# BREAST CANCER DIAGONOSIS IN ANALYSIS OF BRCA GENE USING MACHINE LEARNING ALGORITHMS

J. Sumitha* and T. Devi

Department of Computer Applications, Bharathiar University, India
Email*:sumivenkat2006@gmail.com

## ABSTRACT

Worldwide, the breast cancer is the second leading cancer type which leads to death among women. Breast cancer exists due to the mutation happens in the normal growth of Breast Cancer Gene (BRCA) under certain circumstances. In this paper, we proposed the existing machine learning algorithms for finding the disease – causing BRCA gene. These existing machine learning algorithms are compared with each other to determine the efficiency in detecting the diseases from gene expression value. The results proved that the Hybrid Radial Bias Neural Network (HRBFNN) performs better than Divide and Conquer Kernel Solving Support Vector Machines (DCKSVM) and other classification algorithms than ever before.

**KEYWORDS:** Confusion matrix, DCKSVM, Identification of BRCA gene, HRBFNN, Sequential model.

## INTRODUCTION

Cancer is the second leading cause of death and more than 1,500 people per day are affected by this disease all over the world. Approximately, there are 1479350 new cancer cases were diagnosed all over the world and analyzed that the cancer is the second most common reason of death and it accounts for nearly one of every four deaths. Breast cancer is the most common and rigorous cancer among women and it continuing to be a significant public health problem all over the world. There are about 182,000 new cases of breast cancer are diagnosed and about 46,000 women die of breast cancer each year. Nearly about 192,370 new cases of breast cancer were diagnosed among women and thus, the incidence and mortality of breast cancer are very high [1]. Hence, breast cancer is considered as the most common cause of death among women. The normal gene of BRCA exists in the nature of proto-oncogene in the human body. When the mutation occurs under certain circumstances such as habits, radiation and by hereditary aspects, this proto-oncogene is then translated into oncogene which leads to cancer.

Since breast cancer has very high incidence of death rate, the reason for breast cancer is still mysterious and there is no efficient way to prevent the existence of breast cancer. So, early detection is the first significant step to-wards treating breast cancer. The most common screening methods for earlier detection of breast cancer are mammography and sonography.

Among these two methods, mammography is the most important tool that doctors prefer to diagnose [1,2]. But, mammography has some disadvantages on diagnosing breast cancer. Since it is very sensitive, it is not accurate in detecting breast cancer because mammography can identify an abnormality that looks like cancer in the human body, but it turns out to be normal which is a false positive [3,4] and such a misdiagnosis gives hectic for patients. Hence it is essential to develop software which could give reliable diagnostic results to prevent these drawbacks.

The objective of this paper is to detect BRCA gene with the help of gene expression value using machine learning algorithms. The algorithms which used in this research are the sequential algorithm, DCKSVM and HRB FNN and the prediction is done on the basis of confusion matrix method.

This paper is organized as follows: Section 2 presents proposed methods to solve the task of detecting breast cancer. Section 3 contains results obtained and discussions and section 4 for conclusions.

## METHODS

The methods used for this research is categorized as the sequential algorithm, Divide and Conquer Kernal Solving Support Vector Machine (DCKSVM) and Hybrid Radial Bias Neural Network (HRBFNN).

**Data set Description:** The breast cancer dataset used in this research is taken from the URI

repository having 567 data with gene expression value which is shown in Figure-1. From this data, 380 data are taken as training data and remaining data are for testing data. The number of attributes taken for predicting the efficiency from this dataset is thirty-four and Matlab is the software that has been used for implementing this work. The URI repository link of this prognostic type of breast cancer dataset: http://archive.ics.uci.edu/ml/datasets/breast cancer+win consin+%28.

The parameters used for predicting the efficiency are accuracy, precision, recall and f-measure. But this performance measure for these classifiers in identifying disease is based on confusion matrix method which has been computed from this breast cancer dataset [5]. Algorithms :

Sequential Model It is used to distinguish the differences between the pair of gene expression values in the dataset. For example, gene1, gene2. …..gene10 are genes in the dataset in which the gene pair (gene4 gene8gene2) and the another gene pair (gene8gene4gene2) are found to be similar [6]. Then it shows that this two gene pair stimulates the same disease in the human body. This sequences of each gene in breast cancer dataset are taken as an input for predicting results [7].

