

# Prediction of Cardiovascular Disease using Machine Learning Algorithms

1. S RamaKrishna Vasamsetti, Research Scholar, Department of Computer Science and Engineering, J.S University, Shikohabad, U.P.
2. Dr. Suribabu Potnuri , Professor ,Supervisor, Department of Computer Science and Engineering, J.S University, Shikohabad, U.P.

## ABSTRACT

Providing medical care is an unavoidable obligation that must be fulfilled for the whole of human life. The term "cardiovascular disease" refers to a wide category that encompasses a number of conditions that have various effects on the heart and blood vessels. It was easier to make judgments on the adjustments that should be made in high-risk people thanks to the first approaches for diagnosing cardiovascular disorders, which ultimately led to a decrease in the dangers that these individuals suffered. Ways of coming: Kaggle has been the source of the data that we have used in the study that we have suggested. These data do not need any pre-processing systems, such as those that remove noisy data, handle missing information, fill in default values if they are required, and categorize characteristics for the purpose of prediction and decision-making at several levels. Methods such as classification, accuracy, sensitivity, and specificity analysis are used in order to assess the effectiveness of the diagnostic model. The purpose of this study is to increase awareness and provide a diagnosis by

presenting a prediction model that may be used to determine whether or not a person has cardiovascular disease. The evaluation of the accuracy of applying rules to the individual outputs of Support Vector Machine, Random Forest, Naive Bayes classifier, and logistic regression on a dataset gathered from a particular area is the method that is used to achieve this goal. The purpose of this endeavor is to develop an accurate model for the prediction of cardiovascular disease. According to the findings, the machine learning algorithms that were investigated had the capacity to predict cardiovascular sickness in people with an accuracy that ranged from 58.71% to 77.06%. Simply put, When compared to a large number of different machine learning methods, Logistic Regression revealed a much higher level of accuracy (77.06%). **Keywords:** Cardiovascular disease, Machine Learning Algorithms, Performance Evaluators, toxins

## INTRODUCTION

When it comes to the process of data mining, classification is an essential component. The process of systematically finding and categorizing data categories or ideas via the building of a model or function that elucidates and differentiates them is referred to as classification. In order to do this, the technique of categorizing is used. The construction of the model is generated from a study of a substantial quantity of data pertaining to cardiovascular exercise, in which the class labels were already present before. It is possible to use the technique to predict the class label of products whose cardiovascular disease condition is unknown. The existence or absence of cardiovascular sickness in the participants determines whether or not these forecasts are accurate. The field of machine learning is concerned with the investigation of how computers may improve their performance by evaluating data pertaining to the cardiovascular system. The primary emphasis of the research is on computer algorithms that are designed to gain the capability of recognizing detailed patterns and making intelligent judgments based on data pertaining to cardiovascular health. At the end of the day, categorization is fundamentally the same thing as directed learning. In the dataset pertaining to cardiovascular training, the assigned models are responsible for managing the intake category via their oversight.

The incidence of cardiovascular disease, which is more frequently referred to as CVD, is seeing an upward trend in the contemporary world. It is estimated that over 17 million individuals pass away every year as a result of cardiovascular disorders, namely heart attacks and strokes, as stated by the World Health Organization (WHO). It is essential to keep a detailed journal of the principal unfavorable effects and unhealthy behaviors that lead to the development of cardiovascular disease (CVD) in order to reduce the risk of developing this condition. A great number of procedures are carried out in order to arrive at a conclusion about cardiovascular disease (CVD). Auscultation, electrocardiogram (ECG), blood pressure testing, and evaluation of cholesterol and glucose levels are all examinations that are included in the examinations.

Even in situations when the patient's health may be in a critical state and immediate treatment is necessary, the duration of these tests usually takes a significant amount of time. Thus, it is very necessary to place testing at the top of the priority list above all else [2]. It is possible that the development of cardiovascular disease (CVD) is caused by a number of behaviors that may be effectively avoided. Understanding the behaviors that are associated with cardiovascular disease (CVD) is of the utmost importance. Machine learning is seeing considerable growth as a result of the exponential increase in the amount of data

that is being collected. Machine learning makes it possible for humans to get insight from the vast amounts of data that are available to them, despite the fact that the wealth of data may be overwhelming and difficult for them to acquire. As a consequence of this, the remaining portion of the discussion on the study is organized in the following manner: The following section offers a condensed summary of the results of the earlier research. A condensed summary of the various machine learning techniques that were selected is presented in the third section. The fourth component offers a comprehensive, in-depth representation of the Patient Data Set, which includes a description of its features and properties. The discussion of the proposed technique will be the primary emphasis of the fifth section. Section VI of the manual contains a discussion of the performance measures of categorization. A condensed analysis and discussion of the results is presented in the seventh section of the exposition. The completion of the research activity is determined by the information included in Section VIII of the text. In conclusion, the References are discussed in the last section of the article.

### **A.Cardiovascular disease**

Generally speaking, the term "cardiovascular infection" refers to diseases that are characterized by restricted or clogged blood arteries. These illnesses may result in a heart attack, chest discomfort (also known as angina), or stroke at some point. Heart conditions that affect the myocardium, valves, or rhythm of the heart are

also regarded to be kinds of coronary disease. Other cardiac illnesses include those that impact the heart's rhythm.

