



Lung Cancer Prediction Using Machine Learning Technique Over Big Data

T. Thangarasan^{1*}, R. Keerthana², M. Arunkumar³, S. Ramya⁴,
M. Bommy⁵, V. Surendhiran⁶

Abstract

In current position, cancer disease is substantial menace to human life globally. 32 percent of people worldwide are affected by various types of cancer. But lung cancer depicts the highest ratio. Nowadays peoples are not having awareness about to detect the cancer in early stage. The survival rate of five year for lung cancer disease is 55 percent of the cases are affected most. However, only 14 percent of lung tumor cases are diagnosed at an early stage. For slight tumors the five-year survival rate is simply 3 percent. There are 4 stages in lung cancer. If we predict the disease in I and II stage, it is easy to cure by effortless operations. If it exceeds second stage, it may not be cured. So, diagnosing the cancer in earlier stage is the best solution to predict the patients from death. For that, the system uses the Decision Tree and K-Nearest Neighbor (KNN) Algorithms as preferred classification model. By using these algorithms, it becomes easier to diagnose the cancer in early stage. So, the survival rate of lung cancer patients becomes higher. The propound system analyze, calculate and compares the precision of Random forest, Naive Bayes and KNN and the preliminary result reveals that ID3 furnish better precision for cancer dataset. The input has been accessed only in numeric format. The algorithms also maintain key stuffs of the dataset, which are predominant for extracting performance, and so it may warrant the correct defense and effective preservation. This leads to protection of any extracting works that depends on the sequence of distances between objects, such as Random forest, Naive Bayes -search and classification, as well as many visualization techniques. In particular, it establishes a restricted isometric property, that is the tight leap on the shrinkage and enlargement of the original distances.

910

Key Words: Lung cancer, Decision Tree, KNN, ID3, Naïve Bayes

DOI Number:10.14704/nq.2022.20.8.NQ44098

NeuroQuantology 2022; 20(8):910-917

Introduction

Big data is a term which is used to manage and analyze the large quantity of data which is not able to deal by the traditional software systems. Every day there is billions of data are generated from various factors such as e-commerce websites, social media, hospitals etc. Big data plays an important role

to analyze these data with more effectiveness and provides the best result. Big data techniques are used to work with the effective performance of surgery strategies, other medical tests, and also to discover the relationships among very rushed medical, clinical and diagnosis data. In the field of health sector, the facility for doctors had introduced various data chassis with an enormous amount of

Corresponding author: T. Thangarasan

Address: ¹Assistant Professor, School of Computers, Assistant Professor, Madanapalle Institute of Technology & Science, Madanapalle, Andhra Pradesh, India. ²Assistant Professor, Department of Computer Science and Engineering, K.S.R College of Engineering, Tiruchengode, Namakkal, Tamil Nadu, India. ³Assistant Professor, Department of Electronics and Communication Engineering, Sengunthar Engineering College, Tiruchengode, Namakkal, Tamil Nadu, India. ⁴Assistant Professor, Department of Master of Computer Applications, M.Kumarasamy College of Engineering, Karur, Tamil Nadu, India. ⁵Assistant Professor, School of Computers, Madanapalle Institute of Technology and Science, Madanapalle, Andhra Pradesh, India. ⁶Assistant Professor, Department of Computer Science and Engineering, Excel Engineering College, Komarapalayam, Tamil Nadu, India
E-mail: thangarasant@mits.ac.in



information to regulate the patient information but regrettably, the effectiveness over hidden information discovering is less because of the data are not mined. In medical sector, the information mining is a term which means to overview and inspect the medical information to envision the condition of the patient's health. So the statistical analysis, machine learning, decision making and database techniques are applied to discover desiring pattern from healthcare data.

