

AI Development Workflow Assignment

Course: AI for Software Engineering

Project Title: Predicting Hospital Patient Readmission Risk Using Machine Learning

Team/Author: Yvette Lando

Part 1: Short Answer Questions (30 Points)

1. Problem Definition (6 points)

Hypothetical Problem: Predicting student dropout rates in online learning environments.

Objectives:

1. Detect early signs of disengagement among students.
2. Enable timely academic interventions.
3. Improve overall student retention through data-driven support.

Stakeholders:

- Students
- Educational Administrators/Program Managers

Key Performance Indicator (KPI): F1-Score of the dropout prediction model.

2. Data Collection & Preprocessing (8 points)

Data Sources:

1. Learning Management System (LMS) usage logs
2. Student demographics and academic history

Potential Bias: Students with limited internet access may appear disengaged in LMS data, causing unfair labeling.

Preprocessing Steps:

1. Handle missing data through imputation or exclusion.
2. Normalize continuous features such as time spent online and assignment completion rates.
3. Encode categorical variables like course type and region using one-hot encoding.

3. Model Development (8 points)

Chosen Model: Random Forest

Justification: Handles non-linear data well, provides feature importance, and is less prone to overfitting compared to single decision trees.

Data Splitting Strategy:

- 70% Training
- 15% Validation
- 15% Testing (Stratified sampling due to class imbalance)

Hyperparameters to Tune:

1. `n_estimators`: Controls the number of trees in the forest.
2. `max_depth`: Prevents overfitting by limiting tree complexity.

4. Evaluation & Deployment (8 points)

Evaluation Metrics:

1. F1-Score: Balances precision and recall for imbalanced dropout data.
2. ROC-AUC: Evaluates model performance across classification thresholds.

Concept Drift: When data distribution or relationships change over time, degrading model accuracy.

Monitoring Strategy: Regularly evaluate metrics (e.g., F1, accuracy) over time and schedule periodic retraining.

Deployment Challenge: Scalability of serving predictions in real-time to thousands of students.

Part 2: Case Study Application (40 Points)

1. Problem Scope (5 points)

Problem: Predict whether a hospital patient will be readmitted within 30 days of discharge.

Objectives:

1. Reduce 30-day readmission rates and associated penalties.
2. Improve post-discharge planning and follow-up.
3. Support clinical decision-making through AI-driven insights.

Stakeholders:

- Patients
- Hospital administrators and clinical staff

2. Data Strategy (10 points)

Data Sources:

- Electronic Health Records (EHRs): Medications, lab results, vitals
- Patient demographics and past hospital admission records

Ethical Concerns:

1. Patient privacy and data protection (compliance with HIPAA)
2. Bias in training data may disadvantage underrepresented groups

Preprocessing Pipeline:

1. Impute missing values (e.g., using median for lab results).
2. Normalize continuous variables such as vital signs.
3. Feature engineering: time since last admission, chronic illness indicators, number of previous visits.

3. Model Development (10 points)

Model Chosen: XGBoost

Justification: Strong performance on structured/tabular data, built-in regularization, and supports missing values.

Confusion Matrix (Hypothetical):

	Predicted Yes	Predicted No
Actual Yes	45	15
Actual No	20	120

Precision: $45 / (45 + 20) = 0.692$

Recall: $45 / (45 + 15) = 0.75$

4. Deployment (10 points)

Integration Steps:

1. Export trained model to a .pkl file.
2. Create REST API endpoint using Flask or FastAPI.
3. Integrate API with hospital's EHR system.
4. Add model outputs to discharge summary workflow.

Compliance Measures:

- Use HIPAA-compliant cloud services.
- Apply encryption at rest and in transit.
- Implement audit logs and role-based access control.

5. Optimization (5 points)

Method to Prevent Overfitting: Apply cross-validation and use L2 regularization to penalize complex models.

Part 3: Critical Thinking (20 Points)**1. Ethics & Bias (10 points)****Impact of Biased Data:**

- If training data underrepresents certain groups (e.g., minority populations), the model may misclassify their risk, leading to health disparities.

Mitigation Strategy:

- Apply fairness-aware training methods (e.g., re-weighting) and conduct fairness audits.

2. Trade-offs (10 points)**Interpretability vs. Accuracy:**

- Highly accurate models like deep learning can be black boxes, which may reduce clinician trust.

Solution: Use models like XGBoost with SHAP value explanations to balance performance and interpretability.

Computational Constraints:

- Hospitals with limited infrastructure may need lightweight models or cloud-based inference services.

Part 4: Reflection & Workflow Diagram (10 Points)**1. Reflection (5 points)**

Most Challenging Stage: Deployment, due to complexities in integrating with hospital IT systems and ensuring compliance with privacy laws.

Improvement Plan: Allocate more time for stakeholder engagement and consider using MLOps tools for streamlined deployment and monitoring.

2. Workflow Diagram (5 points)

[Problem Definition]



[Data Collection]



[Preprocessing]



[Model Development]



[Evaluation]



[Deployment]



[Monitoring & Feedback]