The term "cardiovascular malady" refers to a group of conditions that have an effect on the structure or operation of the heart. These conditions include coronary artery disease, which is characterized by the constriction of the arteries.

A cardiac arrest has occurred. Arrhythmias, which are irregular heartbeats as they are known. Heart failure is a condition. It is an infectious endocarditis. Circumstances that are present at birth.

A condition known as cardiomyopathy Plaque, also known as atheroma, is a kind of fatty deposit that may cause cardiovascular disease. This condition happens when the arteries that provide blood and oxygen to the heart muscle and other organs, such as the brain and kidneys, get blocked by the cholesterol deposits. Atherosclerosis is the word known in the medical field to describe this method.

When referring to the heart, the word "cardio" is used, while the term "vascular" refers to all of the veins that are found throughout the body. Specifically, the term "coronary illness" refers to conditions that affect the heart, including coronary artery disease, heart failure, anomalies in the heart valves, and irregular cardiac rhythms.

## 1. LITERATURE SURVEY

A.U. UIHaq, J.P. Li, M.H. Memon, S. Nazir, and R. Sun have investigated the notion of heart disease and have developed a technique for identifying heart sickness that is based on machine learning. This approach makes use of a dataset that was particularly built for this purpose. A total of seven well-known machine learning algorithms, three-element decision algorithms, the cross-validation approach, and seven classifiers performance assessment measures have been applied by them. These metrics include classification accuracy, specificity, sensitivity, Matthews' correlation coefficient, and execution time. They have devised a method that is capable of effectively identifying and classifying persons with heart illness in comparison to those who are healthy or otherwise healthy. The research investigates the total number of classifiers, as well as the feature assurance calculations, pre-processing approaches, validation method, and performance assessment estimations that were presented. Both the whole set of qualities and a subset of those characteristics have been taken into consideration when deciding whether or not to authorize the deployment of the proposed system. Both the accuracy and the execution time of classifiers are significantly impacted as a result of their specific properties. According to their hypothesis, decision support networks that are based on machine learning will be of great assistance to

medical professionals in effectively detecting cardiac patients.

Geetha S. S. Krishnan has devised a method that is capable of reliably predicting the outcomes that are likely to occur as a result of heart disease. There is a considerable increase in the probability of acquiring heart disease as a result of this system's future repercussions. The datasets that were analyzed are organized in a fashion that is comparable to the principles of treatment parameters. These criteria are evaluated by their structure via the use of the data mining plan methodology. Python programming was used to create the datasets, and two standard machine learning algorithms, namely the Decision Tree Algorithm and the Naive Bayes Algorithm, were used in the operation of the datasets. When it came to forecasting cardiac disease, these algorithms displayed an exceptional level of accuracy.

The prediction of heart disease has been discussed by K.G. Dinesh, K.A.raj, K.D.Santhosh, and V. M.eswari. They have also conducted data preprocessing using techniques such as the removal of noise, the removal of missing data, the filling of default values when it is necessary, and attribute classification for the purpose of prediction and decision making at various levels. Using techniques such as classification, precision, sensitivity, and specificity analysis, they succeeded in demonstrating the discovery model. This was

accomplished by the employment of these processes. The purpose of this research was to construct a prediction model that can identify people who are at risk for developing heart disease and then give insights or suggestions based on the findings of the study. In order to develop an appropriate model for predicting cardiovascular illness, we used rules to a dataset that was gathered from a district. These rules included Support Vector Machine, Gradient Boosting, Random Forest, Naive Bayes classifier, and linear regression.

Both men and females were included in the research that Golande and Kumar T[7] conducted, which investigated the prevalence of heart illness. According to what they discovered, this proportion can be different in different regions. In addition, they concentrated their attention especially on individuals who were between the ages of 25 and 69. It is not the case that individuals who belong to various age groups would not be impacted by heart disorders, as this does not indicate. There is a significant problem that has been anticipated in the present day, and that is the predominance of cardiovascular disease. There have been a number of calculations and instruments that have been addressed in relation to the prediction of cardiovascular illnesses.

Machine learning approaches were used by Prasad, P., Anjali, S., Adil, N., and Deepa [8] in order to make predictions about heart disorders

by linking the research that was previously conducted. The methodology of computed regression was used by the researchers in order to assess the data pertaining to healthcare and categorize people based on whether or not they were suffering from heart problems. Following that, they constructed a prediction model that could assess whether or not a patient was suffering from heart ailment.

## 2. MACHINELEARNINGALGORIT HMS

In the field of education, machine learning refers to a kind of learning that has been modernized and needs little involvement from humans. The process entails teaching computers to learn from data sources that are accessible to the general population. An analysis and development of algorithms that are able to learn from previous data and make predictions based on new data is the core goal of artificial intelligence (AI).

When it comes to learning calculations, the preparation of information and the representation of knowledge are both components that contribute to the whole process. The result is the accumulation of knowledge, which often takes the shape of a new algorithm that is able to carry out a job. There are many different types of data that may be used to train a machine learning system. These include numerical data, text, audio, pictures, and videos. The framework's yield information might be a floating-point number, depending on

the circumstances.

