

# Healthy Heart?

## Contents

I.	Project Description.....	2
1.	Project Title and Team Members:.....	2
2.	Goal and Objectives:.....	2
i.	Abstract: .....	2
ii.	Motivation: .....	2
iii.	Significance: .....	2
iv.	Objectives: .....	2
v.	Features:.....	3
II.	Increment - 1 .....	4
1.	Related Work (Background) .....	4
2.	Dataset.....	4
3.	Detail design of Features .....	5
4.	Analysis.....	7
i.	Exploratory Data Analysis : .....	7
ii.	Profile of the Data.....	8
iii.	Outlier detection - sys_bp, dia_bp : .....	8
iv.	IQR filtering - height, weight: .....	9
5.	Implementation .....	11
i.	Scaling the DATA: .....	12
ii.	Performing some standardizing using minmaxscaler .....	12
iii.	Data Preparation for training and testing : .....	13
iv.	Random Forest Classifier: .....	13
6.	Preliminary Results .....	13
7.	Project Management .....	14
i.	Implementation status report.....	14
8.	References/Bibliography References:.....	15

## I. PROJECT DESCRIPTION

### 1. PROJECT TITLE AND TEAM MEMBERS:

**Project Title:** Healthy Heart?

**Team Members:**

1. Neha Goud Baddam
2. Reshmi Chowdary Divi
3. Purandhara Maharshi Chidurala

**GitHub Link:** <https://github.com/nehabaddam/SDAI-Project.git>

### 2. GOAL AND OBJECTIVES:

#### i. **ABSTRACT:**

Artificial intelligence (AI) is taking over the world by providing its services in many ways by making the world a better place to live. We have been using artificial intelligence in everyday life for shopping online, streaming videos, smart houses, automation, agricultural fields, health industries, etc. In America, the leading cause of death is due to heart disease, this could be due to many factors like lifestyle, genetics, food habits, etc. It is important to spread awareness about cardiovascular diseases (CVD) and help people lead better life by predicting their risk of getting CVD. “Healthy Heart?” is a web application that will use AI to predict the risk of a person having a heart disease by taking a few major parameters that cause CVD as an input.

#### ii. **MOTIVATION:**

As the cause of death due to heart disease is increasing due to the sedentary lifestyle and unhealthy habits, it’s important to curb it by helping people self-access their health using AI, at their homes without having to visit any hospital. After knowing their risk of getting CVD, they can change their habits or visit a medical expert for more advice. This is the main purpose of the “Healthy Heart?”.

#### iii. **SIGNIFICANCE:**

There is much research going on in the field of medicine to identify the risk factors of CVD, but few factors have been identified that contribute to the CVDs like diabetes, High blood pressure, food habits, lifestyle changes, obesity, etc. Studies are showing that reducing these risk factors for heart disease can help in preventing heart disease. By analyzing these factors there is still a scope to calculate the risk factor and warn people about their health, which can help them have a healthy life ahead. “Healthy Heart?” plays a significant role by using previous research data and the latest data about the above-mentioned factors to accurately predict the risk factor using AI algorithms.

#### iv. **OBJECTIVES:**

The main objective of “Healthy Heart?” would be to take input (basic questions about the factors that cause heart disease) from the user and calculate the risk of having a heart disease based on the inputs given. Even if a person is healthy now they take this

assessment, they can find out their risk of having heart disease and improve their lifestyle to make themselves healthier.

#### v. FEATURES:

“Healthy Heart” focuses on predicting the risk factor accurately. Firstly, we train the model with the existing data and use some latest data for validation and testing. Firstly. We shall try to implement it using different models and check the model that produces maximum accuracy. We shall be using the model with the maximum accuracy for predicting the risk factor. Once the model is trained, we can host the website to get the input from the user dynamically and produce the risk factor as output.

