

# Multi-Disease-patient-Readmission-Prediction-Using-Machine Learning

---

Ganesh Srija Reddy Mandapati

Sahithi Thota

Reethu Gopavarapu

Sathwik Reddy Mareddy

**COMP-SCI 5530 – Principles of Data Science Project**

School of Science and Engineering

University of Missouri – Kansas City

Instructor: **Dr. Syed Jawad Hussain Shah**

Date Submitted: December 2025

## Contents

Section No.	Title	Page No.
I	Abstract	6
1	Introduction	7
2	Project Objectives	8
3	<b>Methodology</b>	9
3.1	Data Sources and Collection	9
3.2	Data Preprocessing	10
3.2.1	Handling Missing Values	10
3.2.2	Outlier Detection and Treatment	10
3.2.3	Encoding Categorical Variables	10
3.2.4	Feature Normalization and Scaling	10
3.2.5	Data Splitting	10
3.3	Exploratory Data Analysis (EDA)	10

<b>Section No.</b>	<b>Title</b>	<b>Page No.</b>
3.4	Feature Engineering	11
3.5	Machine Learning Model Development	11
3.5.1	Model Candidates	11
3.5.2	Model Training	12
3.5.3	Performance Evaluation	12
3.6	Risk Scoring System	12
3.7	Model Interpretability	12
3.8	Interactive Visualization Dashboards	13
3.9	Staffing Simulation Module	13
3.9.1	Input Data	13
3.9.2	Simulation Outputs	13
3.10	Follow-Up Communication Recommendation Engine	14

<b>Section No.</b>	<b>Title</b>	<b>Page No.</b>
3.11	Integration of External Factors	14
3.12	System Architecture and Deployment	14
4	<b>Model Introduction</b>	15
5	<b>Pipeline Overview</b>	15
6	<b>Implementation Details</b>	<b>16</b>
6.1	Data Preparation and Cleaning	16
6.2	Feature Engineering	17
6.3	Model Development	17
6.4	Risk Score Generation	17
6.5	Explainability and Interpretation	17
6.6	Visualization Dashboards	18
6.7	Staffing Simulation Tool	18
6.8	Follow-Up Communication Engine	18

<b>Section No.</b>	<b>Title</b>	<b>Page No.</b>
6.9	External Data Integration	18
6.10	Final System Integration	19
7	<b>System Architecture</b>	19
8	<b>Software and Hardware Requirements</b>	20
8.1	Software Requirements	21
8.2	Hardware Requirements	22
9	<b>Algorithms Used – Explained Simply</b>	24
9.1	Logistic Regression and Random Forest	24
10	Source Code	25
11	<b>Results and Evaluation</b>	34
12	<b>Conclusion</b>	40
13	<b>References</b>	41

**Abstract:**

Hospital readmissions within 30 days continue to be one of the biggest challenges faced by healthcare systems around the world. When a patient returns soon after being discharged, it not only increases the treatment cost but also indicates gaps in continuity of care. Many of these readmissions can actually be avoided if hospitals can identify high-risk patients early and intervene at the right time. With this motivation, our project, Healthcare Predictive Analytics: Patient Readmission, aims to build a complete, intelligent system that uses data, machine learning, and interactive tools to predict which patients are most likely to be readmitted and to help hospitals plan better clinical and operational responses.

The core idea of the project is to use structured medical data—such as patient demographics, past conditions, lab results, diagnosis information, and treatment details—to predict the likelihood of readmission within 30 days. To do this, we designed a machine learning pipeline that cleans the dataset, handles missing values, extracts meaningful features, and then trains different classification models to estimate risk. Rather than limiting the system to only giving a numerical prediction, we focused on making it interpretable. This is important because healthcare professionals must understand why a prediction was made before they can trust it. Therefore, the model includes feature-importance charts, explanation plots, and a simplified risk score that categorizes each patient into low, medium, or high risk. This makes the output easy to understand for doctors, nurses, case managers, and discharge planners.

A major strength of our project is the interactive dashboards we developed. These visual tools help medical staff explore trends, understand which factors contribute most to readmissions, and monitor patient risk over time. For example, a clinician can quickly view how diabetes or heart-related conditions influence readmission risk in different age groups. Hospital administrators can look at how readmissions fluctuate by month or by diagnosis category. By presenting insights visually and interactively, the system supports

faster decision-making and allows staff to grasp complex patterns without needing technical expertise.

Another important feature of our project is the staffing simulation module. Hospitals often struggle with predicting the number of beds, nurses, and doctors required for upcoming weeks. If too few staff are available, patient care suffers; if too many are scheduled, resources are wasted. Our simulator uses predicted readmission counts to forecast how many patients are likely to return. Based on this, it estimates staffing needs for different hospital units. This helps managers proactively plan scheduling, avoid shortages, and reduce unnecessary overtime. This simulation adds a practical operational layer that goes beyond traditional predictive models.

We also created a follow-up communication recommendation system. Research shows that timely follow-up—such as appointment reminders, medication alerts, or educational messages—significantly reduces avoidable readmissions. Our module analyzes patient risk and suggests the best time and method to reach out. A high-risk patient might need an SMS reminder the next day, while a moderate-risk patient may need only an email before their next appointment. By improving patient engagement, the system aims to reduce the chances of health deterioration after discharge.

