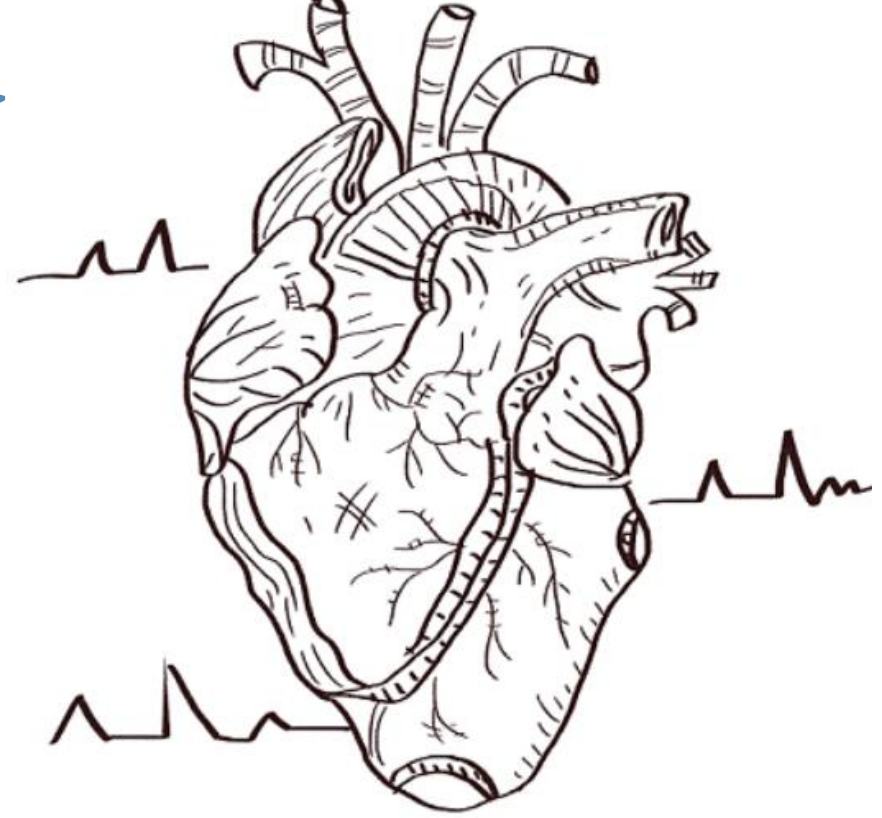


Heart Disease Prediction



GROUP-2

- PRIYANKA AKULA
- SOWJANYA PAMULAPATI
- SOMNA SATTOOR
- SHREYA CHIKATMARLA

Contents

Introduction

Motivations & Objectives

Hybrid Crisp-DM & Waterfall Model

Project Workflow

Data Set Description

Data Cleaning

Visualizations

Algorithms

Conclusion

References

Introduction

Heart disease is a significant health concern worldwide, and early detection plays a crucial role in improving patient outcomes.

In this project, we aim to utilize the power of visualizations and also develop a predictive model that can identify individuals at risk of developing heart disease.

By combining data analysis techniques with effective visual representations, we can enhance our understanding of the factors contributing to heart disease and provide valuable insights for healthcare professionals and patients alike.

