



# An artificial intelligence model for heart disease detection using machine learning algorithms

Victor Chang<sup>a,\*</sup>, Vallabhanent Rupa Bhavani<sup>b</sup>, Ariel Qianwen Xu<sup>b</sup>, MA Hossain<sup>c</sup>

<sup>a</sup> Department of Operations and Information Management, Aston Business School, Aston University, Birmingham, UK

<sup>b</sup> Cybersecurity, Information Systems and AI Research Group, School of Computing, Engineering and Digital Technologies, Teesside University, Middlesbrough, UK

<sup>c</sup> Vice President Office, Cambodia University of Technology and Science, Phnom Penh, Cambodia

## ARTICLE INFO

### Keywords:

Artificial intelligence  
Heart disease detection system  
Machine learning  
Predictive analytics  
Random forest classifier algorithm

## ABSTRACT

The paper focuses on the construction of an artificial intelligence-based heart disease detection system using machine learning algorithms. We show how machine learning can help predict whether a person will develop heart disease. In this paper, a python-based application is developed for healthcare research as it is more reliable and helps track and establish different types of health monitoring applications. We present data processing that entails working with categorical variables and conversion of categorical columns. We describe the main phases of application developments: collecting databases, performing logistic regression, and evaluating the dataset's attributes. A random forest classifier algorithm is developed to identify heart diseases with higher accuracy. Data analysis is needed for this application, which is considered significant according to its approximately 83% accuracy rate over training data. We then discuss the random forest classifier algorithm, including the experiments and the results, which provide better accuracies for research diagnoses. We conclude the paper with objectives, limitations and research contributions.

## 1. Introduction

### 1.1. Introduction

Heart diseases are often used in exchange for cardiovascular diseases. These kinds of diseases mainly refer to the conditions of blocked or narrowed blood vessels, resulting in a stroke, chest pain or angina, and heart attack. Other kinds of heart conditions, such as those affecting the rhythm, valve, or muscle of the heart, are other types of heart diseases. On the other hand, machine learning is crucial for determining whether anyone has suffered from heart disease. In either case, if these are predicted ahead of time, doctors would have a much easier time gaining crucial information for treating and diagnosing patients. Heart disease is mainly an incorrect symptom of coronary artery disease. It is also known as a cardiac disease; therefore, it is not with cardiovascular disease, which is any blood vessel disease.

Python is a programming language with a high level of object-oriented abstraction with a spirited, energetic collection of building options and quick development cycles. As per Loku et al. [1] analysis, it is regarded as one of the safest programming languages with numerous applications in the medical field. Furthermore, it is regarded as a well-liked and well-accepted programming language with applications traversing over AI-based software developments and several other web

applications. As per the suggestion of Mathur [2], the python framework is used easily for creating a desktop or web-based application. As per the depiction of Guleria and Sood [3], with the application of python programming in the health care sectors, especially for detecting heart diseases, clinicians and institutions can provide better and improvised outcomes for the patients through scalable and dynamic applications. However, the coding packages and libraries used in this project are Pandas, Matplotlib, IPython, Numpy, Python, SciPy, and many others.

### 1.2. Problem statement

Currently, the health care sector is generating information from several facilities and patients. By applying the best usage of this data, doctors can easily anticipate superior methods for treatment and enhance the complete delivery system of the health care sectors [4]. One of the most important uses is that the python framework can help make sense and encourage computational facilities in extracting valuable insights from the information over the health care sectors. Moreover, Python is considered to be one of the most renowned programming languages all around the globe. 32% of the UK individuals considered this programming language a secured language for developing healthcare

\* Corresponding author.

E-mail addresses: [victorchang.research@gmail.com](mailto:victorchang.research@gmail.com) (V. Chang), [rupabhavani22@gmail.com](mailto:rupabhavani22@gmail.com) (V.R. Bhavani), [qianwen.ariel.xu@gmail.com](mailto:qianwen.ariel.xu@gmail.com) (A.Q. Xu), [alamgir@camtech.edu.kh](mailto:alamgir@camtech.edu.kh) (MA Hossain).

<https://doi.org/10.1016/j.health.2022.100016>

Received 12 September 2021; Received in revised form 14 November 2021; Accepted 2 January 2022

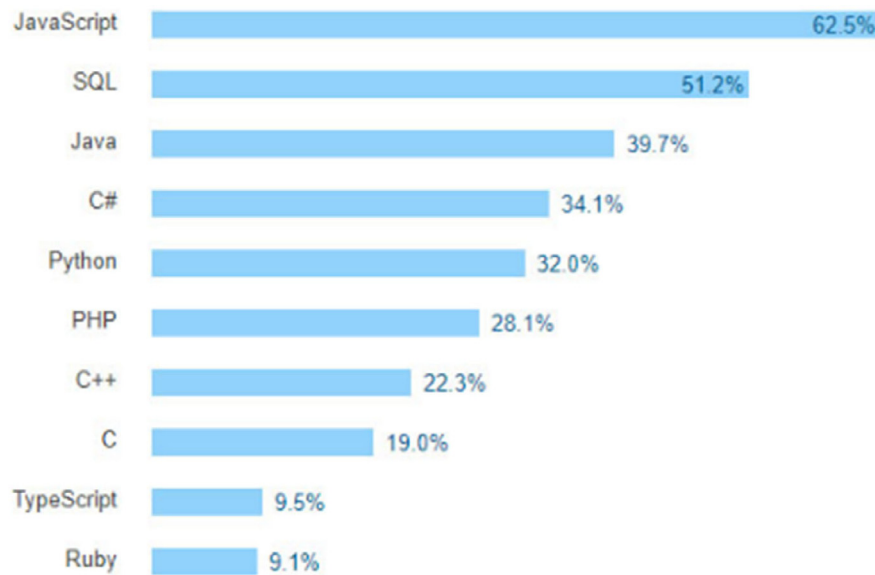


Fig. 1. Percentage of people preferring Python programming language for healthcare applications.  
Source: [5].

applications [5], see Fig. 1. High levels of LDL cholesterol, or “bad” cholesterol, can cause the most common form of heart disease, coronary artery disease (CAD). It is a plaque that has developed up in the arteries of the patient’s heart. CAD has no symptoms in its early stages. Patients can experience symptoms, such as chest pain, shortness of breath, and fatigue when plaque grows large enough to obstruct blood flow.

Additionally, the health care projects made using the Python language must deal with HIPPA (Health Insurance Portability and Accountability Act) requirements for dealing with healthcare records. In this context, as per Nithya and Ilango [6] depiction, Python supports computer security, as it has built-in tools that provide software-defined security. However, according to McPadden et al. [7], Python is currently used in the health care field for data science and machine learning applications that improve patient outcomes. As per the opinion of Panesar [8], the algorithms of machine learning encourage healthcare analytics to use Python, as developers can easily establish tracking and health monitoring applications. Thus, in this case, also, python programming is used for detecting heart disease.

### 1.3. Aims and objectives

#### 1.3.1. Research aim

The research aims are to detect heart disease using the python programming language.

#### 1.3.2. Research objectives

The objectives of the study are as follows:

To critically analyze the ways python language is used to detect heart disease.

To critically investigate the previous activities and apply a suitable methodological approach for superscribing the identified problem.

To critically apply data interpretation strategies in python language for health problem detection.

To critically assess the artifact or product with the help of cybersecurity approaches using appropriate methods and identifying the limitations and strengths of the work.

### 1.4. Research questions

The research questions are —

How would Python language help in detecting heart diseases among the patients?

How can the previous activities be critically investigated for applying appropriate methodological approaches towards addressing the identified issues?

How can the strategies for data interpretation be applied and can the findings be interpreted for achieving rational and logical arguments?

How will the product or the artifact’s insistence assist in evaluating third parties with the assistance of appropriate methods?

### 1.5. Research hypothesis

**H<sub>0</sub>: Python has wide applications in detecting the heart diseases**

**H<sub>1</sub>: Python does not have wide applications in detecting the heart diseases**

### 1.6. Sound justification of evidence

One of the most common diseases is heart disease and the most important reason for death in both developed and developing countries. Davenport and Kalakota’s [9] review looked at several research findings, including the use of the Python programming language for detection and prediction mechanisms for cardiovascular disease. The Python programming language is being used in disease detection systems, especially for heart diseases, to improve other healthcare-related systems.

## 2. Literature review

### 2.1. Introduction

The project comprises of detecting the presence of heart diseases using Python. The dataset comprised several factors, such as Chol, treetops, sex, age, and others. Several other import libraries, such as matplotlib, Numpy, Pandas, warnings, and many others, were used for the project. Correlation matrix, histogram, support vector classifier, K Neighbors Classifier, Random Forest Classifier, and Decision Tree Classifier were used for assessing the outcomes of the specified dataset using a python programming language. Additionally, Python is also considered an open-source language that encourages developing innovative solutions for the health care sectors and supplies better outcomes for the patients, resulting in enhanced care delivery. However, the

language also complies with the HIPAA checklist for assuring the safety of medical information. The major causes of heart disease are diabetes, obesity, unhealthy diet, overweight, excessive alcohol use, and physical inactivity. Therefore, heart disease includes arrhythmia that is considered as atherosclerosis is the hardening of the arteries caused by a heart rhythm abnormality. During a heart attack, some people experience these symptoms. Additionally, pain that spreads to the arm, dizziness or light headedness, throat, snoring, and sweating can occur. Heart attacks, strokes, and coronary heart disease, also known as heart failure and coronary artery disease, are much more common in people over 65 than in younger people.

## 2.2. Demonstration of a deep understanding of an area of an individual interest associated with specialized computing in the health care sectors

One of the most well-known machine learning algorithms tasks is the classification of data. Machine learning tends to be an essential function in this case for extracting knowledge from business activity datasets and transferring it to larger databases. The majority of the machine learning methods rely on a huge number of features that explain the algorithm's behavior, resulting in the model's complexity, indirectly or directly [10]. Many algorithms such as hybrid methods are used in conjunction with logistic regression, naive Bayes, K-nearest neighbor, and neural networks to integrate the heart disease diagnostic algorithms mentioned earlier. Thus, in this case, the system was trained and implemented over the python platform with the help of the UCI (Unique Client Identifier) machine learning deported benchmark dataset.

Coronary artery disease, arrhythmias (heart rhythm problems), heart abnormalities (such as congenital heart defects), and a variety of other disorders are included in the category of heart diseases. Cardiomyopathy and heart infections are among the conditions that fall under this category. The most common measure of heart risk is chest pain, which is a symptom of cardiovascular disease. After that, it has symptoms of Nausea, Indigestion, Heartburn, or Stomach Pain. The paper will exhibit how a program can be created in Python to analyze whether or not an individual is suffering from cardiovascular disease or not [11]. In this paper, the system uses a dataset comprising fourteen characteristics of the test outcomes, carried on around 100 persons. However, the patient suffering from heart disease symptoms will be diagnosed using binary digits, 1 and 0, where 1 will indicate the true value (The patient has heart disease, in other words.) and 0 will indicate the false value (that is, the patient does not have any kind of heart disease). Additionally, co-relation and trends of the obtained features will also be recognized with the help of several features, such as gender, age, cp (chest pain type), chol (cholesterol level), FBS (fasting blood sugar level), exang (exercise-induced angina), thalach (maximum achieved heart rate), old peak (ST depression persuaded by exercise respective to rest), thal (maximum achieved heart rate), ca (number of major vessels).

In this project, initially, the libraries will be imported. Then, the dataset will be loaded, and it will be stored within a variable for printing the information. Finally, the dataset will be imported and the data will be processed. However, after analyzing the outcomes, it is seen that the K-neighbor classifier algorithm showed an 87% score, whereas the support-vector, decision tree, and random forest classifier displayed 83%, 79%, and 84%, respectively [13]. See Fig. 2.

Contrarily, in this case, a co-relation matrix will be used for evaluating the connections within several types of variables. A positive correlation exists between the predictor and the chest pain variable, indicating that the amount of chest pain is directly proportional to the probabilities of suffering from heart diseases. In this case, chest pain is considered a statistical feature with four values: value 1, value 2, value 3, and value 4, referring to atypical angina, typical angina, asymptomatic and non-anginal pain, respectively [14]. A negative correlation among these variables would indicate that more amount of blood is required by the heart.

