

Breast Cancer Prediction Using Clustering Techniques

Priyabrata Karmakar¹, Debolina Dalui², Dharmpal Singh^{3*}, Ira Nath⁴

^{1,2,3,4}Computer Science and Engineering, JIS College of Engineering, JIS University, Kalyani, India

*Corresponding Author: dharmpal.singh@jiscollege.ac.in, Tel.: 7003058392

Available online at: www.ijcseonline.org

Abstract— Data mining is the process of extracting hidden interesting patterns from massive database. It is used to extract the hidden information/knowledge from the real-life database. It also has use in the medical field. Thus want to apply concept of data mining into breast cancer prediction. According to many surveys, it has been observed that all over the world most of the women are dying in breast cancer in recent days. So, we have made an effort to design the knowledge base using clustering technique on the data. We have applied the hierarchical clustering and found the error is 2.13%. The K means and fuzzy C means will be applied in the future to minimize the error further.

Keywords— Breast cancer, Data mining, Clustering, Hierarchical clustering.

I. INTRODUCTION

Cancer is one of the most common diseases in the world that results in majority of death. Cancer is caused by the growth of cells which are not in control in any of the tissues or parts of the body. Cancer may grow up in any part of the body and may spread to other parts. Only early detection of cancer at the benign stage and prevention from spreading to other parts in malignant stage could save a person's life. There are several things that could affect a person's susceptibility for cancer. Developed countries have done a number of surveys & studies that shown cancer varies between people to people with various levels of education. The most of breast cancer has been found among the people with high levels of education whereas an inverse level of education has been found for the incidence of cancers of the lung, stomach and uterine cervix. Such differences in cancer risks are related to the education, life-style factors and disclosure of both environmental and daily-work related carcinogens. This study describes the relation between cancer pattern and risk levels of various factors by generating a risk prediction system for breast cancer which helps in prediction.

The risks related to Breast Cancer can be reduced via early detection of the disease. According to the American Cancer Society (2007) early detection of breast cancer risks can help cutting down the possibility of preventing the full growth of tumors. The various ways of detecting breast cancer includes: clinical examination by a physician, self breast examination and mammography. Clinical examination of breast by a physician is one of the effective ways of reducing breast cancer possibility. A woman should go for clinical examination annually when above 40 years and every 3 years in between 20 to 40 years. Mammography involves the use

of x-rays but with lower radiation; it has a breast cancer detection accuracy of (85 – 90)% where routing mammogram leads to a (25 – 30)% decrease in breast cancer possibility (American Cancer Society, 2007). Self-breast examination requires monthly observation of the breast and underarm by the patient. It makes the patient to be habituated with her breast and easily detect any exception she observes during the exercise. Diagnosis is the process of predicting the presence of breast cancer as either benign or malignant cases.

II. RELATED WORK

SHWETA KHARYA et.al [1] have discussed lots of data mining approaches that have been utilized for breast cancer diagnosis and prognosis. This study paper summarizes various review and technical articles on breast cancer diagnosis and prognosis also we focus on current research being carried out using the data mining techniques to enhance the breast cancer diagnosis and prognosis.

A.PRIYANGA & DR. S.PRAKASAM et.al [2] has implemented the data mining based cancer prediction System (DMBCPS). This system estimates the risk of the breast cancer in the earlier stage. This system is validated by comparing its predicted results with patient's prior medical information and it was analyzed by using weka system.

P.RAMACHANDRAN, N. GIRIJA & T.BHUVANESWARI et.al [3] have proposed a novel multi layered method combining clustering and decision tree techniques to build a cancer risk prediction system which predicts lung, breast, oral, cervix, stomach and blood cancers and is also user friendly, time and cost saving. It uses data mining technology

such as classification, clustering and prediction to identify potential cancer patients.

NIHARIKA SAXENA & PROF. MEETA KUMAR et.al [4] have gathered a comprehensive study of different clinical and non-clinical approaches used for diagnosis and prediction on different dataset of breast cancer. It gives an idea of our proposed hybridized evolutionary model.