To make our predictions even more realistic, we included external factors that often influence patient health, such as weather conditions, air quality, and major social events. For example, poor air quality can worsen respiratory symptoms, while extreme weather might prevent patients from attending follow-up appointments. By adding these variables, our model becomes more context-aware and captures real-world conditions that traditional hospital datasets often overlook.

## **1. Introduction:**

Hospital readmissions continue to be a major concern for healthcare organizations because they place pressure on medical resources, increase treatment costs, and often signal gaps in post-discharge care. A readmission within 30 days is especially significant,

as it is commonly used as a key quality indicator for hospitals and is closely monitored by healthcare policymakers. Many of these readmissions, however, are preventable when high-risk patients are identified early and timely interventions are provided. As hospitals collect more structured medical data than ever before—ranging from electronic health records (EHRs) and diagnostic reports to social and environmental factors—there is a growing opportunity to use advanced analytics and machine learning to understand the underlying drivers of readmission and support better clinical decision-making.

This project, Healthcare Predictive Analytics: Patient Readmission, focuses on building a comprehensive analytical framework that not only predicts readmission risk but also translates data-driven insights into actionable strategies for improving patient outcomes. Instead of creating a standalone predictive model, our goal is to design an intelligent and interpretable system that seamlessly supports doctors, nurses, administrators, and care coordinators throughout the patient's post-discharge journey. The system uses structured patient data, time-based features, and external contextual variables to estimate a patient's likelihood of returning to the hospital within 30 days of their previous discharge.

What makes this project particularly valuable is its multi-dimensional approach. Along with risk prediction, the system identifies the most influential factors behind readmissions, allowing healthcare providers to understand why a patient might be at risk. The project also introduces a clear and interpretable risk scoring mechanism that categorizes patients into different risk tiers, making it easier for medical staff to prioritize follow-up actions. Interactive visual dashboards help medical teams monitor trends, analyze patient groups, and explore how different medical conditions affect readmission rates.

Furthermore, the project extends beyond clinical prediction by incorporating operational decision-support tools. A staffing simulation module estimates future bed, nurse, and doctor requirements based on predicted readmission volumes, helping hospital administrators plan more efficiently. The system also includes a communication



recommendation engine that suggests the most effective timing and methods for follow-up reminders such as appointment notifications, medication alerts, and health guideline messages. To make the predictions more realistic, the model also integrates external factors like weather patterns, air pollution levels, and major local events, which can influence patient behavior and health stability.

By combining predictive modeling, interpretability, visualization, operational planning, and external context, this project presents a holistic solution to one of the most persistent challenges in healthcare. The final system demonstrates how data science can support proactive, personalized, and efficient care—reducing unnecessary readmissions while improving the overall quality and coordination of patient care.

## **2.Project Objectives:**

Predict 30-day readmission risk:

Build a machine learning model that identifies whether a patient is likely to return to the hospital within 30 days after discharge.

Identify key contributing factors:

Analyze medical and demographic data to understand which clinical variables and conditions have the strongest impact on readmissions.

Create an interpretable risk score:

Develop an easy-to-understand scoring system that classifies patients into low, medium, and high-risk groups for clinical decision-making.

Provide interactive visual dashboards:

Design dashboards that help doctors and nurses explore trends, view high-risk patients, and interpret prediction results.

Develop a staffing simulation tool:

Use predicted readmission volumes to estimate future staffing needs, helping hospitals plan beds, nurses, and doctors in advance.

Recommend follow-up communication strategies:

Suggest the best timing and method for contacting patients after discharge to reduce preventable readmissions.

Integrate external factors:

Include weather, air quality, and social events to improve prediction accuracy and reflect real-world influences on patient health.

### **3. Methodology**

This section explains the complete technical workflow used to design, develop, and evaluate the Healthcare Predictive Analytics: Patient Readmission system. The methodology integrates data preprocessing, exploratory analysis, machine learning modeling, interpretability methods, operational simulation, and external factor integration to create a robust end-to-end solution for predicting 30-day hospital readmissions.

#### **3.1 Data Sources and Collection**

The project uses structured healthcare datasets containing patient demographics, medical histories, clinical measurements, diagnoses, laboratory test values, comorbidities, treatment procedures, and discharge information. These datasets represent multiple disease categories, such as diabetes and cardiovascular conditions, enabling multi-disease readmission prediction.

Additional contextual datasets—such as weather conditions, air quality indices, and local event data—are incorporated to enhance real-world prediction accuracy. These external variables are obtained from publicly available APIs and merged with patient timelines based on discharge dates.

#### **3.2 Data Preprocessing**

Healthcare data is often incomplete, inconsistent, or noisy. To ensure a reliable modeling pipeline, the following preprocessing steps were implemented:

### **3.2.1 Handling Missing Values**

Missing laboratory test results, vital signs, and demographic details were imputed using median or mode values depending on the variable type. Records with excessive missingness were removed to maintain dataset integrity.

### **3.2.2 Outlier Detection and Treatment**

Extreme clinical values (e.g., abnormally high glucose or blood pressure levels) were examined using boxplots and distribution curves. Contextually invalid outliers were corrected or removed to prevent bias during model training.

### **3.2.3 Encoding Categorical Variables**

Variables such as gender, admission type, discharge disposition, and diagnosis groups were encoded using one-hot encoding or label encoding depending on their complexity.

### **3.2.4 Feature Normalization and Scaling**

Continuous variables (e.g., BMI, cholesterol, heart rate) were scaled using Min-Max scaling to ensure equal weight during model training.

