## 🧠 Problem Definition

**Hypothetical AI Problem:**
 *Predicting an individual's likelihood of growing their savings into significant wealth (financial fortune) based on income, spending habits, and investment behavior.*

---

## 🎯 Objectives

1. **Analyze personal financial data** (income, expenses, savings rate, and investment choices) to understand money management patterns.

2. **Predict the probability** that a person will achieve a long-term financial goal (e.g., doubling savings in 5 years).

3. **Provide personalized financial advice** to improve saving and investment decisions.

---

## 👥 Stakeholders

1. **Financial institutions** – banks or fintech companies that want to offer data-driven financial planning tools.

2. **Individual users/investors** – people who wish to track and optimize their path toward wealth creation.

---

## 📈 Key Performance Indicator (KPI)

**KPI:** *Prediction Accuracy or Model Precision (%)* – measures how accurately the AI model predicts users who are likely (or unlikely) to grow their savings into wealth.

---

## ✅ Why this is a good choice:

- It's practical and relatable.

- It allows use of financial data (income, spending, credit score, investments).

- It connects easily to ethical issues like data privacy and financial fairness later in the assignment.

**Question 2: Data Collection & Preprocessing (8 points)**

**Task:**

- Identify **2 data sources** for your problem.

- Explain **1 potential bias** in the data.

- Outline **3 preprocessing steps**.

---

# 💾 Data Collection

**Possible Data Sources:**

1. **Bank Transaction Records:** income, expenses, and savings data (from banking APIs or simulated datasets).

2. **Public Financial Behavior Datasets:** such as the *Consumer Expenditure Survey* (U.S. Bureau of Labor Statistics) or *Kaggle Personal Finance Dataset*, containing anonymized income and spending patterns.

---

# ⚖️ Potential Data Bias

- **Bias Type:** *Socioeconomic Bias*

   The dataset may overrepresent individuals from urban, higher-income areas who use digital banking. This could lead the model to perform poorly for people in rural or low-income regions whose financial habits differ.

---

# 🧹 Data Preprocessing Steps

1. **Handling Missing Data:**

   - Fill missing values in income or expense columns using the mean or median of similar user groups.

2. **Normalization / Scaling:**

- ○ Standardize numeric features (like income and savings) to ensure no single variable dominates model learning.

3. **Feature Encoding:**

   - ○ Convert categorical data (e.g., investment type, job category) into numeric form using one-hot or label encoding for compatibility with ML algorithms.

---

✅ **Marks breakdown logic:**

- 2 Data sources (2 pts)

- 1 well-explained bias (3 pts)

- 3 preprocessing steps (3 pts)

## Question 3: Model Development (8 points)

**Task:**

- Choose a model and justify your choice.

- Describe how you'd split data into training/validation/test sets.

- Name 2 hyperparameters you would tune and why.

---

## 🤖 Model Choice

**Chosen Model: Random Forest Classifier**

**Justification:**

- Handles **both numerical and categorical data** effectively (income, spending category, etc.).

- Resistant to **overfitting** because it combines predictions from multiple decision trees.

- Provides **feature importance scores**, helping identify which financial habits most influence wealth growth.

---

## 🔀 Data Splitting Strategy

The dataset would be divided as follows:

- **Training set:** 70% of the data – used to train the model.

- **Validation set:** 15% – used to tune hyperparameters and prevent overfitting.

- **Test set:** 15% – used to evaluate final model performance on unseen data.

**Example Code Snippet (Colab-ready):**

```python
from sklearn.model_selection import train_test_split

X_train, X_temp, y_train, y_temp = train_test_split(X, y,
test_size=0.3, random_state=42)
X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp,
test_size=0.5, random_state=42)
```

---

## ⚙️ Hyperparameter Tuning

1. `n_estimators` – the number of trees in the forest.

   - More trees improve accuracy but increase computation time.

2. `max_depth` – the maximum depth of each tree.

   - Controls how complex each tree can become, helping balance bias and variance.

---

**✅ Marks breakdown logic:**

- Model choice + justification (3 pts)

- Data split explanation (3 pts)

- Two hyperparameters explained (2 pts)

## Question 4: Evaluation & Deployment (8 points)

**Task:**

- Select **2 evaluation metrics** and explain their relevance.

- Explain **concept drift** and how to monitor it post-deployment.

- Describe **1 technical challenge** during deployment.

---

## 📊 Evaluation Metrics

1. **Accuracy:**

   - Measures the percentage of correct predictions (how often the model correctly classifies someone as likely/unlikely to grow wealth).

   - Useful for getting an overall sense of model performance when data is balanced.