## 2.3. Development of an approach for addressing the significant research areas or practices over specialized computing areas in health care sectors

However, a major benefit of Python within the health care sector is that it assists in making sense of the information by working with Machine Learning and AI within the healthcare sectors. As per the analysis of Ozgur et al. [15], the development services of Python is a suitable option for a strong and powerful language to encourage computational abilities in obtaining valuable insights from the information of the patients suffering from heart diseases, that will, in turn, help in supporting healthcare based applications. It is convenient in case one has to deliver the diversity of developing something with the help of an internet connection or has autonomously worked without any internet connection. As per the opinion of Srinath [16], the pliability of running over numbers of operating systems is compounded by a large district and a distinct syntax. Moreover, Python proved to be a suitable language for evaluating huge datasets, with the help of machine learning algorithms in receiving significant insights [17]. The language is also favored by data scientists due to the availability of extensive libraries, such as SciPy, Pandas, Numpy, and many others.

## 2.4. Demonstration of the capability to evaluate, synthesize, and search the information's from the appropriate sources in health care sectors

In this project, the information was gathered from outside databases and a logistic regression was performed during Python. As per the analysis of Jiang et al. [18], several pieces of information are also used for determining the attributes of datasets. For instance, induced angina for the exercise, maximum heart rate, resting blood pressure, resting electrocardiographic measurements, fasting sugar level, thalassemia level, induced depression, number of major vessels, and many others were used for representing the datasets comprising several values. However, the sex of a person can be evaluated using two values, either 0 and 1, where 0 indicates female and 1 refers to male. Contrarily, the chest pain categories will be evaluated with the help of four values, 0, 1, 2 and 3, indicating asymptomatic condition, atypical angina, non-anginal pain, and typical angina, respectively. However, a confusion matrix is also used for generating false positive and negative outcomes. Moreover, as opined by van den Burg et al. [19], the details for the regression analysis are obtained from adequate CSV files. On the other hand, the classification scores for detecting heart disease can also be obtained. In contrast, help vector classifiers, decision tree classifiers, random forest classifiers, and a variety of other machine learning algorithms are only a few examples. However, in this case, the data wrangling procedures will also be used for determining the relation between the negative and positive binary predictor. As per the depiction of Holdgraf [20], this self-service data wrangling equipment helps deal with more complicated information rapidly and generates accurate outcomes to reach superior decisions.

Additionally, the features are also compared with positive and negative heart patients. From investigating all the information, it has been found that the positive patients experienced increased heart rates and displayed around one-third of the ST depression's amount persuaded by exercise associated with old peak [14]. Thus, developers can effectively use Python to build the required models in predicting heart diseases before they become severe.

## 2.5. Critical application of the cybersecurity techniques to ensure conformities with networking configurations and management system of information security within the healthcare sectors

Python programming language is mostly chosen to be used within the health care sectors as the cybersecurity professionals can accomplish the project efficiently. As shown in Fig. 3, as per the opinion of Calix et al. [21], the language is also used for decoding and sending

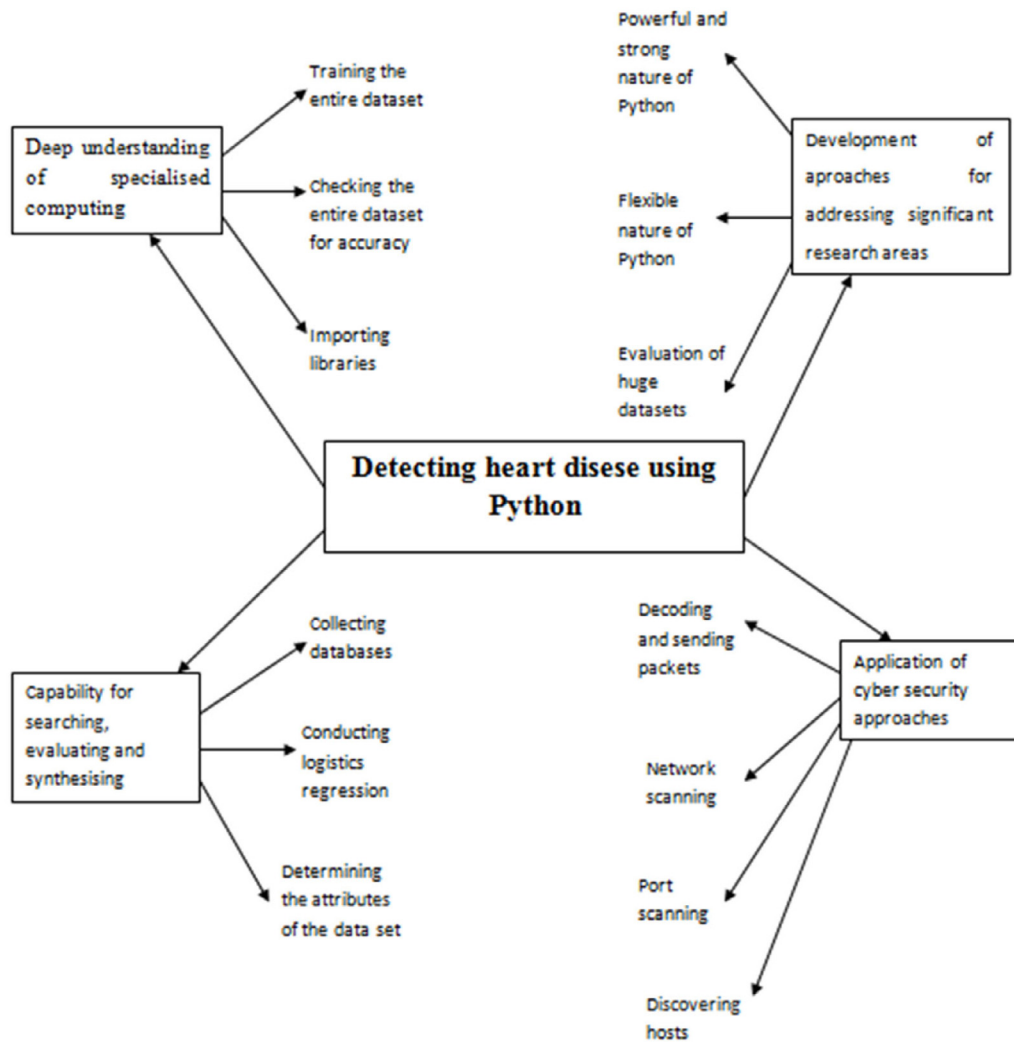


Fig. 2. Conceptual framework.  
Source: [12].

packets, network scanning, port scanning, accessing servers, discovering hosts, and analyzing malware. Additionally, it is also useful to conduct a stream of cybersecurity applications, such as malware analysis and scanning. Moreover, as the health care sector comprises a huge number of confidential information of the patients, this language is well suited for developing an application within this sector [22]. It has an executive and well-defined immaculate method. Moreover, the availability of a huge number of libraries also decreases the amount of effort required to conduct specific tasks, such as cyber threat analysis, detection, penetration testing. The language also has a simple syntax that can easily be picked up by the new developers entering this cybersecurity field [23]. Thus, this language is chosen for developing a system for detecting heart diseases.

## 2.6. Literature gap

The gap in the literature is the availability of minimal information associated with the creation of heart disease detection. Python is a programming language used by the system within the health care systems. As per the depiction of Bau et al. [24], python programming, being a concept-oriented programming language, numerous equipment is involved in developing a prophetic model. However, the selection of the appropriate technology or the appropriate tool can help in the development of a proficient model. Complexity is another issue that

is being faced. Additionally, there might also be a lack of the desired resources.

Python language is the most appropriate language for detecting patients suffering from heart diseases. It is one of the robust languages that foster computational capabilities in gaining valuable insights from the information of the patients suffering from heart diseases. Thus, it is apt for the health care sector. Moreover, it also complies with the HIPAA regulations that ensure medical information safety. However, the project will be using a database from external sources and the libraries will be imported. Loading of the datasets will occur following it, and it will be seized within a variable to reproduce the information. At last, the dataset will be imported with specific import libraries, followed by data processing.

## 3. Research methodology

### 3.1. Overview of methodology

Machine Learning (ML) is important in predicting the existence or absence of heart arrhythmia, locomotor disorders, heart diseases, and other conditions. It was expected well to provide significant insights to physicians, allowing them to adjust their diagnosis and care on a patient-by-patient basis. This project follows the *random forest algorithm* for developing heart disease detection, and this algorithm is used for the methodology to detect heart disease using Python.

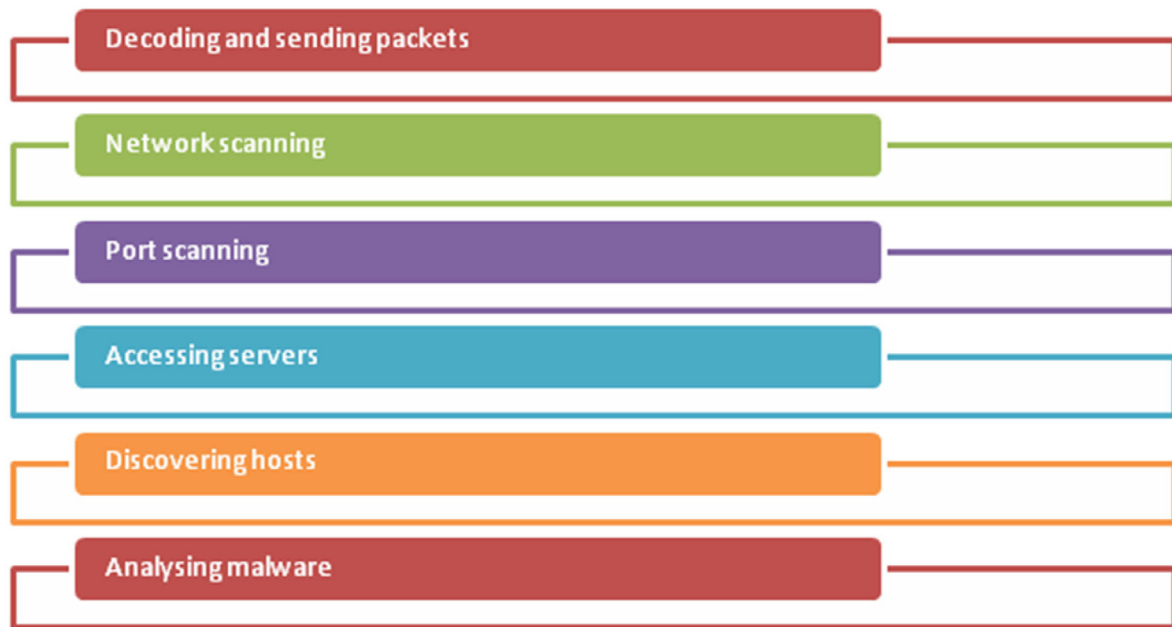


Fig. 3. Advantages of python programming language in cybersecurity scanning.  
Source: Influenced by [21].

A model regarding ML has played a significant role in creating accuracy and determining results with the aid of training data. This model is considered very significant due to its 83% (approximately) accuracy rate over training data [25]. This accuracy level has been significantly highlighted with the aid of a confusion matrix in terms of using accuracy store calculations. Based on the view of Iwendi et al. [26], through the implementation of the confusion matrix, the accuracy is printed as well as displayed on the screen. This model tends to exhibit an approximate 70% of accuracy rate in terms of evaluating test data. It is effective for receiving higher accuracy levels of scores in relation to existing ones. Based on this context, it can be highlighted that selection of “random forest” is implemented based on a particular dataset as well as a decision tree. A vote for the forecasted outcome can be found in this research. Furthermore, the final forecast is based on the highest number of votes that will be presented as the final result. See Fig. 4.