Objectives

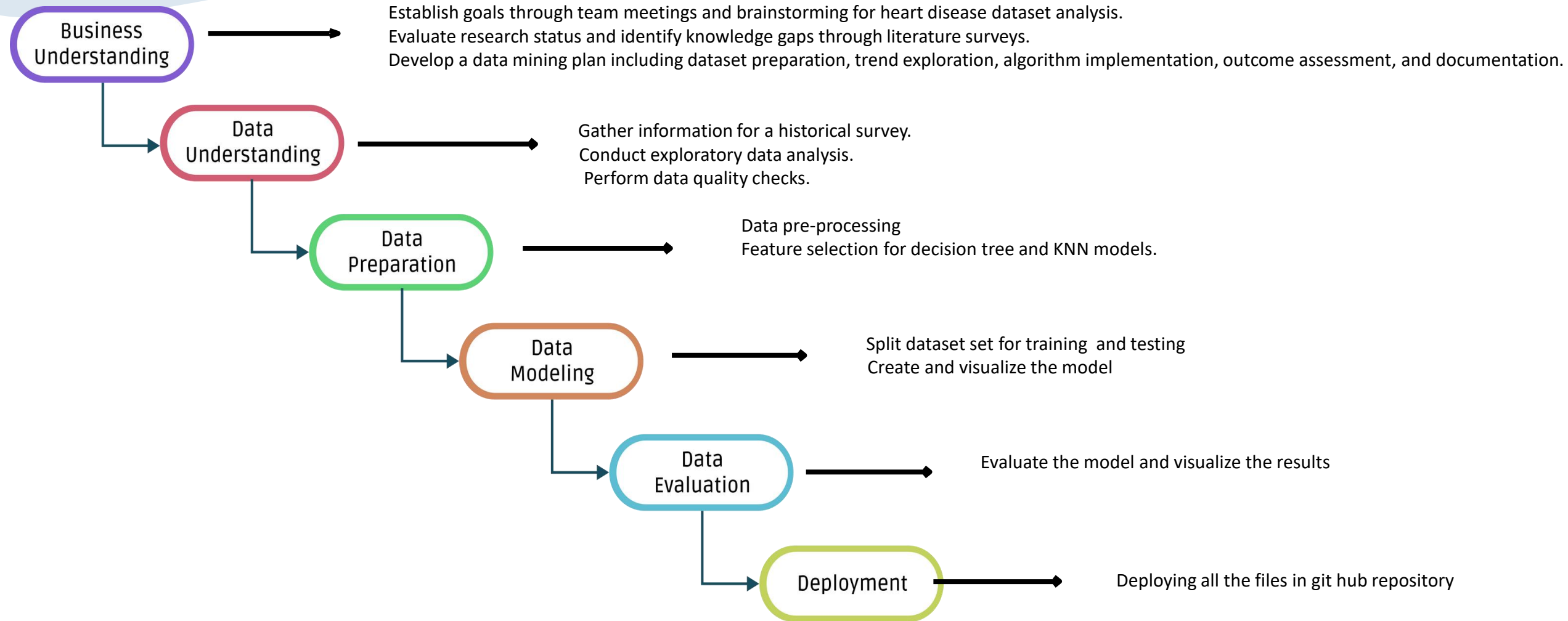


Motivations

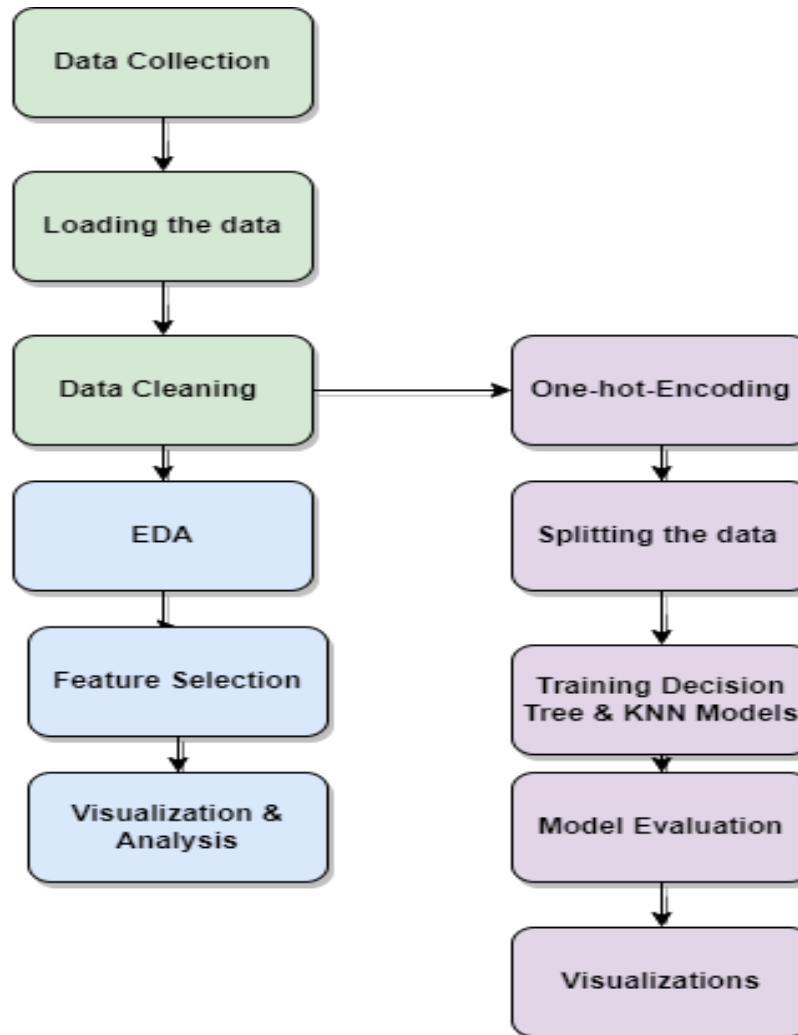
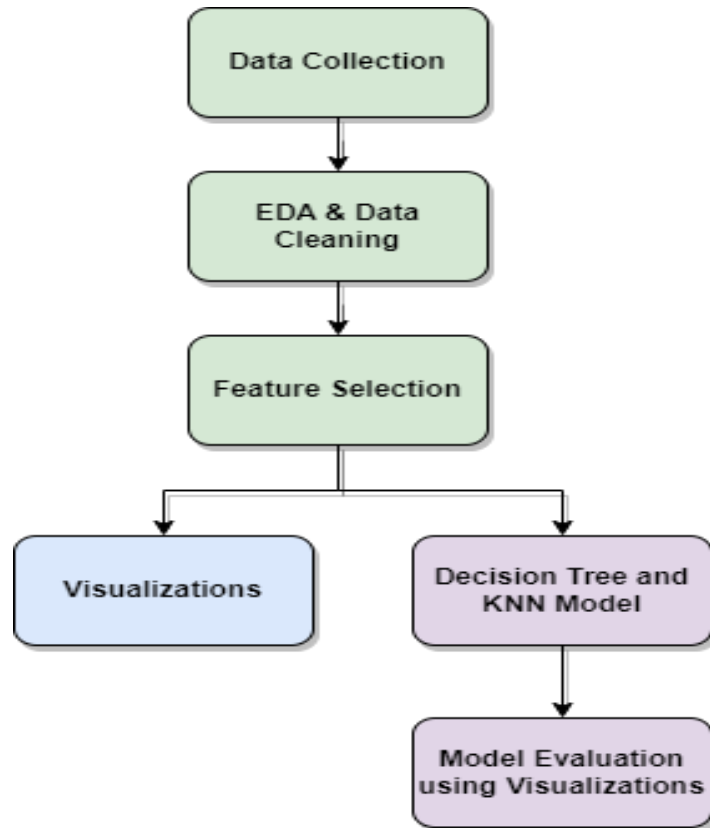
The motivation is using visualizations to better understand heart disease and create a predictive algorithm to detect at-risk individuals.

This can enhance heart disease prevention and management, improving patient outcomes and healthcare expenditures.

Hybrid Crisp-DM & Waterfall Model



Project Workflow



Dataset Description

The Heart Disease Prediction dataset was taken from Kaggle. It includes information from a study conducted on patients with suspected heart disease.

The dataset contains 18 columns of patient data, including age, BMI, blood pressure, cholesterol levels, heart disease, diabetic and other pertinent health indicators.

There are 319796 rows in the dataset, each of which represents a different patient. The target variable is heart disease which represents whether a patient having heart disease or not.

```
df1.describe().T
```

	count	mean	std	min	25%	50%	75%	max
Age	270.0	54.433333	9.109067	29.0	48.0	55.0	61.0	77.0
Sex	270.0	0.677778	0.468195	0.0	0.0	1.0	1.0	1.0
Chest pain type	270.0	3.174074	0.950090	1.0	3.0	3.0	4.0	4.0
BP	270.0	131.344444	17.861608	94.0	120.0	130.0	140.0	200.0
Cholesterol	270.0	249.659259	51.686237	126.0	213.0	245.0	280.0	564.0
FBS over 120	270.0	0.148148	0.355906	0.0	0.0	0.0	0.0	1.0
EKG results	270.0	1.022222	0.997891	0.0	0.0	2.0	2.0	2.0
Max HR	270.0	149.677778	23.165717	71.0	133.0	153.5	166.0	202.0
Exercise angina	270.0	0.329630	0.470952	0.0	0.0	0.0	1.0	1.0
ST depression	270.0	1.050000	1.145210	0.0	0.0	0.8	1.6	6.2
Slope of ST	270.0	1.585185	0.614390	1.0	1.0	2.0	2.0	3.0
Number of vessels fluro	270.0	0.670370	0.943896	0.0	0.0	0.0	1.0	3.0
Thallium	270.0	4.696296	1.940659	3.0	3.0	3.0	7.0	7.0

Data Cleaning

- Checked null values, duplicate entries, outliers and converted the categorical data into numeric and removed the unnecessary data.

```
df[['Diabetic', 'Sex']].T
```

	0	1	2	3	4	5	6	7	8	9	...	319785	319786	319787	319788	319789	319790	319791	319792	...
Diabetic	Yes	No	Yes	No	No	No	No	Yes	No, borderline diabetes	No	...	No	Yes	No	No	No	Yes	No	No	...
Sex	Female	Female	Male	Female	Female	Female	Female	Female	Female	Male	...	Male	Female	Male	Female	Female	Male	Male	Female	Female