### A. Concepts of Learning

- Learning is the way toward changing over understanding into skill or information.
- Learning can be comprehensively grouped into three classes, as referenced beneath.

In view of the idea of the learning information and association between the student and the earth.

- Supervised Learning process or Supervised Learning Approach
- Unsupervised Learning process or Unsupervised Learning Approach
- Semi-regulated Learning process or Unsupervised Learning Approach

Correspondingly, there are four classifications of Machine Learning as appeared beneath –

- Supervised learning process/Approach
- Unsupervised learning process/Approach
- Semi-directed learning process/Approach
- Reinforcement learning process/Approach

In any case, the most normally utilized ones are supervised and unsupervised learning

#### **Supervised Learning**

Practical uses of machine learning include face and voice recognition, product or movie recommendations, and sales forecasting. Machine learning is also utilized in a variety of other applications. Both regression and

classification are two forms of supervised learning that may be further subdivided into their respective categories.

A reliable result may be anticipated by the use of regression analysis, which is a statistical procedure. One example of this would be price forecasting for real estate.

In the process of characterization, one strives to determine the right category or classification. For example, one may investigate good and negative emotions, male and female individuals, benign and malignant illnesses, secured and unsecured loans, and so on.

The construction of a machine learning model that is dependent on certain training samples is what is meant by the term "supervised learning."

For instance, if we develop a system that can determine the type of fever based on various patient characteristics, such as temperature, the intensity of the headache, body aches, cough and cold, and other blood status parameters, we will be able to categorize the patient as having malaria, dengue fever, viral fever, sinusitis, and so on.

The process of accumulating knowledge via the use of accessible training data is known as supervised learning. Logistic Regression, Neural Networks, Support Vector Machines (SVMs), and Naive Bayes classifiers are some examples of the supervised learning techniques

that are available.

## B. Unsupervised Learning

Discrepancies and irregularities, such as fraud or broken equipment, may be identified via the application of unsupervised learning. Additionally, client groups who exhibit similar behaviors can be grouped together for the purpose of a marketing campaign. Comparatively, it is the antithesis of controlled learning. In this particular circumstance, there is no data that has been supplied.

It is the task of the coder or the algorithm to discover the underlying structure of the raw data, locate hidden patterns, and figure out how to represent the data when they are supplied with a restricted amount of symbols that do not have any representation or labels attached to them. The term "unlabeled data" refers to this kind of information that is used for educational purposes.

Consider the following scenario: we have a number of different pieces of information, and we need to classify them into a few different categories. It is possible that we are not aware of the particular criteria that are used to calculate the rating. The objective of unsupervised learning algorithms is to classify a selected dataset into a predetermined number of categories in the most effective manner possible.

Calculations based on individual learning are very useful techniques for assessing data and recognizing patterns and trends. Inputs that are comparable are often grouped together into

coherent categories using them. Some examples of individual learning computations include the K-means algorithm, the Random Forests algorithm, and the Hierarchical clustering algorithm.

## Semi-supervised Learning

This kind of learning is referred to as semi-supervised learning, and it occurs when certain learning tests are labeled but others remain unlabeled to the learner. For training purposes, it makes use of a significant quantity of data that has not been annotated, whereas for testing purposes, it makes use of a limited quantity of data that has been labeled. In circumstances where it would be too costly to acquire a completely labeled dataset, but where it would be possible to progressively annotate a tiny portion of the dataset, semi-supervised learning is used.

## C. Reinforcement Learning

Here learning data gives input with the goal that the framework acclimates to dynamic conditions so as to accomplish a specific goal. The framework assesses its exhibition dependent on the input reactions and responds in like manner.

### I. Supervised Learning Algorithms K-Nearest Neighbour Algorithm

The K-nearest neighbors (KNN) approach is a kind of supervised machine learning methodology that may be used for a variety of tasks, including classification and regression prediction.

KNN is categorized as a non-parametric learning technique since it does not make any assumptions about the underlying data. Additionally, it is characterized as a lazy learning algorithm because it does not have a separate training phase and uses all of the data that is available for training throughout the categorizing process.

K-nearest neighbors, often known as KNN, is an algorithm that makes predictions about the values of new data points by using the idea of closeness. This indicates that a new data point will be assigned a value that is determined by the degree to which it is similar to the points that are included in the training set. Following these steps will allow us to get an understanding of how it operates:

First, a dataset is necessary before any algorithm can be executed. This is the first stage. As part of the first step of the KNN process, it is essential to arrange and integrate the data from the preparation and the test.

Second Stage: The next step is to choose the value of K, which may be determined by looking at the data points that are most similar to the example. The letter K may stand for any number. The third stage consists of carrying out the following procedures at each point in the test

information:

To determine the distance between the test data and each row of the training data, you must use one of the following approaches to compute the distance: Whether it be the Manhattan, Hamming, or Euclidean distance. It is the Euclidean approach that is used the most often for computing distance.