S.SYED SHAJAHAN, S.SHANTHI & V.MANOCHITRA et.al [5] have explored the applicability of decision trees to predict the presence of breast cancer. Also they analyzed the performance of conventional supervised learning algorithms viz. Random tree, ID3, CART, C4.5 and Naive Bayes.

PETER ADEBAYO IDOWU, KEHINDE OLADIPO WILLIAMS, JEREMIAH ADEMOLA BALOGUN AND ADENIRAN ISHOLA OLUWARANTI et.al [6] did a study which is mainly focused at using two data mining techniques to predict breast cancer risks in Nigerian patients using the naïve bayes' and the J48 decision trees algorithms. The J48 decision trees showed a higher accuracy with lower error rates compared to that of the naïve bayes'.

It has been observed that, the authors have used data mining concepts in different real life optimization problem. Therefore in this paper the concept of data mining and soft computing will be used to form the knowledge base. The knowledge base predict the value based on the new data. In section 1, implementation and literature work have been furnished where as in section 2 and 3 methodology and implementation have been discussed. In the section 4 result and conclusion have been discussed.

III. METHODOLOGY

1. Multivariate Data Analysis

Multivariate statistics is a form of statistics encompassing the simultaneous observation and analysis of more than one outcome variable.

Multivariate statistics concerns understanding the different aims and background of each of the different forms of multivariate analysis, and how they relate to each other. The practical implementation of multivariate statistics to a particular problem may involve several types of univariate and multivariate analysis in order to understand the relationship between variables and their relevance to the actual problem being studied.

In addition, multivariate statistics is concerned with multivariate probability distributions, in terms of both:

1. How these can be used to represent the distributions of observed data;
2. How these can be used as part of statistical inference, particularly where several different quantities are of interest to the same analysis.

Certain types of problem involving multivariate data, for example simple linear regression and multiple regression, are not usually considered as special cases of multivariate statistics because the analysis is dealt with by considering the (univariate) conditional distribution of a single outcome variable for the r given other variables.

There are many different models, each with its own type of analysis:

1. Multivariate analysis of variance (MANOVA) extends the analysis of variance to cover cases where there is more than one dependent variable to be analyzed simultaneously.
2. Multivariate regression analysis attempts to determine a formula that can describe how elements in a vector of variables respond simultaneously to changes in others. For linear relations, regression analyses here are based on forms of the general linear model.
3. Principal components analysis (PCA) creates a new set of orthogonal variables that contain the same information as the original set. It rotates the axes of variation to give a new set of orthogonal axes, ordered so that they summarize decreasing proportions of the variation.
4. Factor analysis is similar to PCA but allows the user to extract a specified number of synthetic variables, fewer than the original set, leaving the remaining unexplained variation as error. The extracted variables are known as latent variables or factors; each one may be supposed to account for correlation in a group of observed variables.

1.1 Factor Analysis

Factor analysis is a statistical method used to describe variability among observed, correlated variables in terms of a potentially lower number of unobserved variables called **factors**. In other words, it is possible, that variations in fewer observed variables mainly reflect the variations in total effect. Factor analysis searches for such joint variations in response to unobserved latent variables. The observed variables are modeled as linear combinations of the potential factors, plus "error" terms. The information gained about the interdependencies between observed variables can be used later to reduce the set of variables in a dataset. Computationally this technique is equivalent to low rank approximation of the matrix of observed variables. Factor analysis originated in psychometrics, and is used in behavioral sciences, social sciences, marketing, product management, operations research, and other applied sciences that deal with large quantities of data.

1.2 Clustering

Clustering deals with finding a structure in a collection of unlabeled data. A loose definition of clustering could be "the process of organizing objects into groups whose members

are similar in some way". A cluster is therefore a collection of objects which are "similar" between them and are "dissimilar" to the objects belonging to other clusters. The graphical example of cluster has shown in figure 1.9. In this case it is easily identified the 4 clusters into which the data can be divided; the similarity criterion is distance: two or more objects belong to the same cluster if they are "close" according to a given distance (in this case geometrical distance). This is called distance-based clustering.

