

## Efficient Subspace Clustering of High Dimensional Breast Cancer Gene Set for Multi Variant Gene Population (MVGP) Using Fuzzy Rule Sets and Multi Gene Impact Matrix

<sup>1</sup>N. Magendiran and <sup>2</sup>S. Selvarajan

<sup>1</sup>CSE at Paavai Engineering College, Namakkal, Tamilnadu, India

<sup>2</sup>College of Engineering, Rasipuram, Tamilnadu, India

**Abstract:** The breast cancer is the most threatening factor of women's lifestyle and the reason of the disease has many factors. Still the gene factor has more influence in the generation of breast cancer where the early diagnosis and prevention is essential. There are many approaches has been discussed in the literature, but the identification and selection of a set of genes that influence the disease are still becoming complicated. A multi-variant approach for gene selection is proposed, by performing high dimensional subspace clustering. With the given data set, the method generates a set of rules and unlike generic fuzzy rules the method splits the range values into the number of parts and based on that the rules are generated. Also, according to the different range values, the method generates a multi-gene impact matrix where the frequency of range values of each rule is stored. The data set is clustered according to the generated rules and from the generated rules the gene selection is performed. For the gene selection, the method computes the multi-gene frequency measure that represents how depth the gene has an impact on the classification of disease. The proposed method produces an efficient classification of genes in the influence of breast cancer.

**Key words:** Gene Selection • High-Dimensional Clustering • Multi Gene Impact Matrix • Fuzzy Rule Sets

### INTRODUCTION

The growth of data sets in their dimension increases the challenges in cluster them, where the higher dimensional space requires more sophisticated approaches to clustering the data sets. In any high dimensional space, identifying the subspace is the most important task that has to be performed in an efficient manner. For a breast cancer data set, there are a number of genes influencing or taking part in the appearance of the cancer in the women. To identify them or to cluster such data set the genes are the most important factor that participates in the clustering approach.

Gene selection in high dimensional breast cancer dataset clustering is the most important task and how the gene selection is performed is the big question here. Not all the genes have a great impact, but all the genes have some impact on the cause of breast cancer. To find a strategic approach to selecting the gene selection there must be some efficient approach to be there.

Unfortunately, the existing approach misses the case of gene selection in a modern, sophisticated approach and has no efficient solution to perform the task of gene selection.

The breast cancer can be classified into many cases and to identify the exact subspace we must come up with more efficient measures and gene selection approaches. For example, a subset of genes may be the cause of a specific type of cancer, but they may not have any impact in the presence of another type of cancer. So the gene selection is the most important task that could be used to predict the future appearance of breast cancer. So for the prediction of the breast cancer the gene selection approach can be used which helps early detection and cure of cancer in many ways.

The general fuzzy rule sets are nothing but the range of values for each rule. It contains a number of rules for each case of breast cancer and has a range of values for each gene participated in the appearance of cancer cells. Unlike generic one, the method generates a modern fuzzy

rule set that has many numbers of rules. For each rule, each gene value is split into a number of subdivisions and the method generates many numbers of rules to fill the rule sets.

From the generated fuzzy rule sets, the method computes a multi-gene impact matrix, which represents the impact of genes on the occurrence of cancer. From the rule set with the data set, of each subdivision and rules available, the impact matrix is generated by the values of the data set. The method computes the number of possible occurrences of cancer appearance based on the genes selected or the values of genes. According to the appearance, a single value for the pattern of the genes and values is generated which decides the impact of the gene and their values for the occurrence of cancer in the human body.

**Related Works:** There are many approaches has been discussed for the gene selection of breast cancer and the paper discuss a few of them here in this section for better understanding of breast cancer gene selection.

Stable Gene Selection from Microarray Data via Sample Weighting [1], proposed a framework of sample weighting to improve the stability of feature selection methods under sample variations. Their experiments show that the sample weighting algorithm improves the stability of gene selection. It evaluated with SVM-REF, Relief classifiers.