**DCKSVM:** The most commonly used classification algorithm is the DCKSVM algorithm which applied to the breast cancer dataset for predicting the accuracy [9,10]. DCKSVM algorithm [8] is mainly used to divide the main clusters of data into a number of sub-clusters which makes its predictivity which becomes reliable on that clusters in this breast cancer dataset [8].
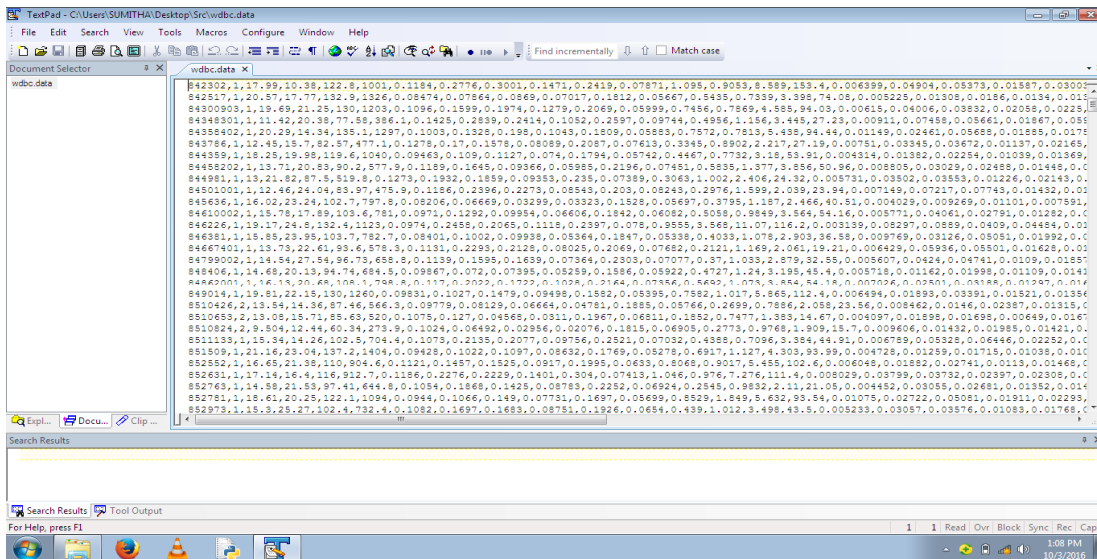


Figure-1:Breast Cancer Dataset

**Algorithm 1:** Divide and Conquer SVM

Input: Training Data

Output: The SVM solution A.

Step 1: Partitioning the variables into k subsets {v1….vk}

Step 2: Time complexity for solving sub-problems reduced to $O(k*n/k)2=O(n2/k)$ with space complexity, where n is the variable and k is the cluster subset.

Step 3: After computation of all subproblem solutions, concatenate them to form solution for the whole problem ab=[a1….ak]

Step 4: A bound is derived on $||ab-a*||2$ where ab is the optimal solution by adding cluster-kernel values.

Step 5: Minimizing the off-diagonal values of the kernel matrix with a balancing normalization.

Step 6: Go to Step 2 for partition the data in each cluster k, and calculate Step 4 for absolute scale.

**HRBFNN:** HRBFNN is applied to improve the capabilities of data granulation and for pre-processing the data in the breast cancer dataset [11]. It incorporates the characteristics of both data granulation and Principal Component Analysis

(PCA) method which is mainly used for prepro-cessing the data or neglecting the missing values in the dataset [12,13].

**Algorithm 2: HRBFNN**

Step 1: Preprocess the data set using PCA for dimentionality reduction.

Step 2: Training and testing data sets are formed.

Step 3: The generic parameters used in this research are decided.

Step 4: Selected inputs are determined.

Step 5: PFNs are designed. For the selecting r inputs, the number of nodes (PFNs) generated in each layer becomes equal to     k =      n!

n!(n−r)! r! , where, n is the number of total inputs and r stands for the number of the chosen input variables and k is the clusters [11].

Step 6: Check the termination criterion.

Step 7: Select the best predictive capability nodes and construct their corresponding layer.

It takes only a few seconds to complete 1000 iterations and generates high performance than other algorithms in terms of accuracy, precision, recall and f-measure.

**RESULTS AND DISCUSSION**

In this research, the machine learning algorithms which are applied to analyze the prognostic breast cancer dataset for detecting diseases are shown in the Table-1. HRBFNN algorithm gives better results when compared to the sequential algorithm and DCKSVM algorithm. The performance can be calculated in terms of accuracy, precision, recall and F-Measure. The criterion for these classifiers in disease detection is based on Confusion matrix. The accuracy percentage of the sequential model, DCKSVM and HRBFNN is 78, 80 and 85 respectively as depicted in Table-1. The results show that the HRBFNN gives 85.18 percentage of accuracy, 0.79 percentage of precision and 0.83 percentage of F-measure which is higher than the sequential model and the DCKSVM algorithms. However, DCKSVM algorithm shows higher recall percen-tage than sequential model and HRBFNN. But, the accuracy value is most significant in pridicting
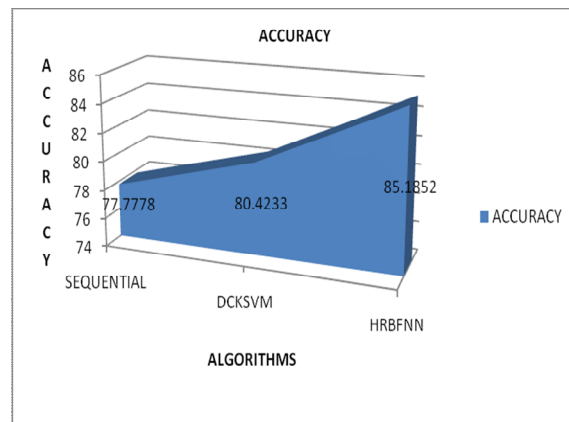
Table-1:Results of sequential, DCKSVM and HRBFNN

