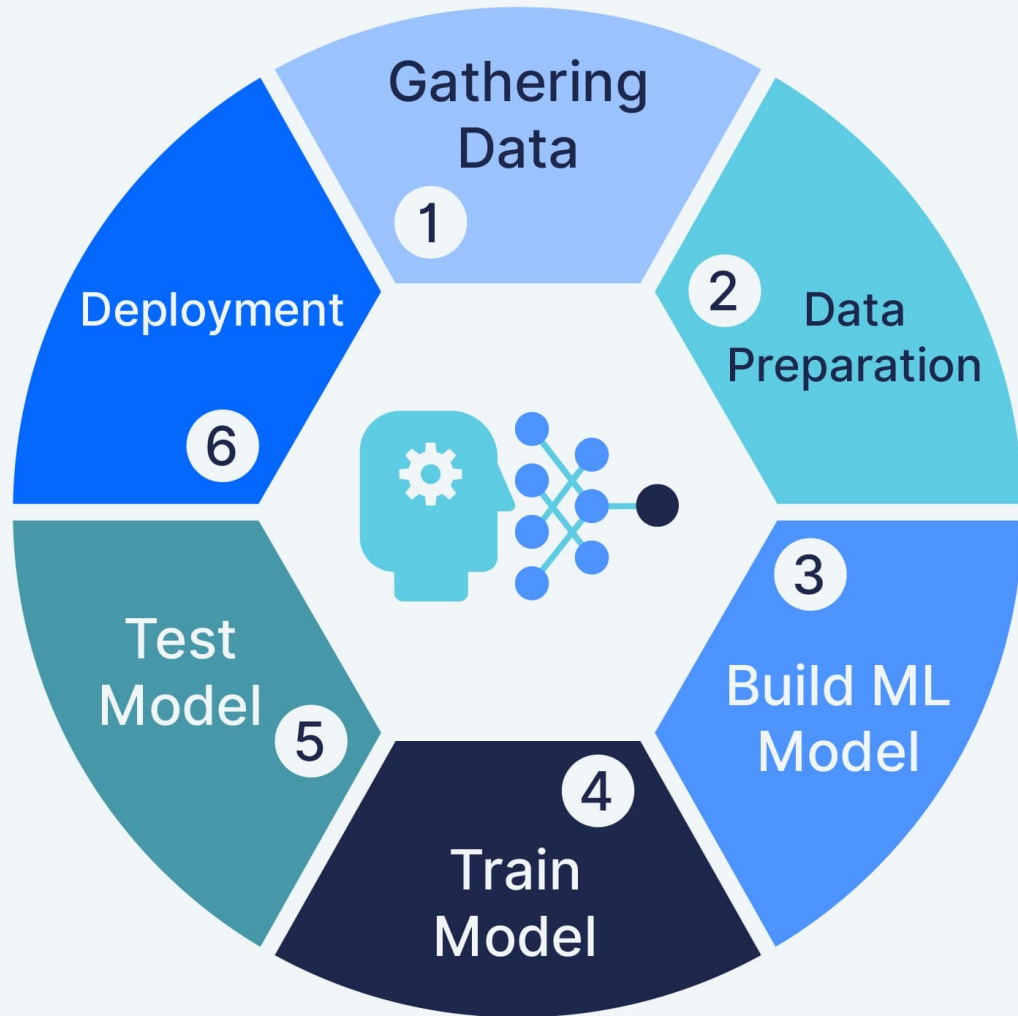


# Machine Learning Life-cycle



# Introduction

**The Selection, Application, and Evaluation of Data Science Methods in University Ranking Prediction**

**Overview:** Using advanced data science methodologies, a comprehensive analysis of the World University Rankings dataset is conducted.

**Objective:** The purpose of this study is to use machine learning models and data science techniques to predict university rankings.

Heatmap showing the correlation matrix for 11 variables. The color scale ranges from -1 (dark red) to 1 (dark blue).

	scores_teaching	scores_research	scores_citations	scores_international_outlook	member_level	stats_number_students	stats_student_staff_ratio	stats_pc_intl_students	overall_score	female_ratio	male_ratio
scores_teaching	1.0	0.85	0.65	0.55	0.15	0.05	-0.15	0.65	0.75	0.05	0.05
scores_research	0.85	1.0	0.75	0.65	0.25	0.15	-0.05	0.75	0.85	0.15	0.15
scores_citations	0.65	0.75	1.0	0.85	0.35	0.25	0.15	0.85	0.95	0.25	0.25
scores_international_outlook	0.55	0.65	0.85	1.0	0.45	0.35	0.25	0.75	0.85	0.35	0.35
member_level	0.15	0.25	0.35	0.45	1.0	0.15	0.05	0.15	0.25	0.05	0.05
stats_number_students	0.05	0.15	0.25	0.35	0.15	1.0	0.15	0.05	0.15	0.05	0.05
stats_student_staff_ratio	-0.15	-0.05	0.15	0.25	0.05	0.15	1.0	0.15	0.05	0.05	0.05
stats_pc_intl_students	0.65	0.75	0.85	0.75	0.15	0.05	0.15	1.0	0.75	0.15	0.15
overall_score	0.75	0.85	0.95	0.85	0.25	0.15	0.05	0.75	1.0	0.15	0.15
female_ratio	0.05	0.15	0.25	0.35	0.05	0.05	0.05	0.15	0.15	1.0	0.85
male_ratio	0.05	0.15	0.25	0.35	0.05	0.05	0.05	0.15	0.15	0.85	1.0