This algorithm establishes the robust as well as high accuracy of the data. This algorithm also allows the number of participative decision trees within the process. It does not lead to over-fitting issues. The reason behind that is that it considers the average of all predictions and cancels out the prejudices. This random forest algorithm is also useful for regression as well as classification issues [28]. Using Machine Learning algorithms, we can predict possible Heart Diseases in people. There are a variety of Machine Learning algorithms, including the K Neighbors Classifier, Support Vector Classifier, Decision Tree Classifier, and Random Forest Classifier. The Random Forest algorithm is a set of algorithm models demonstrating many decision trees using bootstrapping, random subsets of tools, and average voting to make predictions. For a classification issue, Random Forest provides the probability of belonging to the class.

### 3.2. Research method

Here, import all essential libraries used in the project, such as NumPy that works with arrays and the pandas libraries work with the CSV files and data frames. After that, matplotlib creates charts using pyplot. This library is to define parameters using rcParams as well as color them with cm.rainbow. After that, split the dataset into training as well as testing data.

Machine learning is the science of programming a system that can learn from different types of data, according to Larsen et al. [29].

Python is the most commonly used programming language for this type of project. It will also replace lots of languages in the industry and have a huge amount of collection of libraries such as Numpy, scipy, pandas, scikit-learn, matplotlib and many more. When doing exploratory data analysis, Pandas dataframe.info() function is useful for providing a concise overview of the data frame. After that, import the dataset and use read\_csv() that reads the dataset as well as will save it to the dataset variable. Pandas describe() is also used to display certain simple statistical explanations like percentile, mean, standard deviation, and many others. A set of numeric values if not part of a data frame. This method accepts a series of strings and returns a variety of results. After that, use a correlation matrix that understands the data. pyplot was used to display the xticks and yticks in the correlation matrix and the addition of names for the correlation matrix colorbar () displays the matrix's colorbar.

### 3.3. Use of algorithm with justification

This project is based on the random forest algorithm because this algorithm is considered as flexible as well as easy to use in machine learning. Even without hyper-parameter tuning, this algorithm achieves excellent results in the majority of cases. It is also one of the most widely used algorithms due to its versatility and simplicity. Random forests are considered a supervised learning algorithm, according to Amini et al. [30]. This algorithm is used to categorize trustworthy loan applicants, detect fraudulent behavior, and predict diseases.

The random forest algorithm belongs to the machine learning group, is a useful learning method under supervision. The ensemble learning theorem is its foundation, which states that it is a tool for solving a complex problem by integrating several classifiers and improving the model's accuracy. The Random Forest algorithm is a classification algorithm that uses a random forest to classify data which combines the results of several decision trees into a single result. The aim is to apply to different subsets of a dataset to increase the dataset's predictive accuracy. The random forest is formed in two steps: The first is to mix and match to make the random forest, you will need  $N$  decision trees in total, and the second step is to make predictions for each of the trees you made in the first step. The following steps and diagram can be used to illustrate the working process:



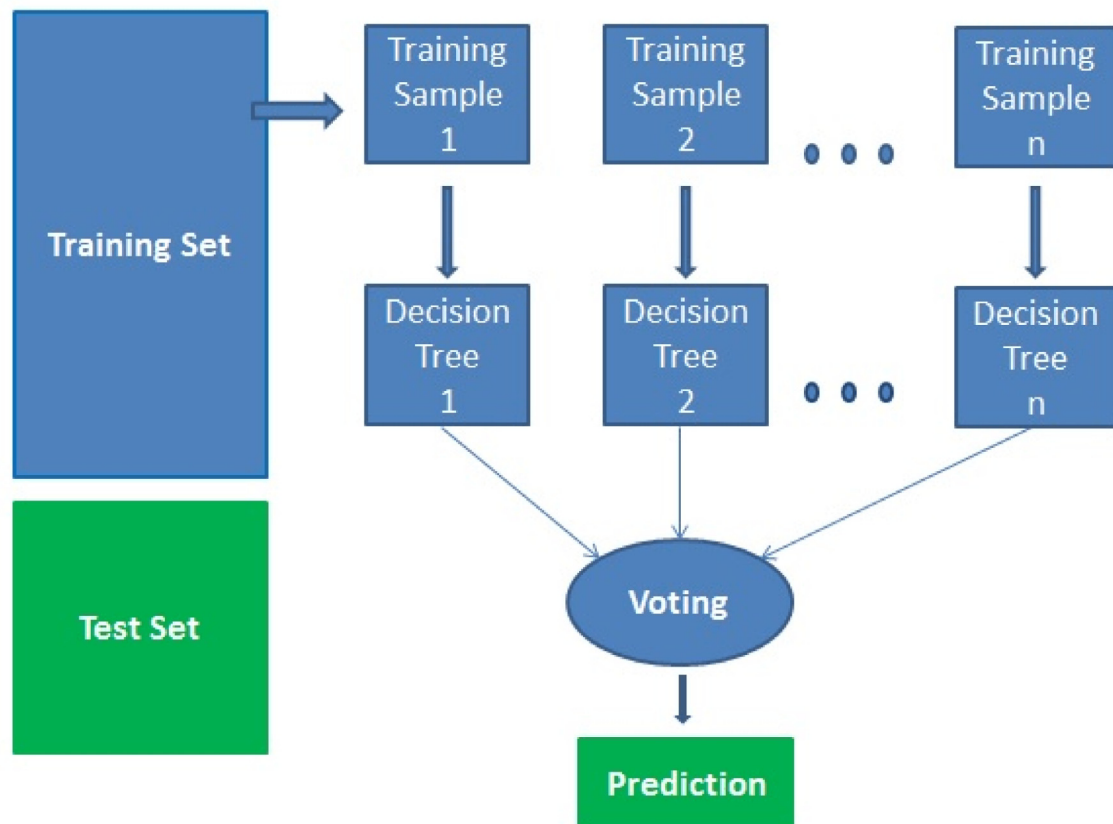


Fig. 4. Algorithm of random forest.

Source: Influenced by [27].

Step 1: Pick  $K$  data points at random from the training collection.

Step 2: Build decision trees for the data points you've chosen (Subsets).

Step 3: Choose the number  $N$  for the number of decision trees you want to build.

Step 4: Repetition of Steps 1 and 2.

Step 5: Find the predictions of each decision tree for new data points, and allocate the new data points to the group with the most votes.

The steps for implementation are listed below.

The first phase is data pre-processing. The training set is then equipped with the Random forest algorithm, and the test result is predicted. The training set should be fitted to the Random forest algorithm. To make it fit, we will use the RandomForestClassifier class from the sklearn.ensemble library. To visualize the training set results, we will plot a graph for the Random forest classifier. Last but not least, assess the accuracy of the results to build the uncertainty matrix and then visualize the test and training set outcomes.

After that, data is preparing for training and scaling. After data scaling training has been generated for random forest algorithm.

```
fromsklearn.ensemble import RandomForestRegressor
regressor = RandomForestRegressor(n_estimators=20, random_state=0)
```

```
regressor.fit(X_train, y_train)
y_pred = regressor.predict(X_test)
```

Muñoz et al. [31] also mentioned that the regression approach is used in the 'sklearn.ensemble' library's RandomForestRegressor class. Additionally, the RandomForestClassifier class of the 'sklearn.ensemble' library was classified. The 'n\_estimators' is a parameter in the RandomForestClassifier class. The number of trees in our random forest is determined by this parameter and the search information for all of the RandomForestClassifier parameters.

After that, classification faced challenges. The metrics used to evaluate an algorithm are accuracy, confusion matrix, precision recall, and F1 values.

```
fromsklearn.metrics import classification_report, confusion_matrix,
accuracy_score
print(confusion_matrix(y_test,y_pred))
print(classification_report(y_test,y_pred))
print(accuracy_score(y_test, y_pred))
```

The accuracy is achieved by our random forest classifier. To enhance the accuracy using parameters of the RandomForestClassifier class as well as to see improvement on the results.

### 3.4. Data analysis

Data analysis works with categorical variables and it will break particular categorical columns 1s and 0s are placed in dummy columns. Furthermore, column Gender has values of 1 for males and 0 for females, which can be divided into two columns, one with the value 1 for true and the other with the value 0 for false. The idea of decision trees is taken to the following level with this random forest classifier and building a forest of trees, each of which is made up of a random selection of features from the total features. Here is a list of how many trees will be used to estimate the class.

### 3.5. Data output

There are many classifiers; therefore, random forest classifiers give better accuracy to this project. For prediction, this project uses a variety of medical parameters such as age, sex, blood pressure, cholesterol, and obesity. Aside from that, the EHDPS predicts a patient's heart disease. Those patients develop heart disease because of their condition. This project work allows for substantial information, health causes, relationships linked to heart disease, and trends to be identified.

### 3.6. Use of the required software with justification

This project required the dataset and jupyter notebook editor that used Python's package manager, pip instead of anaconda. Jupyter notebook is a free, open-source, interactive web component, also known as a computational notebook, as explained by Reich et al. [32]. Software code, numerical performance, explanatory text, and multimedia tools can all be combined into a single document by the researchers. This app allows everyone to view and share documents containing images, live code, calculations, visualizations, narrative text, data cleaning and transformation, numerical simulation, mathematical modeling, data visualization, machine learning, and other features. Normally, Python programming allows building code for multi-purposes. Apart from this, the anaconda is considered most likely preferred for Data science projects. It will also provide pre-built libraries to help projects like machine learning get up and running quickly. Jupyter is a data science platform designed for Python-based data science applications. Besides, it would lower the barriers for data scientists, as Jupyter has simplified documentation, data visualization, and caching a lot easier. Apart from this, Mendez et al. [33] have described that Jupyter helps programmers write code and display the results in real-time without waiting for other parts of the program to finish. If it is a code that is training an ML model or a code that is downloading gigabytes of data from a remote server, Jupyter caches the output of every cell that is running. Jupyter Notebook is a platform-agnostic and language-agnostic programming environment. Jupyter can be interpreted in a variety of languages. Apart from this, Yin et al. [34] have narrated that Jupyter supports visualizations. It also includes rendering some of the datasets such as graphics and charts that are generated from codes with the help of modules such as Matplotlib, Plotly and many more. Sample documentation code is available in Jupyter notebook. The more Jupyter is used, the easier it is for developers to describe their codes line-by-line with suggestions attached all along the way. The better the Jupyter is used, the more interactivity and explanations can be added once the fully functional code.

## 4. Data findings and analysis

### 4.1. Introduction

Python is widely regarded as the most effective and useful programming language. It contains a number of libraries that are used in this machine learning project. A subset of the Artificial Intelligence model is this approach to Machine Learning. Python libraries are used to make predictions with SKLEARN, which is a machine learning prediction tool.

### 4.2. Critical analysis regarding the description of heart disease detection

The healthcare industry generates massive quantities of data, also known as big data. The most common cause of death is heart disease, worldwide and a significant public health issue. In medical science, the detection of heart disease in its early stages has become one of the most serious problems. According to Ramalingam et al. [35], the RR interval, QT interval, and QRS interval are some of the characteristics that are studied in heart disease detection. This method determines whether or not the patient is well. This method determines whether the patient is normal or not.

On the other hand, Subhadra and Vikas [36] have opined that heart disease is considered a fatal human disease that repeatedly enhances the world in both undeveloped and developed countries and increases with consequently and last of all it causes death. This system helps to classify a complex as well as large medical dataset along with that to detect heart disease. Along with that, it detects some steps that use a map-reduce algorithm to both detect the disease and reduce the size of the dataset.

This output represents a histogram of the dataset of heart disease detection. Here, use the code of `df.hist(figsize=(10,10))`. This histogram is used to interpret discrete data visually and to summarize it. By showing the number of data points that fall within a given range of values, necessitates focusing on the most significant points, or facts, of numerical data. See Fig. 5. Exploratory Data Analysis (EDA) finds relevant data, detects the mistake, checks assumptions, and determines the correlation between these explanatory variables. This EDA allows data analysis that excludes statistical models as well as interferences. Apart from this, Ambale-Venkatesh et al. [37] have narrated that the risk factors are because the random forest algorithm is used to consider and forecast heart disease, and the study is done using publicly accessible heart disease data. There are 304 records in this dataset, each with 14 attributes, including age, gender, and more. In order to predict heart disease, a random forest algorithm is used for data visualization as well as data analytics. Along with that, Dogan et al. [38] have described that this research paper discusses classification performances, pre-processing methods, and evaluation metrics. Furthermore, the result of the visualized data shows that the prediction is considered correct. The precision of this method was 83 percent. The ANN-based three-stage method for predicting heart disease. Heart disease prediction using deep neural networks as well as his proposed model performed well. Apart from this, it also produced good outcomes. See Fig. 6.

