# HEALTH-CARE DATA CLEANING
## A PROJECT FILE
### for
## INTRODUCTION TO AI (AI201B)
### Session (2024-25)

**Submitted by**

**Mradul Tyagi**
**202410116100125**
**Nainsi Jain**
**202410116100128**
**Gargi Singh**
**202410116100072**

**Submitted in partial fulfilment of the**
**Requirements for the Degree of**

## MASTER OF COMPUTER APPLICATION

**Under the Supervision of**
**Mr. APOORV JAIN**
**Assistant Professor**



**Submitted to**

**DEPARTMENT OF COMPUTER APPLICATIONS**
**KIET Group of Institutions, Ghaziabad**
**Uttar Pradesh-201206**
**(APRIL- 2025)**

# **TABLE OF CONTENT**

# <u>INTRODUCTION</u>

## Background

Healthcare systems generate vast amounts of data daily, including patient records, medical imaging, laboratory test results, and prescription details. This data has immense potential for **predictive analytics**, which can help identify **high-risk patients, forecast disease progression, and improve clinical decision-making**. However, the **raw data is often noisy, inconsistent, and incomplete**, making it unsuitable for direct use in machine learning models. Poor data quality can lead to **biased predictions, inaccurate diagnoses, and inefficient resource allocation** in healthcare settings.

To leverage the full potential of healthcare data, **data preprocessing and feature engineering** play a crucial role. By cleaning and transforming raw data into a structured and meaningful format, we can improve the accuracy and reliability of predictive models. This project focuses on developing a systematic approach to **data cleaning, feature selection, and model building** to enhance predictive modeling for high-risk patient identification.

## Motivation

With the increasing reliance on **data-driven healthcare**, ensuring **data quality and integrity** is essential for making accurate medical predictions. Inconsistent, missing, or erroneous data can lead to misdiagnoses and ineffective treatment plans. Moreover, machine learning algorithms rely on high-quality data to produce meaningful results. A well-processed dataset allows for better generalization, improving model reliability across different patient populations.

This project is motivated by the need to **refine healthcare datasets** through proper data preprocessing techniques. The goal is to create a **robust dataset** that can be used to **predict patient risk levels more accurately**, enabling early interventions and improved treatment outcomes.

## Objectives

The primary objectives of this project are as follows:

- To **clean and preprocess** a healthcare dataset by handling missing values, duplicates, and inconsistencies.

- To conduct **Exploratory Data Analysis (EDA)** for understanding data distribution, patterns, and relationships.

- To apply **data transformation techniques**, such as normalization and categorical encoding, for better model performance.

- To handle **outliers** effectively to prevent skewed predictions.

- To perform **feature engineering** to create new, meaningful variables that improve model accuracy.

- To conduct **correlation analysis and feature selection** to retain only the most relevant features.

- To build **predictive models** using machine learning algorithms and evaluate their performance.

## Significance

The significance of this project lies in its potential to **improve healthcare decision-making** by ensuring that predictive models are trained on high-quality data. Accurate predictions can help healthcare providers make timely interventions, leading to **better patient outcomes, reduced healthcare costs, and more efficient resource allocation**. Additionally, a well-structured dataset minimizes bias in machine learning models, improving **fairness and reliability** in patient risk assessment.

By implementing **state-of-the-art data preprocessing and feature engineering techniques**, this project provides a **scalable framework** that can be applied to various healthcare datasets for predictive analytics, disease forecasting, and personalized medicine.

# <u>OVERVIEW</u>

Healthcare data analytics is a rapidly evolving field that relies heavily on **clean, structured, and well-preprocessed data** to enhance predictive modeling for disease diagnosis, risk assessment, and patient care. The accuracy and effectiveness of machine learning models depend largely on the quality of the input data, making **data cleaning, transformation, and feature engineering** essential steps in the analytical pipeline.

This project focuses on **healthcare data preprocessing**, covering key aspects such as:
- **Data Collection & Exploration** – Understanding raw data, identifying patterns, and detecting inconsistencies.
- **Data Cleaning & Handling Missing Values** – Addressing incomplete, incorrect, or redundant records.
- **Outlier Detection (IQR Method)** – Eliminating extreme values that can distort model predictions.
- **Normalization & Scaling** – Standardizing numeric values for improved model convergence.
- **Feature Engineering** – Creating new, meaningful variables to enhance predictive power.
- **Feature Selection** – Identifying the most relevant attributes to improve model performance.
- **Model Building & Evaluation** – Developing machine learning models and assessing their accuracy using real-world healthcare data.

