

KUMARAGURU COLLEGE OF TECHNOLOGY

DEPARTMENT OF ARTIFICIAL INTELGINCE AND DATA SCIENCE

24AD1204 – DATA SCIENCE AND VISULIZATION

NAME	ROLLNO
GOKULNAATH M	24BAD026
JAYARAKSHA REGURAJ	24BAD044
KAMALESH N	24BAD054
NISHANTH P	24BAD405

Submitted to
Mrs. E Shriarththy

Week No: 1**Team Formation & Data Hunting****Objectives:**

Selection of dataset. Setting up GitHub. Introduction to Pandas/NumPy.

Team Composition:

NAME	ROLES
GOKULNAATH M	DATA ENGINEER
JAYARAKSHA REGURAJ	STORYTELLER
KAMALESH N	DATA ANALYST
NISHANTH P	DATA VISUALIZATION

Work done:**1. Domain Selection and Problem Statement Abstract****Team Contribution:****1. Dataset Selection & Evaluation**

Gokulnaath M – 24BAD028

Dataset 1**Heart Attack Prediction Dataset**

<https://www.kaggle.com/datasets/iamsouravbanerjee/heart-attack-prediction-dataset>

Characteristics

- Low / limited attributes
- Low object (record) count
- Semi-structured data

Pros

- Uncleaned dataset – ideal for demonstrating data preprocessing, handling missing values, and noise reduction
- Simple structure, easy to understand feature relationships
- Good for baseline model development

Cons

- Limited number of attributes restricts deep risk factor analysis
- Small dataset size reduces model generalization
- Not suitable alone for dimensionality reduction showcase

Jayaraksha Reguraj – 24BAD044**Dataset 2****Heart Disease Dataset**

<https://www.kaggle.com/datasets/johnsmith88/heart-disease-dataset>

Characteristics

- Already cleaned and pre-processed
- Low / limited object count

Pros

- Clean and ready for quick model training
- Good for algorithm comparison
- Minimal preprocessing effort required

Cons

- Already cleaned – not ideal for project requirements that demand full preprocessing pipeline
- Less opportunity to demonstrate data cleaning skills
- Limited data volume

Kamalesh N – 24BAD054**Dataset 3 (Finalized Dataset)****Heart Disease Prediction using Logistic Regression**

<https://www.kaggle.com/datasets/dileep070/heart-disease-prediction-using-logistic-regression>

Characteristics

- Good number of attributes
- Dataset is not pre-processed

Pros

- Uncleaned dataset – perfectly aligned with project requirement to perform:
 - Data cleaning
 - Feature scaling
 - Dimensionality reduction (PCA)
- Balanced attribute set supports risk score generation
- Suitable for end-to-end ML pipeline

Cons

- Requires significant preprocessing effort
- Noise and missing values may affect initial model performance
- Needs careful feature engineering

Nishanth P – 24BAD405

Dataset 4

Framingham Heart Study Dataset

<https://www.kaggle.com/datasets/amanajmera1/framingham-heart-study-dataset>

Characteristics

- Long-term cardiovascular study data
- Moderate to high attribute count
- Includes demographic, behavioural, and medical risk factors

Pros

- Raw / uncleaned dataset – excellent for showcasing:
 - Missing value handling
 - Outlier detection
 - Feature transformation
- Rich clinical depth enables strong risk factor interpretation
- Ideal for dimensionality reduction & risk score modeling

Cons

- Requires extensive preprocessing
- Some attributes may need domain understanding
- Higher complexity compared to smaller datasets

Summary Comparison Table

Dataset	Pre-processed	Attribute Count	Object Count	Suitability
Dataset 1	No	Low	Low	Basic preprocessing demo
Dataset 2	Yes	Low	Low	Model testing only
Dataset 3 (Final)	No	High	Medium	Best fit for project
Dataset 4	No	High	Medium	Advanced analysis

2. Setting up GitHub

All *

Procedure

- Create GitHub Repositories with name “[24ADI204 DSV Team14](#)”
- Add Description of "Heartbeats & Habits" combats the silent threat of cardiovascular disease by applying rigorous cleaning and dimensionality reduction to noisy clinical data. This approach transforms complex biomarker and lifestyle interactions into a unified "Risk Score," visualizing a clear path to improved cardiac health
- Add README.md

3. Introduction to Pandas/NumPy

All *

Procedure

- Research and learn the basic concept of NumPy
- Research and learn the basic concept of Pandas

Resources and References

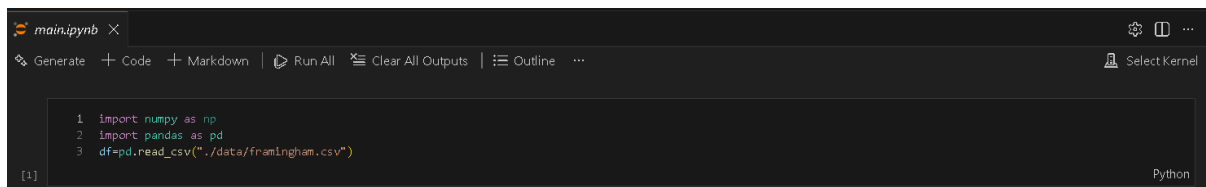
- **Dataset**
 - <https://www.kaggle.com/datasets/iamsouravbanerjee/heart-attack-prediction-dataset>
 - <https://www.kaggle.com/datasets/johnsmith88/heart-disease-dataset>
 - <https://www.kaggle.com/datasets/dileep070/heart-disease-prediction-using-logistic-regression>
 - <https://www.kaggle.com/datasets/amanajmera1/framingham-heart-study-dataset>
- **GitHub**
 - https://github.com/Gokulnaath-gif/24ADI204_DSV_Team14
- **Learns**
 - <https://numpy.org/doc/stable/user/index.html#user>
 - https://www.w3schools.com/python/numpy/numpy_intro.asp
 - https://pandas.pydata.org/docs/user_guide/index.html
 - <https://www.w3schools.com/python/pandas/default.asp>

Week No: 2**Know Your Data****Objectives:**

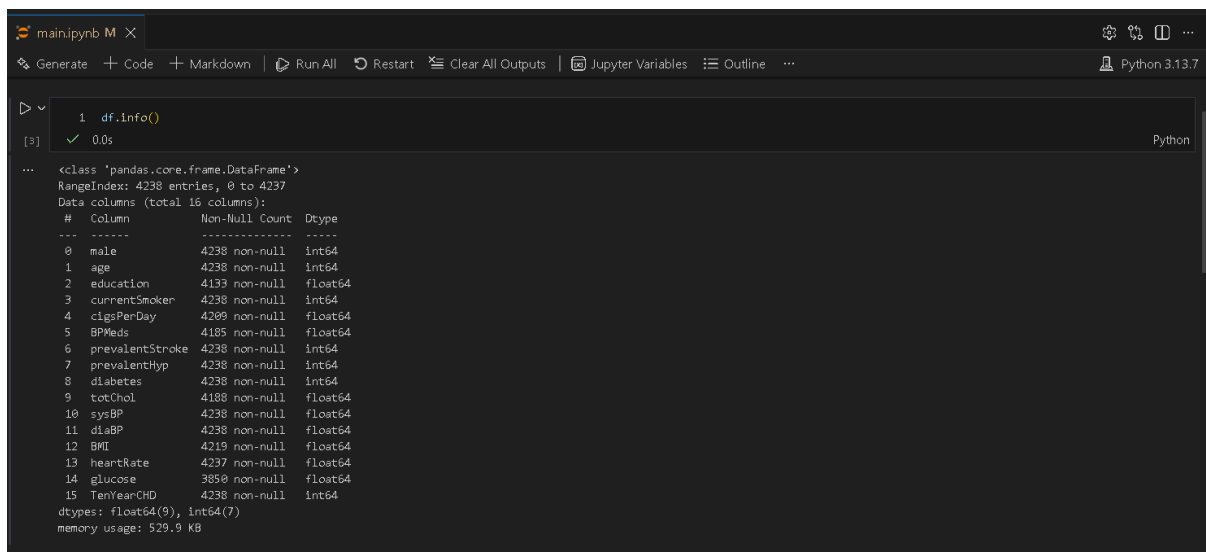
Loading data, checking types, inspecting the structure.

Work done:**1. Load Data and Analyses Dataset**

Gokulnaath M – 24BAD028, Nishanth P – 24BAD405

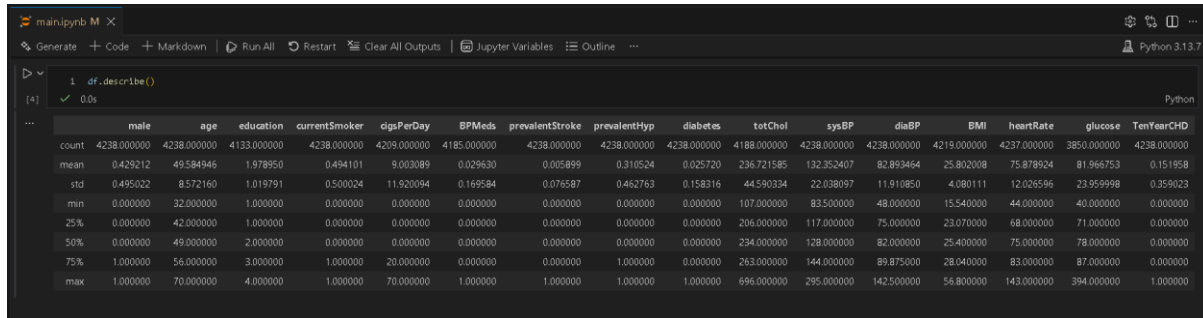
Procedure**1. Load the Dataset**

```
main.ipynb X
Generate + Code + Markdown Run All Clear All Outputs Outline ...
1 import numpy as np
2 import pandas as pd
3 df=pd.read_csv("../data/framlingham.csv")
[1] Python
```

2. Data Structure

```
main.ipynb M X
Generate + Code + Markdown Run All Restart Clear All Outputs Jupyter Variables Outline ...
1 df.info()
[3] ✓ 0.0s Python
...
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4238 entries, 0 to 4237
Data columns (total 16 columns):
#   Column                Non-Null Count  Dtype
---  -
0   male                  4238 non-null  int64
1   age                   4238 non-null  int64
2   education             4133 non-null  float64
3   currentSmoker         4238 non-null  int64
4   cigsPerDay            4209 non-null  float64
5   BPmeds                4185 non-null  float64
6   prevalentStroke       4238 non-null  int64
7   prevalentHyp          4238 non-null  int64
8   diabetes              4238 non-null  int64
9   totChol               4188 non-null  float64
10  sysBP                4238 non-null  float64
11  diaBP                4238 non-null  float64
12  BMI                  4219 non-null  float64
13  heartRate            4237 non-null  float64
14  glucose              3850 non-null  float64
15  TenYearCHD           4238 non-null  int64
dtypes: float64(0), int64(7)
memory usage: 529.9 KB
```