Table 1. Stages And Survival Rates Of Lung Cancer Patients In India

Lung Cancer Stages	Diagnosis Frequency	Five-Year Survival Rate (%)
I	15%	>70%
II	25%	35%-40%
IIIA	20%	25%-40%
IIIB	20%	4%-7%
IV	45%	<2%

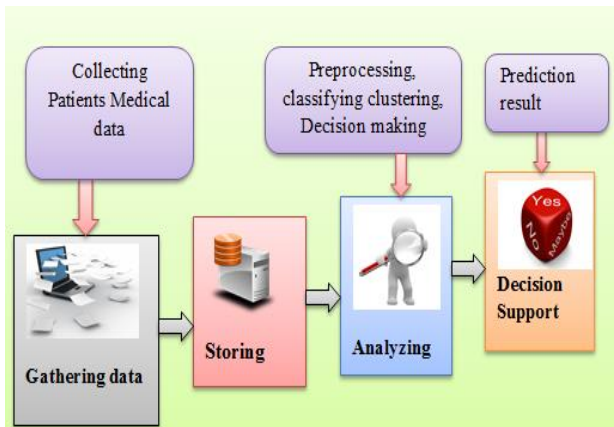


Fig.1. Pipeline process of big data analytics

Fig. 1 explains the schema of disease prediction processing duct which contains various steps of gathering data, storing, reviewing and decision support. Distinct methods of data mining are applied for evaluating these data after gathering enormous amount of health complaint data from various sectors. This analyzing process followed by data preprocessing, then feature selection and certainly it makes the prediction by applying machine learning which may do classification and clustering approaches for those healthcare data.

Lung Cancer Stages And Survival Rates

Lung cancer is the master cause of cancer deaths for both the gender in worldwide. Cigarette smoking, tobacco and genetics is the major risk factor for development of lung cancer. The prognosis of lung cancer is very much least because a doctor not able to find the disease until it reaches severe stage. Because the lung can functioned normally even it is partially damaged. Five-year survival rate in India is around 52% for early stage lung cancer that is localized in the lungs, but only around 3% in advanced, is an inoperable lung cancer.

The stage 1 lung cancer symptoms include lingering or bleeding cough, coughing up or blood, rib ache while breathing deeply, barking cough, raucousness, shortened breathing, wheezing, weakness and dullness, trouble in appetite and weight loss. The stage II lung cancer symptoms are persistent worseness, coughing up blood, shortness in breath, chest pain and back pain, continuous infections like pneumonia and bronchitis. The stage III lung cancer includes pain in the chest, pain when breathing, wheezing, whirling or high pitched sound upon inhalation or exhalation, persistent cough, cough with blood, blood in saliva and mucus, hoarseness or altered voice. The stage IV which is a final serious stage symptom may include a cough that will not get cleared and stacked, out coming of blood or rusty cold while coughing, chest rib ache which may frequently raucous with broad breathing, coughing, raucousness, sudden weight loss and troubles in appetite, compressed breath, tiredness or weak.

Fig. 2 describes the four stages of cancer cell (tumor) development in lung. Stage A is a baseline which describes the minute development of tumor near the lung. In 4 months, the tumor reaches the rib. In 5 months it starts to get grow well and in 8 months it starts to decay the whole lung.

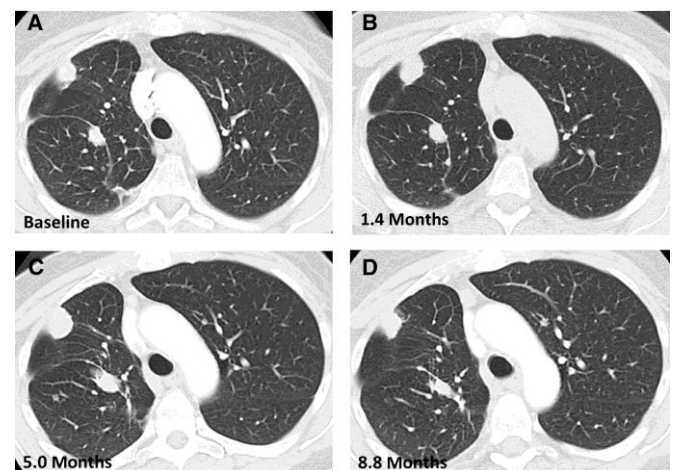


Fig. 2. Lung Cancer Stages



Related Work

There are various exploration that may have benefit of machine learning and extracting techniques are explained beneath.