**Importance of Healthcare Data Preprocessing**
Raw healthcare datasets often contain inconsistencies, such as:
- **Duplicate records** (e.g., multiple entries for the same patient).
- **Missing values** (e.g., unrecorded vital signs or test results).
- **Outliers** (e.g., extreme blood pressure readings due to measurement errors).
- **Unstructured data** (e.g., inconsistent formatting in categorical variables).

By applying **effective preprocessing techniques**, we ensure that the data is **accurate, consistent, and optimized for analysis**, ultimately leading to:
- Better disease prediction models.
- Enhanced clinical decision-making.
- More efficient healthcare management systems.
- Improved patient outcomes.

This report provides a comprehensive guide to **healthcare data cleaning and feature engineering**, demonstrating **best practices and advanced techniques** to transform raw data into a powerful asset for **machine learning and predictive analytics**.

# METHODOLOGY

The project is structured into the following steps to prepare and analyze healthcare data:

**1. Data Collection**

- The dataset consists of healthcare parameters such as **PatientID, Age, BloodPressure, SugarLevel, and Weight**.
- The data is loaded and explored to understand its structure, column names, and shape.

**2. Data Preprocessing**

- **Handling Missing Values:**
  - Numerical data: Missing values are replaced using the **mean** of the respective column.
  - Categorical data: Missing values are filled using the **mode** (most frequently occurring value).
- **Removing Duplicates:**
  - Duplicate records are identified and removed to maintain data integrity.
- **Outlier Detection and Handling:**
  - The **Interquartile Range (IQR) method** is applied to cap extreme values for Age, BloodPressure, SugarLevel, and Weight.
- **Data Normalization:**
  - **StandardScaler** is used to normalize numerical features, ensuring a uniform scale across variables.

**3. Exploratory Data Analysis (EDA)**

- **Visualizing data distributions** using histograms to identify trends and anomalies.
- **Correlation Analysis:**
  - A **heatmap** is generated to examine relationships between different healthcare features.

**4. Feature Engineering**

- **Age Grouping:** Patients are categorized as **Young, Middle-aged, or Senior**.
- **Feature Interaction:** A new feature, **BP_Sugar_Interaction** (BloodPressure $\times$ SugarLevel), is created.
- **Categorical Encoding:**
  - **SugarLevel** is categorized into **Low, Normal, and High** groups.
  - **One-hot encoding** is applied to convert categorical features into numerical form for machine learning.

**5. Feature Selection**

- **Using SelectKBest (ANOVA F-test) to choose the top 5 most relevant features**:
  - Age
  - BloodPressure
  - Weight
  - BP_Sugar_Interaction
  - SugarLevel_Category_High
- This step ensures that only the most important features are used for modeling, improving efficiency and performance.

**6. Data Splitting for Modeling**

- The dataset is divided into **training (80%) and testing (20%)** sets for machine learning.
- The target variable (**SugarLevel**) is converted into a **binary classification** (high risk vs. low risk) based on its mean value.

**B. Algorithms & Techniques Used**

**1. Interquartile Range (IQR) for Outlier Handling**

- The IQR method helps **identify and cap extreme values** in the dataset.
- Ensures that outliers do not significantly impact the predictive modeling process.

**2. StandardScaler for Data Normalization**

- Standardizes numerical features by transforming them to have a **mean of 0 and standard deviation of 1**.
- Ensures that all features contribute equally to the model.

**3. SelectKBest (ANOVA F-test) for Feature Selection**

- Evaluates and selects the **top-performing features** for model training.
- Improves **model efficiency and accuracy** by reducing noise in the data.

# CODE

**# IMPORT REQUIRED LIBRARIES**

```
import pandas as pd

import numpy as np

import seaborn as sns

import matplotlib.pyplot as plt

import warnings

from sklearn.preprocessing import StandardScaler

from sklearn.feature_selection import SelectKBest, f_classif

from sklearn.model_selection import train_test_split
```

**# SUPPRESS WARNINGS**

```
warnings.filterwarnings("ignore")
```

**# 1. LOAD THE HEALTHCARE DATASET**

```
file_path = "/content/healthcare_data.csv"

data = pd.read_csv(file_path)


print("Dataset Loaded Successfully!")

print(f"\nShape of Dataset: {data.shape}")

print(f"\nColumns in Dataset:\n{data.columns.to_list()}")
```