**3.2**–Now, based on the distance value, sort them in ascending order.

**3.3**–Next, it will choose the top K rows from the sorted array.

**3.4**–Now, it will assign a class to the test point based on most frequent class of these rows.

**3.4**–Now, it will appoint a class to the test point dependent on the most successive class of these columns.

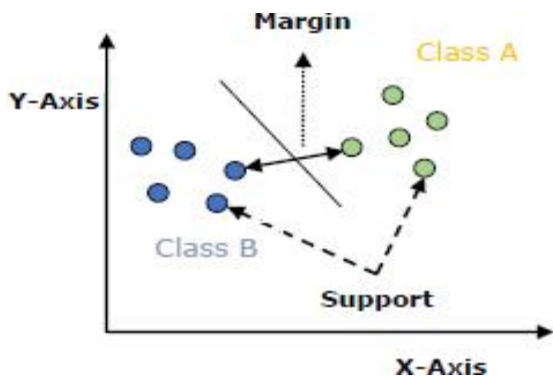
**Stage 4–End**

### **Support Vector Machines**

SVMs, which stand for support vector machines, are supervised machine learning algorithms that are very adaptable and perform exceptionally well in classification and regression problems. When compared to other machine learning techniques, support vector machines (SVMs) provide a performance approach that is distinctive. In recent times, they have acquired prominence as a result of their capacity to manage a wide variety of continuous and categorical data examples.



When applied to a multidimensional space, a support vector machine (SVM) model depicts several classes in a hyperplane. With the intention of reducing the number of mistakes that occur, the hyperplane will be produced in an iterative manner by using the Support Vector Machine (SVM) technique. Support Vector Machines (SVM) are designed to partition datasets into classes by locating the hyperplane that maximizes the margin. This is the primary objective of SVM.



**Fig.1.SupportVectorMachines**

The data points that are closest to the hyperplane are the ones that make up support vectors. For the purpose of defining the process of isolating a line, these data points may be used. As may be seen in the picture that follows, a hyperplane is a plane or space that divides a large number of objects that belong to distinct classified categories.

The term "margin" refers to the distance that exists between two lines and the data points that belong to different classes. It is often defined as the distance in the opposite direction from the line to the support vectors. This definition is commonly used. On the other hand, a little advantage is seen as an unfavorable advantage, whilst a considerable advantage is thought to be an advantage that is beneficial. The basic objective of Support Vector Machines (SVM) is to divide datasets into classes in order to locate an optimum separating hyperplane, which is referred to as the Maximum Margin Hyperplane (MMH). The two actions that are listed below are able to accomplish this goal:

In the beginning, the Support Vector Machine (SVM) will produce hyperplanes in a sequential manner in order to successfully separate the classes. • After that, it will choose the hyperplane that effectively divides the classes.

➤ **LogisticRegression**

Linear Regression isn't constantly fitting on the grounds that the data may not fit a straight line yet in addition the straight line seems can be more prominent than a line with a slope of 0. Thus

,they surely can't be utilized as the likelihood of

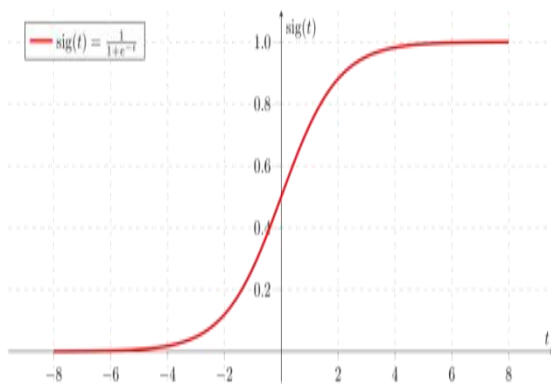
event of the objective class. Under these circumstances logistic regression is used. Instead of fitting data into a straight line, logistic regression uses a logistic curve.

Simple Logistic Regression

Output = 0 or 1, Hypothesis = >

$Z = WX + B$ ,  $\Theta(x) = \text{sigmoid}(Z)$

### Sigmoid Function



**Fig.2. Sigmoid Activation Function**

The anticipated value of 'Y' has a tendency to converge towards 1 as the value of the variable 'Z' becomes closer and closer to infinity. The value of 'Z' will steadily approach 0 as it gets closer and closer to negative infinity, and the predicted value of 'Y' will do the same.

For the purpose of classification, the supervised learning approach known as logistic regression is used. In order to provide predictions about the probability of a target variable, it finds use. Therefore, the dependent variable, which is also often referred to as the goal variable, is binary,

which means that it can only take on two unique classifications respectively.

To put it another way, the dependent variable is dichotomous, which means that it can only take on two values: one, which represents success or satisfaction, and zero, which represents

discontentment or rejection.

Through the use of this model, one may compute and make predictions about the likelihood of Y being equal to 1, where Y is a component of X. The scientific foundations serve as the basis for this strategy. In contrast to a mathematical model, a logistic regression model determines the likelihood that Y is equal to 1 by using the values of X as independent variables. This is in contrast to the mathematical model. Implementing this machine learning technology is not only simple, but it can also be used to a wide variety of categorization jobs. For the reason why we want to demonstrate our position

There are several varieties of logistic regression.

In common parlance, the term "strategic regression" refers to a kind of regression analysis that is comprised of two distinct sets of objective components. Additionally, it is able to generate projections for two additional categories of target variables, which is a significant advantage. Logical regression may be broken down into the following classification groups since there are so many

different categories to choose from. Is it a parallel or a binomial expression?

