# Project Title : Heart Attack Risk Prediction

**Humayera Tabassum Tunan**

# Table of Contents

# Introduction

Globally, heart disease is one of the main causes of death. The therapy and prevention of cardiovascular illnesses greatly depend on the capacity to forecast the risk of heart attacks. In this study, we forecast the risk of heart attack based on multiple lifestyle and health characteristics by using machine learning algorithms.

# Dataset Description

**Source:**

[https://www.kaggle.com/datasets/iamsouravbanerjee/heart-attack-prediction-dataset](https://www.kaggle.com/datasets/iamsouravbanerjee/heart-attack-prediction-dataset)

Acknowledgement:

This dataset is a synthetic creation generated using ChatGPT.

Thumbnail by: vectorjuice on Freepik

**Dataset Description:**

We have 26 features in our dataset.

1. Patient ID - Unique identifier for each patient

2. Age - Age of the patient

3. Sex - Gender of the patient (Male/Female)

4. Cholesterol - Cholesterol levels of the patient

5. Blood Pressure - Blood pressure of the patient (systolic/diastolic)

6. Heart Rate - Heart rate of the patient

7. Diabetes - Whether the patient has diabetes (Yes/No)

8. Family History - Family history of heart-related problems (1: Yes, 0: No)

9. Smoking - Smoking status of the patient (1: Smoker, 0: Non-smoker)

10. Obesity - Obesity status of the patient (1: Obese, 0: Not obese)

11. Alcohol Consumption - Level of alcohol consumption by the patient (None/Light/Moderate/Heavy)

12. Exercise Hours Per Week - Number of exercise hours per week

13. Diet - Dietary habits of the patient (Healthy/Average/Unhealthy)

14. Previous Heart Problems - Previous heart problems of the patient (1: Yes, 0: No)

15. Medication Use - Medication usage by the patient (1: Yes, 0: No)

16. Stress Level - Stress level reported by the patient (1-10)

17. Sedentary Hours Per Day - Hours of sedentary activity per day

18. Income - Income level of the patient

19. BMI - Body Mass Index (BMI) of the patient

20. Triglycerides - Triglyceride levels of the patient

21. Physical Activity Days Per Week - Days of physical activity per week

22. Sleep Hours Per Day - Hours of sleep per day

23. Country - Country of the patient

24. Continent - Continent where the patient resides

25. Hemisphere - Hemisphere where the patient resides

26. Heart Attack Risk - Presence of heart attack risk (1: Yes, 0: No)

**How many features?**

There are 25 features in the dataset.

**Classification or regression problem?**

It is a classification problem because the aim is to predict whether a person is at risk of a heart attack based on various health-related factors. The target variable "Heart Attack Risk" is categorical.

**How many data points?**

There are 8763 data points (patient records) in the dataset.

**What kind of features are in your dataset?**

Numerical Features:

1. Age

2. Cholesterol

3. Blood Pressure

4. Heart Rate

5. Exercise Hours Per Week

6. Previous Heart Problems

7. Stress Level

8. Sedentary Hours Per Day

9. Income

10. BMI

11. Triglycerides

12. Physical Activity Days Per Week

13. Sleep Hours Per Day

Categorical Features:

