

# A Review on Different Classifiers Techniques for Disease Prediction

Rupinder kaur ,student,Punjabi university Patiala,Punjab india  
Dr. Gaurav Gupta , Department of CE, Punjabi university Patiala, Punjab, india.

**Abstract :** A large amount of data is being collected from different locations. Using various data mining techniques useful data can be extracted from the large data repository. Various types of decisions in which this type of information can be used. Because large repository may be having relevant and irrelevant data. In all the cases the usage ability will be minimized. For having better relevant data there requires various types of data mining techniques. Data mining involves classification, it classifies the data into multiple classes. Based on the classification prediction model is developed. This prediction model can be used for the prediction of the heart attack for the patient based on different parameters. These parameters are collected for different patients. Parameters like Pregnancies, Glucose, Blood Pressure, Skin Thickness, Insulin, BMI, Heart attack Pedigree Function, Age, etc. can be useful for the prediction. Different classifiers like CART, RF, SVM, LDA, KNN are being applied on a large dataset.

**keywords** - CART, Random Forest, SVM, Linear discriminate analysis, KNN.

## I. INTRODUCTION

In recent time various types of services are emerging in the society. These services are related to the different fields of the society. Out of those fields major field is medical. India is a large country where large population resides. Various types of organized and unorganized medical facilities are available in the country. But due to the population explosion each facility remain scarce. To overcome and catalyst the growth in this part of the applications. Various researchers are involved which are growing with different researches so that the problem of scarcity of the resources can be catered without increasing the much cost.

Big Data is one of the major thrust area. That can solve the problem of this medical mismanagement. There will be a Big Data of the medical data. That is consisting of integration of various small city level data. Any company will provides processing ability which can process this large integration of the data. So that any patient data can be provided at the required place. Patient case history or we can say medical history can be recorded at each step.

These data are mainly stored and isolated in disparate local systems, and are underutilized in terms of data analysis and knowledge discovery. Data mining techniques have made computational infrastructure capable of handling such enormous information burst in a cost-effective way. Up to the present, most works focus on migrating healthcare IT system and data storage to the cloud platform rather than taking advantage of information hidden in the data. A typical healthcare Data Mining system has a hierarchical structure including system layer, control layer, and service layer. System layer constructs fundamental storage and computing environment using distributed computing resources, storage resources, and network resources. In control layer, system administrators control the load balancing, monitor system performance, and build programming environment. Finally service layer is responsible for providing large-scale healthcare services via real time management, privacy protection, and data analysis.

### 1.1 Types of Heart Diseases

*Type 1 Heart Diseases* is called blockage of the main arteries of the patient. Autoimmune, genetic, and environmental factors are involved in the development of this type of blockage. Diet also plays a major role in having such kind of blockages.

*Type 2 Heart Diseases* is called attacking on the heart. It is the last condition of the heart attack conditions. In the type 2 *Heart Diseases*, the heart usually produces some support the amount produced is not enough for the body's needs, or the body's cells are resistant to it.

### 1.2 Big Data

The amount of data being generated inside and outside each enterprise has exploded. The increasing volume and detail of information, the rise of multimedia and social media, and the Internet of Things are expected to fuel continued exponential data growth for the foreseeable future.

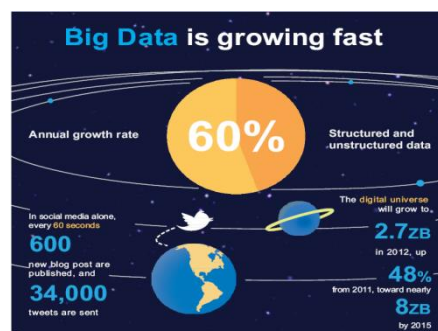


Fig. 1 Big Data Structure[2]

Two common sources of big data exist. First, there is organizational data, which, thanks to improved automation and broader access, is increasingly being shared. Organizational data includes emails, system logs, internal documents, business process events, and other structured, unstructured, and semi-structured data. It also includes social content, such as any blogs and wikis that are available within the organization.

