# Understanding the AI Development Workflow

**Course:** AI for Software Engineering
**Duration:** 7 days
**Total Points:** 100

---

# Part 1: Short Answer Questions (30 points)

### 1. Problem Definition (6 points)

**AI Problem:** Predicting student dropout rates in universities

**Objectives:**

1. Identify students at risk of dropping out early.
2. Recommend interventions to improve retention.
3. Reduce overall dropout rate by targeted support.

**Stakeholders:**

- University administration
- Students

**Key Performance Indicator (KPI):**

- Accuracy of predicting at-risk students (% correctly identified)

---

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

**Data Sources:**

1. Student academic records (grades, attendance)
2. Student demographic surveys (age, financial background, engagement)

**Potential Bias:**

- Urban students may be overrepresented, biasing the model against rural students

**Preprocessing Steps:**

1. Handle missing data: Fill missing grades using median values
2. Normalize numeric features: Scale attendance percentages and grades
3. Encode categorical variables: Convert gender, course, and region into numerical codes

---

## 3. Model Development (8 points)

**Model Choice:** Random Forest – robust to overfitting, interpretable via feature importance

**Data Split:**

- 70% training, 15% validation, 15% test

**Hyperparameters to Tune:**

1. `n_estimators`: Number of trees, impacts accuracy and speed
2. `max_depth`: Prevents overfitting

---

## 4. Evaluation & Deployment (8 points)

**Evaluation Metrics:**

- Accuracy: Overall correctness
- F1-score: Balances precision and recall for imbalanced data

**Concept Drift:**

- When data distribution changes over time (e.g., new curriculum)
- Monitor by retraining periodically

**Technical Challenge:**

- Scalability: Ensuring real-time predictions for thousands of students

---

# Part 2: Case Study – Hospital Readmission Risk (40 points)

## Problem Scope (5 points)

**Problem:** Predict patients at risk of readmission within 30 days post-discharge

**Objectives:**

1. Reduce unnecessary readmissions
2. Improve patient care
3. Optimize hospital resources

**Stakeholders:**

- Hospital administrators
- Doctors and nurses
- Patients

---

## Data Strategy (10 points)

**Data Sources:**

1. Electronic Health Records (EHRs)
2. Patient demographics

**Ethical Concerns:**

- Patient privacy and HIPAA compliance
- Risk of biased predictions against some demographic groups

**Preprocessing Pipeline:**

1. Handle missing lab results with imputation
2. Encode categorical features (gender, disease type)
3. Feature engineering: Create "readmission risk score" from comorbidities and prior admissions

---

## Model Development (10 points)

**Model Choice:** Gradient Boosted Trees (XGBoost) – accurate, handles tabular data, interpretable via SHAP

**Hypothetical Confusion Matrix:**

|                     | Predicted Readmit | Predicted No Readmit |
|---------------------|-------------------|----------------------|
| Actual Readmit      | 80                | 20                   |
| Actual No Readmit   | 10                | 90                   |

**Calculations:**

- **Precision:** $80 / (80 + 10) = 0.888$
- **Recall:** $80 / (80 + 20) = 0.8$

---

## Deployment (10 points)

**Integration Steps:**

1. Convert model to API (Flask/FastAPI)
2. Integrate API with hospital EHR
3. Dashboard for doctors to view risk scores

**Compliance:**

- Encrypt patient data in transit and at rest
- Maintain access logs
- Follow HIPAA regulations

**Optimization (5 points):**

- Use cross-validation and regularization to reduce overfitting

---

# Part 3: Critical Thinking (20 points)

## Ethics & Bias (10 points)

**Impact:** Biased data may misclassify high-risk patients, causing worse outcomes for underrepresented groups

**Mitigation Strategy:** Use re-sampling or weighting to balance training data across demographics

---

## Trade-offs (10 points)

**Interpretability vs Accuracy:** Complex models (XGBoost) are accurate but harder to interpret; simpler models are more understandable but less accurate

**Computational Constraints:** Limited resources may require simpler models (logistic regression, shallow trees)

---

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

## Reflection (5 points)

- **Challenge:** Balancing accuracy, interpretability, and ethical considerations
- **Improvement:** More data and computational resources allow better tuning and fairness adjustments

---

## AI Development Workflow Diagram (5 points)

```
Problem Definition
        ↓
Data Collection
        ↓
Data Preprocessing
        ↓
Feature Engineering
        ↓
Model Development
        ↓
Model Evaluation
        ↓
Deployment
        ↓
Monitoring & Maintenance
```