

An Ensemble Machine Learning Model based on Multiple Filtering and Supervised Attribute Clustering Algorithm for Classifying Cancer Samples

Shilpi Bose^{Corresp., 1}, Chandra Das¹, Abhik Banerjee¹, Kuntal Ghosh², Matangini Chattopadhyay³, Samiran Chattopadhyay⁴, Aishwarya Barik¹

¹ Department of Computer Science and Engineering, Netaji Subhash Engineering College, Kolkata, West Bengal, India

² Machine Intelligence Unit & Center for Soft Computing Research, Indian Statistical Institute, Kolkata, West Bengal, India

³ School of Education Technology, Jadavpur University, Kolkata, West Bengal, India

⁴ Department of Information Technology, Jadavpur University, Kolkata, West Bengal, India

Corresponding Author: Shilpi Bose

Email address: shilpi.bose@nsec.ac.in

Background: Machine learning is one kind of machine intelligence technique that learns from data and detects inherent patterns from large, complex datasets. Due to this capability, machine learning techniques are widely used in medical applications, especially where large-scale genomic and proteomic data are used. Cancer classification based on bio-molecular profiling data is a very important topic for medical applications since it improves the diagnostic accuracy of cancer and enables a successful culmination of cancer treatments. Hence, machine learning techniques are widely used in cancer detection and prognosis.

Methods: In this article, a new ensemble machine learning classification model named Multiple Filtering and Supervised Attribute Clustering algorithm based Ensemble Classification model (MFSAC-EC) is proposed which can handle class imbalance problem and high dimensionality of microarray datasets. This model first generates a number of bootstrapped datasets from the original training data where the oversampling procedure is applied to handle the class imbalance problem. The proposed MFSAC method is then applied to each of these bootstrapped datasets to generate sub-datasets, each of which contains a subset of the most relevant/informative attributes of the original dataset. The MFSAC method is a feature selection technique combining multiple filters with a new supervised attribute clustering algorithm. Then for every sub dataset, a base classifier is constructed separately, and finally, the predictive accuracy of these base classifiers is combined using the majority voting technique forming the MFSAC-based ensemble classifier. Also, a number of most informative attributes are selected as important features based on their frequency of occurrence in these sub-datasets.

Results: To assess the performance of the proposed MFSAC-EC model, it is applied on different high-dimensional microarray gene expression datasets for cancer sample classification. The proposed model is compared with well-known existing models to establish its effectiveness with respect to other models. From the experimental results, it has been found that the generalization performance/testing accuracy of the proposed classifier is significantly better compared to other well-known existing models. Apart from that, it has been also found that the proposed model can identify many important attributes/biomarker genes.

An Ensemble Machine Learning Model based on Multiple Filtering and Supervised Attribute Clustering Algorithm for Classifying Cancer Samples

Shilpi Bose¹, Chandra Das¹, Abhik Banerjee¹, Kuntal Ghosh², Matangini Chattopadhyay³, Samiran Chattopadhyay⁴, Aisarya Barik¹

¹Department of CSE, Netaji Subhash Engineering College, Kolkata, West Bengal, India

²Machine Intelligence Unit & Center for Soft Computing Research, Indian Statistical Institute, Kolkata, West Bengal, India

³School of Education Technology, Jadavpur University, Kolkata, West Bengal, India

⁴Department of Information Technology, Jadavpur University, Kolkata, West Bengal, India

Corresponding Author:

Shilpi Bose

Department of CSE, Netaji Subhash Engineering College, Panchpota, Garia, Kolkata, 700152, West Bengal, India

Email address: shilpi.bose@nsec.ac.in

Abstract

Background: Machine learning is one kind of machine intelligence technique that learns from data and detects inherent patterns from large, complex datasets. Due to this capability, machine learning techniques are widely used in medical applications, especially where large-scale genomic and proteomic data are used. Cancer classification based on bio-molecular profiling data is a very important topic for medical applications since it improves the diagnostic accuracy of cancer and enables a successful culmination of cancer treatments. Hence, machine learning techniques are widely used in cancer detection and prognosis.

Methods: In this article, a new ensemble machine learning classification model named Multiple Filtering and Supervised Attribute Clustering algorithm based Ensemble Classification model (MFSAC-EC) is proposed which can handle class imbalance problem and high dimensionality of microarray datasets. This model first generates a number of bootstrapped datasets from the original training data where the oversampling procedure is applied to handle the class imbalance problem. The proposed MFSAC method is then applied to each of these bootstrapped datasets to generate sub-datasets, each of which contains a subset of the most relevant/informative attributes of the original dataset. The MFSAC method is a feature selection technique combining multiple filters with a new supervised attribute clustering algorithm. Then for every sub dataset, a base classifier is constructed separately, and finally, the predictive accuracy of these base classifiers is combined using the majority voting technique forming the MFSAC-based ensemble classifier. Also, a number of most informative attributes are selected as important features based on their frequency of occurrence in these sub-datasets.

Results: To assess the performance of the proposed MFSAC-EC model, it is applied on different high-dimensional microarray gene expression datasets for cancer sample classification. The proposed model is compared with well-known existing models to establish its effectiveness with respect to other models. From the experimental results, it has been found that the generalization performance/testing accuracy of the proposed classifier is significantly better compared to other well-known existing models. Apart from that, it has been also found that the proposed model can identify many important attributes/biomarker genes.

Introduction

Cancer is one of the most fatal diseases around the globe (Tabares-Soto 2020, Hambali 2020). According to the World Health Organization report, Cancer is marked as the second most deadly disease and an estimated 9.7 million deaths around the world in 2018 have occurred due to this signature disease (Hambali 2020). Generally, one in every 6 deaths all over the world, occurs due to cancer. So, within 2030, the number of new cancer patients per year will increase approximately by 25 million (Hambali 2020, NIH 2019). Although several advanced techniques are already developed for the detection of cancer, the proper prognosis of cancer patients, till date, is very poor and the survival rate is also very low (Tabares-Soto 2020, Hambali 2020, Konstantina 2015). It has been already found that for very accurate cancer sample classification or prediction, adequate information is not available from the clinical, environmental, and behavioral characteristics of patients (Kourou 2015, Hambali 2020, Tabares-Soto Reinel 2020). Recently, due to different types of bio-molecular data analysis, several genetic disorders with different biological characteristics have been revealed which are very helpful for early identification and prognosis of cancer and also to discern the responses for different types of treatment (Colozza2005, Greller 1999, Li 2018, Liu2011, Pilling 2017, Su 2001, Swan 2013). With the rapid advancements in genomic, proteomic, and imaging high-throughput technologies (Colozza2005, Greller 1999, Li 2018, Liu2011, Pilling2017, Su 2001, Swan 2013), now it is possible to accumulate huge amount (in the order of thousands) of different bio-molecular information of patients. Using this huge amount of information, researchers have been trying to develop more advanced techniques for early detection and proper prognosis of cancer, and also to improve cancer therapy for improvement of patients' survival rate. To analyze this huge amount of information, lab-based approaches are not adequate as these methods are costly and time-consuming. So, computational or in-silico methods like statistical methods, machine learning, deep learning, etc. have been being used extensively in this field. It is well-known fact that in cancer-causing cells, gene expression is either overexpressed or under expressed (Tabares-Soto 2020). So, measurement of gene expression in cancer cells can give adequate information to improve cancer diagnostic procedures. Nowadays, different developing countries have been using this procedure for cancer sample detection. It is already known that using DNA microarray technology it is possible to measure the expression level of a numerous number of genes for a single experiment/sample simultaneously. The outcome of DNA microarray technology is a gene expression data matrix. This matrix carries information

about the expression level of a huge number of genes for a limited number of samples (such as diseased patient samples and normal samples). The presence of the limited number of samples in this data matrix is due to the lack of availability of samples. So, based on information of gene expression data matrix, cancer sample classification is one of the essential tasks in the field of cancer research (Chin 2016, Dashtban 2017, Ding 2005, Elyasigomari 2017, Furey2000, Golub 1999, Nada 2019, Tabares-Soto 2020).