Second, data comes from sources outside of the organization. Some of this extra-organizational data is available publicly and at no charge, some of it is obtained through paid subscription, and the rest is selectively made available by specific business partners

or customers. This includes information from social media sites and various other sources, including product literature; corporate information; health, life, and consumer websites; helpful hints from third parties, and customer complaints, such as complaints posted on the websites of regulatory agencies.

### 1.3 Big Data challenges in healthcare

Inferring knowledge from complex heterogeneous patient sources. Leveraging the patient/data correlations in longitudinal records. Understanding unstructured clinical notes in the right context.

- Efficiently handling large volumes of medical imaging data and extracting
- Potentially useful information and biomarkers. Analyzing genomic data is a computationally intensive task and combining
- With standard clinical data adds additional layers of complexity. Capturing the patient's behavioral data through several sensors.
- Various social interactions and communications.
- Take advantage of the massive amounts of data and provide right intervention to the right patient at the right time.
- Personalized care to the patient
- Potentially benefit all the components of a healthcare system Personalized care to the patient. i.e., provider, payer, patient, and management.

## 2. Literature survey

[1] This paper is based on dengue fever analysis by using the data mining techniques. In this model clustering techniques. The main aim to check the dengue patients or non-dengue patients. This model observed symptoms it analyzed the dengue fever mostly various on fever temperature followed by muscle and joint pain and red spots (August, No, & Kaur, 2017).

[2] This model is based on reduce the problem of leprosy spread is to find and predict the cause of leprosy through data mining tools and techniques. we use the data collection, data analysis, data sampling and t-test and clustering techniques. The model predictor importance show that the disease is widely affected according to area and age where gender and occupation does not affect much for the disease (Kaur & Gupta, 2015)

## 3. TECHNIQUES

### 3.1 CART

It is classification and regression tree for the machine learning based technique where various types of data items are being classified based on tree based technique. This type of decision tree based technique can be useful for the prediction based modeling technique. This decision tree represents the root element as the single input variable and a split point on that variable. The leaf node of the tree contain the output variable like yes and no. Further another parent node represents the variable then further the grant children having yes or no.

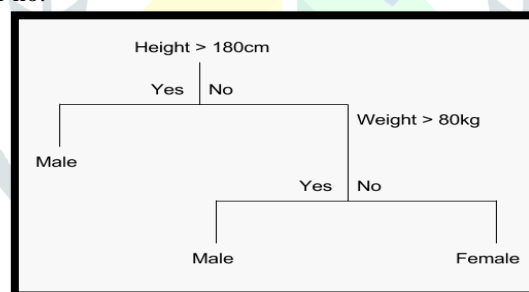


Fig. 2 CART based tree.

The values of the tree are stored into the file. This stored file will be having various rules making the system of prediction quite easy. For a new input the tree will be traversed from the root and move in downward direction. Based on the input appropriate decision can be derived.

### 3.2 RANDOM FOREST

Random forest is another classifier used for the prediction. Currently it is the efficient technique for classification of the heart disease data. It includes various small decision trees. Each decision tree helps in having decision for the specific classification. Higher is the number of tree higher will be the accuracy. Decision tree concept is more rule based system. Given a training dataset the decision tree algorithm will be used for deriving various trees. Once the decision tree are being established, the rules are applied on to the testing set. Suppose if the person want to predict whether person will like the movie or not. For the prediction purpose various types of decision tree will be required. Each decision tree will be having different types of lead actor in the film. Now for the person who has seen the movie and like the most has certain lead actor. Based on the decision tree if this movie has the same lead actor then the person definitely will like the movie else will not like the movie.

### 3.3 SUPPORT VECTOR MACHINE

It is another supervised learning technique. This is used for the classification and regression challenge. It is the classification algorithm. Where each item based on the feature values is put into the 3- dimensional space. N is the number of features. Each value of the feature being the value of particular coordinate. This way the different classes will be prepared on the basis of segregation of the data along the hyper plane.