### 4.3. Data interpretation of the selected dataset

Data interpretation is considered a process that reviews data over some predefined processes. It helps to assign some data. It involves the outcomes of data analysis and makes inferences on the relations. Mahdavejad et al. [39] described that data interpretation also helps with data analysis, data collection, and data presentation in row form. A machine learning approach's interpretation is described as the process of attempting to comprehend a machine learning model's predictions. Python libraries are helping to build the data interpretation through a machine learning approach. To understand why the classifier chooses a specific class, interpret the models using the predictions as well as the parameters. The data analysis of machine learning models facilitates the method of attempting to comprehend a machine learning model's predictions. This involvement of heart disease detection has two significant phases: it has monitored the evaluation metric and tried various ideas of algorithms selection that enhance and develop the more robust approach. It is also essential to interpret the models using the parameters as well as predictions to understand the classification chose the exact class. On the other hand, Tauzin et al. [40] have declared that data interpretation is a part of data analysis in the modern past, cost-effective technologies, and it acts as an essential part of the health sector involves emergency situations as well as outbreaks of disease. Data interpretation by the term "machine learning" refers to a form of data processing. The construction of analytical models is automated using this data analysis. Data analysis is an artificial intelligence specialization focused and based on the assumption that computers can learn and recognize patterns in data. It also makes decisions with little to no human input. See Fig. 7.

### 4.4. Data interpretation strategies using python language in detecting heart problem

Stats models are considered a python model that helps the users explore data perform statistical tests and estimate statistical models. It has a go-to language for data analysts. According to Peters et al. [41], it depends on important data analysis libraries, data pre-processing, and, last of all, exploratory data analysis. Data analysis libraries contain libraries and packages, and those are open-source and widely used to crunch data.

### Fundamental Scientific Computing

Numpy is a numeric library using Python that can perform linear algebra, Fourier transforms, and a random number. Along with that,

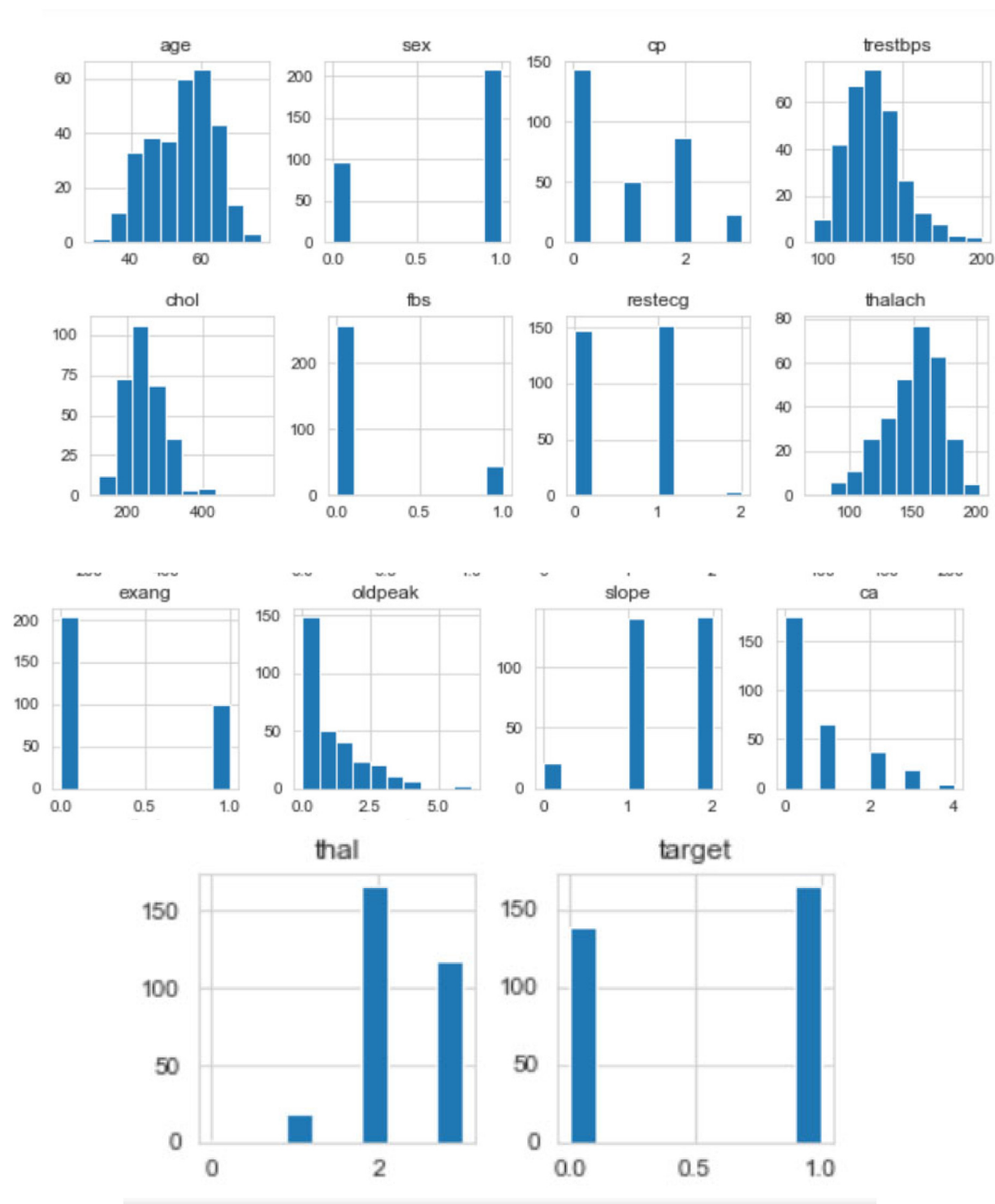


Fig. 5. Bar graph of each data.

SciPy is a scientific library of Python, and this library involves a high-level science module.

The criteria of machine learning approach interpretation processes as well as explore the techniques for interpretation dependent on the scope. This study is focused on an in-depth understanding of current model interpretation approaches and their drawbacks and challenges. It will also go through the age-old trade-off between model accuracy and model interpretability. Finally, consider some of the most popular model interpretation strategies.

#### 4.5. Data manipulation and data visualization

Pandas is used for data analysis as well as machine learning. Pandas are used in the form of data frames. The panda's library also provides data from various file formats, such as CSV, Excel, plain text, JSON, SQL and many more. Peters et al. [41] explained that data manipulation is considered the process of data transformation, formatting, and structuring. Apart from this, the Matplotlib library is used for plotting as well as visualizing data. This library allows plot graphs, Heatmaps, line plots,



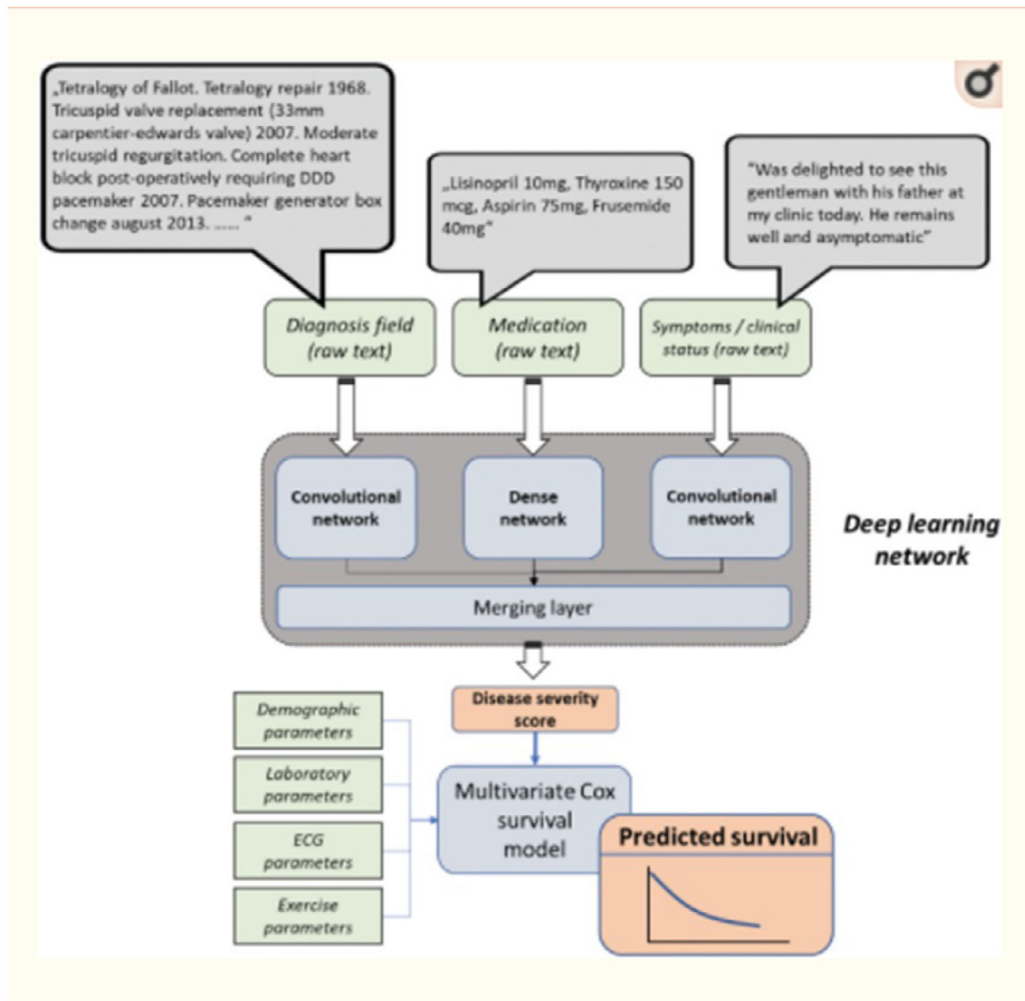


Fig. 6. Heart disease detection using a deep learning approach.  
Source: [38].

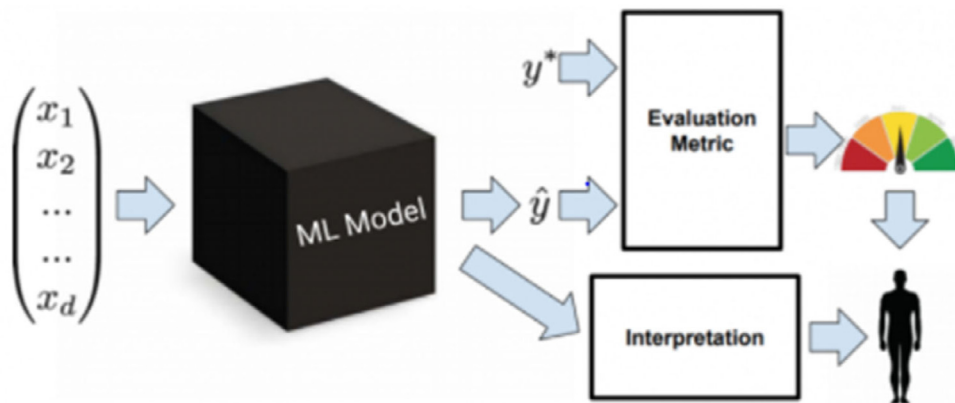


Fig. 7. Data interpretation.  
Source: [40].

histograms as well as a lot more. It is embedded in GUI toolkits. See Fig. 8.

This bar graph of the counterplot represents the target variable of the dataset. It means the sex of the patient who faced much heart disease. 0 represent female and 1 represents male of the patients. Here, used of the code is `sns.set_style('whitegrid')`

```
sns.countplot(x='target',data=df,palette='flare')
```

## Machine Learning

As Carleo et al. [42] explained, Scikit Learn of machine learning is considered a free machine learning library that helps to build on NumPy, SciPy, and Matplotlib. This library contains efficient components for statistical model development. It also runs different classification, regression, as well as clustering algorithms. It is also integrated well with pandas through working on data frames.