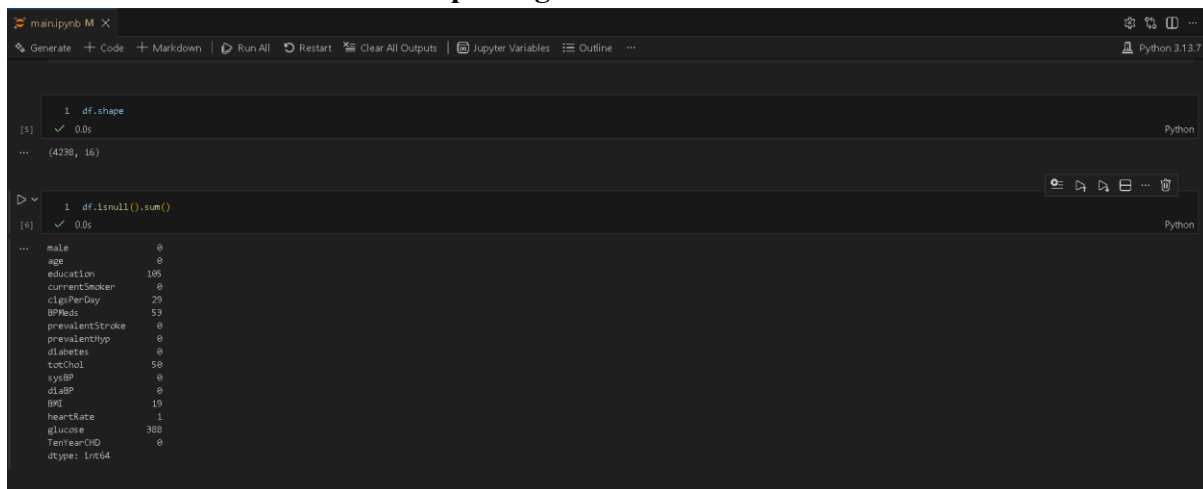
3. Data Statistical Information



A screenshot of a Jupyter Notebook interface. The top bar shows 'mainipynb M X' and various icons. Below the toolbar, the code cell contains 'df.describe()'. The output is a statistical summary table with 16 columns: count, mean, std, min, 25%, 50%, 75%, and max for each feature. The features are: male, age, education, currentSmoker, cigsPerDay, BPMeds, prevalentStroke, prevalentHyp, diabetes, totChol, sysBP, diaBP, BMI, heartRate, glucose, and TenYearCHD.

	male	age	education	currentSmoker	cigsPerDay	BPMeds	prevalentStroke	prevalentHyp	diabetes	totChol	sysBP	diaBP	BMI	heartRate	glucose	TenYearCHD
count	4238.000000	4238.000000	4133.000000	4238.000000	4209.000000	4185.000000	4238.000000	4238.000000	4238.000000	4188.000000	4238.000000	4238.000000	4219.000000	4237.000000	3850.000000	4238.000000
mean	0.429212	49.584946	1.978950	0.494101	9.003089	0.029630	0.005899	0.310524	0.025720	236.721585	132.352407	82.893464	25.802008	75.878924	81.966753	0.151958
std	0.495022	8.572160	1.019791	0.500024	11.920094	0.169584	0.076587	0.462763	0.158316	44.580334	22.038097	11.910850	4.080111	12.026596	23.959998	0.359023
min	0.000000	32.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	107.000000	83.500000	48.000000	15.540000	44.000000	40.000000	0.000000
25%	0.000000	42.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	206.000000	117.000000	75.000000	23.070000	68.000000	71.000000	0.000000
50%	0.000000	49.000000	2.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	234.000000	128.000000	82.000000	25.400000	75.000000	78.000000	0.000000
75%	1.000000	56.000000	3.000000	1.000000	20.000000	0.000000	0.000000	1.000000	0.000000	263.000000	144.000000	89.875000	28.040000	83.000000	87.000000	0.000000
max	1.000000	70.000000	4.000000	1.000000	70.000000	1.000000	1.000000	1.000000	1.000000	696.000000	295.000000	142.500000	56.800000	143.000000	394.000000	1.000000

4. Data Structure and Inspecting Null Values



A screenshot of a Jupyter Notebook interface. The first code cell contains 'df.shape' and the output is '(4238, 16)'. The second code cell contains 'df.isnull().sum()' and the output is a list of features with their respective null counts: male (0), age (0), education (105), currentSmoker (0), cigsPerDay (29), BPMeds (53), prevalentStroke (0), prevalentHyp (0), diabetes (0), totChol (58), sysBP (0), diaBP (0), BMI (19), heartRate (1), glucose (368), and TenYearCHD (0). The dtype is listed as 'int64'.

Feature	Count
male	0
age	0
education	105
currentSmoker	0
cigsPerDay	29
BPMeds	53
prevalentStroke	0
prevalentHyp	0
diabetes	0
totChol	58
sysBP	0
diaBP	0
BMI	19
heartRate	1
glucose	368
TenYearCHD	0

2. “First Impression” report.

Jayaraksha Reguraj – 24BAD044, Kamalesh N – 24BAD054

Report

1. Dataset Overview

The dataset used for this analysis is the Framingham Heart Study dataset, which contains medical and lifestyle information of patients. The primary objective of the dataset is to predict the 10-year risk of coronary heart disease (CHD).

- Number of records: 4,240
- Number of features: 16
- Target variable: TenYearCHD (binary: 0 = No CHD, 1 = CHD)

2. Data Structure

- The dataset consists of numerical features (both integer and float types).
- Variables include:
 - Demographic features: age, sex
 - Behavioural features: smoking status, cigarettes per day
 - Medical history: blood pressure medication, diabetes, stroke history

Clinical measurements: cholesterol, systolic & diastolic BP, BMI, heart rate, glucose

3. Data Types

- Most features are of type int64 or float64.
- The target variable TenYearCHD is an integer categorical variable representing a binary outcome.
- No categorical string variables are present, reducing the need for encoding at this stage.

4. Missing Values

- Several columns contain missing values, including:
 - education
 - cigsPerDay
 - BPMeds
 - totChol
 - BMI
 - heartRate
 - glucose

5. Initial Observations

- The dataset appears realistic and noisy, which is expected in medical data.
- Feature scales vary significantly (e.g., age vs. cholesterol vs. glucose), suggesting the need for feature scaling before modeling.
- The target variable is imbalanced, with fewer positive CHD cases compared to negative ones, which may affect model performance.

6. Data Quality Assessment

- No duplicate rows were observed.
- No obvious data corruption or invalid values were detected during initial inspection.
- Presence of missing values is the main data quality concern.

7. Overall, First Impression

The dataset is:

- Well-structured
- Relevant for classification tasks
- Contains missing values
- Requires preprocessing and scaling
- Potential class imbalance in the target variable

Resources and References

- **Dataset**
 - <https://www.kaggle.com/datasets/dileep070/heart-disease-prediction-using-logistic-regression>
- **GitHub**
 - https://github.com/Gokulnaath-gif/24ADI204_DSV_Team14

Week No: 3

The Cleaning Sprint

Objectives:

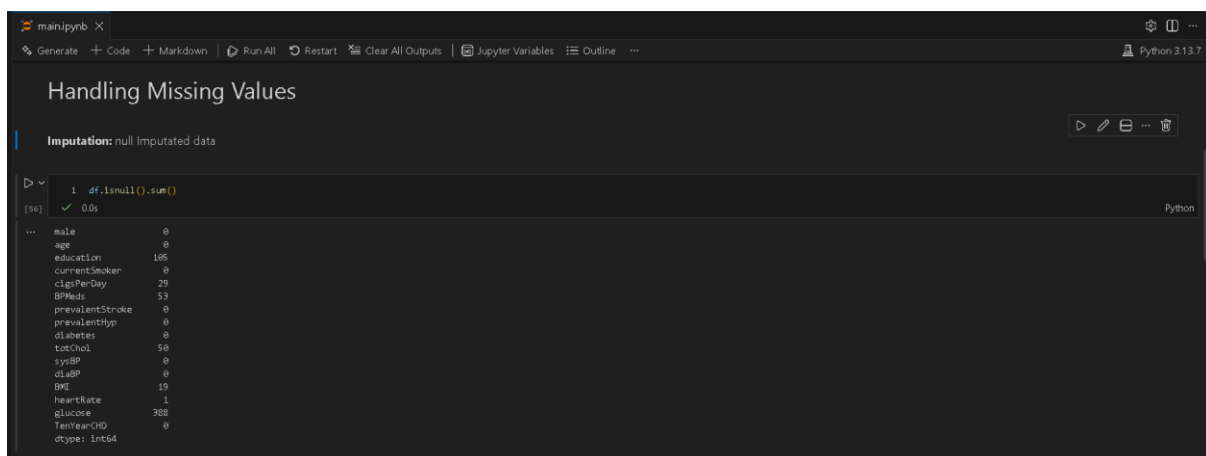
Handling Missing Values (Imputation strategies) and Outlier Detection (Boxplots/Z-score).

Work done:

1. Handling Null Values

Gokulnaath M – 24BAD028, Nishanth P – 24BAD405

1. Attributes Contains Null Values



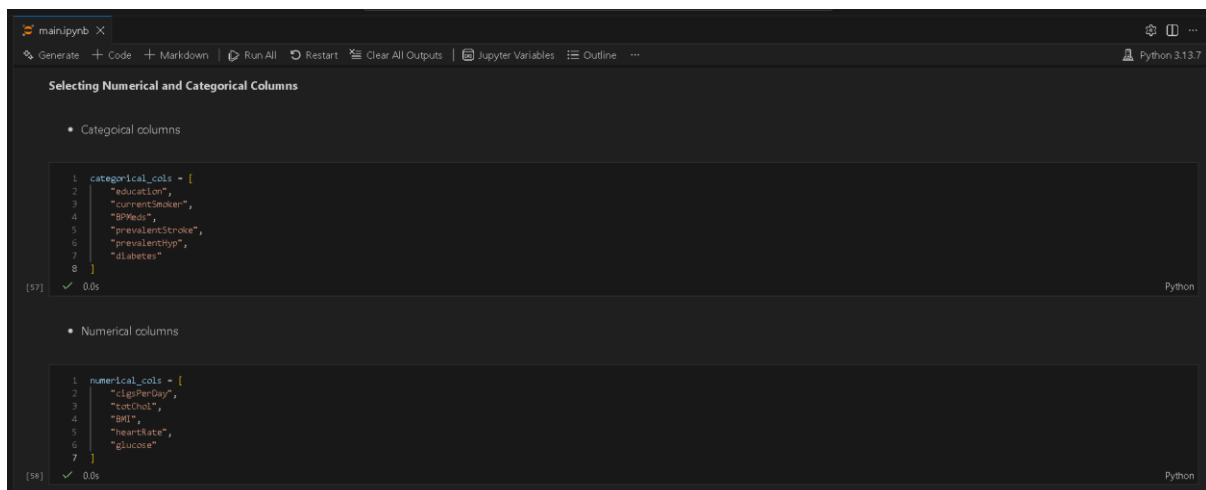
The screenshot shows a Jupyter Notebook titled "Handling Missing Values". The code cell contains the following Python code:

```
1 df.isnull().sum()
```

The output shows the sum of null values for each column in the dataset:

```
male      0
age       0
education 185
currentSmoker  0
cigsPerDay 29
BPmeds    53
prevalentStroke  0
prevalentHyp  0
diabetes   0
totChol   50
sysBP     0
diabP     0
BMI       19
heartRate  1
glucose    388
TenYearCHD  0
dtype: int64
```