We are focusing on using the following classification models:

- Random Forest

**Input:** The webpage will be hosted that consists of many input variables like height, weight, age, diabetes(yes/no), hypertension(yes/no), etc. We are hoping to consider below inputs:

Attribute	Description
age	Age (int)
height	Height (int)
weight	Weight (float)
gender	Gender (categorical code)
sys_bp	Systolic blood pressure (int)
dia_bp	Diastolic blood pressure (int)
cholesterol	Cholesterol (1: normal, 2: above normal, 3: well above normal)
glucose	Glucose (1: normal, 2: above normal, 3: well above normal)
smoke	Smoking (binary)
alco	Alcohol intake (binary)
active	Physical activity (binary)
CVD	Presence or absence of cardiovascular disease (binary)

**Output:** A calculated risk factor.

## II. INCREMENT - 1

### 1. RELATED WORK (BACKGROUND)

Cardiovascular diseases (CVDs) are the leading cause of death globally. According to projections, 17.9 million deaths globally in 2019—or 32% of all fatalities—were caused by CVDs. 85% of these deaths were caused by heart attacks and strokes. Low- and middle-income countries have a high rate of CVD deaths.

The main behavioral risk factors for heart disease and stroke include poor eating habits, inactivity, cigarette use, and alcohol abuse. Many people experience elevated blood pressure, elevated blood glucose, and elevated blood lipids, as well as overweight and obesity due to these factors. These "intermediate risk factors" indicate an increased risk of consequences like heart attack, stroke, and heart failure and can be assessed in primary care settings.

It has been established that lowering the risk of cardiovascular disease entails giving up smoking, reducing salt intake, increasing fruit and vegetable consumption, exercising frequently, and avoiding dangerous alcohol usage. Health policies that support environments where healthy options are both affordable and accessible are essential if people are to acquire and maintain healthy behaviors.

By identifying those who are most vulnerable to CVDs and ensuring they receive the right care, premature deaths can be prevented. Access to noncommunicable disease drugs and core health technology in all primary healthcare facilities is essential to ensuring that people in need receive care and counseling.

We are grateful to the WHO for providing the information we needed for our background investigation.[2]

### 2. DATASET

Cleveland residents were screened between 2000 and 2006 to assess CVD risk, identify high-risk individuals, understand the risk factors associated with CVD, and create a cohort for follow-up. Most of the methods used in the survey were questionnaires, along with physical measures, lab tests, and other techniques. Every two years, blood was drawn, and socioeconomic information was gathered once a year during follow-up.

The 76 attributes in this database have been reduced to the top 14 for consideration.

To this day, ML researchers in particular use the Cleveland database. The "target" field alludes to the patient's having heart illness. It is an integer value between zero and one (presence). Investigations into the Cleveland. [1]

Selected Attributes:

Attribute	Description
age	Age (int)
height	Height (int)
weight	Weight (float)
gender	Gender (categorical code)
sys_bp	Systolic blood pressure (int)
dia_bp	Diastolic blood pressure (int)
cholesterol	Cholesterol (1: normal, 2: above normal, 3: well above normal)
glucose	Glucose (1: normal, 2: above normal, 3: well above normal)
smoke	Smoking (binary)
alco	Alcohol intake (binary)
active	Physical activity (binary)
CVD	Presence or absence of cardiovascular disease (binary)

The selected 14 attributes cover Quantitative data, Categorical data, and Binary data.