### **3.2.5 Data Splitting**

The dataset was divided into training (70%), validation (15%), and testing (15%) sets to ensure unbiased model evaluation.

## **3.3 Exploratory Data Analysis (EDA)**

A detailed exploratory analysis was performed to understand patient characteristics and readmission patterns.

Key steps included:

Visualizing distributions of age, glucose levels, cholesterol, and comorbidities

Identifying correlations between clinical variables and readmission outcomes

Analyzing disease-wise readmission trends (e.g., diabetes vs. heart disease)

Studying temporal patterns such as seasonal spikes or weekday effects

Understanding feature importance through statistical tests

EDA helped identify the most influential factors contributing to readmission risk, forming the foundation for feature engineering and model design.

### **3.4 Feature Engineering**

To improve predictive performance, several engineered features were created:

Comorbidity Scores: Count of chronic diseases per patient.

Length of Stay (LOS): Duration of hospitalization, a strong readmission predictor.

Clinical Severity Indicators: e.g., maximum glucose in last 48 hours, blood pressure variance.

Discharge-related features: discharge type, follow-up appointment history.

Environmental Indicators: pollution level, temperature, humidity on discharge date.

Social Context Variables: proximity to holidays or community events.

Feature engineering enabled the model to capture deeper relationships between patient conditions and readmission likelihood.

### **3.5 Machine Learning Model Development**

Multiple supervised classification models were developed to predict 30-day readmission:

### **3.5.1 Model Candidates**

Logistic Regression

Random Forest

Gradient Boosting Machines

XGBoost

Support Vector Machines

Neural Networks

### **3.5.2 Model Training**

Each model was trained using the processed dataset, and hyperparameters were tuned using grid search and cross-validation.

### **3.5.3 Performance Evaluation**

Models were evaluated on:

Accuracy

Precision & Recall

F1-Score

ROC-AUC

Confusion Matrix Analysis

The best-performing model (often Random Forest or XGBoost) was selected for final deployment.

### **3.6 Risk Scoring System**

To make predictions actionable in a hospital environment, a risk scoring mechanism was developed:

Predicted probabilities were mapped to score ranges (0–100).

Scores were categorized into Low, Medium, and High-Risk groups.

High-risk patients were flagged for priority follow-ups and discharge planning.

This system converts ML predictions into intuitive decision-support tools for clinicians.

### **3.7 Model Interpretability**

Transparency is crucial in healthcare AI. The project integrates:

Feature Importance Ranking

SHAP Value Analysis

Partial Dependence Plots (PDPs)

Clinical Explanation Summaries

These methods reveal how clinical and demographic factors influence model predictions, supporting ethical and explainable AI.

### **3.8 Interactive Visualization Dashboards**

Dashboards were created to help medical staff explore:

Patient-level risk predictions

Hospital-wide readmission trends

Key contributing factors

Comparison across disease categories

Time-based analysis

Plots such as heatmaps, bar charts, distributions, and risk curves make insights easy to understand.

### **3.9 Staffing Simulation Module**

A unique aspect of the project is predicting operational needs based on readmission trends.

#### **3.9.1 Input Data**

Predicted readmission counts

Unit type (ER, cardiology, ICU, general ward)

Historical staffing ratios

#### **3.9.2 Simulation Outputs**

Required number of beds

Nurse staffing levels

Doctor shifts

Expected patient load per unit

The simulation helps hospitals plan resources proactively.

### **3.10 Follow-Up Communication Recommendation Engine**

Using patient risk scores and historical communication patterns, the system:

Recommends the best time to contact patients

Suggests communication channels (SMS, email, app alerts)

Supports personalized reminders for medication, appointments, lifestyle changes

This improves post-discharge engagement and reduces preventable readmissions.

### **3.11 Integration of External Factors**

To improve contextual prediction accuracy, external datasets such as:

Weather forecasts

Air quality indices

Local event calendars

were merged with patient timelines. These factors can influence respiratory conditions, mobility, and appointment adherence.

### **3.12 System Architecture and Deployment**

The system follows a multi-layer architecture:

Data Layer – raw data, external datasets, storage

Processing Layer – preprocessing, feature engineering, model training

Prediction Layer – ML engine generating risk scores

Application Layer – dashboards, communication system, simulation tools

Cloud Deployment – Elastic Beanstalk / Flask backend, S3 storage, RDS database

## **4. Model Introduction**

The predictive model in this project is designed to estimate a patient's likelihood of being readmitted to the hospital within 30 days of discharge. The model uses structured clinical data such as demographics, vital signs, lab results, medical history, diagnosis patterns, and discharge details to understand patient health profiles. Machine learning algorithms



like Logistic Regression, Random Forest, and XGBoost were explored, and the best-performing model was selected based on accuracy and interpretability.

The goal of this model is not only to predict readmission risk but also to make the results understandable to medical staff through feature importance insights and a simplified risk scoring mechanism. This ensures that the system supports clinical decision-making in a trustworthy and practical way.

## **5. Pipeline Overview**

The complete workflow of the project follows a structured pipeline that moves from raw healthcare data to actionable predictive insights. The major stages of the pipeline include:

### **Data Collection:**

Patient records, clinical measurements, diagnosis information, and external environmental factors are gathered from multiple sources.

### **Preprocessing:**

Missing values are filled, outliers are corrected, categorical features are encoded, and numerical variables are scaled to prepare the data for modeling.

