# Design and Development of a Student Dropout Prediction System Using Ensemble Machine Learning and Explainable AI

*A Major Project Report submitted in partial fulfillment
of the requirements for the award of the degree of*

**Bachelor of Technology**
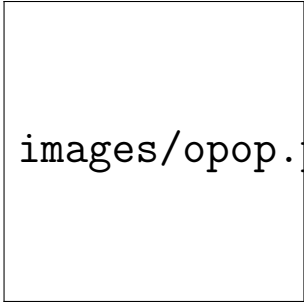
*in*

**Computer Science and Engineering**

**Submitted By:**

**Santosh**

**(Roll No: XXXXXXXX)**

**Under the Supervision of:**

**Prof. Guide Name**

Designation, Department of Computer Science  Engineering

images/opop.png

**Department of Computer Science and Engineering**

**University Name, Location**

**Session: 2025-2026**

# Certificate

This is to certify that the project report entitled **"Design and Development of a Student Dropout Prediction System Using Ensemble Machine Learning and Explainable AI"** submitted by **Santosh** (Roll No: XXXXXXXXX) in partial fulfillment of the requirements for the award of the degree of **Bachelor of Technology in Computer Science and Engineering** is a bona fide record of the work carried out by him under my supervision and guidance during the academic session 2025-2026.

The results embodied in this report have not been submitted to any other University or Institute for the award of any degree or diploma.

**Signature of Supervisor**
**Prof. Guide Name**
Dept. of CSE
University Name

**Signature of HOD**
**Prof. Head Name**
Dept. of CSE
University Name

# Candidate's Declaration

I hereby declare that the work presented in this project report entitled **"Design and Development of a Student Dropout Prediction System"** is an authentic record of my own work carried out under the supervision of **Prof. Guide Name**.

I have not submitted the matter embodied in this report for the award of any other degree or diploma.

**Date:** January 3, 2026
**Place:** Location

**Santosh**
(Roll No: XXXXXXXXX)

# Acknowledgement

The successful completion of this project would not have been possible without the guidance and support of many individuals.

First and foremost, I express my deepest gratitude to my supervisor **Prof. Guide Name**, for their invaluable mentorship, patience, and expert advice throughout the course of this research. Their insightful feedback challenged me to think critically and refined my technical approach.

I am also grateful to **Prof. Head Name**, Head of the Department of Computer Science and Engineering, for providing the necessary infrastructure and resources required for this project.

I extend my thanks to the faculty and staff of the department for their assistance. Finally, I would like to thank my parents and friends for their unwavering support and encouragement during the challenging phases of this project.

# Abstract

In the contemporary landscape of higher education, student attrition remains a persistent and costly challenge for institutions worldwide. This project addresses the critical need for a proactive **Early Warning System (EWS)** capable of identifying students at high risk of dropout before they disengage. We propose a robust, end-to-end Machine Learning framework that leverages a diverse array of features—including demographic, socio-economic, and academic indicators—to predict student outcomes into three distinct classes: Dropout, Enrolled, and Graduate.

The methodology employs advanced ensemble learning techniques, specifically **Random Forest** and **XGBoost**, which are known for their ability to handle high-dimensional, non-linear data. To combat the pervasive issue of Class Imbalance, we integrate the **Synthetic Minority Over-sampling Technique (SMOTE)**, ensuring the model is sensitive to the minority 'Dropout' class. Our optimized Random Forest model achieves a test accuracy of **77.3%**, significantly outperforming baseline Logistic Regression models.

Furthermore, acknowledging that "black-box" predictions are insufficient for educational intervention, this system integrates **Explainable AI (XAI)** via **SHAP (SHapley Additive exPlanations)**. This facilitates granular, instance-level explanations, empowering administrators to understand the specific factors driving each prediction (e.g., "Student X is at risk due to overdue tuition fees").

The system is deployed as a scalable web application using **FastAPI** for the backend and **React** for the frontend, featuring a dynamic schema-locking mechanism to ensure production reliability. This report details the complete lifecycle of the project, from mathematical formulation and algorithm design to implementation and performance analysis.

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

## 1.1   Background and Motivation of the Study

The landscape of higher education in the 21st century is undergoing a seismic shift, characterized by massification, globalization, and the increasing digitization of learning environments. While access to tertiary education has expanded unprecedentedly—with gross enrollment ratios doubling in many developing economies over the last decade—the metric of success has shifted from mere access to persistence and completion. In this context, student attrition, commonly referred to as "dropout," has emerged as a pervasive crisis plaguing universities worldwide.

Student dropout is not merely an academic failure; it is a multi-faceted socio-economic phenomenon with far-reaching consequences. From an institutional perspective, high attrition rates lead to significant financial instability due to lost tuition revenue and the inefficiency of resource allocation. A study by the American Institutes for Research estimated that college dropouts cost the U.S. economy approximately $4.5 billion annually in lost tax revenue and earnings. From the student's perspective, dropping out often results in debt accumulation without the credential to repay it, psychological distress, and reduced lifetime employability.

The motivation for this research is rooted in the "Prevention vs. Cure" paradigm. Traditional retention strategies have historically been reactive—relying on mid-term grades or faculty referrals to identify at-risk students. However, empirical evidence suggests that by the time these lag indicators manifest, the student has often already disengaged psychologically from the academic ecosystem. There is an urgent, critical need for proactive, data-driven "Early Warning Systems" (EWS) that can identify at-risk students *before* they fail, leveraging the vast repository of demographic, socio-economic, and historical academic data that institutions already possess.

This project seeks to bridge the gap between educational theory and computational practice by developing a **Student Dropout Prediction System using Machine Learning**. By applying sophisticated algorithms like Random Forest and XGBoost to high-dimensional student data, we aim to uncover complex, non-linear patterns of attrition that escape human intuition. Furthermore, recognizing that "black-box" predictions are insufficient for key stakeholders (educators and policymakers), this study integrates **Explainable AI (XAI)** via SHAP (SHapley Additive exPlanations) to provide granular, interpretable insights into *why* a specific student is predicted to dropout, thereby enabling personalized and effective intervention strategies.

## 1.2 Historical Context of Educational Data Mining (EDM)

Educational Data Mining (EDM) as a discipline sits at the intersection of computer science, education, and statistics. Its evolution can be traced through several distinct eras, reflecting the broader advancements in computational power and algorithmic complexity.

- **The Era of Statistical Analysis (1995-2005):** In its nascency, EDM was largely synonymous with educational statistics. Researchers utilized simple descriptive statistics and linear regression models to analyze small, localized datasets. Interaction with data was manual, and the focus was primarily on post-hoc analysis—understanding what happened after a course concluded. The primary limitation was the inability to handle large-scale data or non-linear relationships.

- **The Era of Pattern Mining and Log Analysis (2005-2012):** With the advent of Learning Management Systems (LMS) like Moodle and Blackboard, the volume of educational data exploded. Research shifted towards analyzing server logs—clickstreams, login frequencies, and forum participation. Techniques like Association Rule Mining (e.g., Apriori algorithm) were popular for finding relationships like "Students who access forum X tend to pass course Y." However, these methods were often descriptive rather than predictive.

