

AI Project Assignment

Part 1: Short Answer Questions

1. Problem Definition

Hypothetical AI Problem: Predicting student dropout rates in online learning platforms.

Objectives:

- Identify students at high risk of dropping out within the next month.
- Recommend personalized interventions to improve engagement.
- Reduce overall dropout rates by 15% within a semester.

Stakeholders:

- University administration (decision-makers for interventions)
- Students (beneficiaries of early support)

Key Performance Indicator (KPI):

- Accuracy of dropout prediction model $\geq 85\%$
-

2. Data Collection & Preprocessing

Data Sources:

- Learning Management System (LMS) logs (e.g., login frequency, assignment submissions)
- Demographics and academic history (e.g., GPA, prior course performance)

Potential Bias:

- Students with poor internet access may have fewer LMS interactions, causing the model to overestimate dropout risk

Preprocessing Steps:

1. **Handling Missing Data:** Impute missing values using mean/mode or predictive models
 2. **Normalization/Scaling:** Standardize features to have zero mean and unit variance
 3. **Encoding Categorical Variables:** Convert non-numeric features (e.g., gender, major) to one-hot or label encoding
-

3. Model Development

Chosen Model: Random Forest Classifier

Justification: Robust to overfitting, handles mixed data types, interpretable feature importance

Data Split:

- Training: 70%
- Validation: 15%
- Test: 15%

Hyperparameters to Tune:

- Number of trees (**n_estimators**): Affects model performance and variance
 - Maximum tree depth (**max_depth**): Prevents overfitting while maintaining accuracy
-

4. Evaluation & Deployment

Evaluation Metrics:

- **F1-score**: Balances precision and recall, crucial for imbalanced dropout prediction
- **ROC-AUC**: Measures model's ability to distinguish between dropout vs. retention

Concept Drift:

- Concept drift occurs when data distributions change over time (e.g., new student behaviors)
- Monitor via continuous model evaluation on new student data; retrain model periodically

Technical Challenge During Deployment:

- **Scalability**: Handling real-time predictions for thousands of students during peak usage
-

Part 2: Case Study Application

Scenario: Predicting Patient Readmission Risk Within 30 Days

Problem Scope:

Problem: Reduce unplanned patient readmissions within 30 days of hospital discharge

Objectives:

- Predict 30-day readmission probability per patient
- Support targeted post-discharge interventions
- Reduce readmission costs and improve patient outcomes

Stakeholders:

- Hospital management and clinicians
 - Patients receiving care
-

Data Strategy

Data Sources:

- Electronic Health Records (EHRs) (diagnoses, medications, vitals)
- Demographic data (age, gender, socioeconomic status)

Ethical Concerns:

- **Patient privacy:** Data must comply with HIPAA regulations
- **Bias:** Historical disparities may reflect in model predictions

Preprocessing Pipeline:

1. Remove personally identifiable information (PII)
 2. Handle missing lab results (imputation)
 3. Feature engineering: compute comorbidity scores, previous admissions count
-

Model Development

Chosen Model: Gradient Boosted Trees (e.g., XGBoost)
Justification: High accuracy on tabular data; robust to missing values

Hypothetical Confusion Matrix:

	Predicted Yes	Predicted No
Actual Yes	40	10
Actual No	15	85

Precision: $40 / (40 + 15) = 0.727$
Recall: $40 / (40 + 10) = 0.8$

Deployment

Integration Steps:

1. Develop API endpoint for model predictions
2. Connect API with hospital EHR system
3. Implement real-time dashboard for clinicians

Compliance:

- Data encryption in transit and at rest
- Role-based access control for patient data
- HIPAA-compliant logging and auditing

Optimization:

- Use cross-validation and regularization (e.g., L1/L2) to prevent overfitting
-

Part 3: Critical Thinking

Ethics & Bias

Impact of Bias:

- ### Mitigation Strategy:

- ## Trade-offs

Resource Limitation:

- ## Reflection

Workflow Diagram