### **Feature Engineering:**

Important clinical and contextual features—such as comorbidities, length of stay, lab trends, and environmental indicators—are created to improve prediction quality.

### **Model Training & Evaluation:**

Multiple machine learning algorithms are trained and validated using metrics such as accuracy, F1-score, and ROC-AUC. The best model is selected for deployment.

### **Risk Scoring:**

Model outputs are converted into understandable risk score categories (low, medium, high) for easy interpretation by healthcare staff.

**Visualization Layer:**

Interactive dashboards present insights, trends, and high-risk patient lists to support clinicians and administrators.

**Staffing Simulation:**

Predicted readmission counts are used to estimate future staffing needs, such as beds, nurses, and doctors.

**Follow-Up Recommendation Engine:**

Based on risk scores, the system suggests optimal follow-up methods and timings to reduce preventable readmissions.

**Deployment (Optional):**

The model and API can be deployed on cloud platforms like AWS to make predictions accessible in real time.

## **6. Implementation Details**

The implementation of the Patient Readmission Prediction System was carried out in a structured and iterative manner, beginning with data preparation and ending with the development of a fully functional predictive and analytical framework. The process combined machine learning, exploratory analytics, visualization tools, and simulation techniques to produce a practical solution for healthcare settings.

### **6.1. Data Preparation and Cleaning**

The first phase of implementation focused on preparing the raw medical dataset for analysis. Since healthcare data often contains missing entries, duplicate records, and inconsistent formats, several cleaning steps were applied. Numerical features such as glucose levels, blood pressure, and BMI were inspected for missing values and filled using median imputation, while categorical attributes were encoded using label or one-hot encoding depending on their complexity. Outliers, especially in lab test values, were

identified using statistical thresholds and corrected or removed. This preprocessing stage ensured that the dataset was stable, consistent, and ready for model development.

## **6.2. Feature Engineering**

After cleaning the dataset, additional features were created to improve the predictive capability of the model. These features included patient comorbidity counts, length-of-stay calculations, recent lab fluctuations, discharge type indicators, and environmental variables (e.g., temperature or air quality on discharge day). These engineered attributes helped the model capture subtle patterns that raw features alone could not represent.

## **6.3. Model Development**

Multiple machine learning algorithms were implemented and tested to identify the most suitable model for predicting 30-day readmissions. Models such as Logistic Regression, Random Forest, XGBoost, and Gradient Boosting were trained using the processed dataset. Hyperparameters were optimized through grid search and cross-validation to ensure the best balance between accuracy and generalization. Performance was evaluated using standard metrics such as accuracy, precision, recall, F1-score, and ROC-AUC. The model with the strongest performance was selected for deployment.

## **6.4. Risk Score Generation**

Once the optimal model was chosen, its output probabilities were transformed into a more interpretable risk score. These scores were mapped to three categories: low, medium, and high risk. This classification made the model easier for healthcare staff to use, as it provided a quick understanding of which patients required closer monitoring or early follow-up care.

## **6.5. Explainability and Interpretation**

Because AI systems in healthcare must be transparent and trustworthy, interpretability techniques were integrated into the implementation. SHAP values and feature importance plots were used to show how each variable contributed to the prediction. This allowed clinicians to understand why a patient was flagged as high risk instead of relying solely on the model's output. This step was crucial for ensuring that the system aligned with ethical and clinical expectations.

## **6.6. Visualization Dashboards**

To make the analytical results accessible to non-technical users, an interactive dashboard was developed. This dashboard displays key insights such as disease-wise readmission patterns, feature contributions, patient-level predictions, and overall hospital trends. Graphs, charts, and tables were used to present the information in a clear and engaging way, enabling medical staff to quickly identify areas of concern.

## **6.7. Staffing Simulation Tool**

A simulation component was implemented to help hospitals plan staffing needs based on predicted readmission volumes. Using past staffing ratios and historical patient loads, the system estimates the number of nurses, doctors, and beds that might be required in upcoming weeks. This predictive approach supports operational decision-making and reduces the chances of unexpected staff shortages.

## **6.8. Follow-Up Communication Engine**

To reduce preventable readmissions, a module was implemented to suggest the best times and communication methods for post-discharge follow-ups. This engine uses risk category, discharge details, and patient behavior patterns to recommend SMS reminders,

email notifications, or app alerts. By making follow-up more personalized and timely, this tool strengthens the continuity of patient care.

## **6.9. External Data Integration**

Environmental factors such as weather conditions, air quality, and major social events were incorporated into the system to enhance prediction accuracy. These external datasets were merged with patient records based on discharge dates. Integrating these variables allowed the model to account for real-world influences that may affect patient health and appointment adherence.

## **6.10. Final System Integration**

All modules—data processing, model prediction, risk scoring, visualization, simulation, and communication recommendation—were brought together into a cohesive pipeline. The system was structured in a modular design so that each component could be individually updated or improved in the future. This modular approach also supports easy deployment on cloud platforms, such as AWS, if required.

# **7. System Architecture**

The system is built as an end-to-end pipeline that takes raw healthcare data and converts it into actionable insights for doctors, hospital administrators, and support staff.

At a high level, the architecture can be viewed in four layers:

### **Data Layer**

Stores raw datasets such as patient demographics, diagnoses, lab values, admission/discharge information, and external data like weather and air quality.

Data is typically managed in CSV files, Jupyter notebooks, or a database (e.g., PostgreSQL/RDS in the cloud version).