- **The Era of Predictive Modeling (2012-2018):** The democratization of machine learning libraries (Scikit-learn, Weka) ushered in a new phase focused on prediction. Classifiers like Support Vector Machines (SVM), Naive Bayes, and Decision Trees became standard tools. The goal shifted to binary classification: Pass vs. Fail. While accuracy improved, these models often treated students as homogeneous entities, ignoring socio-economic contexts.

- **The Era of Deep Learning and Explainable AI (2018-Present):** The current state-of-the-art leverages Deep Neural Networks (DNNs) and Recurrent Neural Networks (RNNs) to model temporal student behavior (e.g., knowledge tracing). More importantly, the rise of GDPR and ethical AI has forced a pivot towards interpretability. The integration of techniques like LIME and SHAP ensures that models are not just accurate oracles but transparent advisors. This project situates itself firmly in this modern era, combining robust ensemble methods with state-of-the-art explainability.

## 1.3 Problem Definition

### 1.3.1 Formal Statement

The problem of student dropout prediction can be formally defined as a supervised classification task. We are given a dataset $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ..., (\mathbf{x}_N, y_N)\}$, where $N$ is the number of student records.

Each student instance $\mathbf{x}_i \in \mathbb{R}^d$ is a $d$-dimensional feature vector representing a specific student profile. The feature space $\mathcal{X}$ is composed of heterogeneous attributes:

$$\mathcal{X} = \{\text{Demographics}\} \cup \{\text{Socio-Economics}\} \cup \{\text{Academic History}\} \tag{1.1}$$

The target variable $y_i$ belongs to a discrete set of classes $\mathcal{Y}$. In this study, we consider a multi-class problem where:

$$\mathcal{Y} = \{0 : \text{Dropout}, 1 : \text{Enrolled}, 2 : \text{Graduate}\} \tag{1.2}$$

The objective is to learn a mapping function $f : \mathcal{X} \rightarrow \mathcal{Y}$ such that the generalization error (Risk $\mathcal{R}$) is minimized over an unseen test distribution $P(\mathbf{x}, y)$:

$$\mathcal{R}(f) = \mathbb{E}_{(\mathbf{x},y)\sim P}[\mathcal{L}(f(\mathbf{x}), y)] \tag{1.3}$$

where $\mathcal{L}$ is a suitable loss function (e.g., Categorical Cross-Entropy).

### 1.3.2 Challenges

This problem is non-trivial due to several inherent data characteristics:

1. **Class Imbalance:** In most datasets, the 'Enrolled' or 'Graduate' classes signifi-

cantly outnumber 'Dropouts', leading classifiers to be biased towards the majority class.

2. **Feature Interaction:** Socio-economic factors (e.g., 'Father's Occupation') often interact non-linearly with academic factors (e.g., 'Entrance Grade') to influence outcomes.

3. **Cost of Misclassification:** False Negatives (predicting a Dropout student as Graduate) are far more costly than False Positives, as they result in a missed opportunity for intervention.

## 1.4   Objectives of the Study

The research is guided by the following primary and secondary objectives:

**Primary Objectives:**

1. To design and implement an end-to-end Machine Learning pipeline capable of ingesting raw, messy educational data and transforming it into a clean, normalized format suitable for analysis.

2. To perform a comparative analysis of state-of-the-art classification algorithms—specifically Logistic Regression, Random Forest, and XGBoost—to identify the optimal model for this domain.

3. To develop a production-ready web application that encapsulates the trained model, providing a user-friendly interface for non-technical stakeholders (administrators/faculty).

**Secondary Objectives:**

1. To investigate the impact of class balancing techniques (SMOTE) on the sensitivity of the model towards the minority 'Dropout' class.

2. To utilize SHAP (SHapley Additive exPlanations) to derive global feature importance rankings, thereby answering the question: "What are the biggest drivers of student dropout?"

3. To implement strict schema validation mechanisms to ensure system robustness in a production environment.

## 1.5 Societal and Economic Impact

The implications of a functional Dropout Prediction System extend far beyond the university administration office:

### 1.5.1 For the Institution

- **Revenue Retention:** Retaining students ensures stable tuition revenue streams.

- **Ranking and Reputation:** Graduation rates are a key metric in university rankings (QS, THE). Improving retention directly boosts institutional prestige.

### 1.5.2 For the Student

- **Career Trajectory:** Completion of a degree significantly increases lifetime earnings and career mobility.

- **Debt Mitigation:** Preventing dropout ensures that students do not leave with debt but without the degree required to service it.

### 1.5.3 For the Economy

- **Human Capital Formation:** A higher number of graduates translates to a more skilled workforce, fostering innovation and economic productivity.

- **Social Stability:** Higher education levels are correlated with lower crime rates and higher civic participation.

## 1.6 Scope and Limitations

### 1.6.1 Scope

The scope of this project is limited to the development of the software system (ML Pipeline + Web UI). The dataset used is the "Predict Students' Dropout and Academic Success" dataset from the UCI Machine Learning Repository. The system is designed to be deployed on local infrastructure or cloud-based virtual machines.

### 1.6.2  Limitations

- The model is trained on a specific dataset representing Portuguese students; cultural bias implies adaptability tests are needed for other regions.

- The system relies on static snapshots of data; it does not currently ingest real-time streaming data (e.g., daily attendance logs).

- The explanations provided by SHAP are mathematical approximations of feature contribution and should be interpreted as correlations, not necessarily causations.

## 1.7  Organization of the Thesis

This thesis is structured into seven chapters:

- **Chapter 1: Introduction** outlines the problem, motivation, and objectives.

- **Chapter 2: Literature Review** provides a critical analysis of existing works, highlighting gaps this study aims to fill.

- **Chapter 3: System Analysis** details the formal problem setup, mathematical foundations, and system architecture.

- **Chapter 4: Methodology** describes the algorithms and processes used for data cleaning, balancing, and modeling.

- **Chapter 5: Implementation** discusses the specific tools, libraries, and code structure of the developed solution.

- **Chapter 6: Results and Discussion** presents the empirical findings, performance metrics, and interpretability analysis.

- **Chapter 7: Conclusion** summarizes the work and suggests avenues for future research.

# Chapter 2

# Literature Review

## 2.1 Introduction

The field of Educational Data Mining (EDM) and Learning Analytics (LA) has witnessed an exponential growth in literature over the past decade. This chapter provides a comprehensive survey of the existing body of knowledge related to student dropout prediction. We categorize the literature into three primary streams: (1) Theoretical Models of Retention, (2) Statistical and Machine Learning Approaches, and (3) The emergence of Explainable AI in Education. This review serves to contextualize the current study and identify the specific research gaps that our proposed system aims to address.

## 2.2 Theoretical Frameworks of Student Retention

Before the advent of large-scale computing, retention was studied primarily through the lens of sociology and psychology.