Using computational or in-silico approaches, gene expression-based cancer sample classification task has been reviewed extensively in different papers (Chin 2016, Dashtban2017, Ding 2005, Elyasigomari 2017, Furey2000, Golub 1999, Nada 2019, Tabares-Soto 2020). However, the main difficulties in the sample classification task arise due to several factors. Firstly, in these data sets, a substantially small number of samples is available (generally in the order of hundreds) compared to the availability of a huge number of genes (generally in the order of thousands) (Chin 2016, Hambali 2020, Nada 2019). For sample classification, genes are treated as features/attributes. So, the high-dimensional gene space is an overhead for most classification algorithms. Secondly, only a very few genes are informative (differentially expressed) and the rest of the section is non-informative (noisy) (Chin 2016, Hambali 2020, Nada 2019) for sample classification and responsible for degrading the classifier's performance. Gene dimension reduction by identification of informative genes as biomarkers can improve the classification accuracy of classifiers. Apart from the improvement of classification accuracy, the identification of informative biomarkers (here, informative genes) has great prospects from a biomedical point of view. These are beneficial for finding the biological reason for a disorder, assessing disease risk, and developing therapeutic targets. The third problem arises due to the small sample size which creates an overfitting problem in classifier construction. Another problem that degrades classifier performance is the sample class imbalance problem. This problem occurs due to the presence of more instances/samples of one class (majority class) with respect to other class(es) (minority class) in a dataset.

A fairly large number of works have been already developed for sample classification. These works are divided into two categories. In the first category (Chin 2016, Hambali 2020, Nada 2019), the major emphasis is given to the selection of relevant genes for the reduction of feature space. Then based on this reduced feature space, predictive/classification accuracy of the samples is measured using different existing single classification models like naïve Bayes, support vector machine, relevance vector machine, K-nearest neighbor, decision tree, logistic regression, etc. As gene selection is a feature selection task, so based on feature selection techniques, these methods are divided into different categories. These are (1) filter methods (2) wrapper methods (3) embedded methods and (4) hybrid methods. Before we mention the second category of classification methods, let us first elaborate on the first category methods one by one. Filter methods (Chin 2016, Hambali 2020, Nada 2019) select a subset of features without taking any information from any classification model. These methods select features that are differentially expressed with respect to sample class labels. The filter methods rank individual features according to their class discrimination power based on some statistical score function

and then select a number of high-ranked features to form a reduced and relevant feature subset. The popular statistical score functions used in filter methods are Fisher's score, Signal to Noise ratio (SNR), correlation coefficient, mutual information, Relief (Das 2019), etc. Filter methods are computationally simple, fast, and unbiased in favor of any specific classifier as these methods do not consider any knowledge from any classifier at the feature selection phase. The drawback of filter methods is that the number of selected features is based solely on the trial-error method. Wrapper methods (Chin 2016, Hambali 2020, Nada 2019), on the other hand, judge discrimination capability of a feature subset using classification error rate or prediction accuracy of a classifier as the feature evaluation function. It selects the most discriminative feature subset via minimizing the classification error rate or maximizing the classification accuracy of a classifier. The wrapper methods generally achieve better classification accuracy than the filter methods because the selection of feature subset is classifier-dependent. One drawback of these methods is that these are biased to used classifiers and another drawback is that these are computationally more expensive than the filter methods as generation of the best feature subset for the high-dimensional dataset is an NP-complete problem. Due to these reasons, these methods are not applicable for high-dimensional datasets.

In Embedded methods (Chin 2016, Hambali 2020, Nada 2019), the optimal feature subset is selected through the unique learning procedure of a specific classifier at the time of classifier construction. Actually, in these methods, the optimal feature subset selection part is embedded as part of classifier construction. These methods are faster than wrapper methods but are biased to the specific classifier. In embedded approaches, the feature selection process is specific for a particular classifier and is not applicable to other classifiers. These are also computationally expensive. Due to these reasons for high-dimensional datasets, these methods are not applicable. On the other hand, recently hybrid feature selection methods (Chin 2016, Hambali 2020, Nada 2019) are also developed. In hybrid methods, different category-based methods are combined to take advantage of all of these methods for improving classification accuracy.

Apart from these methods, clustering techniques (Chin 2016, Hambali 2020) are also used for feature selection purposes. Clustering techniques divide the data space in such a manner that objects in the same cluster are similar while in different clusters they are dissimilar. For the feature selection task, clustering methods (famous as attribute clustering in feature selection domain) (Au 2005) divide the features into several distinct clusters and then reduce the feature dimension by selecting a small number of significant features from each cluster. A lot of unsupervised gene (attribute) clustering algorithms (Au2005, Chin 2016, Hambali 2020) are already developed for this task. However, these methods are unsuccessful to find informative functional groups of genes for sample classification as in clustering genes, no supervised information from sample classes is considered (Au2005, Chin 2016, Hambali 2020). So, scientists have developed a number of supervised gene (attribute) clustering algorithms (Dettling 2002, Hastie 2000, Hastie 2001, Maji 2012) in which genes are grouped using supervised information from sample classes and a reduced gene set is formed via selecting the most informative genes from each cluster.

All the above-mentioned variants deliver comparable feature selection and classification accuracy. Quite often this type of classification models with only a few genes and with a limited number of training samples can classify the majority of training samples correctly, but the generalization capability of such classification models cannot be guaranteed (Bolo2012, Ghorai 2011, Nagi 2013, Wang 2006, Wang 2012, Yang 2010). So, the most important task for a medical diagnosis system is to improve the classification accuracy of unknown samples (generalization performance) which cannot be solved by this type of classification model. Apart from this problem, the microarray data is related to several uncertainties due to fabrication, hybridization, and image processing procedure in microarray technology. These uncertainties introduce various types of noise in microarray data. Due to the presence of these uncertainties with a limited number of training samples, the conventional machine learning approaches face challenges to develop reliable classification models.

To overcome the above-mentioned problems, it is therefore essential to develop general approaches and robust methods. In this regard, researchers are motivated to develop the second category-based model. These are the different robust ensemble classification models (Bolo2012, Ghorai 2011, Nagi 2013, Osareh 2013, Wang 2006, Wang 2012, Yang 2010) which can overcome small sample size problems and are capable of removing uncertainties of gene expression data.

Ensemble methods (Dietterich 2000) are a class of machine learning technique which combines multiple base learning algorithms to produce one optimal predictive model. Ensemble classification model refers to a group of individual/base classifiers that are trained individually on the trained dataset in a supervised classification system and finally, an aggregation method is used to combine the decisions produced by the base classifiers. These ensemble classification models have the potential to alleviate the small sample size problem by applying multiple classification models on the same training data or on bootstrapped samples (sampling with replacement) of the training data to decrease the chance of overfitting in the training data. In this way, the training dataset is utilized more efficiently, and as a consequence, the generalization ability is improved.

Although different category-based ensemble classification models exist in the literature but these ensemble models are not capable of addressing all the above-mentioned problems (small sample size, high dimensional feature space, and sample class imbalance problem) related to microarray data.

In this regard, here a new Multiple Filtering and Supervised Attribute Clustering algorithm-based ensemble classification model named MFSAC-EC is proposed. In this model, first, a number of bootstrapped versions of the original training dataset are created. At the time of the creation of bootstrapped versions, an oversampling technique (Błaszczyński 2013) is adopted to solve the class imbalance problem. For every bootstrapped dataset a number of sub-datasets (each with a subset of genes) are generated using the proposed MFSAC method. The MFSAC is a hybrid method combining multiple filters with a new supervised attribute clustering method. Then for every sub dataset, a base classifier is constructed. Finally, based on the prediction accuracy of all

these base classifiers of all sub-datasets for all bootstrapped datasets an ensemble classifier (EC) is formed using the majority voting technique. The novelty of the proposed MFSAC-EC model is that here the emphasis is given simultaneously on the high dimensionality problem of gene expression data, small sample size problem as well as the class imbalance problem. All of these problems at the same time are not considered in any existing ensemble classification model. First of all, due to the use of bootstrapping method with a class balancing strategy, the proposed model can handle a small sample size and overfitting problem. Secondly, in MFSAC, different filter methods are used with their unique characteristics. So, different characteristics-based relevant gene subsets are selected via different filters to form different sub-datasets from every bootstrapped dataset. Finally, every gene subset is modified using a supervised attribute clustering algorithm. In this way, the high-dimensionality problem of gene expression data is handled here. Apart from this, from the MFSAC generated sub-datasets, the frequency of occurrence is counted for every gene and informative genes are ranked accordingly. The prediction capability of the proposed model is experimented with over different microarray datasets and compared with the existing well-known models. Experimental outputs demonstrate the superiority of the proposed model over existing models.

Materials & Methods

The proposed MFSAC-EC model is composed of different filter score functions, a new supervised attribute clustering method, and an ensemble classification method. In the following subsections, first, a brief overview is given on different filter score functions and then the proposed MFSAC-EC model is described.

