



EFFICIENT CLASSIFIER ALGORITHM FOR GENE EXPRESSION DATA ANALYSIS

C. Kondalraj¹, Dr. R. Murugesan²

¹Assistant Professor, Department of Computer Science & Information Technology

²Associate Professor of Computer Science

^{1,2}CPA College (Affiliated to Madurai Kamaraj University), Bodinayakanur

¹kondalrajc@gmail.com, ²rmncpa90@gmail.com

Volume 6, Issue 10, May 2024

Received: 09 May 2024

Accepted: 19 June 2024

Published: 20 July 2024

doi:

[10.48047/AFJBS.6.10.2024.7311-7321](https://doi.org/10.48047/AFJBS.6.10.2024.7311-7321)

Abstract: In this paper we propose an efficient and effective classifier algorithm for gene expression data analysis. In this research work, the proposed PCA approach is used on improvisation on feature extraction. Once the feature extraction is completed, using Multi Algorithm Fusion (MAF) method is investigated and performed. Also, the Polynomial Support Vector Machine (PSVM) has been related to MAF which supports feature extraction. From this, the absolute weight of SVM, fisher ratio and PSVM are attained. Finally, random forest classifier was utilized for the effective classification approach and it was proved with efficient outcome of measures such as error rate, accuracy, specificity, precision, recall, F1score, false positive rate.

Index Terms— PCA- Principal Component Analysis, RF- Random Forest Classifier, Polynomial Support Vector Machine (PSVM), MCF-Multi Algorithm Fusion.

INTRODUCTION

The process of accumulating and excavating information from enormous amounts of data to discover remarkable patterns and their correlation with extracted data is known as data mining. It entails smart technologies and amenability to investigate the probability of concealed information inside the data. Nowadays, the tangible world data are immense and of poor quality. Mining

such data will not provide suitable and expedient results because those data are vulnerable to noise. Occasionally, it is observed that some fields are partly vanished. Hence, pre-processing is crucial before data mining which involves cleaning, integrating, normalizing and reducing of data. There are different data sets from which the data are exhumed. Feature selection is a method of finding the most pertinent features from the data set and signifying the high dimensional data to a smaller extent. The chosen features are extricated and categorized for various analyses. The basic steps involved in data mining are illustrated in the below figure [1].

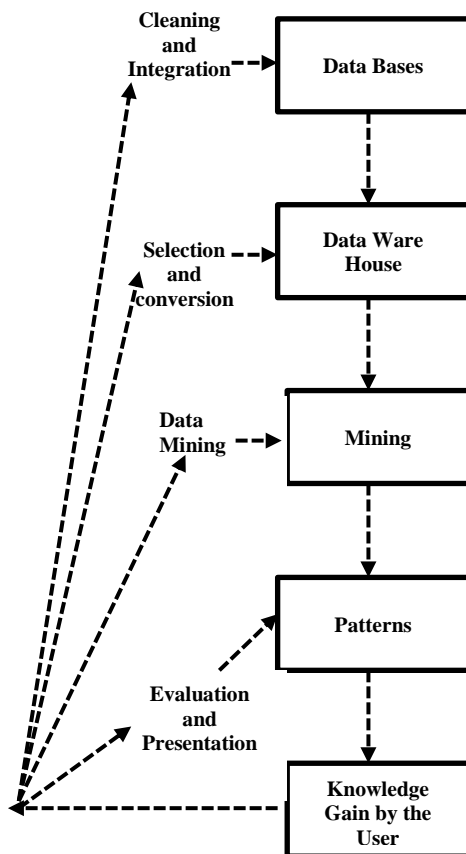


Fig. 1 Steps involved in data mining

The proposed work ponders on an important concept named as Gene selection. It plays a significant role in colon cancer data classification to lessen the number of noisy and inappropriate genes and to choose the allied genes for enhancing the classification results.

The main objectives of this work are:

- To improve feature selection process using PCA method.
- To extract the features by utilizing MAF-PSVM approach.
- To classify genes by means of Random forest classifier.

This paper is organized as follows:

Section I provides the general introduction of data mining, feature selection and gene classification.

Section II reviews the works that are closely related to gene expression data mining.

Section III describes the proposed work of colon and cancer dataset that utilizes data mining algorithms for extraction, selection and classification.

Section IV deals with the performance analysis of various factors such as classification error, gene classification accuracy and feature robustness. And finally,

Section V concludes the paper along with the references to future work.