The research of paper [Andrzej Skalski et al] developed the proposal as segmentation of lung for Computed Tomography data of a particular patient with Lung cancer. Using these regional conclusions, a wholly automated approach of Lung collocation has been introduced. The collocation is based on 73-feature vectors calculated for the lung sector which consist of 3D edge, region, orientation and spatial relational information. For more accurate classification, paper [BendiVenkataRamana et al] has implemented Naïve Bayes classifier and Support Vector Machines for accurate classification. The paper [Cybenko.G et al] has overviewed the classification and regression techniques and makes the comparison over them (SVM provided better accuracy- 79.32%) for blood cancer disease prediction. In paper [Emrana Kabir Hashi et al] system uses the classification models as Decision Tree and K-Nearest Neighbor (KNN) Algorithms. At the end, the proposed system which makes an analysis over it and correlate the certainty between C4.5 and KNN and the overviewed conclusion finalize that the C4.5 gives better certainty for diabetes mellitus diagnosis. For the medical database, the Pima Indians Dataset is used in this exploration. For the automatic prediction, paper [Han Sang Lee et al] proposed an automated process for identifying and fragmenting SRM (small renal mass) in contrast-enhanced CT images by using texture, pixel and context feature collocation and compared their result (accuracy- 77.45%). For classification of CT images of lung cancer, paper [Isabelle Guyon et al] has compared Multi segmented injection Forward Neural Network, J-48, Random Forest and Genetic programs are tested using ILPD (Indian Lung Patient Dataset) Data Set. The paper [Jankishran Pahariyavohra et al] uses, the covering methodology which is recently popularized and it offers an effortless and effective way to find the problem of selecting the variable, nevertheless of the chosen machine learning. Because, the machine learning is referred as a perfect black box and it lend to use of off-the-shelf machine learning software package. This approach consists in using the diagnosis and predicting valuation of a machine learning to evaluate the useful variable subset. In routine, one should define: (i) how to search all possible variable subset space; (ii) how to evaluate the prediction valuation of learning machine to

supervise the search and kept it. To predict the type II blood cancer, paper [John C. Platt et al] has make the comparison over the accuracy of C4.5 used to develop a decision tree. In paper [Rajeswari P et al], it introduces advanced algorithm for coaching the support vector machines: Sequential Minimal Optimization, or SMO.

Proposed Method

The projected KNN based classifier determines nearest values persistently from trained datasets and it works only with numeric and image vector of the Lung cancer dataset. The main advantage of this approach is the perfect predictive working based on symptoms and test diagnostic data of the Lung cancer patients. The proposed approach is used to detect the Lung cancer patients and their stage of cancer affected and the experimental application shows the conclusion of the efficiency of the prospective approach. With the combination to that for reviewing medical care data, important steps of extracting techniques like preprocessing data, change the missing values, feature selection, machine learning and decision support and also in addition naïve Bayes also used for the result accuracy are used in training dataset. Finally the random forest method has been executed on the training dataset of cancer for the classification process.

Step 1: In this proposed system, user (doctor, physician etc.) can insert the attributes like their symptoms, their fantastic doubted symptoms, x-ray images and CT images and also the disease value and these values are added to the system with the help of admin through internet. Then, by applying the machine learning technique, the final result after various decision making and pruning can be shown to user.

Step 2: On the server, there is an allowance for admin to load different disease datasets and suiting various extracting algorithms to train up the set of data. The inputs which are requested by users are assembled and refined on main server to predict the analysis result.