Preliminaries

In this paper, a data set (here, a microarray gene expression data set) is represented by a data matrix, $K_{U \times V}$, with U data objects (samples) and V features (genes). The set of objects or samples is represented as $E = \{E_1, E_2, \dots, E_s, \dots, E_U\}$ while the set of genes is represented as $G = \{G_1, G_2, \dots, G_t, \dots, G_V\}$. Here, each sample is a V -dimensional feature vector containing V number of gene expression values. Similar way, every gene is a U -dimensional vector containing U number of sample values. Here, $C_{U \times 1}$ is a class vector representing the associated class label for every sample. The class label is taken from a set $DC = \{d_1, d_2, \dots, d_j, \dots, d_N\}$ with N distinct class labels.

Brief overview of Filter score functions used in MFSAC

The filter score functions used in the proposed MFSAC-EC model are modified Fisher score (Gu 2011), modified T-test (Zhou 2007), Chi-square (Das 2019), Mutual information (Das 2019), Pearson correlation coefficient (Leung 2010), SNR (Leung 2010) and Relief-F (Das 2019). A summary of these 7 filters used in the MFSAC-EC model is given in the supplemental Table S1.

Proposed MFSAC-EC Model

In the proposed MFSAC-EC model, initially, bootstrapping (sampling with replacements) with a class balancing procedure of samples is applied on training dataset K to create D number of different bootstrapped versions from the training dataset. Here, every bootstrapped dataset with

239 U samples is formed by random sampling with replacements U times from the original dataset K .
 240 After that oversampling procedure is applied to each minority class to achieve data balance.
 241 Oversampling consists of increasing the minority class instances by their random replication to
 242 exactly balance the cardinality of the minority and majority classes in each bootstrapped dataset.
 243 Due to oversampling each bootstrapped dataset will contain more instances than the original
 244 dataset.
 245 The MFSAC method of the MFSAC-EC model, which is an integration of multiple filters and a
 246 new supervised attribute (gene) clustering method, is applied on every newly created
 247 bootstrapped (BK_l) training dataset. The proposed MFSAC method first calculates the class
 248 relevance score of every gene present in the bootstrapped training dataset using each filter score
 249 function (FT_x), $x = 1$ to 7 mentioned above. Then for each filter score function, a sub dataset S
 250 D_{lx} with a gene subset (GS_{lx}) is created by selecting a predefined number (let P) of the most
 251 relevant genes from the full gene set G . So, $|GS_{lx}| = P$. After that on every gene subset (GS_{lx}) of
 252 every sub dataset SD_{lx} , the SAC (Supervised Attribute Clustering) method is applied and a set of
 253 clusters CGS_{lx} and corresponding cluster representatives (considered as modified features) are
 254 formed. Finally, Q numbers of most relevant cluster representatives are selected as modified
 255 features and a reduced sub dataset RSD_{lx} of the sub dataset SD_{lx} is formed. How the SAC
 256 method works on GS_{lx} of every sub dataset SD_{lx} is discussed below.
 257 For any sub dataset SD_{lx} , the SAC method starts by selecting the gene from the subset (GS_{lx})
 258 with the highest FT_x value. Let gene $G_{li} \in GS_{lx}$ with the highest FT_x value be selected as the first
 259 member (let $FT_x(G_{li}, C) = A$) and it also becomes the initial cluster representative $R(R = G_{li})$
 260 of the first cluster C_1GS_{lx} and G_{li} is deleted from GS_{lx} . In effect, $G_{li} \in C_1GS_{lx}$, and $GS_{lx} = GS_{lx}$
 261 $- \{G_{li}\}$ and so $FT_x(R, C) = A$. This cluster is then grown up in parallel with the cluster
 262 representative refinement process which is described next. In this process, the gene (let G_{lm}) with
 263 next highest FT_x value is taken from GS_{lx} subset and is merged with the current cluster
 264 representative R . The merging is done in two ways. Firstly the expression profile of G_{lm} is
 265 directly added with R and a temporary augmented representative TR^+ is formed and its FT_x
 266 value (let B_1) is calculated. The second one is that the sign-flipped value of the expression
 267 profile of G_{lm} is added with R and another temporary augmented representative TR^- is formed
 268 and its FT_x value (let B_2) is calculated. If $FT_x(TR^+, C) \geq FT_x(TR^-, C)$ that is $B_1 \geq B_2$ then
 269 TR^+ is chosen else TR^- is chosen. Let TR^+ is chosen. Now if $FT_x(TR^+, C) > FT_x(R, C)$
 270 then $R = TR^+$ otherwise, R is unaltered. Similar way if TR^- is chosen and if $FT_x(TR^-, C) >$
 271 $FT_x(R, C)$ then $R = TR^-$ otherwise, R remains unchanged. If R is modified then the gene G_{lm} is
 272 included in the cluster and G_{lm} is deleted from GS_{lx} . In effect, $G_{lm} \in C_1GS_{lx}$, and $GS_{lx} = GS_{lx} -$
 273 $\{G_{lm}\}$. So, the next chosen gene is included in the current cluster if it improves the class

relevance value of the current cluster representative. The merging process is described in Figure 1.

Here g_0 represents the current cluster representative (R) and its class relevance score ((FT_x, R) , here Pearson score), is shown. Now among all the genes g_1, g_2, g_3, g_4 , and g_5 , the Pearson score of g_1 is the highest. So, g_1 is chosen for the merging process. Then g_1 is added with R to create the temporary augmented representative ($TR^+ = R + g_1$) and also its sign-flipped value is added with the R to form the temporary augmented representative ($TR^- = R - g_1$). The Pearson score of TR^+ is greater than the Pearson score of TR^- , so TR^+ is chosen. Now the Pearson score of TR^+ is greater than the Pearson score of R , so TR^+ is considered as the current cluster representative and $R = TR^+$. This process is continued for all other genes. Now, g_3 is chosen as it is the gene with the next highest Pearson value. g_3 and its sign-flipped value are added individually with current cluster representative R to form $TR^+ = R + g_3$ and $TR^- = R - g_3$ respectively. In this case, Pearson score of TR^- is greater than the Pearson score of TR^+ . So, TR^- is chosen. Then Pearson score of TR^- is Checked with the Pearson score of R and here Pearson score of TR^- is greater than the Pearson score of R . So, TR^- is considered as current cluster representative and $R = TR^-$. In this way, cluster representative is refined. This process is repeated for every member of GS_{lx} subset.

After the formation of the first cluster and its corresponding augmented representative, R is assigned to AR_{lx1} that means $AR_{lx1} = R$, and the supervised clustering process is repeated to form the second cluster with the gene (let G_{lx}) with next highest FT_x value from GS_{lx} subset. In this way a set of clusters $CGS_{lx} = \{C_1GS_{lx}, C_2GS_{lx}, \dots, C_kGS_{lx}, \dots\}$ and their corresponding augmented cluster representatives $AR_{lx} = \{AR_{lx1}, \dots, AR_{lxk}, \dots\}$ are formed. After that Q number of most powerful augmented cluster representatives are chosen (as modified features) according to their FT_x value from the generated clusters and with these Q number of modified features, a reduced sub dataset RSD_{lx} of sub dataset SD_{lx} is formed.

In this way, for every bootstrapped version (BK_l) of the training dataset, 7 number of RSD_{lx} sub-datasets are created and for every RSD_{lx} an individual classifier is constructed using any existing classifier and finally, an ensemble classifier (EC) is formed by combining all these classifiers of all bootstrapped versions using the majority voting technique. To classify every sample using this ensemble classifier, each classifier votes or classifies the sample for a particular class, and the class for which the highest number of votes is obtained is considered as the output class.

MFSAC method based Informative Attribute Ranking

For every gene (feature/attribute), the frequency of occurrence that means the total number of times it appears in all sub-datasets generated by the MFSAC method for all bootstrapped versions is calculated. Then according to their frequency of occurrence, those genes are ordered or ranked. The top-ranked genes with the highest occurrence frequency are considered the most informative cancer-related genes.

The block diagram of the proposed MFSAC-EC model is shown in Figure 2, while the block diagram of the MFSAC method is shown in Figure 3. The algorithm of the proposed model is described below.

Algorithm: MFSAC-EC

Input: A $K_{U \times V}$ data matrix (here, gene expression data matrix) containing U number of data objects (here, cancer samples) and V number of attributes (here, genes).

Output: An ensemble classifier MFSAC-EC is formed to classify test samples. From MFSAC generated sub-datasets, informative genes are selected according to their rank. Every gene is ranked according to its frequency of occurrence.