A Novel Approach for Single Gene Selection Using Clustering and Dimensionality Reduction [2], proposed Hybrid Fuzzy C Means-Quick Reduct algorithm for single gene selection. Average Correlation Value (ACV) is calculated for the high class discriminated genes. The algorithm is evaluated using WEKA classifier with leukemia cancer data set.

Relevant and Significant Supervised Gene Clusters for Microarray Cancer Classification [3], have developed mutual information based supervised gene clustering (MSG) algorithm to form the reduced gene clusters for cancer classification. The approach has been evaluated using different microarray cancer data sets with different classifiers like Naïve Bayes, K-nearest rule and SVM.

Gene-Expression-Based Cancer Subtypes Prediction through Feature Selection and Transductive SVM [4] proposed combined gene selection and transductive support vector machine (TSVM). Their method identified the potential genes and used the TSVM to improve the prediction accuracy compared to standard inductive SVM.

Experimental results confirm the effectiveness compared to the ISVM and low-density separation method in the area of semi-supervised cancer classification as well as gene-marker identification.

Identifying Gene Pathways Associated with Cancer Characteristics via Sparse Statistical Methods [5] proposed statistical method for uncovering gene pathways that characterize cancer heterogeneity. They define a set of activities of pathways from microarray gene expression data based on the Sparse Probabilistic Principal Component Analysis (SPPCA). It creates a novel gene-gene association relating to the cancer phenotypes. This method analyzes breast cancer gene expression data.

An optimal gene selection based on search mechanism has been adapted for diagnosis of cancer in Optimal Search-Based Gene Subset Selection for Gene Array Cancer Classification [6]. The method introduces tabu search (TS) to gene selection from high-dimensional gene array data. TS were shown to be a promising tool for gene subset selection.

A Family wise error rate based gene selection approach has been discussed in Gene Selection for Sample Classifications in Microarray Experiments” [7], which uses two or multiple samples for the classification of genes. The method reduces the false positive ratio that has been evaluated using colon cancer data set.

Gene Expression Profiles for Predicting Metastasis in Breast Cancer: A Cross-Study Comparison of Classification Methods [8], performs a comparative analysis of different methods. The method uses a voting mechanism for the prediction of metastasis. The voting based method produces efficiency with the breast cancer patients.

Comparison of feature selection methods for cross-laboratory microarray analysis [9], investigate four feature selection methods. They are t-Test, Significance Analysis of Microarrays (SAM), Rank Products (RP) and Random Forest (RF) across breast cancer and lung cancer microarray data which consists of three cross lab data sets each. Their results show that SAM has the best classification performance. RF also gets high classification accuracy, but it is not as stable as SAM. The Test performance is the worst of the four methods.

Data Mining Techniques for the Identification of Genes with Expression Levels Related to Breast Cancer Prognosis [10], applied data mining techniques to study the gene expression values of breast cancer patients with

known clinical outcomes. Created the classification models for clinical practice to support the rapid prescription. With nine algorithms of feature selection, they extracted a group of subsamples of data, which was analyzed with different classification algorithms for comparison purpose. They used five learning algorithms implemented in YaLE or WEKA. Classifying a patient as “good prognosis” when she is in a state that will develop metastasis (i.e., an FN error) is much more serious than classifying a patient as “poor prognosis” when she is not in a state that will develop metastasis (i.e., an FP error). The algorithm classified ill patients more accurately (lower FN and higher TP) at the expense of the classification of healthy patients (higher FP and lower TN).

All the above-discussed methods have the problem of gene selection efficiency and could not perform any prediction about the breast cancer identification and the reason for them.

**Preliminaries in High-dimensional Clustering:** Let the data set  $D_s$  has  $N$  number of data points from  $\{Dp1, Dp2, \dots, Dpn\}$  where each data point has  $Q$  dimensions from  $\{D1, D2, D3, \dots, Dq\}$  and the number of attribute types from  $A1$  to  $Am$ . Clustering such high dimensional data set can be performed based on the similarity measure  $HSim$ .