### 3. DETAIL DESIGN OF FEATURES

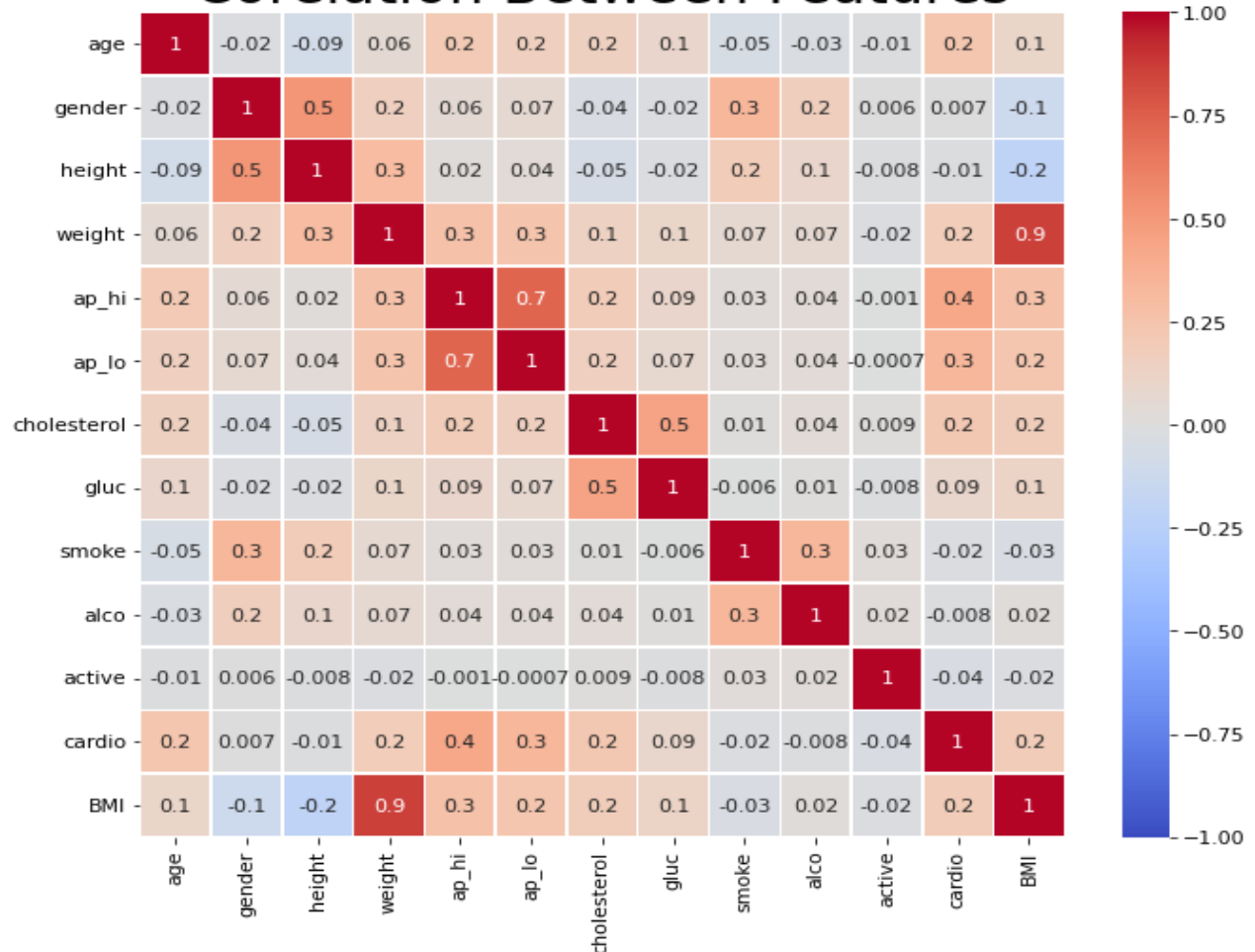
Here in this project, we are trying to implement a Random Forest algorithm. As tree splits are unaffected by scaling, the Random Forest Classifier is unaffected.

With cross-validation, hyperparameter tuning was carried out. A compromise between consistency and runtime, the number of trees produced stable results at around 300 trees, and 500 trees were ultimately chosen.

Although more tests (for tree size and parameters) were run than are currently in the code, they were omitted for readability and runtime considerations.

With the help of the random forest algorithm, the current study forecasts a patient's risk of developing heart disease. If the risk factor indicates 1, it is concluded that the patient is likely to develop CVD. If the risk factor indicates 0, it is concluded that the patient has a healthy heart. 14 attributes totaling 69000 data samples were collected to predict heart disease. Eighty percent of the dataset was used for training, and twenty percent was used for testing.