Definitions:

$E = \{E_1, E_2, \dots, E_s, \dots, E_U\}$ is the set of objects or samples of $K_{U \times V}$ data matrix. Every sample E_s is a V dimensional vector.

$G = \{G_1, G_2, \dots, G_t, \dots, G_V\}$ is the set of features or genes of $K_{U \times V}$ data matrix. Every gene G_t is a U dimensional vector.

$BK = \{BK_1, BK_2, \dots, BK_l, \dots, BK_D\}$ is a set of the bootstrapped version of the original training dataset. In every bootstrapped dataset the number of samples varies from the original dataset but the number of features is the same as the original dataset.

$C_{U \times 1}$ is a class vector representing the associated class label for every sample. For a data matrix N distinct class labels exist and class labels are taken from a set $DC = \{d_1, d_2, \dots, d_k, \dots, d_N\}$.

$FT_x(G_t, C)$ is x^{th} filter score function which returns the class relevance value of G_t gene with respect to class vector C using FT_x score function, for $x = 1$ to 7 as 7 represents the total number of filtering score functions used here.

$GS_{lx} (GS_{lx} = P)$ is a set of top-ranked genes of G selected using FT_x score function and SD_{lx} is corresponding sub dataset of BK_l . Here SD_{lx} is a data matrix containing P number of genes.

$CGS_{lx} = \{C_1GS_{lx}, C_2GS_{lx}, \dots, C_kGS_{lx}, \dots\}$ and $AR_{lx} = \{AR_{lx1}, \dots, AR_{lxk}, \dots\}$ are the set of clusters and corresponding cluster representatives respectively generated from the corresponding subset GS_{lx} of SD_{lx} . Here every AR_{lxk} is a vector.

TR^+ , TR^- , R are vectors similar to a gene vector.

$RSD_l = \{RSD_{l1}, RSD_{l2}, \dots, RSD_{lx}, \dots, RSD_{l7}\}$ is a set of sub-datasets each containing Q number of most relevant cluster representatives formed for every bootstrapped dataset BK_l .

$CF_l = \{IC_{l1}, IC_{l2}, \dots, IC_{lx}, \dots, IC_{l7}\}$ is a set of classifiers formed for every bootstrapped dataset.

1. Create D number bootstrapped version of training dataset K .

2. For Every bootstrapped dataset BK_l repeat step 3

3. Repeat for $x = 1$ to 7

- A. Repeat for $t = 1$ to V

- a) Calculate class relevance score $FT_x(G_t, C)$ of G_t gene, where $G_t \in G$, with respect to class vector C
- B. Select P number of top-ranked genes from G based on FT_x score function and form GS_{lx} gene subset with corresponding SD_{lx} sub dataset
- C. Set $k = 0$
- D. Repeat until $GS_{lx} = \emptyset$
 - a) Set $k = k + 1$
 - b) Set $AR_{lxk} = 0$, $R = 0$, and $i = 0$
 - c) Select the gene (let G_{li}) whose FT_x score value is maximum among all genes of GS_{lx} and set $R = G_{li}$
 - d) Add G_{li} to $C_k GS_{lx}$, and delete G_{li} from GS_{lx}
 - e) Set count = 1
 - f) Repeat for $j = 1$ to $|GS_{lx}|$
 - I. Compute first augmented representatives TR^+ by adding $G_{lj} \in GS_{lx}$ with R that means $TR^+ = R + G_{lj}$
 - II. Compute second augmented representatives TR^- by adding sign-flipped version of $G_{lj} \in GS_{lx}$ with R that means $TR^- = R - G_{lj}$
 - III. Compute class relevance value $FT_x(TR^+, C)$ and $FT_x(TR^-, C)$ using FT_x score function
 - IV. If $FT_x(TR^+, C) \geq FT_x(TR^-, C)$ then
 - If $FT_x(TR^+, C) > FT_x(R, C)$ then
 - Set $R = R + G_{lj}$ and add G_{lj} to $C_k GS_{lx}$ and delete G_{lj} from GS_{lx}
 - count = count + 1
 - V. If $FT_x(TR^-, C) > FT_x(TR^+, C)$ then
 - If $FT_x(TR^-, C) > FT_x(R, C)$ then
 - Set $R = R - G_{lj}$ and add G_{lj} to $C_k GS_{lx}$ and delete G_{lj} from GS_{lx}
 - count = count + 1
 - g) Set $R = R/\text{count}$
 - h) Set $AR_{lxk} = R$
- E. Select Q number of most relevant cluster representatives according to FT_x score from AR_{lx} set and form RSD_{lx} sub data set.
- F. Construct a classifier C_{lx} for RSD_{lx} sub data set

4. Apply a test sample over all the classifiers of all bootstrapped dataset and calculate the prediction accuracy of each classifier
5. Apply simple voting over all predictions to form an ensemble classifier *EC* and get final prediction.
6. Calculate number of occurrences for every gene for all RSD_{lx} sub datasets across all bootstrapped versions and rank them according to their count.
7. Select a number of top-ranked genes as informative genes.
8. End

Description and Preprocessing of the Datasets

The experimentation has been carried out over ten publicly available different gene expression binary class and multi-class datasets. Among these datasets, eight datasets are cancer datasets and two arthritis datasets. The eight cancer datasets are Leukemia (Golub1999), Colon (Alon 1999), Prostate (Singh 2002), Lung (Gordon 2002), RBreast (Veer 2002), Breast (West 2001), MLL (Armstrong 2001), and SRBCT (Khan 2001). To show the accuracy of the proposed model with respect to other than cancer datasets here two arthritis datasets RAHC (Pouw Kraan 2003) and RAOA (Pouw Kraan 2007) are also considered. The summary of the datasets is represented in Table 1.

In the Leukemia dataset (Golub1999), the gene expression data matrix is prepared using Affymetrix oligonucleotide arrays. The original dataset consists of two datasets: the training dataset and the testing dataset. The training dataset consists of 38 samples (27 Acute Lymphoblastic Leukemia (ALL) and 11 Acute Myeloid Leukemia (AML)) while the test dataset consists of 34 samples (20 Acute Lymphoblastic Leukemia (ALL) and 14 Acute Myeloid Leukemia (AML)), each with 7129 probes from 6817 genes. For the Leukemia dataset, training and test datasets are merged here and genes with missing values are removed and finally, the dataset with 7070 genes and 72 samples is prepared.

In the Colon cancer dataset (Alon 1999), gene expression of 6500 genes for 62 samples is measured using Affymetrix oligonucleotide arrays. Among these 62 samples, 40 are Colon cancer samples and 22 are normal samples. Among these 6500 genes, 2000 genes are selected based on the confidence of measured expression levels.

Prostate cancer dataset (Singh 2002) also consists of training and testing datasets. In the training dataset, among 102 samples, 50 are normal samples and 52 are prostate cancer samples. In the test dataset among 34 samples, 25 are prostate cancer samples and 9 are normal prostate samples. Gene expression of every sample is measured with respect to 12600 genes using Affymetrix chips. Here, training and test datasets are merged, and a dataset with 12600 genes and 136 samples is formed.

The Lung cancer dataset (Gordon 2002) consists of 181 samples. Among these samples, 31 are malignant pleural mesothelioma and rest 150 adenocarcinoma of lung cancer. Each sample is

represented by 12533 genes and the gene expression of every sample is measured using Affymetrix human U95A oligonucleotide probe arrays.

In Rbreast data set (Veer 2002), the patients, who are considered as breast cancer patients after 5 years intervals of initial diagnosis, fall under the category of relapse and rest as no relapse of metastases. 97 samples have been provided in which 46 patients developed distance metastases within 5 years and they are considered as relapse while the remaining remained healthy and are labeled as non-relapse. This dataset comprises 24481 genes and among them, 293 are removed.

In the Breast cancer dataset (West 2001), the gene expression of 49 samples is measured using HuGeneFL Affymetrix microarray arrays. Breast tumors are positive or negative in the presence or absence of estrogen receptors (ER). In this dataset, 25 samples are ER+ tumors and 24 samples are ER- tumors.

MLL (Armstrong 2001) is a type of dataset which comprises of training data set of 57 leukemia samples including 20 ALL, 17 MLL, and 20 AML and the test dataset including 4 ALL, 3 MLL, and 8 AML samples. For MLL cancer dataset training and test, datasets are merged here and finally, the dataset with 12582 genes and 72 samples are prepared.

