# Performance Analysis of Supervised Learning Techniques on Heart Disease Prediction

**Ronakkumar Ashokbhai Modi, S Govinda Rao**

*Abstract-- In this Era, Internet is utilized on substantial scale and each field, for example, Health care, Economic, Feedback accumulation and different applications. In these estimation investigations, prior individuals used to give their criticism about the films, item, administrations and so forth, these things they have referred. This input freely accessible for upcoming orientations. In this work, Prediction of coronary illness is testing viewpoint looked by specialist and particularly in medical clinics they gather input from their patients. The Performance of conclusion examination is essential errand for machine to get yield in type of criticism for example Positive or Negative input. Sentiment Analysis and forecast of coronary illness is primary rule of medication and emergency clinics just as specialists. Machine learning calculations assume imperative job in this area.in this zone of research work to build up a product which helps the machine learning calculations to take choice with respect to both forecast and evaluate the assumption examination. The principle goal of this examination is foreseeing coronary illness of a patient utilizing machine learning calculations. Coronary illness is significant malady in social insurance industry. Along these lines, it is troublesome errand to foresee infection, in this work, execute managed machine learning strategies which gives better comprehension of heterogeneous forecast model and help to discover best symptomatic for medicinal services framework.*

*Index Terms: Sentiment Analysis, Machine Learning, NLTK, Naïve Bayes, Logistic Regression, Linear Model.*

## I. INTRODUCTION

Machine Learning is a field of programming building. It gives the "PC system to ability to "Learn" (for instance legitimately upgrade execution on a specific endeavor) with data, without being unequivocally adjusted. The name machine learning was established in 1959 by Arthur Samuel. It is resolved of AI. Advance the examination in machine learning of "Precedent Reorganization" and "Computational Learning" Theory in AI. It explores the examination and improvement of estimations that can be gain from and make conjecture of data, for instance, figuring's beat following completely "Static Program"(Static models are less requesting to build and test).

.

**Revised Manuscript Received on April 05, 2019.**
  **Ronakkumar Ashokbhai Modi,** Department of CSE, Gokaraju Rangaraju Institute of Engineering and Technology (GRIET), Hyderabad, India
  **Dr S Govinda Rao**, Department of CSE, Gokaraju Rangaraju Institute of Engineering and Technology (GRIET), Hyderabad, India

The system incorporates for progressing in the direction of "Data Breach", "optical character acknowledgment",

"Figuring out how to rank" and "Computer vision". It firmly identified with "computational Statistics". (Which likewise center on forecast making using PC).

The system of estimation of perceiving and organizing suppositions conveyed in a touch of substance, especially in order to choose if the writer's attitude in the direction of a particular point, things, certain, undesirable, or reasonable. Communicating the feelings and emotions with the assistance of words makes people one of a kind. These emotions are known as the opinions and the way toward breaking down these announcements is known as the Sentiment Analysis. [1]

Coronary illness is the real reasons for death among the general population now-a-days. As per World Health Organization (WHO) estimation about 12 million of death happens because of coronary illness in request to decrease the danger of coronary illness, expectation ought to be finished. A portion of the properties which anticipate the coronary illness are age, sex, chest torment type, family ancestry, ECG perusing, cholesterol, glucose level.

### A. Supervised machine Learning Techniques

In this system comprehensively exploited to order reason. Here the procedure of representative and beginners are initially prepared through example information. That data just remained relegated for classes, prototypes are tried via giving few example information contributions doing for characterization dependent upon earlier preparation. That execution of the representative is estimated by precision of an order.

In SML the learning calculation is given marked model data sources, where the names show the ideal yield. [2] SML itself is composed of

- **Classification: -** where the output is qualitative, and
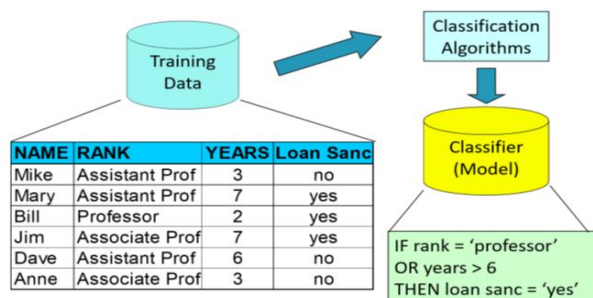- **Regression/Prediction**: - where the output is quantitative.



Fig 1: Classification of Training dataset

Typical Applications

Medical diagnosis, Credit approval, Target marketing, Treatment effectiveness analysis, Treatment viability investigation, Climate, entertainment, sports, etc.