I. RELATED WORKS

The gene expression data analysis is an important aspect for interpreting diseases and plays vital role in understanding diseases and discerning medications for them. At all times, gene expression data was prone to high dimensionality. Due to this concern, feature selection became the essential mechanism for classifying cancer. [2]J.c.Ang suggested a Semi-Supervised SVM based Feature Selection (S³VM-FS) that utilized information from labelled and un-labelled gene expression data. The lung cancer data were

experimented and it was observed that this method attained a maximum accuracy but the processing time was longer than the familiar methods like SVM-RFE and S³VM-RFE.

The cancer clusters were demarcated and patient genes were classified for predicting cancer by means of micro-array systems. [3]J.X.Liu offered an innovative technique for organizing the tumour samples from gene expression data based on Robust Principal Component Analysis (RPCA). Initially, the RPCA emphasized typical genes related to a distinctive biological method. After that, RPCA and RPCA combined with Linear Discriminant Analysis (LDA) were employed to ascertain their characteristics. As a final point, Support Vector Machine (SVM) was pertained to categorize the samples based on the applied to classify the tumor samples of gene expression data based on the realized features. The techniques were efficient and viable in tumour grouping. Robust feature selection detected the appropriate genes for coherent sample classification. The prevailing selection methods were affluent and lay-off in gene expression data resulted in poor accuracy that deprived multi-class classification. [4]J.Bennet introduced an ensemble feature selection incorporated with RFE and BBF for selecting genes and utilized SVM algorithm for classifying them. The ensemble method relented an improved performance than the other classifiers and afforded a novel perception in feature selection.[5]T.Latkowski presented the data mining algorithms to distinguish the substantial genes and their structures from micro-array gene expression data set of autism. Certain feature selection techniques were chosen and their respective outcomes were assimilated. The numerical results of gene selection and classification were deliberated from which a synthesis of different selection methodologies concomitant with autism. An exceptional process to calculate the number of top-ranked genes that were availed in classification was also included.

Bi-clustering procedures evaluated the gene expression data proficiently by realizing a gene group having stabilized expression configurations with definite conditions. [6]Z.Wang initiated an uncomplicated method in which the bi-clusters were found from large, noisy and intricate data. Longest Common Subsequence (LCS) configuration nominated row duos in a key matrix obtained from a data matrix to uncover the kernel of each bi-cluster. Several bi-clustering algorithms were verified on gene expression dataset and on comparison, it was proven that the UniBic algorithm outclassed all other earlier processes in discovering the constricted bi-clusters. The surviving Rotation Forest Algorithm (RFA) based methods concentrated only on the accuracy of gene classification and ignored the cost of classification. [7] H.Lu recommended cost-intuitive RFA for arranging gene expression data. The types of classification costs such as misclassification, test and rejection were analysed in this work. The tentative results confirmed that the proposed algorithm minimized the classification cost efficiently.[8]S.Cogill proposed a classification strategy to recognize the autism related disorders making use of expression patterns in relation to brain development. The acquaintance of familiar protein-coding genes was swapped with all gene varieties. The method exactly classified and ordered the genes in terms of risks associated with autism. This gene classification was reliable than the former techniques.[9]H.Lu put forward a composite feature selection algorithm by merging AGA and MIM. The investigational fallouts indicated the supremacy of MIM-AGA selection in decreasing the facet of gene data and eliminating classification redundancy. The compact dataset achieved greatest accuracy when compared to the traditional selection algorithms. Four classifiers were smeared with abridged dataset to exhibit the vigour of projected scheme. A huge quantity of gene expression data were collected from biological experimentations by applying the microarray technology. The collected data were then explored successfully with the help of some clustering algorithms. In the next step, co clustering was employed for recovering the co-clusters present in the gene expression data and with this the group were categorized further into sub groups. Subspace Weighting Co-Clustering (SWCC) was proposed for the purpose of co-clustering the high dimensional gene expression data [10].X.chen Then a novel co-clustering function was framed with the objective of recovering the co-clusters where the subspace weight matrix was employed. The work [11] focused on reducing the dimensioning of features present in the categorization by proposing a novel strategy of weighted feature selection that was based on bacterial algorithm. It differentiates the features and classifies them based on their performances and occurrence frequency of population. The efficiency of the proposed work was analyzed by 4 bacterial dependent techniques. In addition to that, 3 famous algorithm that depends on population which utilizes 15 different datasets of micro-array that possess unique classes and features. As a result, the embedded and weighted feature selection strategies have enhanced the capacity of feature selection of some bacterial algorithms.

