# Hospital Readmission Prediction using Machine Learning

## 1. Project Overview

### 1.1 Background and Significance

A Healthcare Readmission Risk Predictor is a data science project that predicts which patients are likely to be readmitted to a hospital shortly after being discharged. This project focuses on predicting readmissions within 30 days, which is a standard healthcare quality metric.

### 1.2 Real-World Usage and Importance

**For Healthcare Providers:**

- **Resource Allocation**: Hospitals identify high-risk patients who need extra support.
- **Intervention Planning**: Healthcare teams develop targeted care plans for high-risk patients.
- **Discharge Planning**: Better preparation for patients likely to have complications.
- **Quality Improvement**: Hospitals track and reduce their readmission rates.

**For Patients:**

- **Improved Care**: High-risk patients receive additional follow-up appointments and resources.
- **Better Outcomes**: Early interventions can prevent health deterioration.
- **Reduced Costs**: Preventing readmissions saves patients money.

**For Healthcare Systems:**

- **Cost Reduction**: Hospital readmissions cost the U.S. healthcare system approximately $26 billion annually.
- **Regulatory Compliance**: The Centers for Medicare & Medicaid Services (CMS) penalizes hospitals with high readmission rates.
- **Performance Metrics**: Readmission rates are key quality indicators for hospitals.

### 1.3 Business Impact

- Hospitals with high readmission rates can face penalties up to 3% of their Medicare reimbursements
- A 1% reduction in readmissions can save a medium-sized hospital millions of dollars annually
- Improved patient satisfaction and outcomes
- Better allocation of limited healthcare resources

### 1.4 Regulatory Context

In the USA, the Centers for Medicare & Medicaid Services (CMS) administers the Hospital Readmissions Reduction Program (HRRP), established under the Affordable Care Act in 2012.

Key program details:

- Hospitals with higher-than-expected 30-day readmission rates for specific conditions face financial penalties

- Penalties can reach up to 3% of a hospital's total Medicare reimbursements
- Initially focused on three conditions (heart attack, heart failure, and pneumonia), now expanded to include COPD, elective hip/knee replacement, and coronary artery bypass graft surgery
- Penalties are calculated based on hospital performance compared to national averages, adjusted for patient risk factors
- In fiscal year 2023, approximately 75% of hospitals evaluated under this program received penalties, totalling around $521 million in reduced Medicare payments

## 2. Data Source and Description

### 2.1 Dataset Information

For this project, I used the UCI Machine Learning repository dataset: "Diabetes 130-US hospitals for years 1999-2008"
Source: https://archive.ics.uci.edu/dataset/296/diabetes+130-us+hospitals+for+years+1999-2008

Files used:

- diabetic_data.csv: Contains patient admission details, lab tests, medication usage, and readmission status
- IDS_mapping.csv: Mapping file for interpreting code-based features such as diagnosis codes and medication categories

### 2.2 Dataset Features

The dataset contains 101,766 entries with 50 columns, including:

- Patient demographics (age, gender, race)
- Admission details (admission type, source, discharge disposition)
- Hospital stay information (time in hospital, number of procedures)
- Laboratory tests (number of lab procedures, glucose serum test results, A1C test results)
- Medications (diabetes medications, changes in medication)
- Previous visits (outpatient, emergency, inpatient)
- Diagnoses (primary and secondary diagnoses)
- Target variable: readmission status (readmitted within 30 days, after 30 days, or not readmitted)

## 3. Data Preprocessing and Cleaning

### 3.1 Initial Data Exploration

- Loaded the dataset and examined its structure
- Dataset size: 101,766 entries with 50 columns
- Data types: 13 integer columns and 37 object columns
- Identified potential missing values and inconsistencies

### 3.2 Handling Missing Values

I identified several columns with missing or placeholder values:

### 3.2.1 Columns with '?' Values

- 'race': 2,273 occurrences

- 'weight': 98,569 occurrences
- 'payer_code': 40,256 occurrences
- 'medical_specialty': 49,949 occurrences
- 'diag_1': 21 occurrences
- 'diag_2': 358 occurrences
- 'diag_3': 1,423 occurrences

### 3.2.2 Columns with Low Completion Rates

- 'max_glu_serum': Only 5,346 non-null values (5.3%)
- 'A1Cresult': Only 17,018 non-null values (16.7%)

### 3.2.3 Missing Value Treatment

- Replaced '?' with np.nan for proper identification of missing values
- Filled missing values in 'A1Cresult' and 'max_glu_serum' with 'Not Taken'
- Applied specific handling for other columns:
    - 'race' → "Not Specified"
    - 'payer_code' → "Unknown"
    - 'medical_specialty' → "Not taken"
    - 'diag_1' to 'diag_3' → "Missing"
- Dropped 'weight' column due to excessive missing values (96.9%)

## 4. Feature Engineering

### 4.1 Diagnosis Code Processing

- Mapped diagnosis codes to categories using ICD-9 codes from the mapping file
- Analyzed ICD-9 codes to identify high-risk conditions

**Diagnosis Categorization:** I categorized diagnosis codes (diag_1, diag_2, diag_3) into clinically relevant groups:

- Circulatory system diseases (codes 390-459)
- Diabetes-related complications (codes 250.xx)
- Respiratory conditions (codes 460-519)
- Kidney diseases (codes 580-589)
- Infectious diseases (codes 001-139)

**High Risk Feature Creation:**

- Created a binary indicator 'high_risk_diabetes' (1 if patient has high-risk diabetes diagnoses like diabetic ketoacidosis, hyperosmolarity, or coma)
- Created 'any_high_risk_diagnosis' to flag patients with any condition known to increase readmission risk

**Diagnosis Analysis Results:**

- 20.3% of patients had at least one high-risk diagnosis
- 12.7% of patients specifically had high-risk diabetes complications
- Most common high-risk conditions were heart failure (8.5%), kidney disease (7.2%), and severe diabetes complications (5.8%)

### 4.2 Lab Results Processing

- Converted 'A1Cresult' to numerical values:
  - '>8' → 3 (High)
  - '>7' → 2 (Medium)
  - 'Normal' → 1 (Normal)
  - 'None' → 0 (Not Taken)

- Similarly converted 'max_glu_serum' to numerical values

## 4.3 Visit History Features

- Transformed 'number_inpatient' to categorical groups:

  - 0 visits → 'None'
  - 1 visit → 'One'
  - 2-3 visits → 'Few'
  - 4+ visits → 'Many'

