

## AI Development Workflow Assignment Report

**Title:** Predicting Student Dropout Rates Using Machine Learning

**Author:** Caren Rayon

**Date:** July 2025

---

### Part 1: Short Answer Questions (30 points)

#### 1. Problem Definition (6 points)

**Problem:** Predicting the likelihood of a student dropping out before completing their academic program.

**Objectives:**

- Identify at-risk students early.
- Enable targeted interventions by educators.
- Improve overall retention rates.

**Stakeholders:**

- School administrators
- Students and their families

**KPI:** Model accuracy or F1-score in identifying dropouts correctly.

---

#### 2. Data Collection & Preprocessing (8 points)

**Data Sources:**

- Student Information Systems (SIS)
- Learning Management Systems (LMS) like Moodle or Canvas

**Potential Bias:**

- Students from low-income backgrounds may be underrepresented or inaccurately labeled.

**Preprocessing Steps:**

1. Handle missing attendance/grade entries (e.g., fill or remove).
2. Normalize continuous variables like GPA and login hours.

3. Encode categorical variables (e.g., program type).
- 

### 3. Model Development (8 points)

**Chosen Model:** Random Forest – It's interpretable and performs well on tabular data.

**Data Splitting:**

- 70% training
- 15% validation
- 15% testing

**Hyperparameters to Tune:**

- `n_estimators`: Number of trees in the forest
  - `max_depth`: Prevents overfitting by limiting tree growth
- 

### 4. Evaluation & Deployment (8 points)

**Evaluation Metrics:**

- Accuracy: Overall correct predictions.
- F1-Score: Balance between precision and recall for dropout prediction.

**Concept Drift:**

When model accuracy degrades due to changes in student behavior or academic policies.

**Monitoring:** Compare recent predictions vs real outcomes regularly.

**Deployment Challenge:**

Scalability — Ensuring the model performs across various departments and programs.

---

### Part 2: Case Study Application (40 points)

**Problem Scope (5 points)**

**Problem:** Predict 30-day readmission risk after hospital discharge.

**Objectives:**

- Reduce hospital costs.

- Improve patient care and follow-up.

**Stakeholders:**

- Hospital IT team
  - Doctors and healthcare planners
- 

**Data Strategy (10 points)**

**Data Sources:**

- Electronic Health Records (EHRs)
- Demographic data

**Ethical Concerns:**

- Privacy of patient data
- Biased treatment recommendations for minority groups

**Preprocessing Pipeline:**

- Fill missing lab values
  - Encode gender, diagnosis
  - Normalize age, time in hospital
  - Feature engineering: e.g., count of past readmissions
- 

**Model Development (10 points)**

**Model:** Gradient Boosting Classifier – balances accuracy and interpretability.

**Confusion Matrix (Hypothetical):**

	Predicted Yes	Predicted No
Actual Yes	80	20
Actual No	10	90

**Precision:**  $80 / (80 + 10) = 0.89$

**Recall:**  $80 / (80 + 20) = 0.80$

---

### **Deployment (10 points)**

#### **Integration Steps:**

- Connect model to EHR via API
- Trigger predictions at discharge

#### **Compliance:**

- Use data encryption and de-identification
- Adhere to HIPAA laws

#### **Optimization:**

Use cross-validation or dropout regularization to avoid overfitting.

---

### **Part 3: Critical Thinking (20 points)**

#### **Ethics & Bias (10 points)**

**Risk:** Biased data may mislabel vulnerable patients, denying them follow-up care.

#### **Mitigation:**

- Use IBM AI Fairness 360 to test for bias
  - Adjust thresholds or resample training data
- 

#### **Trade-offs (10 points)**

##### **Interpretability vs Accuracy:**

Doctors may prefer simpler models they can understand, even if slightly less accurate.

##### **Resource Constraints:**

Low compute may force use of logistic regression over deep neural nets.

---

### **Part 4: Reflection & Workflow Diagram (10 points)**

#### **Reflection (5 points)**

**Challenge:** Balancing model performance with fairness.

**Improvement:** With more time, I would collect more diverse and recent data.

**Workflow Diagram (5 points)**

plaintext

CopyEdit

Problem Definition → Data Collection → Preprocessing → Model Training →

Evaluation → Deployment → Monitoring & Maintenance

---

**References**

- IBM AI Fairness 360 Toolkit
  - Kaggle: Breast Cancer Dataset
  - Scikit-learn Documentation
-