II. PROPOSED WORK

The proposed work in this section deliberates the methodologies involved in the proposed work and the flow diagram as well. The intelligent technologies are required for exploring the likelihood of hidden knowledge existing over the data. Initially, preprocessing step is utilized for eliminating the unwanted data. There may be a noise presence in some of the data, so preprocessing step is essentially required to remove all the unwanted data. The elimination of unwanted data processing has been accomplished with preprocessing techniques, which also reduces computation time. Feature selection is the significant step that has to be concerned more for efficient gene selection process. The selection of appropriate features are the essential step which is performed using an operative PCA approach. After completing the selection process, best features has to be extracted for undertaking classification process. In the microarray gene expression data analysis two feature extraction approaches are employed i.e. ranking-based feature extraction and set-based feature extraction. Under this ranking based approach, the features are selected based on the individual term not concerning inter relationship among the features. The set-based feature extraction is the feature extraction process, where the dependency among the features present in feature set is considered. The robustness is the significant problem often occurred from the generalization ability of feature extraction in learning methods. To improve the extraction efficiency, a multi algorithm fusion technique has been implemented. The polynomial support vector machine related to multi-algorithm fusion is considered for the feature extraction process. The absolute weight of SVM, fisher ratio and Polynomial Support Vector Machine (PSVM) were utilized in this fusion

approach. An effective feature extraction is accomplished with this multi algorithm fusion technique. The multi algorithm fusion based PSVM algorithm is chosen for the efficient outcome.

A. DATA PREPROCESSING

In this preprocessing stage, an unwanted data is removed from the preferred dataset. The removal of noise and white space is achieved and the null values are also removed to retrieve an effective significant information from the dataset.

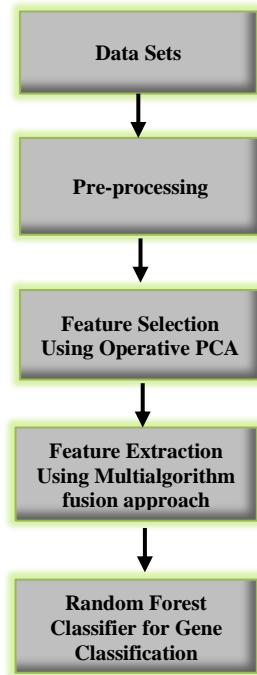


Figure 1: Proposed methodology

B. FEATURE SELECTION USING OPERATIVE PCA APPROACH

In this gene expression approach, a Principal Component Analysis (PCA) based dimensionality reduction approach is proposed. In this operative PCA approach the information is plotted linearly to the lower-dimensional space in such a manner the maximum data variance is caught over the low-dimensional depiction. Generally, the eigenvector matrix is being computed through generating the correlation or the co-variance data matrix. From this the reconstruction of large amount of original data variance is completed which is corresponding to higher eigenvalues. Additionally, certain initial eigenvectors can frequently be understood in relations of the significant physical performance of the system. The new space is abridged together with the lost data where the most significant variance traversed by some eigenvectors are retained. The following algorithm is utilized for the creating input data for PCA.

Operative PCA-Based Feature Selection

Input: PCA Input Data

Output: Feature Selected Attributes

Procedure:

Initialize Variables:

All variables as double.

Calculate Mean:

III.

IV. READ FILE AND CALCULATE MEAN:

V. $OUT += D / ENTRIES.LENGTH$ FOR EACH D IN ENTRIES.

VI. RETURN OUT.

VII. CALCULATE COVARIANCE:

VIII.

IX. INITIALIZE VARIABLES: $SUM = 0$, $AM = MEAN(A)$, $BM = MEAN(B)$, $DV = A.LENGTH - 1$.

X. LOOP THROUGH OUT.LENGTH TO CALCULATE SUM:

XI. $SUM += A(I) + AM * B(I) * BM$.

XII. IF $SUM == 0.0$, ASSIGN VAL TO SUM.

XIII. RETURN SUM / DV .

XIV. CALCULATE COVARIANCE MATRIX:

XV.

XVI. INITIALIZE MATRIX $OUT = [MAT.LEN][MAT.LEN]$.

XVII. LOOP THROUGH OUT.LENGTH TO FILL MATRIX:

XVIII. $DTA = MAT[I]$, $DTB = MAT[J]$.

XIX. $OUT[I][J] = COV(DTA, DTB)$.

XX. RETURN OUT.

XXI. CALCULATE EIGEN SET:

XXII.