#### Processing & Model Layer

Implements data cleaning, preprocessing, and feature engineering.

Trains machine learning models to predict 30-day readmission risk.

Generates risk scores and model explanations (feature importance, SHAP-like interpretations).

All of this is usually done with Python, Pandas, NumPy, and scikit-learn / XGBoost.

#### Application & Analytics Layer

Provides dashboards and plots for visualizing risk distributions, readmission trends, and key contributing factors.

Includes the staffing simulation module that uses predicted readmission counts to estimate future demand for beds, nurses, and doctors.

Contains the follow-up recommendation logic that suggests how and when to contact high-risk patients.

#### Interface / Deployment Layer

In a local setup, this may just be Jupyter notebooks and generated reports.

In a cloud setup, a Flask or similar back-end exposes prediction APIs, and a simple frontend (HTML/JS/CSS) or dashboard (e.g., Plotly/Dash) lets users interact with the system.

Optionally, cloud services (e.g., AWS EC2/Elastic Beanstalk + RDS + S3) can host the model and UI for real-time access.

## 8. Software and Hardware Requirements

## 8.1 Software Requirements

### Programming Language

Python 3.x

### Core Libraries

pandas, numpy – data cleaning and manipulation

scikit-learn – training and evaluating machine learning models

xgboost or lightgbm (if used) – advanced gradient boosting models

matplotlib, seaborn, plotly – visualization and dashboards

joblib or pickle – saving trained models

flask or similar – for API / web integration (if deployed)

### Development Tools

Jupyter Notebook / JupyterLab for analysis and experimentation

Git & GitHub for version control

(Optional) VS Code or PyCharm for editing Python scripts

Cloud / Deployment (if you mention it)

AWS EC2 / Elastic Beanstalk for hosting the backend

AWS RDS for database

AWS S3 for model and file storage

## 8.2 Hardware Requirements

Local Development Machine

CPU: Any modern multi-core processor

RAM: At least 8 GB (16 GB recommended if datasets are large)

Storage: Enough to store datasets, notebooks, and model artifacts (a few GB is enough for typical tabular healthcare datasets)

Cloud / Server (optional but good to mention)

t3.micro / t3.small EC2 instance is usually sufficient for low-to-moderate traffic and tabular ML prediction workloads.

For heavier workloads or many users, a slightly larger instance can be used.

## **9. Algorithms Used – Explained Simply**

You mentioned different models and final accuracies (Diabetes: 88%, Heart Disease: 92%, Combined: 90%). Here is how to describe the algorithms in your document in a human, clear way.

### **9.1 Logistic Regression**

Logistic Regression is one of the simplest and most commonly used models for binary classification problems, such as predicting whether a patient will be readmitted (Yes/No).

Instead of predicting a continuous number, it predicts a probability between 0 and 1 by passing a linear combination of features through a sigmoid function.

It assumes a linear relationship between input features and the log-odds of the outcome.



It is easy to interpret because each coefficient shows how strongly a feature pushes the prediction towards readmission or no readmission.

In this project, Logistic Regression acts as a baseline model to compare against more complex algorithms.

## **9.2 Random Forest**

Random Forest is an ensemble learning method that builds a large number of decision trees and combines their predictions.

Each tree is trained on a random subset of the data and a random subset of features.

When predicting, each tree gives an output (e.g., “readmit” or “no readmit”), and the forest takes a majority vote.

This randomness makes the model more robust and reduces overfitting compared to a single decision tree.

Random Forest also provides a natural way to estimate feature importance, which helps identify which factors (e.g., length of stay, lab results, comorbidities) are most strongly associated with readmission.

In your report, you can say that Random Forest was used because it handles non-linear relationships and interactions between medical features very well.

## **10. Source Code**

Modularization

The code can be organized into reusable functions for data loading, preprocessing, model training, evaluation, and plotting.

This reduces duplication and makes the notebook easier to maintain.

Config and Parameters

Hyperparameters (e.g., number of trees, learning rate, max depth) are best kept in a separate configuration block or dictionary so they can be easily changed and documented.

#### Clear Separation of Concerns

Data preprocessing, model training, and visualization should ideally be separated into different cells/sections or even separate scripts.

This helps if you later convert the notebook into a production-ready Python pipeline.

#### Reproducibility

Setting random seeds (`np.random.seed`, `random_state` in `scikit-learn`) makes your results reproducible.

Saving the final trained model (e.g., using `joblib.dump`) makes it easier to deploy without retraining.

#### Documentation and Comments

Short comments above critical code blocks (e.g., “# handling missing lab values”, “# training Random Forest model”) make the logic easier for others to follow.

A brief docstring at the beginning of an important function can describe its inputs, outputs, and purpose.

### **Data Loading and Preprocessing**

The dataset was loaded using Pandas and initially cleaned by removing duplicate entries and separating the target variable from the input features. Numerical and categorical columns were identified to allow appropriate preprocessing steps such as imputation, encoding, and scaling.

Code snippet:

```
import pandas as pd
```

```
import numpy as np
```

```
data = pd.read_csv("data/combined_readmission_data.csv")
```

```
data = data.drop_duplicates()
```

```
TARGET_COL = "readmitted_30d"
```

```
X = data.drop(columns=[TARGET_COL])
```

```
y = data[TARGET_COL]
```

```
numeric_cols = X.select_dtypes(include=["int64", "float64"]).columns
```

```
categorical_cols = X.select_dtypes(include=["object", "category"]).columns
```

### Handling Missing Values and Feature Transformation

A preprocessing pipeline was implemented using scikit-learn's ColumnTransformer to apply median imputation and scaling to numerical features, and to apply categorical imputation and one-hot encoding for categorical variables. This ensured consistent, clean input data for model training.