Step 3: For reviewing medical care data, the important steps of big data accessions like preprocess data, replace missing values, feature selection, machine learning, KNN and decision support are evolved for training the data. On the server there are numerous algorithms like ID3 and KNN gets accomplished over trained dataset and gets qualified to segment the data which are given for testing.



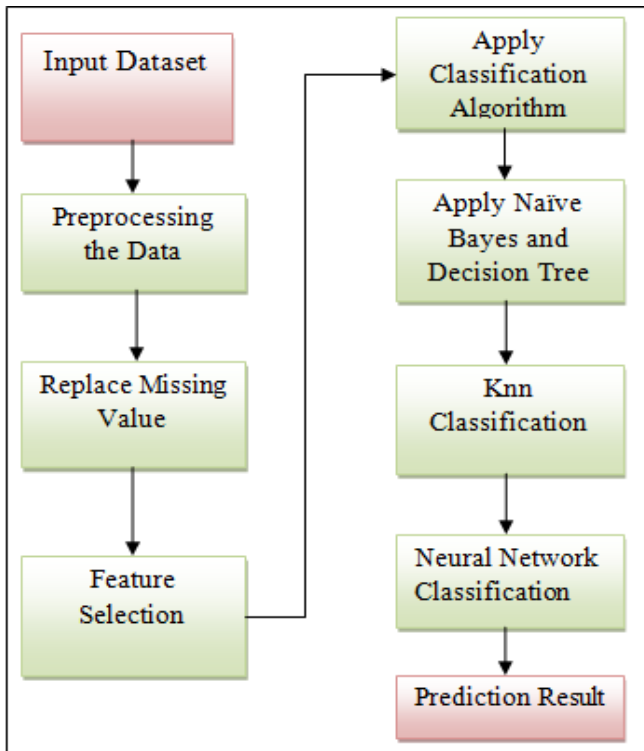


Fig. 3. Architecture Diagram

Methodology

In advanced system, the medical prediction has created on the basis of big data analyzing technique to detect the lung cancer disease. Decision tree, Random forest, Naïve Bayes and K-Nearest Neighbor (KNN) classifiers are used to make the system to train the dataset.

A. Naive Bayesian algorithm

On the basis of Bayes’ theorem the self-subsistence estimate between the predictors done by Bayesian classifier. Naive Bayes classifiers are a family of elementary feasibility classifiers on the basis of applying naïve Bayesian theorem. The anterior possibility $P(c/x)$, from $P(c)$, $P(x)$, and $P(x/c)$ has been thought out with the Bayes theorem. A self-subsistence value of another predictor has been generated by assumption of predictor(x) on the liable class(c) with the help of naïve Bayes classifier. The formula for predicting the tuple ‘x’ related to the class ‘c’ by naïve Bayesian classification is

$$P(c/x) = \frac{P(x/c) \cdot P(c)}{P(x)}$$

$P(c/x)$ is the rear possibility of class target concluder.

$P(c)$ is the earlier likelihood of a class.

$P(x/c)$ is the probability which is the concluder given class.

$P(x)$ is the prior probability of predictor.

B. Decision Tree

[Badr HSSINA et al] Decision tree is a more strong categorization technique. There are various techniques like ID3, C4.5, C5, J48, CART and CHAID algorithms are obtainable to estimate the dataset. Here, there is a continual value in this dataset, so we have chosen ID3 as our perfect classifier. By calculating the topmost information profit in all accredits we can find the decision points.

Pseudo-code for the Decision Tree (T) is declared beneath:

- Step 1: Compute Class Frequency(T);
- Step 2: if One Class or Few Cases
Return a leaf; Create a decision node N;
- Step 3: For Each Attribute A Compute Gain(A);
- Step 4: N. test = Attribute With Best Gain;
- Step 5: if N. test is continuous Find Threshold;
- Step 6: For Each T’ in the splitting of T
- Step 7: if T’ is Empty Child of N is a leaf else Child of N = Decision Tree(T’);
- Step 8: Compute Errors of N; Return N.