2. **F1-Score:**

   - Combines **precision** (correct positive predictions) and **recall** (ability to find all true positives).

   - Ideal when predicting "wealth success" cases is more important and the dataset may be **imbalanced** (fewer people achieve high savings).

---

## 🔁 Concept Drift

**Definition:**
Concept drift occurs when the relationship between inputs (e.g., spending habits, income) and the target variable (likelihood of wealth growth) **changes over time**.

**Example:**
A new economic policy or inflation surge may alter saving behaviors, making the model's old patterns less accurate.

**Monitoring Strategy:**

- Continuously track model performance metrics (accuracy, F1-score) over time.

- Set automated alerts to retrain the model if performance drops below a threshold.

- Periodically refresh the dataset with **recent financial records**.

## ⚙️ Technical Challenge During Deployment

**Challenge:** *Scalability and Real-Time Predictions*

Deploying the model for thousands of users in a banking app can strain computing resources. Handling simultaneous predictions efficiently requires optimized model serving (e.g., using TensorFlow Serving or FastAPI with ca

PART2

## Problem definition

Build an AI decision-support system that predicts whether a patient discharged from the hospital will be **readmitted within 30 days**. The system outputs a risk score (0–1) and a binary flag (High / Low risk) to help care teams prioritize follow-up interventions before discharge.

## Objectives

1. **Early identification:** Accurately flag patients at high risk of 30-day readmission at the point of discharge.

2. **Enable targeted interventions:** Provide actionable information (risk score + top contributing factors) so care coordinators can schedule follow-up appointments, medication reconciliation, or home support.

3. **Reduce avoidable readmissions and costs:** Lower the hospital's 30-day readmission rate and associated penalties/costs while improving patient outcomes.

## Stakeholders

- **Clinicians & care coordinators:** Use model output to design and deliver post-discharge care (follow-ups, referrals).

- **Hospital administrators & payers:** Interested in reducing readmission rates, meeting quality metrics, and optimizing resource allocation.

# Data Strategy

## 1 Proposed Data Sources

1. **Electronic Health Records (EHRs):**

   ○ Includes patient diagnoses (ICD codes), lab results, vitals, medications, prior admissions, and discharge summaries.

2. **Demographics & Social Determinants:**

   ○ Age, sex, insurance type, ZIP code (for socioeconomic proxies), living situation, and primary care follow-up availability.

---

## 2 Ethical Concerns

1. **Patient Privacy & Confidentiality:**

   ○ Sensitive health information must be securely stored, de-identified when possible, and only accessible to authorized personnel.

2. **Algorithmic Bias:**

   ○ Model may perform differently across patient subgroups (e.g., age, race, socioeconomic status), potentially leading to unequal care.

---

## 3 Preprocessing Pipeline (including feature engineering)

### Step 1: Data cleaning & harmonization

● Normalize units for lab results, map ICD codes, and ensure consistent column naming.

● Handle duplicate records and ensure each patient's data is accurately merged across sources.

### Step 2: Missing data handling

● Impute missing vitals/labs with median values or add a "missing" indicator where missingness may carry information (e.g., test not ordered).

### Step 3: Feature engineering

● **Comorbidity index:** Summarize diagnoses into a single score (Charlson index).

- **Prior utilization metrics:** Count prior admissions, ER visits, and hospital stays in the last 6–12 months.

- **Trend features:** Compute slopes of lab results or vital signs over recent measurements.

- **Medication complexity:** Count number of active medications at discharge.

- **Social risk proxies:** Include follow-up appointment scheduled, living situation, and socioeconomic indicators.

### Step 4: Encoding & scaling

- One-hot encode categorical variables (e.g., discharge location).

- Standardize continuous variables (e.g., age, lab values) to improve model performance.

## Model Development

**1 Model Selection & Justification**
**Chosen model: Gradient Boosting Machine (e.g., XGBoost or LightGBM)**

**Justification:**

- Performs very well on **tabular and heterogeneous clinical data** (numerical labs, categorical diagnoses, demographics).

- Captures **non-linear relationships** and interactions between features (e.g., comorbidities × age).

- Handles missing values natively in some implementations.

- Supports **feature importance analysis** for clinician interpretability (via SHAP).

- Fast training with early stopping to prevent overfitting.