The high-dimensional similarity measure  $Hsim$  is computed by computing the similarity or closure value of data points at each dimension as follows.

$$Hsim = \{S1(Dpx, Dpy), S2(Dpx, Dpy), S3(Dpx, Dpy), \dots, Sq(Dpx, Dpy)\} \quad (1)$$

From the equation (1), the variables  $S1, S2, S3$  represent the similarity measure on a particular dimension. To identify the cumulative similarity measure, the method can perform the averaging method or standard deviation or any other mathematical approach.

Finally, the cumulative similarity measure can be computed using the equation (2).

$$CSim = \frac{\sum_{i=1}^Q HSim(i)}{Q} \quad (2)$$

Equation (2) shows how to compute the similarity between two data points. However, when the method cluster data points with high dimensions, in order to assign a data point to a group of data points  $C_s$ , the similarity measure is performed with all the data points.

For example, there exist  $C$  number of clusters with each cluster has variable number of data points  $Dp$ . To identify the group of data point to which the testing sample  $T_s$  belongs to, the method must compute the similarity measure  $k$  as follows:

$$K = \frac{\sum_{i=1}^{Size(C)} HSim(i)}{Q} \quad (3)$$

Here  $K$  is the similarity measure that represents closeness value towards the cluster  $C_i$ ; similarly the closeness of the data point towards each cluster  $C_x$  can be used to identify the final class of the data point.

**Gene Selection from High Dimensional Space:** The breast cancer dataset is categorized into many classes, as of the data points are high dimensional they are clustered based on any similarity measure. Even though the clustering is performed based on all the dimensional similarity, there is only little number of dimensions and their values which decides the closure of the data point.

If each dimension is about a gene and the data point has numerous numbers of genes, among them only a small set of gene values which decides the category of the data point. When the method classifies the data sample towards number of classes, in order to come into a class, the gene values at a data point has to be close to the gene values of the data points available in the cluster. By identifying such small set of genes, which has more influence in getting assigned to the class, the causes and the specific small set of genes can be monitored in disease prediction.

**Multi Variant Gene Population Technique:** For any given gene set  $G_s$ , with dimension  $N$ , belongs to  $K$  number of diseases, the problem of identifying a small set of genes called impact genes  $I_g$ , can be performed by using any population technique. The population technique is one, which choose a small set of genes from a large set of genes which has more impact or more influence over a class of disease. A multi-variant gene population technique is proposed which selects a small set of genes based on their influence on any disease. The proposed multi-variant gene selection using fuzzy sets and impact matrix approach has the following stages, namely preprocessing, rule generation, Multi Gene Impact Matrix Generation, Multi Variant Gene Selection. We discuss each of them in detail in this section.

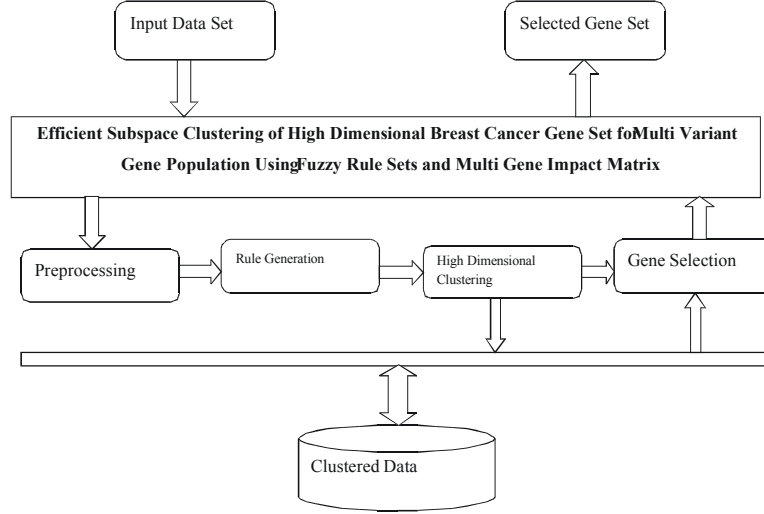


Fig. 1: Proposed System Architecture

Figure 1 show the architecture of the proposed multi-variant gene selection approach and its functional components which will be discussed in detail in this section.

