# Disease Prediction Using Data mining Classification Algorithm

[1]Loda Prasanna Kumar,[2]T Naga Raju
[1]Student of M.Tech (CSE) and Department of Computer Science Engineering,
[2] Asst.Prof, Department of Computer Science and Engineering,
Kakinada Institute of Engineering and Technology,AP.

## Abstract

Knowledge discovery in databases has built up its prosperity rate in different noticeable fields, for example, e-business, advertising, retail and medical. Medical Data mining has extraordinary intensity for investigating the outside of anyone's ability to see designs in the separate medical Data collections. This venture plans to arrange the disease Data indexes and deliver the reports in light of their manifestations. The Data collections are ordered by utilizing multilayered encourage forward neural systems. The datasets for the diseases are procured from UCI, an online vault of vast Data collections.

**Keywords:** Clustering, Classification, Data mining tools, Disease prediction, Health care.

## I. Introduction

Data Mining is the discovery of knowledge in databases. Strategies of Data mining help to process the Data and transform them into valuable data. Prediction results from Data mining are valuable in different fields like Business Intelligence, Bioinformatics, Healthcare Management, Finance and so forth. Medical field has wide sum and also assortment of Data for handling and there exist numerous testing assignments. This field requires exact and auspicious mannered diagnosis which can spare numerous patients life. Data mining methods assumes an imperative part in medicinal services investigation. Early location and exact results are achievable by doctors utilizing Data mining calculations. Diverse calculations will be utilized for shifted disease diagnosis. In light of the Data utilized the exactness and execution additionally shift. We are focusing on Healthcare databases, which have a tremendous measure of Data yet in any case, there is an absence of successful investigation apparatuses to find the concealed knowledge. In this study we display an outline of the ebb and flow look into being done utilizing the DM systems for the diagnosis and guess of different diseases, featuring basic issues and compressing the methodologies in an arrangement of educated lessons. Data Mining utilized as a part of the field of medical application can misuse the shrouded designs display in voluminous medical Data which generally is left unfamiliar. Data mining systems which are connected to medical Data incorporate affiliation manage mining for finding incessant examples, prediction and order.

## II. Related work

Utilization of Predictive Data Mining in Clinical Diagnosis In disease diagnosis, [7] contrasted control based Repeated Incremental Pruning with Produce Error Reduction (RIPPER), Decision Tree (DT), Artificial Neural Networks (ANN) and Support Vector Machine (SVM) based on Sensitivity, Specificity, Accuracy, Error Rate, True Positive Rate and False Positive Rate, and 10-overlay cross approval to gauge the unprejudiced gauge of these prediction models. SVM display was pronounced the best classifier for Cardiovascular Disease (CVD) diagnosis. [8], looked at DT, k-Nearest Neighbor (kNN), ANN, and grouping in light of bunching in the prediction of coronary illness, and hereditary calculation (GA) to enhance the prescient capacities of the picked models. [1], showed how to execute a confirmation based clinical master arrangement of a Bayesian model to identify coronary corridor disease. The Bayesian was considered to have extensive favorable position in managing a few missing factors contrasted with coordinations and straight relapse models. In the diagnosis of Asthma with master framework, [10] completed a similar examination of machine learning calculations, for example, Auto-affiliated Memory Neural Networks (AMNN), Bayesian systems, ID3 and C4.5 and discovered