2. Selecting Categorical and Numerical Attributes



The screenshot shows a Jupyter Notebook titled "Selecting Numerical and Categorical Columns". The code cell contains the following Python code:

```
1 categorical_cols = [
2     "education",
3     "currentSmoker",
4     "BPmeds",
5     "prevalentStroke",
6     "prevalentHyp",
7     "diabetes"
8 ]
```

The output shows the list of categorical columns:

```
1 categorical_cols = [
2     "education",
3     "currentSmoker",
4     "BPmeds",
5     "prevalentStroke",
6     "prevalentHyp",
7     "diabetes"
8 ]
```

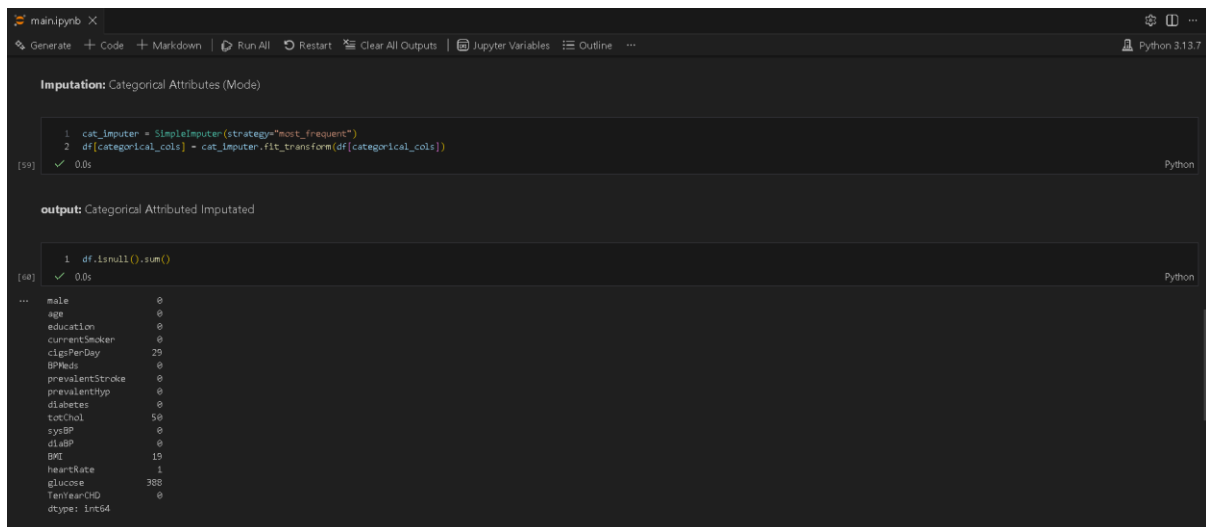
The code cell also contains the following Python code:

```
1 numerical_cols = [
2     "cigsPerDay",
3     "totChol",
4     "BMI",
5     "heartRate",
6     "glucose"
7 ]
```

The output shows the list of numerical columns:

```
1 numerical_cols = [
2     "cigsPerDay",
3     "totChol",
4     "BMI",
5     "heartRate",
6     "glucose"
7 ]
```

3. Imputation of Categorical Attributes (Mode)



The screenshot shows a Jupyter Notebook interface with a dark theme. The top bar includes a file explorer, a toolbar with icons for Generate, Code, Markdown, Run All, Restart, Clear All Outputs, Jupyter Variables, and Outline, and a Python version indicator (Python 3.13.7). The notebook has two cells. The first cell, titled 'Imputation: Categorical Attributes (Mode)', contains two lines of Python code: `1 cat_imputer = SimpleImputer(strategy="most_frequent")` and `2 df[categorical_cols] = cat_imputer.fit_transform(df[categorical_cols])`. The second cell, titled 'output: Categorical Attributed Imputed', contains a line of code: `1 df.isnull().sum()`. The output of the second cell is a list of 17 categorical variables and their corresponding null counts, all of which are 0. The variables are: male, age, education, currentSmoker, cigsPerDay, BPmeds, prevalentStroke, prevalentHyp, diabetes, totChol, sysBP, diaBP, BGL, heartRate, glucose, TenYearCHD, and dtype: int64.

```
mainipynb X
Generate + Code + Markdown Run All Restart Clear All Outputs Jupyter Variables Outline Python 3.13.7

Imputation: Categorical Attributes (Mode)

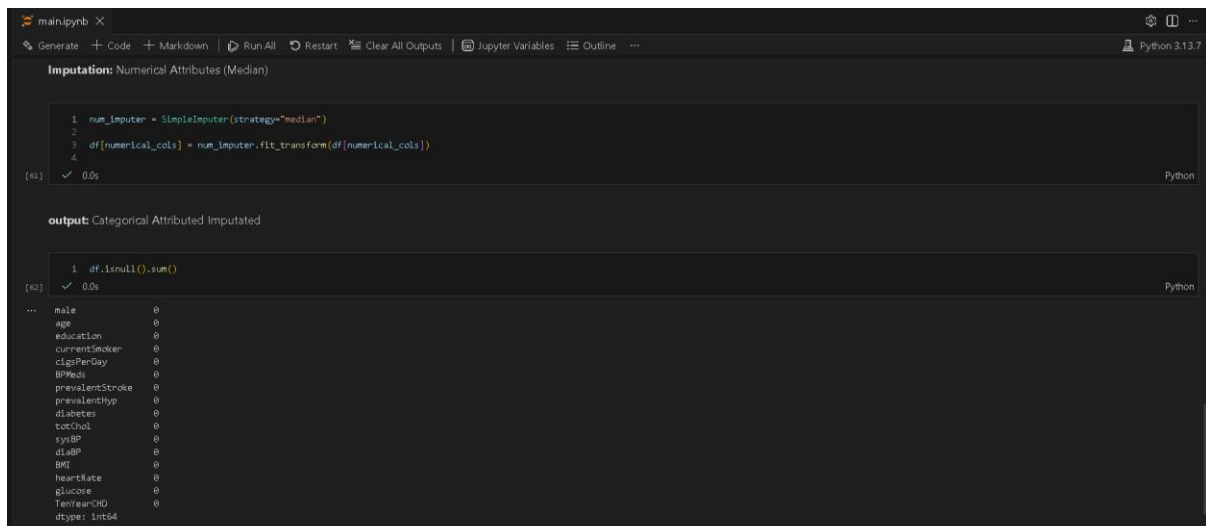
1 cat_imputer = SimpleImputer(strategy="most_frequent")
2 df[categorical_cols] = cat_imputer.fit_transform(df[categorical_cols])
[59] ✓ 0.0s Python

output: Categorical Attributed Imputed

1 df.isnull().sum()
[60] ✓ 0.0s Python

...
male 0
age 0
education 0
currentSmoker 0
cigsPerDay 29
BPmeds 0
prevalentStroke 0
prevalentHyp 0
diabetes 0
totChol 50
sysBP 0
diaBP 0
BGL 15
heartRate 1
glucose 388
TenYearCHD 0
dtype: int64
```

4. Imputation of Numerical Attributes (Median)



The screenshot shows a Jupyter Notebook interface with a dark theme. The top bar includes a file explorer, a toolbar with icons for Generate, Code, Markdown, Run All, Restart, Clear All Outputs, Jupyter Variables, and Outline, and a Python version indicator (Python 3.13.7). The notebook has two cells. The first cell, titled 'Imputation: Numerical Attributes (Median)', contains three lines of Python code: `1 num_imputer = SimpleImputer(strategy="median")`, `2` (blank line), and `3 df[numerical_cols] = num_imputer.fit_transform(df[numerical_cols])`. The second cell, titled 'output: Categorical Attributed Imputed', contains a line of code: `1 df.isnull().sum()`. The output of the second cell is a list of 17 numerical variables and their corresponding null counts, all of which are 0. The variables are: male, age, education, currentSmoker, cigsPerDay, BPmeds, prevalentStroke, prevalentHyp, diabetes, totChol, sysBP, diaBP, BGL, heartRate, glucose, TenYearCHD, and dtype: int64.

```
mainipynb X
Generate + Code + Markdown Run All Restart Clear All Outputs Jupyter Variables Outline Python 3.13.7

Imputation: Numerical Attributes (Median)

1 num_imputer = SimpleImputer(strategy="median")
2
3 df[numerical_cols] = num_imputer.fit_transform(df[numerical_cols])
[61] ✓ 0.0s Python

output: Categorical Attributed Imputed

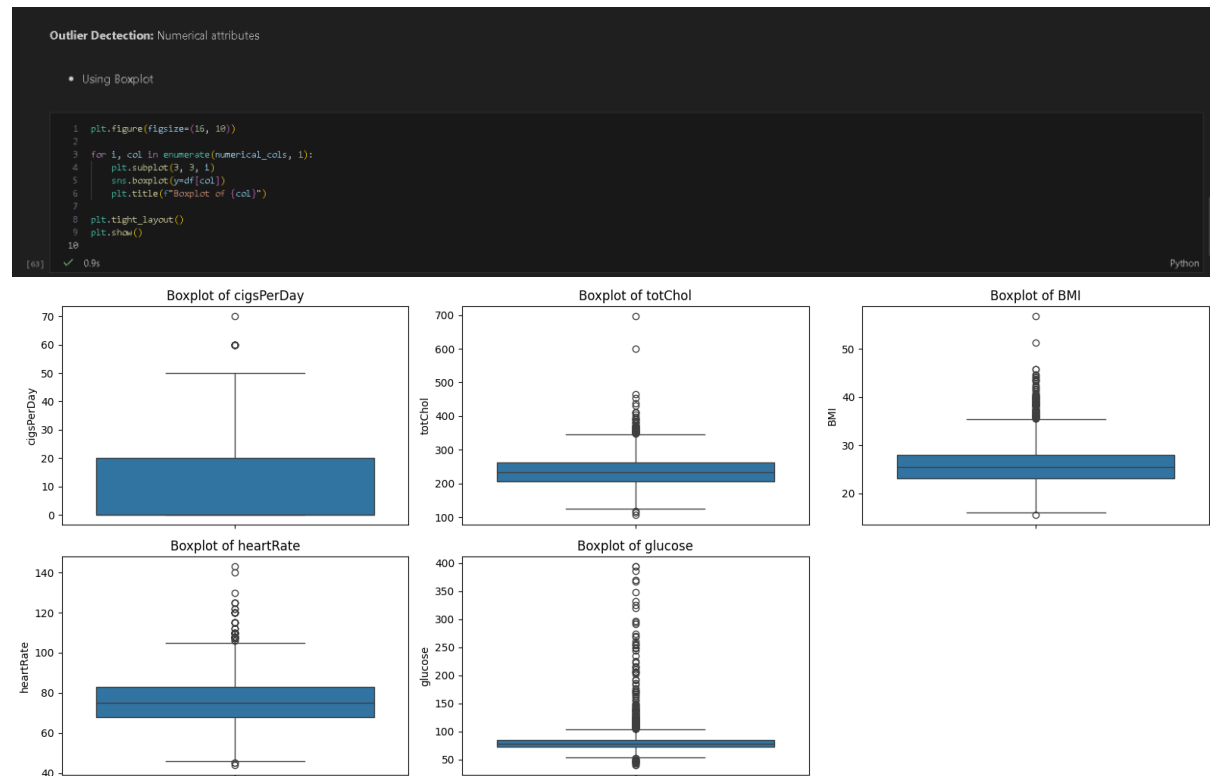
1 df.isnull().sum()
[62] ✓ 0.0s Python

...
male 0
age 0
education 0
currentSmoker 0
cigsPerDay 0
BPmeds 0
prevalentStroke 0
prevalentHyp 0
diabetes 0
totChol 0
sysBP 0
diaBP 0
BGL 0
heartRate 0
glucose 0
TenYearCHD 0
dtype: int64
```