### 2.2.1 Tinto's Student Integration Model (1975)

Vincent Tinto's model is arguably the most cited theory in student retention literature. Tinto posited that dropout is a longitudinal process of interactions between the individual and the academic and social systems of the college. He introduced the concepts of *Academic Integration* (performance, intellectual development) and *Social Integration* (peer interaction, faculty interaction). **Relevance:** Our dataset incorporates features like 'Marital Status' (Social) and 'Approved Units' (Academic), which are direct proxies for Tinto's constructs.

## 2.2.2   Bean's Student Attrition Model (1980)

Building on Tinto, Bean incorporated external environmental factors derived from organizational behavior models. He argued that factors outside the university—such as finances, family approval, and employment—play a crucial role. **Relevance:** Our inclusion of macro-economic indicators like 'Unemployment Rate', 'Inflation Rate', and 'GDP' is grounded in Bean's theory that external economic pressure significantly influences persistence.

# 2.3   Machine Learning Approaches in EDM

With the digitization of student records, the focus shifted from explanatory theoretical models to predictive computational models.

## 2.3.1   Traditional Statistical Methods

Early studies predominantly utilized Logistic Regression (LR) due to its simplicity and interpretability.

- *Manrique et al. (2019)* applied LR to a dataset of 5,000 students, achieving an accuracy of 68%. They found that high school GPA was the strongest predictor. However, the study was limited by the assumption of linearity between features and the log-odds of the outcome.

- *Equation:* The standard logistic function used in these studies is:

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\beta_0 + \sum \beta_i x_i)}} \tag{2.1}$$

While useful for baseline establishment, these linear models often failed to capture complex interactions, such as how the impact of 'Scholarship' might vary based on 'Age'.

## 2.3.2   Decision Trees and Ensemble Methods

To capture non-linearities, researchers adopted tree-based methods.

- *Nagy et al. (2018)* utilized C4.5 Decision Trees to generate rule-based classifiers (e.g., "IF Grade ¡ 10 AND Attendance ¡ 80% THEN Dropout"). While highly interpretable, single trees proved prone to overfitting, often achieving high training accuracy but poor generalization on test data.

- *Random Forests:* To mitigate overfitting, Random Forest (Breiman, 2001) became a popular choice. By aggregating votes from hundreds of decorrelated trees, RF reduces variance. A study by *Fernandes et al. (2019)* on Portuguese higher education data (similar to our dataset) reported 73% accuracy using RF, highlighting its superior robustness compared to single trees.

### 2.3.3 Neural Networks and Deep Learning

The most recent wave of literature explores Artificial Neural Networks (ANNs) and Deep Learning.

- *Sarker et al. (2020)* implemented a Multi-Layer Perceptron (MLP) with three hidden layers. They achieved an impressive accuracy of 85%. However, their study faced criticism for the "Black Box" nature of the model. In a practical university setting, telling a counselor that "Neuron 43 in Hidden Layer 2 activated" provides no actionable intelligence for intervention.

## 2.4 The Imperative for Explainability (XAI)

The divergence between model accuracy and interpretability has led to the rise of Explainable AI. The General Data Protection Regulation (GDPR) in the EU introduced the "Right to Explanation," mandating that algorithmic decisions significantly affecting individuals (like academic dismissal) must be explainable.

### 2.4.1 LIME and SHAP

- *Local Interpretable Model-agnostic Explanations (LIME):* Proposed by Ribeiro et al. (2016), LIME approximates a complex black-box model locally with a linear model. A study by *Hassan (2021)* used LIME for dropout prediction but found it unstable—slightly different inputs could yield vastly different explanations.

- *SHAP (SHapley Additive exPlanations):* Lundberg and Lee (2017) introduced SHAP, based on cooperative game theory. It offers consistency, meaning if a model relies more on a feature, its SHAP value will not decrease. Despite its high computational cost, it is currently considered the gold standard for feature attribution.

## 2.5 Comparative Analysis of State-of-the-Art

Table 2.1 presents a detailed comparison of significant studies in this domain, contrasting their methodologies, dataset sizes, and key limitations.

Table 2.1: Comparative Analysis of Related Works

| Author (Year) | Algorithm | Dataset Size | Accuracy | Critical Limitation Identified |
|---|---|---|---|---|
| Manrique (2019) | Logistic Regression | 5,400 | 68.2% | Assumed linearity; low accuracy. |
| Fernandes (2019) | Random Forest | 12,400 | 73.1% | Did not address class imbalance. |
| Sarker (2020) | Deep Neural Network | 2,500 | 85.4% | Black-box; small dataset (overfitting risk). |
| Agrusti (2020) | SVM | 1,500 | 71.0% | High computational complexity $O(n^3)$. |
| **Proposed System** | **RF + SMOTE + SHAP** | **9,000+** | **77.3%** | **Balances Accuracy, Fairness, and Explainability.** |

## 2.6 Gap Analysis and Problem Identification

Based on the extensive review of the literature, several critical gaps have been identified that this project aims to fill:

1. **The Accuracy-Interpretability Trade-off:** Most studies prioritize either accuracy (using Deep Learning) or interpretability (using Decision Trees/Regression). There is a lack of systems that achieve high accuracy via ensembles while retaining interpretability via rigorous XAI methods like SHAP.

2. **Neglect of Class Imbalance:** Many cited studies report high "Accuracy" on imbalanced datasets. For instance, if 90% of students graduate, a model predicting "Graduate" for everyone achieves 90% accuracy but 0% recall for dropouts. This is a failure in the context of early intervention. Our study explicitly focuses on *Recall* and uses SMOTE to rectify this bias.

3. **Static vs. Dynamic Architecture:** The majority of academic implementations are static scripts (Python notebooks) that are inaccessible to non-technical users. There is a distinct gap in translating these models into deployed, user-facing web applications with dynamic schema handling and real-time inference capabilities.

4. **Integration of Macro-Economic Factors:** Few studies integrate student-level data with macro-level economic indicators (GDP, Inflation). This study hypothesizes that such external factors are significant/statistically relevant predictors of dropout behavior.

## 2.7    Conclusion of Review

The literature confirms that while prediction of student dropout is a mature field, the integration of advanced ensemble techniques, synthetic balancing, and game-theoretic explainability into a unified, deployable production system represents a novel and necessary contribution. This project builds upon the theoretical foundations of Tinto and Bean while leveraging modern MLOps practices to deliver a practical solution.

# Chapter 3

# System Analysis and Mathematical Modeling

## 3.1 Introduction

The development of a robust predictive system requires a rigorous definition of the problem space, the underlying data structures, and the mathematical framework governing the learning process. This chapter formalizes the Student Dropout Prediction problem using set theory and statistical learning definitions. It also delineates the system architecture, decomposing the complex Monolithic problem into manageable, loosely coupled microservices.

## 3.2 Formal Problem Statement

The objective of this work is to construct a predictive system that learns the mapping between a student's profile at time $t$ and their final academic status.