**Preprocessing:** The given data set  $D_s$ , has  $N$  number of genes with  $M$  number of instances. At this stage, the method identifies the number of genes present in the data set and for each of the gene from each of the data points in the dataset, the presence of value is identified. If there is, any data point gene with missing values are identified and the data point will be removed from the data set. The noise removed data set is given for the rule set generation that will be used for further processes.

#### Preprocessing Algorithm:

**Input:** Data Set  $D_s$ .

**Output:** Noise Removed Data Set  $NDS$ .

Identify Number of genes available in the data set  $D_s$ .

$$\text{Gene set } GS = \sum_{i=1}^{Size(D_s)} \sum Gene(D_s(i)) \neq GS \quad (4)$$

The equation (4) identifies the distinct genes available in the data set  $D_s$ .

Perform Noise Removal

For each gene  $G_i$  from  $G_s$

For each data point  $D_{s_i}$  from  $D_s$

If  $D_{s_i}(G_i) == 0$  then //the gene value is 0 then remove the data point

$$D_s = \sum D_s - D_{s_i} \quad (5)$$

The equation (5) removes the noisy data point from the data set  $D_s$ .

Else

End

End

$NDS = ?NDS + D_{s_i}$  //Add the correct data point to the data set.

End

The above algorithm performs the preprocessing of input data set by identifying the unique genes and for each data point, the presence of all the gene values are verified. If there is any missing cases or has no value, then they are removed from the data set.

**Rule Generation:** The rule generation is performed based the values of genes present in the data set. Rule generation is performed using the preprocessed data, for each gene value the method identifies the maximum and minimum values and from them number of range values are generated by splitting them. Based on generated multiple range values, the method generates a number of rules using which the multi-gene impact matrix is generated and clustering is performed.

#### Rule Generation Algorithm:

**Input:** Input Data Set  $D_s$ .

**Output:** Rule Set  $R_s$ .

Identify set of genes from data set  $D_s$ .

$$\text{Gene set } GS = \sum_{i=1}^{Size(D_s)} \sum Gene(D_s(i)) \neq GS \quad (6)$$

The equation (6) identifies the set of distinct genes from the data set.

For each of the gene identified, the minimum and maximum values are identified as follows.

For each gene  $g_i$  from  $G_s$

Compute minimum and maximum values of gene  $G_i$ .

$$G_{min} = \min(\sum_{i=1}^{size(Ds)} Ds(i)(G_i)) \quad (7)$$

The equation (7) computes the minimum value of each gene  $G_i$  from the data set  $D_s$ .

$$G_{max} = \max(\sum_{i=1}^{size(Ds)} Ds(i)(G_i)) \quad (8)$$

The equation (8) computes the maximum value of the gene  $G_i$  from the data set  $D_s$ .

End

For each gene  $G_i$  generate set of range values.

$$RV = \sum_{i=1}^{size(Gs)} \forall (Range\ Values) \quad (9)$$

The equation (9), generates set of range values between minimum and maximum values of each gene  $G_i$  present in the gene set  $G_s$ .

End

For each Gene  $G_i$  and range values

Generate rules  $R_s$ .

$$R_s = \sum_{i=1}^{size(Gi)} \Sigma Rv(u) \quad (10)$$

Equation (10), generates number of rules for each gene  $G_i$  from gene set  $G_s$  and the number of rule for each gene  $G_i$  depends on the number of range values generated by the equation (9).

End

The above algorithm first computes the minimum and maximum values for each of the gene present in the data point. With the computed minimum and maximum values, the method generate set of range values for each gene attribute. Based on computed range values, number of rules is generated to perform clustering. The total number of rules generated by the rule generation approach could be computed as follows:

$$\text{Total number of rules } X = M \times 2^N \quad (11)$$

The equation (11) computes the total number of rules computed using the gene set  $G_s$ , where  $M$  represent the number of range values  $N$  represent the number of the combination.

