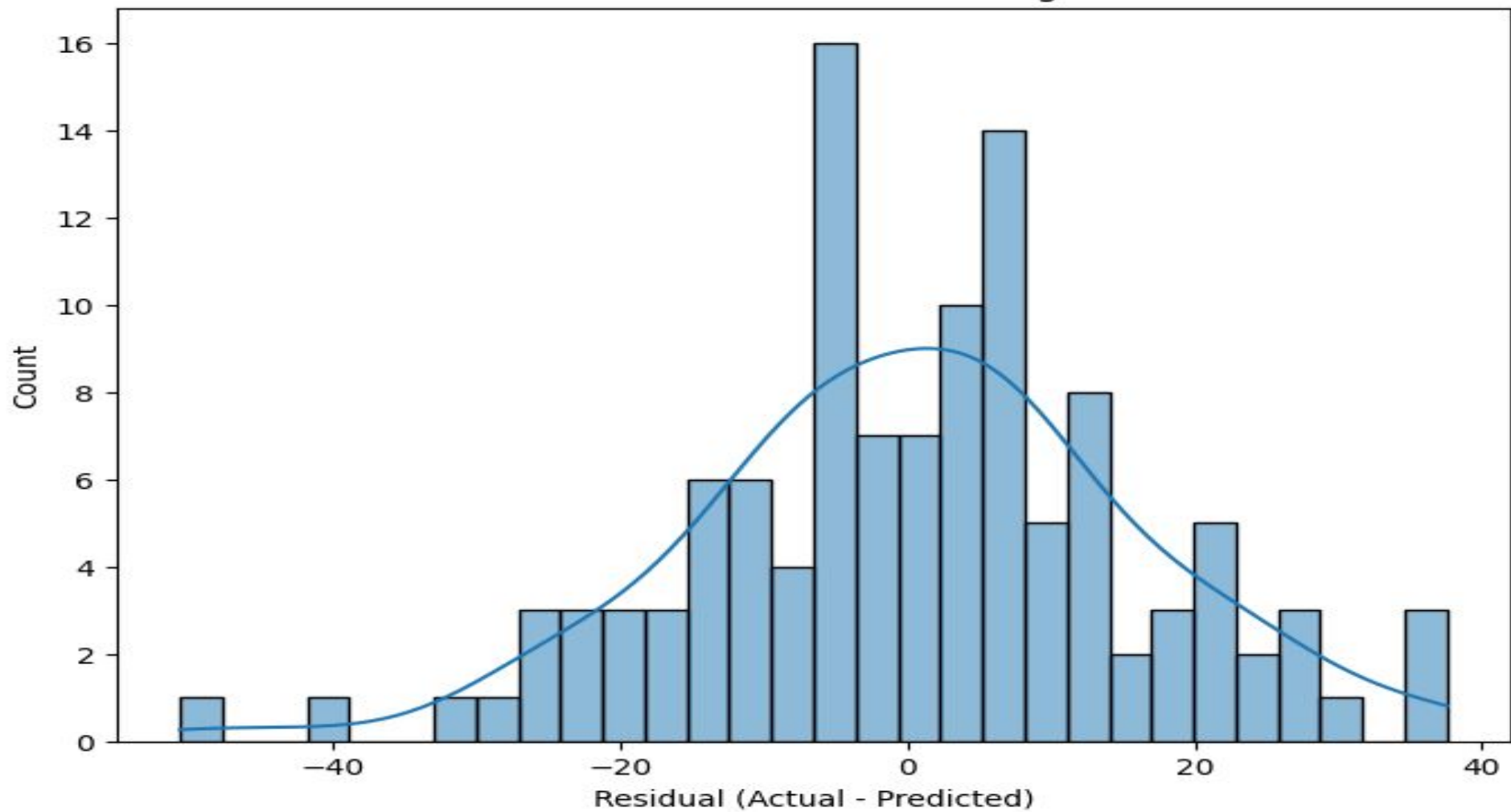
Actual vs Predicted - Linear Regression