AMNN to perform best as far as calculations productivity and exactness of disease diagnosis. In an investigation of Phospholipidosis, [1] utilized structure-action connections (SAR) to analyze kNN, DT, SVM and simulated safe frameworks calculations prepared to distinguish drugs with Phospholipidosis possibilities and SVM created the best predictions took after by a Multilayer Perceptron fake neural system, strategic relapse, and kNN. In the diagnosis of Chronic Obstructive Pulmonary and Pneumonia diseases (COPPD), [2] analyzed neural systems and simulated insusceptible frameworks. Likewise, [3] utilized DT, Naïve Bayes, and Neural Networks to broke down coronary illness. The neural system calculation was found to anticipate coronary illness with most noteworthy precision. 2.2 Erythemato-Squamous Diseases (ESD) Diagnosis AnExtreme learning machine and Artificial Neural Network was proposed by [4] to better distinguish the differential issue of Erythemato-Squamous skin diseases. A help vector machine (SVM) in light of arbitrary subspace (RS) and highlight determination was utilized to tackle differential issue of Erythemato-Squamous diseases [5]. A Catfish Binary Particle Swarm Optimization (CatfishBPSO), Kernelized Support Vector Machines (KSVM), and Association Rule include determination technique was utilized to analyze Erythemato-squamous diseases [6]. The AR-CatfishBPSO-KSVM display picked up gauge execution precision of 99.09% contrasted with Support vector machines and affiliation rules - multilayer perceptron (AR-MLP). Another model in light of Adaptive Neuro-Fuzzy Inference System (ANFIS) was utilized as a part of the identification of Erythemato-squamous diseases [6]. They guarantee ANFIS has a few possibilities in identifying the Erythemato-squamous diseases. Another diagnosis show in light of Support Vector Machine (SVM) with a novel cross breed highlight choice strategy - Improved F-score and Sequential Forward (IFSFFS) was utilized on five irregular preparing test parcels of the Erythemato-squamous diseases datasets from University of California Irvine (UCI) machine learning archive database [7]. The outcome demonstrates that SVM-based model with IFSFFS accomplished the ideal characterization precision. [8], utilized Rough-Neuro half breed technique to accomplish precise Erythemato-squamous diseases diagnosis. The technique joins

Rough sets Johnson Reducer for diminishment of important properties and fake neural system Levenberg-Marquardt calculation for the arrangement of the diseases. The model was asserted to have analyze the diseases at a precision of 98.8%. Twostage cross breed include choice calculations for diagnosing Erythemato-squamous diseases. It consolidate Support Vector Machines (SVM) as a grouping device, and the broadened Sequential Forward Search (SFS), Sequential Forward Floating Search (SFFS), and Sequential Backward Floating Search (SBFS), as inquiry methodologies, and the summed up F-score (GF) to assess the significance of each element. The two-arrange cross breed show was asserted to have accomplished better order precision when contrasted with accessible calculations for Erythemato-squamous diseases [9].

## III. Data Mining Techniques

Classification Techniques

Decision Tree:

Flow chart like tree structure, where each inward hub, signifies a test on a quality, each branch speaks to a result of the test, and each leaf hub holds a class mark. These territories basic and quick to comprehend and can without much of a stretch converted into arrangement rules.

Decision tree are developed by ID3, C4.5, CART (Classification And Regression Tree).

Bayesian Classification:

Bayesian classifiers are factual classifiers. They can anticipate class participation probabilities, for example, the likelihood that a given tuple has a place with a specific class.

Bayesian classification depends on Bayes' hypothesis

$$
\begin{aligned}
P(X|C_i) &= \prod_{k=1}^{n} P(x_k|C_i) \\
&= P(x_1|C_i) \times P(x_2|C_i) \times \cdots \times P(x_n|C_i).
\end{aligned}
$$

Fig. 1 Bayesian classification recipe

Run Based Classification:

The rulebased classifier learned model is spoken to as an arrangement of IFTHEN rules. An IFTHEN lead is a declaration of the frame:

A lead R can be evaluated by its scope and precision. Given a tuple, X, from a classlabeled informational collection, D, let ncovers be the quantity of tuples secured by R; ncorrect be the quantity of tuples effectively classified by R; and |D| be the quantity of tuples in D. We can define scope exactness R as

$$coverage(R) = \frac{n_{covers}}{|D|}$$

$$accuracy(R) = \frac{n_{correct}}{n_{covers}}.$$

Fig 1.2 Rule based order

Nearest Neighbor Method:

A method that characterizes each record in a dataset in view of a mix of the classes of the k record(s) most like it in a verifiable dataset (where k 1), Sometimes called the knearest neighbor procedure.