**Multi Variant Impact Matrix:** The multi-variant impact matrix is generated using the rule set generated and based on the rule generated with the data set available the method compute the multi-variant impact matrix. For each data point of the data set, the method performs matching of a rule with the other gene values. Based on the matching rule, the method computes the multi-variant support values and with the total number of data points, the method computes the impact factor. Similarly for each of the rules identified the methods generate the impact matrix and the generated impact matrix is used to cluster the data points of the data set.

**Impact Matrix Generation Algorithm:**

**Input:** Rule Set  $R_s$ .

**Output:** Multi variant impact matrix MVIM.

For each rule  $R_i$  from  $R_s$

Compute number of data points match with the range values of genes.

Gene count  $G_{sup}$ .

$$G_{sup} = \sum_{i=1}^{size(Ds)} Ds(i).Gi.value > Ri.min \& \& Ds(i).Gi.value < Ri.max \quad (12)$$

The equation (12) computes the gene support value that represent the total number of data points has the value falls within the range of the gene value of the rule.

$$\text{Compute impact value } G_{imp} = \frac{G_{sup}}{size(Ds)} \quad (13).$$

The equation (13) computes the impact value of each gene using computed gene support value and the total number of data set.

End

$$MVIM(i) = R_i + G_{imp} \quad (14).$$

The equation (14) adds the impact value to the static variable and adds to the multi-variant impact matrix, which contains the impact value of all the genes.

The above algorithm computes the impact matrix, where for each data point, the impact matrix is computed using the rule set available.

**High-Dimensional Clustering:** The clustering is performed based on the multi-variant impact matrix generated and with the input data set. For each data point, from the rule set available, the matching rule is identified and from the impact matrix the values of genes are identified. For each level, the method finds the matching cluster or subspace to identify the cluster to which the data point belongs. With the data points of each cluster or subspace, the input data point is computed with the multi-variant similarity measure and the deviation with the multi-variant impact factor is performed. If the deviation is more than the process will look for the next cluster and finally a more closure data cluster is identified and the data point is assigned to the selected class label.

#### High-Dimensional Clustering Algorithm:

**Input:** Rule Set Rs, Dataset Ds, MVIM.

**Output:** Cluster Cs.

For each data point Di

    Identify the matching rule Ri.

    For each level of the gene value

        Identify the cluster Cs.

        Compute multi-variant similarity MVS.

$$Mvs = \sum_{i=1}^{size(Cs)} \forall (dp(csi)) \square Di \quad (15)$$

The equation (15) computes the multi-variant similarity value for each data point and the matching rule Ri.

$$\text{Compute impact deviation } Idev = Ri(Imp) - Cs(Imp) \quad (16)$$

The equation (16) computes the impact deviation between the rule and the cluster.

If Idev < Th then

    Assign di with the class Cs.

    End

    End.

End.

The above discussed high dimensional clustering algorithm performs the grouping of similar gene data points into a class. For each data point available, for each of the gene present in the data point, a multi-variant similarity is computed and based on that an impact deviation is computed. Finally, if the impact deviation is less than the threshold, then the data point is assigned to the class.

**Gene Selection Approach:** With the available information, the other gene values with the data points available in the other cluster are identified. The gene values that are not more similar to the other cluster data point are identified. For each gene present in the data point, the method identifies the impact of the gene in the impact matrix and using the values of genes the impact of the gene in the other cluster is identified. For the number of data points from the each cluster, the data points closure to the available gene value is identified and the number of data points with the same gene value of another cluster is identified. These details are used to compute the impact of the gene value.