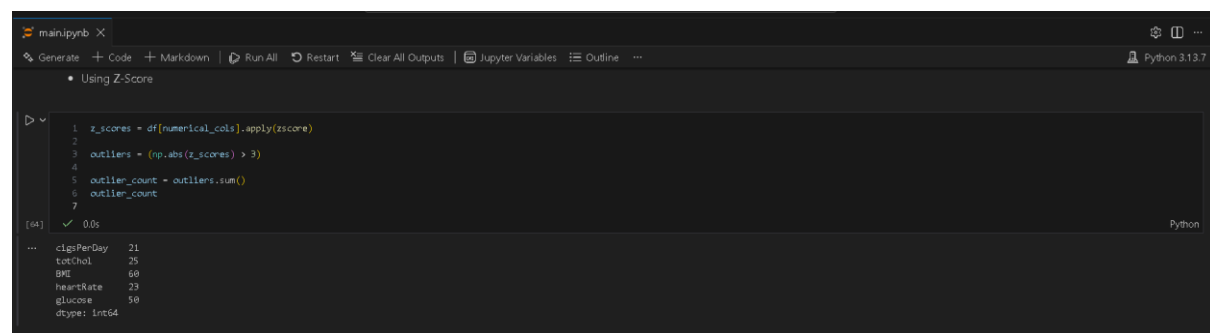
2. Handling Outliers values

Jayaraksha Reguraj – 24BAD044, Kamalesh N – 24BAD054

1. Outlier Detections using Boxplot

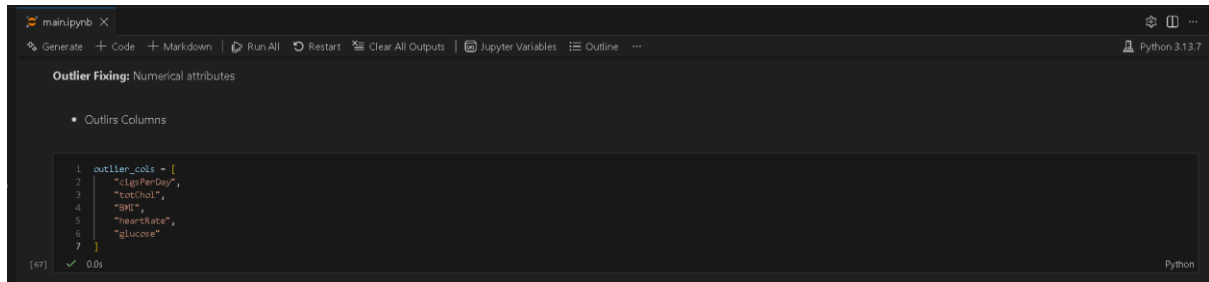


2. Outlier Detections using Z-Score



3. Outlier Fixing

- Select Outlier Columns

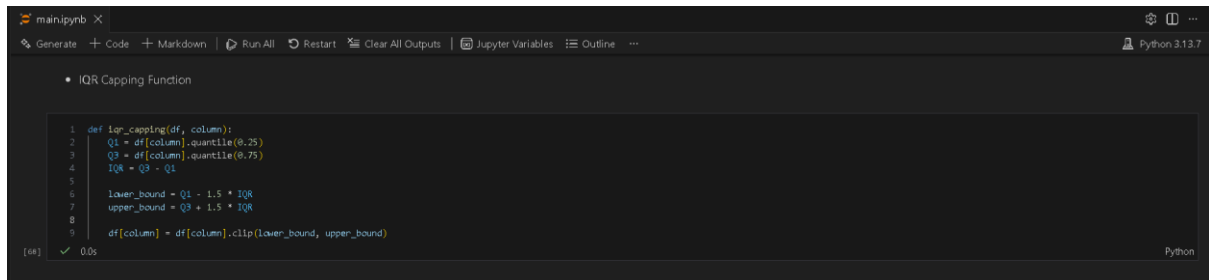


A Jupyter Notebook interface with a dark theme. The title bar shows 'mainipynb' and various icons. The menu bar includes 'Generate', '+ Code', '+ Markdown', 'Run All', 'Restart', 'Clear All Outputs', 'Jupyter Variables', 'Outline', and a settings icon. The notebook content is titled 'Outlier Fixing: Numerical attributes' and contains a bullet point 'Outliers Columns'. Below this, a code cell [67] contains the following Python code:

```
1 outlier_cols = [  
2     "cigsPerDay",  
3     "totChol",  
4     "BMI",  
5     "heartRate",  
6     "glucose"  
7 ]
```

The code cell is executed, showing a green checkmark and '0.0s'.

- IQR Capping Function

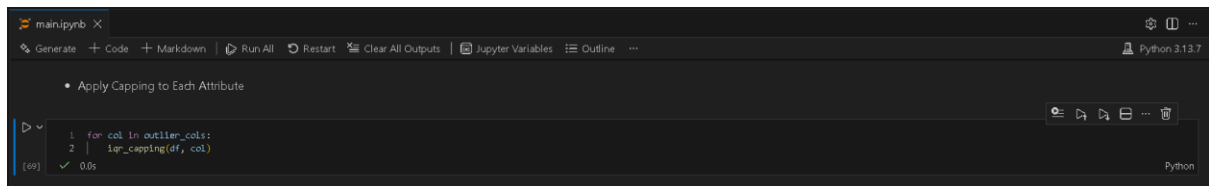


A Jupyter Notebook interface with a dark theme. The title bar shows 'mainipynb' and various icons. The menu bar includes 'Generate', '+ Code', '+ Markdown', 'Run All', 'Restart', 'Clear All Outputs', 'Jupyter Variables', 'Outline', and a settings icon. The notebook content is titled 'IQR Capping Function' and contains a code cell [68] with the following Python code:

```
1 def iqr_capping(df, column):  
2     Q1 = df[column].quantile(0.25)  
3     Q3 = df[column].quantile(0.75)  
4     IQR = Q3 - Q1  
5  
6     lower_bound = Q1 - 1.5 * IQR  
7     upper_bound = Q3 + 1.5 * IQR  
8  
9     df[column] = df[column].clip(lower_bound, upper_bound)
```

The code cell is executed, showing a green checkmark and '0.0s'.

- Apply Capping to Each Attribute



A Jupyter Notebook interface with a dark theme. The title bar shows 'mainipynb' and various icons. The menu bar includes 'Generate', '+ Code', '+ Markdown', 'Run All', 'Restart', 'Clear All Outputs', 'Jupyter Variables', 'Outline', and a settings icon. The notebook content is titled 'Apply Capping to Each Attribute' and contains a code cell [69] with the following Python code:

```
1 for col in outlier_cols:  
2     iqr_capping(df, col)
```

The code cell is executed, showing a green checkmark and '0.0s'.

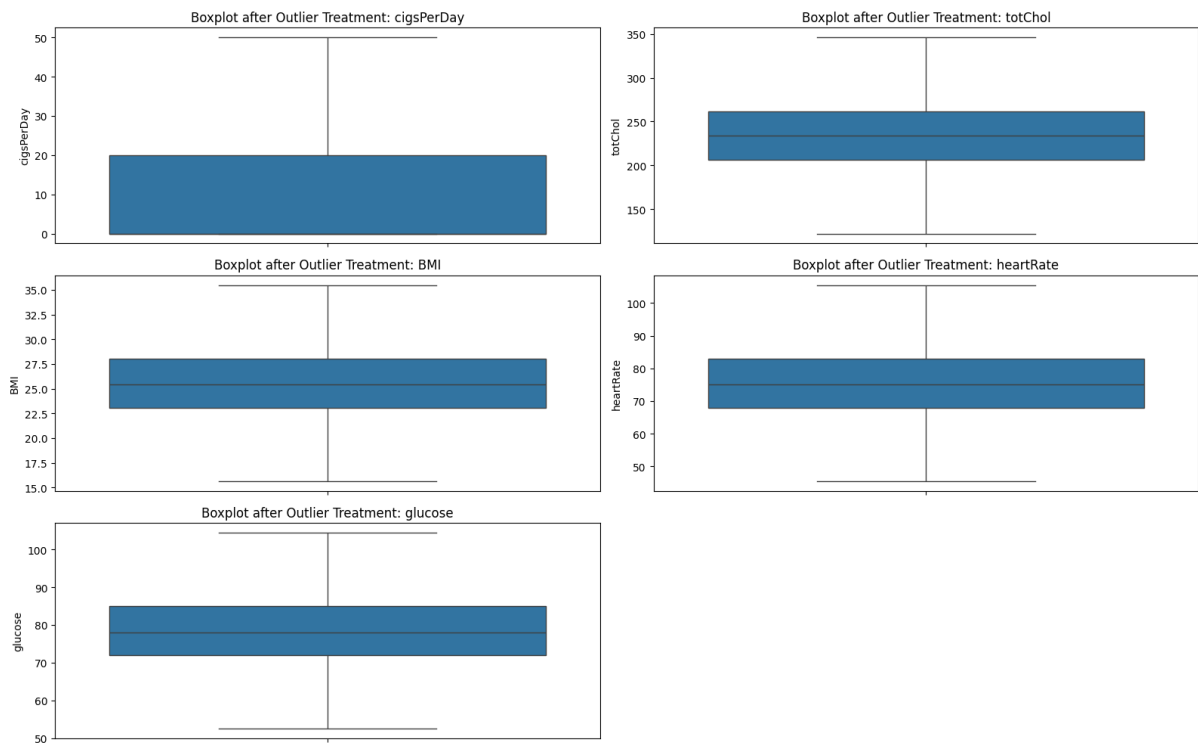
- Verify Outliers After Treatment

```
mainipynb X
Generate + Code + Markdown Run All Restart Clear All Outputs Jupyter Variables Outline Python 3.13.7

• Verify Outliers After Treatment

1 plt.figure(figsize=(16, 10))
2
3 for i, col in enumerate(outlier_cols, 1):
4     plt.subplot(3, 2, i)
5     sns.boxplot(y=df[col])
6     plt.title("Boxplot after Outlier Treatment: {col}")
7
8 plt.tight_layout()
9 plt.show()

[78] ✓ 0.8s
```



Resources and References

- **Dataset**

- <https://www.kaggle.com/datasets/dileep070/heart-disease-prediction-using-logistic-regression>