## Corelation Between Features

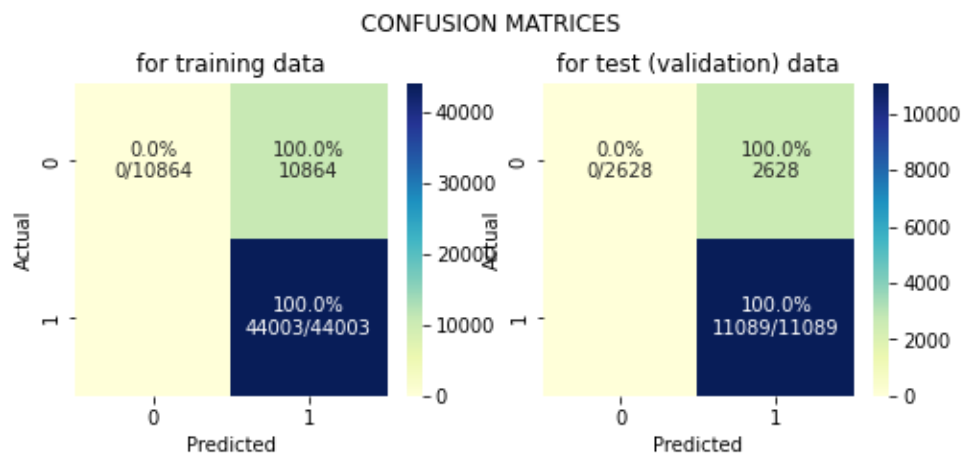


	id	age	gender	height	weight	sys_bp	dia_bp	cholesterol	glucose	smoke	alco	active	target
0	0	18393	2	168	62.0	110	80	1	1	0	0	1	0
1	1	20228	1	156	85.0	140	90	3	1	0	0	1	1
2	2	18857	1	165	64.0	130	70	3	1	0	0	0	1
3	3	17623	2	169	82.0	150	100	1	1	0	0	1	1
4	4	17474	1	156	56.0	100	60	1	1	0	0	0	0
5	8	21914	1	151	67.0	120	80	2	2	0	0	0	0
6	9	22113	1	157	93.0	130	80	3	1	0	0	1	0
7	12	22584	2	178	95.0	130	90	3	3	0	0	1	1
8	13	17668	1	158	71.0	110	70	1	1	0	0	1	0
9	14	19834	1	164	68.0	110	60	1	1	0	0	0	0

We obtained the correlation matrix after performing exploratory data analysis, which correlates the attributes of the data set.

On the testing data set, we are using the random forest algorithm to generate a confusion matrix.

More sophisticated metrics, such as sensitivity, specificity, and AUC, are obtained from the confusion matrix and can aid in our decision-making during the classification process.



## 4. ANALYSIS

### i. EXPLORATORY DATA ANALYSIS :

Performing the Analysis on the raw data that is exported from the data source. We have observed that there are few duplicate rows detected in the analysis, which is around 25000. Finally, there are no duplicates.

```
[5] # printing the length of the dataframe
length = df.shape[0]*df.shape[1]
print('Length of the data frame : ', length)

# printing missing values
missing_vals = df.isna().sum().sum()
print('Missing values: ', missing_vals, '\n')

# checking for duplicates
df_dup = df.duplicated().sum()
if df_dup:
    print('Duplicates Rows : {}'.format(df_dup))
else:
    print('No duplicates')
```

```
Length of the data frame : 909974
Missing values: 0
```

```
No duplicates
```

## ii. PROFILE OF THE DATA

Observing that the age is being considered concerning the number of days. That must be changed into years by dividing the value by 365. Systolic blood pressure "sys\_bp" and Diastolic blood pressure "dia\_bp" cannot be negative

	age	gender	height	weight	sys_bp	dia_bp	cholesterol	glucose	smoke	alco	active	target
count	69998.000000	69998.000000	69998.000000	69998.000000	69998.000000	69998.000000	69998.000000	69998.000000	69998.000000	69998.000000	69998.000000	69998.000000
mean	19468.892454	1.349567	164.359296	74.205639	128.817109	96.630747	1.366839	1.226464	0.088131	0.053773	0.803737	0.499700
std	2467.272520	0.476837	8.210060	14.395955	154.013595	188.475211	0.680228	0.572277	0.283487	0.225571	0.397172	0.500003
min	10798.000000	1.000000	55.000000	10.000000	-150.000000	-70.000000	1.000000	1.000000	0.000000	0.000000	0.000000	0.000000
25%	17664.000000	1.000000	159.000000	65.000000	120.000000	80.000000	1.000000	1.000000	0.000000	0.000000	1.000000	0.000000
50%	19703.000000	1.000000	165.000000	72.000000	120.000000	80.000000	1.000000	1.000000	0.000000	0.000000	1.000000	0.000000
75%	21327.000000	2.000000	170.000000	82.000000	140.000000	90.000000	2.000000	1.000000	0.000000	0.000000	1.000000	1.000000
max	23713.000000	2.000000	250.000000	200.000000	16020.000000	11000.000000	3.000000	3.000000	1.000000	1.000000	1.000000	1.000000

If sys\_bp and dia\_bp are more than 180 and 120 mmHg respectively then it is a hypertensive crisis, which is an emergency case. Therefore, max values are not realistic.