Another kind of clustering is conceptual clustering: two or more objects belong to the same cluster if one object defines a concept common to all that objects. In other words, objects are grouped according to their fit to descriptive concepts, not according to simple similarity measures.

Clustering Algorithms

Clustering algorithms may be classified as listed below:

1. Exclusive Clustering
2. Overlapping Clustering
3. Hierarchical Clustering
4. Probabilistic Clustering

In the first case data are grouped in an exclusive way, so that if a certain datum belongs to a definite cluster then it could not be included in another cluster. A simple example of that is shown in the figure 1.10, where the separation of points is achieved by a straight line on a bi-dimensional plane. On the contrary the second type, the overlapping clustering, uses fuzzy sets to cluster data, so that each point may belong to two or more clusters with different degrees of membership. In this case, data will be associated to an appropriate membership value.

Instead, a hierarchical clustering algorithm is based on the union between the two nearest clusters. The beginning condition is realized by setting every datum as a cluster. After a few iterations it reaches the final clusters. Finally, the last kind of clustering uses a completely probabilistic approach. The mostly used clustering algorithms are:

1. K-means
2. Fuzzy C-means
3. Hierarchical clustering

K-means is an exclusive clustering algorithm, Fuzzy C-means is an example of overlapping clustering algorithm, hierarchical clustering is obvious and lastly. Mixture of Gaussian is a probabilistic clustering algorithm. The distance measure plays important role in clustering which have been discussed in next section.

Distance Measure

An important component of a clustering algorithm is the distance measure between data points. If the components of the data instance vectors are in the same physical units it is possible that the simple Euclidean distance metric is sufficient to successfully group similar data instances.

Clustering algorithm

K-Means Clustering

K-means (MacQueen, 1967) clustering algorithm is one of the simplest unsupervised learning algorithms that solve the well known clustering problem. The procedure follows a simple and easy way to classify a given data set through a certain number of clusters (assume k clusters) with a priori. The main idea is to define k centroids, one for each cluster. These centroids should be placed in a cunning way because of different location causes different result. So, the better choice is to place them as much as possible far away from each other. The next step is to take each point belonging to a given data set and associate it to the nearest centroid. When no point is pending, the first step is completed and an early groupage is done. At this point it is needed to re-calculate k new centroids as barycenters of the clusters resulting from the previous step. After these k new centroids have been obtained a new binding has to be done between the same data set points and the nearest new centroid. A loop has been formed. As a result of this loop it may be noticed that the k centroids change their location step by step until no more changes are done. In other words centroids do not move any more.

Finally, this algorithm aims at minimizing an objective function, generally a squared error function. The objective function as

$$J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2, \text{ where } \|x_i^{(j)} - c_j\|^2 \text{ is a chosen distance measure between a data point } x_i^{(j)} \text{ and the cluster centre } c_j.$$

K-Means Algorithm

1. Place K points into the space represented by the objects that are being clustered. These points represent initial group of centroids.
2. Assign each object to the group that has the closest centroid.
3. When all objects have been assigned, recalculate the positions of the K centroids.
4. Repeat Steps 2 and 3 until the centroids no longer move. This produces a separation of the objects into groups from which the metric to be minimized can be calculated.

Although it can be proved that the procedure will always terminate, the k-means algorithm does not necessarily find the most optimal configuration, corresponding to the global objective function minimum. The algorithm is also significantly sensitive to the initial randomly selected cluster centers. The k-means algorithm can be run multiple times to reduce this effect. K-means algorithm is a simple algorithm that has been adapted to many problem domains.

Hierarchical Clustering Algorithm

Hierarchical clustering technique is one of the important clustering methods. There are two basic approaches for generating a hierarchical clustering.

Agglomerative: Start with the points as individual clusters and, at each step, merge the closest pair of clusters. This requires defining a notion of cluster proximity.