XXIII. INITIALIZE $COPY = INPUT$, $Q = [COPY.LENGTH][COPY.LENGTH]$.

XXIV. LOOP THROUGH Q.LENGTH TO SET $Q[I][J] = 1$.

XXV. INITIALIZE $DONE = FALSE$.

XXVI. CREATE $nMAT = MULTIPLY(Q[I], Q[J])$.

XXVII. IF $nMAT - COPY > 0.0000001$, UPDATE $COPY = nMAT$.

XXVIII. CALCULATE EIGEN VALUES AND VECTORS:

XXIX.

XXX. OBTAIN COVARIANCE MATRIX: $DATA = COVMAT()$.

XXXI. RETURN EIGEN.VALUES.

XXXII. LOOP THROUGH VALS[0].LENGTH TO CHECK ISNAN AND SET VALS[J][I] = 1.0.

XXXIII. CALCULATE FEATURE VALUES: FV = PC(10, EIGEN).

XXXIV. CALCULATE TRANSPOSE MATRIX: K = TRANSPOSE(MUL(PC, INTRANSPOSE)).

XXXV. PRINT MATRIX:

XXXVI.

XXXVII. CONVERT K TO STRING: DT = K.TOSTRING().

XXXVIII. LOOP THROUGH LENGTH(K) TO PRINT TT.

OUTPUT: RELEVANT FEATURE SET ACHIEVED

C. FEATURE EXTRACTION USING MULTI ALGORITHM FUSION APPROACH

The steps involved in the implementation of MAF-PSVM are announced in detail. The output of the PCA based feature extraction is given as an input for feature extraction process. The feature extraction algorithm such as Fisher Ratio, AW-SVM and PSVM are fused together, however selection of the kernel function is the challenging task for this fused feature extraction process. Different studies have exposed that the linear kernel function is an appropriate for the linearly separate data, therefore polynomial kernel was designated. In future, other kernel function will be studied. The polynomial kernel functions have different generalization ability. The optimum classifier performance is reached in many cases when parameter d has the range of 1. Thus, polynomial kernel order value d is engaged as 1 and feature extraction has been conducted for clustered samples. The fusion method content utilized in this approach is introduced below,

The score related multi algorithm is utilized in this study. Initially, all the feature scores were produced from the score vector with the each basis criteria. Next, the multiple score vectors were aggregated as single consent score vector using the score combination algorithm. Lastly, ranking is done based on consensus scores through employing feature ranking procedure. The score-based multi-algorithm fusion procedure is exemplified in Fig 2.

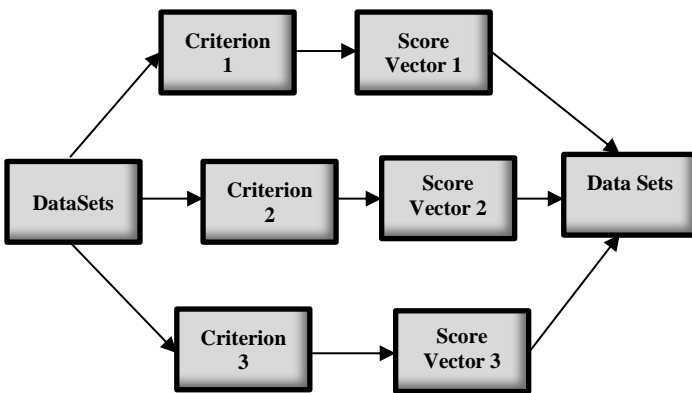


Figure 2: Score based multi algorithm

At distinct criteria, scores are produced under score aggregating will be analogous. The score normalization should be undertaken firstly then the score combination is performed, the normalized range is [0, 1] in this approach. Let's take u_i as the score vector formed by basis criterion i, the score normalization is accomplished as follows:

$$u'_i = \frac{u_i - u_{i \min}}{u_{i \max} - u_{i \min}} \quad (1)$$

Where $u_{i \min}$ and $u_{i \max}$ are min and max values of vector u .

The larger score in all the basis criteria achieves better features using the formula which is given below

$$u = \frac{1}{m} \sum_{i=1}^m u'_i \quad (2)$$

Where m is the basis criterion in fusion.