## iii. OUTLIER DETECTION - SYS BP, DIA BP :

Outlier Detection in systoli blood pressure and diastolic blood pressure columns

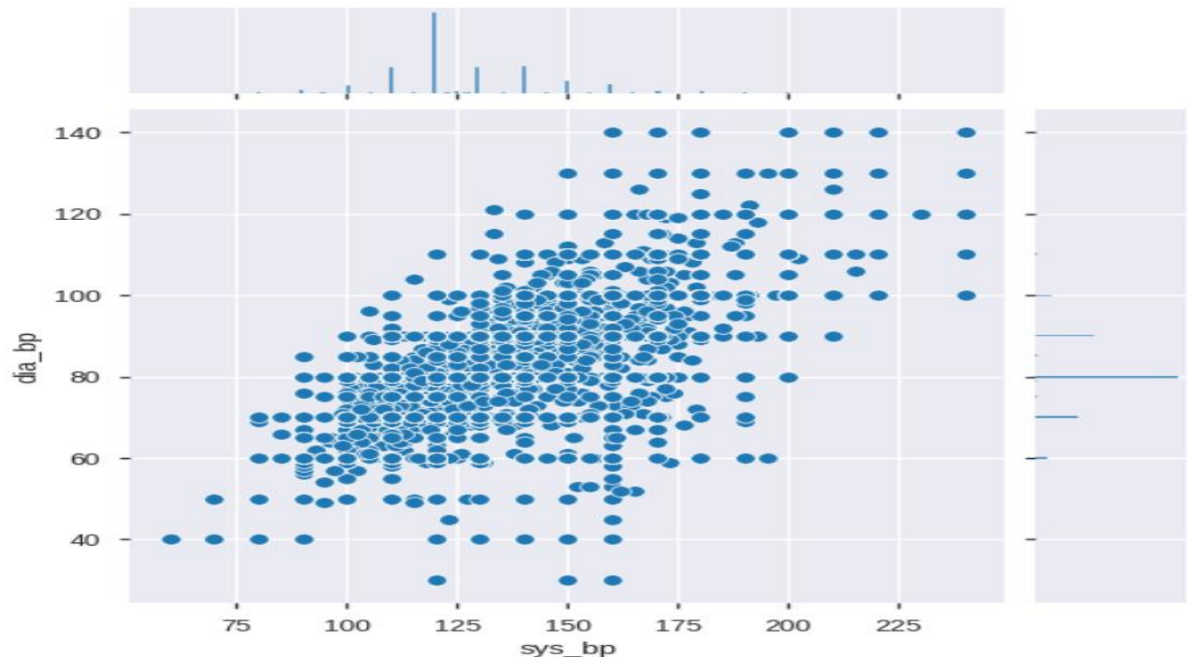
```
[ ] # Determining the outliers
outliers = len(df[(df["sys_bp"]>=280) | (df["dia_bp"]>=220) | (df["dia_bp"] < 0) | (df["sys_bp"] < 0) | (df["sys_bp"]<df["dia_bp"])])

print(f'total {outliers} outliers')
print(f'percent missing: {round(outliers/len(df)*100,1)}%')

total 1275 outliers
percent missing: 1.8%
```

According to the graph below, it can be assumed that values for sys bp and dia bp that are greater than 280 mm Hg and 120 mm Hg, respectively, will be eliminated as outliers. They cannot have negative values, so they also have positive values.