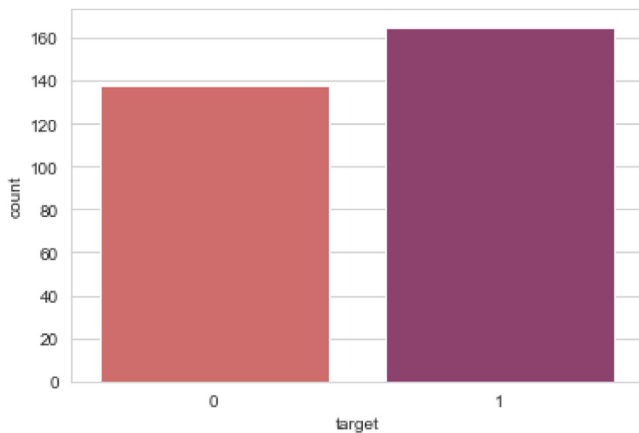


Fig. 8. Count Plot the sex.

### Importing Libraries and Loading the Data

This approach imports libraries as NumPy and pandas and then loads the dataset. The CSV is the most generally used format for that machine learning data is considered as presented. This CSV file is used for automatically assigning names. Otherwise, it is labeled if the file does not have a header, each column of the dataset manually names the attributes.

#### 4.6. Exploratory data analysis

The EDA data analysis is used to get a better understanding of data and look for the data. For statisticians, it is similar to a kind of storytelling. It allows for the discovery of trends and observations within data using visual methods. Aside from that, EDA is frequently used as the first step in the data modeling process. It will explore the dataset as well as perform the exploratory data analysis. It follows handling missing value, outlier treatment, encoding categorical variables and normalizing, and finally, scaling and removing duplicates.

#### 4.7. Evaluation of data

Evaluation of data introduces Python, and it is a mathematical data processing language that uses Pandas Data Frame objects to store data. Importing, washing, and converting data in preparation for review is the work involved in data analysis.

#### 4.8. Data cleaning

Data cleaning is considered as preparing data for analysis by cleaning the raw data that prepare the data visualizing data and predicting data. It is the method of removing false, corrupted, improperly formatted, and redundant data from a dataset that would otherwise be incorrect, corrupted, and incomplete.

### Correlation Matrix Plot

According to Harper et al. [43], a correlation Matrix Plot is a covariance matrix that is a metric called the correlation that defines the strength of the linear association. The Correlation matrix sums up the power and direction of a linear relationship between two variables, and it allows values between  $-1$  as well as  $+1$ . The feature of the correlation matrix displays the correlation between the coefficients. A particular random variable is considered correlated with each of its other values. This represents an excellent way to check correlations among features by visualizing the correlation matrix as a heat map. See Fig. 9.

The relationship between age, sex, cp, trestbps, chol, FBS, restecg, thalach, exang, oldpeak, slope, and ca is expressed as a graph. The linear relationship between two continuous variables is defined using the dataset's correlation.

#### 4.9. Assessment of artifact regarding cybersecurity approach

The cybersecurity consultants are deploying a revolution as users can transition from managing the perimeter that is extracting as well as analyzing any residue left by cyber thieves on every endpoint device such as laptop, desktop and many more. As narrated by Jia et al. [44], they focus on finding artifacts that convey every user as well as applications. These applications are ever interacted with by the system and find these artifacts deep in the OS system files, memory, file systems, as well as more systems. Artifacts can reveal evidence even when the perpetrators proclaim innocence. It can also show the cyber criminal's intent by showing their Internet searches and websites visited.

On the other hand, Galinec et al. [45] have narrated that concern about cybersecurity services and artifacts provides important clues about unauthorized access by unauthorized entities. Cybersecurity services include investigative activities, and when assessments are drawn, the artifacts help corroborate the findings. Mainly, the root cause of a cyber-attack is considered as never discovered, nor are the threat actors ever found.

#### 4.10. Different methods for cybersecurity

Computers, mobile devices, networks, servers, and data can be protected by cybersecurity systems, which are similar to security gatekeepers to ensure that the information on the computers and the data stored on them are not vulnerable to external threats. According to Fan et al. [46], cybersecurity has various security methods, such as cloud security based on the cloud. Because of its increased privacy, its storage has been a common choice over the last decade. Data is stored in the cloud, which is considered more reliable because it is protected by a software program that tracks activity and can warn users if anything unusual occurs with their cloud accounts. Following that, Network Protection protects an internal network from external threats by enhancing network security. Firewalls supervised internet access, anti-spyware, and antivirus programs are all part of network security. Apart from this, Jin et al. [47] have opined that application Security protects the data. This information is kept on the applications that are used to operate the business. Applications are more accessible across various networks that are particularly vulnerable to cyber-attacks, and the safety of applications with cybersecurity antivirus programs, firewalls, and encryption services is critical. It is also essential for the control access to carefully manage the physical access to that premises and the computer network. Control access restricts access to outside threats as well as unauthorized users. Apart from this, Anees and Hussain [48] have explained that it has limited access to the data. Otherwise, it also has limited access to the services or application controls. Control access is limited to sending as well as receiving certain kinds of email attachments. Hence, firewalls are effectively gatekeepers between the system as well as the internet. It has one of the most effective defenses against cyber threats, including malware and viruses. The firewall of the device properly checks them regularly and it ensures that it is up to date in terms of applications and firmware. Otherwise, it will not be completely successful. Apart from this, use for security software such as antivirus, anti-malware, and anti-spyware that help detect the malicious code after detection will be removed from it. Along with that, Monitor for intrusion that detects unusual network activity as well as monitor the system. If a monitoring system detects a possible security breach, it may issue a warning based on its discovered features, such as an email alert. Apart from that, raising knowledge has to assist in the security of the company and ensure that understanding of the user's role and any relevant policies and procedures that provide them with regular cybersecurity awareness and training. Besides this, many important security enhancements are included in updates to help defend against various threats. This is referred to as bugs or glitches, and it ensures that apps and systems are up to date to avoid being targeted by criminals. See Fig. 10.

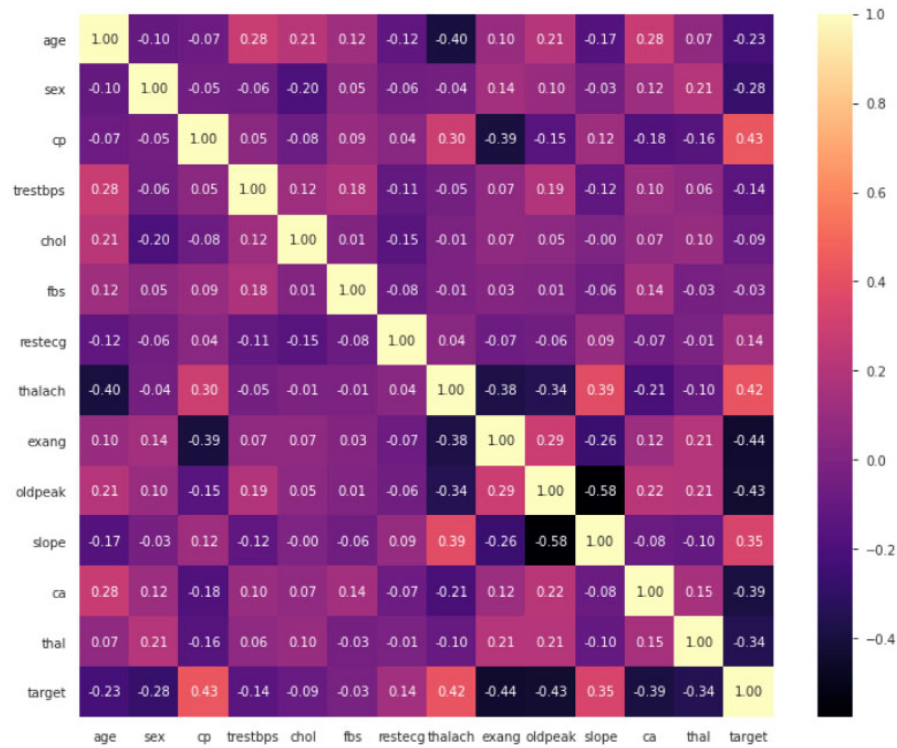


Fig. 9. Correlation of each feature in the dataset using the heatmap.

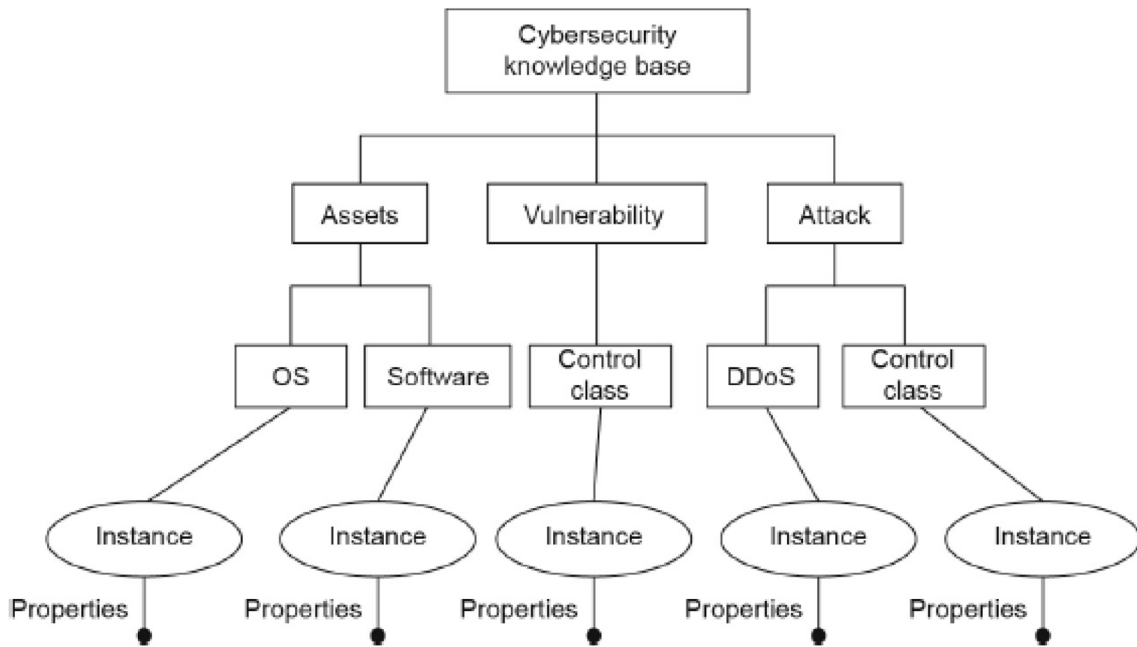


Fig. 10. The architecture of the cybersecurity knowledge base.  
Source: [47].

#### 4.11. Critical evaluation regarding the identified methods of cybersecurity

The traffic passing through cybersecurity is monitored by a firewall. This network allows information to flow in the form of packets. Apart from that, the firewall examines each of these packets separately for

any potentially dangerous attacks. If the firewall comes across these threats by chance, it will automatically block them. Firewalls also come with an access policy. Apart from this, Anees and Hussain [48] have narrated that it can support development for certain hosts as well as services. Following that, attackers exploit certain hosts, and the best

course of action is to prevent those hosts from gaining access to the device. If a user believes that they need protection from this type of unauthorized access, an access policy may be implemented. One of a user's most important characteristics is privacy. Hackers are constantly on the lookout for data privacy to gather knowledge about the user. As a result of the use of a firewall, various services are provided where the domain name service and the finger, which are used by a website, are blocked. On the other hand, Jin et al. [47] have described that the hackers have no chance of getting a privacy description. Furthermore, firewalls will prevent the site system's DNS information from being accessed. According to, the names and the IP address are not available to the attackers. Firewalls, on the other hand, require an expenditure that varies depending on the type of cost. Hardware firewalls are believed to be more costly than software firewalls. Hardware firewalls often necessitate installation as well as maintenance, which can be expensive.