SRBCT dataset (Khan 2001) is introduced as a dataset comprising of gene-expression for identifying small round blue-cell tumors of childhood SRBCT and samples of this dataset are further divided into four class which are neuroblastoma, rhabdomyosarcoma, non-Hodgkin lymphoma, and Ewing family of tumors and they are obtained from cDNA microarrays. A training set consisting of 63 SRBCT tissues, a test set consisting of 20 SRBCT and 5 non-SRBCT samples are available. Here we have considered only the training dataset. Each tissue sample is already standardized to zero mean value and has a unit variance across the genes.

RAHC commonly known as Rheumatoid Arthritis versus Healthy Controls is a data set (Pouw Kraan 2003) which comprises of gene expression characterizing as peripheral blood cells of 32 patients with RA, 3 patients with probable RA, and 15 age with sex-matched healthy controls performed under microarrays with a complexity of 26000 unique genes of 46000 elements.

RAOA commonly known as Rheumatoid Arthritis versus Osteoarthritis is a dataset (Pouw Kraan 2007) that includes the gene expression of thirty patients in which 21 of them are with RA and the remaining 9 of them are with OA. The Cy5 labeled experimental cDNA and Cy3 labeled common reference sample were pooled and hybridized to the lymphochips (consisting of 18000 cDNA spots which symbolize immunology in the genes of relevance).

Results

To assess the performance of the proposed MFSAC-EC model, four well-known existing classifiers named K-Nearest Neighbor (Duda 1999), Naive Bayes (Duda 1999), Support vector machine (Vapnik 1995), and Decision tree(c4.5) (Duda 1999) are applied independently in this model and four different ensemble classification models are formed. To prove the superiority of the proposed model, it is compared with existing well-known filter methods (used here) and existing recognized gene selection methods (Ding 2005, Au2005, Maji2005) and also with different existing ensemble classifiers (Bolo2012, Nagi 2013, Osareh 2013, Wang 2006, Wang 2012). To analyze the performance, the methods are applied to different publicly available

cancer and other disease-related gene expression datasets. The major metrics used here for evaluations of the performance of the proposed classifier are the Cross-validation method (LOOCV, 5-fold, and 10-fold), ROC Curve, and Heat map.

Tools Used

The algorithms are implemented using Python programming language and Scikit-learn libraries (Pedregosa 2011) which are explained in (Komer 2014) for ML algorithms. The programs are executed on an online Colab platform with 12 GB RAM and Intel(R) Xeon(R) processor available in the "CPU" Runtime Type at the time of writing. Figures and tables are generated in the Matplotlib library (Hunter 2007) and also in Microsoft Excel. The python codes used here are available at https://github.com/NSECRsearchCD-SLB/PEERJ_MFSAC_EC.

In the following subsections, first, the different types of metrics used here are discussed, and then the performance of the proposed MFSAC-EC model is verified with respect to these metrics.

This is followed by comparing the classification performance of the proposed model with different existing methods in terms of 10-fold cross-validation. The proposed model does not only perform the task of classification but also ranks every attribute or gene in descending order based on its information present in the dataset. To show the effectiveness of this ranking procedure topmost eight genes from Colon cancer and Leukemia cancer datasets are represented with their corresponding names, symbols, and references in significant cancer-related journals to demonstrate their significant roles in these cancers.

Evaluation Metrics

The performance of the proposed MFSAC-EC classifier is established with respect to the following measures.

Cross-Validation method

The first well-known metric used here to evaluate the classification model performance is the k -fold cross-validation method (Wang 2012). In the k -fold cross-validation method, the dataset is randomly divided into k number of folds and $k-1$ folds are used for training and one fold is used for testing. The process is repeated for k number of times and average classification accuracy is taken. When k is set at 1 that means the fold size is equal to the size of the dataset (training dataset size is equal to one less than the number of samples in the dataset and validation is done using the remaining sample) then it is considered as Leave one out cross-validation method (LOOCV). For k is equal to 2, the cross-validation method is named the household method. It has been found that when k is set at a very small value that means the fold size is large then the accuracy of the classification model is affected by low bias and high variance problems. On the other hand, if k is set at a high value that means the fold size is not so large then the classification accuracy of the classification model has a high bias but low variance. It has been found that 10-fold cross-validation method outperforms the LOOCV method (Breiman 1992, Ambroise 2002, Asyali 2006) and it has been also endorsed that the 10-fold cross-validation method as a better measure for classification.

In training-testing random splitting the dataset is initially randomly partitioned into training set ($2/3^{\text{rd}}$ of the dataset) and testing set ($1/3^{\text{rd}}$ of the dataset) with 50 runs.

ROC curve analysis

The performance of the proposed classifier for two-class datasets is also judged using Receiver Operator Characteristic (ROC) analysis (Wang 2012). It is a visual method for evaluating binary classification models. Under this analysis, the following measures are considered to judge the binary classification model.

Classification accuracy (Acc) is defined as,

$$Acc = \frac{TP + TN}{TP + FP + TN + FN} \quad 0 \leq Acc \leq 1$$

The sensitivity (SN) or True Positive Rate (TPR) can be defined as,

$$SN = TPR = \frac{TP}{TP + FN}$$

The specificity (SP) or True Negative Rate (TNR) can be defined as,

$$SP = TNR = \frac{TN}{TN + FP}$$

The False Positive Rate (FPR) can be defined as:

$$FPR = (1 - specificity) = \frac{FP}{FP + TN}$$

The Positive Predicted Value (PPV) can be defined as:

$$PPV = \frac{TP}{TP + FP}$$

The Negative Predicted Value (NPV) can be defined as:

$$NPV = \frac{TN}{TN + FN}$$

Where TP, TN, FP, FN are true positive, true negative, false positive, and false negative respectively.

The ROC curve is plotted considering TPR along the y-axis and FPR along the x-axis. The area under the ROC curve (AUC) is used to represent the performance of the binary classification model. The higher AUC value of a ROC curve for a particular classification model signifies the better performance of the classification model in differentiating positive and negative examples. The range of AUC value is $0 \leq AUC \leq 1$.

Heat map analysis

A heatmap is a data representation diagram in which the values for a variable of interest are portrayed using a data matrix. In this data matrix, the values of the variable are represented across two-axis variables as a grid of colored squares. The axis variables are divided into ranges and each cell's color represents the intensity of that variable for the particular ranges of values of axis variables.

Here, the performance of the proposed classifier for multi-class datasets is judged using Heat map representation of confusion matrix (Liu 2014), where a confusion matrix is a tabular representation to visualize the performance of a classification model in terms of true positive, true negative, false positive and false negative.

Parameter Estimation

Before running the MFSAC-EC, the parameters are settled down. In MFSAC-EC the input training dataset is bootstrapped. The proposed MFSAC-EC model is run here varying the number of bootstrapped datasets (D) from 5 to 30 and the classification accuracy of this model is more or less the same from 10 to the rest of the range. So, the number of bootstrapped datasets for every training dataset for this model is set at 10.

In MFSAC method, initially P number of genes is selected by each filter method. Here in Table 2, the classification accuracy of the proposed model is shown with respect to different values of P . From Table 2, it has been found that the proposed model gives the best result for $P=100$ for RAOA and RAHC datasets. In case of Breast cancer, Lung cancer, MLL and SRBCT datasets it gives the best result at $P = 200$. For Leukemia datasets it gives the best result at $P= 500$. So, it can be said that MFSAC-EC gives best result for P value within 200 to 500 in all cases for all datasets except Colon and Prostate. In Colon and Prostate, it shows the best result for $P = 1500$. Here we have used SVM, DT (C4.5), NB, and KNN classifiers individually for forming different ensemble classification models. All the classifiers are implemented using Scikit-learn libraries of Python. For all classifiers, we have set parameters with default parameter values. For DT as default setting we have used splitting function = Gini, Splitting criterion = best, height = none (that means for every sample it reaches a leaf/class node). For SVM, we have used the RBF kernel function. For KNN we have chosen K (number of nearest neighbor) value from 3 to 7. The overall execution time of a single run of the MFSAC-EC model (considering bootstrapped dataset creation, feature selection using MFSAC, and then generating classification accuracy of test samples using LOOCV, 5-fold, 10-fold, and random splitting) and testing time using only 10-fold are shown for different datasets in Table 3.

Classification Performance of the Proposed MFSAC-EC Classifier

In Table 4, using the LOOCV method, the classification accuracy of our proposed MFSAC-EC model is 100% for different datasets (Leukemia, Breast, RBreast, Lung, RAOA, and RAHC) for all cases. In the Prostate dataset, we did not get 100% accuracy using our model with respect to any type of existing classifier. In MLL, Colon, and SRBCT it also gives 100% accuracy using all types of ensemble classifiers.