Artificial Neural Network:

Artificial neural network (ANN) is a machine learning approach that models human cerebrum and comprises of various counterfeit neurons.

Different Methods:

Languid learning, fuzzy logic and hereditary algorithms

**IV. Proposed System:**

This task means to arrange the disease informational collections and create the reports in light of their side effects. The informational collections are arranged by utilizing single class order procedure. The datasets for the diseases are gained from UCI, an online storehouse of expansive informational collections.

*MULTILAYERED FEED FORWARD NEURAL NETWORK:*

The most mainstream managed prepared counterfeit neural network is the multilayer perceptron network (MLP). It comprises of a few layers of computational

components. The non-intermittent adaptation of the network is additionally called feed forward neural network, since a yield can't impact its neuron esteems.

A feed forward network has a layered structure. Each layer comprises of units which get their contribution from the past layer units by assessing the underneath recipe in which the information units are increased by the synaptic weights. MLP comprises of info layer, shrouded layers and yield layer in which the qualities are passed from the information layer through concealed layers lastly at the yield layer required yield is gotten.

$$y = \varphi(\sum_{i=1}^{n} w_i x_i + b) = \varphi(\mathbf{w}^T \mathbf{x} + b)$$

Fig 2.calculation work

Where indicates the vector of weights, is the vector of data sources, is the predisposition and is the actuation work.

Network Design Parameters:

1)  Number of Input Nodes:

Info layer comprises of set of information units which is totally relying upon the application or the current information. Qualities typically falls under various classifications like parallel, ostensible and ordinal information.

2)  Number of Output Nodes:

The yield layer comprises of set of yield neurons which is totally in view of the application and furthermore acclimated to fall inside the range in light of the kind of initiation work picked.

3)  Number of center or shrouded layers:

The shrouded layers permit various possibly extraordinary mixes of sources of info that may results in high (or low) yields. Each progressive shrouded layer speaks to the likelihood of perceiving the significance of mixes of mixes.

4)  Number of Nodes for every Hidden Layers:

The more hubs there are the more noteworthy the quantity of various information mixes that the

network can perceive. It depends on the general guideline as take after:

a.  Number of shrouded neurons must be in the middle of information and yield neurons.

b.  Number of shrouded neurons = (input neurons + yield neurons)*2/3.

c.  Hidden neurons ought not be more noteworthy than 2(input neurons).

5)  Initial Connection Weights:

The weights on the information joins are instated to some irregular potential arrangement. Since the preparation of the network relies upon the underlying beginning arrangement, it can be imperative to prepare the network a few times utilizing diverse beginning stages. The underlying weights must be in the scope of 1 to 1.

6)  Initial Node Biases:

Hub inclination esteems confer a hugeness of the info mixes feeding into that hub. All in all hub predispositions are permitted to be adjusted amid preparing, yet can be set to specific esteems at network instatement time. Adjustment of the hub inclinations can be additionally permitted or refused.

7)  Learning Rate:

At each preparation step the network processes the bearing in which each predisposition and connection esteem can be changed to compute a more right yield. The rate of change at that arrangement state is likewise known. A learning rate is client assigned keeping in mind the end goal to decide how much the connection weights and hub predispositions can be altered in view of the alter course and change rate. The higher the learning rate (max. of 1.0) the speedier the network is prepared. In any case, the network has a superior shot of being prepared to a nearby least arrangement. A nearby least is a time when the network balances out on an answer which isn't the most ideal worldwide arrangement.

8)  Momentum Rate:

To help abstain from sinking into a nearby least, a force rate enables the network to conceivably skip through neighborhood minima. A background

marked by alter rate and course are kept up and utilized, to some extent, to push the arrangement past nearby minima. An energy rate set at the most extreme of 1.0 may bring about preparing which is very flimsy and accordingly may not accomplish even neighborhood minima, or the network may take an over the top measure of preparing time. In the event that set at a low of 0.0, energy isn't considered and the network will probably sink into a neighborhood least.