1. Sex

2. Diabetes

3. Family History

4. Smoking

5. Obesity
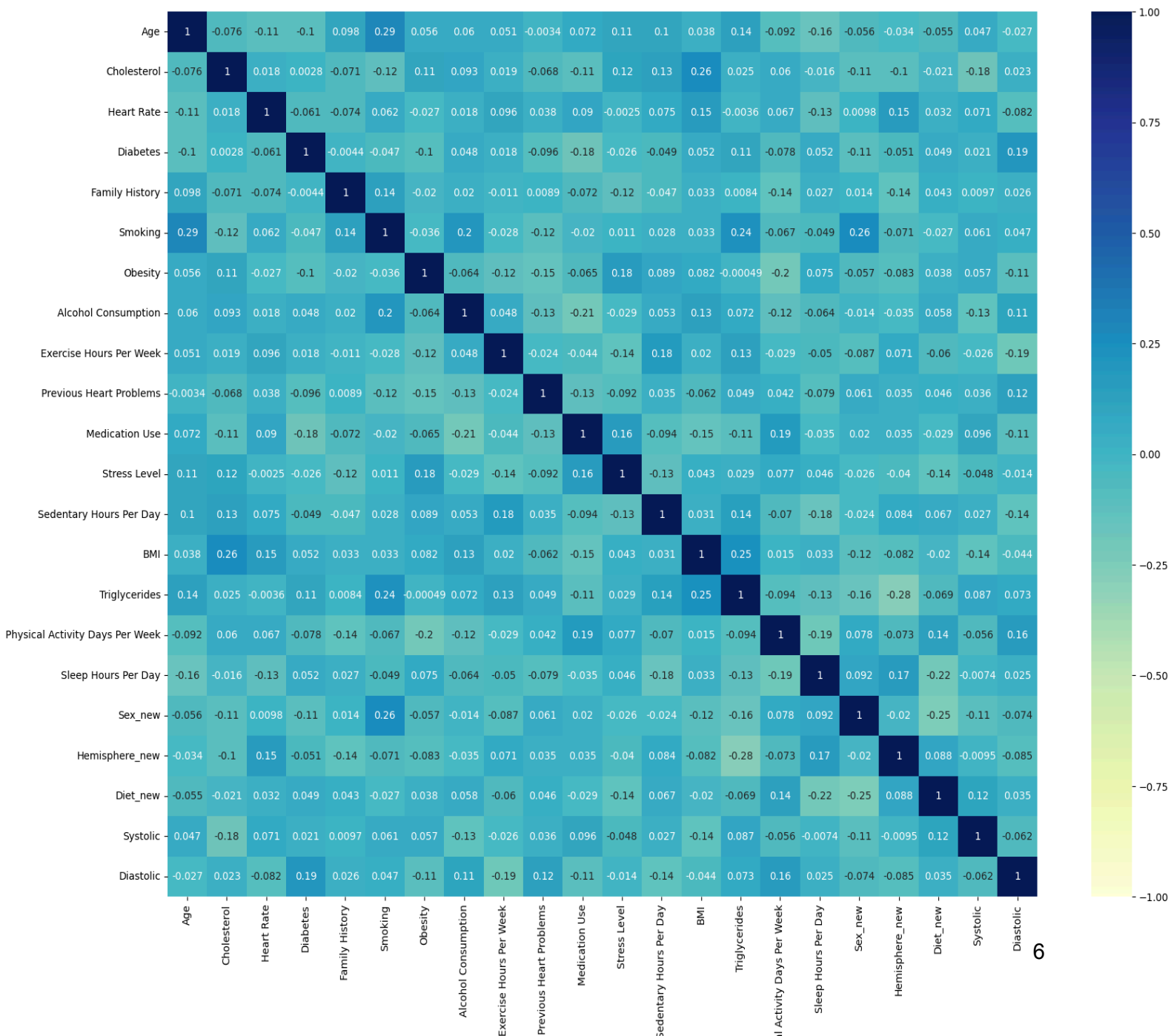
6. Alcohol Consumption

7. Diet

8. Medication Use

9. Country

10. Continent

11. Hemisphere

## Correlation of all features:

# Data pre- processing

**Problem 1: Null Values.**

[The dataset did not have any null values so I inserted 16,000 random null values for demonstration]

Solution: Removing Null Values

In order to remove null values from the dataset, I found some data points where the highly important features had null values. These data points will be of no use to our project, that's why I dropped them using dropna().

**Problem 2: Irrelevant features.**

Solution: Feature selection

There were some features['Patient ID', 'Continent', 'Country', 'Income'] that were irrelevant to our project, which will be of no use. I dropped those features using drop().