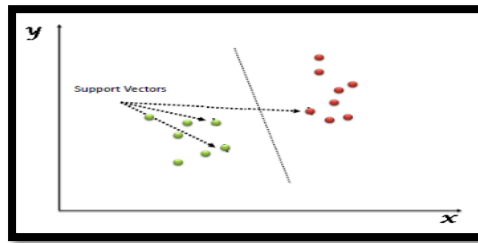


Fig. 3 SVM Plane

### 3.4 LINEAR DISCRIMINANT ANALYSIS

It is the supervised learning technique. it simply classify the whole dataset into two or more classes. The resulting combination may be used for linear classifier or may for the dimensionality reduction. LDA is closely related is ANOVA (analysis of variance), which also attempt to express one dependent variable as a linear combination of other features or measurements. LDA is also closely related to principal component analysis (PCA) and factor analysis in that they both look for linear combinations of variables which best explain the data. LDA explicitly attempts to model the difference between the classes of data. PCA on the other hand does not take into account any difference in class, and factor analysis builds the feature combinations based on differences rather than similarities.

### 3.5 KNN

KNN assumes that the data is in a feature space. More exactly, the data points are in a metric space. The data can be scalars or possibly even multidimensional vectors. Since the points are in feature space, they have a notion of distance – This need not necessarily be Euclidean distance although it is the one commonly used. Each of the training data consists of a set of vectors and class label associated with each vector. In the simplest case, it will be either + or – (for positive or negative classes). But KNN, can work equally well with arbitrary number of classes. We are also given a single number "k". This number decides how many neighbors (where neighbors is defined based on the distance metric) influence the classification. This is usually a odd number if the number of classes is 2. If  $k=1$ , then the algorithm is simply called the nearest neighbor algorithm.

### 4. DATASETS

Different parameters are considered for the prediction of the heart disease patients. The attributes of dataset are either related directly or indirectly. It is PIMA dataset.

Interactions	Glucose	BloodPressure	SkinThickness	Heart Rate	BMI	heartPedigreeFunction	Age	Outcome
6	148	72	35	0	33.6	0.627	50	1
1	85	66	29	0	26.6	0.351	31	0
8	183	64	0	0	23.3	0.672	32	1
1	89	66	23	94	28.1	0.167	21	0
0	137	40	35	168	43.1	2.288	33	1
5	116	74	0	0	25.6	0.201	30	0
3	78	50	32	88	31	0.248	26	1
10	115	0	0	0	35.3	0.134	29	0
2	197	70	45	543	30.5	0.158	53	1

Table.1 Attributes for Heart attack prediction

### 5. COMPARISON RESULTS

Collected dataset of the patients on different health parameters named as PIMA dataset has been under gone through different classifiers. The CART based technique performance in comparison to other techniques is better. So Cart based technique performance is comparatively efficient for the classification of the PIMA dataset into positive and negative classes

#### 5.1 Table for Accuracy

Accuracy							
	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
CART	0.666667	0.759358	0.787879	0.805184	0.848485	0.969697	0
LDA	0.617647	0.757576	0.787879	0.782217	0.812166	0.909091	0
SVM	0.65625	0.759358	0.787879	0.790027	0.842246	0.911765	0
KNN	0.617647	0.709099	0.757576	0.758183	0.792558	0.909091	0
RF	0.65625	0.735294	0.787879	0.790116	0.842246	0.939394	0

Table

represents the value for the CART, LDA, SVM KNN, and RF at the Confidence level of 0.95. CART based technique has optimal accuracy and the Random forest has least accuracy.

Table 3.kappa

Cohen's kappa coefficient ( $\kappa$ ) is a statistic which measures inter-rater agreement for qualitative (categorical) items. It is generally

Kappa						
Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
<b>CART</b> <b>0.15384615</b>	0.456294	0.514923	0.536164	0.647947	0.930233	0
<b>LDA</b> <b>0.05990783</b>	0.382867	0.477086	0.475126	0.565906	0.790698	0
<b>SVM</b> <b>0.11374408</b>	0.380488	0.474672	0.477213	0.619801	0.793522	0
<b>KNN</b> <b>0.14671815</b>	0.30273	0.441558	0.429303	0.529555	0.790698	0
<b>RF</b> <b>0.22123894</b>	0.349121	0.511628	0.503959	0.608541	0.863636	0

thought to be a more robust measure than simple percent agreement calculation, as  $\kappa$  takes into account the possibility of the agreement occurring by chance.