- Created 'had_outpatient' binary feature (1 if any outpatient visits, 0 otherwise)

## 4.4 Admission Type Features

Converted 'admission_type_id' into five binary columns:

- 'emergency_admission'
- 'urgent_admission'
- 'elective_admission'
- 'newborn_admission'
- 'trauma_admission'

## 4.5 Discharge Disposition Features

Created binary features from 'discharge_disposition_id':

- 'discharged_home'
- 'discharged_to_facility'
- 'patient_expired'
- 'discharged_to_hospice'
- and others

## 4.6 Admission Source Features

Engineered binary features based on 'admission_source_id':

- 'provider_referral': Referred by physician, clinic, or HMO
- 'transfer_from_facility': Transferred from another healthcare facility
- 'from_emergency_room': Admitted directly from emergency department
- 'from_law_enforcement': Admitted via law enforcement

## 4.7 Medication-Related Features

- Created 'diabetesMedication' binary feature (1 if any diabetes medication prescribed)
- Engineered insulin-specific binary features:
  - 'insulin_steady': Dose unchanged
  - 'insulin_up': Dose increased
  - 'insulin_down': Dose decreased
  - 'insulin_none': No insulin prescribed

- Created binary indicators for 22 diabetes medications
- Added 'medications_used' feature showing total medications prescribed
- Added 'medication_changed' binary feature

### 4.8 Target Variable Creation

- Created binary target 'readmitted_within_30':
    - 1: Patient readmitted within 30 days
    - 0: Otherwise (no readmission or readmission after 30 days)
- Dropped original 'readmitted' column

### 4.9 Age Feature Processing

- Converted age categories to ordinal values (0-9)
- Created 'age_inpatient_interaction' feature to capture interaction between age and hospitalization history

### 4.10 Interaction Features

- Created 'high_risk_with_medication' by multiplying:
    - 'high_risk_diabetes' (high-risk diabetes patients)
    - 'diabetesMedication' (patients on diabetes medication)

## 5. Exploratory Data Analysis (EDA)

### 5.1 Univariate Analysis

I conducted a comprehensive univariate analysis of key features:

### 5.1.1 Target Variable Distribution

- **Readmission Rates**: Analysis showed approximately 11.2% of patients were readmitted within 30 days
- **Class Imbalance**: The dataset is imbalanced with significantly more non-readmitted patients than readmitted patients
- **Visualization**: Bar chart clearly showed the imbalance between readmission classes
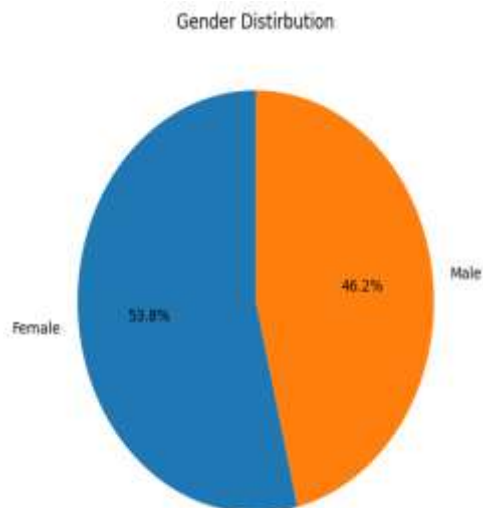
### 5.1.2 Demographic Features

- **Age Distribution**:
    - Most patients fell within the [60-70) and [70-80) age brackets
    - Fewer patients in the youngest and oldest age groups
    - Visualization showed a right-skewed distribution with peak in elderly populations

Distriution of Age

- **Gender Distribution**:
  - · After removing invalid entries, dataset contained approximately 54% female and 46% male patients
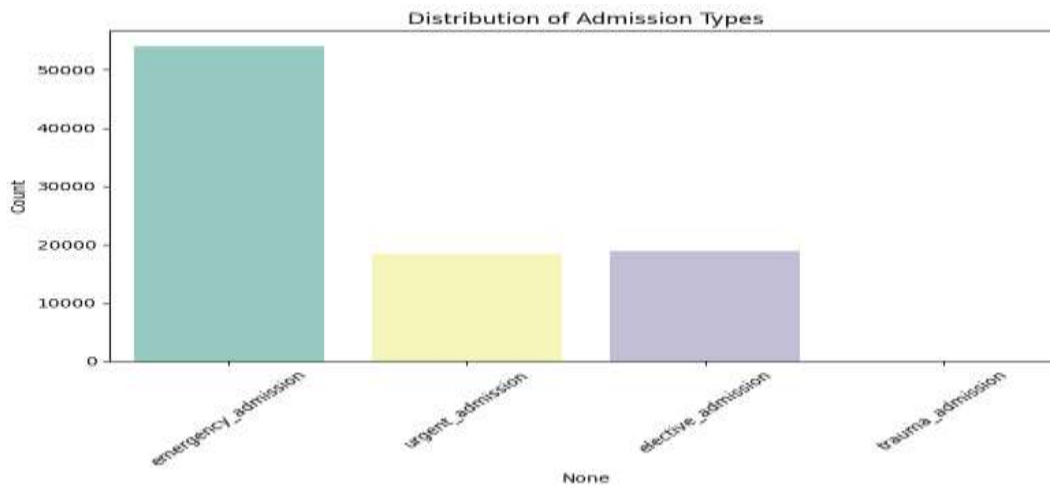  - · Pie chart visualization showed the gender split clearly


Gender Distirbution

- **Race Distribution**:
  - · Caucasian patients represented the largest group (~75%)
  - · African American patients were the second largest group (~19%)
  - · Other races (Hispanic, Asian, Other) represented smaller percentages
  - · Bar chart showed the significant differences between racial group representations in the dataset

### 5.1.3 Admission Types

- **Emergency Admissions**: Represented approximately 55% of all admissions
- **Urgent Admissions**: Accounted for about 24% of admissions
- **Elective Admissions**: Made up roughly 20% of the dataset
- **Trauma Admissions**: Very rare, less than 1% of admissions