**Problem 3: Missing values.**

Solution: Imputing Missing Values

Some of the numerical features['Sedentary Hours Per Day','Physical Activity Days Per Week', 'Sleep Hours Per Day', 'Income', 'Exercise Hours Per Week'] had some missing values. I replaced those missing values with the average of all values in those features by using mean(). For the binary features['Family History', 'Smoking', 'Obesity', 'Alcohol Consumption'] I replaced the missing values with 0 by using fillna(0).

**Problem 4: Categorical Features.**

Solution: Labeling/One hot Encoding

There were three categorical features: "Diet", "Sex" and "Hemisphere". I encoded "Sex" and "Hemisphere" into binary values using LabelEncoder(). For "Diet", I used map() function and mapped 'Healthy' as 2, 'Unhealthy' as 0, and 'Average' as 1. We can also solve this using One Hot encoding, as shown for demonstration.

**Problem 5: Features with multiple values.**

Solution: Feature Engineering.

I separated the feature "Blood Pressure" into two different features Systolic and Diastolic.

# Feature Scaling

Features that we edited were removed, also unnecessary features were removed. No further Feature scaling was required.
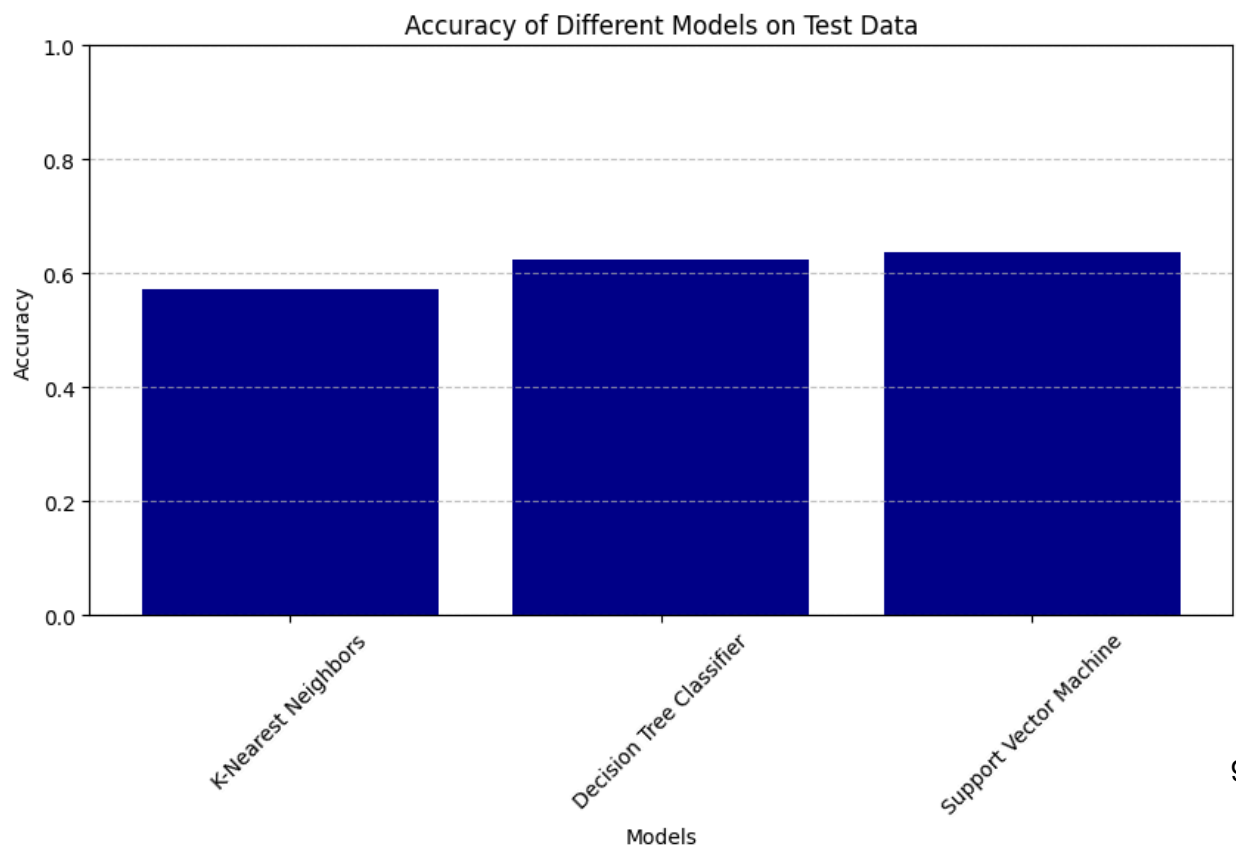
# Dataset Splitting