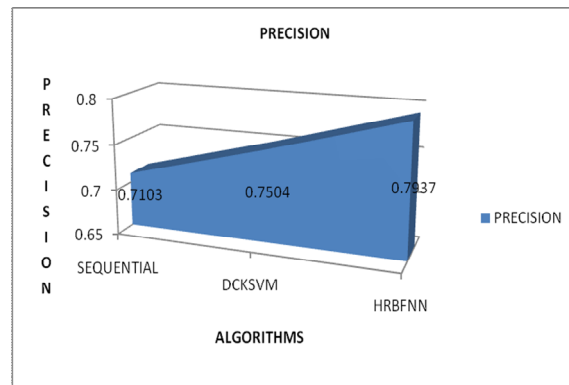| Algorithms<br><br>Parameters | Sequential Model(%) | DCKSVM(%) | HRBFNN (%) |
|---|---|---|---|
| Accuracy | 77.7778 | 80.4233 | 85.1852 |
| Precision | 0.7103 | 0.7504 | 0.7937 |
| Recall | 0.7659 | 0.8923 | 0.8713 |
| F-measure | 0.7371 | 0.7892 | 0.8307 |

the disease – causing gene. Hence, HRBFNN is best in predicting the BRCA gene in the dataset. The graph plotted for these results are shown in Figure-2.

Accuracy is the percentage of correct predictions from the dataset. Accuracy is the Proportion of the total number of predictions that were correct and it can be calculated using the equation,
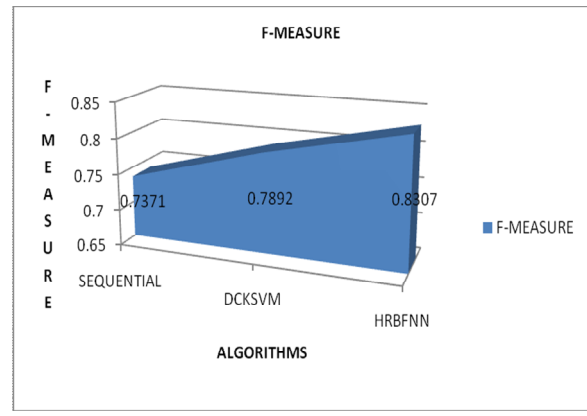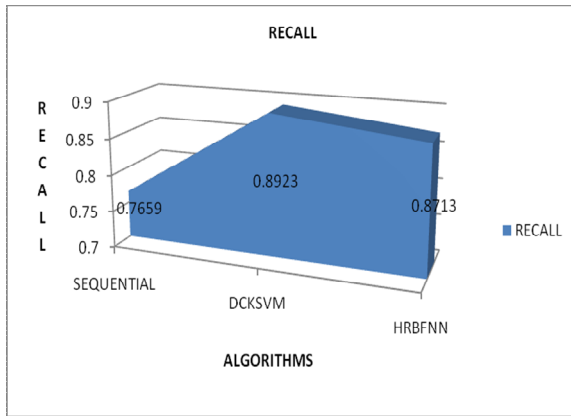
Accuracy= p+s/(p+q+r+s) …..................(1)



(a)



(b)

(c )



(d)

**Figure-2:** (a) Accuracy graph of Sequential, DCKSVM and HRBFNN Algorithms  (b) Precision graph of Sequential, DCKSVM and HRBFNN Algorithms (c) Recall graph of Sequential, DCKSVM and HRBFNN Algorithms (d) F-measure graph of Sequential, DCKSVM and HRBFNN Algorithms

Recall is the proposition of positive cases that are correctly identified and it can be calculated using the equation,

Recall=s/(r+s)             ……………...(2)

Precision is the proposition of the predicted positive cases that were correct and it can be calculated using the equation,

Precision=s/(q+s)          ……………(3)

F-measure (F) can be calculated using the equation,

F=2* (Precision*Recall)/(Precision+Recall) .... (4)

Where p is the number of correctly predicted genes which are not disease-causing gene, q is the number of incorrectly predicted genes which predicted as diseases causing gene, r is the number of incorrect prediction of undiseased gene and s is the number correctly predictions of disease-causing genes in the dataset.