In Table 5 and Table 6, it has been shown that using 5 -fold and 10-fold cross-validation, MFSAC-EC does not provide 100% accuracy only for Colon and Prostate cancer datasets. For other datasets, it provides 100% accuracy with respect to all types of ensemble classifiers. To show the generalization property of the proposed ensemble classifiers, the classification accuracy of these classifiers is also measured repeatedly with respect to the random splitting of the dataset into a training set (2/3 data of original dataset) and test set (1/3 data of original dataset). Random splitting is done with care such that class proportion is alike in the training set and test set. In Table 7, the classification accuracy of the above mentioned four different types of ensemble classifiers for the different number of cluster representatives is shown in different datasets which are based on the best result of 50 random splitting of the dataset into a training set (2/3 data of original dataset) and test set (1/3 data of original dataset).

From the results of Table 4 to Table 7, it has been observed that classification accuracy in the LOOCV method, 5-fold cross-validation, and 10-fold cross-validation methods is higher than the random splitting of the dataset, and the overall generalization performance of the proposed classification model is also good.

The performance of the proposed model for different two-class datasets with respect to different parameters like SN, SP, PPV, NPV, FPR is shown in Table 8. From this table, it is found that the performance of the proposed model is very good with respect to all these parameters for all two-class datasets.

In Figures 4, the ROC curve is shown for different two-class datasets. In Figures 4a, 4b, and 4c, the ROC curves are shown for Breast cancer using LOOCV, for Colon cancer using 5-fold cross validation, and for RAHC dataset using 10-fold cross-validation respectively. The ROC curves for Leukemia Cancer, and Lung cancer datasets using LOOCV are given in Supplemental Figures F1a and F1b respectively. For Breast cancer, Leukemia cancer, and Lung cancer, the AUC value is equal to 1.0 in every case. The ROC curves are shown for RAOA, and RBreast cancer datasets using 5-fold cross-validation in Supplemental Figures F2a, and F2b respectively. For these datasets also the prediction accuracy using 5-fold cross validation is very high according to the AUC value. In Supplemental Figures F2c, the ROC curves are shown for Prostate cancer using 10-fold cross-validation. From these curves of 10-fold cross validation, it may be seen that except for Prostate cancer, for all other datasets the AUC value is 1 and for Prostate cancer, the AUC value is close to 1.

In Figures 5a and 5b, heatmap representation of the confusion matrix are shown for multi-class datasets: SRBCT and MLL with respect to 5-fold cross-validation, and 10-fold cross-validation respectively. From these figures, it is clear that for the proposed model prediction accuracy is accurate in most cases.

Comparison of MFSAC-EC Model with Well-Known Existing Filter Methods used in this model

In Supplemental Figure F3, the proposed MFSAC-EC model in combination with different existing classifiers is compared with different filter methods used in this model with respect to SRBCT, RAHC, Prostate, and Colon datasets in terms of 10-fold cross-validation. In all cases, the performance of the proposed model is significantly better with respect to all filters.

Comparison of MFSAC-EC Model with Well-Known Existing Gene Selection Methods

In Figure 6, the MFSAC-EC model with different existing classifiers as base classifiers are compared with existing well-known supervised gene selection methods named mRMR (minimum redundancy maximum relevance framework) (Ding 2005), MSG (mutual information based supervised gene clustering algorithm) (Maji 2012), CFS (Correlation-based Feature Selection) (Ruiz 2006), and FCBF (Fast Correlation-Based Filter) (Ruiz 2006) with respect to

different classifiers using 10-fold cross-validation method. From these results, it has been found that the proposed model outperforms in most of the cases.

In Figure 7, the MFSAC-EC model is compared with well-known existing unsupervised gene selection methods named MGSACO (Tabakhi 2015), UFSACO (Tabakhi 2014), RSM(Lai 2006), MC (Haindl 2006), RRFS (Ferreira 2012), TV (Theodoridis 2008), and LS (Liao 2014) with respect to DT, SVM, NB classifiers using random splitting method. From these results, it can be said that the MFSAC-EC model outperforms in all cases.

Comparison of MFSAC-EC Model with Well-Known Existing Ensemble Classification and DEEP learning Models

In Table 9, the proposed MFSAC-EC model using the DT classifier is compared with well-known existing ensemble classification models with respect to 10-fold cross-validation. These models are PCA-basedRotBoost (Osareh 2013), ICA-based RotBoost (Osareh 2013), AdaBoost (Osareh 2013), Bagging (Osareh 2013), Arcing (Osareh 2013), Rotation Forest (Osareh 2013), EN-NEW1 (Wang 2006), and EN-NEW2 (Wang 2006). From Table 9, it is clear that the proposed model using DT classifier outperforms in all cases.

In Table 10, the proposed MFSAC-EC model using DT, NB, KNN as base classifiers are compared with different existing ensemble classifiers with respect to 10-fold cross-validation. These classifiers are Bagging based ensemble classifier (Nagi 2013), Boosting based ensemble classifier (Nagi 2013), Stacking based ensemble classifier (Nagi 2013), Heuristic breadth-first search-based ensemble classifier (HBSA) (Wang 2012), Sd_Ens (Nagi 2013), and Meta_Ens (Nagi 2013). In Table 11 our model using SVM and KNN as base classifiers is compared with auto-encoder-based deep learning models (Nabendu 2020) in terms of random splitting. Here, results are shown only for the datasets for which results are available in the literature, and all other fields are marked as “Not Found”. In all cases, the MFSAC-EC model outperforms all the well-known existing ensemble models (except for the Colon cancer dataset) and deep learning models which in turn validates the usefulness of the proposed model.

Biological Significance Analysis

The top 8 genes selected by the MFSAC-EC model for Colon cancer and Leukemia are listed in Table 12. For every gene, the name and symbol of the gene as well as the Accession number of the Affymetrix chip are listed. Apart from this information, to validate those genes, biomedical literature of the genes is searched and for every gene, the corresponding reference about its role and significance for a particular disease is provided.

Discussion

In this paper, a new Multiple Filtering and Supervised Attribute Clustering algorithm-based ensemble classification model named MFSAC-EC is proposed. The main motivation behind this work is to develop a machine learning-based ensemble classification model to overcome the over-fitting problem which arises due to the presence of sample class imbalance problem, small

sample size problem, and also high dimensional feature set problem in the microarray gene expression dataset, to enhance the prediction capability of the proposed model.

Nowadays, in designing machine learning models, the use of ensemble methodology has been increasing day by day as it incorporates multiple learning algorithms and also training datasets in different efficient manners to improve the overall prediction accuracy of the model. Due to the inclusion of prediction accuracy of multiple learning models and also the use of different bootstrapping datasets, the chances of potential overfitting in training data is greatly reduced in the ensemble models, and as a consequence the prediction accuracy increases. One necessary condition of the superior performance of an ensemble classifier with respect to its individual member/base classifier is that every base classifier should be very accurate and diverse (Osareh 2013). A classifier is considered accurate if its generalization capability is high and two classifiers satisfy diverse property if their prediction in classifying the same unknown samples vary from each other. The general principle of ensemble methods is to rearrange training datasets in different ways (either by resampling or reweighting) and build an ensemble of base classifiers by applying a base classifier on every rearranged training dataset (Osareh 2013).

In our proposed ensemble model, at first, a number of bootstrapped datasets of the original training dataset is created. In every bootstrapped dataset, the class imbalance problem is solved using the oversampling method. Then for every bootstrapped dataset, a number of sub-datasets are created using the MFSAC method (which is a hybrid method combining multi-filters and a new supervised attribute/gene clustering method), and then for every generated sub dataset, a base classifier is constructed using any existing classification model. After that, a new ensemble classifier (EC) is formed using the majority voting scheme by combining the prediction accuracy of all those base classifiers.