2 rows × 319795 columns

```
df = df[df.columns].replace({'Yes':1, 'No':0, 'Male':1, 'Female':0, 'No, borderline diabetes':'0', 'Yes (during pregnancy)':1 })
df['Diabetic'] = df['Diabetic'].astype(int)
```

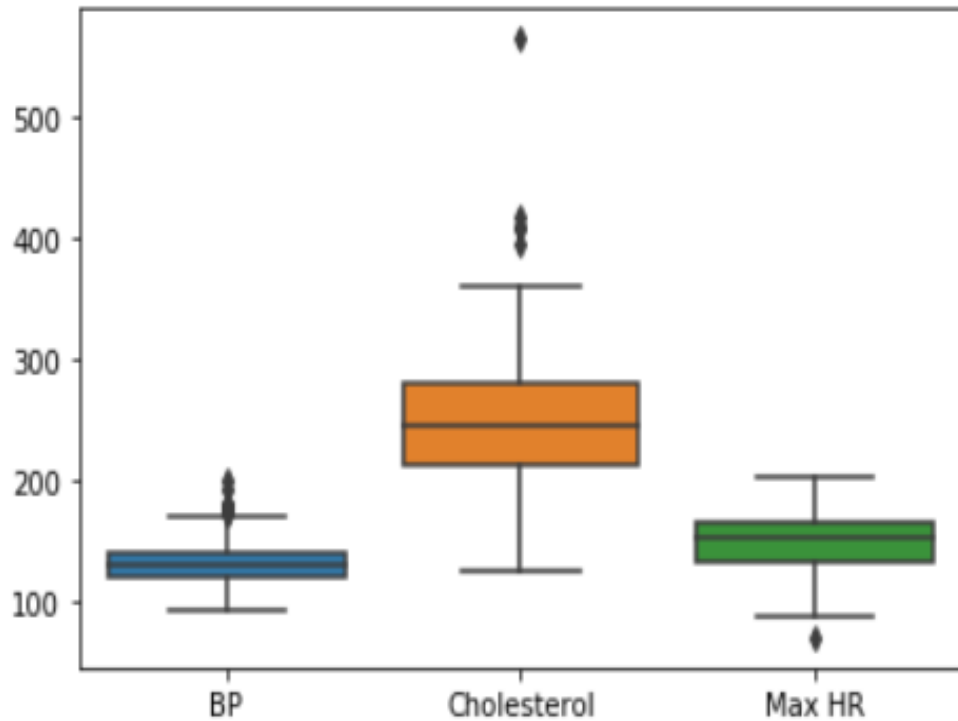
```
df[['Diabetic', 'Sex']].T
```

	0	1	2	3	4	5	6	7	8	9	...	319785	319786	319787	319788	319789	319790	319791	319792	319793	319794
Diabetic	1	0	1	0	0	0	0	1	0	0	...	0	1	0	0	0	1	0	0	0	0
Sex	0	0	1	0	0	0	0	0	0	1	...	1	0	1	0	0	1	1	0	0	0