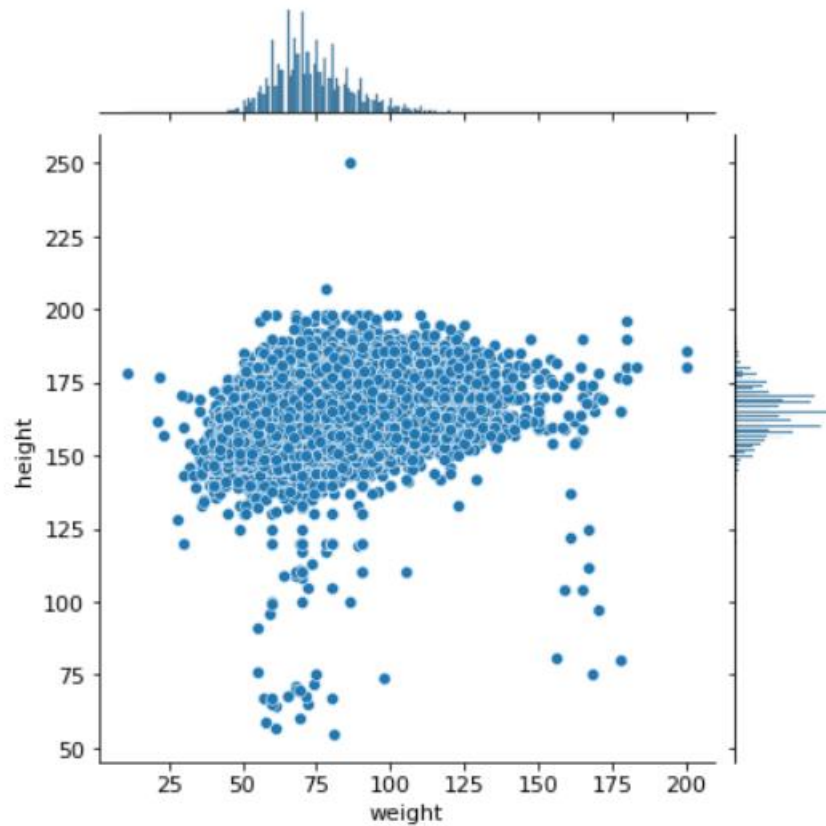
Detected Outliers

**Note:**

- Pulse pressure can not be negative which is the difference between the Systolic and diastolic blood pressures.
- Systolic blood pressure and Diastolic blood pressure cannot be negative
- Values larger than 180 mm Hg and 120 mm Hg for ap\_hi and ap\_lo respectively are an hypertensive crisis, which is in an emergency case

1. **Systolic blood pressure** is the pressure when the heart beats – while the heart muscle is contracting (squeezing) and pumping oxygen-rich blood into the blood vessels.
2. **Diastolic blood pressure** is the pressure on the blood vessels when the heart muscle relaxes.

**iv. IQR FILTERING - HEIGHT, WEIGHT:**



From the plot, it is observed that the smallest person has been recorded at 54cm and the tallest at 251cm

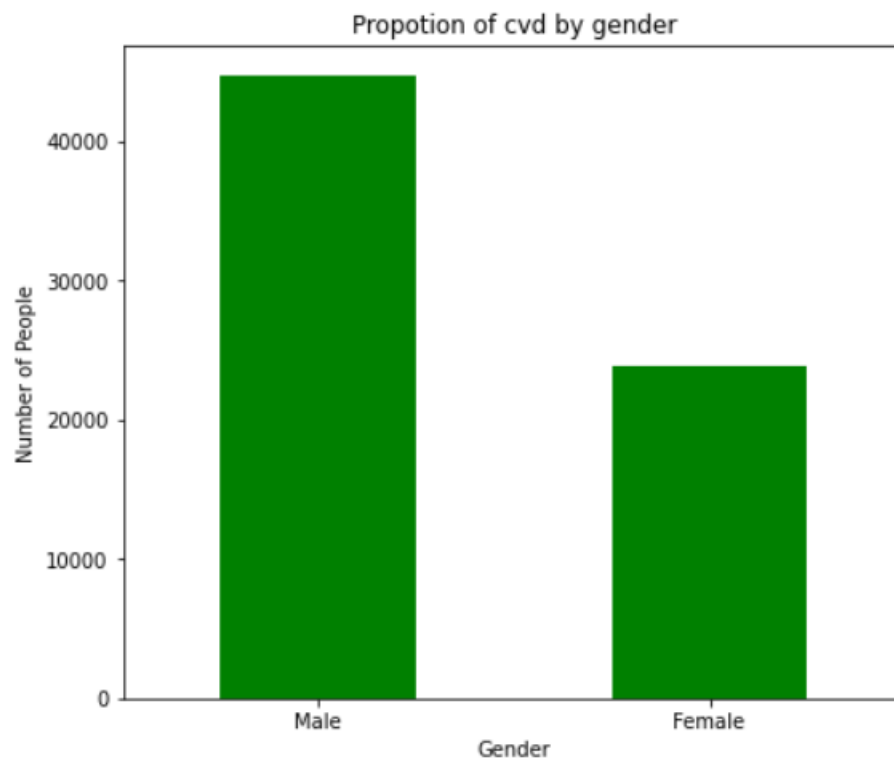
Feature: height  
Percentiles: 5th=152.0, 95th=178.0, IQR=26.0  
Identified outliers: 46