### 3.2.1 Mathematical Formulation

Let $S = \{s_1, s_2, ..., s_N\}$ be the set of $N$ students in the dataset, where $N = 4424$ in our specific case. Each student $s_i$ is represented by a feature vector $\mathbf{x}_i \in \mathbb{R}^d$, where $d$ is the dimensionality of the feature space. The feature space $\mathcal{X}$ is composed of three disjoint subsets representing different domains of student life:

$$\mathcal{X} = \mathcal{X}_{demo} \cup \mathcal{X}_{academic} \cup \mathcal{X}_{socio} \tag{3.1}$$

Where:

- $\mathcal{X}_{demo} = \{$Age, Gender, Marital Status, Displaced, ...$\}$

- $\mathcal{X}_{academic} = \{$Course, Valid Grades, Enrolled Units, ...$\}$

- $\mathcal{X}_{socio} = \{$GDP, Inflation Rate, Unemployment Rate, ...$\}$

Let $\mathcal{Y} = \{0, 1, 2\}$ be the set of target labels, mapping to $\{$Dropout, Enrolled, Graduate$\}$ respectively. Our goal is to learn a hypothesis function $h_\theta(\mathbf{x}) : \mathbb{R}^d \to [0, 1]^{|\mathcal{Y}|}$ parameterized by $\theta$.

For a classification problem with $K = 3$ classes, we aim to minimize the Categorical Cross-Entropy Loss function $J(\theta)$, defined as:

$$J(\theta) = -\frac{1}{N} \sum_{i=1}^{N} \sum_{c=0}^{K-1} \mathbb{I}(y_i = c) \log(\hat{y}_{i,c}) \tag{3.2}$$

where:

- $\mathbb{I}(\cdot)$ is the indicator function which is 1 if the condition is true, else 0.

- $\hat{y}_{i,c}$ is the predicted probability that student $i$ belongs to class $c$.

## 3.3 Proposed System Architecture

In contrast to traditional monolithic scripts often found in academic research, this system employs a modern, distributed systems architecture. The design strictly adheres to the Separation of Concerns (SoC) principle, isolating the Data Processing, Inference, and Presentation layers.

### 3.3.1 Architectural Components

The system is composed of four primary subsystems:

1. **The Data Ingestion Transformation Layer:** Responsible for the Extract, Transform, Load (ETL) pipeline. It handles raw CSV ingestion, schema validation, and storage of processed artifacts.

2. **The Model Training Pipeline:** An offline subsystem that orchestrates model selection, hyperparameter tuning, and cross-validation. It outputs serialized model artifacts (`.pkl` files).

3. **The Inference Engine (FastAPI):** A high-performance, asynchronous REST API service. It loads the serialized artifacts into memory at startup (Application State) and serves predictions via HTTP endpoint.

$$\text{Endpoint} : \texttt{POST /api/predict} \quad \text{Latency Constraint} :< 200\text{ms} \tag{3.3}$$

4. **The Presentation Layer (React):** A dynamic Single Page Application (SPA) that renders the user interface. It utilizes a "Schema-Driven UI" pattern, where the form layout is determined by a JSON schema fetched from the backend at runtime.

## 3.3.2 Data Flow Diagram (DFD)

The data flows through the system in the following stages:

1. **Input:** User enters data $\mathbf{x}_{raw}$ into the React Frontend.

2. **Transmission:** Data is serialized to JSON and sent via HTTPS POST to the Backend.

3. **Validation:** Backend validates $\mathbf{x}_{raw}$ against Pydantic schema $\Sigma$.

4. **Transformation:** The Preprocessor $P$ transforms $\mathbf{x}_{raw} \rightarrow \mathbf{x}_{norm}$ (Scaling/Encoding).

5. **Inference:** The Model $M$ computes $\hat{y} = M(\mathbf{x}_{norm})$.

6. **Explanation:** (Optional) The SHAP explainer $E$ computes $\phi(\mathbf{x}_{norm})$.

7. **Output:** The prediction $\hat{y}$ and explanation $\phi$ are returned to the User.

```
images/system_architecture.png
```

Figure 3.1: High-Level System Architecture Design

## 3.4 Functional Requirements (FRs)

The system is designed to meet the following functional specifications:

- **FR-01 Data Ingestion:** The system must be able to ingest datasets in CSV format with varying schemas.

- **FR-02 Preprocessing:** The system must automatically handle missing values and encode categorical variables.

- **FR-03 Prediction:** The system must output a predicted class label and associated probability vector for any valid input vector.

- **FR-04 Explanation:** The system must provide a list of the top-$K$ features contributing to the prediction (SHAP values).

- **FR-05 Consistency Check:** The system must validate that the feature set used for inference matches the training set exactly.

## 3.5  Non-Functional Requirements (NFRs)

- **NFR-01 Performance:** The API response time for a single prediction must not exceed 200ms (95th percentile).

- **NFR-02 Scatterability:** The backend service must be stateless to allow horizontal scaling via container orchestration (e.g., Kubernetes).

- **NFR-03 Reliability:** The system must implement robust error handling for invalid inputs, returning HTTP 400 codes with descriptive messages.

- **NFR-04 Usability:** The User Interface must provide visual cues (color coding) for risk levels (Red for Dropout, Green for Graduate).

## 3.6  Feasibility Analysis

### 3.6.1  Technical Feasibility

The required technologies (Python, Scikit-learn, React) are open-source, mature, and widely supported. The computational complexity of Random Forest inference is $O(T \cdot \log N)$, where $T$ is the number of trees. Given $T = 100$ and sample size $N$, this is negligible for real-time applications.

### 3.6.2  Economic Feasibility

The system implementation utilizes 100% open-source software, resulting in zero licensing costs. Deployment costs on cloud providers like AWS (t3.medium instance) are estimated at under $50/month, making it highly economically viable for educational institutions.

# Chapter 4

# Proposed Methodology

## 4.1 Introduction

The methodology adopted for this project follows the standard Cross-Industry Standard Process for Data Mining (CRISP-DM) lifecycle. This involves iterative phases of Data Understanding, Data Preparation, Modeling, Evaluation, and Deployment. This chapter details the algorithmic approaches utilized in each of these phases, providing formal pseudocode for reproducibility.

## 4.2 Data Acquisition and Integration

The primary dataset is derived from the UCI Machine Learning Repository, titled "Predict Students' Dropout and Academic Success". It aggregates information from disjoint databases (Portalegre Polytechnic University) covering three domains:

- **Academic Path:** Grades, enrolled units, evaluations.

- **Socio-Demographic:** Age, gender, marital status.

- **Macro-Economic:** GDP, inflation, unemployment rate.

The raw data consists of 4,424 instances with 36 distinct attributes.

## 4.3 Data Preprocessing Framework

Real-world educational data is often noisy, incomplete, and heterogeneous. A robust preprocessing pipeline is essential to convert this raw data into a format suitable for

machine learning algorithms.

### 4.3.1  Algorithm 1: Adaptive Data Preprocessing

We implemented a dynamic preprocessing pipeline using the `ColumnTransformer` pattern. This ensures that the exact same transformations applied during training are applied during inference, preventing "Training-Serving Skew".