Random Split is done on the dataset using the train_test_split function from SK learn library. Test_size = 0.3 represents the 30% data used for training and the remaining 70% were used for testing .
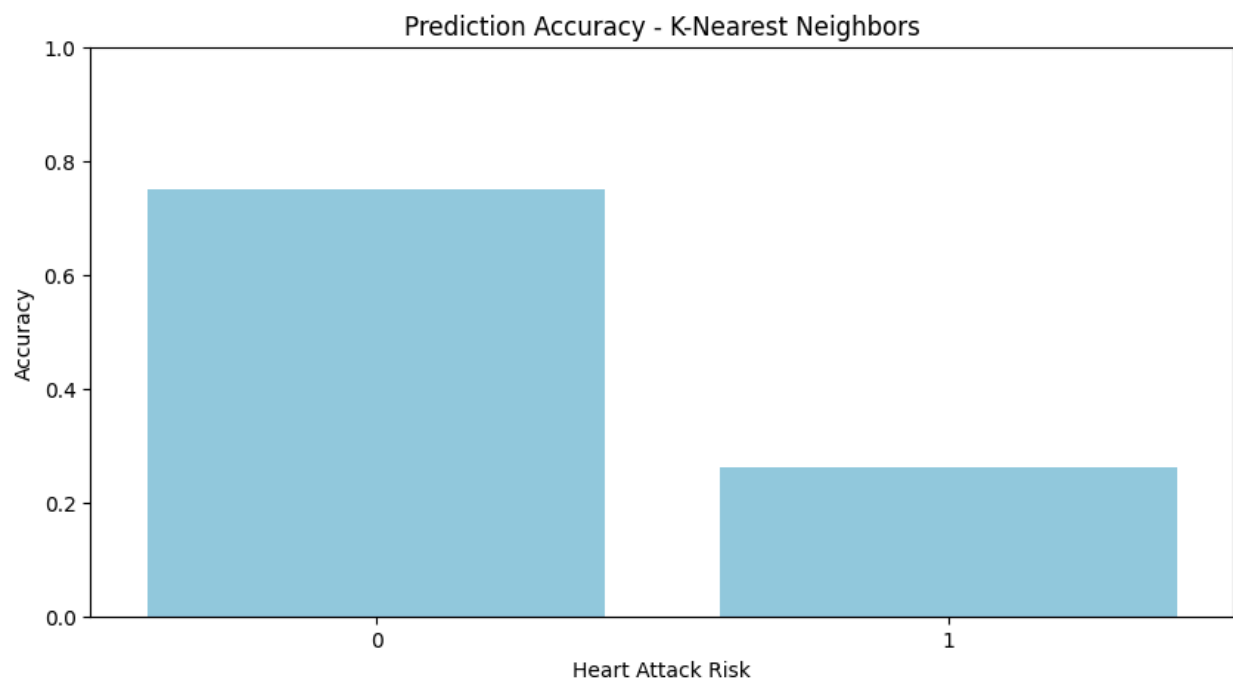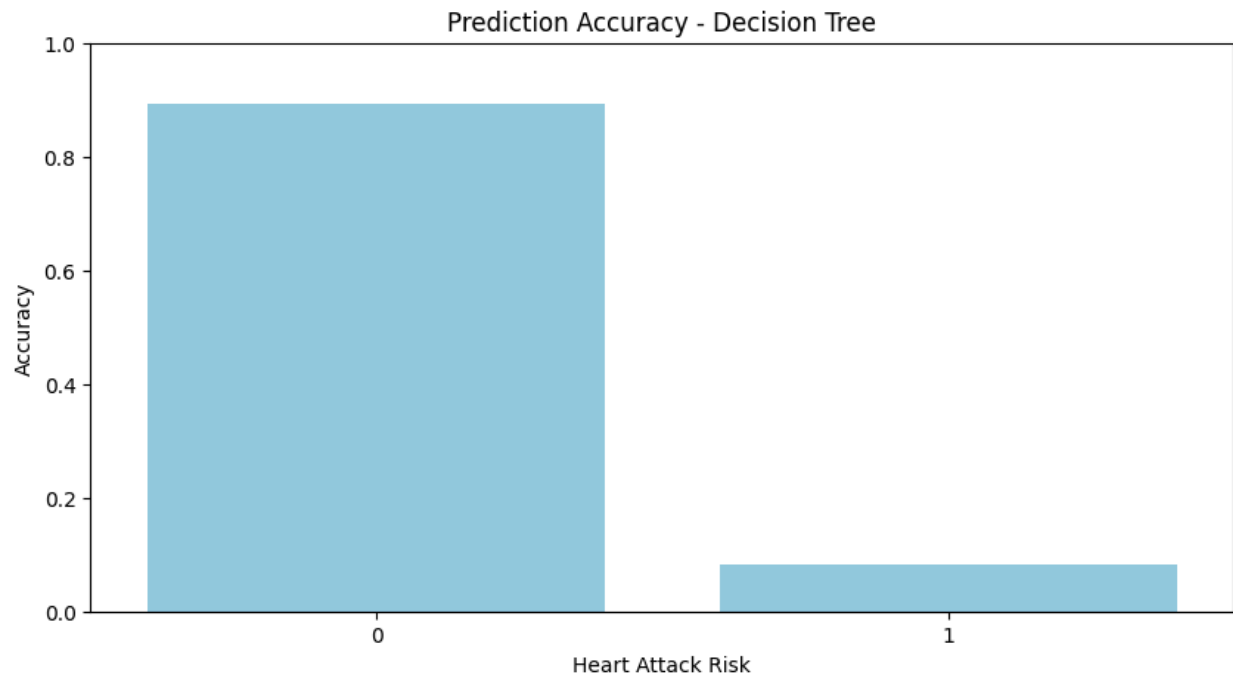
# Model Training and Testing