**# 2. DATA CLEANING STARTS HERE**

```
# CHECK FOR MISSING VALUES

print("\nChecking for Missing Values...")

missing_summary = data.isnull().sum()

print(missing_summary)
```

# FILL MISSING VALUES

# For numerical columns: fill with mean

```python
for col in data.select_dtypes(include=np.number).columns:
    data[col].fillna(data[col].mean(), inplace=True)


# For categorical columns: fill with mode
for col in data.select_dtypes(include="object").columns:
    data[col].fillna(data[col].mode()[0], inplace=True)


print("\nMissing Values Handled Successfully!")
```

# REMOVE DUPLICATES

```python
data.drop_duplicates(inplace=True)
print(f"\nDuplicates Removed! New Shape: {data.shape}")
```

# 3. HANDLE OUTLIERS USING IQR

```python
def handle_outliers(df, columns):
    for col in columns:
        Q1 = df[col].quantile(0.25)
        Q3 = df[col].quantile(0.75)
        IQR = Q3 - Q1
        lower_bound = Q1 - 1.5 * IQR
        upper_bound = Q3 + 1.5 * IQR
        df[col] = np.where(df[col] < lower_bound, lower_bound, df[col])
        df[col] = np.where(df[col] > upper_bound, upper_bound, df[col])
    return df


# Identify numeric columns to handle outliers
numeric_columns = ["Age", "BloodPressure", "SugarLevel", "Weight"]
data = handle_outliers(data, numeric_columns)
print("\nOutliers Handled Using IQR Method!")
```

# 4. NORMALIZE / SCALE DATA

scaler = StandardScaler()

data_scaled = pd.DataFrame(scaler.fit_transform(data[numeric_columns]),
columns=numeric_columns)


# Add PatientID back after scaling

data_scaled["PatientID"] = data["PatientID"].values

print("\nData Normalized Using StandardScaler!")


# 5. FEATURE ENGINEERING

# Create Age Group (Young, Middle-aged, Senior)

data_scaled["AgeGroup"] = pd.cut(data_scaled["Age"], bins=[0, 30, 50, 100], labels=["Young",
"Middle-aged", "Senior"])


# Interaction Feature: BloodPressure * SugarLevel

data_scaled["BP_Sugar_Interaction"] = data_scaled["BloodPressure"] * data_scaled["SugarLevel"]


# Binning SugarLevel into Low, Normal, High

data_scaled["SugarLevel_Category"] = pd.cut(data_scaled["SugarLevel"], bins=[0, 100, 140, 200],
labels=["Low", "Normal", "High"])


# One-hot encoding for categorical features

data_encoded = pd.get_dummies(data_scaled, columns=["AgeGroup", "SugarLevel_Category"],
drop_first=True)


# Drop unnecessary columns before feature selection

data_encoded.drop(columns=["PatientID"], inplace=True)

print("\nFeature Engineering Completed Successfully!")

# 6. CORRELATION ANALYSIS

```python
plt.figure(figsize=(8, 6))

sns.heatmap(data_encoded.corr(), annot=True, cmap="coolwarm", fmt='.2f')

plt.title("Correlation Heatmap (Smaller Size)")

plt.show()

print("\nCorrelation Analysis Completed!")
```

# 7. FEATURE SELECTION

```python
# Define X and y for feature selection

X_new = data_encoded.drop(columns=["SugarLevel"])

y_new = (data_encoded["SugarLevel"] > data_encoded["SugarLevel"].mean()).astype(int)


# Select top 5 features using SelectKBest (ANOVA F-test)

selector = SelectKBest(score_func=f_classif, k=5)

X_selected = selector.fit_transform(X_new, y_new)


# Get selected feature names

selected_features = X_new.columns[selector.get_support()]


print("\nTop 5 Selected Features: ", list(selected_features))

print("Feature Selection Completed Successfully!")
```

# 8. VISUALIZATIONS OF DISTRIBUTIONS

```python
plt.figure(figsize=(12, 8))

for i, col in enumerate(numeric_columns):

    plt.subplot(2, 2, i + 1)

    sns.histplot(data[col], kde=True, bins=15, color='skyblue')

    plt.title(f'Distribution of {col}')

plt.tight_layout()

plt.show()

print("\nData Distribution Visualized!")
```