#### Gene Selection Algorithm:

**Input:** MVIM, Data point Dp, Cluster Cs.

**Output:** Gene selection Gs.

For each cluster Cs

    For each gene value Gv of Dp

        For each cluster cs

            Compute number of data points present within the range.

$$Nodp = \sum_{i=1}^{size(Cs)} Gs == C(i)(gv)$$

    Compute number of data points present within range in another cluster.

$$Nodpo = \sum_{i=1}^{size(OCs)} Gs == OCs(i)(gv)$$

    If Nodpo > Nodp

        Else

        Add to list Gs.

        End

    End

End

The above-discussed procedure shows how the gene selection is performed. In this procedure, for each cluster available and for each gene value, the method identify the number of data points falls within the range values and identify the number of data points falls within the range of another cluster. Using both the values, that if the gene has number of data samples in the current category than the other category then the gene is selected [11-14].

## RESULTS AND DISCUSSION

The proposed multi-variant gene selection approach has been implemented and tested in Matlab environment and the proposed method has produced efficient results in clustering and the method has produced efficient results in time complexity and accuracy of clustering. Also, the method has produced efficient results in gene selection and reduced the time complexity also.

Table 1 shows the details of the data set being used to evaluate the performance of multi-variant gene selection approach. The data sets listed in Table 1 has number of data points and each has a different number of classes. The data point has number of dimensions and each dimension has different attribute type. The method has been evaluated with each of the data set where each has various dimensions and the number of samples with data types is also different.

Each data set has been used to evaluate the performance of the proposed method for clustering

accuracy. From the data set 80 percent of the data set has been used as training set and remaining 20 percent of data points are used to perform the evaluation. The method produces efficient results with each of the data sets and the gene selection approach has produced less time complexity with higher detection accuracy.

The clustering efficiency shows the accuracy of the cluster being produced by any clustering algorithm that can be measured by computing the true positive and true negative, false positive and false negative values.

The true positive is the value of clustering that is computed using the number of samples given and the number of correct label assignment performed. Similarly, the true negative value is computed using the number of samples of a class that is assigned with other class labels. The False positive is the ratio computed using the number of other class samples assigned to a specific class from given set of samples.

$$\text{True Positive } Tp = \frac{\text{Number of samples indexed correctly}}{\text{Number of samples given for testing}} \times 100$$

$$\text{True Negative } TN = \frac{\text{Number of samples indexed incorrectly}}{\text{Number of samples given for testing}} \times 100$$

$$\text{False Positives } Fp = \frac{\text{Number of samples of other class assigned with the current class}}{\text{Number of samples given for testing}} \times 100$$

$$\text{False Negative Ratio } Fn = \frac{\text{Number of samples of current class assigned with the other class}}{\text{Number of samples given for testing}} \times 100$$

Table 1: Details of the data set being used

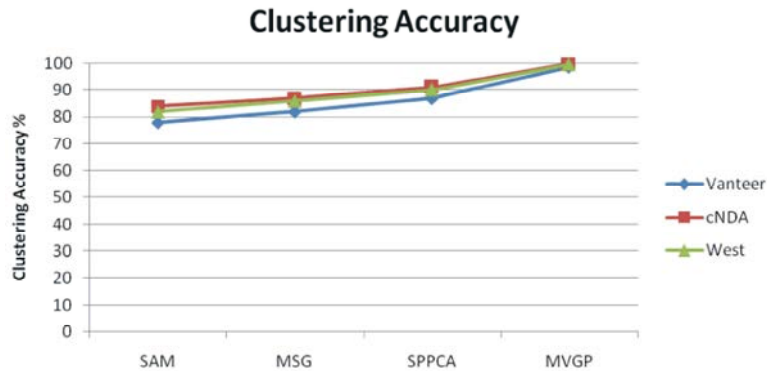
Dataset	Number of Data Points (N)	Attributes (d)	Attribute Values (AA)	Classes (K)
Vanteer	19	24481	1	2
cDNA	34	176	1	2
West <i>et al.</i> ,	49	7129	1	2