Heatmap showing the correlation matrix for 11 variables. The color scale ranges from -1 (dark red) to 1 (dark blue).

	scores_teaching	scores_research	scores_citations	scores_international_outlook	member_level	stats_number_students	stats_student_staff_ratio	stats_pc_intl_students	overall_score	female_ratio	male_ratio
scores_teaching	1.0	0.85	0.65	0.55	0.15	0.05	-0.15	0.65	0.75	0.05	0.05
scores_research	0.85	1.0	0.75	0.65	0.25	0.15	-0.05	0.75	0.85	0.15	0.15
scores_citations	0.65	0.75	1.0	0.85	0.35	0.25	0.15	0.85	0.95	0.25	0.25
scores_international_outlook	0.55	0.65	0.85	1.0	0.45	0.35	0.25	0.75	0.85	0.35	0.35
member_level	0.15	0.25	0.35	0.45	1.0	0.15	0.05	0.15	0.25	0.05	0.05
stats_number_students	0.05	0.15	0.25	0.35	0.15	1.0	0.15	0.05	0.15	0.05	0.05
stats_student_staff_ratio	-0.15	-0.05	0.15	0.25	0.05	0.15	1.0	0.15	0.05	0.05	0.05
stats_pc_intl_students	0.65	0.75	0.85	0.75	0.15	0.05	0.15	1.0	0.75	0.15	0.15
overall_score	0.75	0.85	0.95	0.85	0.25	0.15	0.05	0.75	1.0	0.15	0.15
female_ratio	0.05	0.15	0.25	0.35	0.05	0.05	0.05	0.15	0.15	1.0	0.85
male_ratio	0.05	0.15	0.25	0.35	0.05	0.05	0.05	0.15	0.15	0.85	1.0

Heatmap showing the correlation matrix for 11 variables. The color scale ranges from -1 (dark red) to 1 (dark blue).