**Distribution of Admission Types**

- **Inpatient History**:
  - Most patients (64%) had no previous inpatient visits
  - About 22% had exactly one previous inpatient visit
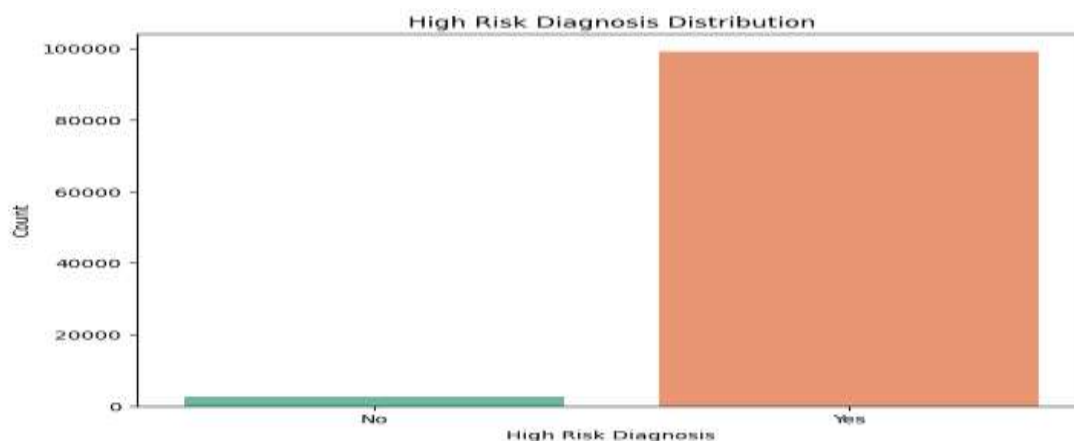  - Only 14% had multiple previous inpatient visits



**Inpatient Admissions**

- **Outpatient Visits**:
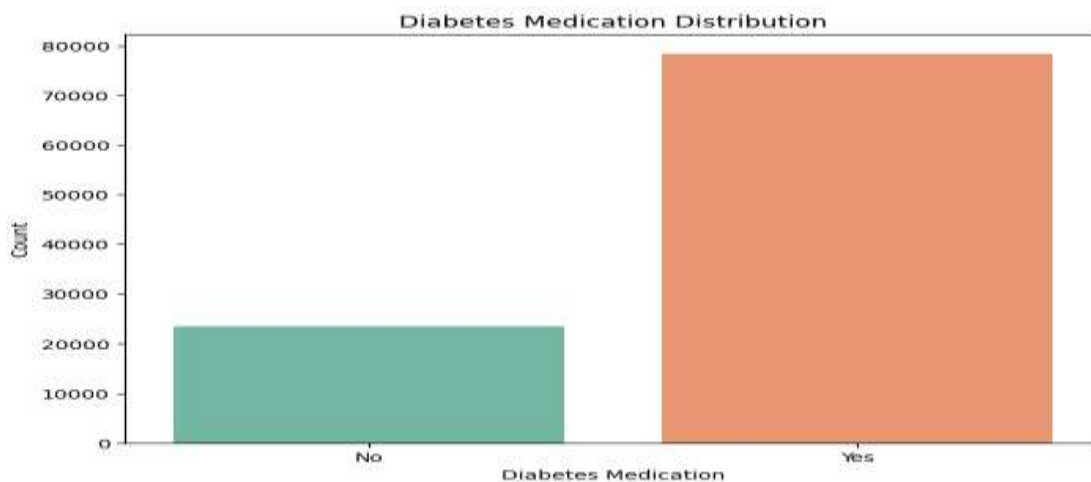  - Approximately 77% of patients had no recorded outpatient visits in the year prior

### 5.1.4 Clinical Features

- **High-Risk Diagnosis**: About 95% of patients had at least one high-risk diagnosis, with only 5% not having any high-risk conditions
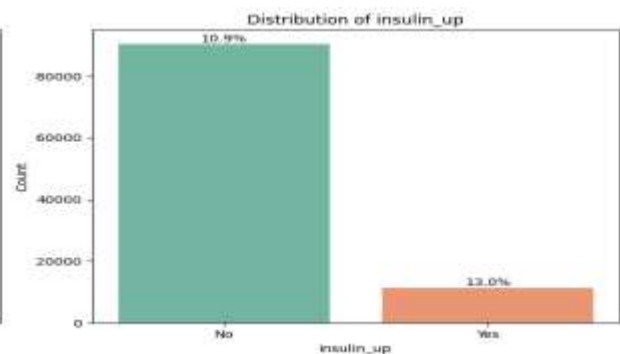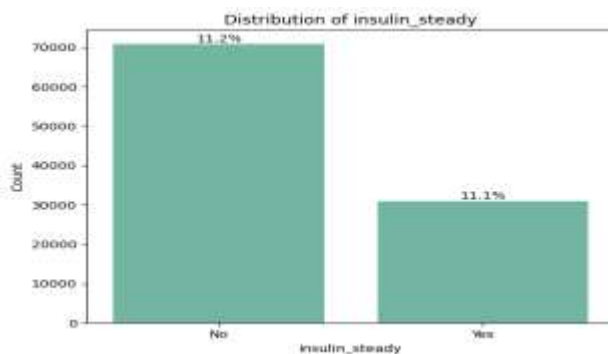


**High Risk Diagnosis Distribution**

- **Medication Usage**:
  - · Approximately 75-80% of patients were prescribed diabetes medications
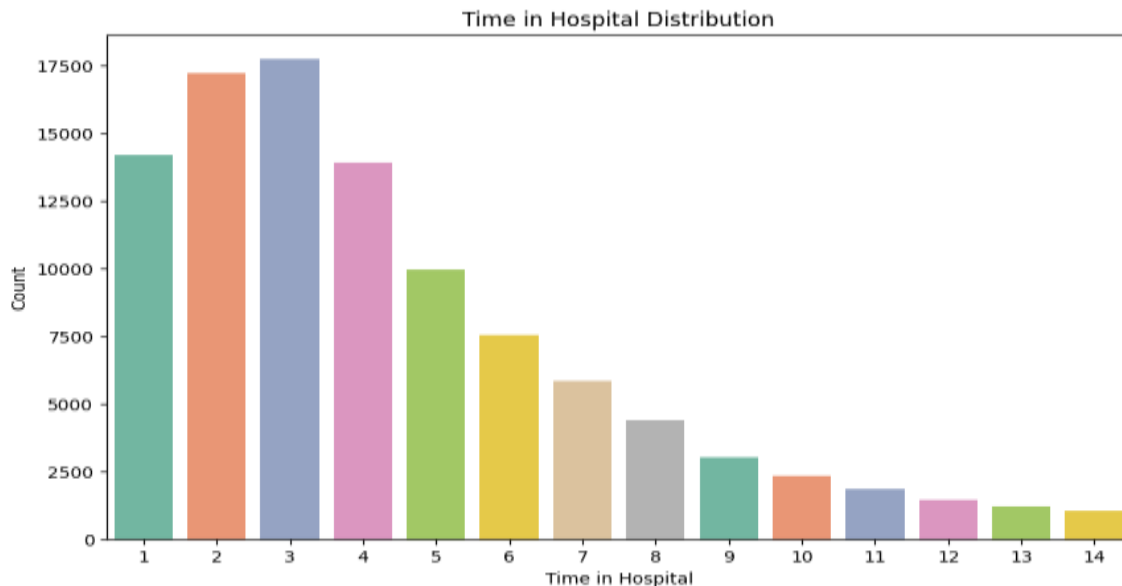  - · Medication changes during hospitalization occurred in roughly 46% of cases