Methods in Classification and choosing the best.

Here, approximately used classification algorithms such as,

- K—Nearest Neighbor
- Decision Trees
- Naïve Bayes
- Support Vector Machines

**B. Supervised Learning with Python**

In Supervised learning, we begin with bringing in dataset containing preparing traits and the objective characteristics. The Supervised Learning calculation will gain proficiency with the connection between preparing models and their related target factors and apply that educated relationship to arrange completely new contributions (without targets).

To delineate how administered learning functions, we should consider a case of anticipating the characteristics of an understudy dependent on the quantity of hours he contemplated. [4]

Mathematically,

$$Y = f(X) + C$$

Here,

**F** will be the connection between the imprints and number of hours the understudy arranged for a test.

**X** is the INPUT (Number of hours he arranged)

**Y** is the yield

**C** will be arbitrary mistake**.**

A definitive objective of the supervised learning calculation is to foresee Y with the greatest exactness for a given new information X. There are a few different ways to actualize Supervised Learning Techniques.

In view of the given datasets the machine learning issue is arranged into two sorts, Classification, and Regression.

"On the off chance that the given information has both info (preparing) qualities and yield (target) values, at that point it is a Classification issue". "If the dataset has continuous numerical values of attributes without any target labels, then it comes under Regression problem".

**C. Classification**

Consider the case of a therapeutic scientist who needs to investigate Heart sickness information to foresee which one of every three explicit medications a patient ought to get. This information examination assignment is called Classification, where a model or classifier is built to foresee class marks, for example, "Treatment A," "Treatment B" or "Treatment C."

Characterization is a forecast issue that predicts the clear-cut class marks which are discrete and unordered. It is a two-advance procedure, comprising of a learning step and a characterization step [3].

**D. Regression**

Regression is ordinarily named as choosing associations between something like two elements. e.g., deliberate the anticipate compensation and perspective on the specified data X. [5]

Regression Models:

Here, approximately used Regression models such as,

- Linear Regression
- Logistic Regression
- Polynomial Regression

## II. LITERATURE SURVEY

### A. Performance Analysis of Supervised Machine learning Techniques for Sentiment analysis

The purpose of this paper is loped at the accompanying: to recognize the qualities of the film surveys and classify them as Positive or Negative. This is on the grounds that the watchers dependably hope to invest their valuable energy and cash for viewing the great nature of motion pictures. That is the reason they used to take assistance from the input frameworks by experiencing the criticisms posted by the watchers who as of now have watched that films, and in the event that they get progressively positive criticisms about that motion picture, they like to watch that motion picture[1].

All over it furthermore happens that because of such countless concerning a film, the watcher might be skip scrutinizing all of the explanations, they want to display up or propose the watchers the dimension of negative and positive contributions around the relating film. So, the people can save their hour and energy of the watcher and besides takes the fitting result. In this investigation, research the film analyses and endeavors to group of the overview category. With the objective of the proposed method, the customers about the overview categories and at whatever points need to exhibit the number negative and positive contribution to respects to the picked movie(s). [3]

### B. Heart Disease Prediction System Using Supervised Learning Classifier

Cardiovascular contamination remains the best purpose behind passings worldwide and the Heart Disease Prediction toward the starting time is criticalness.

The examination of affliction is a basic work in medication. The social protection industry accumulates colossal proportion of human administrations data and after that they are mined to discover covered information for reasonable essential authority. Cardiovascular ailment is a kind of veritable prosperity gambling and visit happening disease. Cardiovascular ailments suggest any contamination that impacts the cardiovascular structure [5].

### C. Performance Analysis of unsupervised machine learning Techniques for network traffic classification

K-Means Clustering Technique:

K-means calculation is notable unsupervised grouping calculation. In this preparation stage it arbitrarily produces K group. These bunches speak to a given convention. It segments the N entity enter into K bunches. Contribution of calculation is N stream of numbers. Additionally, calculation is an esteem C show sum of groups. [6] In the arrangement stage this haphazardly produced

bunches further order the obscure traffic. As per the similitude (Euclidean separation) all new article are doled out to their most comparative groups.

After task of new item again focal point of all groups is determined. The procedure proceeds until every one of the items are allocated to their closest Neighbor bunches. The nearest centroid of info entity and produced items are determined by Euclidean separation as appeared.[6]

$$Dist(i,c) = \sqrt{\sum(ij - cj)2}$$

From above formula, i is an information variable of n measurements and C is group of n measurements. Yield k implies make it the centroid of C bunches and every group displays the convention Means is straightforward and preparation period is quick.