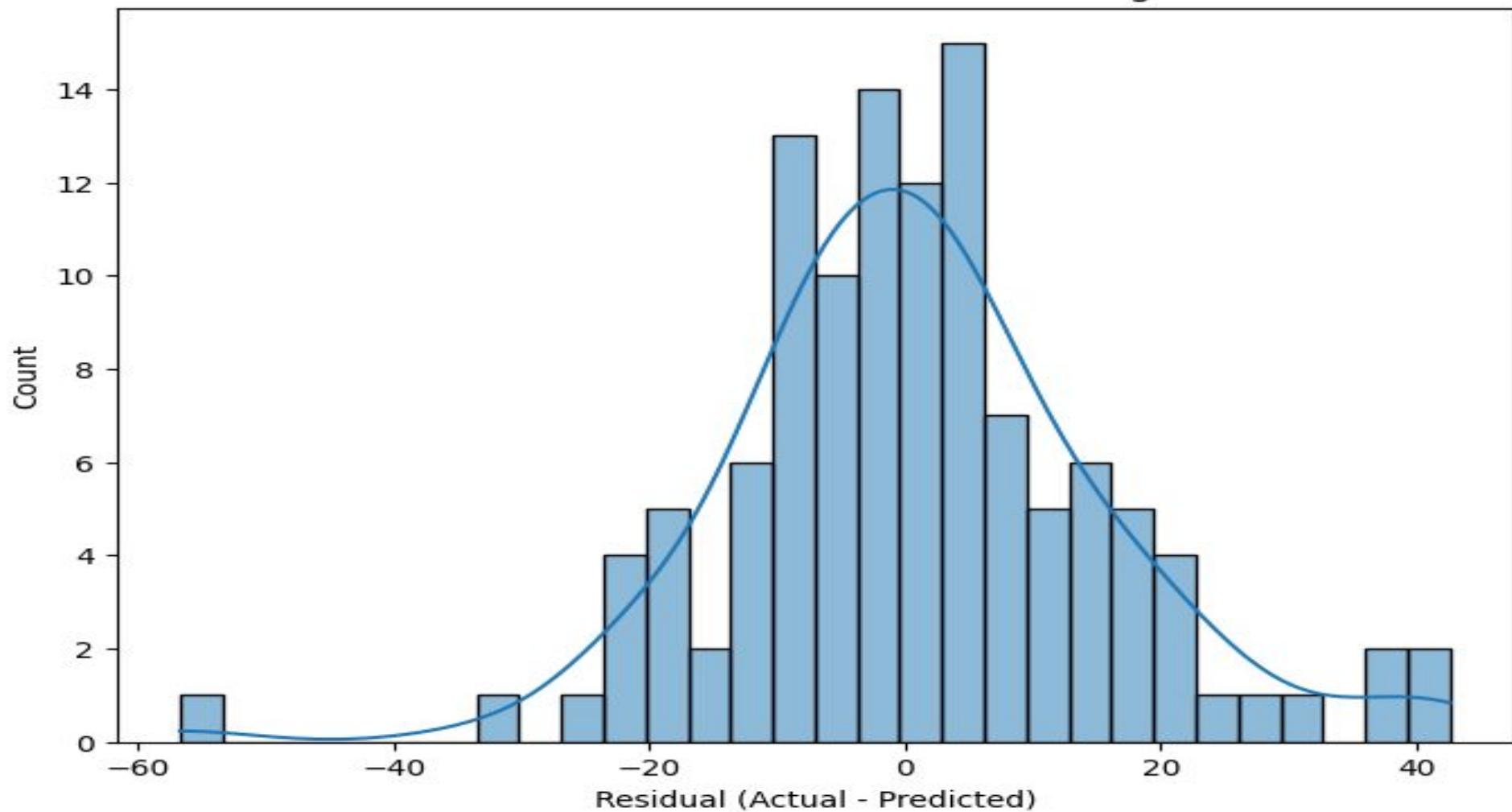
Actual vs Predicted - Random Forest Regressor