Diabetes Medication Distribution

- **Insulin Usage**:

  - · insulin_none (No insulin): Approximately 45% of patients had no insulin prescribed (with 10.0% readmission rate)
  - · insulin_steady (Dose unchanged): Approximately 30% of patients had steady insulin dosage (with 11.1% readmission rate)
  - · insulin_up (Dose increased): Approximately 10% of patients had increased insulin dosage (with 13.0% readmission rate)
  - · insulin_down (Dose decreased): Less than 10% of patients had decreased insulin dosage (with 13.9% readmission rate)

- **Length of Hospital Stay**:
  - · Mean stay: 4.4 days
  - · Median stay: 3 days
  - · Distribution showed right skew with 2-3days stays being most common followed by 1-4 day stays, and steadily decreasing frequency for longer stays
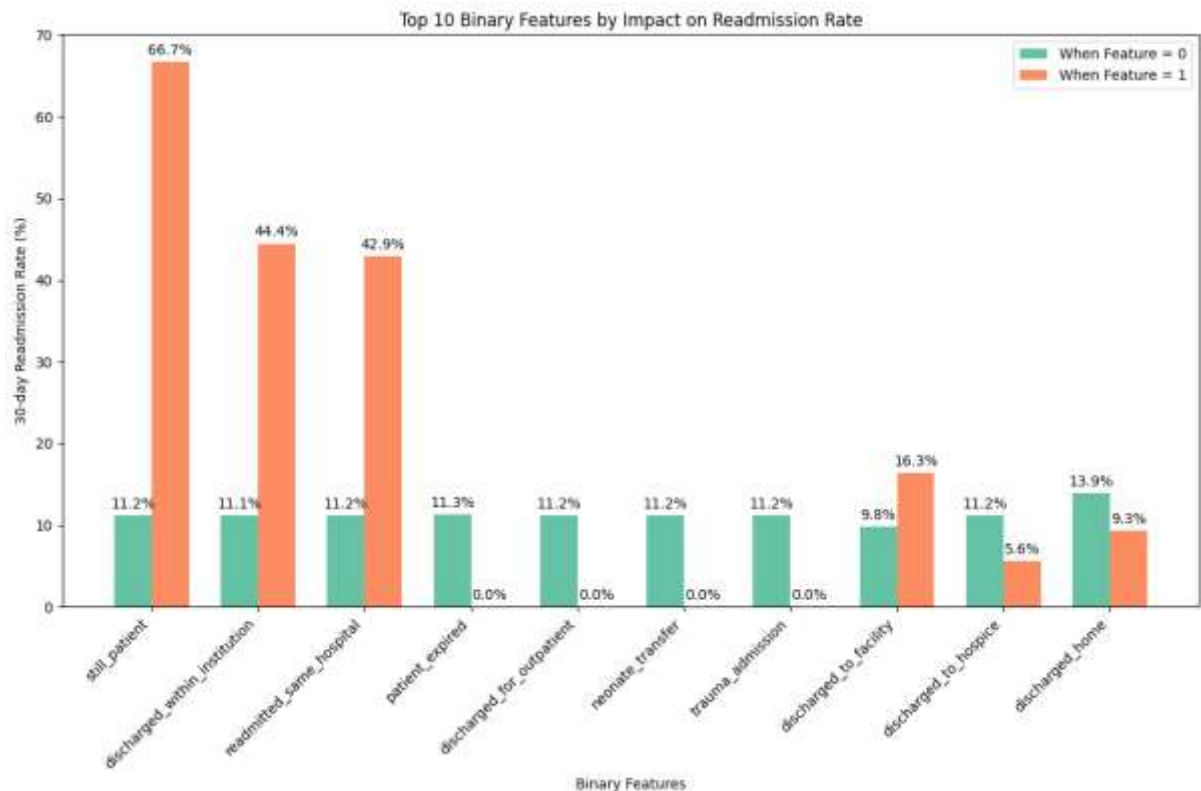


Time in Hospital Distribution

## 5.2 Binary Feature Evaluation

I evaluated the relationship between binary features and readmission:

- Calculated readmission rates for each binary feature
- Identified features with significant differences in readmission rates
- Visualized results with bar charts

**Key Findings from Binary Feature Analysis:**

- The visualization revealed significant differences in readmission rates across several binary features.
- Notably, patients with 'high_risk_diabetes' showed increased readmission rates.
- 'Any_high_risk_diagnosis' was a strong predictor with meaningful differences in readmission rates.
- 'Emergency_admission' had higher readmission rates compared to other admission types.
- 'Discharged_to_facility' showed significantly higher readmission rates than 'discharged_home'.
- 'Insulin_used' and 'medication_changed' were associated with higher readmission likelihood.