Divisive: Start with one, all-inclusive cluster and, at each step, split a cluster until only singleton clusters of individual points remain. In this case, it is needed to decide which cluster to split at each step and how the splitting to be done.

Agglomerative hierarchical clustering techniques are by far the most common. A hierarchical clustering algorithm is often displayed graphically using a tree-like diagram called a dendrogram, which displays both the cluster, sub cluster relationships and the order in which the clusters are merged (agglomerative view) or splitted (divisive view). For sets of two-dimensional points, a hierarchical clustering can also be graphically represented using a nested cluster diagram..

Defining Proximity between Clusters

The cluster proximity defines the difference among the various agglomerative hierarchical techniques. Cluster proximity is typically defined with a particular type of cluster in mind, for an example, many agglomerative hierarchical clustering techniques, such as MIN, MAX, and Group Average which have come from a graph-based view of clusters. Min defines cluster proximity as the proximity between the closest two points that are in different clusters, or using graph terms, the shortest edge between two nodes in different subsets of nodes. Alternatively, max, takes the proximity between the farthest two points in different clusters to be the cluster proximity, or using graph terms, the longest edge between two nodes in different subsets of nodes

For similarities, however, where higher values indicate closer points, the names seem reversed. For that reason it is usually preferred to use the alternative names, single link and complete linkage respectively. Another graph-based approach, the group average technique, defines cluster proximity to be the average pair wise proximities (average length of edges) of all pairs of points from different clusters. Figure 1.13 illustrates these three approaches

For a given set of N items to be clustered, and an N*N distance (or similarity) matrix is to be constructed, the basic process of hierarchical clustering is as follows:-

1. Start by assigning each item to a cluster, so that there N items, N clusters are available, each containing just one item. Let the distances (similarities) between the clusters be the same as the distances (similarities) between the items they contain.

2. Find the closest (most similar) pair of clusters and merge them into a single cluster, so that one less number of clusters is available..
3. Compute distances (similarities) between the new cluster and each of the old clusters.
4. Repeat steps 2 and 3 until all items are clustered into a single cluster of size N. (*)

IV. IMPLEMENTATION

The available data contains the Breast Cancer based on the following attributes viz. ID No, Clump Thickness, Uniformity of Cell Size, Uniformity of Cell Shape, Marginal Adhesion, Single Epithelial Cell Size, Bare Nuclei, Bland Chromatin, Normal Nucleoli, Mitoses, Class etc. The purpose of this work is to correlate the items, so that based on any value of attributes, the value of disease can be estimated. Now it is necessary to check whether the data items are proper or not. If the data items are proper, extraction of information is possible otherwise the data items are not suitable for the extraction of knowledge. In that case preprocessing of data is necessary for getting the proper data.

Data Cleansing

The data cleansing techniques include the filling in the missing value, correct or rectify the inconsistent data and identify the outlier of the data which have been applied on the available data.

It has been observed that each data set does not contain any missing value. The said data item does not contain any inconsistent data i.e. any abnormally low or any abnormally high value. All the data values are regularly distributed within the range of that data items. Therefore the data cleansing techniques are not applicable for the available data.

Data Integration

The data integration technique has to be applied if data has been collected from different sources. The available data have been taken from a single source therefore the said technique is not applicable here.

Data transformation

The data transformation techniques such as smoothing, aggregation, normalization, decimal scaling have to be applied to get the data in proper form. Smoothing technique has to be applied on the data to remove the noise from the data. Aggregation technique has to be applied to summarize the data. To make the data within specific range smoothing and normalization technique have to be applied. Decimal scaling technique has to be applied to move the decimal point to particular position of all data values. Out of these data transformation techniques, smoothing and decimal scaling techniques have been applied on the data

The methods of factor analysis has been applied on the available input data to make a decision to select the optimal model for the formation of association rule because Apriori Algorithm is not applicable for the numeric data.