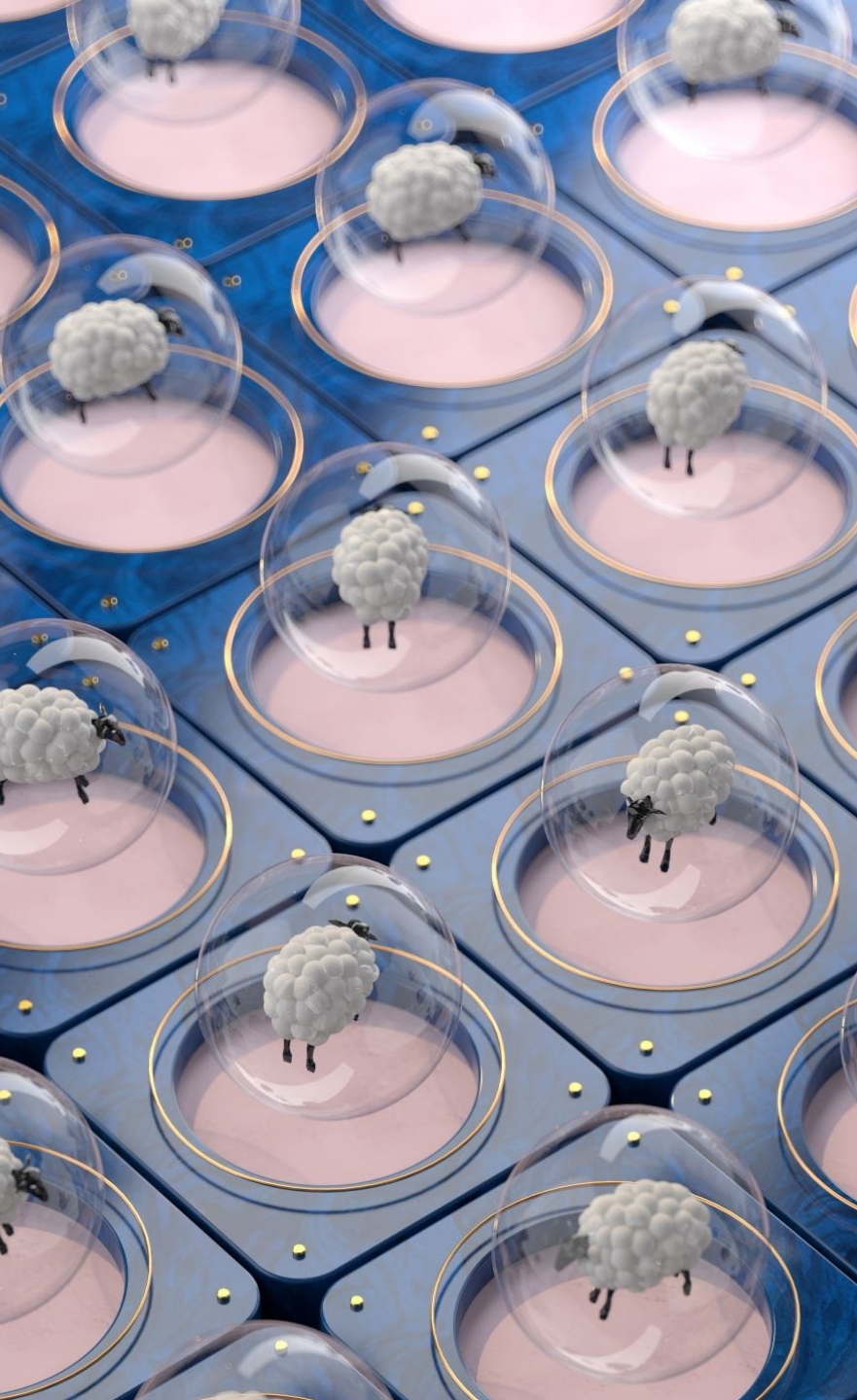
	scores_teaching	scores_research	scores_citations	scores_international_outlook	member_level	stats_number_students	stats_student_staff_ratio	stats_pc_intl_students	overall_score	female_ratio	male_ratio
scores_teaching	1.0	0.85	0.65	0.55	0.15	0.05	-0.15	0.65	0.75	0.05	0.05
scores_research	0.85	1.0	0.75	0.65	0.25	0.15	-0.05	0.75	0.85	0.15	0.15
scores_citations	0.65	0.75	1.0	0.85	0.35	0.25	0.15	0.85	0.95	0.25	0.25
scores_international_outlook	0.55	0.65	0.85	1.0	0.45	0.35	0.25	0.75	0.85	0.35	0.35
member_level	0.15	0.25	0.35	0.45	1.0	0.15	0.05	0.15	0.25	0.05	0.05
stats_number_students	0.05	0.15	0.25	0.35	0.15	1.0	0.15	0.05	0.15	0.05	0.05
stats_student_staff_ratio	-0.15	-0.05	0.15	0.25	0.05	0.15	1.0	0.15	0.05	0.05	0.05
stats_pc_intl_students	0.65	0.75	0.85	0.75	0.15	0.05	0.15	1.0	0.75	0.15	0.15
overall_score	0.75	0.85	0.95	0.85	0.25	0.15	0.05	0.75	1.0	0.15	0.15
female_ratio	0.05	0.15	0.25	0.35	0.05	0.05	0.05	0.15	0.15	1.0	0.85
male_ratio	0.05	0.15	0.25	0.35	0.05	0.05	0.05	0.15	0.15	0.85	1.0



# Key Findings and Analysis

**Evaluation:** It was determined that Random Forest had the lowest Mean Squared Error (MSE) and the highest R<sup>2</sup> score, thus being the most effective model.

- For Linear Regression:
  - The Mean Squared Error (MSE) is: 4.82
  - The R<sup>2</sup> Score is: 0.98
- For Decision Tree:
  - The Mean Squared Error (MSE) is: 25.69
  - The R<sup>2</sup> Score is: 0.91
- For Random Forest:
  - The Mean Squared Error (MSE) is: 4.04
  - The R<sup>2</sup> Score is: 0.98



# Critical Evaluation and Reflection

**Model Evaluation:** Predicting university rankings with the Random Forest model was highly accurate, reflecting the success of this project.

**Reflection:** An important aspect from the project that feature selection and the balance between the complexity of the model and its interpretation were important considerations.

**Real-World Application:** The findings can be used to guide universities in improving their ranking metrics.

**Personal Learning:** Developed an understanding of data science applications in educational contexts, and improved proficiency in machine learning techniques.

# Conclusion and Ethical Considerations

**Conclusion:** Data science is successfully demonstrated by the analysis, with Random Forest providing the most accurate predictions.

Scores Research: (importance 34.93%); Scores Citations: (importance 30.10%); Scores Teaching: (importance 17.60%); Scores International Outlook: (importance 10.57%). These features have the highest impact on the overall score.

**Ethical Considerations:** Ensured that data was used responsibly, that privacy was preserved and that fairness was maintained in the interpretation and analysis of data.