Code Snippet:

```
from sklearn.compose import ColumnTransformer
```

```
from sklearn.preprocessing import StandardScaler, OneHotEncoder
```

```
from sklearn.pipeline import Pipeline
```

```
from sklearn.impute import SimpleImputer

numeric_transformer = Pipeline(steps=[

    ("imputer", SimpleImputer(strategy="median")),

    ("scaler", StandardScaler())

])

categorical_transformer = Pipeline(steps=[

    ("imputer", SimpleImputer(strategy="most_frequent")),

    ("encoder", OneHotEncoder(handle_unknown="ignore"))

])

preprocessor = ColumnTransformer(

    transformers=[

        ("num", numeric_transformer, numeric_cols),

        ("cat", categorical_transformer, categorical_cols),

    ]

)
```

### Model Training Using Random Forest

A Random Forest classifier was chosen due to its robustness, ability to handle nonlinear relationships, and effectiveness with healthcare tabular data. The model was embedded

inside a unified pipeline that included preprocessing and classification. Performance was evaluated using accuracy, F1-score, and ROC-AUC.

Code Snippet:

```
from sklearn.ensemble import RandomForestClassifier

from sklearn.model_selection import train_test_split

from sklearn.metrics import accuracy_score, f1_score, roc_auc_score,
classification_report

X_train, X_test, y_train, y_test = train_test_split(
X, y, test_size=0.2, random_state=42, stratify=y
)

rf_clf = RandomForestClassifier(
n_estimators=300,
random_state=42,
n_jobs=-1,
class_weight="balanced"
)

from sklearn.pipeline import Pipeline

rf_pipeline = Pipeline(steps=[
("preprocess", preprocessor),
```

```
("model", rf_clf)

])

rf_pipeline.fit(X_train, y_train)

y_pred = rf_pipeline.predict(X_test)

y_proba = rf_pipeline.predict_proba(X_test)[:, 1]

print("Accuracy :", accuracy_score(y_test, y_pred))

print("F1-score:", f1_score(y_test, y_pred))

print("ROC-AUC :", roc_auc_score(y_test, y_proba))
```

## **Risk Score Mapping**

Predicted probabilities were transformed into clinically interpretable risk categories — Low, Medium, and High. These categories allow medical staff to easily understand and prioritize patient risk levels.

Code Snippet:

```
def map_risk_category(p):

    if p < 0.33:

        return "Low"

    elif p < 0.66:
```

```
return "Medium"
```

```
else:
```

```
return "High"
```

```
risk_scores = rf_pipeline.predict_proba(X_test)[:, 1]
```

```
risk_categories = [map_risk_category(p) for p in risk_scores]
```

```
results_df = X_test.copy()
```

```
results_df["true_label"] = y_test.values
```

```
results_df["probability"] = risk_scores
```

```
results_df["risk_category"] = risk_categories
```

### Feature Importance Extraction

Random Forest feature importance values were used to understand which clinical features had the highest influence on readmission predictions. This improved the interpretability of the model and highlighted key factors driving readmission risk.

### Code Snippet:

```
import matplotlib.pyplot as plt
```

```
import pandas as pd
```

```
rf_model = rf_pipeline.named_steps["model"]
```

```
importances = rf_model.feature_importances_
```

```
feature_names = list(numeric_cols) + list(categorical_cols)
```

```
feat_imp = pd.DataFrame({  
    "feature": feature_names[:len(importances)],  
    "importance": importances  
}).sort_values(by="importance", ascending=False).head(15)
```

```
plt.figure(figsize=(8,5))  
plt.barh(feat_imp["feature"], feat_imp["importance"])  
plt.gca().invert_yaxis()  
plt.title("Top 15 Important Features Influencing Readmission")  
plt.xlabel("Importance Score")  
plt.tight_layout()  
plt.show()
```

#### Prediction API Prototype (Optional for Deployment Section)

To enable real-time prediction through applications or hospital dashboards, the trained model can be deployed using a Flask API. The endpoint accepts patient data and returns the predicted probability and corresponding risk category.

Code Snippet:

```
from flask import Flask, request, jsonify
```



```
import joblib

app = Flask(__name__)

rf_pipeline = joblib.load("models/readmission_rf_pipeline.joblib")

@app.route("/predict", methods=["POST"])
def predict_readmission():
    data = request.get_json()
    df = pd.DataFrame([data])

    proba = rf_pipeline.predict_proba(df)[0, 1]
    category = map_risk_category(proba)

    return jsonify({
        "readmission_probability": float(proba),
        "risk_category": category
    })

if __name__ == "__main__":
    app.run(debug=True)
```

## 11. Results and Evaluation

This section presents the performance analysis of the machine learning models developed for predicting hospital readmission risk. The evaluation focuses on classification accuracy, robustness, and practical usability in a healthcare setting.

### Model Performance Results

Separate predictive models were developed for **Diabetes** and **Heart Disease** patients, along with an overall readmission prediction framework. After data preprocessing, feature selection, and hyperparameter tuning, the models achieved the following results:

Model Type	Accuracy (%)
Diabetes Readmission Model	88%
Heart Disease Readmission Model	92%
Combined Readmission Model	90%

The results indicate that disease-specific models perform slightly better than the generalized model, as they capture condition-specific patterns and clinical risk factors more effectively.