D. EFFECTIVE RANDOM FOREST BOOTSTRAP ALGORITHM:

In random forest algorithm, a group of randomized trees are constructed for classification. The optimal variable selection of nodes and the splitting of that nodes for generating the pure child node has to be concerned while constructing CART. A node impurity is measured utilizing Gini index of classification algorithm. Assume node s and the assessed class probabilities, $q(c/s)$ ($C=1\dots C$). The Gini index of node s is demarcated as

$$N(s) = \sum_{c_1 \neq c_2} q\left(\frac{c_1}{s}\right) q\left(\frac{c_2}{s}\right) = 1 - \sum_{c=1}^C q^2(c/s)$$

Allow o be the splitting point of node s , where the node gets separated into two portions. In that the proportion q_R , t samples are allocated to s_R and the s_l acquired from the q_l proportion. i.e. $s_R + s_l = 1$. Therefore the decrease in Gini index is shown below.

$$\Delta N(o, t) = N(s) - q_R N_{s_R} - q_l N_{s_l}$$

The optimal splitting point s^* and optimal variable j^* yields highest decrease in Gini impurity gotten and it is shown below
 $o^*, j^* = \text{argmax}_{o, j} \Delta N(o, s)$

The RF algorithm is split into three functions named as bootstrap sample, Tree Generation and Randomized Tree Learn. An efficient random forest bootstrap algorithm is designated in order to achieve effective gene classification with great classification accuracy and less computation time.

E. CLASSIFICATION:

After extracting the features, the classification technique is employed to classify the breast cancer diseases. In this work, for classification, Random forest classification is used. The novel in this algorithm is mentioned below,

Algorithm: Effective Random Forest Bootstrap Algorithm

Input: Prioritized cluster.

Output: Classified output.

Procedure:

1. Initialize Random Forest Parameters:

- inter: Number of seeds.
- numTree: Number of trees.
- data: Seed content.
- tdata: Test data.
- tree[]: Array of numTree.
- update: 100 / numTree.
- start: StartTimer().
- k: $\log(\text{data length} - 1) / \log(2) + 1$.
- done: estimateOOB().

2. Create StartTimer:

- treePool: Executor with NumThreads.
- Loop from 0 to numTrees:
 - treePool.create Tree(data, tdata).
- End loop.
- Call TestForest().

3. Create TestForest:

- R: Random values for trees.
- correctness: 0.
- Loop through tdata.size:
 - Loop through tree.size:
 - Add to R.
- End loops.
- pred: ModeOf(R).
- If pred == R, increment correctness.

4. Create ModeOf:

- max, maxclass: 1.
- Loop through tree.predict.size:
 - count: 0.
 - Loop again through tree.predict.size:
 - If $i == j$, increment count.
 - End inner loop.
 - If count > max, update maxclass and max.
- End outer loop.

5. **Create estimateOOB:**
 - If estimateOOB(R) == null:
 - map = C.
 - Call estimateOOB(R, map).
6. **Novelty:**
 - str: "File name".
 - Loop through str.size:
 - If str.get(IN)[str.get(IN).length] < categ, continue.
 - Else, Categ = str.get(IN)[str.get(IN).length].
 - End loop.
7. **Random Forest Function (RaF):**
 - Function RaF:
 - RaF.C = categ.
 - RaF.M = Input.get(0).length.
 - RaF.Ms = Math.log(RaF.M) / Math.log(2) + 1.
 - Call RaF.Start().
 - End Function.

Thus the Classification accuracy is achieved using Novel Random Forest Algorithm. Finally, the performance evaluation was made and the performance metrics were estimated to prove the effectiveness of this proposed system.

IV.RESULTS AND DISCUSSION

Colon cancer dataset is utilized in this approach to calculating the efficiency of the proposed method. The dataset <http://biogps.org/dataset/tag/colon%20cancer/> link is given here for your reference. The performance analysis of proposed work is described in this section. The performance measures such as error rate, accuracy, specificity, precision, recall, F1-score, False positive rate etc. is calculated for our proposed PCA based random forest classifier.