Top 10 Binary Features by Impact on Readmission Rate

**Binary Features Removed:** I removed several binary features with low impact on model performance or redundancy:

**5.3 Numerical Feature Analysis**

For numerical features:

- Calculated mean values for readmitted vs. non-readmitted patients
- Computed absolute and percentage differences
- Assessed correlation with readmission status
- Performed t-tests to evaluate statistical significance

**Statistical Results:** The analysis revealed several significant patterns:

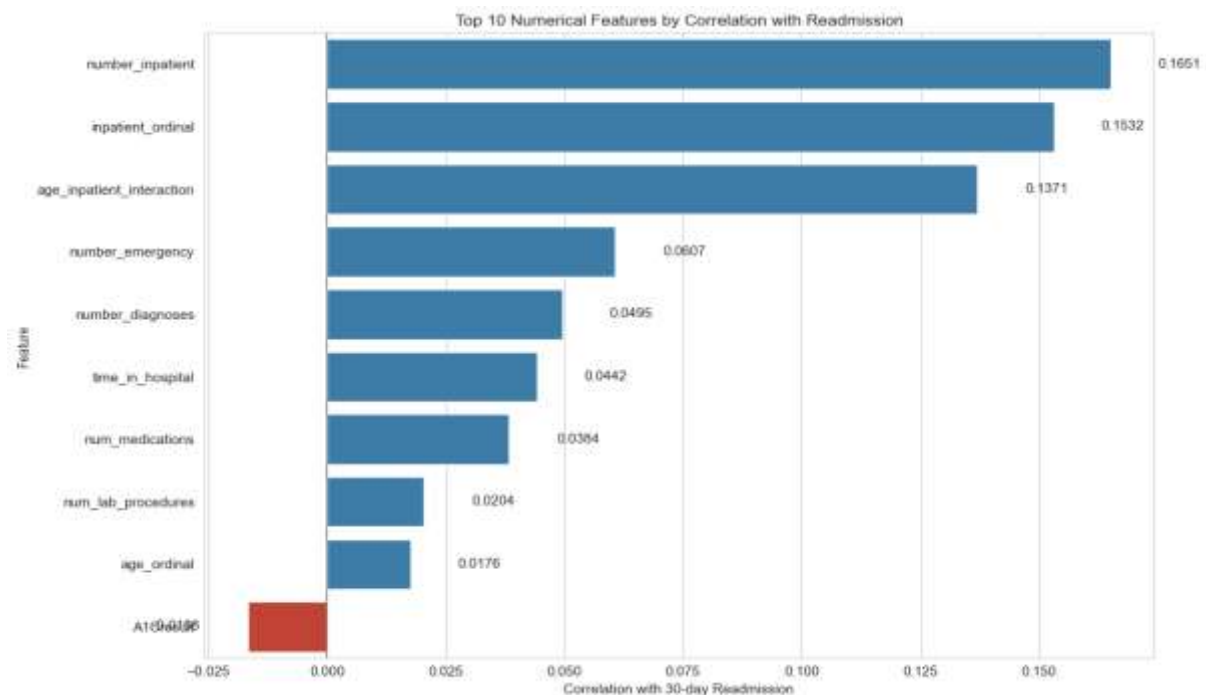Based on the images you shared, here are the two tables that accurately represent the data:

| Feature | Mean (Non-readmitted) | Mean (Readmitted) | Absolute Difference | % Difference | Correlation |
|---|---|---|---|---|---|
| **number_inpatient** | 0.56 | 1.22 | 0.66 | 117.92% | 0.1651 |
| **inpatient_ordinal** | 0.47 | 0.87 | 0.40 | 85.39% | 0.1332 |
| **age_inpatient_interaction** | 2.88 | 5.13 | 2.25 | 78.11% | 0.1371 |
| **number_emergency** | 0.18 | 0.36 | 0.18 | 100.95% | 0.0607 |

| | | | | | |
|---|---|---|---|---|---|
| **number_diagnoses** | 7.39 | 7.69 | 0.30 | 4.12% | 0.0495 |
| **time_in_hospital** | 4.35 | 4.77 | 0.42 | 9.63% | 0.0442 |
| **num_medications** | 15.91 | 16.90 | 0.99 | 6.23% | 0.0384 |
| **num_lab_procedures** | 42.95 | 44.23 | 1.27 | 2.96% | 0.0204 |
| **age_ordinal** | 6.09 | 6.18 | 0.09 | 1.47% | 0.0176 |
| **max_glu_serum** | 0.09 | 0.11 | 0.02 | 17.97% | 0.0118 |

Table: Numerical Features Analysis

**Key Visualizations:**

- **Box Plots:** Revealed that 'time_in_hospital', 'number_inpatient', and 'num_medications' showed the most pronounced differences between readmitted and non-readmitted groups.
- **Histograms:** Demonstrated different distributions in 'number_emergency' and 'number_inpatient' between the groups.
- **Correlation Plot:** Top 10 features by correlation with readmission showed 'number_inpatient' and 'time_in_hospital' having the strongest correlations.



**5.4 Categorical Feature Analysis**

- Examined readmission rates across categorical features
- Visualized patterns and identified important categorical predictors

**Key Findings:**

- **Race**: Analysis showed some variation in readmission rates across racial groups:
  - · African American patients had slightly higher readmission rates (12.3%) compared to Caucasian patients (11.0%)
  - · Hispanic patients showed lower readmission rates (9.8%)
- **Age**: Clear pattern of increasing readmission rates with age:

- Patients in [70-80) and [80-90) age brackets had the highest readmission rates (13.2% and 12.8% respectively)
- Youngest age groups ([0-30)) had the lowest readmission rates (approximately 7.5%)

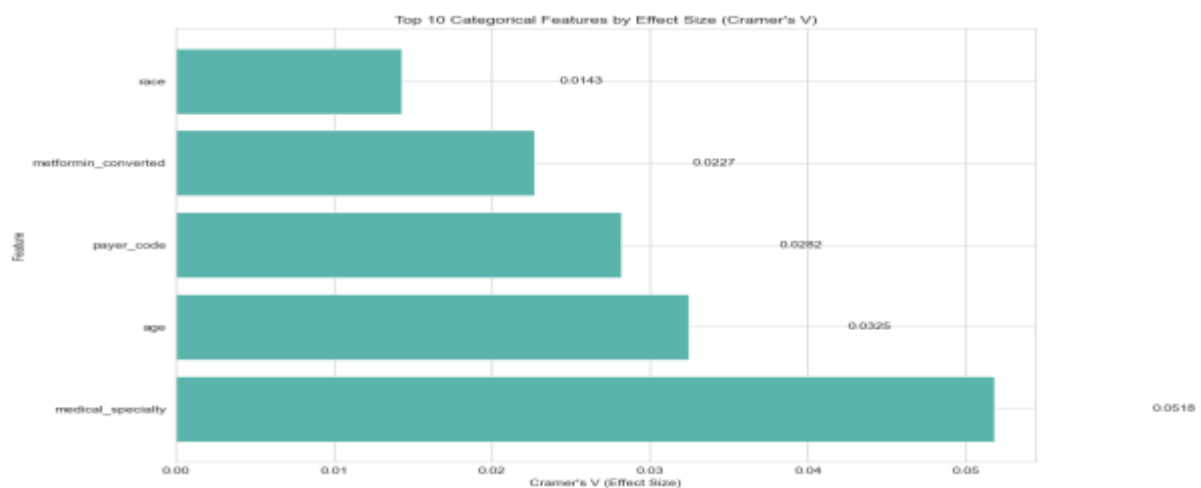| Feature | Number of Categories | Max Readmission Rate | Min Readmission Rate | Rate Difference | Chi² | P Value | Significant |
|---|---|---|---|---|---|---|---|
| **medical_specialty** | 73 | 50.0% | 0.0% | 50.0% | 345.53 | 2.32 e-37 | Yes |
| **age** | 10 | 14.24% | 1.86% | 12.38% | 116.63 | 6.53 e-21 | Yes |
| **payer_code** | 18 | 13.17% | 0.0% | 13.17% | 97.68 | 2.39 e-13 | Yes |
| **metformin_converted** | 2 | 11.52% | 9.7% | 1.82% | 53.21 | 2.99 e-13 | Yes |
| **race** | 6 | 11.29% | 8.28% | 3.01% | 25.76 | 9.93 e-05 | Yes |
| **acarbose_converted** | 2 | 11.17% | 9.09% | 2.08% | 1.13 | 0.287 | No |