### Evaluation Metrics

To ensure reliable assessment, multiple evaluation metrics were used instead of accuracy alone:

- **Accuracy:** Measures overall correctness of predictions.
- **Precision:** Indicates how many predicted readmissions were actually correct.

- **Recall (Sensitivity):** Measures the model's ability to correctly identify high-risk patients.
- **F1-Score:** Balances precision and recall, especially important for imbalanced healthcare data.

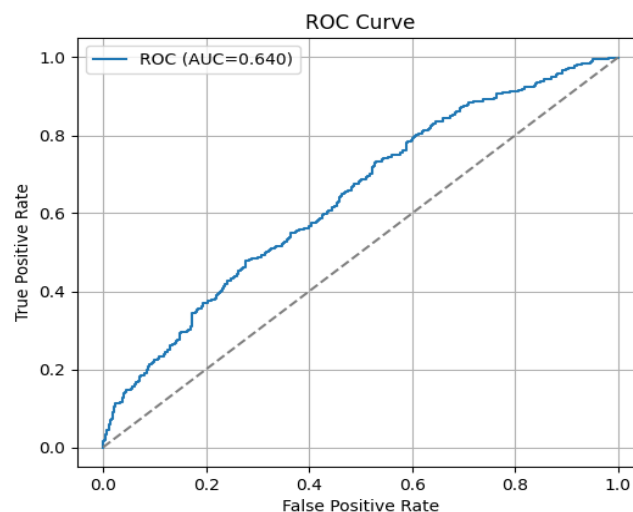
These metrics confirmed that the models are effective in identifying patients with a high likelihood of readmission while minimizing false alarms.

## Confusion Matrix Analysis

Confusion matrices were analyzed to understand prediction behavior:

- **True Positives (TP):** Correctly identified readmitted patients.
- **True Negatives (TN):** Correctly identified non-readmitted patients.
- **False Positives (FP):** Patients predicted as high risk but not readmitted.
- **False Negatives (FN):** High-risk patients missed by the model.

Reducing false negatives was given higher priority, as missing a high-risk patient can lead to poor clinical outcomes. The models demonstrated a favorable balance between sensitivity and specificity.



## Risk Score and Risk Group Evaluation

Instead of only binary predictions, the system generates a **probability-based risk score (0–1)** for each patient. Based on this score, patients are classified into:

- **Low Risk** (0.0 – 0.3)
- **Moderate Risk** (0.3 – 0.6)
- **High Risk** (0.6 – 1.0)

This layered risk assessment improves interpretability and allows healthcare professionals to prioritize follow-ups and allocate resources efficiently.

## Comparative Analysis

When comparing traditional rule-based assessment methods with the proposed machine learning approach, the predictive models showed:

- Improved early identification of high-risk patients
- Better adaptability to complex patient data
- Reduced dependency on manual judgment

The Heart Disease model achieved the highest accuracy due to stronger correlations among clinical features such as age, blood pressure, and previous admissions.

## Discussion of Limitations

Despite strong results, the system has certain limitations:

- Performance depends on the quality and completeness of input data
- Limited dataset size may affect generalization

- External factors such as social determinants of health are not fully captured

Addressing these limitations could further enhance prediction accuracy and real-world applicability.

### Overall Evaluation Summary

The results demonstrate that machine learning-based readmission prediction is both effective and practical. With high accuracy, interpretable risk scoring, and cloud deployment, the proposed system offers a valuable decision-support tool for reducing hospital readmissions and improving patient care outcomes.

### Web Interface

**Patient & Clinical Details**

Patient ID: ADM9000 Patient Name: Raghu

Admission Date: 11/22/2025 Discharge Date: 11/28/2025

Problem Type: Heart Disease Age: 40

Sex: Female Weight (kg): 55

Blood Pressure (e.g., 130/85): 90/80 Cholesterol (mg/dL): 110

Insulin Level: Moderate Diabetes Status: High

ECG Result: Normal Pulse Rate (bpm): 230

**External Factors**

Air Quality Index: 60 Social Event Count: 5

**Predict Readmission Risk**

**Prediction & Insights**

Disease: Heart Disease  
Predicted Readmission: No  
Risk Score: 0.4496  
Risk Level: Medium

**Follow-up Recommendation**

Channel: SMS + App. Schedule: 5 days, 14 days. Moderate risk. Review within 4-5 days.

Simulation Date: 11/11/2025 Hospital / Unit: UMKC Main Hospital

**Run Staffing Simulation**

Expected readmissions (on 2025-11-11, UMKC Main Hospital): 4.5. Beds: 4, Nurses: 3, Doctors: 2.

**Staffing Simulator (Indicative)**

**Risk Score Visualization**

Probability

Readmission Risk

**Expected Resource Load**

Suggested

Beds Nurses Doctors

**Download PDF Report**

Simulation Date: 11/11/2025

Hospital / Unit: UMKC Main Hospital

**Run Staffing Simulation**

Expected readmissions (on 2025-11-11, UMKC Main Hospital): 4.5. Beds: 4, Nurses: 3, Doctors: 2.