- **GitHub**

- https://github.com/Gokulnaath-gif/24ADI204_DSV_Team14

- **Resources**

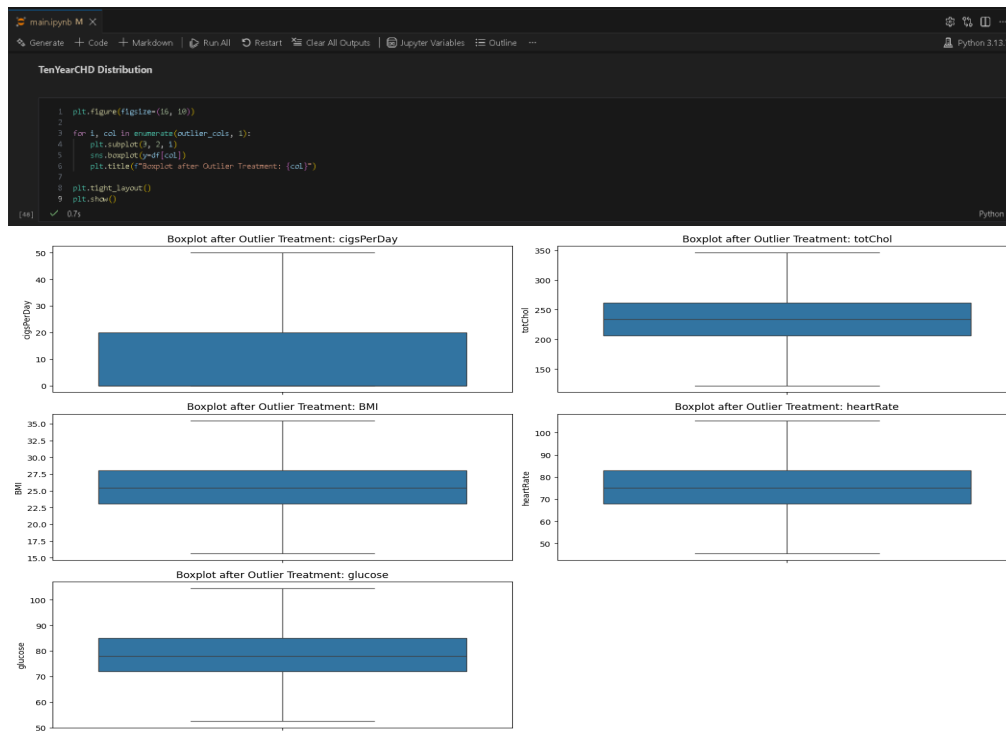
- <https://www.geeksforgeeks.org/data-analysis/working-with-missing-data-in-pandas/>
- https://medium.com/@punya8147_26846/dataframes-handling-missing-values-in-pandas-11f7702afaf7
- <https://www.geeksforgeeks.org/data-science/detect-and-remove-the-outliers-using-python/>
- <https://www.geeksforgeeks.org/pandas/handling-outliers-with-pandas/>
- <https://llego.dev/posts/outlier-detection-handling-python-guide/>
- <https://x.ai/grok>
- <https://openai.com/chatgpt>

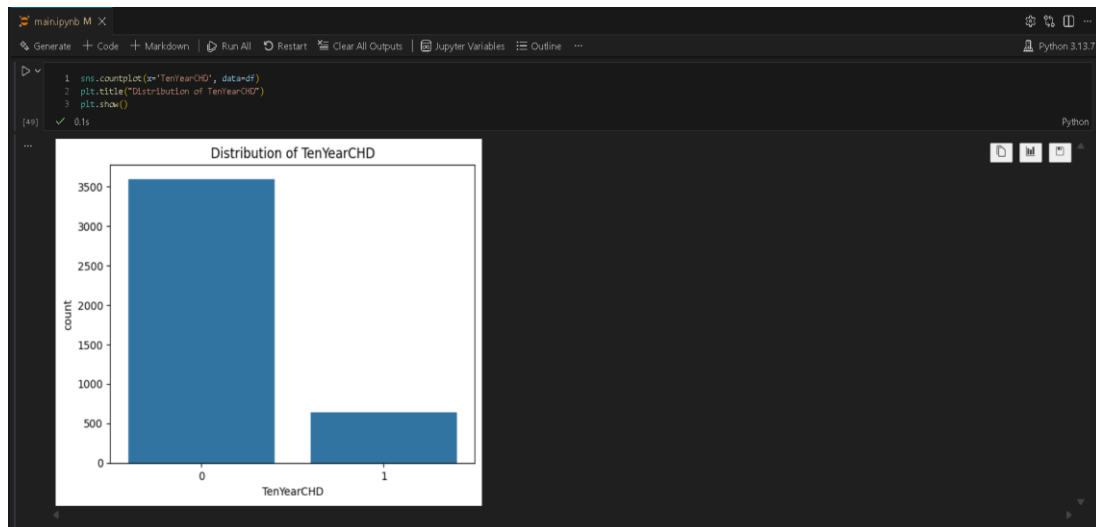
Week No: 4

EDA Deep Dive

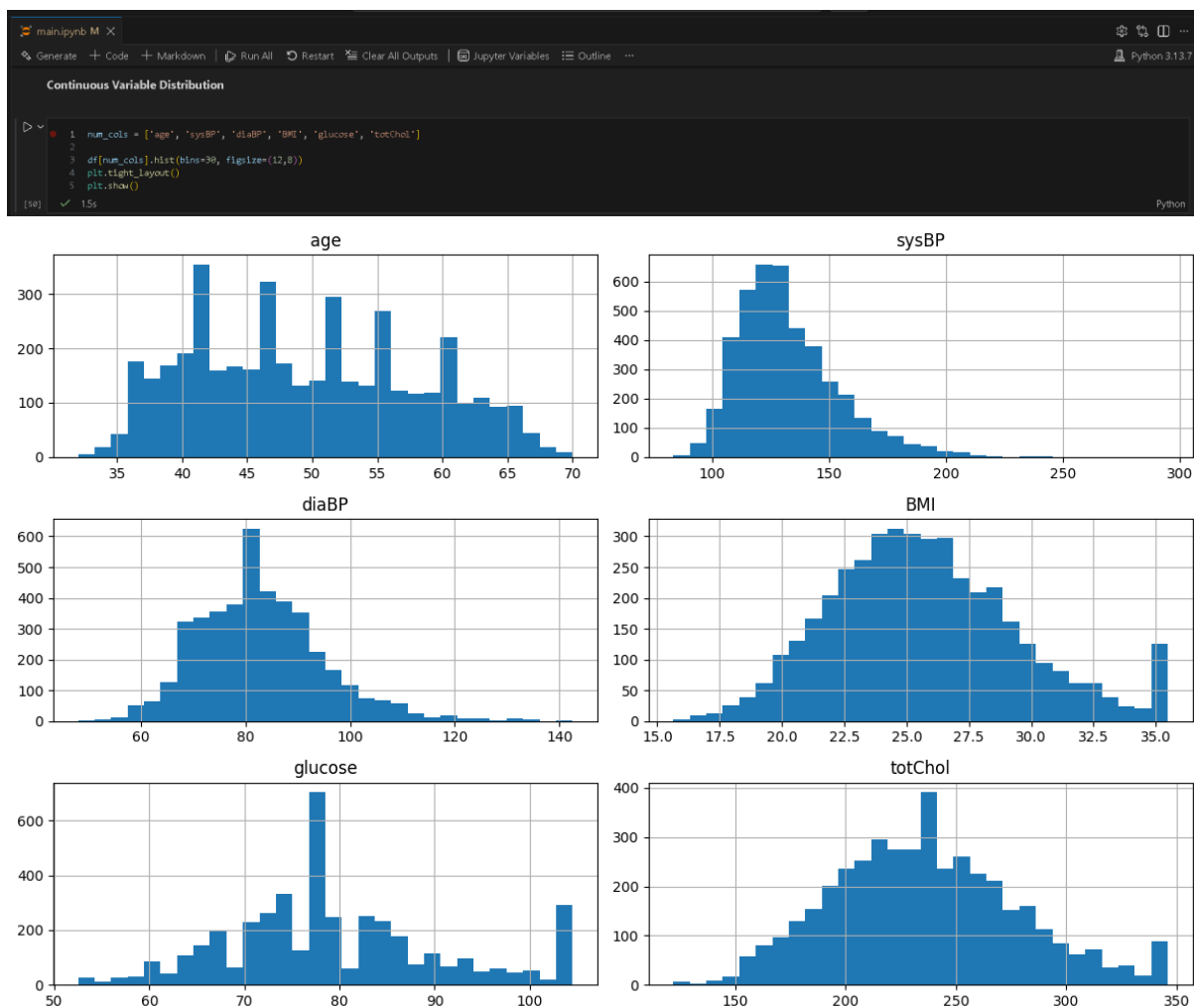
Objectives:

Variant analysis and Visualizing distributions.

Work done:**1. Variant Analysis****Kamalesh N 24BAD054 - Nishanth P 24BAD405****1. Univariate****1. TenYearCHD Distribution**



2. Continuous Variable Distribution



Why Histogram?

- Shows distribution shape (Normal / Skewed)
- Helps identify skewness

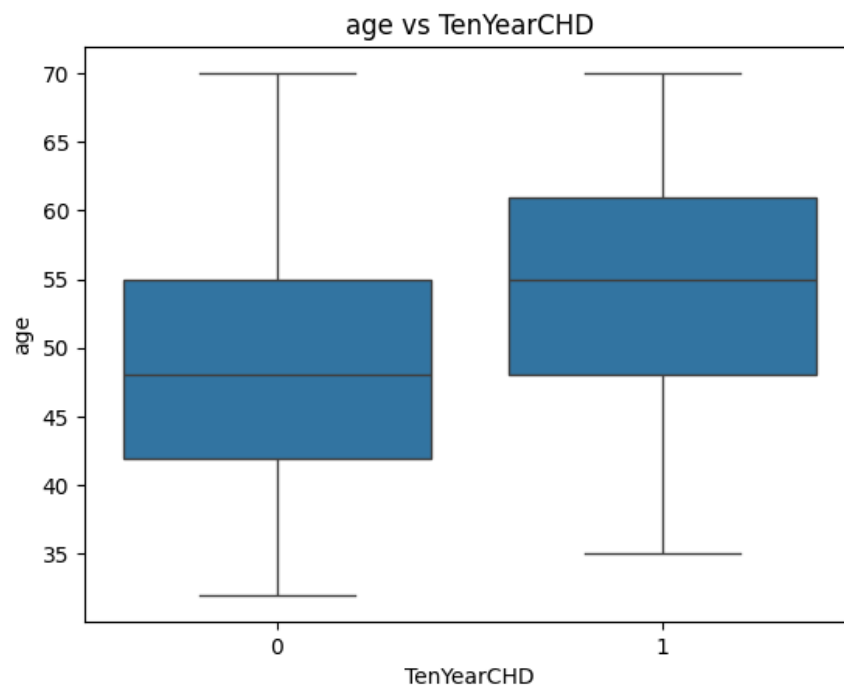
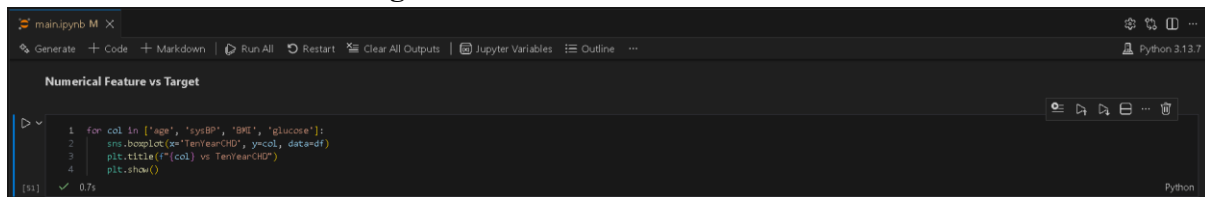
- Helps check extreme values