Cho-Jui et al., [8] calculated accuracy and time of DCKSVM with LIBSVM in covtype dataset. But it is compared within the SVM classifiers itself and predicted the results in abse-nce of recall value. In this research, DCKSVM is compared with other classifiers rather than SVM classifiers and predicted the accuracy as 80.4 percentage and recall as 0.89 percentage. Wei et al., [11] estimated the performance of HRBFNN with the use of PCA and without the use of PCA for determining the predicting capabilities which show that the use of PCA with HRBFNN gives better performance than the without use of PCA experimented on using a series of numerical data. Hence this algorithm is selected for this research and predicted the accuracy of 85.18 percentage

which is the highest performance than other classifiers in diagnosing of breast cancer.

Machine learning algorithm is most com-monly used in all the research related to this work, but particularly in SVM classifiers, DCK SVM is implemented for the purpose of dividing the problem into sub-problems which aid for independent solutions with high efficiency. Table-1 concludes that HRBFNN predicts high efficiency than other classifiers. Based on the confusion matrix calculation method, the para-meters used for predicting the results are shown in equation (1) to (4).

CONCLUSION

In this research, HRBFNN algorithm proved that it is effective in analyzing the BRCA gene from the breast cancer dataset. The performance analysis of the sequential algorithm with DC KSVM and HRBFNN shows that the HRBFNN gives better efficiency than the DCKSVM and sequential algorithm. It can also be further enhanced with some other computer-assisted algorithms to improve the efficiency.

**REFERENCES**
[1]. Erin L. Linnenbringer, Social Constructions, Biological Implications: A Structural Exa-mination of Racial Disparities In Breast Cancer Subtype, University of Michigan, (2014)

[2].  Laura J.Van,  M.Mao, L.Hans, M.J. Marton, P.S. Linsley,   Gene Expression Profiling Predicts Clinical Outcome of Breast Can-cer. *Nature* 15: 415 (2012) www.nature.com,

[3]. Guha S., A.Meyerson, N.Mishra, R.Motwani and O'Callaghan, Clustering Data Streams: Theory and Practice. *IEEE Transactions*

*on Knowledge and Data Engineering* 15: 515-528 (2003)        http://dx. doi. org/ 10.1109/TKDE.2003.119  8387.

[4]. Ingrid A. Hedenfalk, Markus Ringner, Jerey M. Trent and Ake Borg, Gene  Expression in Inherited Breast Cancer.  *Advances in Cancer Research* 4: 1-34, (2002)

[5].   Isabelle G.,  W.Jason,  B.Stephen  and  V. Vladimir, Gene Selection for Cancer Classi-fication using Support Vector Machines. *Machine Learning* 16(1):  389-422 (2014)

[6].   Ken Kaneiwa, A sequential pattern mining algorithm  using  rough  set  theory. *Inter-national    Journal    of    Approximate Reasoning* 52: 881–893 (2011)

[7]. Carl H. Mooney and  John F. Roddick, Sequ-ential  Pattern  Mining – Approaches  and Algorithms. *ACM  Computing  Surveys* (2013)

[8].   Cho-Jui, S. Si and Inderjit S. Dhillon, A Divide  and  Conquer  Solver  for  Kernel Support Vector Machines. *Proceedings of the* 31st *International  Conference  on Machine Learning,* Beijing, China (2014)

[9]. Minseung K. and K.Sung-Hou, Empirical prediction of genomic susceptibilities for multiple  cancer  classes.  *PNAS*  11(5): 1921–1926 (2015)

[10]. Oana F., I.Diana and T.Thomas, A Machine Learning Approach for Identifying Disease-Treatment Relations in Short Texts. *IEEE Transactions  on  Knowledge  and  Data Engineerin* 23(6):    (2016)

[11]. Wei H., Sung-Kwun and P. Witold,  Design of hybrid radial basis function neural net-works (HRBFNNs) realized with the aid of hybridization of fuzzy clustering method (FCM)  and  polynomial  neural  networks (PNNs). *Neural  Networks*  60:  166–181 (2015)

[12]. Hilmi Berk Celikoglu, Application of radial basis functions and generalized regression neural networks in non-linear utility fun-ction specification for travel mode choice modeling. *Mathematical  and  Computer Modelling*  44: 640–658 (2006)

[13].  Seema  Singh  and  H.  Sushmitha,  An Efficient  Neural  Network  Based  System for Diagnosis of Breast Cancer. *Interna-tional Journal of Computer Science and Information Technologies* 5(3): .4354-4360 (2014)