Enactment work:

Most units in neural network change their net contributions by utilizing a scalartoscalar work called an enactment work, yielding an esteem called the unit's initiation. But potentially for yield units, the actuation esteem is nourished to at least one different units. Initiation capacities with a limited range are frequently called squashing capacities. Probably the most regularly utilized enactment capacities are

1) Identity work

Clearly the information units utilize the personality work. Some of the time a consistent is increased by the net contribution to frame a straight capacity.
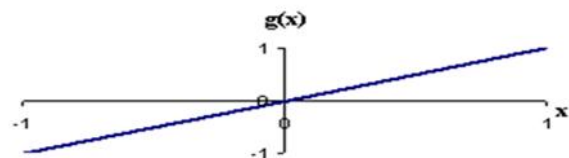


Fig 3. Identity work

2) Binary advance capacity Otherwise called edge capacity or Heaviside work. The yield of this capacity is constrained to one of the two esteems:

$$g(x) = \begin{cases} 1 & \text{if } ( x \geq \theta ) \\ 0 & \text{if } ( x < \theta ) \end{cases}$$

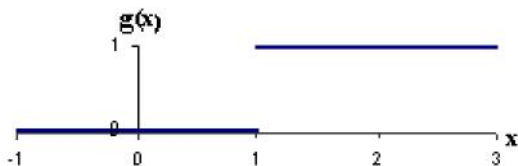This kind of function is often used in single layer networks.

Fig 4. Binary step function

3) Sigmoid function (Fig 2.4.4)

$$g(x) = \frac{1}{1+e^{-x}}$$

This function is especially advantageous for use in neural networks trained by back propagation; because it is easy to differentiate, and thus can dramatically reduce the computation burden for training and its output values are between 0 and 1
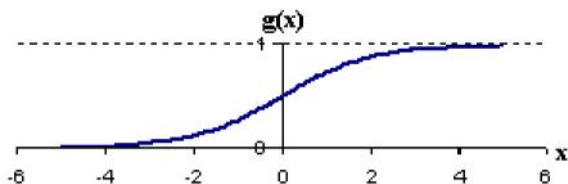


Figure 5. sigmoid function

4) Bipolar sigmoid function (Fig 2.4.5)

$$g(x) = \frac{1-e^{-x}}{1+e^{-x}}$$

This function has similar properties with the *sigmoid function.* It works well for applications that yield output values in the range of [1,1].
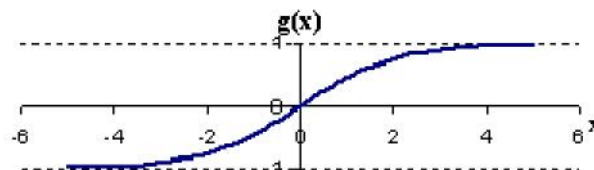


Fig 6. Bipolar sigmoid function

Initiation capacities for the concealed units are expected to bring nonlinearity into the networks. The reason is that an arrangement of straight capacities is again a direct capacity. Be that as it may, the nonlinearity makes multilayer networks so effective. The sigmoid capacities are the most widely recognized decision

5) Hyperbolic digression work:

The hyperbolic digression actuation work (TANH) enactment work is a typical initiation work for neural networks. The hyperbolic digression capacity will deliver positive numbers in the vicinity of 1 and 1. The hyperbolic digression initiation work is most helpful for preparing information that is likewise in the vicinity of 0 and 1. Since the hyperbolic digression enactment work has a subordinate, it can be utilized with inclination based preparing strategies. The hyperbolic digression initiation work is maybe the most well-known actuation work utilized for neural networks.
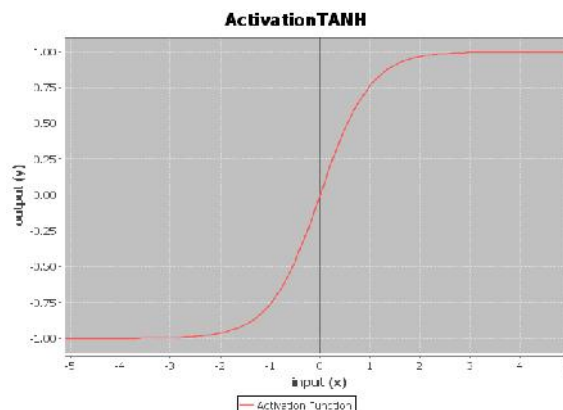