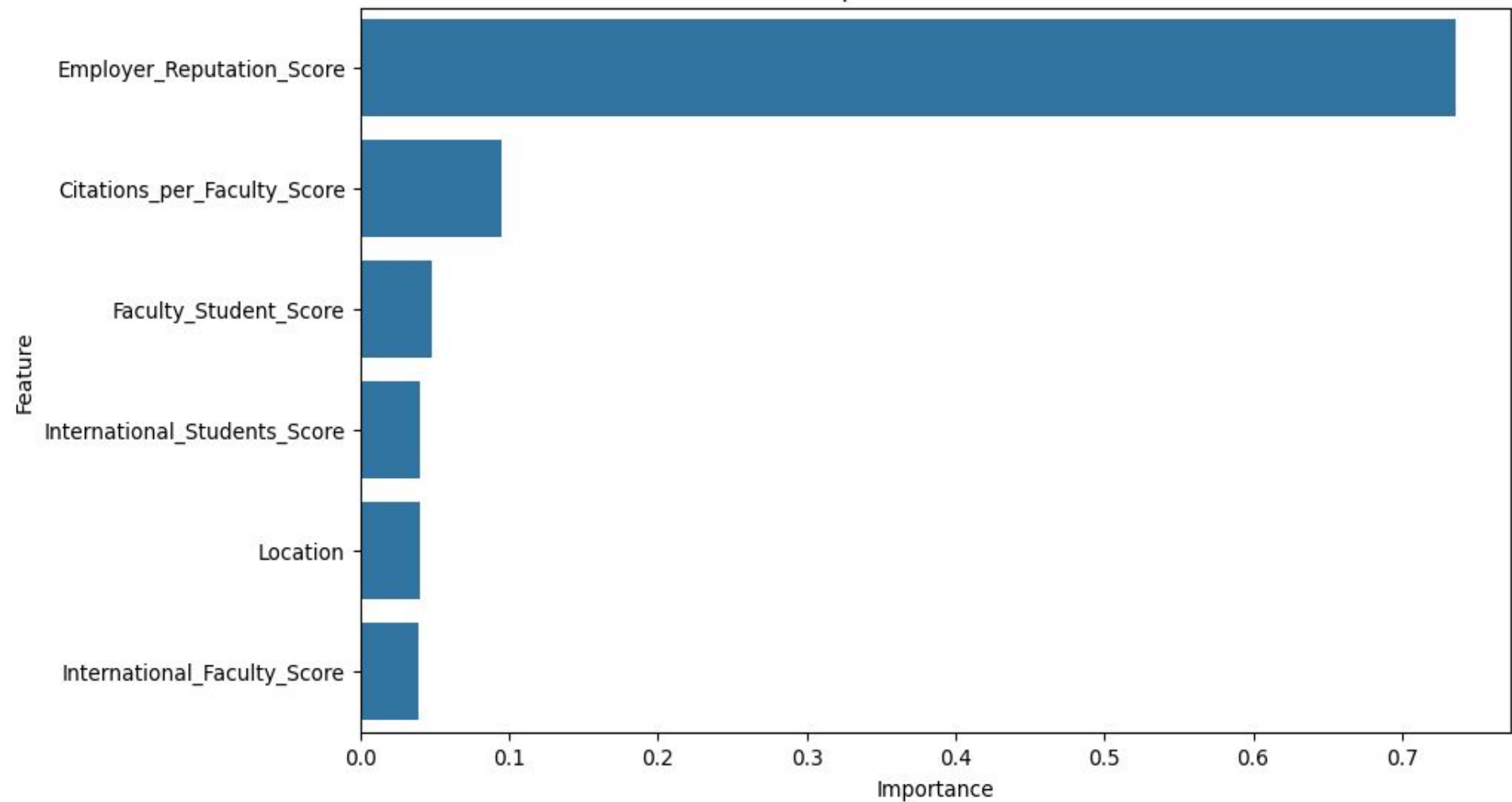
Residual Distribution - Linear Regression



Residual Distribution - Random Forest Regressor



Feature Importance - Random Forest



Key Insights

- Research impact (citations) is a critical determinant of academic reputation.
- Employer reputation plays a significant supporting role.
- Random Forest captures hidden patterns better than linear models.
- High-quality data preprocessing is crucial for model accuracy.
- Institutional strategies should prioritize research output and industry engagement.

Challenges Faced

Data Issues:

- Missing or corrupted entries
- Inconsistent data formatting

Modeling Issues:

- Linear Regression limited in capturing complex patterns

Mitigation:

- Careful data cleaning, choosing an ensemble model

Conclusion

This project analyzed the QS World University Rankings 2025 dataset to predict Academic Reputation Scores based on institutional features. Through a structured machine learning workflow, several important insights emerged:

- Random Forest outperformed Linear Regression
- Focus on research and employer links
- Machine learning provides strategic insights

Thank you