Using the values of antecedent item, the correlation of coefficient of the items have been computed and these coefficients have been stored in table 1 termed as correlation matrix.

Table 1: Correlation matrix

1	0.64 54	0.65 51	0.48 67	0.52 21	0.58 89	0.55 85	0.53 62	0.34 98
0.64 54	1	0.90 69	0.70 56	0.75 18	0.68 73	0.75 6	0.72 29	0.45 89
0.65 51	0.90 69	1	0.68 31	0.71 97	0.70 97	0.73 62	0.71 94	0.43 91
0.48 67	0.70 56	0.68 31	1	0.59 98	0.66 6	0.66 7	0.60 34	0.41 83
0.52 21	0.75 18	0.71 97	0.59 98	1	0.58 22	0.61 64	0.62 9	0.47 97
0.58 89	0.68 73	0.70 97	0.66 6	0.58 22	1	0.67 63	0.57 72	0.33 98
0.55 85	0.75 6	0.73 62	0.66 7	0.61 64	0.67 63	1	0.66 61	0.34 49
0.53 62	0.72 29	0.71 94	0.60 34	0.62 9	0.57 72	0.66 61	1	0.42 88
0.34 98	0.45 89	0.43 91	0.41 83	0.47 97	0.33 98	0.34 49	0.42 88	1

The eigen value and eigne vector (table 3) of the elements using the above correlation matrix have been calculated using Matlab Tools. The contribution of eigen value of each item among all other items has been calculated and have been furnished in table 2.

Table 2 Eigen Value

0.014531	0.014531	0.014531	0.014531	0.014531	0.014531	0.014531	0.014531
0.052256	0.052256	0.052256	0.052256	0.052256	0.052256	0.052256	0.052256
0.09077	0.09077	0.09077	0.09077	0.09077	0.09077	0.09077	0.09077
0.002872	0.002872	0.002872	0.002872	0.002872	0.002872	0.002872	0.002872
0.053528	0.053528	0.053528	0.053528	0.053528	0.053528	0.053528	0.053528
0.001753	0.001753	0.001753	0.001753	0.001753	0.001753	0.001753	0.001753
0.043806	0.043806	0.043806	0.043806	0.043806	0.043806	0.043806	0.043806
0.002038	0.002038	0.002038	0.002038	0.002038	0.002038	0.002038	0.002038
0.000453	0.000453	0.000453	0.000453	0.000453	0.000453	0.000453	0.000453

From the table2, it has been observed that percentage contribution of the item which have less than 2.5 % as compared to other items therefore that items have been ignored. The eigen vectors of the items have been furnished in table 3

Table 3: Eigen Vectors

0.2355	0.2355	0.2355	0.2355	0.2355	0.2355	0.2355	0.2355
0.4466	0.4466	0.4466	0.4466	0.4466	0.4466	0.4466	0.4466
0.5886	0.5886	0.5886	0.5886	0.5886	0.5886	0.5886	0.5886
-0.1047	-0.1047	-0.1047	-0.1047	-0.1047	-0.1047	-0.1047	-0.1047
-0.452	-0.452	-0.452	-0.452	-0.452	-0.452	-0.452	-0.452
0.0818	0.0818	0.0818	0.0818	0.0818	0.0818	0.0818	0.0818
-0.4089	-0.4089	-0.4089	-0.4089	-0.4089	-0.4089	-0.4089	-0.4089
-0.0882	-0.0882	-0.0882	-0.0882	-0.0882	-0.0882	-0.0882	-0.0882

The major factors have been calculated as per the formula ($\sqrt{\text{eigen value}} \times \text{eigen vector}$) using the selected eigen value as furnished in table 2 and eigen vector as furnished in table 3. The major factors have been furnished in table 4

Table 4: Cumulative Effect Value of Items

Data Attributes	Cumulative Effect Values
A	1.000209
B	0.952342
C	0.96052
D	0.999662
E	0.999728
F	0.9993
G	0.999708
H	0.999772
I	0.999974

Now a relation has been formed by using the cumulative effect value of all the elements to produce total effect value. Total effect value = $(1.00) \times A + (0.952) \times B + (0.999) \times C$Now using the relation, a resultant total effect value formed (due to size of data it has not shown here)

Implementation by Clustering:

The Concept of hierarchical clustering has been used using the Matlab command to form the cluster on the total effect to select the optimal cluster for the formation of knowledge base.