Feature: weight  
Percentiles: 5th=55.0, 95th=100.0, IQR=45.0  
Identified outliers: 20

After cleaning the outliers, each person's BMI is calculated

	age	gender	height	weight	sys_bp	dia_bp	cholesterol	glucose	smoke	alco	active	target	BMI
0	0.138060	2	168	62.0	110	80	1	1	0	0	1	0	22.0
1	0.151833	1	156	85.0	140	90	3	1	0	0	1	1	34.9
2	0.141543	1	165	64.0	130	70	3	1	0	0	0	1	23.5
3	0.132280	2	169	82.0	150	100	1	1	0	0	1	1	28.7
4	0.131162	1	156	56.0	100	60	1	1	0	0	0	0	23.0

Distribution of CVD by gender:



## 5. IMPLEMENTATION

After acknowledging the advantages and the accuracy of the Random Forest Classifier, we are implementing this as our ML model.

## i. SCALING THE DATA:

We have several features with various scales. Here, the data will be transformed using the StandardScaler library so that the mean is 0 and the standard deviation is 1. It essentially organizes the data.

```
[ ] #we perform some Standardization
df_scaled=df_cleaned.copy()

columns_to_scale = ['age', 'weight', 'sys_bp', 'dia_bp','cholesterol','gender','BMI','height']

scaler = StandardScaler()
df_scaled[columns_to_scale] = scaler.fit_transform(df_cleaned[columns_to_scale])

df_scaled.head()
```

	age	gender	height	weight	sys_bp	dia_bp	cholesterol	glucose	smoke	alco	active	target	BMI
0	-0.434234	1.36696	0.453577	-0.853013	-0.999463	-0.137129	-0.537112	1	0	0	1	0	-1.047946
1	0.309309	-0.73155	-1.063547	0.772519	0.799812	0.923668	2.409079	1	0	0	1	1	1.437583
2	-0.246221	-0.73155	0.074296	-0.711663	0.200053	-1.197926	2.409079	1	0	0	0	1	-0.758931
3	-0.746238	1.36696	0.580004	0.560493	1.399570	1.984465	-0.537112	1	0	0	1	1	0.242988
4	-0.806613	-0.73155	-1.063547	-1.277065	-1.599221	-2.258723	-0.537112	1	0	0	0	0	-0.855269

## ii. PERFORMING SOME STANDARDIZING USING MINMAXSCALER

```
#we perform some Standardization using minmaxscaler
df_scaled_mm=df_cleaned.copy()

columns_to_scale_mm = ['age', 'weight', 'sys_bp', 'dia_bp','cholesterol','gender','BMI','height']

mmScaler = MinMaxScaler()
df_scaled_mm[columns_to_scale_mm] = mmScaler.fit_transform(df_cleaned[columns_to_scale_mm])

df_scaled_mm.head()
```

	age	gender	height	weight	sys_bp	dia_bp	cholesterol	glucose	smoke	alco	active	target	BMI
0	0.588076	1.0	0.585106	0.246377	0.277778	0.454545	0.0	1	0	0	1	0	0.231557
1	0.730159	0.0	0.457447	0.413043	0.444444	0.545455	1.0	1	0	0	1	1	0.495902
2	0.624003	0.0	0.553191	0.260870	0.388889	0.363636	1.0	1	0	0	0	1	0.262295
3	0.528455	1.0	0.595745	0.391304	0.500000	0.636364	0.0	1	0	0	1	1	0.368852
4	0.516918	0.0	0.457447	0.202899	0.222222	0.272727	0.0	1	0	0	0	0	0.252049

### iii. DATA PREPARATION FOR TRAINING AND TESTING :

Using `train_test_split` we are dividing the data accordingly, using different parameters like `random_state`, and `test_size`. Here, in this dataset, we are considering the target column to be predicted by the model. In this case, we'll use the 80:20 split ratio. Thus, 80% of the dataset is used for the

```
x_train shape is (54865, 12)
x_test shape is (13717, 12)
y_train shape is (54865,)
y_test shape is (13717,)
```

### iv. RANDOM FOREST CLASSIFIER:

By Importing the `sklearn` library we are importing a random forest classifier.