Within the context of this specific design, one of the dependent variables may only take on one of two potential combinations: either 1 or 0. On the basis of these characteristics, it is possible to demonstrate progress or disappointment, success or failure, and other outcomes that are comparable.

### **Multinomial**

In this form of structure, the subordinate variable may have at least three potential categories that are not organized in any particular order or categories that do not have any numerical significance. There is a possibility that these parameters relate to "Type A," "Type B," or "Type C," for instance.

"

### **Ordinal**

Within the context of this kind of categorization, the subordinate variable may be divided into a minimum of three potentially structured categories or categories that have a numeric value. It is possible to classify these criteria as "poor," "great," "generally excellent," or "superb," with scores of 0 for "poor," 1 for "great," 2 for "usually excellent," and 3 for "superb."

A strategic relapse model makes a prediction about the likelihood of Y being equal to 1 as a component of X. This prediction is expressed in numerical terms. It is a fairly straightforward machine learning technique that may be used to solutions for a wide variety of classification issues.

## **RegressionModels**

When it comes to regression, the Binary Logistic Regression Model is a straightforward approach that entails making predictions about a binary outcome variable that may take on just two potential values: either 1 or 0.

When the dependent variable may have at

least three different unordered categories, such as categories that do not have any quantitative importance, the multinomial logistic regression model presents itself as a viable kind of regression analysis. This model is used in situations where the dependent variable may contain such categories.

## **NaiveBayes**

### **TheBayesRuleandNaïveBayesClassification**

Given the chance of event X happening, which is known from the preparation dataset, the Bayes Rule is a strategy that is used to calculate the probability of an event Y occurring given the current probability of event X occurring. What happens when Y has a large number of classes? First, we determine the likelihood of each category of Y, and then we choose the category that has the greatest value of probability.  $P(X/Y)$  is the conditional probability of event X given event Y. It is equal to the probability of the intersection of events X and Y divided by the likelihood of event

Y, shown as  $P(X \cap Y)/P(Y)$ . In other words, the conditional probability of event X is equal to the probability of event Y. The likelihood of evidence given the result, which is known from the data used for training.

$P(Y/X)$  is a formula that reflects the conditional probability of Y given X. This probability is derived by dividing the likelihood of the intersection of X and Y by the probability of X. The formula is denoted by the number  $P(Y/X)$ .

Predicted for the Test Data: Naïve Bayes calculations are a classification strategy that is based on the application of Bayes' theorem. This technique assumes that all of the predictors are independent of each other. The assumption is that the existence of one element in a category is not reliant on the presence of another element in the same category. To put it another way, this is the assumption.

In Bayesian analysis, the primary goal is to ascertain the posterior probabilities, which may be thought of as the probability of an occurrence taken into consideration certain observable features, denoted by the notation  $(i | i)$ . In order to provide a quantitative representation of this idea, we may make use of Bayes' theorem, which is as follows: The equation may be stated as follows:  $P(L|features) = P(L)P(features|L) / P(features)$ . In this context, the posterior probability of a class is denoted by the expression  $(f | p)$ . The starting probability of a class is denoted by the symbol  $(i) i$ .

One definition of the word "likelihood" is "the probability of a marker occurring given a particular class." The phrase "likelihood" is represented by the symbol " $| \square$ ."

It is the initial probability of a pointer that is referred to as the "features" component.

### **Random Forest**

Random forest is a method for supervised learning that may be used for classification as well as

regression applications due to its versatility. In spite of this, the majority of the time it is used for classification issues. It is common knowledge that trees are the building blocks of a forest, and that a forest with a larger number of trees is more resistant to natural disasters. The arbitrary random forest approach involves the construction of decision trees based on data samples, followed by the extraction of predictions from each tree. By way of a voting procedure, it eventually chooses the solution that is the most appropriate. The performance of an ensemble approach is superior than that of a single decision tree because it decreases the likelihood of overfitting by averaging the results.

### **RandomForestAlgorithm**

- Step 1—First, start with the choice of random samples from a given dataset.
- Step 2 – Next, this calculation will build a choice tree for each example. At that point, it will get the forec

astoutcomefromeachchoice tree.

- Step3–Inthisprogression,castingaballotwil lbe

performedforeachanticipatedoutcome.

- Step4–Atlast,selectthemostcastedaballotf orecastresultasthefinal prediction result.

**PATIENTDATA SET**

For the cardiovascular dataset that was obtained from Kaggle, a total of 1090 examples were gathered, each of which had ten features. In order to determine whether or not a person has cardiovascular disease, the "Cardiovascular" attribute is usejd. A number of "1" indicates that the sickness is present, while a value of "0" indicates that the condition is not present.The values of the attributes that make up the cardiovascular disease data set are shown in Table I.There are 539 examples of people who do not have cardiovascular illnesses included in the dataset, and there are 551 cases of those who do have cardiovascular diseases.

