

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

**Scenario:** A hospital wants an AI system to predict patient readmission risk within 30 days of discharge.

### Problem Scope (5 points)

- **Problem:** To accurately predict which patients are at high risk of readmission within 30 days of discharge, enabling targeted interventions to improve patient outcomes and reduce healthcare costs.
- **Objectives:**
  - Reduce 30-day patient readmission rates.
  - Improve patient care coordination and post-discharge planning.
  - Optimize hospital resource allocation.
- **Stakeholders:**
  - **Patients:** Benefit from improved health outcomes and reduced readmissions.
  - **Hospital Administration:** Aims to reduce costs, improve quality metrics, and optimize resource use.
  - **Clinicians (Doctors, Nurses):** Use predictions to guide care decisions and interventions.
  - **Data Scientists/AI Developers:** Responsible for building and maintaining the system.
  - **Regulators:** Ensure compliance with healthcare data privacy and safety standards.

### Data Strategy (10 points)

- **Proposed Data Sources:**
  - **Electronic Health Records (EHRs):** Patient demographics, diagnoses (ICD-10 codes), procedures (CPT codes), medications, lab results, vital signs, discharge summaries, clinical notes.

- **Demographics:** Age, gender, socioeconomic status, insurance type, residential address (for social determinants of health).
- **Past Readmission History:** Previous hospitalizations and readmission events.
- **Social Determinants of Health (SDOH) Data:** Information on housing stability, food security, access to transportation, and social support (if available and ethically sourced).

- **Ethical Concerns:**

1. **Patient Privacy and Data Security (HIPAA Compliance):** Ensuring that sensitive patient health information (PHI) is protected from unauthorized access, use, or disclosure. This includes anonymization/de-identification techniques.
2. **Algorithmic Bias and Fairness:** The model might inadvertently discriminate against certain demographic groups (e.g., based on race, socioeconomic status) if the training data is unrepresentative or contains historical biases, leading to unequal access to interventions.

- **Preprocessing Pipeline:**

1. **Data Cleaning:**

- Handle missing values (e.g., imputation for lab results, dropping features with excessive missingness).
- Correct inconsistencies and errors (e.g., typos in diagnoses, out-of-range values).
- Remove duplicate records.

2. **Data Transformation:**

- **Categorical Encoding:** One-hot encoding for nominal variables (e.g., gender, insurance type), ordinal encoding for ordinal variables (e.g., severity scores).
- **Numerical Scaling:** Standardization or normalization for numerical features (e.g., lab values, age) to ensure features contribute equally to the model.

3. **Feature Engineering:**

- **Comorbidity Scores:** Calculate a Charlson Comorbidity Index or Elixhauser Comorbidity Index based on patient diagnoses to quantify overall health burden.
- **Length of Stay (LOS):** Calculate the duration of the current hospitalization.
- **Medication Adherence Indicators:** Create features from prescription refill history or medication reconciliation at discharge.
- **Frequency of Past Hospitalizations:** Number of admissions in the last 6 months or year.
- **Time Since Last Admission:** Days since the most recent previous hospitalization.
- **Discharge Disposition:** Categorical feature indicating where the patient is discharged to (e.g., home, skilled nursing facility, hospice).

### Model Development (10 points)

- **Model Selection and Justification:**
  - **Model:** Gradient Boosting Machine (e.g., XGBoost or LightGBM).
  - **Justification:**
    - **High Performance:** GBMs are known for their strong predictive accuracy on tabular data, often outperforming simpler models.
    - **Handles Mixed Data Types:** Can naturally handle both numerical and categorical features.
    - **Feature Importance:** Provides insights into which factors are most influential in predicting readmission, aiding clinical interpretability.
    - **Robustness to Outliers:** Less sensitive to outliers compared to some other models.
    - **Scalability:** Efficient for large datasets, which are common in healthcare.
- **Confusion Matrix and Precision/Recall (Hypothetical Data):**
  - Let's assume our model predicts readmission for 100 patients.

- **Confusion Matrix:**

	Predicted Readmission	Predicted No Readmission	
Actual Readmission	20 (TP)	5 (FN)	25
Actual No Readmission	10 (FP)	65 (TN)	75
	30	70	100

True Positives (TP) = 20 | False Negatives (FN) = 5 | **Actual Readmission** |  
False Positives (FP) = 10 | True Negatives (TN) = 65 |
- **Calculation:**
  - **Precision:** Of all patients predicted to be readmitted, what proportion actually were?  $\text{Precision} = \frac{TP}{TP+FP} = \frac{20}{20+10} = \frac{2}{3} \approx 0.67$  (or 67%)  
*Interpretation:* When the model predicts readmission, it is correct 67% of the time.
  - **Recall:** Of all patients who actually were readmitted, what proportion did the model correctly identify?  $\text{Recall} = \frac{TP}{TP+FN} = \frac{20}{20+5} = \frac{4}{5} = 0.80$  (or 80%)  
*Interpretation:* The model identifies 80% of all actual readmissions.

## Deployment (10 points)

- **Integration Steps:**
  1. **API Development:** Create a RESTful API endpoint for the trained model, allowing other hospital systems to send patient data and receive readmission risk predictions.
  2. **Data Ingestion Pipeline:** Establish a secure and automated pipeline to extract relevant patient data from EHRs in near real-time or batch, transform it into the features required by the model, and feed it to the API.
  3. **Integration with Clinical Workflow:**
    - Display predictions within the EHR system (e.g., as a risk score on a patient's dashboard).
    - Trigger alerts or notifications to care managers or discharge planners for high-risk patients.
    - Integrate with existing patient management or care coordination platforms.
  4. **Monitoring and Alerting:** Implement a system to continuously monitor model performance (e.g., drift in predictions, data quality issues) and alert relevant teams if performance degrades.

5. **User Interface (UI) Development:** Develop user-friendly dashboards for clinicians to view risk scores, contributing factors, and track intervention effectiveness.
- **Ensuring Compliance with Healthcare Regulations (e.g., HIPAA):**
    1. **Data De-identification/Anonymization:** Implement robust techniques to de-identify PHI before it is used for model training, testing, or even prediction where possible, adhering to HIPAA's Safe Harbor or Expert Determination methods.
    2. **Access Control and Encryption:** Implement strict role-based access control (RBAC) to the AI system and underlying data. Ensure all data (at rest and in transit) is encrypted using industry-standard protocols.
    3. **Audit Trails:** Maintain comprehensive audit logs of all data access, model predictions, and system changes to demonstrate accountability and compliance.
    4. **Data Use Agreements (DUAs) and Business Associate Agreements (BAAs):** Establish formal agreements with all third-party vendors or internal departments involved in data handling to ensure they comply with HIPAA regulations.
    5. **Regular Security Audits and Penetration Testing:** Conduct periodic assessments to identify and remediate vulnerabilities in the system.
    6. **Transparency and Explainability:** While not directly a HIPAA requirement, providing explainable AI (XAI) insights can help clinicians understand predictions, which is crucial for ethical use and demonstrating due diligence in patient care.

### Optimization (5 points)

- **Method to Address Overfitting:**
  - **Regularization:** Apply L1 (Lasso) or L2 (Ridge) regularization during model training. This adds a penalty to the model's loss function based on the magnitude of the coefficients (L2) or the absolute values of the coefficients (L1), discouraging overly complex models that fit the training data too closely. For tree-based models like GBMs, regularization can be applied by controlling parameters like `max_depth`, `min_child_weight`, `subsample`, `colsample_bytree`, and `lambda/alpha` (L1/L2 regularization terms).