### D. Heart disease prediction system based on hidden naïve Bayes classifier

Coronary ailment is condition that impacts the heart. As the "coronary supply courses slim, circulatory system to the heart can back off or quit, initiating chest torment, heart attack". Detecting coronary sickness need to significantly talented and qualified specialists. There is a mix of various requests is to isolate capable information from epic proportions of data is called "Data Mining". There are many uses to remove human administrations learning for clinical essential authority and make hypothesis of immense restorative data.

Classification is an unavoidable issue which is utilized some operations as well as discover obscure examples. Analysts are concentrating on planning proficient arrangement calculations for substantial informational collections. Characterization framework will help doctors to inspect a patient, regardless of whether quiet is probably going to show some kindness ailment or not.[7]

Langseth and Nielsen projected Hidden naive Bayes in 2005. Shrouded parent is made and concealed gullible Bayes for every element that joins the impacts from every single component. Concealed nave Bayes exhibits noteworthy execution than further conventional arrangement calculations.

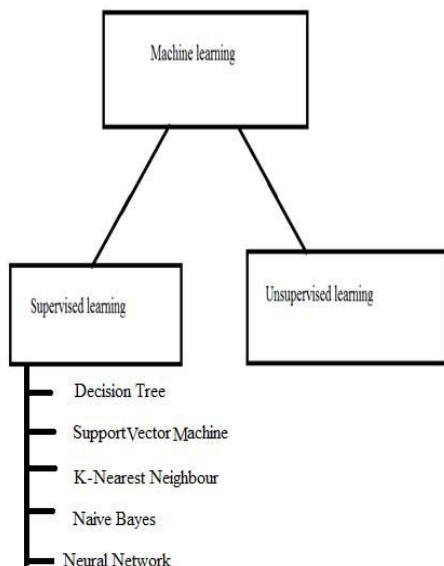### E. Performance Analysis of Various Supervised Algorithms on Big Data

*Fig 2. Algorithms of ML*

## III. PROPOSED METHODOLOGY

### A. Gathering Health Disease Data Sets

Accumulation of informational indexes is the essential occupation for any sort of estimation investigation inquire about, luckily there are some Health Disease informational indexes openly accessible over web. [1]

### B. Cleaning the Data Sets

Heart Disease informational collection comprises of alphanumeric, extraordinary characters and unknown characters. Which is helpful for our classifier, that is the reason subsequent to gathering the informational collections, we have attempted the informational index cleaning technique. Where we used to make the informational collections free of every single undesirable substance. Also, presently the cleaned informational collections are prepared for the subsequent stage which is grouping the audits accessible in the informational collections. [1]

### C. Information Classification

Supervised machine learning strategies used to take named information, which as of now have classified in accessible modules. That is the reason we have to dole out name, for example, "Positive" or "Negative" to the surveys as per their attributes.

### D. K-nearest Neighbor classifier (KNN)

K-Nearest Neighbor (KNN) is basic, languid and nonparametric classifier. KNN is favored when every one of the highlights are nonstop. KNN is likewise called as case-based thinking and has been utilized in numerous applications like example acknowledgment, measurable estimation. Characterization is gotten by recognizing the closest Neighbor to decide the class of an obscure example. KNN is favored over other order calculations because of its high intermingling rate and straightforwardness. [10] Figure 1 show closest Neighbor characterization. KNN arrangement has two phases

1) Discover the k number of cases in the dataset that is closest to case S

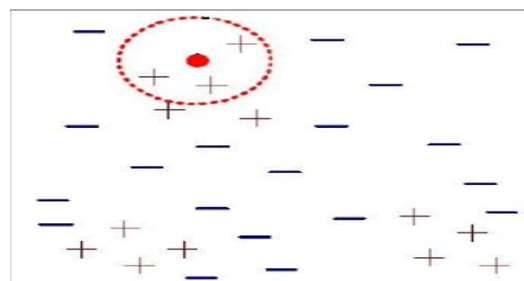2) This k number of cases at that point vote to decide the class of case S

*Fig 3 Clustering of dataset*

## IV. ABOUT THE DATASET

The data set utilized for this work is from "Data. world" store in which the coronary illness dataset is utilized. The dataset has 210 occurrence and 8 properties. These 8

properties are the generally consider factors for the coronary illness forecast. Despite the fact that it has 210 cases of which just 10 are finished and the rest of the lines contain missing qualities and expelled from the analysis.