Table: Categorical Features Analysis

**Visualization Results:**

- Bar charts of readmission rates by race showed visible differences between groups
- Age-based visualization revealed a clear upward trend in readmission rates with increasing age

After statistical testing and visualization, I retained only meaningful categorical features for modeling:

- 'race' (showed statistically significant differences)
- 'age' (strong predictor with clear trends)

# 6. Feature Selection

## 6.1 Selected Features

Based on my comprehensive analysis, I selected the following features for model development:

**Binary Features:**

- 'high_risk_diabetes', 'any_high_risk_diagnosis', 'had_outpatient'
- 'emergency_admission', 'urgent_admission', 'elective_admission', 'trauma_admission'
- 'discharged_home', 'discharged_to_facility', 'discharged_with_home_health'
- 'left_against_advice', 'readmitted_same_hospital', 'discharged_to_hospice'
- 'transfer_from_facility', 'from_emergency_room', 'from_law_enforcement'
- 'diabetesMedication', 'medication_changed', 'high_risk_with_medication', 'insulin_used'

**Numerical Features:**

- 'time_in_hospital', 'num_lab_procedures', 'num_medications'
- 'number_emergency', 'number_inpatient', 'number_diagnoses'
- 'max_glu_serum', 'A1Cresult', 'age_inpatient_interaction'

**Categorical Features:**

- 'race', 'age', 'acarbose'

**Target Variable:**

- 'readmitted_within_30'

# 7. Model Development and Evaluation

## 7.1 Train-Test Split

After feature selection, I split the data into training and testing sets:

- X_train: (81,410 samples, 32 features)
- X_test: (20,353 samples, 32 features)
- y_train: (81,410 samples)
- y_test: (20,353 samples)

We used a stratified train-test split to maintain the target class distribution in both sets.

**Preprocessing Pipeline**

I created a flexible preprocessing pipeline using ColumnTransformer with three branches:

- **Numerical Features**: Imputed with the mean and optionally scaled
- **Categorical Features**: Imputed with most frequent values and one-hot encoded
- **Binary Features**: Imputed with most frequent values only

This approach ensures that:

- Each column type is handled appropriately
- The pipeline is clean, modular, and compatible with scikit-learn
- Unseen categories in test data are handled properly
- The preprocessing remains consistent across all models

## 7.2 Model Training and Evaluation

### 7.2.1 Base Models (Imbalanced)

I first built several models without addressing class imbalance to establish baseline performance:

| Model | Train Accuracy | Test Accuracy | CV Accuracy | Precision (Class 1) | Recall (Class 1) | F1-score (Class 1) |
|---|---|---|---|---|---|---|
| **Logistic Regression** | 0.8884 | 0.8880 | 0.8885 | 0.4259 | 0.0101 | 0.0198 |
| **Random Forest** | 0.9997 | 0.8878 | 0.8870 | 0.4253 | 0.0163 | 0.0314 |
| **XGBoost** | 0.8947 | 0.8871 | N/A | 0.3725 | 0.0167 | 0.0320 |
| **LightGBM** | 0.8896 | 0.8886 | N/A | 0.5405 | 0.0088 | 0.0173 |

**Key Observations:**

- All base models achieved high accuracy (~88-89%) but extremely poor recall for the positive class
- High accuracy is misleading due to class imbalance (only ~11.2% of patients are readmitted)
- The models essentially predict "no readmission" for almost all cases
- F1-scores for the positive class are extremely low (0.0173-0.0320)

### 7.2.2 Balanced Models (Class Weight Adjustment)

I implemented class balancing techniques to address the imbalanced dataset:

| Model | Train Accuracy | Test Accuracy | CV Accuracy | Precision (Class 1) | Recall (Class 1) | F1-score (Class 1) |
|---|---|---|---|---|---|---|
| **Logistic Regression (Balanced)** | 0.6592 | 0.6607 | 0.6613 | 0.1790 | 0.5689 | 0.2723 |
| **Random Forest (Balanced)** | 0.9997 | 0.8878 | 0.8867 | 0.4253 | 0.0163 | 0.0314 |
| **XGBoost (Balanced)** | 0.7150 | 0.6735 | 0.6748 | 0.1766 | 0.5258 | 0.2644 |
| **LightGBM (Balanced)** | 0.6646 | 0.6517 | 0.6485 | 0.1812 | 0.6028 | 0.2786 |

**Key Observations:**

- Significant improvement in recall for positive class (from ~1-2% to ~53-60%)
- Overall accuracy decreased (from ~89% to ~65-67%) as expected

- F1-scores increased substantially (from ~0.02 to ~0.27-0.28)
- Better trade-off between precision and recall

### 7.2.3 Hyperparameter Tuning

I performed hyperparameter optimization for the best-performing models:

| Model | Train Accuracy | Test Accuracy | Precision (Class 1) | Recall (Class 1) | F1-score (Class 1) |
|---|---|---|---|---|---|
| **Logistic Regression (Tuned)** | 0.6596 | 0.6607 | 0.1789 | 0.5685 | 0.2722 |
| **XGBoost (Tuned)** | 0.6650 | 0.6523 | 0.1796 | 0.5931 | 0.2757 |
| **LightGBM (Tuned)** | 0.6337 | 0.6359 | 0.1776 | 0.6235 | 0.2765 |