Table 2: Comparison of gene selection produced

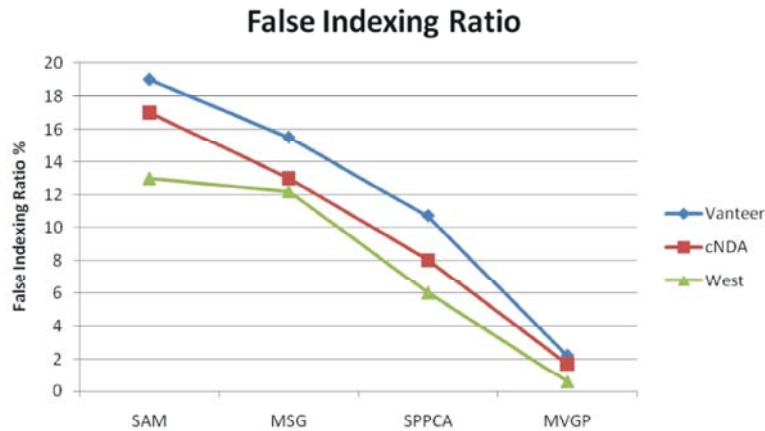
Methods	Number of genes selected from different data set		
	Vanteer	cDNA	West
SAM	241	22	87
MSG	198	19	43
SPPCA	113	14	26
MVGP	77	3	14

Table 3: Comparative analysis results on different data set

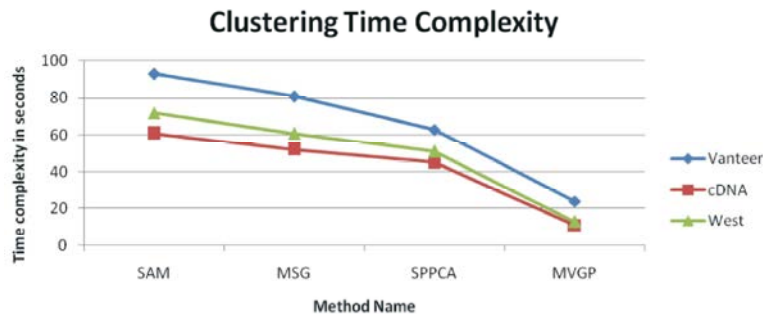
Methods	Accuracy %			False Indexing %			Time Complexity in seconds		
	Vanteer	cDNA	West	Vanteer	cDNA	West	Vanteer	cDNA	West
SAM	81	83	87	19	17	13	93	61	72
MSG	84.5	87	87.2	15.5	13	12.2	81	52	61
SPPCA	89.3	92	94	10.7	8	6	63	45	51
MVGP	97.8	98.3	99.4	2.2	1.7	0.6	24	11	13



Graph 1: Comparison of clustering accuracy



Graph 2: Comparison of false indexing ratio



Graph 3: Comparison of time complexity

Classification ratio is computed using all the above-discussed measures like  $T_p$ ,  $T_n$ ,  $F_p$  values. The clustering efficiency of any algorithm can be computed as follows:

$$\text{Clustering efficiency CE} = \frac{TP}{TP + TN} \times 100$$

The Graph 1 shows the comparison of clustering accuracy produced by different methods while using two



different data sets. It shows clearly that the proposed method has produced efficient clustering than other approaches.

The false indexing ratio is computed using the values of false positive and the false negative ratio that is computed using the formula described below.

$$\text{False Indexing Ratio FIR} = \frac{(FP + FN)}{\text{Total number of samples}} \times 100$$

The Graph 2, shows the false indexing ratio produced by different methods and it shows clearly that the proposed method has produced less false ratio than other methods.

The Graph 3 shows the time complexity produced by the proposed method of clustering different data sets. It shows clearly that the proposed approach has produced less time complexity in all the data sets where each of them varies with the dimensions and number of samples.

Table 2 shows the number of genes selected as influencing by all the method being considered. It shows clearly that the proposed method has produced good accuracy and number of genes selected is less.

Table 3 shows the comparative results produced by the different methods of clustering accuracy, false indexing and time complexity while using the all the three data set being considered for the evaluation. It shows clearly that the proposed method has produced efficient results than other methods.