Access control allows an effective path of the unwanted control entry of the logical as well as physical assets. On the other hand, access control is more damaging; it is caused when the key is compromised, which is the disadvantage.

Client security gives them peace of mind by making them aware of the security system. In addition, cybersecurity is becoming increasingly important, with the advantage of potentially increasing revenue and marketability. On the other hand, Security flaws in software are commonly caused by bugs. Many types of bugs can be found in software. Minor issues, such as incorrect print output rendering, are among them. Otherwise, it is an error message that has not been formatted properly. Some bugs are security flaws that could allow unauthorized access to sensitive information. As a result of these vulnerabilities, attackers will take advantage of security flaws. Apart from that, Fan et al. [46] have mentioned that the protection and safety of their confidential personal details is one of the users' main concerns, which makes them reluctant to share their information and, as a result, make online transactions.

On the other hand, authentication allows the process of ascertaining. User and session authentication in software that has been improperly configured is extremely vulnerable. Aside from that, it has security solutions to ensure that customers' information is secure in the users' system. The best business practices can increase the number of buyers, increase revenue, and create a positive consumer reputation. In the software development industry, security misconfiguration is a common issue.

#### 4.12. Role of python in detecting heart disease

This project predicts the heart disease of the patient by extracting the patient's medical history. It leads to fatal heart disease from a dataset. It contains the patient's medical history, such as sugar level, blood pressure, chest pain, and age and gender. Based on the given scenario, Python provides an accuracy of the heart disease detection and achieves 83% accuracy using the random forest classifier. This classifier shows a collection of decision trees drawn randomly from the training set described by Mehmood et al. [49]. It combines the votes of different decision trees are used to determine the final class of the test object. It works by building a large number of decision trees during training and then generating the class that represents the mode of the classes; otherwise, the average of heart disease detection prediction mean is calculated. Therefore, there is not much point in enhancing this number of estimators to enhance the accuracy further [Refers to appendix].

#### 4.13. Conclusion

The other library used in this prediction is SEABORN, which links all of the attributes together. Last but not least, the confusion matrix

determines accuracy flawlessly after importing CONFUSION MATRIX. The random forest algorithm creates and merges many decision trees based on the given scenario to generate a more accurate and stable prediction. The Random forest classifier is a bagging classifier with hyperparameters, much like a decision tree.

## 5. Discussion

### 5.1. Introduction

In this section, discuss the method that is developed using Python. This language allows lots of packages and libraries that are dependent on "machine learning algorithms". These "machine learning algorithms" as well as their outcomes are compared with the proposed model.

### 5.2. K-Neighbors classifier

The k-Neighbors classifier recognizes a vector based on the plurality vote of its neighbors. The classification route of this classifier is very important. According to Tjahjadi and Ramli [50], this method is used based on the neighbor's majority vote. Along with that, the k-neighbor is given a weight of age. It is much more distant than the others. This distance between the vector and the neighbor is represented by d, and the weight assigned is constituted by 1/d. Apart from this; For instance-based learning, this classification technique is used. According to, the classification of heart disease detection is when the estimation also occurs.

On the other hand, Pedrozo et al. [51] have opined that the K-Nearest Neighbor is considered a machine learning algorithm that can solve issues of classification and regression. The K-Nearest Classifier is used to determine the data's accuracy. As a result, it can determine which model should be the most suitable among the following models for usage in the future. In Eq. (1), the Euclidean distance is symbolized as

$$d(p, q) = d(q, p) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2} \quad (1)$$

### 5.3. Decision tree classifier

The decision tree classifier creates a tree from data item observations. This data item is made up of branches that connect the branches, which are made up of target values that make up the leaves. Assegie and Nair [52] narrated that the decision tree is used to classify heart disease detection after that. The leaves show the class labels and branches, which show the features that are the leading features of the labels for the classes. It is used to identify data objects using visual components. The decision tree model's accuracy value is calculated by inputting the XTrain parameters as well as the YTrain, which is the fit shape The Decision tree model's score is then discovered, bypassing the XTest and YTest parameters to the system score() function, which searches for the Decision tree model's score. Apart from this, Herbold [53] has stated that the decision tree algorithm is considered as a flowchart such as a tree structure. It has an internal mode that displays the outcome and the branches that constitute a decision rule. Along with that, every leaf of the node includes the outcome. Apart from this, the top mode node is known as the root node of the decision tree algorithm. See Fig. 11.

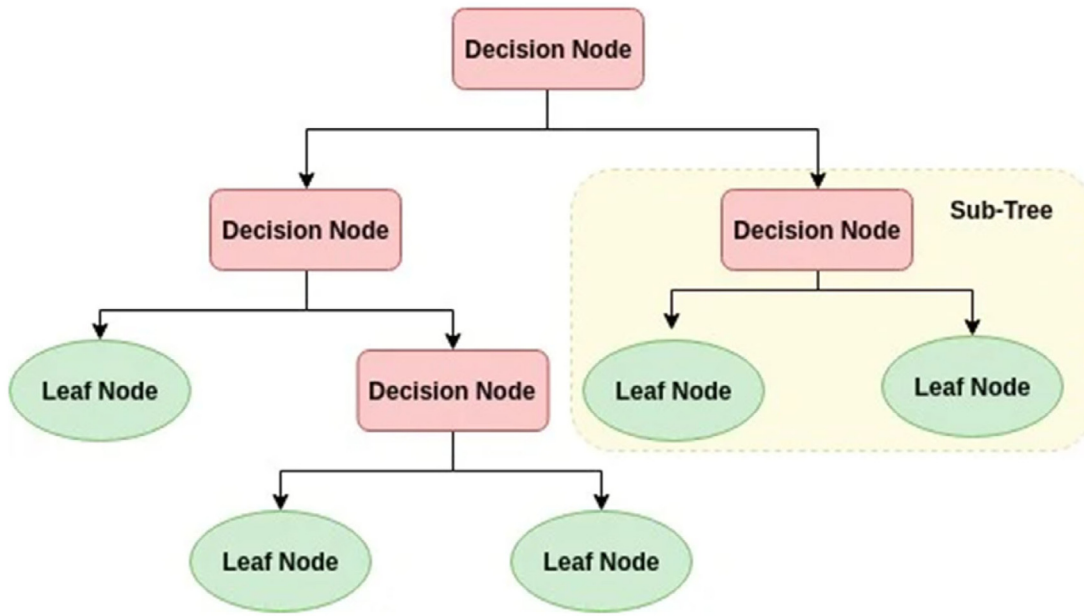


Fig. 11. Decision tree classifier.  
Source: [53].

#### 5.4. Support vector machines

A Support Vector Machine (SVM) algorithm generates a model that classifies new examples into one of several categories. Aside from that, it is divided into two groups. The SVM model can be referred to as a “non-probabilistic classifier”, since it can be thought of as a space-based characterization of instances. A hyper-plane is considered a building that divides the various examples into various categories. SVM supervises the machine learning algorithm used for both classifications and regression problems, as Zahariev et al. [54] define. The Support Vector simplifies the coordinates of separate observations. After that, the SVM classifier is considered as a frontier that best segregates the two classes. The performance is called the SVM model’s level of accuracy that inputs the parameters of the XTrain as well as the YTrain that is the method of fit, according to Wang and Liang [55]. After that, the model’s score is input into the system of the score using the parameters XTest and YTest. The score of the SVM model can be found using this tool.

#### 5.5. Random forest classifier

Random forest classification is considered an ensemble learning approach used for solving machine learning challenges such as classification and regression. This random forest classification of heart disease detection algorithm works by constructing multiple decision trees. Apart from this, this classifier uses a technique called “bootstrap aggregation”.

The likelihood of a reduction in node impurity is weighted, and it can be improved by hitting the node to determine feature significance. With only two child nodes, Scikit-learn calculates a node’s importance using Gini Importance. Between Eqs. (2) and (6), we show key formulas for random forest classifiers.

$$ni_j = w_j C_j - w_{\text{left}(j)} C_{\text{left}(j)} - w_{\text{right}(j)} C_{\text{right}(j)} \quad (2)$$

ni sub (j) = node j significance

W sub (j) = weighted number of samples hitting node j

C sub (j) = node j impurity value j

left (j) = child node from left split on node j

right (j) = child node from right split on node j

sub () is used because subscript is not available in medium.

After that, the value of each function on a decision tree is determined as follows:

$$f_{i_i} = \frac{\sum_{j: \text{node } j \text{ splits on feature } i} ni_j}{\sum_{k \in \text{all nodes}} ni_k} \quad (3)$$

normfi sub i = the feature’s normalized importance i

fi sub i = the significance of i

By dividing by the number of all function importance values, these can be normalized to a value between 0 and 1.

$$\text{norm} f_{i_i} = \frac{f_{i_i}}{\sum_{j \in \text{all features}} f_{i_j}} \quad (4)$$

At the Random Forest stage, the final feature importance is the average of all the trees. The total number of trees is divided by the amount of the feature’s significance value on each tree.

$$RF f_{i_i} = \frac{\sum_{j \in \text{all trees}} \text{norm} f_{i_{ij}}}{T} \quad (5)$$

RFfi sub (i) = the importance of feature i calculated from all trees in the Random Forest model

normfi sub (ij) = the normalized feature importance for i in tree j

T = total number of trees

As Huljanah et al. [56] stated, the classification has produced the class that the majority of the decision trees predicted in the forest. On the other hand, For regression outputs, the class with mean predictions of the individual trees. Based on the given scenario, trees and entropy criteria are used for developing the classifier. The entropy criterion is represented by

$$\text{Entropy} = - \sum_{j=1}^c p_j \log 2(p_j) \quad (6)$$

A “Random forest classifier”, according to Mehrang et al. [57], shows a collection of decision trees built from randomly selected subsets of the training set. The sepsis data is then aggregated from various decision trees, determining the final class of the test object. The Random Forest Prediction depicts the model’s pictorial representation as



well as its precision for the dataset's data collection. Using the Random Forest, the accuracy of the heart disease detection system is 0.821875.

Furthermore, by inputting the XTrain and YTrain parameters required for the method of fit, the accuracy value of the "Random Forest Score model" is obtained. Then, to determine the score of this model, use the parameters XTest and YTest in the method score() to determine the Random Forest Score model's score.

### 5.6. Logistic regression

A widely used classification methodology is the logistic regression model. The variables of the logistic regression show the class. This class is a variable of categorical dependency. Ge et al. [58] explained that the logistic regression includes the logistic equation, with the model as the dependent variable. Along with that, according to Saif et al. [59], a supervised learning classification algorithm called logistic regression is used to estimate the likelihood of a target variable. This algorithm of the target variable is dependent on the categorical method that is used. Aside from that, the dependent variable has a dichotomous nature, which means it can be classified into two groups.

### 5.7. Summary

These models are compared depending on precision, specificity, f-measure, accuracy as well as sensitivity. The VLRAKN achieved the highest accuracy when it is compared with algorithms such as naive Bayes, random forest classification, logistic regression and decision trees. With the observation of this approach, Random forest classification gives the parameters for the algorithms are based on the best accuracy score of the data reading.

## 6. Conclusion

### 6.1. Conclusion

Based on the given scenario, the first section discusses heart disease prediction using Python. Python is object-oriented as well as it is also a high-level programming language that has quick development cycles and spirited, energetic building options. This language helps better to be able to predict the heart disease pathway accurately. The heart care industry is generating the data from several facilities as well as patients over applying the best use of the strategy of this data. Apart from this, doctors are easily demonstrating this superior model for treatments and it will be improving the complete delivery system of the healthcare sector. This prediction model of heart disease is especially used in heart diseases, clinicians, and institutions that can provide better along with that improvised outcomes of the patients over scalable and dynamic applications and conclude the problem of this model. Apart from this, chapter two discussed detecting the presence of heart diseases using Python. This application depends on the heart disease dataset that involves data of the patients, which are age, sex, chol, treetops, and many more.