A random forest is a meta-estimator that fits several decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree

Using a pre-defined function for this object called `GridSearchCV`, which is useful to find the best fit combination out of a given, we are performing the model train against it.

Here are the parameters that are given to find the best.

```
RandomForestClassifier(max_depth=20, max_features='sqrt', max_leaf_nodes=400,
                        n_estimators=500)
```

## 6. PRELIMINARY RESULTS

Here is an outcome that `GridSearch` has thrown for the best estimator. These Parameter values will be used to train the actual model to get the best output.

```

Accuracy: 0.7333236130349202
True Positive Rate: 0.6963490650044524
Report:
           precision    recall  f1-score   support

      0       0.72       0.77       0.75       6979
      1       0.74       0.70       0.72       6738

 accuracy          0.73          0.73          0.73       13717
  macro avg       0.73       0.73       0.73       13717
 weighted avg     0.73       0.73       0.73       13717

True Positive : 4692
True Negative: 5367

```

After performing the training with the outcome, the above picture shows the results. We can find the best output would occur at max-leaf nodes = 400

## 7. PROJECT MANAGEMENT

### i. IMPLEMENTATION STATUS REPORT

#### Work completed:

**Description:** We have completed gathering data for training and testing purposes. We have cleaned the data by removing all missing and duplicate values. We have detected outliers and removed them. We have also standardized the cleansed data. We have used the standardized data to train the model using different nodes and select the random forest classifier with a node that generates the maximum accuracy. Once trained, the test data is used for predictions. All necessary model metrics have been measured.

#### Responsibility (Task, Person):

- Neha Goud Baddam: Worked on model generation and model metrics calculation. Also worked on documentation.
- Reshmi Chowdary Divi: Worked on data gathering and data cleaning.
- Purandhara Maharshi Chidurala: Worked on outlier detection and removal.

#### Contributions (members/percentage):

- Neha Goud Baddam: 40%
- Reshmi Chowdary Divi: 30%
- Purandhara Maharshi Chidurala: 30%

**Work to be completed:**

**Description:** We shall be using the already-created model in Increment-1 to create a website, that will take input from the users. We also need to host the website using Apache tomcat/Xampp. Once done with coding, we also need to test the application to predict accurate risk factors.

**Responsibility (Task, Person):**

- a. Neha Goud Baddam: Designing a web application using a flask that will take user inputs (all 14 attributes) and generate risk factor as output.
- b. Reshmi Chowdary Divi: Testing the web application.
- c. Purandhara Maharshi Chidurala: Documentation and hosting of the website.

**Issues/Concerns:** There are no specific issues as of now, the issues could arise in Increment-2, while testing the web application. Also, the accuracy of the model is 73%, we may try to increase the accuracy to some extent.

**8. REFERENCES/BIBLIOGRAPHY REFERENCES:**

1. **Cleveland Database:** Dua, D. and Graff, C. (2019). UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: the University of California, School of Information and Computer Science.
2. **Cardiovascular Diseases – WHO:** cardiovascular diseases World Health Organization. World Health Organization. Available at: <https://www.who.int/health-topics/cardiovascular-diseases/> (Accessed: November 20, 2022).
3. **Heart Disease Prediction using Artificial Intelligence:** Zaibunnisa L. H. Malik, Momin Fatema, Nikam Pooja, Gawandar Ankita, 2021, Heart Disease Prediction using Artificial Intelligence, INTERNATIONAL JOURNAL OF ENGINEERING RESEARCH & TECHNOLOGY (IJERT) NREST – 2021 (Volume 09 – Issue 04)
4. **Mayo Clinic Health System:** Created by Mayo Foundation for Medical Education and Research using content from Framingham Heart Study Cardiovascular Disease 10-Year BMI-Based Risk Score Calculator, Framingham Heart Study General Cardiovascular Disease 30-Year Lipid-Based and BMI-Based Calculators, and ACC/AHA Pooled Cohort Equations CV Risk Calculator.