## CONCLUSION

The proposed multi-variant gene selection approach works based on multi-variant impact matrix to develop the clustering accuracy. The method preprocesses the data points and generates the rule sets according to number of range values available at each data point. With the help of the data points and the rule sets, the method generates the multi-factor impact matrix that is used to identify the class of data points. Similarly, based on the multi-variant impact matrix and rule sets, a set of the gene is selected based on computed similarity measure values. The similarity value is computed using all the gene values, but the similarity value becomes higher only when the data points have more similar values for each dimension. However, in case of gene selection, the attribute gene values are get selected only when they are deviate with others. The gene selection approach proposed populate minimum set of genes which has more impact on the occurrence of breast cancer and helps monitoring, the diagnosis of breast cancer. The proposed method

increases the accuracy of clustering and reduces the time complexity highly. Also, the proposed method reduces the false indexing ratio produced by other methods also. The gene selection approach can be further improved by adapting other relational measures that could be computed using the gene values.

## REFERENCES

1. Lei Yu, Yue Han and Michael E. Berens, 2012. Stable Gene Selection from Microarray Data via Sample Weighting, *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 9(1): 262-272.
2. Sathishkumar E.N., K. Thangavel and T. Chandrasekha, 2013. A Novel Approach for Single Gene Selection Using Clustering and Dimensionality Reduction, *International Journal of Scientific & Engineering Research*, 4(5): 1540-1545.
3. Pradipta, Maji and Chandra Das, 2012. Relevant and Significant Supervised Gene Clusters for Microarray Cancer Classification, *IEEE Transactions on Nanobioscience*, 11(2): 161-168.
4. Ujjwal Maulik, Anirban Mukhopadhyay and Debasis Chakraborty, 2013. Gene-Expression-Based Cancer Subtypes Prediction through Feature Selection and Transductive SVM, *IEEE Transactions on Biomedical Engineering*, 60(4): 1111-1117.
5. Shuichi Kawano, 2012. Identifying Gene Pathways Associated with Cancer Characteristics via Sparse Statistical Methods, *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 9(4): 966-972.
6. Li Jiexun, Hua Su, Hsinchun Chen and Bernard W. Futscher, 2007. Optimal Search-Based Gene Subset Selection for Gene Array Cancer Classification, *IEEE Transactions on Information Technology in Biomedicine*, 11(4): 398-405.
7. CHEN-A TSAI, 2004. Gene Selection for Sample Classifications in Microarray Experiments, *DNA and Cell Biology*, 23(10): 607-614.
8. Burton Mark, Mads Thomassen, Qihua Tan and Torben A. Kruse, 2012. Gene Expression Profiles for Predicting Metastasis in Breast Cancer: A Cross-Study Comparison of Classification Methods, *The Scientific World Journal Volume 2012*, Article ID 380495, pp: 11.
9. Liu Hsi-Che, 2013. Comparison of feature selection methods for cross-laboratory microarrayanalysis, *IEEE/ACM Transactions on Computational Biology and Bioinformatics*.

10. Giarratana Gabriele, 2009. Data Mining Techniques for the Identification of Genes with Expression Levels Related to Breast Cancer Prognosis, Ninth IEEE International Conference on Bioinformatics and Bioengineering, pp: 295-300.
11. Bessarabova Marina, 2010. Bimodal gene expression patterns in breast cancer, BMC Genomics, Supplementary, pp: 1.
12. Bevilacqua Vitoantonio, 2012. Comparison of data-merging methods with SVM attributes selection and classification in breast cancer gene expression, BMC Bioinformatics.
13. Tyrer Jonathan, Stephen W. Dey and Jack Cuzick, 2004. A breast cancer prediction model incorporating familial and personal risk factors, Statistics in Medicine Statist. Med., pp: 1111-1130.
14. Jorgen Aaroe, 2010. Gene expression profiling of peripheral blood cells for early detection of breast cancer, Breast Cancer Research.