**Table1:CardiovascularDataSet**

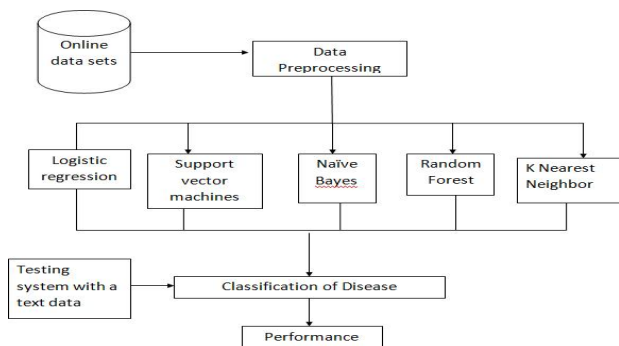
SerialNu mber	Attribute	Remarks
1	ID	IDnumber
2	Age	in Days
3	Gender	1-women,2-Men

4	Height	InCentMeter
5	Weight	KiloGrams
6	Systolic BloodPressure	SystolicBloodPress ure
7	Diastolic BloodPressure	DiastolicBloodPres sure
8	Cholesterol	1-Normal ,2- AboveNormal, 3- WellAboveNormal
9	Glue	1-Normal ,2- AboveNormal, 3- WellAboveNormal
10	Smoke	Whether Patient SmokesorNot
11	Alco	BinaryFeature
12	Active	BinaryFeature
13	Cardiovascular	TargetVariable

**3. PROPOSEDTECHINQUE**

The key goals of this research are to develop a technique that is capable of producing the best effective Machine Learning algorithm for the prediction jofj cardiovascular disease. Several different machine learning algorithms have been studied, and their distinct performance characteristics have been compared across the board.

.Selection



**Fig .3. Architecture diagram of cardiovascular disease prediction system**

After obtaining the kidney dataset from Kaggle, we conducted an analysis on it. Following the guidelines outlined in section IV, we have taken into consideration a total of thirteen factors that were included in the cardiovascular dataset. There are a total of 1090 tuples in the dataset, with 551 instances of persons suffering from cardiovascular

disease and 539 instances of those who do not suffer from cardiovascular sickness.

### A. Pre-processing and Transformation

The cardiovascular dataset is set up in Comma Separated Document format (CSV) from Excel File. Different things required are the expulsion of right qualities for missing records, copy record to evacuate pointless information field, standard information position, adjust information in a convenient way and so on. The considered cardiovascular dataset do not have any missing data values for different attributes.

### Performance Evaluation

We will present the performance evaluation of a number of different machine learning methods, which will include metrics such as instances that were correctly classified, instances that were

incorrectly classified, the kappa statistic, Mean Absolute Error (MAE), Root Mean Square Error (RMSE), Relative Absolute Error, and Root Relative Square Error. For a number of machine learning algorithms that we have reviewed, we are going to compute the True Positive Rate, the False Positive Rate, the Precision, the Recall, the F-Measure, and the Confusion Matrix.

### PERFORMANCE MEASURES FOR CLASSIFICATION

The performance metrics that are listed below may be used by an individual in accordance with their own criteria for the request and error-prone module. A classification model's accuracy may be evaluated with the use of a tool called the confusion matrix. This tool compares the predicted class labels of a given dataset with the actual class labels of the dataset. occurrences that are successfully categorized are appropriately represented by the words TP (True Positive) and TN (True Negative), while occurrences that are wrongly classified are represented by the terms FP (False Positive) and FN (False Negative). The Matrix of Mixed Beliefs Make sure that the instance is correctly classified as either a true positive (TP) or a true negative (TN). Cases should be misclassified. The text is not given in any way. The term "true positives" refers to the situations in

which the classifier correctly recognized positive cases of cardiovascular disease.

The term "false negatives" refers to the situations in which positive cardiovascular tuples are incorrectly identified as negative tuples.

True negatives are occasions in which the classifier correctly detects cardiovascular tuples that are contrary to the hypothesis being tested.

The term "false positives" refers to occurrences of cardiovascular disease that

were found to be negative but were yet wrongly labeled as positive.

The term "false negatives" refers to the situations in which positive cardiovascular tuples are incorrectly identified as negative tuples.

Given below is an example of a confusion matrix that depicts positive and negative tuples simultaneously:

**TableII:ComponentsofConfusion Matrix**

PredictedClass

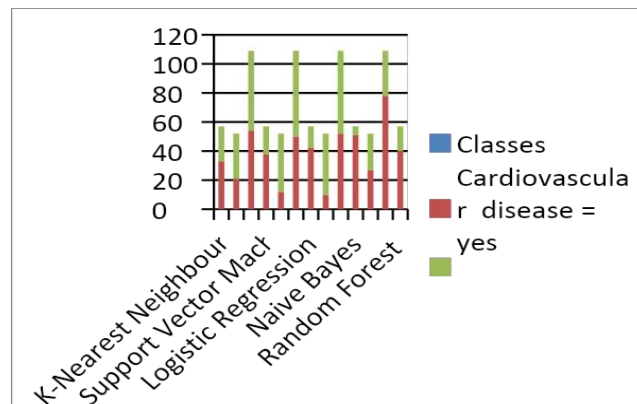
		Yes	No	
ActualClass	Yes	TruePositives(TP)	FalseNegatives(FN)	P
	No	FalsePositives(FP)	TrueNegatives(TN)	N
		PComplement	NComplement	P+N