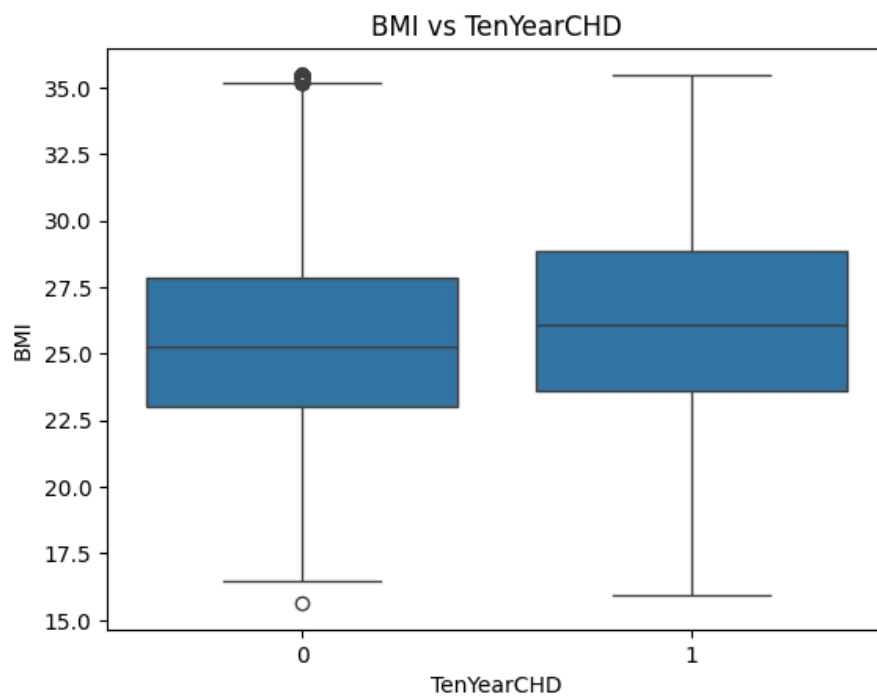
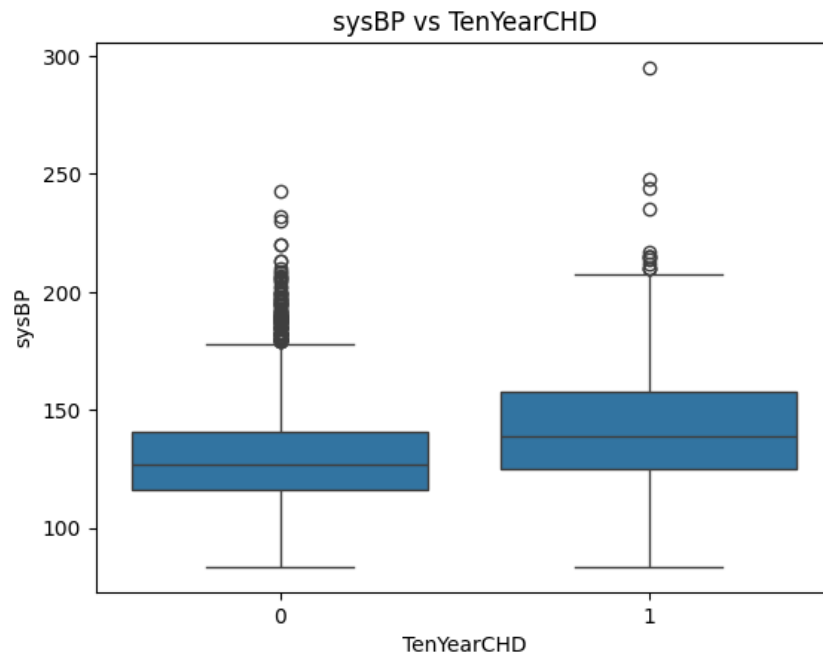
Univariate Observations

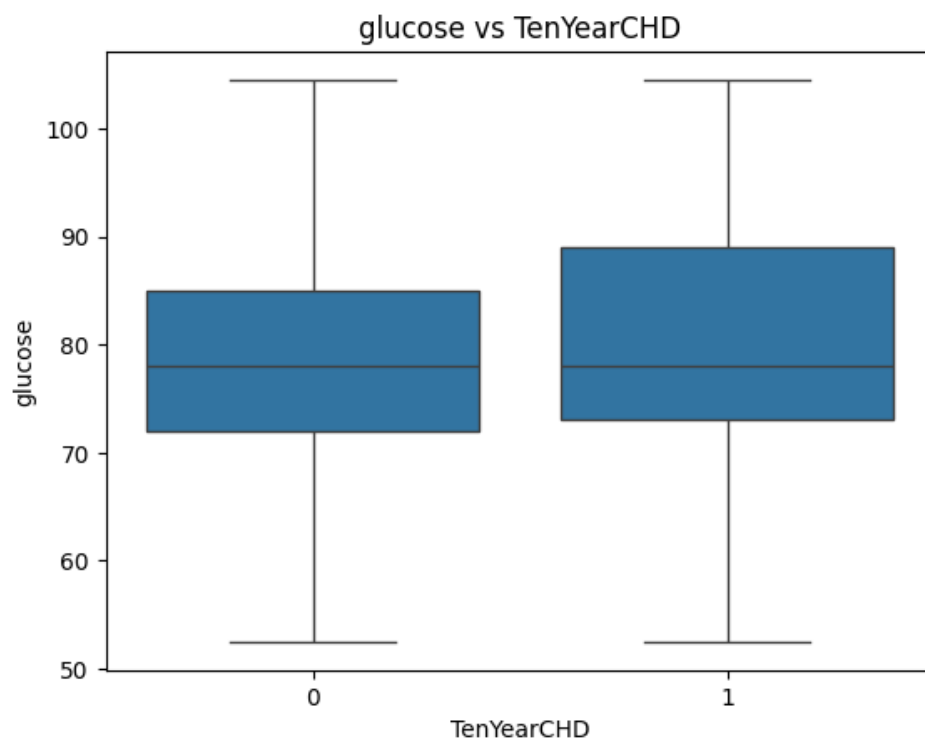
- age, BMI, glucose show right-skewed distributions
- sysBP and diaBP are approximately normal
- Outliers present in glucose and cholesterol

2. Bivariate

1. Numerical Feature vs Target





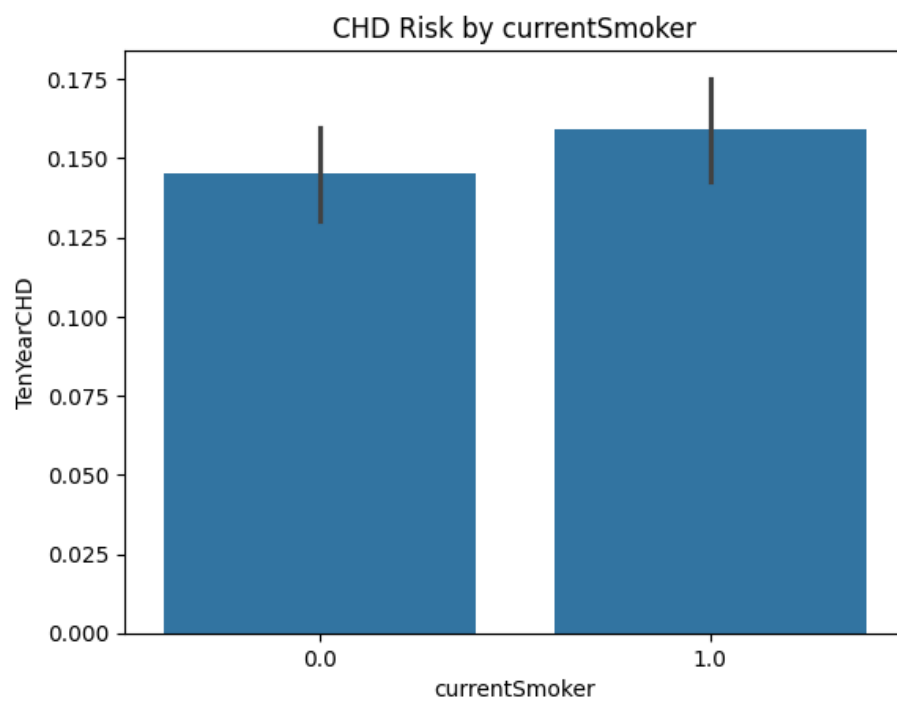
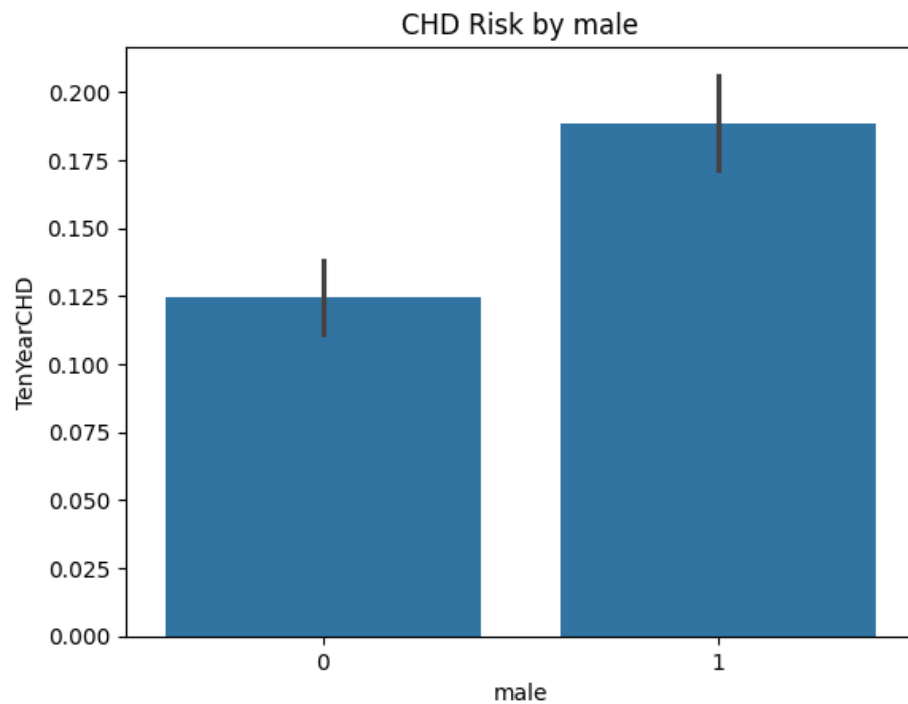