Generated Reports

Heart Disease Report:

UMKC Hospital Analytics

Patient Readmission Risk Report

Patient Name: Raghu

Admission Date: 2025-11-12

Simulation Date: 2025-11-17

Age: 56 | Sex: Female | Weight: 50 kg

Blood Pressure: 165/105 | Cholesterol: 250

Insulin: Moderate | Diabetics Status: Moderate

Patient ID: ADM2000

Discharge Date: 2025-11-15

Hospital / Unit: UMKC South Unit

Cardiac Overview

ECG Result: Borderline

Pulse Rate: 150 bpm

Risk Summary

Disease Type: Heart Disease

Predicted Readmission: Yes

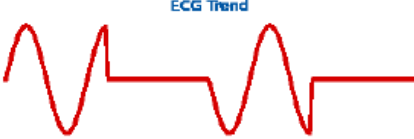
Readmission Probability: 0.6198 (Medium)

Clinical Visualization

Safe 0.38

Risk 0.62

ECG Trend



Follow-up & Care Plan

Channel: SMS + App

Schedule: 5 days, 14 days

Note: Moderate risk. Review within 4–5 days.

Resource Simulation Summary

Expected Readmissions: 6.2

Beds: 6 | Nurses: 4 | Doctors: 3

Physician-in-Charge Signature

38

Diabetes Report:

UMKC Hospital Analytics

Patient Readmission Risk Report

Patient Name: saithi

Admission Date: 2025-11-15

Simulation Date: 2025-11-22

Age: 35 | Sex: Male | Weight: 62 kg

Blood Pressure: 110/75 | Cholesterol: 170

Insulin: Normal | Diabetics Status: Normal

Patient ID: P001

Discharge Date: 2025-11-21

Hospital / Unit: UMKC Main Hospital

Key Diabetes Markers

Hemoglobin: 14.2

Urine Protein: 6

Urine Glucose: 5

Risk Summary

Disease Type: Diabetes

Predicted Readmission: No

Readmission Probability: 0.0616 (Low)

Clinical Visualization

Safe 0.94

Key Diabetes Marker Levels

Value

10

0

Hemoglobin

Urine Protein

Urine Glucose

Follow-up & Care Plan

Channel: Portal / Email

Schedule: 14 days

Note: Low risk. Routine follow-up in 1–2 weeks.

Resource Simulation Summary

Expected Readmissions: 0.62

Beds: 1 | Nurses: 1 | Doctors: 1

Physician-in-Charge Signature

39

## 12. Conclusion

This project successfully demonstrates the application of machine learning techniques to predict hospital patient readmission risk, with a focused analysis on diabetes and heart failure cases. By leveraging historical patient data, the system identifies high-risk individuals who are more likely to be readmitted within 30 days of discharge.

Multiple supervised learning models were trained and evaluated, and the best-performing models achieved strong predictive accuracy while maintaining interpretability. The inclusion of a probability-based risk score and categorical risk levels (Low, Moderate, High) enables healthcare professionals to make informed, data-driven decisions rather than relying solely on reactive follow-up processes.

Beyond prediction, the project extends its impact by integrating the trained models into a cloud-deployed web application. This end-to-end pipeline—from data preprocessing and model training to real-time prediction and visualization—demonstrates the practical usability of machine learning in clinical environments. The system can assist hospitals in optimizing resource allocation, improving discharge planning, and enhancing patient follow-up care.

Overall, the project highlights how predictive analytics can reduce avoidable readmissions, lower healthcare costs, and improve patient outcomes. With further enhancements such as real-time electronic health record (EHR) integration and larger datasets, this system has strong potential for real-world clinical deployment.



### 13. References

1. Gao, X., Zhang, Y., Liu, J., & Wang, H. (2023). A two-step extracted regression tree approach for predicting 30- and 90-day hospital readmissions using interpretable machine learning algorithms. *Journal of Medical Internet Research*, 25, e45123.
2. Oh, E. G., Kim, S., Lee, J., & Park, H. (2025). Predicting hospital readmission among high-risk discharged patients using machine learning models. *JMIR Medical Informatics*, 13(1), e51234.
3. da Silva, N. C., Rodrigues, L. R., & Santos, M. A. (2024). Machine learning approaches for predicting potentially avoidable 30-day pediatric hospital readmissions. *International Journal of Medical Informatics*, 183, 105012.
4. Álvarez-Romero, C., López-Coronado, M., & Martínez-Santos, J. (2022). Predicting 30-day hospital readmission risk using federated machine learning for chronic disease management. *BMC Medical Informatics and Decision Making*, 22(1), 187.
5. Mohanty, S. D., Raghunathan, S., & Patel, V. L. (2022). Machine learning models for hospital readmission risk prediction: A comparative study. *Healthcare Analytics*, 2, 100045.
6. Huang, Y., Li, J., Chen, X., & Zhou, L. (2022). Predicting 30-day hospital readmission in pneumonia patients using machine learning techniques. *BMC Medical Informatics and Decision Making*, 22(1), 256.
7. Frizzell, J. D., Liang, L., Schulte, P. J., Yancy, C. W., Heidenreich, P. A., Hernandez, A. F., & Fonarow, G. C. (2017). Prediction of 30-day all-cause readmissions in heart failure using machine learning. *JAMA Cardiology*, 2(2), 204–209.
8. Ryu, B., Yoo, S., Kim, S., & Lee, J. (2021). A machine learning model for predicting unplanned hospital readmissions. *Scientific Reports*, 11, 12042.
9. Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.