The command is furnished as follows:

Step 1 :

$X = [\text{total effect value}]$;

Step 2 :

Use this we can define the whole data set as a matrix. Then pass it to pdist.

$Y = \text{pdist}(X)$

Step 3:

The pdist is a function which calculates the distance between every object.

$Z = \text{linkage}(Y)$

The linkage function is a function which generates a hierarchical cluster tree, then the linkage information is returning in a matrix Z. After that when hierarchical cluster tree is created we will partition the data into the clusters.

Step 4 :

At first, we create two cluster, then 4,8,16 onwards.

$T = \text{cluster}(Z,2)$

$T = \text{cluster}(Z,4)$

$T = \text{cluster}(Z,8)$

V. RESULTS AND DISCUSSION

Then we calculated the total distance of each cluster and we found that the total distance for the 4th cluster ($T = \text{cluster}(Z,16)$) is less as furnished in table 5.

Table 5. Distance of cluster based on number of cluster

Cluster Number	Total Distance
2	6222.2152
4	3333.782
8	1092.03
16	796.383

After that based on the cluster number 16 , we have related the item with the sixteen cluster and then predict the output based on it. The error has been found as 2.13 %

VI. CONCLUSION AND FUTURE SCOPE

This research work is based on clustering in health care which has been developed by using data mining techniques and clustering techniques. The concept of hierarchical Clustering for prediction of breast cancer has been used here and prediction has been done based on it. It has been observed that , hierarchical clustering has given 2.13 % as error.

The concept of K-means and fuzzy c means will be used to minimize the error in our future work.

ACKNOWLEDGMENT (HEADING 5)

Authors are grateful towards JIS College of Engineering for providing lab and related facilities for carrying out the research activities.