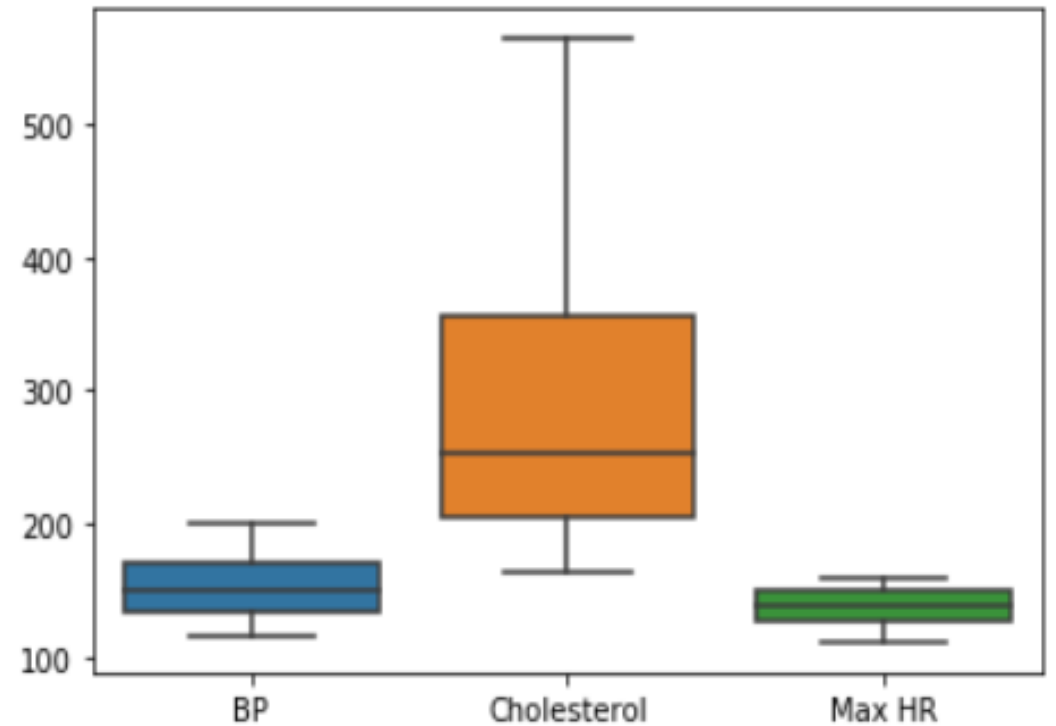
2 rows × 319795 columns

Data Cleaning

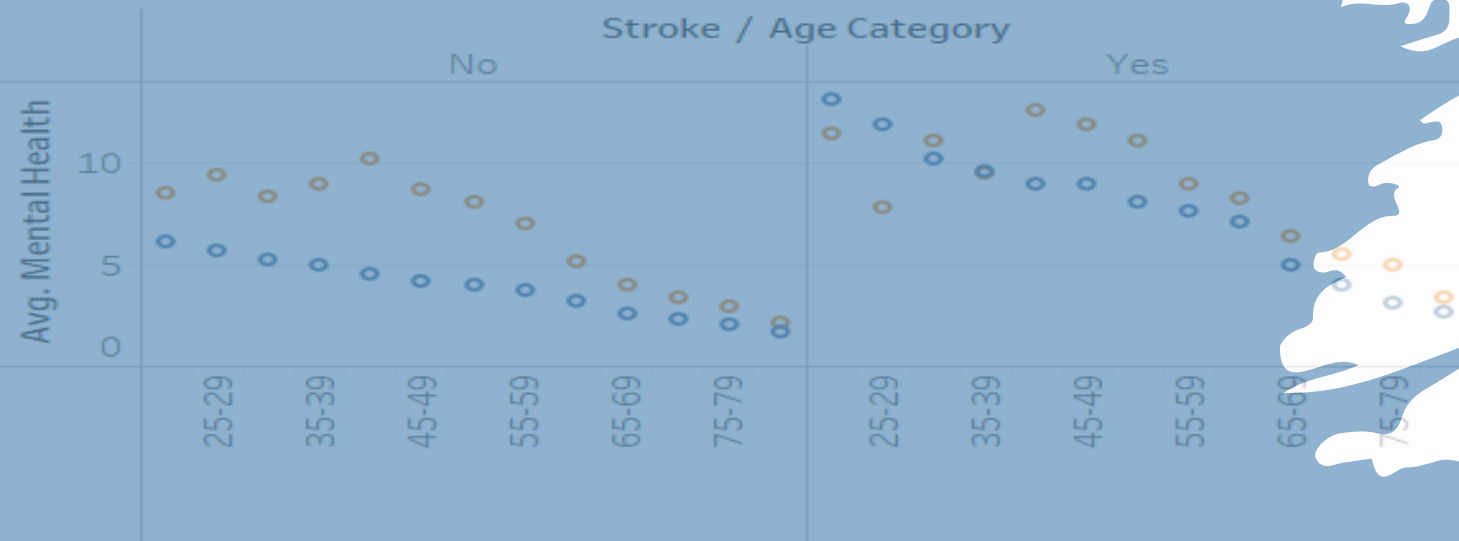
Box plot with outliers



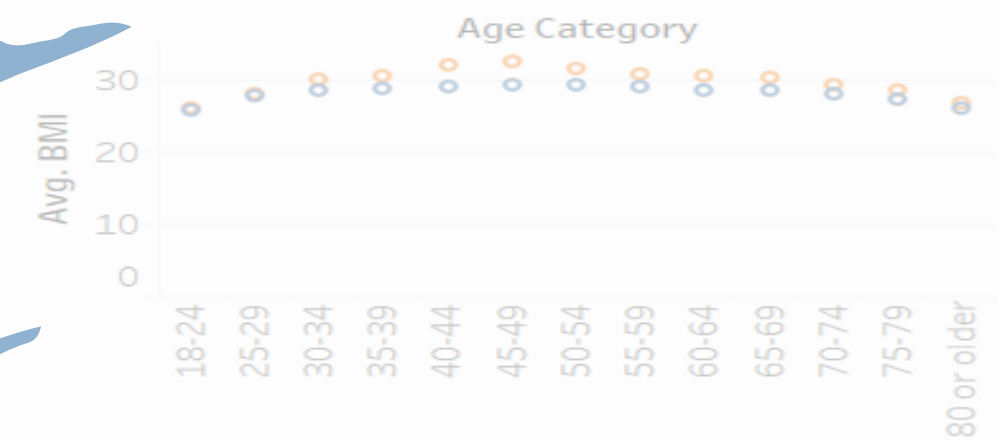
Box plot after removing outliers



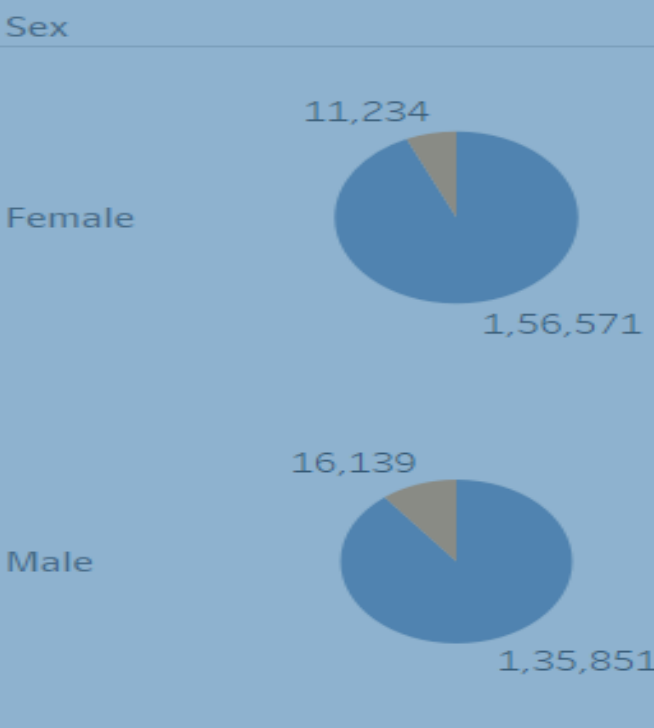
Heart disease based on age and stroke with average mental health.



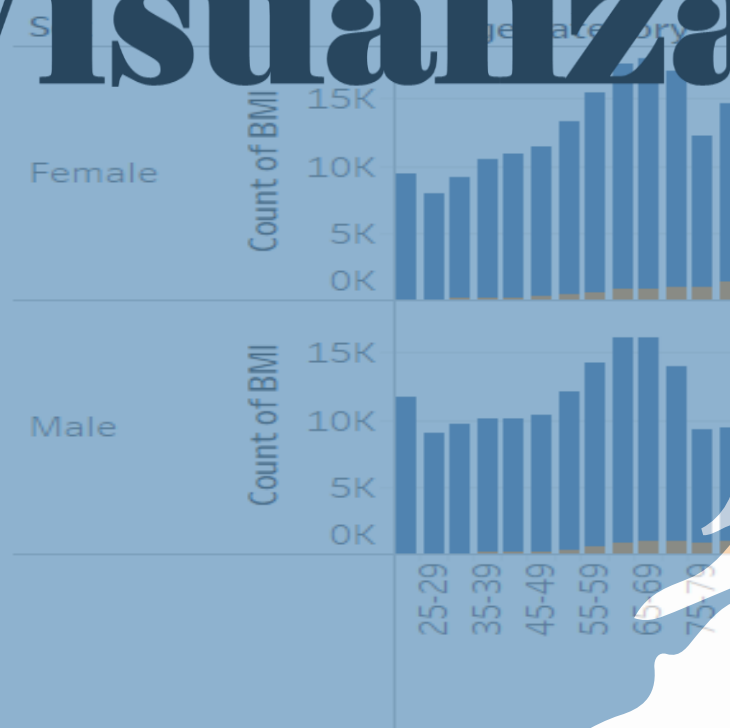
Heart disease based on age and average BMI value