### 5.2 Table of kappa

In current table the Kappa parameter is compared for all the classifiers like CART, LDA, SVM, KNN, RF. MAX kappa value is for CART based technique. the least value for the Kappa technique is for the LDA based technique

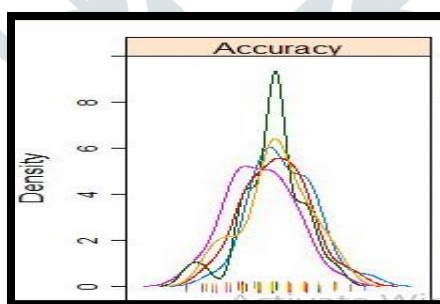
Kappa						
Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
<b>RF</b> <b>0.22123894</b>	0.349121	0.511628	0.503959	0.608541	0.863636	0
<b>SVM</b> <b>0.11374408</b>	0.380488	0.474672	0.477213	0.619801	0.793522	0
<b>KNN</b> <b>0.14671815</b>	0.30273	0.441558	0.429303	0.529555	0.790698	0

Table 4.

kappa

accuracy

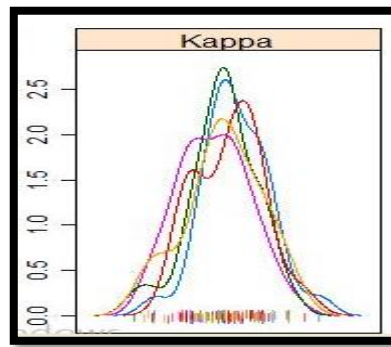
### 5.3 Graph for accuracy comparison



Graph 1 comparative Accuracy

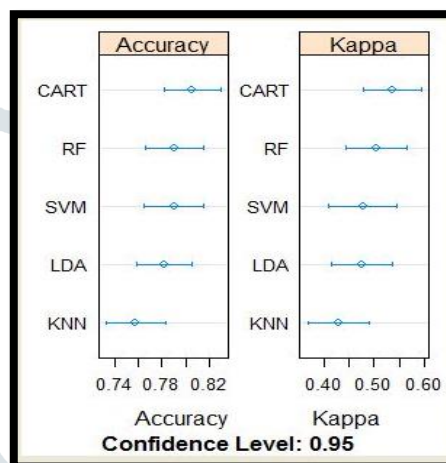
This graph shows the accuracy for the CART based technique as the maximum. And the LDA based technique has lowest accuracy.

#### 5.4 Graph for the Kappa



Graph 2 Kappa

This graph shows the Kappa comparison for different techniques. These technique has best technique as CART based technique. it has max. kappa value.



Graph 3 Line graph for accuracy and Kappa

This graph shows that the accuracy and Kappa value are maximum for CART based technique. the evaluation is done at the confidence value of 0.95.

#### 6. CONCLUSION

Data mining is the most prevalent field in the area of the extraction of the relevant data. There are various applications which are generating data in large amount. These data items are in billions of bytes on every day basis. Traditional data processing applications are not enough which can process these data items on larger scale. There are different classification techniques which are used for the classification of the large dataset into two or more different classes. Later on certain prediction based system will be required for generating accurate prediction. In current research there are different classifiers are being used which has been applied on the PIMA dataset for the purpose of prediction. So that the patient whether will be having Heart attack or not. Various classifiers like LDA, Random Forest, SVM, CART, KNN has been applied. The result for the accuracy and Kappa factor has been evaluated at the confidence level of 0.95. From all the classifier CART based technique has been best in both the factors like accuracy and Kappa.

#### 7. FUTURE WORK

Mining is the most important aspect as far as system of decision making is concerned. There are various applications which are used by different businesses for the enhancement of the quality of the decision. Various classification techniques are there which can be used for the purpose of prediction. In future the accuracy can be enhanced further by hybridization of classifiers.