TABLE: 1 Performance Measures

PERFORMANCE MEASURES	VALUES
Error Rate	0.011461318
Accuracy	0.938538682
Specificity	0.971223022
Precision	0.981308411
Recall	1
F1-Score	0.990566038
False Positive Rate	0.028776978
False Negative Rate	0
Negative Predict Value	1
Mathews Correlation Coefficient	0.976252693

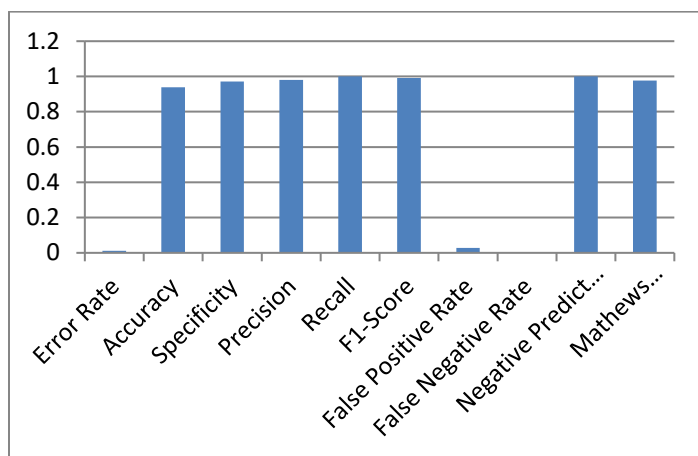


Figure 3: Performance measures

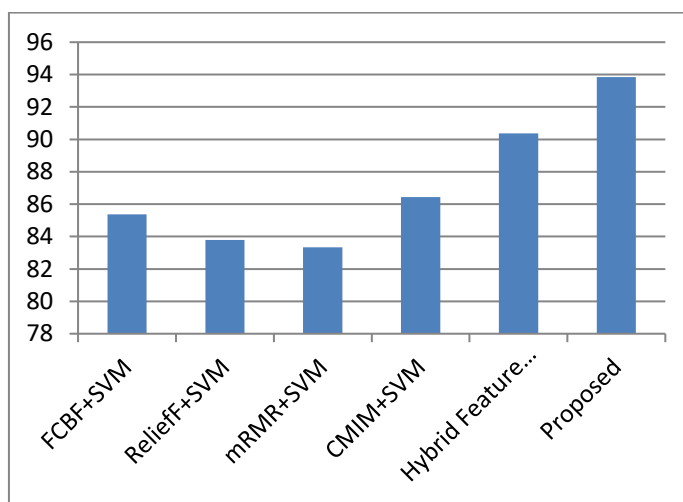


Figure 4: Comparison of gene classification accuracy

In figure 5, classification error for feature extraction of colon dataset for proposed MAF-PSVM, AW-SVM, PSVM are represented graphically, where x-axis represents number of features and y-axis denotes classification error rate. For instance, when considering 50 features, the classification error rate is 0.161 for proposed, 0.165 for AW-SVM, 0.172 for PSVM. From that graph, it is noted that, the proposed work has been recorded with low value, when compared to other existing methods.