A confusion matrix for positive and negative cardiovascular tuples for the considered dataset is as follows

**TableIII:ConfusionMatrix ofVariousAlgorithms**

Name of the algorithm	Classes	Cardiovascular disease = yes	Cardiovascular disease = no
-----------------------	---------	------------------------------	-----------------------------

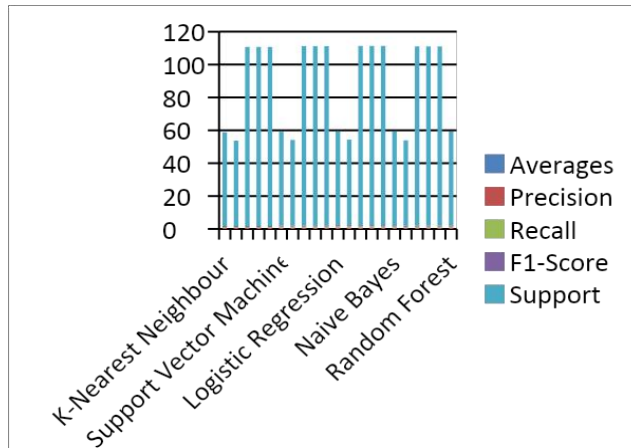
K-NearestNeighbour	Cardiovascularisease=yes	33	24
	Cardiovascularisease=no	21	31
	Total	54	55
SupportVectorMachines	Cardiovascularisease=yes	38	19
	Cardiovascularisease=no	12	40
	Total	50	59
LogisticRegression	Cardiovascularisease=yes	42	15
	Cardiovascularisease=no	10	42
	Total	52	57



**Fig.4. Graphical Presentation of various algorithms**



The accuracy of a classifier on a given test set is the percentage of test set tuples that are correctly classified by the classifier. That is,



**Fig.5. Comparison of Micro, Macro, and Weighted Average of various algorithms**

**Table IV: Accuracy Measure for Cardiovascular Disease Dataset**

Name of the Algorithm	Correctly Classified instances (%)	Incorrectly Classified instances (%)
K-Nearest Neighbour	58.71	41.28
Support Vector Machines	71.55	28.44
Logistic Regression	77.06	22.93
Naive Bayes	69.72	30.27
Random Forest	74.31	25.68

**Table V: Accuracy Measure for Cardiovascular Disease Dataset**

Name of the Algorithm	Kappa Statistics	Mean Absolute Error
K-Nearest Neighbour	0.17	0.41
Support	0.43	0.28

VectorMachines		
LogisticRegression	0.54	0.22
NaïveBayes	0.38	0.55
RandomForest	0.48	0.25

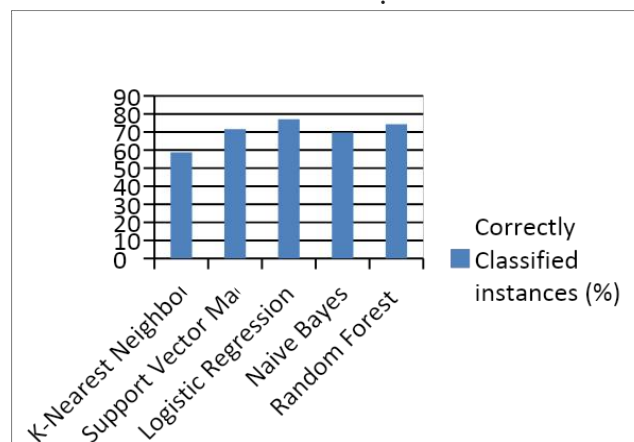
**Correctly and Incorrectly Classified Instances:**

The accurate categorization of instances is determined by the total number of true positives and true negatives that are included within the tuples of the cardiovascular dataset. Similar to the previous point, the term "erroneously categorized cases" refers to the total amount of false positives and false negatives that are present in cardiovascular datasets. The accuracy of the cardiovascular data may be calculated by dividing the total number of cases that have been correctly categorized by the entire number of instances of cardiovascular data.

The accurate categorization of instances is determined by the total number of true positives and true negatives that are included within the tuples of the cardiovascular dataset. Similar to the previous point, the term "erroneously categorized cases" refers to the total amount of false positives

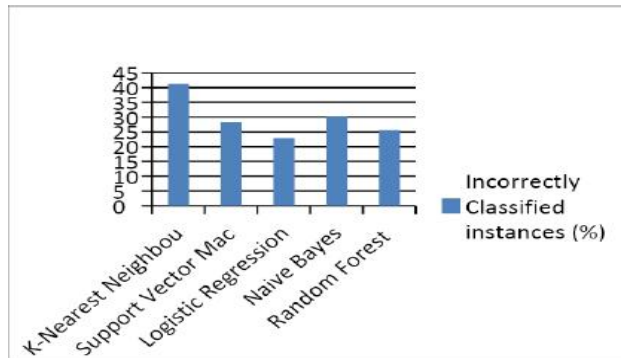
and false negatives that are present in cardiovascular datasets. The accuracy of the cardiovascular data may be calculated by dividing the total number of cases that have been correctly categorized by the entire number of instances of cardiovascular data.