#### References

- [1] August, J., No, I., & Kaur, R. (2017). Available Online at [www.ijarcs.info](http://www.ijarcs.info) International Journal of Advanced Research in Computer Science Check the result is Dengue or No Analysis of Dengue Data Dataset Indicate which symptom is more. 8(7), 1233–1237.
- [2] Kaur, N. K., & Gupta, G. (2015). PREDICTING THE VARIOUS RISK FACTORS OF LEPROSY. (3), 79–87.
- [3] Takuo Emoto, Tomoya Yamashita, Toshio Kobayashi, Naoto Sasaki, Yushi Hirota, Tomohiro Hayashi, Anna So, Kazuyuki Sahara, Keiko Yodoi, Takuya Matsumoto, Taiji Mizoguchi, Wataru Ogawa, Ken-ichi Hirata, " Characterization of gut microbiota profiles in coronary artery disease patients using data mining analysis of terminal restriction fragment length polymorphism: gut microbiota could be a diagnostic marker of coronary artery disease", Springer, vol. 3, pp: 89-100, 2016.
- [4] Dr. M. Renuka Devi and Shyal, "Exploring various data mining techniques", vol. 43, pp: 189-196, 2016.
- [5] ROSA MEO, GIUSEPPE PSAILA and STEFANO CERI "An Extension to SQL for Mining Association Rules", Data Mining and Knowledge Discovery, 2, 195–224 (1998)
- [6] Charu C. Aggarwal and Philip S. Yu, "Data Mining Techniques for Associations, Clustering and Classification", c\_Springer-Verlag Berlin Heidelberg, N. Zhong and L. Zhou (Eds.): PAKDD'99, LNAI 1574, pp. 13–23, 1999.



- [7] Andreas Buja , Young-Seop Lee, “Data Mining Criteria for Tree-Based Regression and Classification” ACM , 2001 1-5811-391-x/01/08
- [8] Dursun Delen \*, Christie Fuller, Charles McCann, Deepa Ray, “Analysis of healthcare coverage: A data mining approach”, Elsevier Expert Systems with Applications 36 (2009) 995–1003
- [9] Hai-Dong Meng · Yu-Chen Song · Fei-Yan Song · Hai-Tao Shen , “Research and application of cluster and association analysis in geochemical data processing”, Springer Comput Geosci (2011) 15:87–98
- [10] Gao Huang, Shiji Song, Jatinder N. D. Gupta, and Cheng Wu , “Semi-Supervised and Unsupervised Extreme Learning Machines” , 2168-2267 c ( 2014) IEEE
- [11] Aliza Ahmad and Aida Mustapha H, Comparison between Neural Networks against Decision Tree in Improving Prediction Accuracy for Heart attack Mellitus. ICDIPC 2011, Part I, CCIS 188, pp. 537–545, 2011.
- [12] Alexis Marcano-Cedeño, Joaquín Torres, and Diego Andina, A Prediction Model to Heart attack Using Artificial Metaplasticity. IWINAC 2011, Part II, LNCS 6687, pp. 418–425, 2011.
- [13] Humar, K. and Novruz, A., Design of a hybrid system for the Heart attack and heart diseases. Expert Systems with Applications, 35, 82–89, 2008.
- [14] Rebecca Schnall and Marlene Rojas, A user-centered model for designing consumer mobile health (mHealth) applications (apps). Journal of Biomedical Informatics 60 (2016) 243–251.
- [15] Ki-Hyun Kim, Ehsanul Kabir and Shamin Ara Jahan, The use of cell phone and insight into its potential human health impacts, Environ Monit Assess (2016)
- [16] Ioannis Kavakiotis , Olga Tsave, Athanasios Salifoglou , Nicos Maglaveras , Ioannis Vlahavas , Ioanna Chouvarda , “ Machine Learning and Data Mining Methods in Heart attack Research” , Elsevier, Computational and Structural Biotechnology Journal 15 (2017) 104–116.

