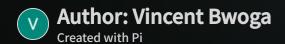


# Predicting Student Dropout Rates for Early Intervention



# CONTENTS

- 1. Objectives
- 3. Key Performance Indicator (KPI)
- 5. Model Development
- 7. Case Study Application
- 9. Model Development
- 11. Optimization
- 13. Reflection & Workflow Diagram

- 2. Stakeholders
- 4. Data Collection & Preprocessing
- 6. Evaluation & Deployment
- 8. Data Strategy
- 10. Deployment
- 12. Critical Thinking

Objectives



# **Objective Details**

- 1 Identify students at high risk of dropping out by the end of the first semester.
- 2 Provide actionable insights to academic advisors for targeted interventions.
- 3 Improve overall student retention rates by 10% within two years.

# Stakeholders

## Stakeholder List



University administration (to allocate resources and monitor program success).



Academic advisors (to implement interventions based on predictions).

**Key Performance Indicator (KPI)** 

# **KPI Details**

Retention Rate Improvement: Percentage increase in student retention after implementing interventions, targeting a 10% improvement.



# Data Collection & Preprocessing

### **Data Sources**



Student Information System (SIS): Academic records, grades, and enrollment status.



Learning Management System (LMS): Engagement metrics like assignment submissions and login frequency.



## **Potential Bias**

Historical academic records may reflect socioeconomic biases, as students from lower-income backgrounds may have lower grades due to external factors (e.g., part-time work), skewing risk predictions.

# **Preprocessing Steps**

- 1 Handling Missing Data:
  Impute missing grades or
  attendance records using
  median values for numerical
  data or mode for categorical
  data to maintain dataset
  integrity.
- Normalization: Scale numerical features (e.g., GPA, attendance) to a 0-1 range to ensure equal weighting in model training.
- Feature Encoding: Convert categorical variables (e.g., major, enrollment status) into numerical formats using one-hot encoding to make them model-compatible.

Model Development

# Model Development



#### **Model Choice**

- Random Forest
- Justification: Random Forest handles non-linear relationships and mixed data types (numerical and categorical) well, is robust to outliers, and provides feature importance for interpretability, which is valuable for understanding dropout factors.

#### **Data Splitting**

• Split data into 70% training, 15% validation, and 15% test sets. The training set trains the model, the validation set tunes hyperparameters, and the test set evaluates final performance to ensure unbiased assessment.



#### Hyper parameters to Tune

Number of Trees
 (n\_estimators): Increasing