The prediction accuracy of the proposed model is verified by applying it to high-dimensional microarray gene expression data. From Figure 6, and Figure 7 it has been found that the classification accuracy of the MFSAC-EC model is much better than the well-known existing gene selection methods. From Table 9, Table 10, and Table 11, it has been also found that the proposed MFSAC-EC classification model is superior to the existing ensemble classification models in almost every case. The superior performance of the proposed model is due to the following reasons:

- The generation of the different bootstrapped versions of training data and also the use of the oversampling procedure to balance the cardinality of majority class and minority class in every bootstrapped dataset reduces the chances of the overfitting problem of a classifier.
- Different types of filter methods are used in the MFSAC method. It has been already observed that one filter gives better performance for one dataset while the same gives poor results for other datasets. This is because every filter uses separate metrics and so the choice for a filter for a specific dataset is a very complex task. As different filter methods are used in the MFSAC method, so different sub-datasets with different characteristics-based attributes/genes are formed from each dataset. This is shown using

Venn diagram in Supplemental Figures F4a and F4b. Here for Leukemia and Prostate cancer datasets, the first twenty genes, selected by each filter are shown. In case of Leukemia dataset, Relief measure generates non-overlapping gene subset while using other filter metrics presence of a small number of overlapping genes in different gene subsets are observed. In Prostate cancer dataset, Relief generates non-overlapping gene subset and also maximum number of genes are non-overlapping in gene subsets formed by Fisher score, MI (mutual information). From these figures, it is clear that using different filter methods different subsets of genes are selected and different sub datasets are formed. It shows diversity of those filter methods. As a consequence, the base classifiers prepared on these diverse datasets are become diverse. This diversity increases the power of ensemble classifier.

- Moreover, the genes selected by different filter methods are good biomarker also. In Table 12, the top ranked 8 genes selected by MFSAC-EC model are shown for Leukemia and Colon cancer datasets. Among these genes, gene MPO (with column number 1720), CST3 (with column number 1823), ZYX (with column number 4788), CTSD (with column number 2062), CD79A/MB-1 (with column number 2583), LYZ (with column number 6738) in Leukemia dataset are important biomarkers as these are selected by different filter methods mentioned in Supplemental Figure F4.
- In MFSAC, at first, a sub dataset of the most relevant genes is selected by each filter method. Then on each sub dataset, the proposed supervised gene clustering algorithm is applied and a reduced sub dataset of modified attributes/features in the form of augmented cluster representatives is generated. In this method, at the time of cluster formation, genes are augmented based on their supervised information. In other words, such augmentation is considered where it increases the class discrimination power. Thus effectively, the class relevance of any augmented cluster representative is greater than that of any single gene involved in that process. So, this modified sub dataset containing a reduced feature set in the form of augmented cluster representatives is more powerful according to class discrimination power than the sub dataset containing a subset of the most relevant genes. Apart from this, it is well known fact in gene expression data that two genes are functionally similar if they are pattern-based similar (either positively co-expressed or negatively co-expressed) (Das 2016). So, at the time of the augmentation procedure, two types of augmentations are considered here. One is that a gene is added with its original value with the current cluster representative and another one is that the gene is added with its sign-flipped value with the current cluster representative. This is because if the current cluster representative and a gene are positively co-expressed then normal addition is considered but if they are negatively co-expressed then normal addition will hamper the addition process and in that case, sign-flipping of that gene will give proper result. The effect of augmentation with respect to every filter method is shown in Figure 8. In Figure 8, for the Breast cancer dataset, at the time of supervised cluster formation from each filter generated subset, the original gene, and its

corresponding class relevance value, and also augmented gene and its corresponding class relevance are shown. From Figure 8, it is clear that for every filter method the class relevance score of every original gene is increased with respect to that filter after augmentation. In Figure 8, different class labels are distinguished by different colors.

- Finally, for each sub dataset with modified attributes in the form of augmented cluster representatives, a classifier is constructed using any existing classifier, and these classifiers are combined using the majority voting technique to form an ensemble classifier (EC). The use of different sub-datasets with optimal gene subsets in the form of augmented cluster representatives and the formation of a classifier for every sub dataset can solve the overfitting problem of any single classifier. This is due to the reason that not all sub-datasets can consistently perform well on all types of cancer datasets (due to inherent characteristics of the datasets), but due to the use of majority voting in ensemble classifiers, this problem can be solved or reduced.

Another outcome of our proposed model is to rank informative genes for every cancer dataset. For this task, the frequency of occurrence of each gene present in the form of augmented cluster representatives in every sub dataset is counted and these genes are ranked according to the counted value to measure the importance of those genes for any specific disease, here cancer. To establish the biological significance of those selected genes for every cancer dataset, their contribution has been confirmed by other existing studies where they are referred already. From these existing studies, it is clear that the selected genes are important for cancer class discrimination and also are important as cancer biomarkers for molecular treatment targets.

Conclusions

Many machine learning and statistical learning-based classifiers for sample classification already exist in the literature, but these methods are prone to suffer from overfitting due to small sample size problems, class imbalance problems, and the curse of the high dimensionality of microarray data. Although some of the existing methods can mitigate these issues to quite an extent, the problems have still not been satisfactorily overcome. Due to this reason, here a novel feature selection-based ensemble classification model named MFSAC-EC is proposed. It has been shown that the proposed model can handle the above-mentioned issues present in existing models. To check the performance of the proposed MFSAC-EC model, this classifier is applied to test sample classification accuracy in high dimensional microarray gene expression data, a domain that will be beneficial in the field of cancer research. From the experimental results, it has been found that the proposed model outperforms all other well-known existing classification models combined with the different recognized feature selection methods and also the newly developed ensemble classifiers for all types of cancer datasets mentioned here. Apart from this classification task, the proposed model can also rank informative attributes according to their importance. The efficiency of the proposed model in this task is vindicated by finding the most informative genes for Colon cancer and Leukemia cancer datasets using this model. These genes are biologically validated based on other well-known existing studies. Consequently, it is clear

that the selected genes are vital for sample class discrimination and are also important biomarkers for molecular treatment targets of deadly diseases.

Acknowledgments

References

- Alon U, Barkai N, Notterman D. A., Gish K., Ybarra S., Mack D., and Levine A. J. 1999. Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by Oligonucleotide arrays. *Proceedings of National Academy of Sciences, USA.* 96(12): 6745–6750.
- Ambroise C, McLachlan GJ. 2002. Selection bias in gene extraction on the basis of microarray gene-expression data. *Proceedings of National Academy of Sciences, USA.* 99(10):6562–6566.
- Armstrong S, Staunton J, Silverman L, Pieters R, den Boer M, Minden M, Sallan S, Lander E, Golub T and Korsmeyer S. 2001. MLL translocations specify a distinct gene expression profile that distinguishes a unique leukemia. *Nature Genetics.* 30:41–47.
- Asyali MH, Colak D, Demirkaya O, Inan MS. 2006. Gene expression profile classification: A review. *Current Bioinformatics.* 1(1):55–73.
- Austin H Chen, Yin-Wu Tsau and Ching-Heng Lin. 2010. Novel methods to identify biologically relevant genes for leukemia and prostate cancer from gene expression profiles. *BMC Genomics.*
- Au W. H., Chan K. C. C., Wong A.K.C., and Wang Y. 2005. Attribute clustering for grouping, selection, classification of gene expression data. *IEEE/ACM Transactions of Computational Biology and Bioinformatics.* 2(2):83–101.
- Bai YX, Yi Ji-Lin, Li Jian-Feng, Hong Sui. 2007 Clinicopathologic significance of BAG1 and TIMP3 expression in colon carcinoma. *World Journal of Gastroenterology.* 13(28): 3883–3885.
- Błaszczyński Jerzy, Stefanowski Ł Jerzy, Idkowiak ukasz. 2013. Extending Bagging for Imbalanced Data. *Proceedings of the 8th International Conference on Computer Recognition Systems (CORES 2013).*
- Bolo' n-Canedo V, Sánchez-Marño Noelia, Alonso-Betanzos Amparo. 2012. An ensemble of filters and classifiers for microarray data classification. *Pattern Recognition.* 45(1):531–539.
- Botchkina Inna L, Rowe Rebecca A, Rivadeneira David E, Karpeh Martin S Jr, Crawford Howard, Dufour Antoine, Ju Jingfang, Wang Yuan, Leyfman Yan, Botchkina Galina I. 2009. Phenotypic Subpopulations of Metastatic Colon Cancer Stem Cells: Genomic Analysis. *Cancer Genomics & Proteomics.* 6(1): 19–30.
- Breiman L, Spector P. 1992. Submodel selection and evaluation regression – the X-random case. *International Statistical Review.* 60(3):291–319.
- Brown Adam R, Simmen C.M. Rosalia, Raj Vinay R., Van Trang T., MacLeod Stewart L, Simmen Frank A. 2015. Krüppel-like factor 9 (KLF9) prevents colorectal cancer through inhibition of interferon-related signaling. *Carcinogenesis.* 36(9): 946–955.
- Chen Xiao, Gao Bing, Ponnusamy Murugavel, Lin Zhijuan, Liu Jia. 2017. MEF2 signaling and human diseases. *Oncotarget.* 8(67):112152–112165.
- Chin A J, Mirzal A, Haron H, Hamed H N A. 2016. Supervised, Unsupervised and Semi-supervised Feature Selection: A Review on Gene Selection. *IEEE Transactions on Computational Biology and Bioinformatics.* 13(5):971–989.
- Colozza M, Cardoso F, Sotiriou C, Larismont D, Piccart M J. 2005. Bringing molecular prognosis and prediction to the clinic. *Clin Breast Cancer.* 6(1):61–76.
- Das Chandra, Bose Shilpi, Chattopadhyay Matangini, Chattopadhyay Samiran. 2016. A novel distance-based iterative sequential KNN algorithm for estimation of missing values in microarray gene expression data. *IJBRA.* 12(4):312.
- Das Chandra, Bose Shilpi, Banerjee Abhik, Dutta Sourav, Ghosh Kuntal, Chattopadhyay Matangini. 2019. Comparative Performance Analysis of Different Measures to Select Disease Related Informative Genes from

Microarray Gene Expression Data. International Conference on Innovation in Modern Science and Technology (ICIMSAT-2019), Springer.