The Data of patients recorded by number of times of chest agony and age in years. In dataset there are 8 qualities utilized in this framework, with 5 alphabetic and 3 numerical traits. [10]

Approximately, the essential parameters are variable parameters and that should be analyzed for at regular intervals the lock (greatest pulse accomplished), circulatory strain (mm Hg), serum cholesterol in (mg/dl), and electrocardiographic outcome. In genuine world, information isn't constantly finished and on account of the medicinal information, it is in every case genuine. To evacuate the quantity of irregularities which are related with information we use Data pre-handling.

## V. ALGORITHMS/TECHNIQUES

### A. K—Nearest Neighbor (KNN)

In the learning step, the characterization show assembles the classifier by breaking down the training set. In the characterization step, the class names for given information is anticipated. The dataset attributes and their related class marks inside investigation part into a training set and test set. [10]
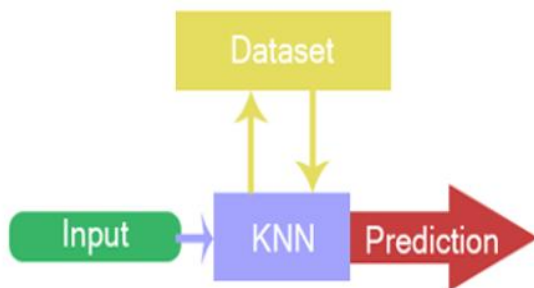


*Fig 4: Flow chart of KNN*

Our proposed technique means to improve the execution of KNN classifier for malady forecast. Calculation for our proposed technique is appeared as Algorithm 1. [7]

Stage 1: Input: Heart Disease informational index

Stage 2: Output: Classification of informational collection into patients with coronary illness and typical.

Stage 3: Input the informational collection

Stage 4: Apply pre-preparing strategies Fill in missing qualities

Stage 5: Discard excess highlights

Stage 6: Apply (KNN) on Predominant highlights

Stage 7: Measure the execution of the KNN demonstrate.

**Executing KNN in Scikit-Learn on Health dataset to order the kind of Heart Disease dependent on the given information.**
Initial steps, to apply our machine learning calculation we have to comprehend and investigate the given dataset. In this precedent, we use Heart illness dataset which is imported from the scikit-learn bundle. [8]

Pip install ~m pandas
pip install ~m matplotlib
pip install ~m scikit-learn
### K-Nearest Neighbors in scikit-learn
KNN depends on knowledge by communication affair, so that, by brightness of a submitted test attributes and creating attributes those are equal to that. The created attributes are depicted by n attributes. Every attribute addresses a point in a n-dimensional space. Consequently, every planning attributes are secured in n-dimensional precedent space.
At that moment when we give an ambiguous attribute, a k-Nearest Neighbor classifier checks the model space for the k created attributes those are closest to the ambiguous attributes.
These k created attributes, that are the k nearest Neighbors of the ambiguous attributes.
Presently, we give import KNN classifier from sklearn and apply to our information which at that point groups the Heart Disease dataset. [8]

## VI. ANALYSIS AND EXPERIMENTAL RESULT

Every one of the tests are directed in this work by using CORE I5 processor with 2 GHz processor with one Tera Bite Hard disk with minimum eight GB RAM in MS-Windows OS. Python 3.7 with pythons Scikit-Learn library is used in the Python Program for finishing our future work. [1]

Examination outcome acquire by different classifiers by various datasets and present the classified result in figure 5. from fig 5 it very well may be plainly seen that execution of KNN classifier from sklearn and apply to our data which by then gatherings the Heart Disease dataset.

*Table I: Heart Disease Dataset*

|       | age       | rest_bpress | max_heart_rate |
|-------|-----------|-------------|----------------|
| Count | 209.00000 | 209.000000  | 209.000000     |
| Mean  | 47.966507 | 133.660287  | 137.574163     |
| Std   | 8.054133  | 17.432960   | 23.876797      |
| Min   | 28.000000 | 92.000000   | 82.000000      |
| 25%   | 42.000000 | 120.000000  | 120.000000     |
| 50%   | 49.000000 | 130.000000  | 140.000000     |
| 75%   | 54.000000 | 140.000000  | 153.000000     |
| max   | 66.000000 | 200.000000  | 188.000000     |