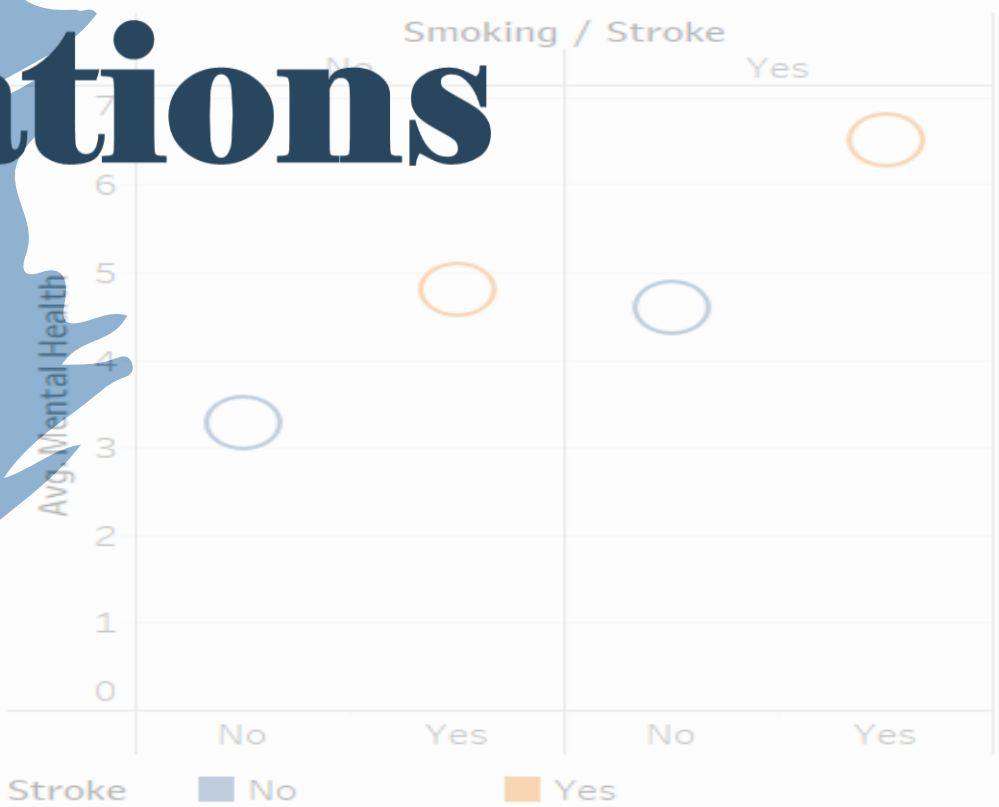
Heart Disease based on sex and BMI value



Heart stroke based on age, sex and BMI value

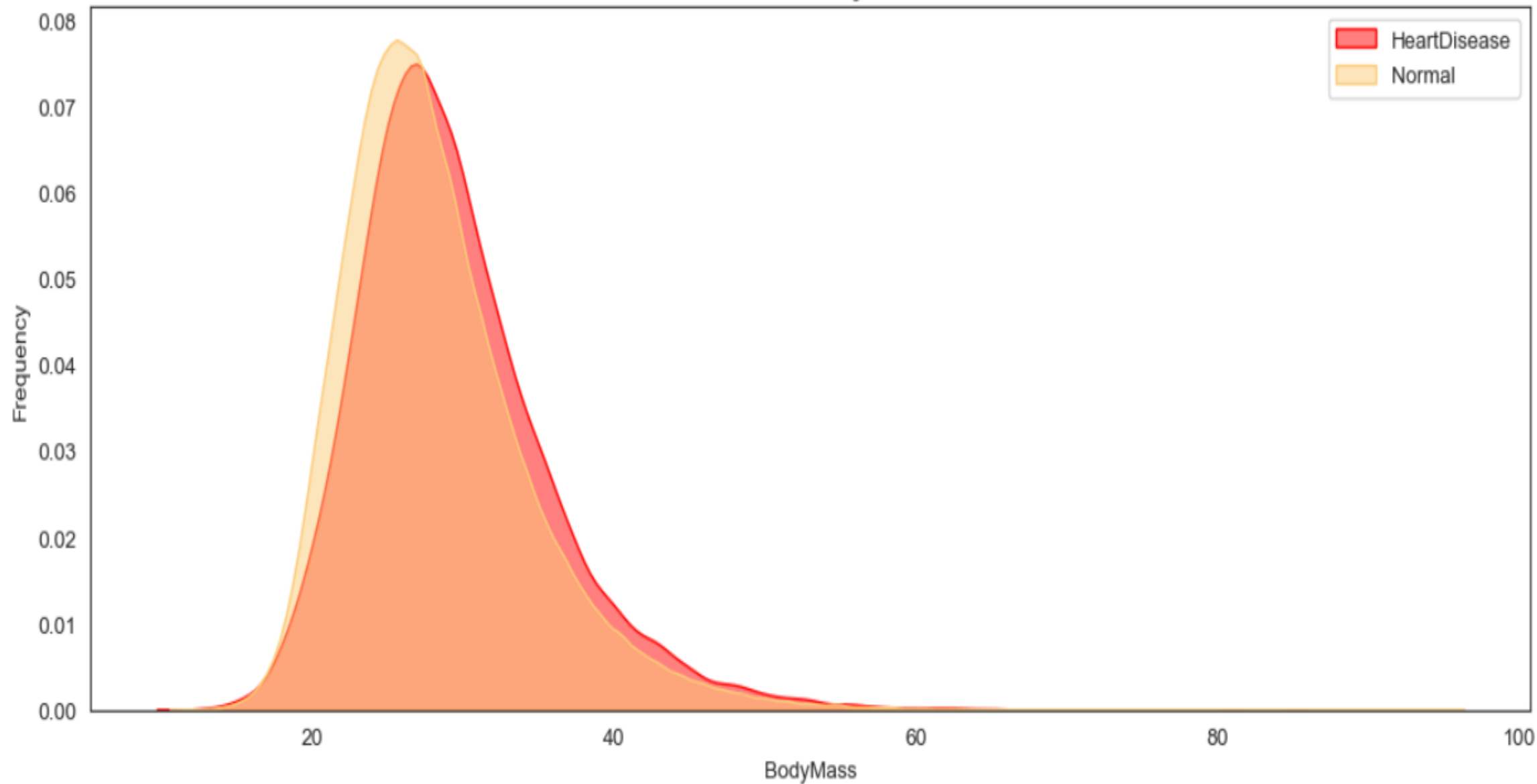


Heart Stroke for people with average mental health and smoking habit.