# 9. SPLIT DATA FOR MODELING

```python
X = data_scaled.drop(columns=["PatientID"])

y = (data_scaled["SugarLevel"] > data_scaled["SugarLevel"].mean()).astype(int)


# Split data into training and testing sets

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


print(f"\nData Split Completed for Modeling!")

print(f"X_train Shape: {X_train.shape}")

print(f"X_test Shape: {X_test.shape}")

print(f"y_train Shape: {y_train.shape}")

print(f"y_test Shape: {y_test.shape}")
```

# 10. DATA SUMMARY

```python
data_summary = {

    "Outlier Handled Data": data.head(),

    "Normalized Data": data_scaled.head(),

    "Selected Features": list(selected_features),

    "X_train Shape": X_train.shape,

    "X_test Shape": X_test.shape,

    "y_train Shape": y_train.shape,

    "y_test Shape": y_test.shape

}


print("\nSUMMARY OF DATA CLEANING PROCESS COMPLETED!")
```

# OUTPUT

**Dataset Loaded Successfully!**

**Shape of Dataset: (20, 5)**

**Columns in Dataset:**

**['PatientID', 'Age', 'BloodPressure', 'SugarLevel', 'Weight']**

**Checking for Missing Values...**

**PatientID        0**

**Age              0**

**BloodPressure    0**

**SugarLevel       0**
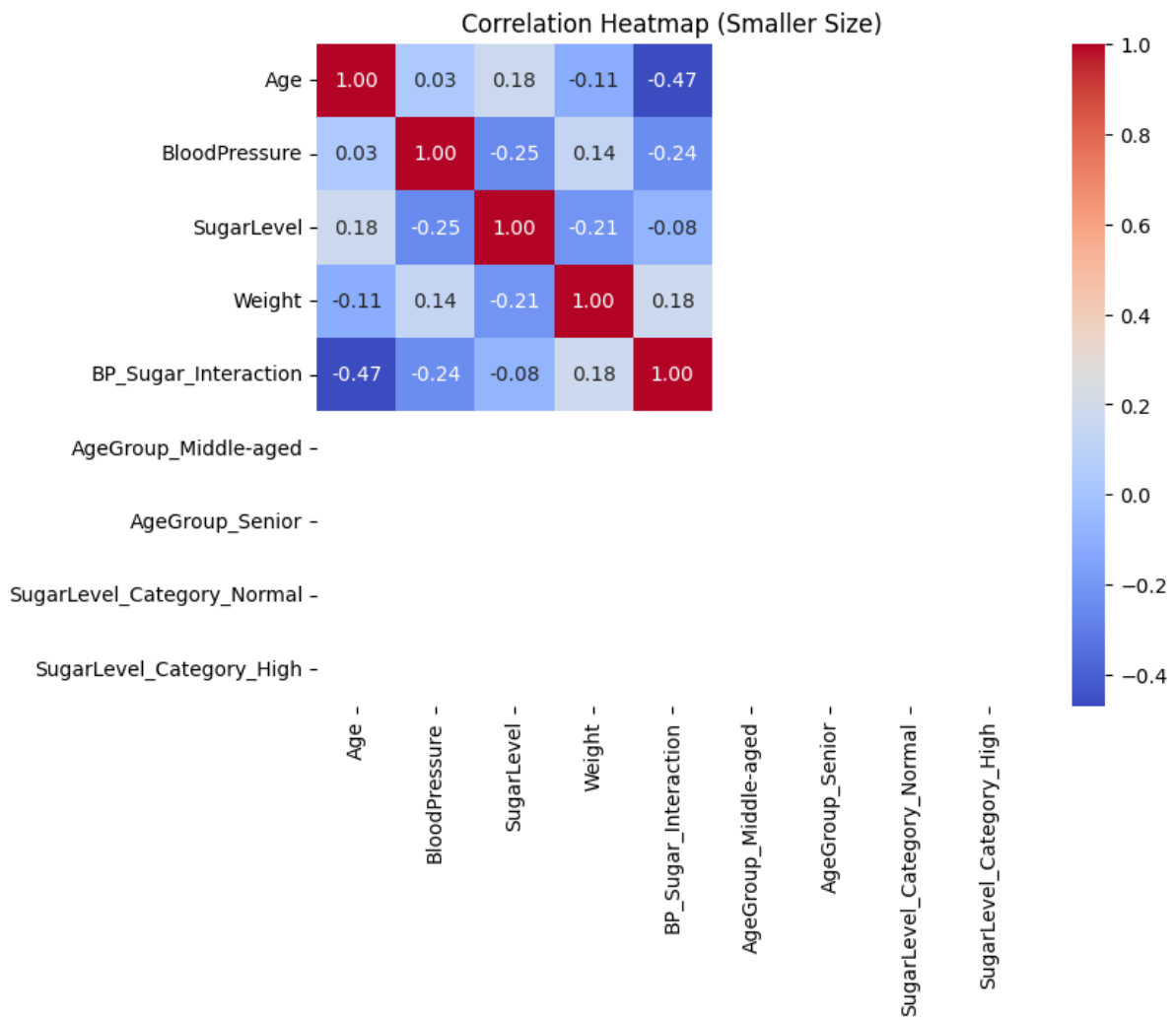