C. K-Nearest Neighbor (Knn)

KNN is superintend learning algorithm that finds and recollect radical data based on marginal distance from radical data to the K nearest neighbor. The Euclidean distance is used as a proposal to explain the togetherness.

Pseudo-code for the KNN classifier is stated below:

- Step 1: Input: $D = \{(x_1, c_1), \dots, (x_n, c_n)\}$ $x = (x_1, \dots, x_n)$ radical occurrence to be classified
- Step 2: For each labeled occurrence (x_i, c_i) Calculated $d(x_i, x)$
- Step 3: Order $d(x_i, x)$ from lowest to highest, $(i=1, \dots, N)$
- Step 4: Select the K nearest instances to x : $D_x K$
- Step 5: Assign to x the most frequent class in $D_x K$

D. Neural Network

Mostly the regression type classification and grouping in data science are done perfectly by artificial neural networks. The feature vectors are grouped into classes and it gets ready for receiving the input from the user and find the label which suits most. This can be worked with its efficiency to label the things, like customer types, images and music genres. In medical science department, these classifiers are frequently used to explore the health of the patients, establish them as regular, unsure, or defective.



The network is a consulting firm of artificial neurons, as same as neuron connection in human brain. It gets trained by experiencing various instances of every class and learns the likeness and contrast by making comparison. This is a routine which is similar as how the brain learnt, with repeated pattern that forms a better association over time. Visual inspection and k-means clustering of data from four aspects proposed seven various means of health. We can take color as an example for our prediction result.

Grey: "early" (starting run-in of the aspects)

Green: "normal"

Yellow: "suspect" (health seems to be crumble)

Black: "failure.S2", "failure. inner" (S3), or "failure.

roller" (S1 and S4)

Violet: "stage2" (secondary failure of S4)

Results And Discussions

A. Dataset

The major Centre for machine learning and artificial intelligent systems from university of California (UCI), Irvine used this dataset for predicting the cancer. There are 600 samples with 8 numerical valued attribute and 400 which gives the negative test and remaining 100 possess the positive instances in this dataset. The attributes which have been selected are explained briefly for cancer data analysis in TABLE II.

Table 2. Training Dataset

S.No	Training Dataset	Normal Image	Disease Image
1	150	55	95
2	300	74	226
3	450	116	334
4	600	124	476

Table 3. Test Dataset

S.No	Test Dataset	Normal Image	Disease Image
1	100	39	61
2	200	65	135
3	300	92	208
4	400	115	285

In d detection, sensitivity refers the fraud detection rate and it is defined as

$$\text{Sensitivity} = \text{TP}/\text{N}$$

False Alarm Rate: False alarm rate refers the section of actual negative instances which are predicted as positive instances and it is explained as

$$\text{False Alarm Rate} = \text{FP}/\text{N}$$

Result Comparison

The outcomes of final prediction using machine learning for lung cancer disease are discussed in this part. Here 600 instances of random chosen authentic set of data are used to train the machine

and other 100 precedents are used for testing the data. This prediction system gets matured using three familiar algorithm i.e. Random forest, ID3 and KNN algorithm. The accomplishment of these three algorithms is discussed beneath.

1) Performance Assessment of Random Forest algorithm: While Applying Random Forest classifier, TABLE IV has designed as the preliminary result of analysis. Random Forest classifier makes correct classification over 136 instances and made incorrect classification over 74 precedents. The certainty of perfectly classified precedent is 69.43% and imperfectly classified precedent is 30.57%.