Here, the table 1 demonstrates the separate age, rest_bpress, max_heart_rate and their precise outcomes. Fig 5 likewise can demonstrate the diagram of positive and negative infection. Likewise, it delineates the contrast between two diagrams and plotting the dataset values on x-pivot and y-hub.
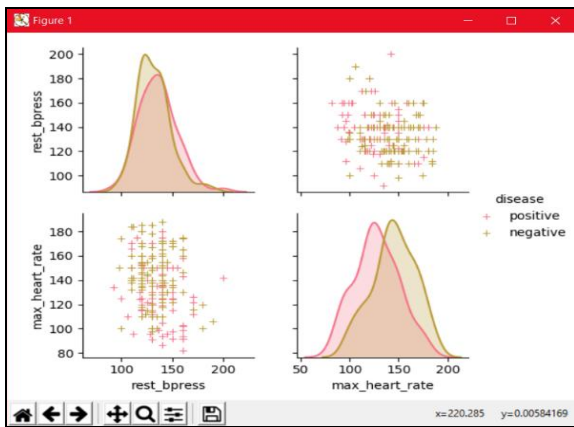
*Fig 5: Result of classification Dataset in Positive & Negative*

## VII. CONCLUSION

This paper tended to the forecast of coronary illness dependent on KNN. Our methodology utilizes KNN as a classifier to decrease the misclassification rate. This paper additionally explores KNN based element choice measure to choose few highlights and to enhance the order execution. The outcomes recommend that proposed methodology can altogether enhance the learning exactness. From reproduction results, it is inferred that KNN based component determination is vital for arrangement of coronary illness. This model helps the doctors in a productive forecast of illnesses with dominating highlights. In future, we need to incorporate outfit classifiers with KNN to build up a choice emotionally supportive network for early finding of coronary illness and furthermore might want to think about GA and PSO for coronary illness set.

## REFERENCES

1. BiswaRanjanSamal, Mrutyunjaya Panda, Human Being Character Analysis from Their Social Networking Profiles. A Semi supervised Machine Learning Approach, (IJCSIS) International Journal of Computer Science and Information Security, Vol. 14, No.5, May 2016
2. Scikit-learn: Machine Learning in Python, Pedregosa et al., JMLR 12, pp. 2825-2830, 2011.
3. TaiwoOladipupoAyodele Types of Machine Learning Algorithms, New Advances in Machine Learning, Yagang Zhang (Ed.), InTech,2010, DOI: 10.5772/9385.Availablefrom:
4. M.A. Jabbar, "Heart Disease Prediction System using Associative Classification and Genetic Algorithm", ICECIT, pp 183-192, Elsevier, vol 1(2012)
5. R. Chitra and Dr.V. Seenivasagam "Heart Disease Prediction System Using Supervised Learning Classifier" Bonfring International Journal of Software Engineering and Soft Computing, Vol. 3, No. 1, March 2013 1
6. Hardeep Singh "Performance Analysis of unsupervised machine learning Techniques for network traffic classification" Assistant Professor Lovely Professional University.2015 Fifth International Conference on Advanced Computing & Communication Technologies.
7. MA Jabbar and Shirina samreen "Heart disease prediction system based on hidden naïve bayes classifier" Vardhaman college of Engineering Hyderabad, India.
8. Athira Unnikrishnan "Performance analysis of various Supervised Algorithms on Big Data" International Conference on Energy, Communication, Data Analysis and Soft Computing (ICECDS-2017).
9. Siqian Chen, Jie Yang, Yun Gu "IMAGE SENTIMENT ANALYSIS USING SUPERVISED COLLECTIVE MATRIX FACTORIZATION" Institute of Image Processing and Pattern Recognition Shanghai Jiao Tong University, Shanghai China
10. Prediction of heart disease using k-nearest neighbor and particle Swarm optimization. Jabbar MA*Vardhaman College of Engineering, Hyderabad, India ISSN 0970-938X.

## UTHORS PROFILE

**Ronakkumar Modi,** pursuing the M. TECH in Software Engineering from Gokaraju Rangaraju Institute of Engineering & Technology(GRIET), Telangana, India. I have done B.E in Information Technology from Gujarat Technological University. I have research interests including Machine Learning, Networking, and Programing Languages. I won first prize in cloud computing seminar which was organized by CSI (Computer Society of India).

**Dr S Govinda Rao** is working as Professor in CSE, Gokaraju Rangaraju Institute of Engineering and Technology (GRIET), Hyderabad. Completed his PhD From JNTUK and M.Tech From Andhra University. He is around 14 years of teaching experience. He published around 35 research publications. He is life member of MIE