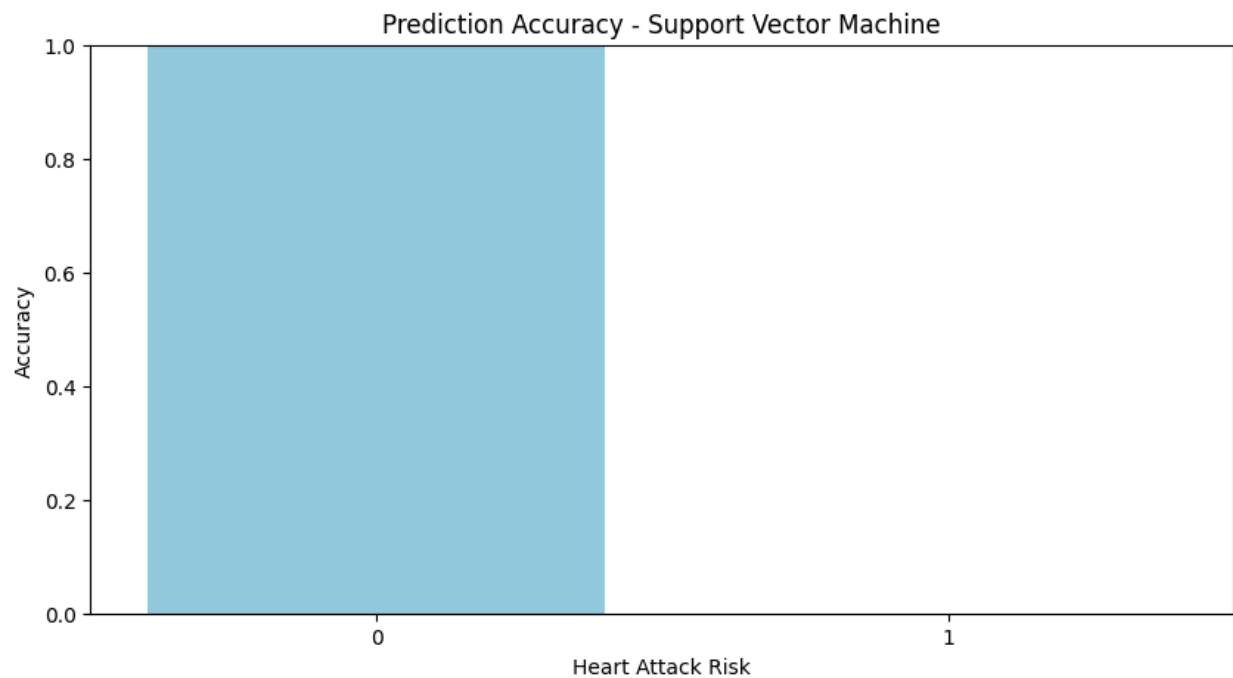I used the following three Machine Learning mopeds to train and test the dataset.