Table 4. Random Forest Classifier

Random Forest	Prediction result		Accuracy	
	perfect classified precedent	Imperfect classified precedent	perfect classified precedent	Imperfect classified precedent
Training instances: 600 Testing instances: 210	136	74	69.43%	30.57%

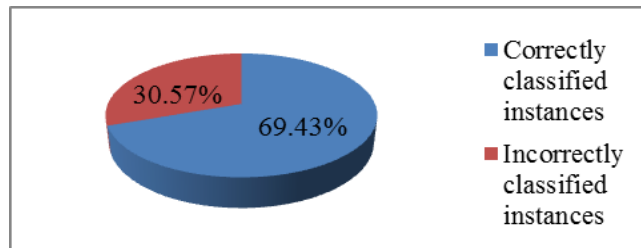


Fig.4. Accuracy chart of Random Forest Classifier

2) Performance Assessment of KNN algorithm: classification over 48 instances. The certainty of While applying KNN classifier, the prevision result perfectly classified precedent is 72.65% and provides, KNN classifier have the correct imperfectly classified precedent is 27.45%. classification over 162 instances and incorrect

Table 5. Knn Classifier

KNN	Prediction result		Accuracy	
	perfect classified precedent	Imperfect classified precedent	perfect classified precedent	Imperfect classified precedent
Training instances: 600 Testing instances: 210	162	48	72.65%	27.45%

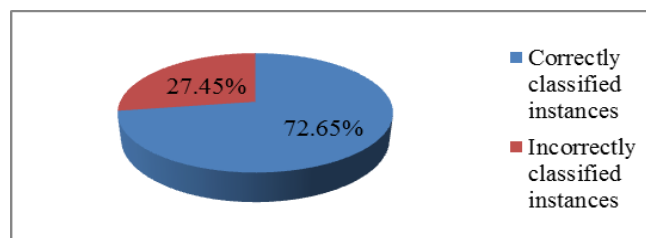


Fig.5. Accuracy chart of KNN Classifier

3) Performance Assessment of Decision tree (ID3) algorithm: After execution of system using ID3 classifier, the machine has classified 198 instances in a correct manner and 12 instances are classified in incorrect manner. Finally, the certainty of perfectly classified precedent is 96.62% and imperfectly classified precedent is 3.38%.



Table 6. Id3 Classifier

ID3	Prediction result		Accuracy	
	perfect classified precedent	Imperfect classified precedent	perfect classified precedent	Imperfect classified precedent
Training instances: 600 Testing instances: 210	198	12	96.62%	3.38%

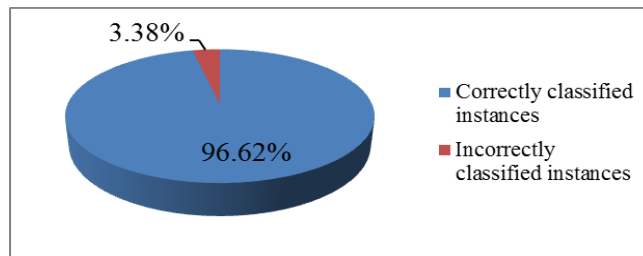


Fig.6. Accuracy chart of ID3 Classifier

4) Allegory of Classification certainty: The consummation of Machine learning system was overviewed with Lung cancer dataset using Random forest, KNN and ID3. In consonance with

demonstration results, TABLE IV denotes the consummation allegory of random forest, KNN and ID3 classifier upon its allegory split (80:20) model

Table 7. Comparison Of Random Forest, Knn And Id3.

Classifier	No. of instances		Accuracy
Random Forest	perfectly classified	136	69.43%
	Imperfectly classified	74	30.57%
KNN	Perfectly classified	162	72.65%
	Imperfectly classified	48	27.45%
ID3	Perfectly classified	198	96.62%
	Imperfectly classified	12	3.38%