REFERENCES

- [1] Shweta kharya, ‘Using data mining techniques for diagnosis and prognosis of cancer disease’, INTERNATIONAL JOURNAL OF COMPUTER SCIENCE, ENGINEERING AND INFORMATION TECHNOLOGY (IJCEIT), vol.2, no.2, april 2012.
- [2] A.Priyanga, Dr.S.Prakasam, ‘ The role of data mining-based cancer prediction system (dmbcps) in cancer awareness’, INTERNATIONAL JOURNAL OF COMPUTER SCIENCE AND ENGINEERING COMMUNICATIONS- IJCSEC, vol.1 issue.1, december 2013.
- [3] P.ramachandran, N. Girija, T.bhuvanawari, ‘Early detection and prevention of cancer using data mining techniques’, INTERNATIONAL JOURNAL OF COMPUTER APPLICATIONS (0975 – 8887) volume 97– no.13, july 2014.
- [4] Niharika saxena, prof. Meeta kumar, ‘Comprehensive study on data clustering for breast cancer prognosis and risk exposure’, INTERNATIONAL JOURNAL OF PURE AND APPLIED MATHEMATICS, volume 118 no. 24 2018.
- [5] S.Syed Shajahaan, S.Shanthi , V.Manochitra, ‘Application of data mining techniques to model breast cancer data’, INTERNATIONAL JOURNAL OF EMERGING TECHNOLOGY AND ADVANCED ENGINEERING, ISSN 2250-2459, iso 9001:2008 certified journal, volume 3, issue 11, november 2013.
- [6] Peter Adebayo Idowu, Kehinde Oladipo williams, Jeremiah Ademola Balogun, Adeniran Ishola Oluwaranti, ‘Breast cancer risk prediction using data mining classification techniques’, SOCIETY FOR SCIENCE AND EDUCATION, UNITED KINGDOM, volume-3, issue-2, ISSN : 2054-7420.
- [7] G . Purkait and D.P Singh, “An effort to optimize the error using statistical and soft computing methodologies” , Journal of Applied Computer Science & Artificial Intelligence Vol .1 No. 1, pp. 15-20, 2017.
- [8] D. P., J. P. Choudhury and M. De, “An Enhance DE algorithm for analysis in data set”, International Journal of Data Science (IJDS), 2016 (Accepted).
- [9] Dharmal Singh, Abhishek Banerjee, Gopal Purkait , “Assessment of Heart DiseaseTypes Using Clustering Techniques Under the Domain of Data Mining”, International Journal of Artificial Intelligence and Knowledge Discovery, Vol.6 Issue 03 pp. 10-16, Print ISSN: 2231-2021 e-ISSN: 2231-0312 .
- [10] D. P. Singh, S. Sahana, SK Saddam Ahmed , "A Comparative Study to Assess the Crohn's Diseasetype using Statistical and Fuzzy Logic Methodology", International Journal of Computer Sciences and Engineering, Volume-04, Issue-06, Page No (41-46), Aug -2016, E-ISSN: 2347-2693
- [11] D. P. Singh, J. P. Choudhury and M. De, “An effort to select a preferable meta heuristic model for knowledge discovery in Data mining”, International Journal of mataheuristics, Vol. 4 No. 1, pp.57-90, September, 2015. 1755-2184, ISSN print: 1755-2176 DBLP Index
- [12] Sk. Saddam and D. P.Singh, “Genotype based classification of Crohn's disease using various BPN training algorithms”, International Journal of Artificial Intelligence and Knowledge Discovery, Vol. 5 No. 2, pp.22-29, July, 2015, Print ISSN: 2231-2021 e-ISSN: 2231-0312.
- [13] D. P. Singh, J. P. Choudhury and M. De, “An Effort to Compare the Clustering Technique on Different Data Set Based On Distance Measure Function in the Domain of Data Mining”, International Journal of Artificial Intelligence and Knowledge Discovery, Vol. 5, No. 1, pp. 1-8, January, 2015, Print ISSN: 2231-2021 e-ISSN: 2231-0312.
- [14] D. P. Singh, J. P. Choudhury and M. De, “A Comparative Study to Select a Soft Computing Model for Developing the Knowledge Base of Data mining with Association Rule Formation by Factor Analysis”, International Journal of Artificial Intelligence and Knowledge Discovery, Vol. 3, No. 3, pp.18-23, October, 2013, Print ISSN: 2231-2021 e-ISSN: 2231-0312.
- [15] D. P. Singh, J. P. Choudhury and M. De, “An Effort to Developing the Knowledge Base in Data Mining by Factor Analysis and Soft Computing Methodology”, Internat