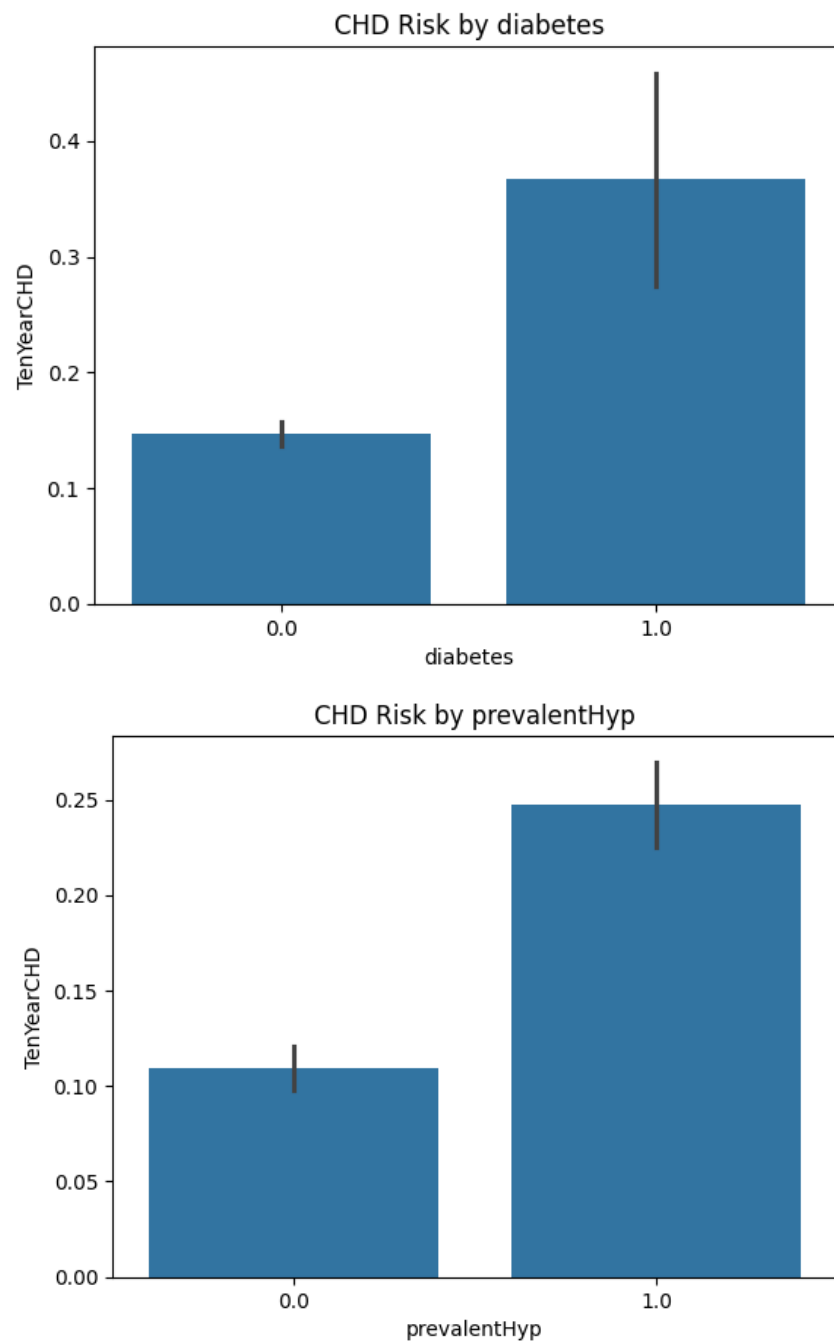
2. Categorical Features vs Target

```
mainipynb M X
Generate + Code + Markdown Run All Restart Clear All Outputs Jupyter Variables Outline Python 3.13.7

Categorical Feature vs Target

1 cat_cols = ['male', 'currentSmoker', 'diabetes', 'prevalentHyp']
2
3 for col in cat_cols:
4     sns.barplot(x=col, y='TenYearCHD', data=df,)
5     plt.title(f"CHD Risk by {col}")
6     plt.show()
```





Bivariate Observations

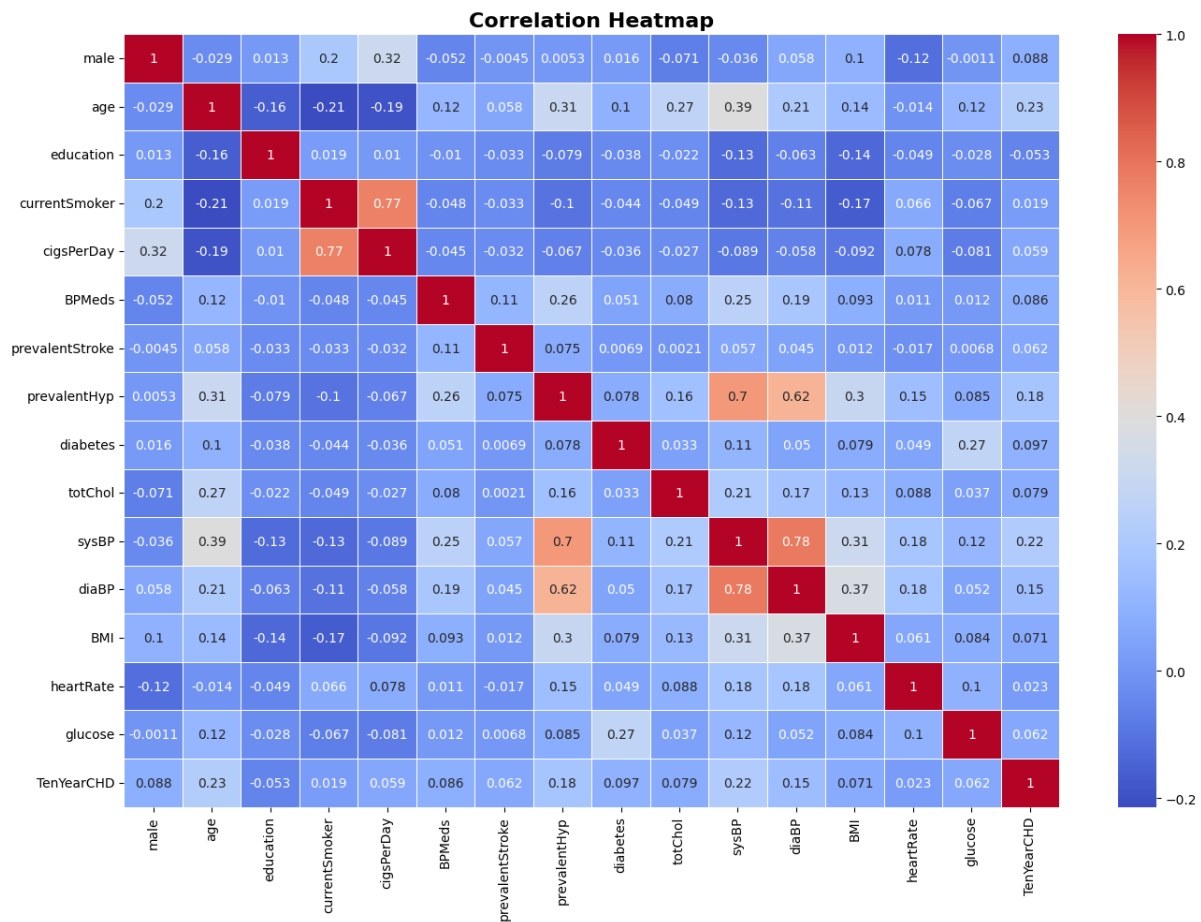
- CHD patients have higher median age and blood pressure
- Diabetes and hypertension show strong association with CHD
- Smoking increases CHD probability

4. Multivariate Analysis

1. Correlation Matrix

```
mainipynb M X
Generate + Code + Markdown | Interrupt Restart Clear All Outputs Go To | Jupyter Variables Outline ... Python 3.13.7

1 plt.figure(figsize=(14, 10))
2 sns.heatmap(
3     df.corr(),
4     mask=None,
5     cmap="coolwarm",
6     linewidths=0.5,
7     annot=True
8 )
9 plt.title("Correlation Heatmap", fontsize=16, fontweight='bold')
10 plt.tight_layout()
11 plt.show()
12
```



Multivariate Observations

- Strong correlation between sysBP and diaBP
- BMI and glucose are positively correlated
- Multicollinearity exists and affects linear models