1. Decision Tree (RandomForestClassifier)
2. K-Nearest Neighbors (KNeighborsClassifier)
3. Support Vector Machine (SVC)
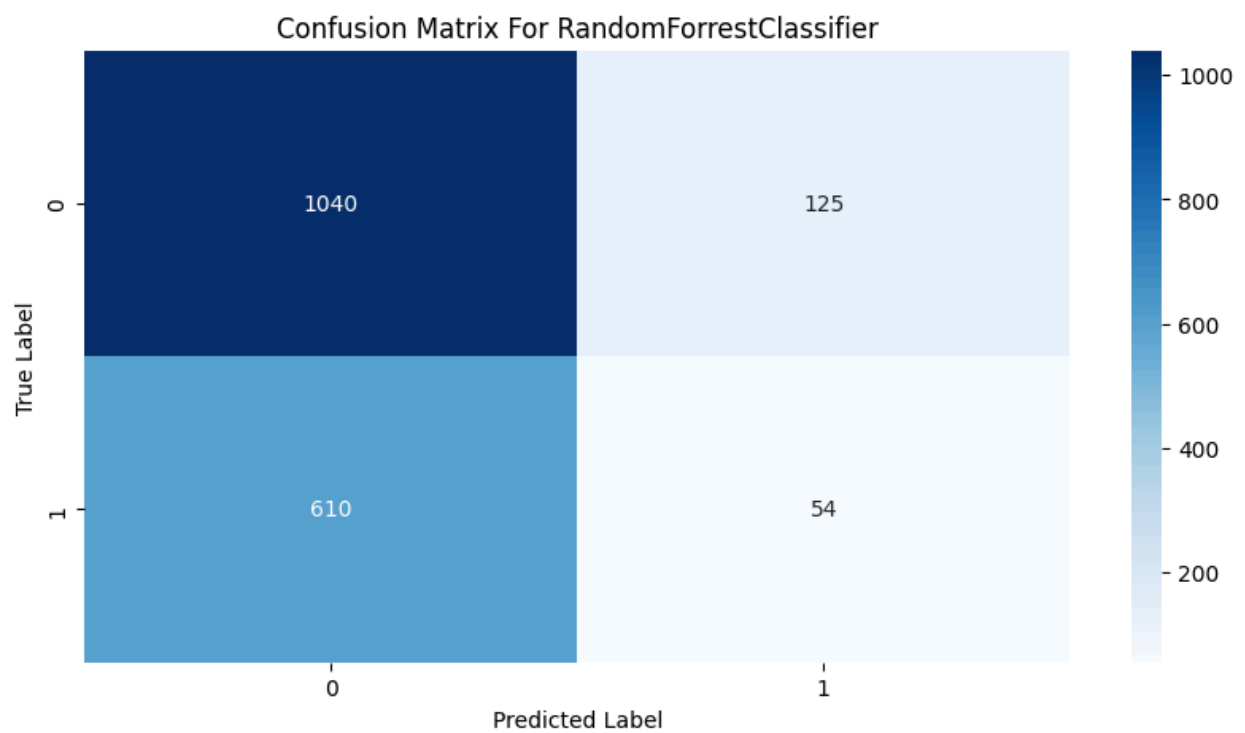
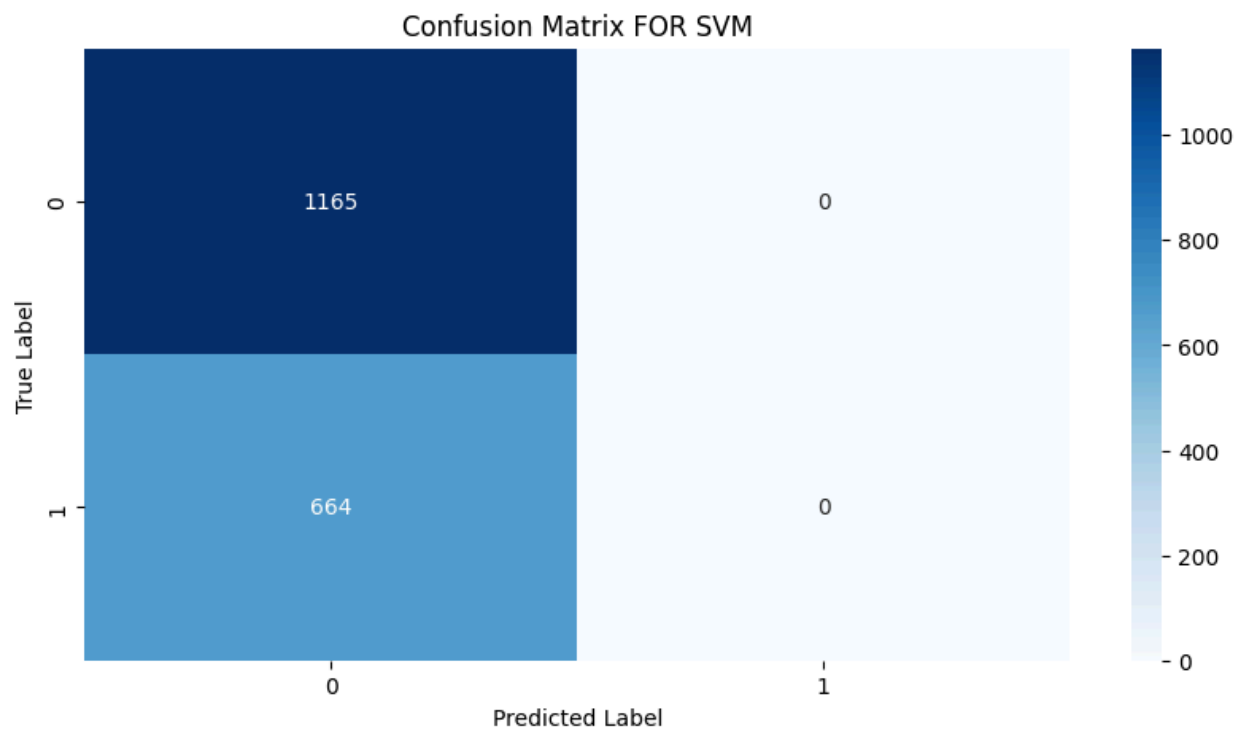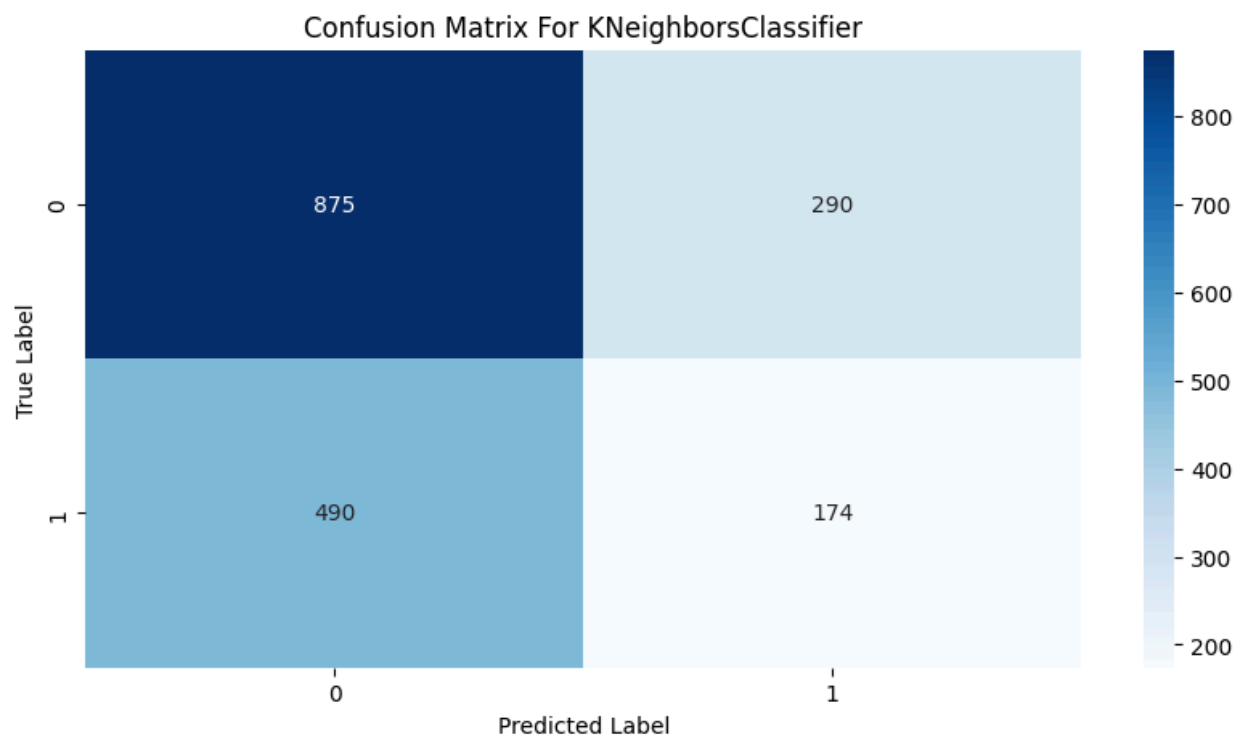# Model Selection / Comparison Analysis

**Bar chart showcasing prediction accuracy of all models:**



Prediction Accuracy - Decision Tree



Prediction Accuracy - K-Nearest Neighbors

Prediction Accuracy - Support Vector Machine

**Confusion Matrix:**



Confusion Matrix For RandomForrestClassifier

Confusion Matrix For KNeighborsClassifier



Confusion Matrix FOR SVM

# Conclusion

The analysis conducted for the project showed that a number of health indicators significantly affect the chance of having a heart attack. The project's goal is to analyze these data in order to gain important knowledge about possible risk factors that will enable early intervention and prevention. The "Heart Attack Risk Prediction" project is a big step in the right direction toward creating a useful tool for heart attack risk prediction. This initiative intends to help with early diagnosis of heart-related conditions and better healthcare outcomes by utilizing machine learning techniques.