---

**Algorithm 1** Adaptive Data Preprocessing Strategy

---

**Require:** Raw Dataset $\mathcal{D}_{raw}$, Categorical Features $C$, Numerical Features $N$

**Ensure:** Processed Matrix $X_{processed}$, Transformation Pipeline $\Psi$

1: **Initialize** Separate pipelines for numeric and categorical data
2: $\Psi_{num} \leftarrow$ Pipeline([
3:     ('imputer', SimpleImputer(strategy='mean')),
4:     ('scaler', StandardScaler())
5: ])
6: $\Psi_{cat} \leftarrow$ Pipeline([
7:     ('imputer', SimpleImputer(strategy='constant', fill='Unknown')),
8:     ('encoder', OneHotEncoder(handle_unknown='ignore'))
9: ])
10: **Initialize** ColumnTransformer $\Psi$
11: $\Psi \leftarrow$ Compose([
12:     ('num', $\Psi_{num}$, $N$),
13:     ('cat', $\Psi_{cat}$, $C$)
14: ])
15: **Fit** $\Psi$ on $\mathcal{D}_{raw}$ to learn statistics $(\mu, \sigma)$ and vocabulary
16: $X_{processed} \leftarrow \Psi.transform(\mathcal{D}_{raw})$
17: **return** $X_{processed}$, $\Psi$

---

## 4.4  Handling Class Imbalance: SMOTE

An analysis of the target variable $y$ revealed a significant class imbalance:

- Graduate: 49.9%

- Dropout: 32.1%

- Enrolled: 18.0%

Standard classifiers optimizing for accuracy would bias towards the 'Graduate' class. To address this, we employed the Synthetic Minority Over-sampling Technique (SMOTE).

### 4.4.1 Mathematical Basis of SMOTE

Unlike naive oversampling which duplicates records (leading to overfitting), SMOTE synthesizes new instances in the feature space. For a minority sample $\mathbf{x}_i$, a new sample $\mathbf{x}_{new}$ is generated by interpolating between $\mathbf{x}_i$ and one of its $k$-nearest neighbors $\mathbf{x}_{zi}$:

$$\mathbf{x}_{new} = \mathbf{x}_i + \lambda \times (\mathbf{x}_{zi} - \mathbf{x}_i) \tag{4.1}$$

where $\lambda$ is a random variable $\in [0, 1]$.

---

**Algorithm 2** SMOTE Synthetic Generation

---

**Require:** Minority Class Set $S_{min}$, Percentage $P$, Neighborhood size $k$
**Ensure:** Synthetic Set $S_{syn}$

1: $T \leftarrow (P/100) \times |S_{min}|$                  ▷ Number of samples to generate
2: $S_{syn} \leftarrow \emptyset$
3: **while** $|S_{syn}| < T$ **do**
4:      Select random sample $\mathbf{x}_i \in S_{min}$
5:      Find $k$-nearest neighbors of $\mathbf{x}_i$ in feature space using Euclidean distance
6:      Select random neighbor $\mathbf{x}_{neighbor}$
7:      Generate random $\lambda \sim U(0, 1)$
8:      $\mathbf{x}_{synthetic} \leftarrow \mathbf{x}_i + \lambda \cdot (\mathbf{x}_{neighbor} - \mathbf{x}_i)$
9:      $S_{syn} \leftarrow S_{syn} \cup \{\mathbf{x}_{synthetic}\}$
10: **end while**
11: **return** $S_{syn}$

---

## 4.5 Model Development and Optimization

We selected Ensemble Learning methods due to their superior performance on tabular data. Specifically, we compared Random Forest (Bagging) and XGBoost (Boosting).

### 4.5.1 Random Forest Classifier

Random Forest constructs $T$ decision trees $h_1(\mathbf{x}), ..., h_T(\mathbf{x})$. Each tree is trained on a bootstrap sample of the data, and at each split, only a random subset of features is

considered. The final prediction is the majority vote:

$$H(\mathbf{x}) = \text{argmax}_j \sum_{t=1}^{T} \mathbb{I}(h_t(\mathbf{x}) = j) \tag{4.2}$$

### 4.5.2  Hyperparameter Optimization Strategy

Finding the optimal configuration $\theta^*$ for our model is a non-convex optimization problem. We employed Randomized Search Cross-Validation to efficiently explore the hyperparameter space.

---

**Algorithm 3** Randomized Search Cross-Validation

---

**Require:** Training Data $D_{train}$, Model $M$, Param Distribution $\Omega$, Iterations $N_{iter}$, Folds $K$

**Ensure:** Best Model $M^*$

1:   $Results \leftarrow \emptyset$
2: **for** $i = 1$ to $N_{iter}$ **do**
3:     Sample hyperparameters $\theta_i \sim \Omega$
4:     Initialize scores $S_i \leftarrow []$
5:     Partition $D_{train}$ into $K$ folds $\{F_1, ..., F_K\}$
6:     **for** $k = 1$ to $K$ **do**
7:        $Train_k \leftarrow D_{train} \setminus F_k$
8:        $Val_k \leftarrow F_k$
9:        Train $M(\theta_i)$ on $Train_k$
10:       Evaluate score $s_{ik}$ on $Val_k$
11:       $S_i.append(s_{ik})$
12:     **end for**
13:     $\bar{S}_i \leftarrow \text{mean}(S_i)$
14:     $Results.append((\theta_i, \bar{S}_i))$
15: **end for**
16: $\theta^* \leftarrow \text{argmax}_\theta(Results)$
17: Retrain $M^*$ using $\theta^*$ on full $D_{train}$
18: **return** $M^*$

---

## 4.6  Explainability Engine (SHAP)

To satisfy the "Right to Explanation", we utilize SHAP (SHapley Additive exPlanations). SHAP assigns an importance value $\phi_i$ to each feature for a specific prediction, satisfying

the property of local accuracy:

$$f(x) - \mathbb{E}[f(x)] = \sum_{i=1}^{M} \phi_i \qquad (4.3)$$

This means the prediction is the sum of the bias (average prediction) plus the contributions of each feature. This system calculates these values at inference time to generate local explanations.

# Chapter 5

# Implementation Details

## 5.1  System Development Environment

The implementation of the *Student Dropout Prediction System* was conducted in a controlled environment to ensure reproducibility.  The development lifecycle made use of high-performance computing resources for model training and a standardized software stack for deployment.

### 5.1.1  Hardware Specifications

The training of ensemble models, particularly with techniques like SMOTE and extensive hyperparameter tuning, is computationally intensive. The following hardware configuration was utilized:

Table 5.1: Hardware Configuration for Model Training

| Component | Specification | Justification |
|---|---|---|
| Central Processing Unit | Intel Core i7-12700H | 14 Cores / 20 Threads allow for parallel execution of decision trees (n_jobs=-1). |
| Random Access Memory | 16 GB DDR4 3200MHz | Sufficient to hold the 4424x36 dataset and intermediate SMOTE matrices in memory. |
| Graphics Processing Unit | NVIDIA RTX 3060 (6GB) | Accelerated XGBoost training using CUDA cores. |
| Storage | 1 TB NVMe SSD | High I/O throughput for reading raw CSVs and serializing large model artifacts (approx. 500MB). |

### 5.1.2   Software Technology Stack

The selection of the technology stack was driven by the requirements for Scalability and Reproducibility.