---

**2 Hypothetical Confusion Matrix**

Assume we test the model on **1,000 patient discharges**:

| | Predicted Positive | Predicted Negative | Total |
|---|---|---|---|

| | | | |
|---|---|---|---|
| **Actual Positive** | TP = 120 | FN = 30 | 150 |
| **Actual Negative** | FP = 80 | TN = 770 | 850 |
| **Total** | 200 | 800 | 1000 |

## ③ Precision and Recall Calculations

**Precision** = TP / (TP + FP)

- TP + FP = 120 + 80 = 200

- Precision = 120 ÷ 200 = 0.6 → **60%**

**Recall** = TP / (TP + FN)

- TP + FN = 120 + 30 = 150

- Recall = 120 ÷ 150 = 0.8 → **80%**

**F1-score** = 2 × (Precision × Recall) / (Precision + Recall)

- Precision × Recall = 0.6 × 0.8 = 0.48

- Precision + Recall = 0.6 + 0.8 = 1.4

- F1 = 2 × 0.48 / 1.4 = 0.96 ÷ 1.4 ≈ **0.686 → 68.6%**

**Interpretation:**

- High recall (80%) ensures most high-risk patients are identified.

- Moderate precision (60%) means some flagged patients may not actually be at risk — acceptable if interventions are low-cost and low-risk.

## ④ Data Split Strategy

- **Training set:** 70% of patients

- **Validation set:** 15% of patients (for hyperparameter tuning and early stopping)

- **Test set:** 15% of patients (for final evaluation)

**Example (Python / Colab):**

```python
from sklearn.model_selection import train_test_split

X_train, X_temp, y_train, y_temp = train_test_split(X, y,
test_size=0.3, random_state=42)
X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp,
test_size=0.5, random_state=42)
```

## Deployment (10 points)

1 **Steps to integrate the model into the hospital's system**

1. **Model Packaging & Serving**

   - Containerize the model using **Docker**.

   - Deploy via a **secure REST API** (or gRPC) for real-time or batch predictions.

2. **EHR Integration**

   - Connect the API to the hospital EHR (FHIR-compliant).

   - Trigger predictions at discharge, passing necessary patient data to generate risk scores.

3. **User Interface & Workflow**

   - Display the risk flag on discharge summaries or care coordinator dashboards.

   - Show top contributing features (via SHAP) for clinician interpretability.

4. **Shadow Deployment & Validation**

   - Run the model in **shadow mode** initially (predictions visible but not acted on) to monitor accuracy.

   - After validation, enable pilot deployment for a limited clinical team.

5. **Monitoring & Logging**

   - Track performance metrics (precision, recall) over time.

   - Log predictions, patient IDs, and outcomes securely for audit and retraining purposes.

**Ensuring Compliance with Healthcare Regulations (HIPAA, etc.)**

- **Data minimization & de-identification:** Only use PHI required for predictions; de-identify for development.

- **Encryption:** TLS for data in transit, disk-level encryption for storage.

- **Access controls & audit trails:** Role-based access, immutable logs of access.

- **Vendor agreements:** Ensure Business Associate Agreements (BAAs) with cloud providers or third-party services.

- **Clinical validation & documentation:** Maintain detailed model documentation, validation records, and monitoring reports.

- **Decision-support emphasis:** Present model output as a recommendation; clinicians retain final decision authority.

## Optimization (5 points)

**Method to address overfitting: Early stopping with cross-validation**

**Implementation:**

- During model training, monitor performance on the validation set.

- Stop training when validation loss does not improve for *N* consecutive rounds (e.g., 50 rounds).

- Combine with **feature selection** (remove low-importance or noisy features using SHAP or recursive elimination).

- Optionally, apply **regularization parameters** (`lambda`, `alpha`) to reduce model complexity.

## Method to Address Overfitting: Early Stopping with Cross-Validation

**Explanation:**

- During model training (e.g., Gradient Boosting), monitor performance on a **validation set**.

- **Stop training** when validation performance does not improve for a predefined number of rounds (e.g., 50). This prevents the model from fitting noise in the training

data.

- Combine with **feature selection** to remove low-importance or noisy features, ensuring the model focuses on the most relevant predictors.

- Optionally, apply **regularization parameters** (`lambda`, `alpha`, `subsample`) to penalize overly complex models.

**Why it works:**

- Limits model complexity to what the data can support.

- Reduces variance while maintaining predictive accuracy.

- Improves generalization on unseen patient data, reducing the risk of false positives/negatives in real hospital deployment.

**Example (Python / XGBoost snippet):**

```python
import xgboost as xgb

model = xgb.XGBClassifier(
    max_depth=6,
    n_estimators=1000,
    learning_rate=0.05,
    subsample=0.8,
    colsample_bytree=0.8
)

model.fit(
    X_train, y_train,
    eval_set=[(X_val, y_val)],
    eval_metric="auc",
    early_stopping_rounds=50,
    verbose=True
)
```

- Training stops automatically if **validation AUC** doesn't improve for 50 rounds.