On the other hand, individually import libraries such as matplotlib, Numpy, Pandas, warnings, and many more are used in this application. This python language is one of the robust languages that foster computational capabilities in enhancing valuable insights from the information of the patients suffering from heart diseases. However, it also complies with the HIPAA regulations that ensure medical information safety. Aside from that, chapter three delves into the technique for detecting heart disease using machine learning, including several algorithms. Machine learning is important for predicting the existence of a threat. This project follows the random forest algorithm for developing heart disease detection. This algorithm is used for the methodology

to detect heart disease using Python. This application is considered significant according to its 83% (approximately) accuracy rate over training data.

Along with that, a model regarding ML has played a significant role in creating accuracy and determining results with the aid of training data. According to the scenario, it can be highlighted that the selection of random forest classification is developed based on the exact dataset as well as the decision tree. In this section, the emphasis for discussion is on data analysis. It works with categorical variables along with that; it will break particular categorical columns into dummy columns with 1s and 0s. Apart from this, this part also mentioned that the data output of this application is used for many medical parameters such as age, gender, blood pressure, cholesterol, and obesity for prediction and requirement of the software used in the development application. Besides this, the Machine learning application using Python is a subset of the Artificial Intelligence model and the python libraries are the prerequisites for making predictions that SKLEARN is normally used in machine learning prediction. Except for the Decision Tree, the best values provided by the ML model are provided by Random Forest. This is the simplest method for predicting heart disease to produce precise results.

### 6.2. Recommendation

**Recommendation 1:** The aim is to introduce a dataset auditing technological setup for removing issues within the structure, as shown in [Table 1](#).

**Recommendation 2:** The aim is to require proper training regarding this machine learning approach for users, see [Table 2](#).

### 6.3. Linking to objectives

#### Objective 1:

The first goal defines critical analysis of the python language used to detect heart disease.

In the literature review, Section 2.3 has highlighted the development approach of addressing the significant research in the health care sector. The major beneficial factor of Python within the health care sector is that it assists in making sense of the information by working with Machine Learning and AI within the healthcare sectors. Thus it is related to the first objective.

Section 4.2 has described the analysis according to the details of heart disease detection. Therefore, this objective is successfully met in the research.

#### Objective 2:

The second objective defines investigating the previous activities critically and applying a suitable methodological approach for super-scribing the identified problem.

Section 2.4 has highlighted the capability to evaluate, synthesize, and search the information from the appropriate sources in health care sectors in the literature review. Thus, it is related to the second objective. Section 4.3 has described the data interpretation of the selected dataset. Therefore, the objective is successfully met.

#### Objective 3:

The third objective defines critically applied data interpretation strategies in python language for health issue detection.

In the literature review, Section 2.2 highlighted a deep learning understanding of individual interest associated with specialized computing in healthcare. Thus it is related to the third objective.

Section 4.4 has described different data interpretation strategies based on python language.

#### Objective 4:

The fourth objective defines critically assessing the artifact or product with the help of cybersecurity approaches using appropriate methods and identifying the limitations and strengths of the work.

**Table 1**  
Summary of recommendation 1.

S-specific	To introduce a data auditing technological setup
M-measurable	It can be measured through the success rate of doctors' treatment in detecting heart disease.
A-attainable	A daily examination of the selected dataset will be required to detect heart disease without any error.
R-realistic	It will be effective for removing barriers in the dataset and reducing the chances of error within it.
T-time plan	5 months

**Table 2**  
Summary of recommendation 2.

S-specific	To introduce a machine learning training
M-measurable	It can be measured over success rate according to the patients
A-attainable	Arrange the weekly training for the users
R-realistic	It will be effective for removing the boundary in training as well as decreases the error in the training
T-time plan	6 months

In the literature review, Section 2.5 has highlighted the Critical application of the cybersecurity techniques that conform to networking configurations within the healthcare industries and the information security management system. Section 4.5 has described various methods of cybersecurity. Thus, it is related to the fourth objective.

#### 6.4. Limitations

Traditional invasive-based approaches and angiography are also considered well-known appropriate techniques for diagnosing cardiac conditions in this study. It is a drawback in heart disease prediction. Apart from this, intelligent learning is dependent on computational techniques that are upright and efficient in predicting the occurrence of heart disease. This prediction method is presented here for the purposes of heart disease prediction and diagnosis. This classification technique depends on different pruning as well as data cleaning techniques. This technique is prepared as well as developed a dataset that is suitable for data mining and selects the proper technique, which provides much better accuracy of this application.

#### 6.5. Research contributions

In this research, we develop a healthcare application to help detect heart diseases among patients and those with symptoms. Based on the random forest algorithm, our work provides better accuracies. It has a very low cost for development. Additionally, research outputs can be used as valuable data for further analyses. To achieve rational and logical arguments, we can then develop better diagnoses to detect heart diseases, including findings and interpretations. The artifact developed from this research will assist in evaluating third parties with the assistance of this appropriate method. Our research is aimed at offering both theoretical and practical contributions to healthcare.

#### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

#### Acknowledgment

This work is partly supported by VC Research (VCR 0000154).

#### Appendix. Codes

```
import numpy as np
import pandas as pd
df = pd.read_csv('dataset.csv')
df.describe()
import matplotlib.pyplot as plt
from matplotlib import rcParams
from matplotlib.cm import rainbow
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')
import seaborn as sns
# get correlations of each features in dataset
corrmat = df.corr()
uniform_data = corrmat.index
plt.figure(figsize=(20,20))
# plot heat map
g=sns.heatmap(df[uniform_data].corr(),annot=True,cmap="rocket")
df.hist(figsize=(20,20))
sns.set_style('whitegrid')
sns.countplot(x='target',data=df,palette='flare')
data = pd.get_dummies(df, columns=['sex', 'cp', 'fbs', 'restecg', 'exang', 'slope', 'ca', 'thal'])
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
standardScaler = StandardScaler()
columns_to_scale = ['age', 'trestbps', 'chol', 'thalach', 'oldpeak']
data[columns_to_scale] = standardScaler.fit_transform(data[columns_to_scale])
data.head(10)
y = data['target']
X = data.drop(['target'], axis = 1)
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import cross_val_score
knn_scores = []
for k in range(1,11):
    knn_classifier = KNeighborsClassifier(n_neighbors = k)
    score=cross_val_score(knn_classifier,X,y,cv=10)
    knn_scores.append(score.mean())
plt.plot([k for k in range(1, 11)], knn_scores, color = 'red')
for i in range(1,11):
    plt.text(i, knn_scores[i-1], (i, knn_scores[i-1]))
plt.xticks([i for i in range(1, 11)])
plt.xlabel('Number of Neighbors (K)')
plt.ylabel('Scores')
plt.title('K Neighbors Classifier scores for different K values')
from sklearn.ensemble import RandomForestClassifier
randomforest_classifier= RandomForestClassifier(n_estimators=10)
score=cross_val_score(randomforest_classifier,X,y,cv=10)
score.mean()
```

#### References

- [1] L. Loku, B. Fetaji, A. Krstev, M. Fetaji, Z. Zdravev, Using python programming for assessing and solving health management issues, *South East Eur. J. Sustain. Dev.* 4 (1) (2020).