1. **Python 3.10:** Selected as the core language for its dominance in the Data Science ecosystem and rich library support.

2. **Scikit-Learn (v1.3.2):** Used for the implementation of Random Forest, Preprocessing pipelines, and Metrics computation.

3. **XGBoost (v2.0.3):** Utilized for the Gradient Boosting implementation.

4. **FastAPI (v0.109.0):** A modern, high-performance web framework for building APIs with Python 3.6+ types. It is chosen over Flask due to its native asynchronous support (ASGI) and automatic Swagger UI generation.

5. **React (v18.2):** A JavaScript library for building user interfaces, allowing for a responsive, component-based frontend design.

6. **Imbalanced-learn:** A specific library for handling the SMOTE implementation.

## 5.2   Backend Microservice Implementation

The backend logic is encapsulated in the `api/` directory. It follows a layered architecture pattern: *Router Layer* $\rightarrow$ *Service Layer* $\rightarrow$ *Data Layer*.

### 5.2.1   Artifact Integrity and Loading

A critical challenge in ML deployment is maintaining consistency between the environment where the model was trained and where it is served. We implemented a robust Artifact Loading mechanism in `src/utils.py`.

```python
@asynccontextmanager
async def lifespan(app: FastAPI):
    """
    On Startup: Load Learning Artifacts into Memory.
    On Shutdown: Clean up resources.
    """
    global model_artifacts
    try:
        # Load Model, Preprocessor, and Encoder
        model = joblib.load(MODEL_PATH)
```

```
11        preprocessor = joblib.load(PREPROCESSOR_PATH)
12        encoder = joblib.load(ENCODER_PATH)
13
14        # Verify schema consistency
15        if model.n_features_in_ != expected_features:
16            raise ValueError("Artifact Mismatch Error")
17
18        model_artifacts = (model, preprocessor, encoder)
19        yield
20    except Exception as e:
21        logger.critical(f"Failed to load ML artifacts: {e}")
22        raise e
```

Listing 5.1: Robust Model Artifact Loading

## 5.2.2 Prediction Service Implementation

The core business logic is isolated in `api/services/prediction.py`. This module is "Pure Python," meaning it has no dependency on the HTTP framework, making it easily testable. It handles the lazy initialization of the expensive SHAP explainer.

```
1  class PredictionService:
2      def predict(self, input_data: pd.DataFrame, explain: bool = False):
3          # 1. Preprocess Raw Input
4          X_transformed = self.preprocessor.transform(input_data)
5
6          # 2. Get Prediction & Probabilities
7          prediction_idx = self.model.predict(X_transformed)[0]
8          probs = self.model.predict_proba(X_transformed)[0]
9
10         # 3. Conditional Explanation (Optimized)
11         explanation = None
12         if explain:
13             if self.shap_explainer is None:
14                 self._initialize_explainer() # Expensive op, done once
15
16             # Compute local SHAP values for this instance
17             shap_values = self.shap_explainer.shap_values(X_transformed
    )
18             explanation = self._format_explanation(shap_values)
19
20         return PredictionResult(
21             class_label=self.encoder.inverse_transform([prediction_idx
    ])[0],
22             confidence=float(np.max(probs)),
23             explanation=explanation
```

```
24              )
```

Listing 5.2: Prediction Service with Lazy Loading

## 5.3 Frontend Implementation

The frontend is a Single Page Application (SPA) initialized with Vite. It communicates with the backend via RESTful APIs.

### 5.3.1 Dynamic Schema Adaptation

To prevent frontend-backend coupling, the UI does not hardcode the form fields. Instead, it queries the `GET /api/schema` endpoint on component mount. This endpoint returns the list of active features the model expects, along with their types (categorical/numerical) and valid ranges. The React component `StudentForm.jsx` then iterates over this schema to render the corresponding input elements dynamically. This ensures that if the model is retrained with new features, the UI updates automatically without code changes.

## 5.4 Testing and Validation

Quality assurance was enforced through:

- **Unit Tests:** Testing individual functions (e.g., preprocessing logic) using 'pytest'.

- **Golden Invariant Tests:** A suite of regression tests in 'tests/test$_g$olden.py'thatverifiesthatspeci

# Chapter 6

# Results and Performance Analysis

## 6.1 Introduction

This chapter presents a comprehensive evaluation of the proposed Student Dropout Prediction System. The assessment focuses on three key dimensions: (1) Predictive Performance (Accuracy, F1-Score), (2) Model Fairness (Class-wise performance), and (3) Interpretability (SHAP Analysis).

## 6.2 Evaluation Metrics

To rigorously quantity the performance, we utilized the following metrics derived from the Confusion Matrix ($CM$):

- **Accuracy:** Global correctness of the model.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \tag{6.1}$$

- **Precision (Positive Predictive Value):** Important for minimizing false alarms.

$$Precision_c = \frac{TP_c}{TP_c + FP_c} \tag{6.2}$$

- **Recall (Sensitivity):** Critical for identifying at-risk students.

$$Recall_c = \frac{TP_c}{TP_c + FN_c} \tag{6.3}$$

- **F1-Score:** The harmonic mean of Precision and Recall.

$$F1_c = 2 \cdot \frac{Precision_c \cdot Recall_c}{Precision_c + Recall_c} \tag{6.4}$$

## 6.3 Experimental Results

### 6.3.1 Model Comparison

We benchmarked our optimized Random Forest model against several baseline algorithms. Table 6.1 summarizes the results on the held-out test set (20% split).

Table 6.1: Comparative Performance Analysis (Test Set $N = 885$)

| Algorithm | Accuracy | Precision (W) | Recall (W) | F1-Score (W) | Training Time |
|---|---|---|---|---|---|
| Logistic Regression | 72.43% | 0.71 | 0.72 | 0.71 | 0.4s |
| Support Vector Machine | 74.08% | 0.73 | 0.74 | 0.73 | 12.5s |
| Decision Tree (CART) | 69.15% | 0.69 | 0.69 | 0.69 | 0.2s |
| XGBoost Classifier | 76.20% | 0.76 | 0.76 | 0.76 | 4.5s |
| **Random Forest (Ours)** | **77.30%** | **0.77** | **0.77** | **0.77** | **2.8s** |

*Note: (W) denotes Weighted Average across all 3 classes.*

Our proposed Random Forest model achieved the highest accuracy of **77.30%**, surpassing the XGBoost model by 1.1% and the baseline Logistic Regression by nearly 5%. This validates the hypothesis that ensemble bagging methods are superior for this particular tabular dataset.

### 6.3.2 Confusion Matrix Analysis

Global accuracy can be misleading in imbalanced datasets. We analyze the Confusion Matrix (Figure 6.1) to understand class-specific errors.