- ional Journal of Scientific & Engineering Research (IJSER), Vol. 4, No. 9, pp. 1912-1923, September, 2013, ISSN 2229-5518.
- [16] D. P. Singh, J. P. Choudhury and M. De, "A comparative study on the performance of Fuzzy Logic, Bayesian Logic and neural network towards Decision Making" International Journal of Data Analysis Techniques and Strategies (IJDATS), Vol. 4, No. 2, pp. 205-216, April, 2012. SSN online: 1755-8069 ISSN print: 1755-8050, Scopus Index
- [17] D. P. Singh, J. P. Choudhury and M. De, "A Comparative Study to Select a Soft Computing Model for Knowledge Discovery in Data Mining", International Journal of Artificial Intelligence and Knowledge Discovery, Vol. 2, No. 2, pp. 6-19, April, 2012.
- [18] D. P. Singh, J. P. Choudhury and M. De, "A Comparative Study on the performance of Soft Computing models in the domain of Data Mining," International Journal of Advancements in Computer Science and Information Technology, Vol. 1, No. 1, pp. 35-49, September, 2011, ISSN 2277-9140.
- [19] D. P. Singh, J. P. Choudhury and M. De, "Optimization of fruit quantity by different types of cluster technique", Journal of Computer Sciences, Punjab, Vol. 9, No.1, pp. 17-28, June-July, 2011, ISSN 0973-4058.
- [20] D. P. Singh, J. P. Choudhury and M. De, "Optimization of Fruit Quantity by comparison between Statistical Model and Fuzzy Logic by Bayesian Net work", PCTE, Journal of Computer Sciences, Vol. 8, No.1, Punjab, pp. 91-95, June-July, 2010.
- [21] D. P. Singh, J. P. Choudhury and M. De, "Prediction Based on Statistical and Fuzzy Logic Membership Function", PCTE, Journal of Computer Sciences, Vol. 8, No. 1, Punjab, pp. 86-90, June-July, 2010.
- [22] D.P. Singh, J. P. Choudhury and M. De, "Performance Measurement of Neural Net work Model Considering Various Membership Functions under Fuzzy Logic", International Journal of Computer and Engineering, Vol. 1, No. 2, pp. 1-5, 2010, ISSN-0976-9587.
- [23] D. P. Singh, J. P. Choudhury, "Assessment of Exported Mango Quantity by Soft Computing Model", International Journal of Information Technology and Knowledge Management, Kurukshetra University, Vol. 2, No. 2, pp. 393-395, June-July, 2009, ISSN : 0973-4414.
- [24] D.P Singh, "A Modified Bio Inspired BAT algorithm," International Journal of Applied Metaheuristic Computing (IJAMC), Vol. 9 No. 1, pp. 60-77, 2018. Scopus Index.
- [25] D. P. Singh, J. P. Choudhury and M. De, "A modified ACO for classification on different data set" International Journal of Computer Application, Vol.123, No. 6, pp-39-52, August 2015, ISSN 0975-8887.
- [26] D. P. Singh, J. P. Choudhury and M. De, "Performance measurement of Soft Computing models based on Residual Analysis", International Journal for Applied Engineering and Research, Vol. 6, No. 5, Delhi, India, pp. 823-832, Jan-July, 2011, Print ISSN 0973-4562. Online ISSN 1087-1090. Scopus Index.

Authors Profile

Priyabrata Karmakar is a 2nd year student, persuing Master of Technology in the department of Computer Science and Engineering, from a renowned engineering college of West bengal.



Dharmal Singh received his Bachelor of Computer Science and Engineering and Master of Computer Science and Engineering from West Bengal University of Technology. He has done his Ph.D in year 2015. He has about 12 years of experience in teaching and research. At present, he is with JIS College of Engineering, Kalyani, and West Bengal, India as an Associate Professor and Head of the department. He has published 32 papers in referred journal and conferences index by **Scopus, DBLP and Google Scholar and editorial team and senior member** of many reputed journal index by **SCI, Scopus, DBLP and Google Scholar**.



He has organized **seven national levels Seminar/Workshop, published three patent** and has applied for the **AICTE Project in year of 2019**.

He is a member of the Computer Science Teachers Association (CSTA), Computer Society of India (CSI) and also a member of International Association of Computer Science and Information Technology (IACSIT).

Debolina Dalui is a 2nd year student, persuing Master of Technology in the department of Computer Science and Engineering, from a renowned engineering college of West bengal.



Mrs Ira Nath is presently working as an Assistant Professor in the Department of Computer Science and Engineering of JIS College of Engineering, India. She received the Master of Technology (M.Tech.) degree in Software Engineering from the Maulana Abul Kalam Azad University of Technology, India formerly West Bengal University of Technology, India in 2008. She also received the degree of Bachelor of Technology (B.Tech.) in Computer Science and Engineering from the same university in 2005. She is presently pursuing her Ph.D in Computer Science & Technology at Indian Institute of Engineering Science and Technology (IEST), Shibpur, India. Her research interests include Network Security regenerator placement, survivability and routing and wavelength assignment in translucent WDM optical Networks.