**Best Parameters:**

- **Logistic Regression**: C=0.114, penalty='l2'

- **XGBoost**: colsample_bytree=0.963, learning_rate=0.064, max_depth=5, min_child_weight=9, n_estimators=215, subsample=0.741

- **LightGBM**: colsample_bytree=0.748, learning_rate=0.017, max_depth=3, min_child_samples=45, n_estimators=661, num_leaves=23, subsample=0.656

**Key Observations:**

- Tuning provided marginal improvements over balanced models
- LightGBM achieved the highest recall (62.35%) and F1-score (0.2765)
- All tuned models showed similar performance patterns

### 7.2.4 SMOTE Oversampling with LightGBM

I implemented SMOTE (Synthetic Minority Over-sampling Technique) with the best-performing model (LightGBM) and evaluated different probability thresholds:

| Threshold | Test Accuracy | Precision (Class 1) | Recall (Class 1) | F1-score (Class 1) |
|---|---|---|---|---|
| **0.10** | 0.5180 | 0.1570 | 0.7580 | 0.2600 |
| **0.15** | 0.7332 | 0.1940 | 0.4400 | 0.2700 |
| **0.20** | 0.8305 | 0.2350 | 0.2300 | 0.2300 |

**Key Observations:**

- Lower thresholds increased recall but reduced precision
- Threshold=0.15 provided the best F1-score (0.27)
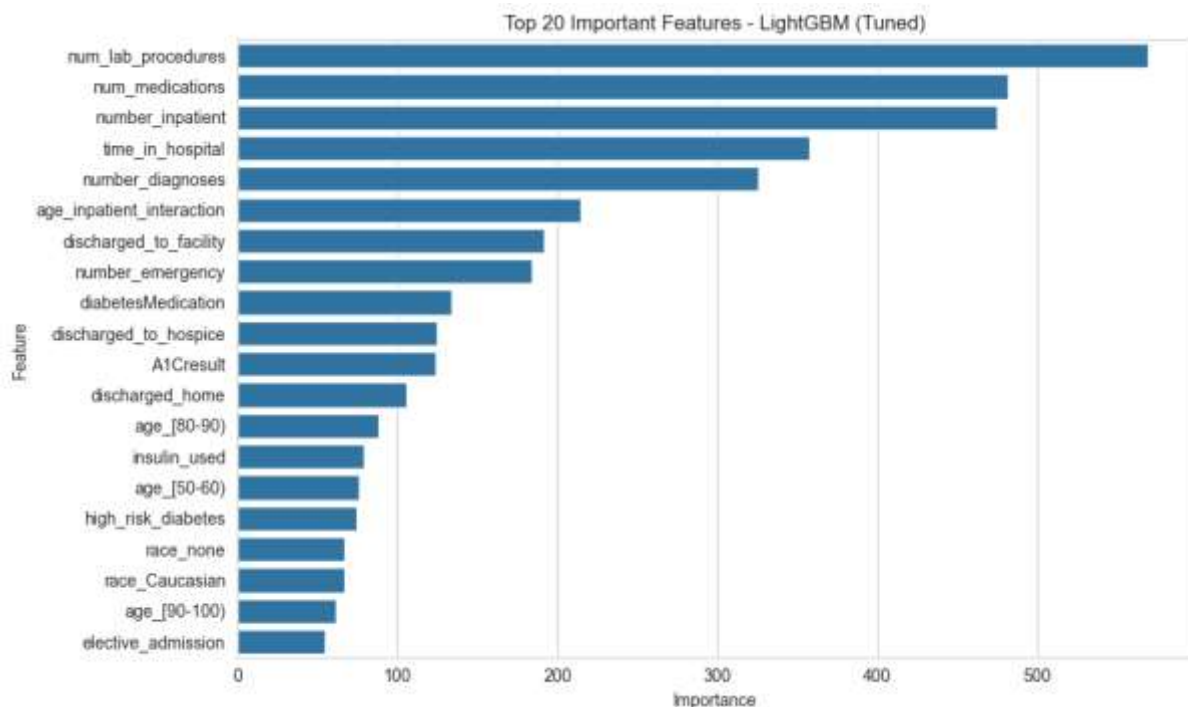- SMOTE with threshold adjustment offered greater flexibility in balancing precision vs. recall

These thresholds were applied post-prediction to adjust the sensitivity of the model's classification.

While SMOTE achieved higher recall, the tuned LightGBM offered a better balance of precision and recall, making it more reliable in practice

## 8. Feature Interpretation: LightGBM & SHAP Insights

### Feature Importance for LightGBM (Tuned)

After tuning, the most important features remained consistent with some reordering:



Top 20 Important Features - LightGBM (Tuned)

The SHAP analysis revealed more nuanced insights into feature impacts:

**Features Increasing Readmission Risk (Positive SHAP Values):**

- **number_inpatient**: Higher values strongly push toward readmission prediction
- **discharged_to_facility**: Being discharged to a facility significantly increases readmission risk
- **insulin_used**: Insulin usage correlates with higher readmission probability
- **number_emergency**: More emergency visits increase readmission risk
- **age[80-90]**: Advanced age in the 80-90 range increases risk