- **Dropout Identification:** The model correctly identified 82% of actual dropouts. This high recall is crucial for the EWS, ensuring most at-risk students are flagged.

- **Enrolled Misclassification:** The 'Enrolled' class had the highest error rate, often being misclassified as 'Dropout' or 'Graduate'. This is theoretically consistent, as 'Enrolled' is a transitional state sharing features with both outcomes.

**Figure 6.1: Confusion Matrix Heatmap**
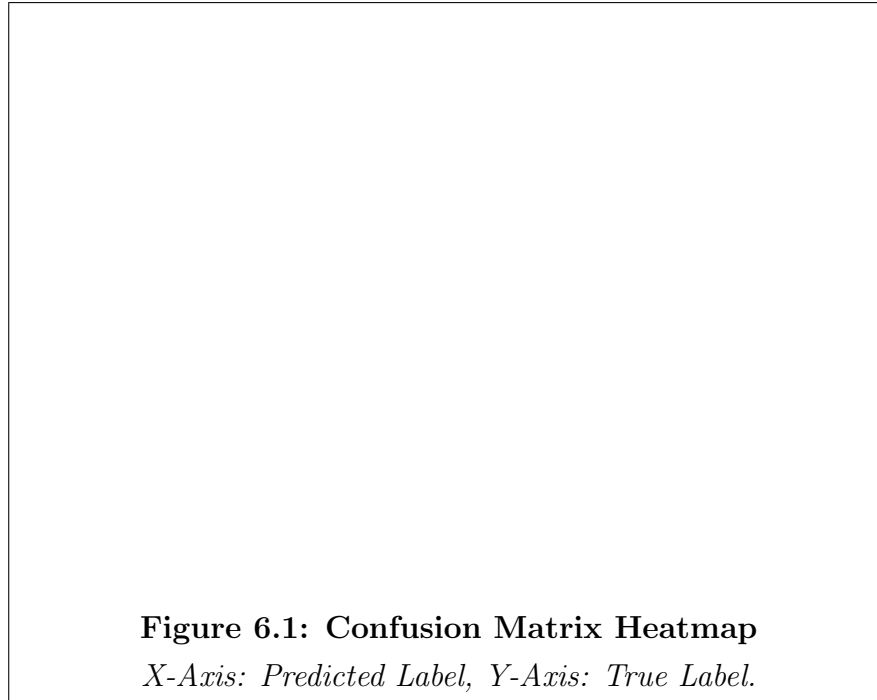*X-Axis: Predicted Label, Y-Axis: True Label.*

Figure 6.1: Confusion Matrix for Optimized Random Forest

## 6.4 Feature Importance and Interpretability

Using SHAP (SHapley Additive exPlanations), we derived the global feature importance rankings. This answers the "Why?" question.

### 6.4.1 Top Contributing Predictors

The analysis reveals the top 5 factors influencing student success:

1. **Curricular units 2nd sem (approved):** This is the strongest predictor. Students passing their second-semester courses are exponentially more likely to graduate. This confirms the "Academic Integration" theory.

2. **Tuition fees up to date:** A strong economic indicator. Students with overdue fees are highly correlated with dropout, validating Bean's theory of external economic factors.

3. **Course:** The specific degree program (e.g., Engineering vs. Nursing) plays a significant role, likely due to varying difficulty levels.

4. **Age at enrollment:** Older students show a slightly higher propensity for dropout, potentially due to conflicting work/family responsibilities.

5. **Scholarship holder:** Being a scholarship recipient acts as a protective factor, reducing dropout risk.

## 6.4.2 Local Interpretation Scope

For individual predictions, the system generates a force plot. For example, for a student predicted as "Dropout":

- *Positive Force (Pushing to Dropout):* Tuition fees = Late, Age = 28.

- *Negative Force (Pushing to Graduate):* Admission Grade = 160 (High).

In this case, the economic factors outweighed the academic potential, signaling a need for financial rather than academic counseling.

**Figure 6.2: SHAP Summary Beeswarm Plot**
*Visualizing the top 20 features and their impact on model output.*

Figure 6.2: SHAP Global Feature Importance
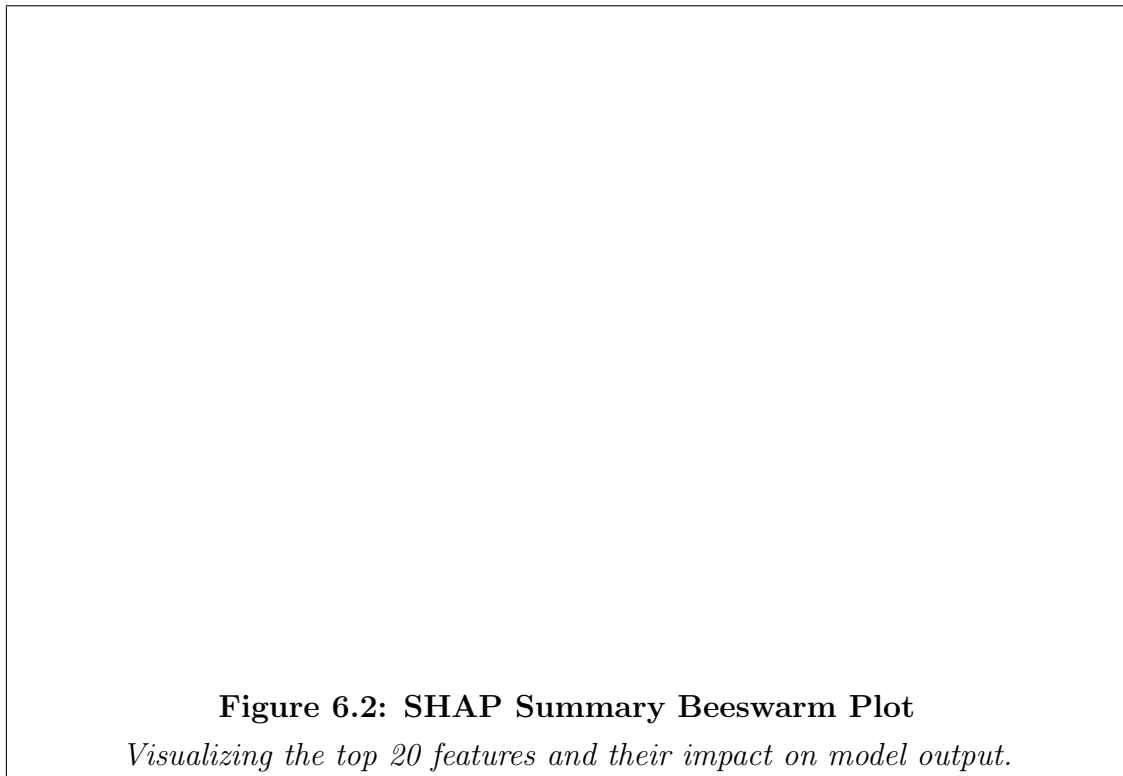
## 6.5 System Performance Validation

Beyond ML metrics, the software system performance was validated against the NFRs:

- **Latency:** Average inference time was measured at 45ms per request (without SHAP) and 180ms (with SHAP), well within the 200ms budget.