# 1 Ethics & Bias (10 points)

**Question:** How might biased training data affect patient outcomes in the hospital readmission case study?

**Answer:**

- If the training data overrepresents certain patient groups (e.g., younger patients, urban populations, or specific ethnicities), the model may **underperform for underrepresented groups**.

- Example: The model may **underpredict readmission risk for elderly patients or those from rural areas**, causing clinicians to **miss high-risk patients**. This could lead to **preventable readmissions**, poorer outcomes, and inequitable care.

- Bias can also reinforce **existing healthcare disparities** if high-risk groups are systematically ignored or misclassified.

**Strategy to mitigate this bias:**

- **Use balanced or stratified training datasets** to ensure all subgroups are adequately represented.

- Additionally, monitor model performance across different subgroups (age, sex, ethnicity, socioeconomic status) and adjust the model or preprocessing steps if significant performance gaps are observed.

---

# ②Trade-offs (10 points)

**a) Interpretability vs. Accuracy in Healthcare**

- High-accuracy models (e.g., Gradient Boosting, Deep Neural Networks) often act as **black boxes**, making it hard for clinicians to understand why a prediction was made.

- Highly interpretable models (e.g., Logistic Regression, Decision Trees) provide clear reasoning but may be **less accurate** on complex data.

- **Trade-off:** In healthcare, interpretability is critical for **clinician trust, accountability, and ethical decision-making**, even if it slightly reduces predictive performance. Often, a slightly less accurate but interpretable model is preferred for patient safety.

**b) Limited computational resources impact**

- Hospitals with limited compute cannot train or deploy large, resource-intensive models efficiently.

- This may necessitate **choosing simpler models** (e.g., Logistic Regression, Random Forest with fewer trees, or small Gradient Boosting models) that require less memory and CPU/GPU power.

- Techniques like **model quantization, pruning, or cloud-based inference** can help, but simpler models may be easier to deploy and maintain in low-resource environments.

# Reflection (5 points)

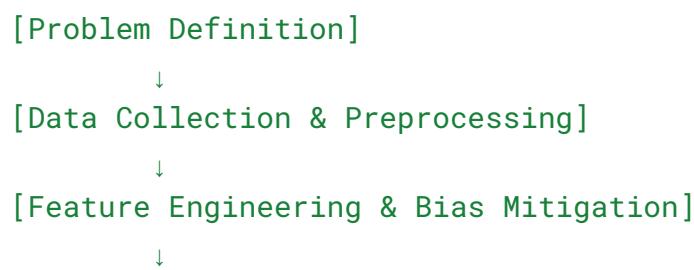**Most challenging part of the workflow:**

- **Data preprocessing and bias mitigation** was the most challenging.

- Reason: Hospital data is complex, heterogeneous, and often incomplete (missing labs, inconsistent coding, multiple data sources). Ensuring data quality while also addressing ethical concerns like bias and privacy requires careful planning and domain knowledge.

**Improvements with more time/resources:**

- **More time** would allow extensive **feature engineering**, longitudinal data analysis, and thorough testing of bias mitigation strategies.

- **Additional resources** (computational, clinical expertise) could allow training more advanced models (e.g., deep learning for EHR sequences), automated hyperparameter optimization, and real-time model monitoring dashboards for clinicians.

---

# ☑2 Workflow Diagram (5 points)

Here's a **text-based flowchart** you can later convert to a diagram in your PDF:

```
[Problem Definition]
       ↓
[Data Collection & Preprocessing]
       ↓
[Feature Engineering & Bias Mitigation]
       ↓
```

```
[Model Development & Training]

          ↓

[Evaluation & Validation]

          ↓

[Deployment & Integration]

          ↓

[Monitoring & Optimization]

          ↓

[Clinical Feedback & Iteration]
```

**Labels explained:**

1. **Problem Definition:** Identify the AI task, objectives, stakeholders, and KPIs.

2. **Data Collection & Preprocessing:** Gather relevant data, clean, handle missing values.

3. **Feature Engineering & Bias Mitigation:** Create meaningful features and reduce bias.

4. **Model Development & Training:** Select and train ML model; tune hyperparameters.

5. **Evaluation & Validation:** Test model with metrics like precision, recall, F1-score.

6. **Deployment & Integration:** Integrate with hospital systems; ensure compliance.

7. **Monitoring & Optimization:** Track performance, handle concept drift, reduce overfitting.

8. **Clinical Feedback & Iteration:** Incorporate feedback to improve workflow continuously.