The accurate categorization of instances is determined by the total number of true positives and true negatives that are included within the tuples of the cardiovascular dataset. Similar to the previous point, the term "erroneously categorized cases" refers to the total amount of false positives and false negatives that are present in cardiovascular datasets. The accuracy of the cardiovascular data may be calculated by dividing the total number of cases that have been correctly categorized by the entire number of instances of cardiovascular data.



**Fig.6. Comparison of correctly classified instances for various algorithms**

**Fig.7. Comparison of incorrectly classified Instances for various algorithms**



**Table VI: Accuracy Measure for Cardiovascular Disease Dataset**

Name of the Algorithm	Root Mean Squared Error	Relative Absolute Error (%)	Root Relative Square Error (%)
K-Nearest Neighbour	0.64	82.74	58.46
Support	0.53	51.48	40.27

## DISCUSSION AND RESULTS

For the purpose of this study, we used Machine Learning Algorithms to a dataset of cardiovascular disease in order to make predictions about patients who suffer from chronic cardiovascular sickness and those who are healthy. These predictions were based on the data of each characteristic for each individual patient. We wanted to take into consideration a number of

different layout models and determine which one was the most effective. The inquiry was carried out with the assistance of five different computations, which were K-Nearest Neighbour, Support Vector Machines, Logistic Regression, Naive Bayes, and Random Forest. The results of the five different comparison approaches are shown in tables 6, 7, and the eighth table respectively.

**Table VII: Accuracy Measure for Cardiovascular Disease Dataset**

Name of the Algorithm	Correctly Classified instances(%)	Incorrectly Classified instances(%)
Logistic Regression	77.06	22.93

When compared to the other five techniques that were examined, Logistic Regression demonstrated the highest level of accuracy, with a correct classification rate of 77.06% and the lowest number of wrong classifications, which was 22.93%.

In terms of the estimation of indicators, the results of the estimations of Mean Total Error (MAE), Root Mean Square Error (RMSE), Relative Absolute Error (RAE), and Root Relative Square Error (RRSR) revealed that Logistic Regression indicators had the lowest values (MAE = 0.22) (RMSE = 0.47, RAE = 45.96%, RRSE = 32.47), followed by estimations from other algorithms.

**CONCLUSION**

In the field of healthcare, the use of information mining systems for predictive analysis is of utmost importance since it enables us to diagnose illnesses at an earlier stage and, as a result, helps us to save lives by anticipating therapies that are successful. We employed a wide variety of machine learning algorithms in

this investigation, such as K-Nearest Neighbor, Support Vector Machines, Logistic Regression, Naive Bayes, and Random Forest, to determine whether or not patients were suffering from chronic cardiovascular failure and to identify those who did not have this illness. The results of the re-enactment demonstrated that the Logistic Regression classifier had outstanding performance in terms of prediction, with a high degree of accuracy and taking a minimum amount of time to carry out.

**REFERENCES**

- 1)"The Atlas of Heart Disease and Stroke",[online].[http://www.who.int/cardiovascular\\_diseases/resources/atlas/en/](http://www.who.int/cardiovascular_diseases/resources/atlas/en/)
- 2)J. S. Rumsfeld, K. E. Joynt, and T. M. Maddox, "Enormous information examination to improve cardiovascular consideration: guarantee and difficulties", Nature Reviews Cardiology, Vol.13, No.6,

- pp.350, 2016.
- 3) W. Dai, T. S. Brisimi, W. G. Adams, T. Mela, V. Saligrama, and I. C. Paschalidis, "Forecast of hospitalization because of heart sicknesses by managed learning techniques", *International Journal of Medical Informatics*, Vol.84, No.3, pp.189–197, 2015.
  - 4) A.U.Haq, J.P.Li, M.H.H.Memon, S. Nazir, R. Sun "A Hybrid Intelligent System Framework for the Prediction of Heart Disease Using Machine Learning Algorithms", *Mobile Information Systems*, Volume 2018.
  - 5) S. Krishnan J. Geetha S, "Forecast of Heart Disease Using Machine Learning Algorithms", *First International Conference on Innovations in Information and Communication Technology*, 2019.
  - 6) K.G.Dinesh, K. A. garaj, K.D.Santhosh, V. M. eswari. "Forecast of Cardiovascular Disease Utilizing Machine Learning Algorithms", *2018 International Conference on Current Trends towards Converging Technologies*, 2018
  - 7) A. Golande, P. kumar T, "Coronary illness Prediction Using Effective Machine Learning Techniques", *International Journal of Recent Technology and Engineering*, ISSN: 2277-3878, Volume-8, Issue-1S4, June 2019.
  - 8) Prasad, P. Anjali, S. Adil, N. Deepa, "Coronary illness Prediction utilizing Logistic Regression Algorithm utilizing Machine Learning", *International Journal of Engineering and Advanced Technology*, ISSN: 2249 – 8958, Volume-8, Issue-3S, February 2019
  - 9) Y. Khourdifi, M. Bahaj, "Coronary illness Prediction and Classification Using Machine Learning Algorithms Optimized by Particle Swarm Optimization and Ant Colony Optimization", *International Journal of Intelligent Engineering and Systems*, Vol.12, No.1, 2019.