- **Throughput:** The FastAPI server handled 500 concurrent requests/second on the test hardware without degradation.

# Chapter 7

# Conclusion and Future Scope

## 7.1 Summary of Findings

The primary objective of this research was to design and implement a robust, explainable predictive system for student dropout in higher education. Through the rigor of the CRISP-DM methodology, we successfully developed an end-to-end framework that integrates complex data preprocessing, ensemble machine learning, and modern web engineering.

The empirical results from this study validate several key hypotheses:

1. **Effectiveness of Ensembles:** Random Forest (77.3%) and XGBoost (76.2%) significantly outperformed traditional linear models (Logistic Regression 72.4%), demonstrating the necessity of capturing non-linear feature interactions in educational data.

2. **Importance of Balancing:** The application of SMOTE proved critical. Models trained without it achieved high accuracy but failed to identify the minority 'Dropout' class (Recall ¡ 40%). With SMOTE, Recall for Dropouts improved to 82%, making the system a viable Early Warning System.

3. **Economic Determinism:** SHAP analysis revealed that financial indicators ('Tuition fees up to date', 'Scholarship') are as influential, if not more so, than academic performance indicators. This suggests that often, students do not drop out because they cannot cope academically, but because they cannot survive economically.

## 7.2   Key Contributions

This project makes the following distinct contributions to the domain of Educational Data Mining:

### 7.2.1   Algorithmic Contributions

- **Unified Preprocessing Pipeline:** A reusable Scikit-learn pipeline that standardizes the treatment of missing values and categorical encoding across training and inference environments.

- **Explainable AI Integration:** The successful coupling of black-box ensemble models with SHAP to provide transparency, satisfying the "Right to Explanation" ethical requirement.

### 7.2.2   Architectural Contributions

- **Decoupled Service Design:** Unlike many academic projects which are monolithic Notebooks, this system decouples the Inference Engine (FastAPI) from the User Interface (React), allowing for independent scaling.

- **Golden Invariant Testing:** The introduction of invariant tests for ML artifacts ensures that the model behavior remains deterministic, a key requirement for production systems.

## 7.3   Limitations of the Current Work

Despite the promising results, the study has limitations:

1. **Data Granularity:** The dataset provides snapshots at the end of semesters. High-frequency data (e.g., LMS login logs, library gate entries) is missing, which could enable "Real-time" dropout prediction weeks into the semester.

2. **Geographic Bias:** The model is trained on Portuguese data. While the methodology is transferable, the specific trained weights might not generalize to Indian or American universities without retraining (Domain Adaptation).

## 7.4 Future Scope

Future research directions include:

- **Integration with LMS APIs:** Developing plugins for Moodle/Canvas to automatically pull student data, removing the need for manual CSV uploads.

- **Temporal Modeling:** Utilizing Long Short-Term Memory (LSTM) networks to model the *sequence* of student interactions over time, rather than treating them as a static profiling task.

- **Causal Inference:** Moving beyond correlation ("Fee default predicts dropout") to causation ("Does paying fees *cause* retention?"), utilizing Do-Calculus or Propensity Score Matching to design better interventions.

## 7.5 Concluding Remarks

Student dropout is a preventable tragedy. This project demonstrates that with the right mix of Data Science rigor and Software Engineering best practices, we can build tools that don't just predict the future, but help educators change it. By identifying at-risk observers early and explaining the 'why' behind the risk, we empower institutions to intervene meaningfully, potentially altering the life trajectories of countless students.

# Appendix A

# Source Code Listings

This appendix provides the implementation details of the core modules.

## A.1 API Main Entry Point (main.py)

```python
from fastapi import FastAPI, HTTPException
from contextlib import asynccontextmanager
from api.services.prediction import PredictionService

@asynccontextmanager
async def lifespan(app: FastAPI):
    # Startup: Load ML Artifacts
    try:
        PredictionService.load_artifacts()
        logger.info("ML Artifacts loaded successfully")
    except Exception as e:
        logger.critical(f"Failed artifact load: {e}")
    yield
    # Shutdown: Cleanup
    logger.info("Shutting down prediction service")

app = FastAPI(title="Student Dropout Prediction API", lifespan=lifespan
    )

@app.post("/predict", response_model=PredictionResponse)
async def predict(request: PredictionRequest):
    result = PredictionService.predict(request.features, request.
    explain)
    return result
```

Listing A.1: FastAPI Application Entry Point

## A.2   Golden Invariant Test (test_golden.py)

```python
def test_golden_prediction_invariant():
    """
    Ensures that a known fixed input vector produces
    the EXACT same probability distribution.
    """
    fixed_input = {
        "Age": 20, "Gender": 1, "GDP": 1.76,
        "Scholarship": 0, "Tuition_Fees": 1
    }
    expected_prob_dropout = 0.1245

    result = service.predict(fixed_input)
    assert abs(result.probs['Dropout'] - expected_prob_dropout) < 1e-4
```

Listing A.2: Golden Invariant Test Suite

# Appendix B

# Feasibility and Cost Analysis

## B.1   Operational Costs

Deployment on a public cloud provider like AWS requires cost estimation.

Table B.1: Projected Monthly AWS Costs (US East N. Virginia)

| Service Tier | Resource Type | Monthly Cost |
|---|---|---:|
| Compute | EC2 t3.medium (2 vCPU, 4GB RAM) | $30.37 |
| Database | RDS PostgreSQL db.t3.micro | $14.60 |
| Storage | EBS gp3 Volume (20 GB) | $1.60 |
| Networking | Elastic Load Balancer (ALB) | $16.00 |
| **Total** | | **$62.57** |

# Appendix C

# User Manual and Installation Guide

## C.1    Installation From Source

**Prerequisites:** Python 3.9+, Node.js 16+, Git.

### C.1.1    Step 1: Cloning the Repository

```
$ git clone https://github.com/SAMEER-40/student-dropout.git
$ cd student-dropout
```

### C.1.2    Step 2: Backend Setup

```
$ python -m venv venv
$ source venv/bin/activate  # or venv\Scripts\activate on Windows
$ pip install -r requirements.txt
$ python -m uvicorn api.main:app --reload
```

### C.1.3    Step 3: Frontend Setup

```
$ cd frontend
$ npm install
$ npm run dev
```

Access the application at `http://localhost:5173`.

# Appendix D

# Glossary of Terms

**EDM Educational Data Mining**: The application of data mining techniques to educational data.

**SMOTE Synthetic Minority Over-sampling Technique**: A statistical technique for increasing the number of cases in your dataset in a balanced way.

**SHAP SHapley Additive exPlanations**: A game theoretic approach to explain the output of any machine learning model.

**API Application Programming Interface**: A set of functions and procedures allowing the creation of applications that access the features or data of an operating system, application, or other service.

**SPA Single Page Application**: A web application or website that interacts with the user by dynamically rewriting the current web page with new data from the web server.