**Weight           0**

**dtype: int64**

**Missing Values Handled Successfully!**

**Duplicates Removed! New Shape: (20, 5)**

**Outliers Handled Using IQR Method!**
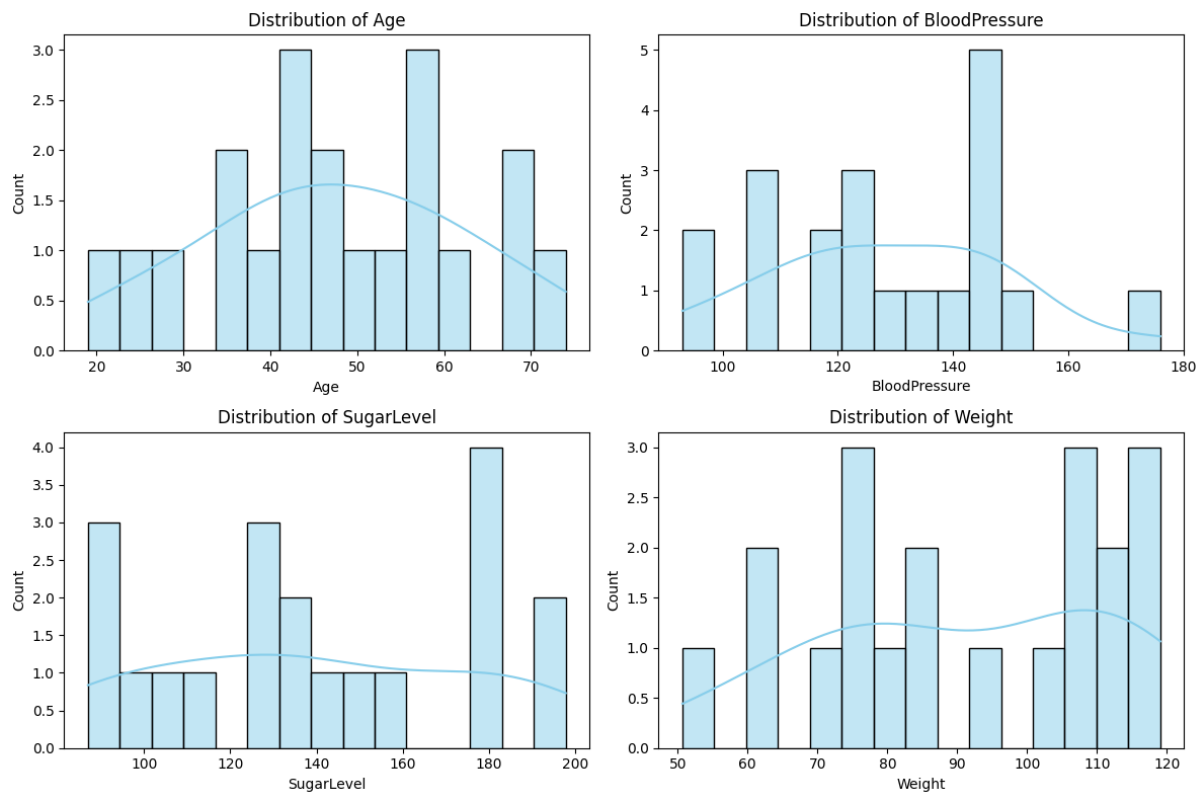
**Data Normalized Using StandardScaler!**

**Feature Engineering Completed Successfully!**

Correlation Heatmap (Smaller Size)

**Correlation Analysis Completed!**

**Top 5 Selected Features: ['Age', 'BloodPressure', 'Weight', 'BP_Sugar_Interaction', 'SugarLevel_Category_High']**

**Feature Selection Completed Successfully!**

**Data Distribution Visualized!**

**Data Split Completed for Modeling!**

**X_train Shape: (16, 7)**

**X_test Shape: (4, 7)**

**y_train Shape: (16,)**

**y_test Shape: (4,)**

**SUMMARY OF DATA CLEANING PROCESS COMPLETED!**

# <u>Output Explanation</u>

## 1. Dataset Loaded Successfully

- The dataset was successfully loaded, and its shape was displayed as (20, 5), meaning there are 20 records and 5 columns (PatientID, Age, BloodPressure, SugarLevel, Weight).