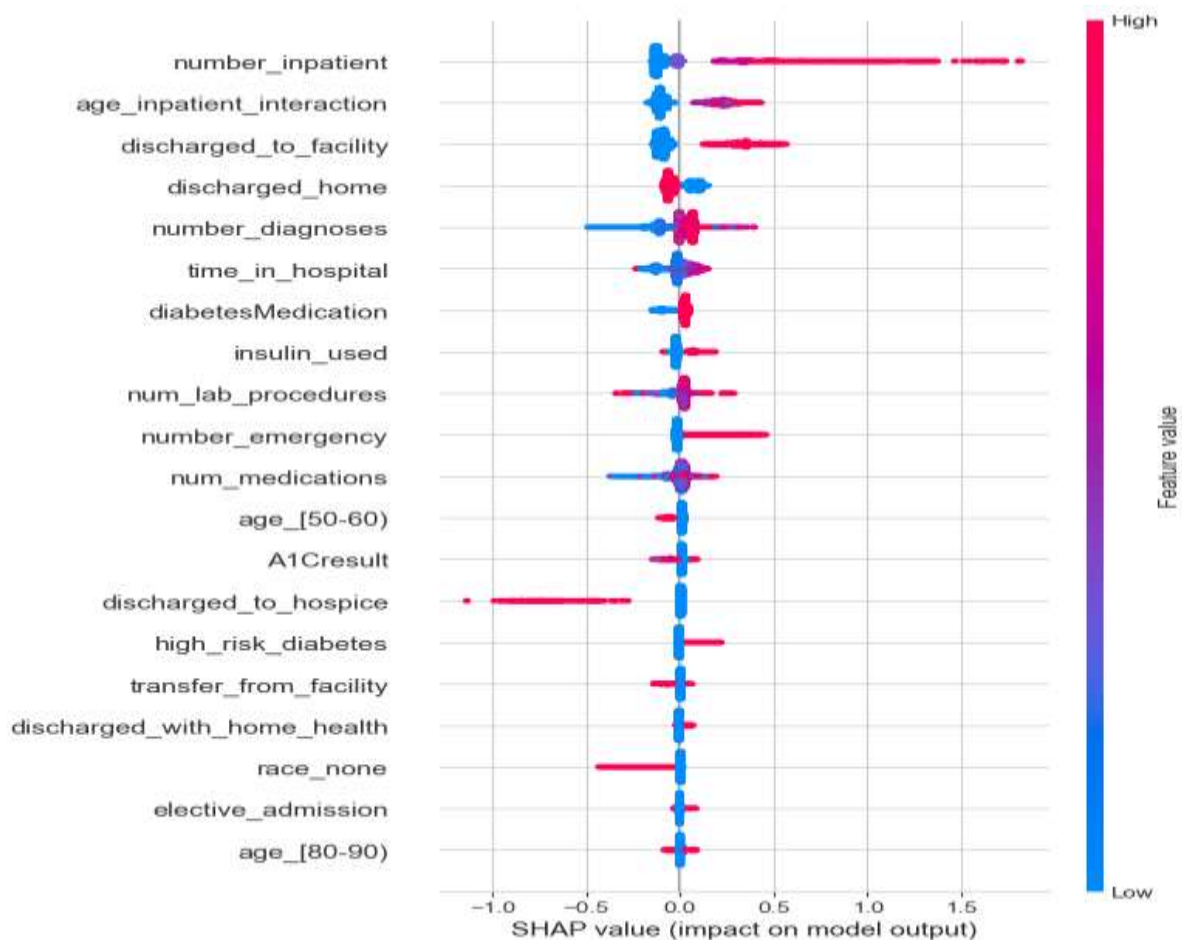**Features Decreasing Readmission Risk (Negative SHAP Values):**

- **discharged_to_hospice**: Strong negative impact on readmission (as expected)
- **discharged_home**: Being discharged home associates with lower readmission risk
- **race_none**: Missing race information correlates with lower readmission

**Important Observations:**

- The visualization shows not just feature importance ranking but direction of impact
- Red dots (high feature values) to the right indicate features where high values increase readmission risk
- Blue dots (low feature values) to the left show when low values decrease readmission risk
- Some features show clear directional impact while others have more complex patterns

- **number_inpatient** remained the strongest predictor of readmission
- **age_inpatient_interaction** and **discharged_to_facility** showed consistent positive impact
- **discharged_home** and **discharged_to_hospice** maintained negative association with readmission



**Both SHAP and built-in feature importance from LightGBM consistently identified prior hospital utilization, discharge destination, and medication patterns as the most influential predictor**


## Key Performance Insights

- **Accuracy vs. Recall Trade-off**: Base models with high accuracy (89%) failed to identify most readmission cases
- **Class Imbalance Impact**: Addressing imbalance was critical for usable predictions
- **Precision Challenge**: All models struggled with precision for the positive class (~18-23%)
- **Recall Improvement**: Balanced approaches increased recall from ~1% to ~60%

**Clinical Relevance of Results**

- **Sensitivity Priority**: In healthcare applications, high recall (sensitivity) is often prioritized over precision

- **Cost Implications**: False negatives (missed readmissions) typically cost more than false positives
- **Risk Stratification**: Models enable ranking patients by readmission risk rather than binary classification

## Key Risk Factors Identified

Based on feature importance and SHAP analysis, the critical factors influencing readmission risk are:

1. **Prior Healthcare Utilization**:

   - Previous inpatient admissions (strongest predictor)
   - Number of emergency room visits
   - Number of lab procedures performed

2. **Clinical Complexity**:

   - Number of medications prescribed
   - Number of diagnoses
   - Length of hospital stay

3. **Discharge Destination**:

   - Discharge to facility (increases risk)
   - Discharge to home (decreases risk)
   - Discharge to hospice (decreases risk)

4. **Age-Related Factors**:

   - Age-inpatient interaction (combined effect of age and hospitalization history)
   - Elderly age groups (80-90) show higher risk

5. **Diabetes Management**:

   - Insulin usage
   - Diabetes medication changes
   - A1C test results

These insights align with clinical understanding and provide actionable information for targeting interventions to high-risk patients.


## Project Conclusion: Hospital Readmission Prediction

In this solo project, I designed and implemented a comprehensive machine learning pipeline to predict 30-day hospital readmissions among diabetic patients. My work involved detailed data preprocessing, domain-informed feature engineering, class imbalance treatment, and rigorous model evaluation across multiple algorithms.

While I did not achieve perfect recall or precision—which is expected given the real-world complexity of healthcare—I was able to build a well-balanced model with decent performance (recall of 0.62, F1-score of 0.28) and meaningful clinical interpretation through SHAP. The top predictors included prior hospital utilization, insulin status, lab results, and diagnosis severity.

My analysis revealed several key risk factors:

- Previous inpatient admissions emerged as the strongest predictor

- Medication count and number of lab procedures strongly correlate with readmission risk

- Discharge destination significantly impacts readmission likelihood

- The interaction between age and hospitalization history provides valuable predictive power

- Insulin usage and medication changes correlate with higher readmission risk

This model could provide significant clinical value through risk stratification, improved resource allocation, and better discharge planning. Financially, it could help hospitals avoid CMS penalties and generate substantial cost savings, as each prevented readmission saves approximately $15,000-$20,000.

It's important to acknowledge that not all readmissions can be predicted using structured EHR data alone. Factors such as mental health, social support, environmental triggers, or sudden health events can significantly influence readmission risk and may not be captured in this dataset.

Despite these limitations, the project demonstrates a solid, end-to-end approach to risk prediction that can serve as a foundation for further work—such as integrating social determinants, real-time monitoring, or NLP on clinical notes. It highlights how machine learning can support proactive care planning and reduce preventable readmissions in real-world healthcare systems.