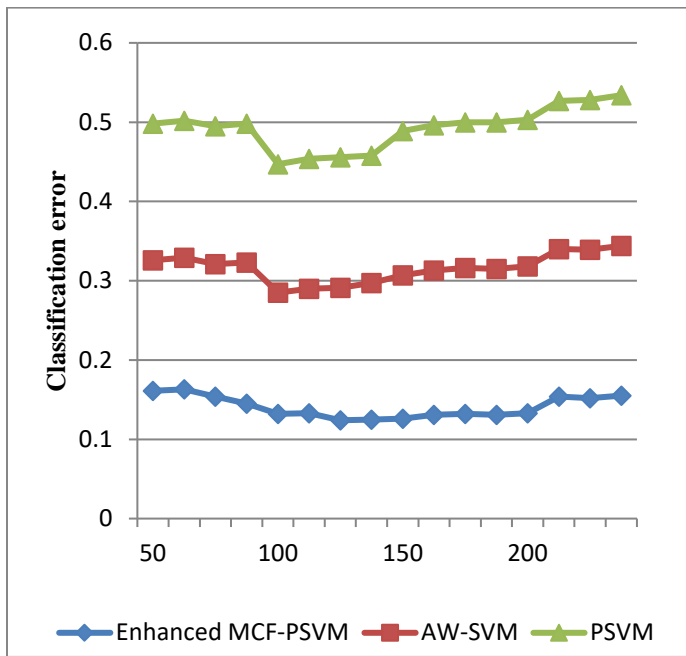


Figure 5: Classification error for colon dataset with enhanced MCF-PSVM

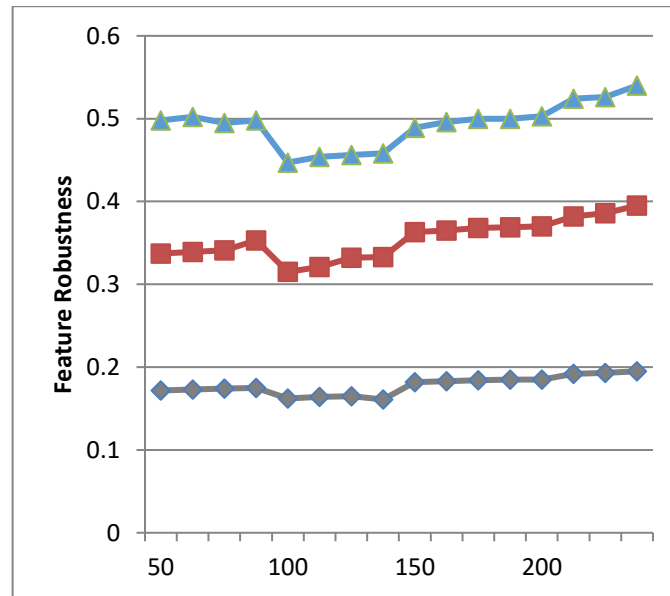


Figure 6: Feature robustness for colon dataset with enhanced MCF-PSVM

In figure 6, feature robustness of colon dataset for proposed MAF-PSVM, AW-SVM, PSVM are represented graphically, where x-axis represents number of features and y-axis denotes feature robustness. For instance, when considering 50 features, the feature robustness is 0.5 for proposed, 0.32 for PSVM, 0.172 for AW-SVM. From that graph, it is noted that, the proposed work has been recorded with high value, when compared to other existing methods

V. CONCLUSION

This method has been proposed to acquire an effective gene classification approach using effective PCA based random forest classifier algorithm. The feature extraction with multiple fusion algorithm has been accomplished using enhanced MAF-PSVM. Finally, random forest classifier was utilized for the effective classification approach and it was proved with efficient outcome of several performance measures.

REFERENCES

- [1] J. Han, J. Pei, and M. Kamber, *Data mining: concepts and techniques*: Elsevier, 2011.
- [2] J. C. Ang, H. Haron, and H. N. A. Hamed, "Semi-supervised SVM-based feature selection for cancer classification using microarray gene expression data," in *International Conference on Industrial, Engineering and Other Applications of Applied Intelligent Systems*, 2015, pp. 468-477.
- [3] J.-X. Liu, Y. Xu, C.-H. Zheng, H. Kong, and Z.-H. Lai, "RPCA-based tumor classification using gene expression data," *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, vol. 12, pp. 964-970, 2015.
- [4] J. Bennet, C. Ganaprakasam, and N. Kumar, "A Hybrid Approach for Gene Selection and Classification using Support Vector Machine," *International Arab Journal of Information Technology (IAJIT)*, vol. 12, 2015.
- [5] T. Latkowski and S. Osowski, "Data mining for feature selection in gene expression autism data," *Expert Systems with Applications*, vol. 42, pp. 864-872, 2015.
- [6] Z. Wang, G. Li, R. W. Robinson, and X. Huang, "UniBic: Sequential row-based biclustering algorithm for analysis of gene expression data," *Scientific reports*, vol. 6, p. 23466, 2016.
- [7] H. Lu, L. Yang, K. Yan, Y. Xue, and Z. Gao, "A cost-sensitive rotation forest algorithm for gene expression data classification," *Neurocomputing*, vol. 228, pp. 270-276, 2017.
- [8] S. Cogill and L. Wang, "Support vector machine model of developmental brain gene expression data for prioritization of Autism risk gene candidates," *Bioinformatics*, vol. 32, pp. 3611-3618, 2016.
- [9] H. Lu, J. Chen, K. Yan, Q. Jin, Y. Xue, and Z. Gao, "A hybrid feature selection algorithm for gene expression data classification," *Neurocomputing*, vol. 256, pp. 56-62, 2017.
- [10] X. Chen, J. Z. Huang, Q. Wu, and M. Yang, "Subspace weighting co-clustering of gene expression data," *IEEE/ACM transactions on computational biology and bioinformatics*, vol. 16, pp. 352-364, 2017.
- [11] H. Wang, X. Jing, and B. Niu, "A discrete bacterial algorithm for feature selection in classification of microarray gene expression cancer data," *Knowledge-Based Systems*, vol. 126, pp. 8-19, 2017.