2. Visualizing distributions

Gokulnaath M 24BAD028 - Jayaraksha Reguraj 24BAD044

1. KDE Distribution by Target

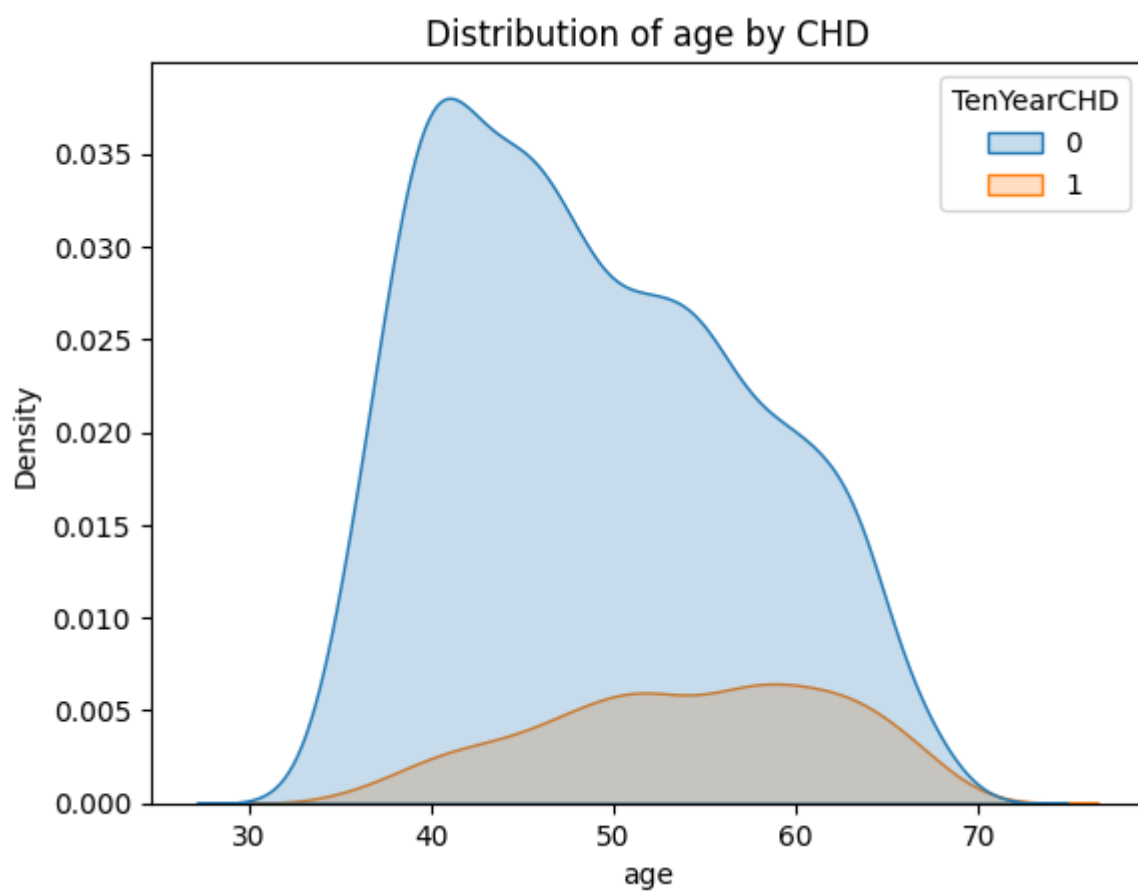
```
main.py M X
Generate + Code + Markdown | Interrupt Restart Clear All Outputs Go To | Jupyter Variables Outline ... Python 3.13.7

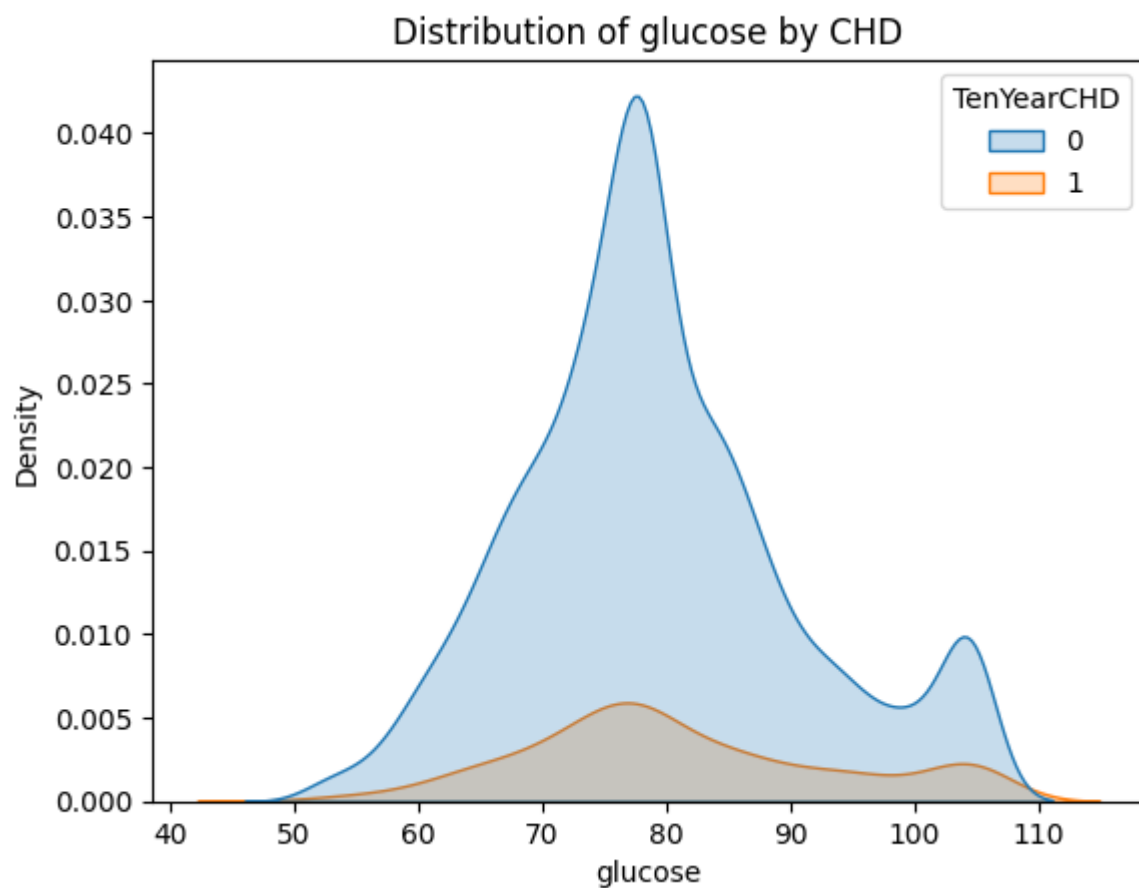
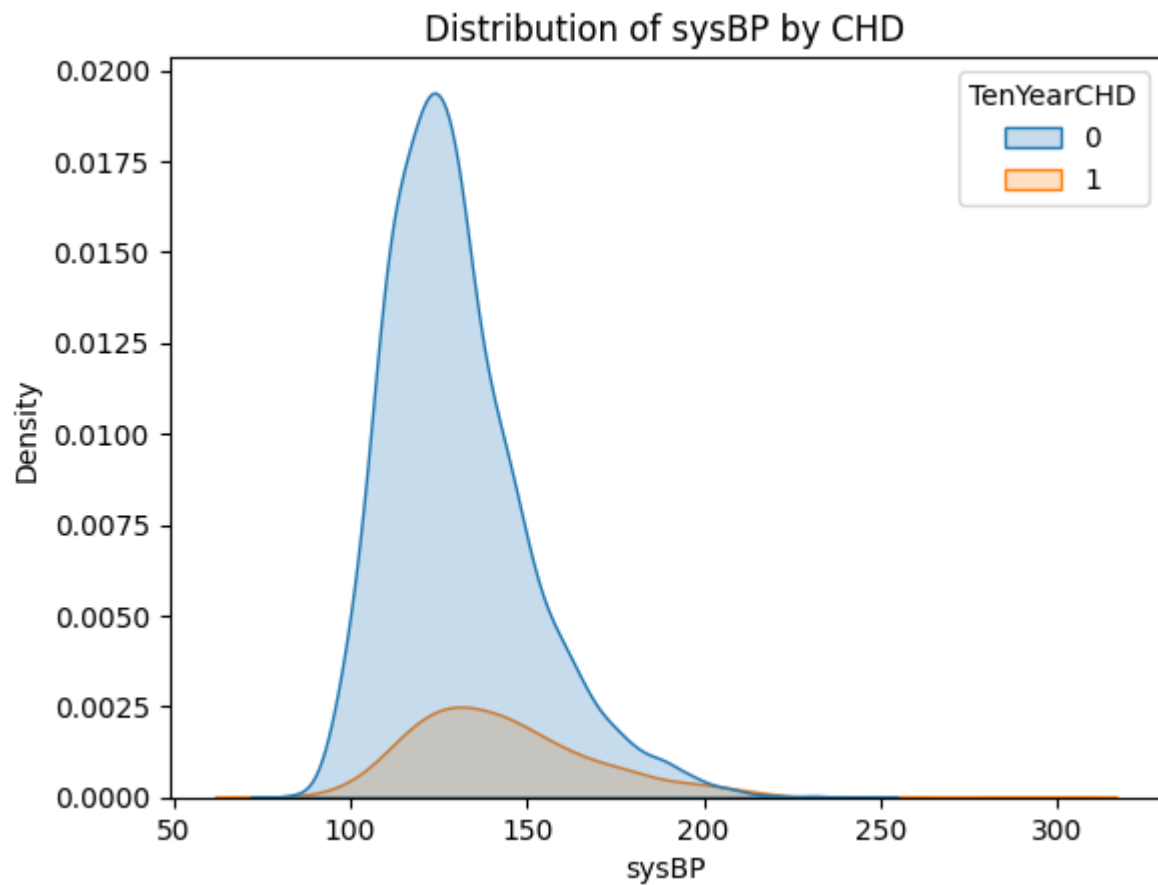
Distribution Analysis

KDE Distribution by Target

1 for col in ['age', 'sysBP', 'glucose']:
2     sns.kdeplot(data=df, x=col, hue='TenYearCHD', fill=True)
3     plt.title(f'Distribution of {col} by CHD')
4     plt.show()

[54] Python
```





2. PCA Distribution by Target

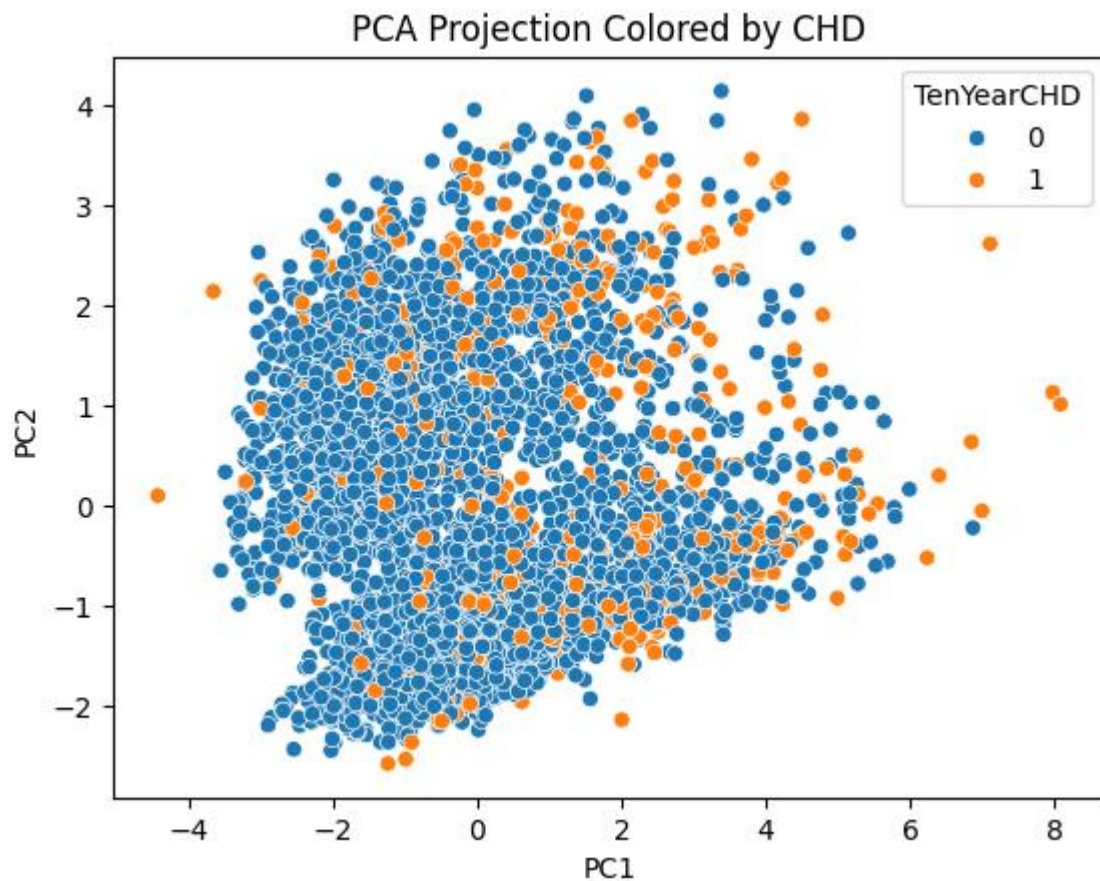
```
mainipynb M X
Generate + Code + Markdown Interrupt Restart Clear All Outputs Go To Jupyter Variables Outline Python 3.13.7

Dimensionality Reduction (PCA)

PCA Implementation

1 from sklearn.preprocessing import StandardScaler
2 from sklearn.decomposition import PCA
3
4 X = df.drop("TenYearCHD", axis=1)
5 y = df["TenYearCHD"]
6
7 X_scaled = StandardScaler().fit_transform(X)
8
9 pca = PCA(n_components=2)
10 pca_result = pca.fit_transform(X_scaled)
11
12 pca_df = pd.DataFrame(pca_result, columns=['PC1', 'PC2'])
13 pca_df["TenYearCHD"] = y
[55] Python

1 sns.scatterplot(x='PC1', y='PC2', hue='TenYearCHD', data=pca_df)
2 plt.title("PCA Projection Colored by CHD")
3 plt.show()
[56] Python
```



Resources and References

- **Dataset**
 - <https://www.kaggle.com/datasets/dileep070/heart-disease-prediction-using-logistic-regression>
- **GitHub**
 - https://github.com/Gokulnaath-gif/24ADI204_DSV_Team14
- **Resources**
 - <https://www.geeksforgeeks.org/data-analysis/what-is-exploratory-data-analysis/>
 - <https://www.ibm.com/think/topics/exploratory-data-analysis>
 - <https://www.kaggle.com/code/imoore/intro-to-exploratory-data-analysis-eda-in-python>
 - <https://x.ai/grok>
 - <https://openai.com/chatgpt>
 - <https://google.com/gemini>