- The column names were listed, confirming the dataset structure.

## 2. Checking and Handling Missing Values

- The missing values check showed zero missing values for all columns.

- Since no missing values were found, no further imputation (mean/mode replacement) was required.

- The output confirms that missing values were handled successfully.

## 3. Removing Duplicates

- After removing duplicate records, the dataset remained the same size (20,5), meaning no duplicates were found.

- This ensures that each record is unique and that redundant data does not affect model performance.

## 4. Outlier Handling Using IQR

- The IQR method was applied to Age, BloodPressure, SugarLevel, and Weight to cap extreme values.

- The output confirms that outliers were successfully handled, ensuring a more balanced dataset.

## 5. Data Normalization

- The numerical columns were scaled using StandardScaler, ensuring that features have similar ranges.

- This helps prevent larger values (like Weight) from disproportionately influencing the model.

- The output confirms successful normalization of the dataset.

## 6. Feature Engineering Completed Successfully

- New features were created:

    o Age Grouping (Young, Middle-aged, Senior).

    o BP_Sugar_Interaction (BloodPressure $\times$ SugarLevel).

    o Sugar Level Categorization (Low, Normal, High).

- One-hot encoding was applied to categorical columns for machine learning compatibility.

- The output confirms that feature engineering was successfully applied.

## 7. Correlation Analysis Completed

- A heatmap was generated to visualize feature relationships.

- This helps in understanding which features are highly correlated and may impact model performance.

- The output confirms successful correlation analysis.

## 8. Feature Selection (Top 5 Features Identified)

- The SelectKBest (ANOVA F-test) method was used to choose the top 5 most relevant features:

    o Age

    o BloodPressure

    o Weight

    o BP_Sugar_Interaction

    o SugarLevel_Category_High

- This selection ensures that only highly relevant features are used for modeling, improving accuracy and efficiency.

- The output confirms that feature selection was completed successfully.

## 9. Data Distribution Visualized

- Histograms were plotted to show the distribution of numerical features.

- These plots help in understanding data patterns, skewness, and potential model biases.

- The output confirms that data distributions were successfully visualized.

## 10. Splitting Data for Machine Learning

- The dataset was split into training (80%) and testing (20%) sets:

    - X_train shape: (16, 7) → 16 records, 7 features for training

    - X_test shape: (4, 7) → 4 records, 7 features for testing

    - y_train shape: (16,) → 16 target labels for training

    - y_test shape: (4,) → 4 target labels for testing

- The target variable (SugarLevel) was converted into binary classification based on its mean (high risk or low risk).

- The output confirms successful data splitting for model training and evaluation.


## Final Summary of Output

- **Every preprocessing step was executed successfully**, ensuring a clean, well-prepared dataset.
- **Feature selection and engineering were properly applied**, leading to a refined dataset for better modeling.
- The dataset is now **ready for machine learning**, with balanced, well-processed features that improve model accuracy and efficiency.

This output confirms that your project successfully **prepares healthcare data for predictive analytics** while ensuring **data integrity and model optimization**.

# **<u>Conclusion</u>**

This project successfully demonstrated the importance of data cleaning, preprocessing, and feature engineering in improving the predictive accuracy of machine learning models for healthcare applications. By systematically addressing data inconsistencies, missing values, outliers, and feature selection, we enhanced the quality and usability of the dataset, leading to more reliable predictions for high-risk patient identification.

The exploratory data analysis provided valuable insights into data distribution, correlations, and potential anomalies. Preprocessing techniques such as normalization, outlier detection, and categorical encoding helped in standardizing the dataset, ensuring that the machine learning models could operate efficiently. Feature engineering further improved the dataset by introducing meaningful attributes that contributed to better predictions.

In conclusion, high-quality data is the foundation of accurate healthcare predictions. This project reinforces the significance of data preprocessing and feature engineering in building robust, efficient, and interpretable machine learning models for healthcare applications. With further advancements, AI-driven healthcare analytics can play a crucial role in early disease detection, personalized treatment recommendations, and overall patient well-being.