- [2] P. Mathur, Overview of machine learning in healthcare, in: *Machine Learning Applications using Python*, A Press, Berkeley, CA, 2019, pp. 1–11.
- [3] P. Guleria, M. Sood, Intelligent learning analytics in healthcare sector using machine learning, in: *Machine Learning with Health Care Perspective*, Springer, Cham, 2020, pp. 39–55.
- [4] V.V. Kumar, Healthcare Analytics Made Simple: Techniques in Healthcare Computing using Machine Learning and Python, Packt Publishing Ltd., 2018, Available at: [https://books.google.com/books?hl=en&lr=&id=nwZnDwAAQBAJ&oi=fnd&pg=PP1&dq=application+of+python+programming+language+in+health+care+sectors&ots=BjzOQe\\_09q&sig=Bhv\\_ixOuZoKJu-WVEKPD6B9y0](https://books.google.com/books?hl=en&lr=&id=nwZnDwAAQBAJ&oi=fnd&pg=PP1&dq=application+of+python+programming+language+in+health+care+sectors&ots=BjzOQe_09q&sig=Bhv_ixOuZoKJu-WVEKPD6B9y0), [Accessed on 5th March, 2021].
- [5] BelltSoft, Python in healthcare, BelltSoft (2017) Available at: <https://beltssoft.com/custom-application-development-services/healthcare-software-development/python-healthcare>, [Accessed on 5th March, 2021].
- [6] B. Nithya, V. Ilango, Predictive analytics in health care using machine learning tools and techniques, in: 2017 International Conference on Intelligent Computing and Control Systems (ICICCS), IEEE, 2017, pp. 492–499.
- [7] J. McPadden, T.J. Durant, D.R. Bunch, A. Coppi, N. Price, K. Rodgerson, C.J. Torre Jr., W. Byron, A.L. Hsiao, H.M. Krumholz, W.L. Schulz, Health care and precision medicine research: analysis of a scalable data science platform, *J. Med. Internet Res.* 21 (4) (2019) e13043–e13045.
- [8] A. Panesar, Machine Learning and AI for Healthcare (1–73), Apress, Coventry, UK, 2019, Available at [https://iedu.us/wp-content/uploads/edd/2020/01/Arjun\\_Panesar\\_Machine\\_Learning\\_and\\_AI\\_for\\_Health-iedu.us.pdf](https://iedu.us/wp-content/uploads/edd/2020/01/Arjun_Panesar_Machine_Learning_and_AI_for_Health-iedu.us.pdf), [Accessed on 5th March, 2021].
- [9] T. Davenport, R. Kalakota, The potential for artificial intelligence in healthcare, *Future Healthcare J.* 6 (2) (2019) 94.
- [10] Analytics Vidhya, Commonly used machine learning algorithms (with python and R codes), 2021, Analytics Vidhya. Available at: <https://www.analyticsvidhya.com/blog/2017/09/common-machine-learning-algorithms/>, [Accessed on 5th March, 2021].
- [11] H. Mayfield, C. Smith, M. Gallagher, M. Hockings, Use of freely available datasets and machine learning methods in predicting deforestation, *Environ. Model. Softw.* 87 (2017) 17–28.
- [12] P. Barot, Why use python in healthcare applications? BoTree Technologies. (2020) Available at: <https://www.botreetechnologies.com/blog/python-in-healthcare-application/>, [Accessed on 5th March, 2021].
- [13] K. Bhanot, Predicting presence of heart diseases using machine learning, Towards Data Sci. (2019) Available at: <https://towardsdatascience.com/predicting-presence-of-heart-diseases-using-machine-learning-36f00f3edb2c>, [Accessed on 5th March, 2021].
- [14] J. Zaidi, Project: Predicting heart disease with classification machine learning algorithms, Towards Data Sci. (2020) Available at: <https://towardsdatascience.com/project-predicting-heart-disease-with-classification-machine-learning-algorithms-fd69e6fd9d6>, [Accessed on 5th March, 2021].
- [15] C. Ozgur, T. Colliau, G. Rogers, Z. Hughes, Matlab vs. Python vs. R, *J. Data Sci.* 15 (3) (2017) 355–371.
- [16] K.R. Srinath, Python—the fastest growing programming language, *Int. Res. J. Eng. Technol.* 4 (12) (2017) 354–357.
- [17] B. Copeland, Advancing opportunities in health care with python-based machine learning, Chief HealthCare Executive (2019) Available at: <https://www.chiefhealthcareexecutive.com/view/advancing-opportunities-in-healthcare-with-python-based-machine-learning>, [Accessed on 5th March, 2021].
- [18] W. Jiang, M. Zhuang, C. Xie, J. Wu, Sensing attribute weights: A novel basic belief assignment method, *Sensors* 17 (4) (2017) 721.
- [19] G.J. van den Burg, A. Nazábal, C. Sutton, Wrangling messy CSV files by detecting row and type patterns, *Data Min. Knowl. Discov.* 33 (6) (2019) 1799–1820.
- [20] C. Holdgraf, Case study 7: Feature extraction and data wrangling for predictive models of the brain in python, in: *The Practice of Reproducible Research*, University of California Press, 2017, pp. 139–148.
- [21] R.A. Calix, S.B. Singh, T. Chen, D. Zhang, M. Tu, Cyber security tool kit (CyberSecTK): A python library for machine learning and cyber security, *Information* 11 (2) (2020) 100.
- [22] H. Cui, F. Li, Andes: A python-based cyber-physical power system simulation tool, in: 2018 North American Power Symposium (NAPS), IEEE, 2018, pp. 1–6.
- [23] Python.org, What Is Python? Executive Summary, Python.org, 2020, <https://www.python.org/doc/essays/blurb/#:text=Python%20is%20an%20interpreted%2C%20object,programming%20language%20with%20dynamic%20semantics.&text=Python's%20simple%2C%20easy%20to%20learn,program%20modularity%20and%20code%20reuse>, [Accessed on 5th March, 2021].
- [24] D. Bau, J. Gray, C. Kelleher, J. Sheldon, F. Turbak, Learnable programming: blocks and beyond, *Commun. ACM* 60 (6) (2017) 72–80.
- [25] Medium, Heart disease detection using machine learning in python, Medium (2021) Available at: <https://randerson112358.medium.com/heart-disease-detection-using-machine-learning-python-a701f39396cb>, [Accessed on 5th March, 2021].
- [26] C. Iwendi, A.K. Bashir, A. Peshkar, R. Sujatha, J.M. Chatterjee, S. Pasupuleti, R. Mishra, S. Pillai, O. Jo, COVID-19 patient health prediction using boosted random forest algorithm, *Front. Public Health* 8 (2020) 357.
- [27] A. Navlani, Understanding random forests classifier in python, DataCamp (2018) Available at: <https://www.datacamp.com/community/tutorials/random-forests-classifier-python>, [Accessed on 5th March, 2021].
- [28] Z. Noshad, N. Javaid, T. Saba, Z. Wadud, M.Q. Saleem, M.E. Alzahrani, O.E. Sheta, Fault detection in wireless sensor networks through the random forest classifier, *Sensors* 19 (7) (2019) 1–21.
- [29] A.H. Larsen, J.J. Mortensen, J. Blomqvist, I.E. Castelli, R. Christensen, M. Duřak, J. Friis, M.N. Groves, B. Hammer, C. Hargus, E.D. Hermes, The atomic simulation environment—a python library for working with atoms, *J. Phys.: Condens. Matter* 29 (27) (2017) 1–58.
- [30] S. Amini, S. Homayouni, A. Safari, A.A. Darvishsefat, Object-based classification of hyperspectral data using random forest algorithm, *Geo-Spat. Inf. Sci.* 21 (2) (2018) 127–138.
- [31] P. Muñoz, J. Orellana-Alvear, P. Willems, R. Céleri, Flash-flood forecasting in an andean mountain catchment—Development of a step-wise methodology based on the random forest algorithm, *Water* 10 (11) (2018) 1–18.
- [32] M. Reich, T. Tabor, T. Liefeld, H. Thorvaldsdóttir, B. Hill, P. Tamayo, J.P. Mesirov, The GenePattern notebook environment, *Cell Syst.* 5 (2) (2017) 149–151.
- [33] K.M. Mendez, L. Pritchard, S.N. Reinke, D.I. Broadhurst, Toward collaborative open data science in metabolomics using Jupyter Notebooks and cloud computing, *Metabolomics* 15 (10) (2019) 1–16.
- [34] D. Yin, Y. Liu, H. Hu, J. Terstriep, X. Hong, A. Padmanabhan, S. Wang, Cybergis-jupyter for reproducible and scalable geospatial analytics, *Concurr. Comput.: Pract. Exper.* 31 (11) (2019) 1–14.
- [35] V.V. Ramalingam, A. Dandapath, M.K. Raja, Heart disease prediction using machine learning techniques: a survey, *Int. J. Eng. Technol.* 7 (2.8) (2018) 684–687.
- [36] K. Subhadra, B. Vikas, Neural network based intelligent system for predicting heart disease, *Int. J. Innov. Technol. Explor. Eng.* 8 (5) (2019) 484–487.
- [37] B. Ambale-Venkatesh, X. Yang, C.O. Wu, K. Liu, W.G. Hundley, R. McClelland, A.S. Gomes, A.R. Folsom, S. Shea, E. Guallar, D.A. Bluemke, Cardiovascular event prediction by machine learning: the multi-ethnic study of atherosclerosis, *Circ. Res.* 121 (9) (2017) 1092–1101.
- [38] M.V. Dogan, I.M. Grumbach, J.J. Michaelson, R.A. Philibert, Integrated genetic and epigenetic prediction of coronary heart disease in the framingham heart study, *PLoS One* 13 (1) (2018) 1–18.
- [39] M.S. Mahdavinjad, M. Rezvan, M. Barekatain, P. Adibi, P. Barnaghi, A.P. Sheth, Machine learning for internet of things data analysis: A survey, *Digit. Commun. Netw.* 4 (3) (2018) 161–175.
- [40] G. Tazuin, U. Lupo, L. Tunstall, J.B. Pérez, M. Caorsi, A.M. Medina-Mardones, A. Dassatti, K. Hess, Giotto-tda: A topological data analysis toolkit for machine learning and data exploration, *J. Mach. Learn. Res.* 22 (39) (2021) 1–6.
- [41] B. Peters, E. Haber, J. Granek, Neural networks for geophysicists and their application to seismic data interpretation, *The Leading Edge* 38 (7) (2019) 534–540.
- [42] G. Carleo, I. Cirac, K. Cranmer, L. Daudet, M. Schuld, N. Tishby, L. Vogt-Maranto, L. Zdeborová, Machine learning and the physical sciences, *Rev. Modern Phys.* 91 (4) (2019) 1–47.
- [43] R. Harper, S.T. Flammia, J.J. Wallman, Efficient learning of quantum noise, *Nat. Phys.* 16 (12) (2020) 1184–1188.
- [44] Y. Jia, Y. Qi, H. Shang, R. Jiang, A. Li, A practical approach to constructing a knowledge graph for cybersecurity, *Engineering* 4 (1) (2018) 53–60.
- [45] D. Galinec, D. Možnik, B. Guberina, Cybersecurity and cyber defence: national level strategic approach, *Automatika: Časopis za automatiku, Mjerenje, Elektroniku, Računarstvo I Komunikacije* (2017).
- [46] Y. Fan, J. Li, D. Zhang, J. Pi, J. Song, G. Zhao, Supporting sustainable maintenance of substations under cyber-threats: An evaluation method of cybersecurity risk for power CPS, *Sustainability* 11 (4) (2019) 1–30.
- [47] G. Jin, M. Tu, T.H. Kim, J. Heffron, J. White, Evaluation of game-based learning in cybersecurity education for high school students, *J. Educ. Learn.* 12 (1) (2018) 150–158.
- [48] A. Anees, I. Hussain, A novel method to identify initial values of chaotic maps in cybersecurity, *Symmetry* 11 (2) (2019) 140.
- [49] F. Mehmood, H.U. Rashidkayani, F. Hussain, Chronic diseases modelling—python environment, *FUJAST J. Biol.* 10 (1) (2020) 31–38.
- [50] H. Tjahjadi, K. Ramli, Noninvasive blood pressure classification based on photoplethysmography using K-nearest neighbors algorithm: A feasibility study, *Information* 11 (2) (2020) 1–18.
- [51] D. Pedrozo, F. Barajas, A. Estupiñán, K.L. Cristiano, D.A. Triana, Data analysis for a set of university student lists using the k-Nearest Neighbors machine learning method, *J. Phys. Conf. Ser.* 1514 (1) (2020) 1–8.
- [52] T.A. Assegie, P.S. Nair, Handwritten digits recognition with decision tree classification: a machine learning approach, *Int. J. Electr. Comput. Eng.* 9 (5) (2019) 1–4.
- [53] S. Herbold, Autorank: A python package for automated ranking of classifiers, *J. Open Source Softw.* 5 (48) (2020) 1–4.

- [54] A. Zahariev, M. Zveryakov, S. Prodanov, G. Zaharieva, P. Angelov, S. Zarkova, M. Petrova, Debt management evaluation through support vector machines: on the example of Italy and Greece, *Entrepreneurship Sustain. Issues* 7 (3) (2020) 1–12.
- [55] C. Wang, C. Liang, Msipred: a python package for tumor microsatellite instability classification from tumor mutation annotation data using a support vector machine, *Sci. Rep.* 8 (1) (2018) 1–10.
- [56] M. Huljanah, Z. Rustam, S. Utama, T. Siswantining, Feature selection using random forest classifier for predicting prostate cancer, *IOP Conf. Ser.: Mater. Sci. Eng.* 546 (5) (2019) 1–9.
- [57] S. Mehrang, J. Pietilä, I. Korhonen, An activity recognition framework deploying the random forest classifier and a single optical heart rate monitoring and triaxial accelerometer wrist-band, *Sensors* 18 (2) (2018) 613.
- [58] J. Ge, X. Li, H. Jiang, H. Liu, T. Zhang, M. Wang, T. Zhao, Picasso: A sparse learning library for high dimensional data analysis in R and python, *J. Mach. Learn. Res.* 20 (44) (2019) 1–5.
- [59] M.A. Saif, A.N. Medvedev, M.A. Medvedev, T. Atanasova, Classification of online toxic comments using the logistic regression and neural networks models, in: *AIP Conference Proceedings*, Vol. 2048 (1) 2018, pp. 1–6.

# Heart Disease Prediction Using Machine Learning Algorithms

## Abstract

Heart disease remains a leading cause of mortality worldwide. Early and accurate detection is crucial for effective treatment and improved patient outcomes. This study explores the application of machine learning algorithms to predict heart disease using Python. Leveraging the UCI Heart Disease dataset, we implemented and evaluated models including Logistic Regression, K-Nearest Neighbors (KNN), Support Vector Machine (SVM), Decision Tree, Random Forest, and Gradient Boosting. Our findings indicate that the Random Forest classifier achieved the highest accuracy at 94%, suggesting its potential as a robust tool for heart disease prediction.

## 1. Introduction

Cardiovascular diseases (CVDs) are a significant global health concern, accounting for approximately 17.9 million deaths annually. Early detection is vital to reduce mortality rates and enhance patient care. Traditional diagnostic methods can be subjective and time-consuming, highlighting the need for automated, data-driven approaches. Machine learning offers promising solutions by analyzing complex medical data to identify patterns indicative of heart disease.

## 2. Methodology

### 2.1 Dataset

We utilized the UCI Heart Disease dataset, comprising 303 instances with 14 attributes, including age, sex, blood pressure, cholesterol levels, and other relevant medical indicators.

### 2.2 Data Preprocessing

Data preprocessing involved handling missing values, normalizing numerical features, and encoding categorical variables. The dataset was split into training (75%) and testing (25%) subsets to



evaluate model performance.

## 2.3 Feature Selection

Recursive Feature Elimination (RFE) was employed to identify the most significant features contributing to heart disease prediction. This process enhances model performance by eliminating irrelevant or redundant features.

## 2.4 Model Implementation

We implemented the following machine learning algorithms using Python's scikit-learn library:

- Logistic Regression
- K-Nearest Neighbors (KNN)
- Support Vector Machine (SVM)
- Decision Tree
- Random Forest
- Gradient Boosting

## 2.5 Model Evaluation

Models were evaluated using metrics such as accuracy, precision, recall, F1-score, and AUC-ROC.

## 3. Results

Performance Summary:

- Logistic Regression: 85.2%
- K-Nearest Neighbors: 86.89%
- Support Vector Machine: 85.2%
- Decision Tree: 99.9%
- Random Forest: 94%
- Gradient Boosting: 73%

#### 4. Discussion

The results underscore the efficacy of ensemble methods, particularly Random Forest, in predicting heart disease. The high accuracy achieved aligns with findings from previous studies. Feature selection through RFE contributed to enhancing model performance by focusing on the most relevant attributes.

#### 5. Conclusion

This study demonstrates the potential of machine learning algorithms, especially the Random Forest classifier, in accurately predicting heart disease. Implementing such models in clinical settings can aid healthcare professionals in early diagnosis and intervention, ultimately improving patient outcomes.

#### References:

- [ijert.org](http://ijert.org)
- [arxiv.org](http://arxiv.org)
- [researchgate.net](http://researchgate.net)