Fig 7.Tanh function graph

Neural network researches grow rapidly with the introduction of the back propagation (BP) algorithm for training the feed forward neural network. This training algorithm is applied for the research presented in this thesis, and the algorithm will be thoroughly explained.

## V. CONCLUSION

In this paper, we have proposed multilayered feed forward neural network to classify the data sets. According to our thesis back propagation method is suitable to classify the disease data sets and we can predict the disease based on the trained network. Our application can be used as a public application by deploying in network and also used in the hospital environment to provide external services to the patient.

Further enhancement can be done by adding the range of inputs and outputs. We can use encoding techniques to add more number of inputs. So that application will be able to classify more number of diseases.

### References

[1]Sara Khalid, David A. Clifton, Lei Clifton, and Lionel Tarassenko- A Two-Class Approach to the Detection of

Physiological Deterioration in Patient VitalSigns, With Clinical Label Refinement IEEE TRANSACTIONS ON INFORMATION TECHNOLOGY IN BIOMEDICINE, VOL. 16, NO. 6, MARCH 2016 .

[2] Dhanya P Varghese, Tintu P B,   A Survey on Health Data using Data Mining Techniques , International Research Journal of Engineering and Technology (IRJET), Volume: 02 Issue: 07, e-ISSN: 2395-0056, p-ISSN: 2395-0072, Oct- 2015.

[3] VahidRafe, RoghayehHashemiFarhoud,   A Survey on Data Mining Approaches in Medicine , International Research Journal of Applied and Basic Sciences, Vol 4 (1), ISSN 2251-838X, 2013.

[4] Lambodar Jena, Narendra Ku. Kamila, "Distributed Data Mining Classification Algorithms for Prediction of ChronicKidney-Disease", International Journal of Emerging Research in Management &Technology, Volume-4, Issue-11, and ISSN: 2278-9359, November 2015.

[5] Dr. S. Vijayarani, Mr.S.Dhayanand, "Kidney Disease Prediction using SVM and ANN algorithms", International Journal of Computing and Business Research (IJCBR), Volume 6, Issue 2, ISSN (Online):2229-6166, March 2015.

[6] T. Revathi, S. Jeevitha,   Comparative Study on Heart Disease Prediction System Using Data Mining Techniques , Volume 4 Issue 7, ISSN (Online): 2319-7064, July 2015.

[7] DevendraRatnaparkhi, Tushar Mahajan, Vishal Jadhav,   Heart Disease Prediction System Using Data Mining Technique , International Research Journal of Engineering and Technology (IRJET), Volume: 02 Issue: 08, e-ISSN: 2395 -0056, p-ISSN: 2395-0072, Nov-2015.

[8] K.Manimekalai,   Prediction of Heart Diseases using Data Mining Techniques , International Journal of Innovative Research in Computer andCommunication Engineering, Vol. 4, Issue 2, ISSN(Online):2320-9801, ISSN (Print):2320- 9798, February 2016.

[9] JyotiRohilla, PreetiGulia,   Analysis of Data Mining Techniques for Diagnosing Heart Disease", International Journal of Advanced Research in Computer Science and Software Engineering, Volume 5, Issue 7, ISSN: 2277 128X, July 2015.

[10] A. Oztekin, D. Delen and Z. J. Kong, "Predicting the graft survival for heart–lung transplantation patients: An integrated data mining methodology," International Journal of Medical Informatics (IJMI), vol. 78, no. 12, pp. e84-e96., 2009.