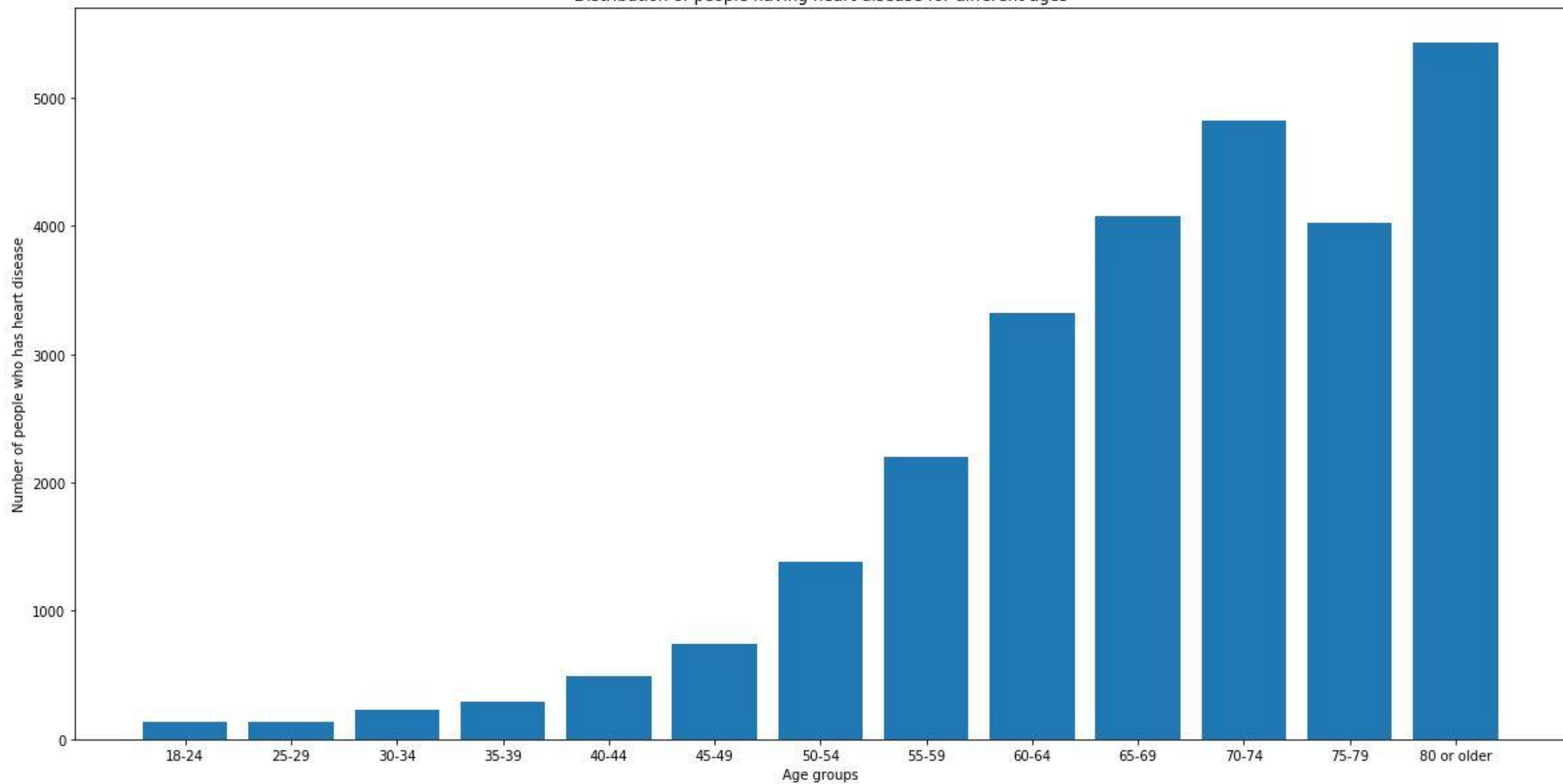


Visualizations

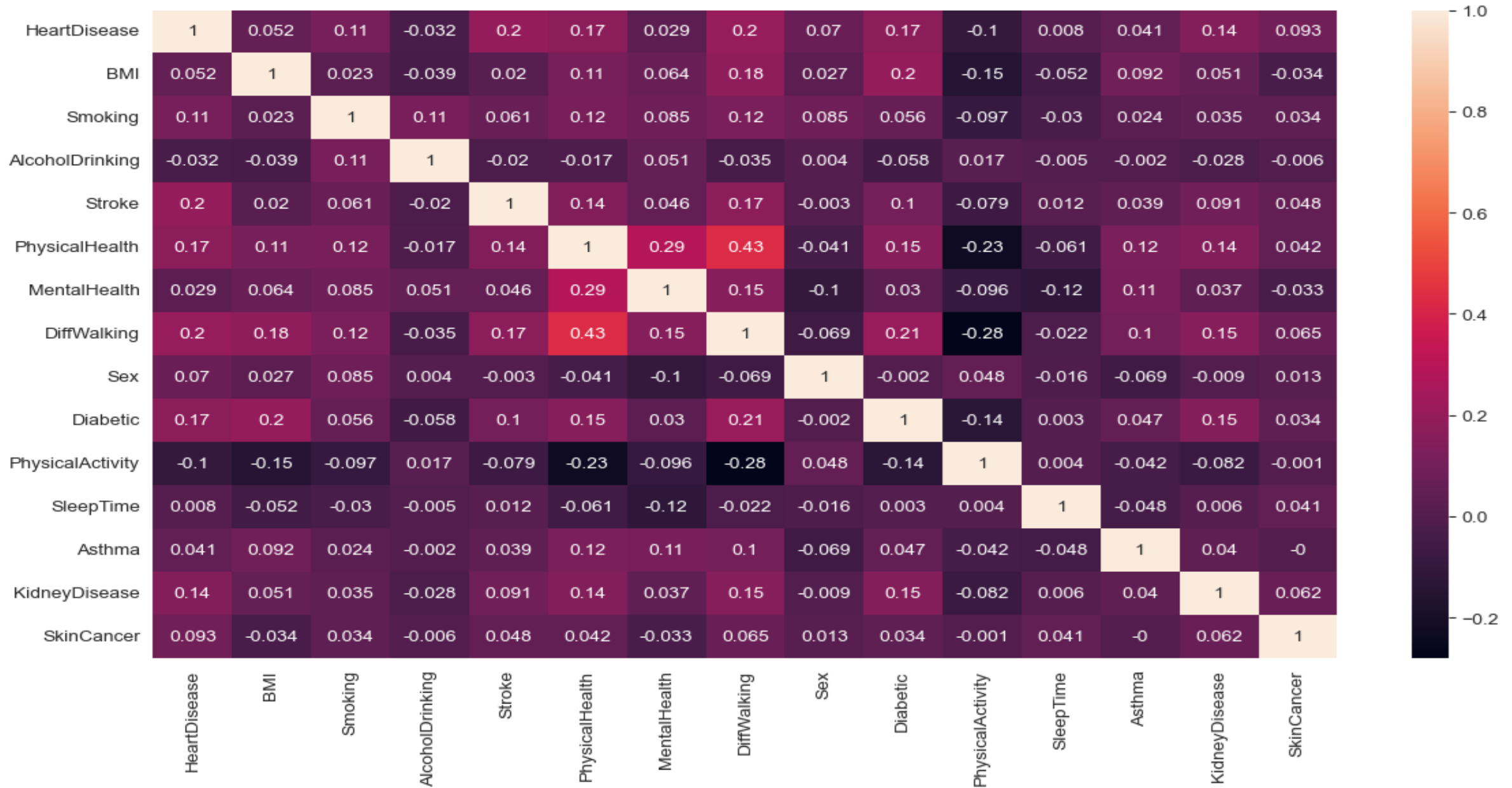
Distribution of Body Mass Index



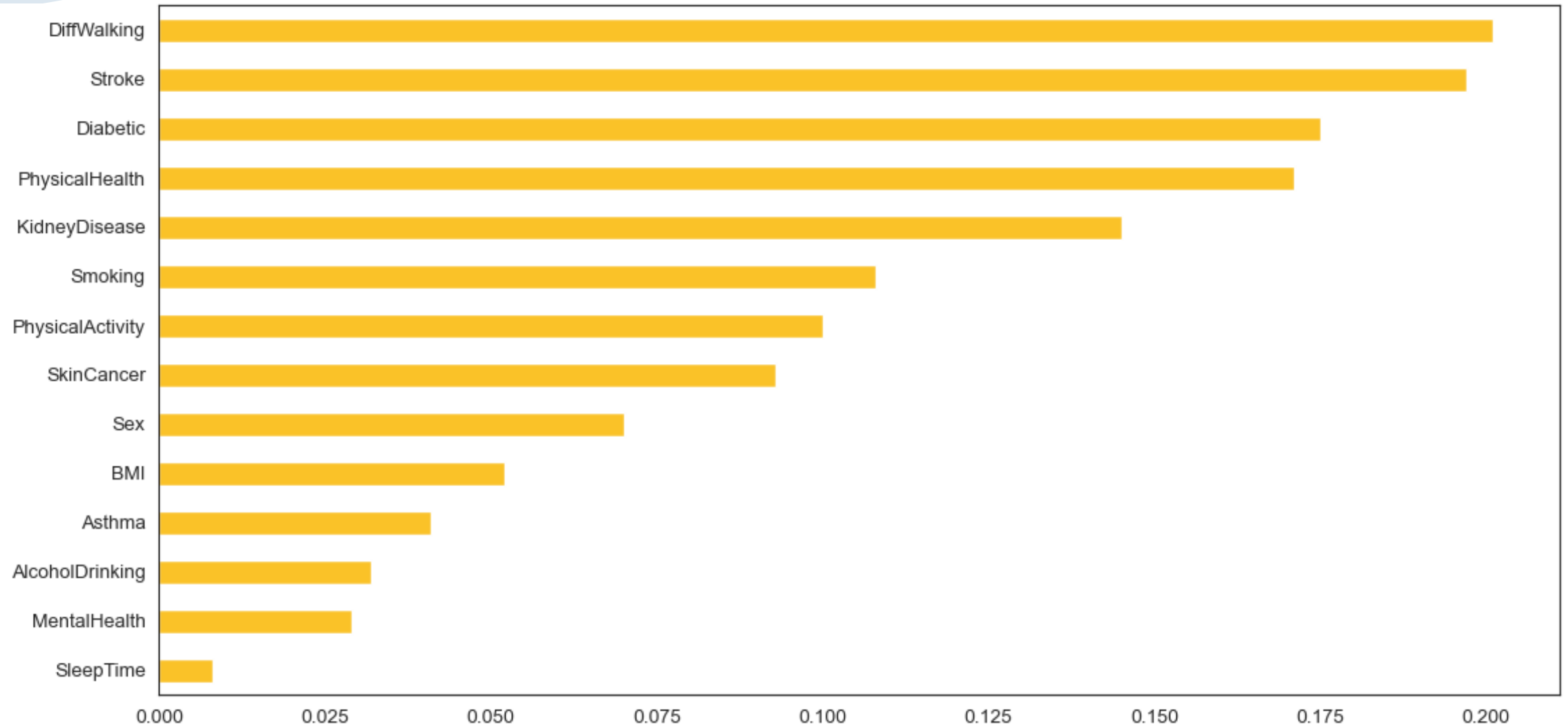
Distribution of people having heart disease for different ages



Correlation between all set of features



Distribution of correlation between the features



Algorithms

The data is spilt into train and test sets with a ratio of 70:30.

Train set is used for both models while test set is used for calculating evaluation metrics.