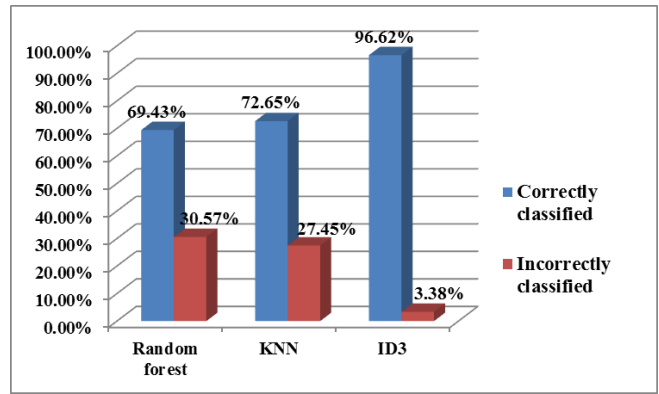


Fig.7. Graph measure Representation of consummation allegory

In Fig. 7, the graph measure representation of consummation allegory is explained from the graphical model which has been proved that ID3 algorithm performs better with higher allegory.

Conclusion

In this project, a novel system is proposed for predicting the lung cancer using big data classification technique such as Random forest, Naive Bayes search and ID3 classification. The outcome of this system is beneficial to the doctors, medical learners and also the patients to take right resolution about the diagnosis and prediction of the Lung cancer disease. Wide use of neural network as an imaging approach enables early detection of tumor formation in the lung. So, this prediction provides great convenience for eliminate only the cancer tumors with a necessary deadline promoting conservation of healthy lung parenchyma volume as possible. Finally, the project uses the K-Nearest Neighbor (KNN) and decision tree (ID3) Algorithms as supervised classification and result prediction model. And the calculation accuracy of Random forest, Naive Bayes and KNN are compared and finalized the accurate prediction.

References

Andrzej Skalski ;Jacek Jakubowski ; Tomasz Drewniak, "Lung tumor segmentation and detection on Computed Tomography data", Imaging Systems and Techniques (IST), 2016 IEEE International Conference on 4-6 Oct. 2016

Aneeshkumarand.A.S C.JothiVenkateswaran,"Estimating the Surveillance of Lung Disorder using Classification Algorithms", International Journal of Computer Applications (0975 - 8887), Volume 57- No.6, November 2012.

Badr HSSINA, Abdelkarim MERBOUHA," A comparative study of decision tree ID3 and C4.5", International Journal of Computer Applications Beni-Mellal, BP: 523

BendiVenkataRamana, Prof. M.Surendra Prasad Babu and Prof. N. B. Venkateswarlu, "A Critical Study of Selected Classification Algorithms for Lung Disease Diagnosis", International Journal of Database Management Systems (IJDMS), Vol.3, No.2 (2011), PP.101-11.

Cybenko.G, "Approximation by superpositions of a sigmoidal function", Mathematics of Control, Signals, and Systems, Vol.2 (1989), PP. 303-314.

Emrana Kabir Hashi, Md. Shahid Uz Zaman and Md. Rokibul Hasan, "An Expert Clinical Decision Support System to Predict Disease Using Classification Techniques", International Conference on Electrical, Computer and Communication Engineering (ECCE), February 16-18, 2017, Cox's Bazar, Bangladesh (2017).

Han Sang Lee ; Helen Hong ; Junmo Kim, "Detection and segmentation of small renal masses in contrast-enhanced CT images using texture and context feature classification", Biomedical Imaging (ISBI 2017), 2017 IEEE 14th International Symposium on 18-21 April 2017, ISSN: 1945-8452

Isabelle Guyon and Andr'eElisseeff," An Introduction to Variable and Feature Selection", Journal of Machine Learning Research 3 (2013) 1157-1182.

Jankishran Pahariyavohra, Jagdeesh makhijani and sanjay patsariya, "Lung patient classification using intelligence techniques ", International journal of advanced research in computer science and software engineering, Volume 4, Issue 2, Pages 295-299, 2013.

John C. Platt," Sequential Minimal Optimization: A Fast Algorithm for Training Support Vector Machines", Technical Report ,April 21, 2010.

Rajeswari P and SophiaReena.G," Analysis Of Lung Disorder Using Data Mining Algorithm", Global Journal Of Computer Science And Technology, Vol. 10 Issue 14 (Ver. 1.0) November 2010.