Decision Tree

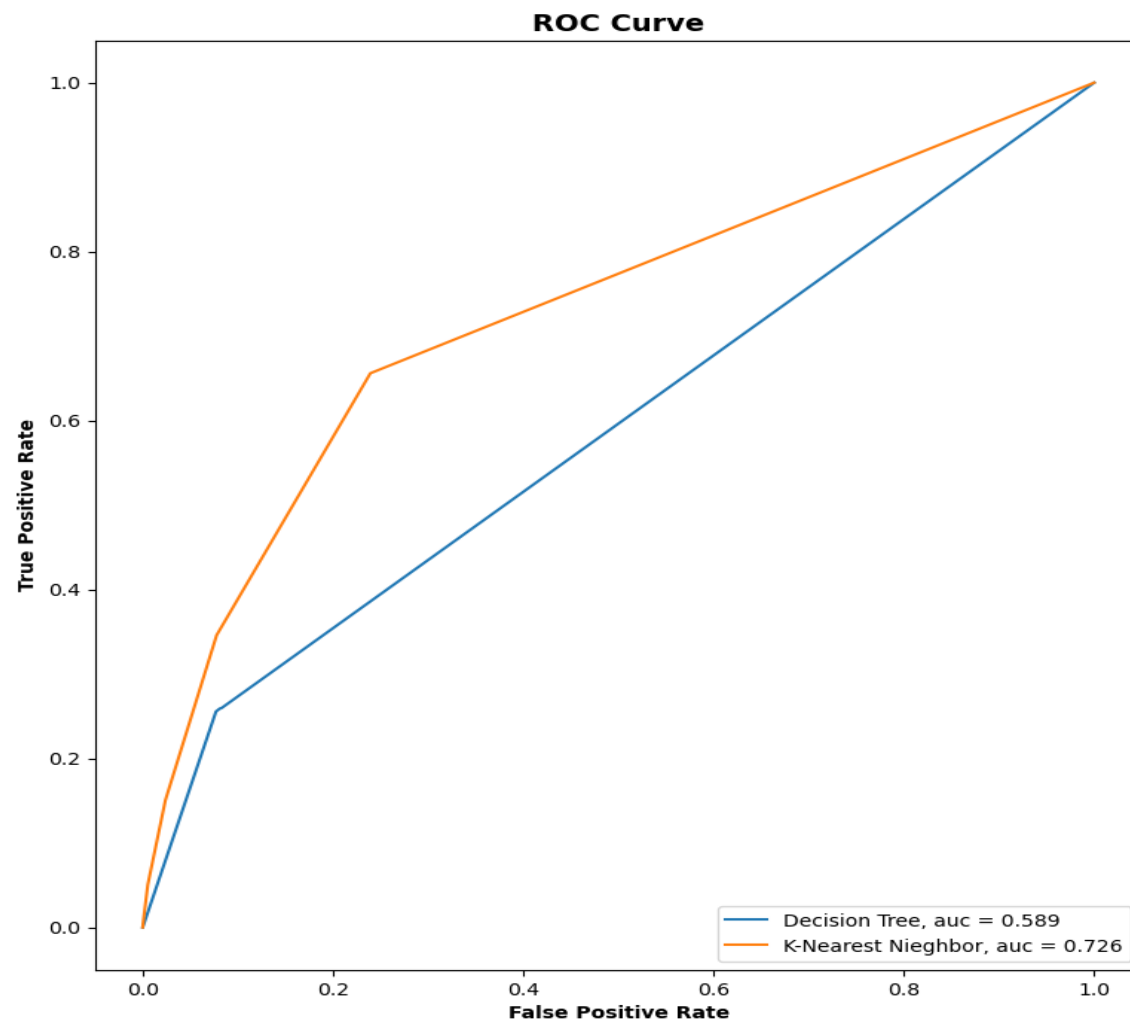
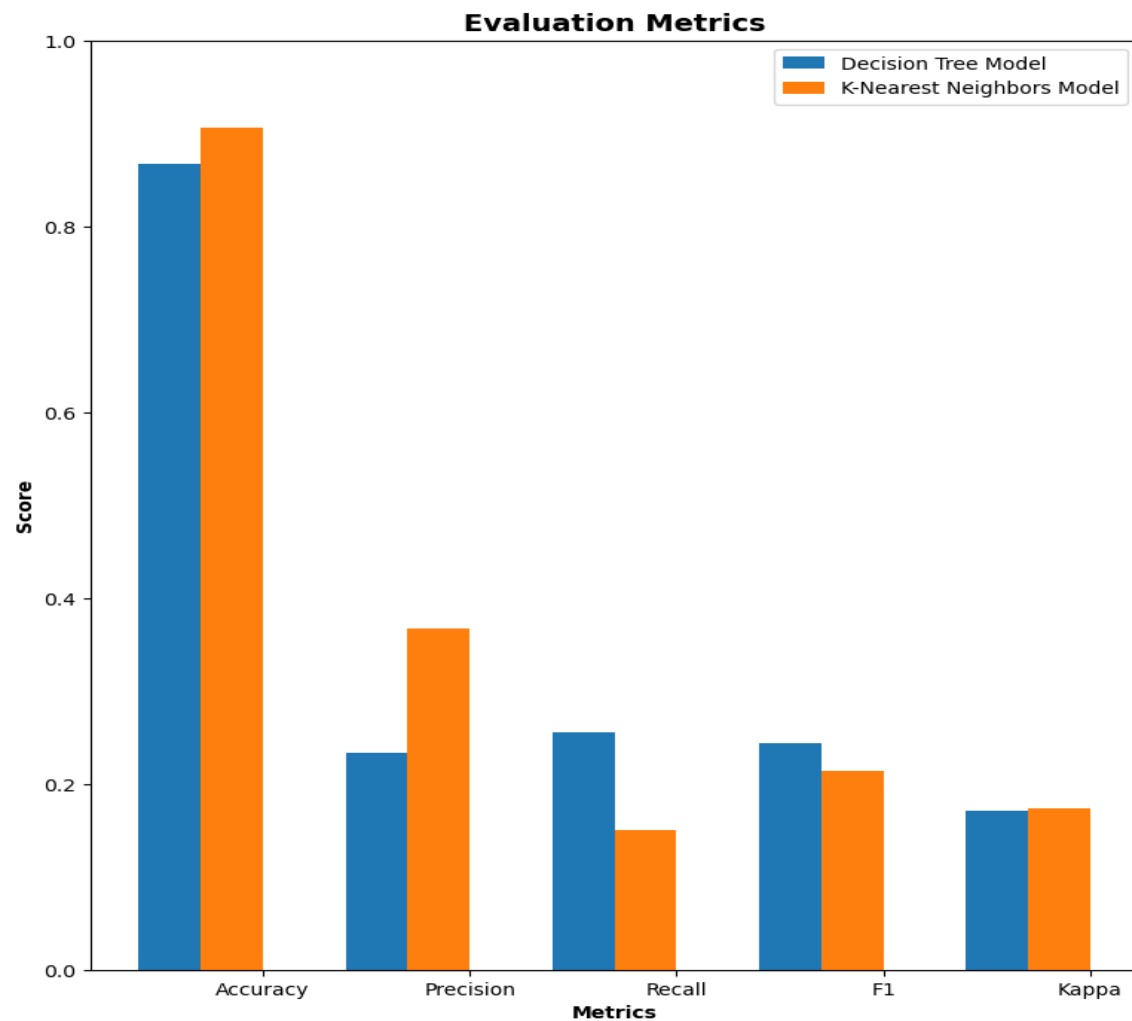
- Since the data set is having both categorical and numerical values where the categorical columns are containing less unique values like yes, no and female, male, we are training decision tree classification model to classify data into yes or no.
- F1 Score = 0.244
- Accuracy = 0.86

K- nearest Neighbors Model

- Another classification model that we taught would work best for this data is KNN since there is good correlation between different features and heart diseases.
- F1 Score = 0.213
- Accuracy = 0.90

Model Evaluation

Comparing the models



Conclusion

In conclusion, the heart disease study project yielded substantial insights regarding the incidence, causes, and mitigation strategies linked with heart disease.

It is possible to discover essential factors and develop efficient strategies for timely diagnosis and prevention using data analysis and statistical modeling.

It is critical to ensure the quality of data in order to generate reliable results.

In general, this also improves comprehension and results in the treatment of heart disease.

References

- <https://github.com/Priyankaakula/DATA230-PROJECT>
- <https://www.kaggle.com/code/andls555/heart-disease-prediction/notebook>
- <https://link.springer.com/article/10.1007/s42979-020-00365-y>

A large, horizontal, blue brushstroke shape with irregular, feathered edges, resembling a paint stroke. The color is a medium blue. In the center of this shape, the words "Thank you" are written in a white, rounded, sans-serif font.

Thank you