Dashtban M. and Balafar M. 2017. Gene selection for microarray cancer classification using a new evolutionary method employing artificial intelligence concepts. *Genomics*.109(2): 91-107.

Dettling M and Buhlmann P. 2002. Supervised Clustering of Genes. *Genome Biology*. 3(12): 0069.1-0069.15.

Dietterich T. G.. 2000. An experimental comparison of three methods for constructing ensembles of decision trees: bagging, boosting, and randomization. *Machine Learning*. 40(2): 139–157.

Ding C. and Peng H. 2005. Minimum redundancy feature selection from microarray gene expression data. *Journal of Bioinformatics and Computational Biology*. 3(2): 185–205.

Duda R. O., Hart P. E., and Stork D. G.1999. *Pattern Classification and Scene Analysis*. NewYork:Wiley.

Durai Rajaraman, Yang Y Shi, Seifalian M Alexander, Goldspink Geoffrey, Winslet C Marc. 2007. Role of insulin-like growth factor binding protein-4 in prevention of colon cancer. *World Journal of Surgical Oncology*. 5:128.

Elyasigomari V, Lee D. A., Screen H. R.C, Shaheed M.H. 2017. Development of a two-stage gene selection method that incorporates a novel hybrid approach using the cuckoo optimization algorithm and harmony search for cancer classification. *Journal of Biomedical Informatics*. 67:11-20.

Feng Hailiang, Liu Yanyan, Bian Xiaocui, Zhou Fangying, Liu Yuqin. 2018. ALDH1A3 affects colon cancer in vitro proliferation and invasion depending on CXCR4 status. *British Journal of Cancer*. 118: 224–232.

Ferreira A.J., Figueiredo M.A.T. 2012. An unsupervised approach to feature discretization and selection, *Pattern Recognition*. 45: 3048-3060.

Furey T.S., Cristianini N, Duffy N, Bednarski D.W, Schummer M, and Haussler D. 2000. Support Vector Machine Classification and Validation of Cancer Tissue Samples Using Microarray Expression Data. *Bioinformatics*. 16(10): 906-914.

Gardina Paul J, Clark Tyson A, Shimada Brian, Staples Michelle K, Yang Qing, Veitch James, Schweitzer Anthony, Awad Tarif, Sugnet Charles, Dee Suzanne, Davies Christopher, Williams Alan, Turpaz Yaron. 2006. Alternative splicing and differential gene expression in colon cancer detected by a whole genome exon array. *BMC GENOMICS*. 7: 325

Ghorai Santanu, Mukherjee Anirban, Sengupta Sanghamitra, Dutta Pranab K.. 2011. Cancer Classification from Gene Expression Data by NPPC Ensemble. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*. 8(3).

Giorgio Eros Di, Hancock Wayne W, Brancolini Claudio. 2018. MEF2 and the tumorigenic process, hic sunt leones. *BBA - Reviews on Cancer*. 1870(2):261-273.

Golub T. R., Slonim D. K., Tamayo P, Huard C, Gaasenbeek M, Mesirov J.P., Coller H, Loh H.L., Downing J. R., Caligiuri M. A., Bloomfield C. D., and Lander E. S.. 1999. Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring. *Science*. 286(5439):531-537.

Gordon G. J., Jensen R. V., Hsiao L.-L., Gullans S. R., Blumenstock J. E., Ramaswamy S., Richards W. G., Sugarbaker D. J., and Bueno R. 2002. Translation of microarray data into clinically relevant cancer diagnostic tests using gene expression ratios in lung cancer and mesothelioma. *Cancer Research*. 62: 4963–4967.

Greller LD, Tobin FL. 1999. Detecting selective expression of genes and proteins. *Genome Research*. 9:282–296.

Gu Q., Li Z., Han J.. Generalized Fisher Score for Feature Selection. 2011. *UAI'11: Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*.

Hai Wang, Hu Haiyan, Zhang Qian, Yang Yadong, Li Yanming, Hu Yang, Ruan Xiuyan, Yang Yaran, Zhang Zhaojun, Shu Chang, Yan Jiangwei, Wakeland Edward K, Li Quanzhen, Hu Songnian, Fang Xiangdong. 2013. Dynamic transcriptomes of human myeloid leukemia cells. *Genomics*. 102:250–256.

Haindl M, Somol P, Ververidis D, Kotropoulos C.2006. Feature Selection Based on Mutual Correlation. *Pattern Recognition, Image Analysis and Applications*. Springer Berlin Heidelberg. 569-577.

Hambali Moshood A., Oladele Tinuke O., Adewole Kayode S. 2020. Microarray cancer feature selection: Review, challenges and research directions. *International Journal of Cognitive Computing in Engineering*.1: 78-97.

872 Handschuh Luiza. 2019. Not Only Mutations Matter: Molecular Picture of Acute Myeloid Leukemia Emerging from
873 Transcriptome Studies. *Journal of Oncology*.
874 Hastie T, Tibshirani R, Eisen M.B., Alizadeh A, Levy R, Staudt L, Chan W. C., Botstein D, and Brown P. 2000.
875 'Gene Shaving' as a Method for Identifying Distinct Sets of Genes with Similar Expression Patterns. *Genome*
876 *Biology*. 1(2):1-21.
877 Hastie T, Tibshirani R, Botstein D, and Brown P. 2001. Supervised Harvesting of Expression Trees. *Genome*
878 *Biology*. 1: 1-12.
879 Hunter JD. 2007. Matplotlib: a 2D graphics environment. *Computing In Science & Engineering* 9(3):90–95 DOI
880 10.1109/MCSE.2007.
881 Kamal Amany M, El-Hefny Nadia H, Hegab Hany M, El-Mesallamy Hala O. 2016. Expression of thioredoxin-1
882 (TXN) and its relation with oxidative DNA damage and treatment outcome in adult AML and ALL: A comparative
883 study. *Hematology*. 21(10):567-575.
884 Karlenius Therese Christina, Tonissen Kathryn Fay. 2010. Thioredoxin and Cancer: A Role for Thioredoxin in all
885 States of Tumor Oxygenation. *Cancers (Basel)*. 2(2): 209–232.
886 Kim Yundeok , Yoon Sulhee, Kim Soo Jeong, Kim Jin Seok, Cheong Jun-Won, Hong Yoo Min. 2012.
887 Myeloperoxidase Expression in Acute Myeloid Leukemia Helps Identifying Patients to Benefit from Transplant.
888 *Yonsei Med J*. 53(3): 530–536.
889 Khan J, Wei JS, Ringner M, Saal LH, Ladanyi M, Westermann F, Berthold F, Schwab M, Antonescu CR, Peterson
890 C, Meltzer S Paul. 2001. Classification and diagnostic prediction of cancers using gene expression profiling and
891 artificial neural networks. *Nat Medicine*. 7(6):673–679.
892 Klimiankou Maksim, Uenalalan Murat, Kandabaraui Siarhei, Nustede Rainer, Steiert Ingeborg, Mellor-Heineke
893 Sabine, Zeidler Cornelia, Skokowa Julia, Welte Karl. 2019. Ultra-Sensitive CSF3R Deep Sequencing in Patients
894 With Severe Congenital Neutropenia. *Front. Immunol*.
895 Komer B, Bergstra J, Eliasmith C. 2014. Hyperopt-sklearn: automatic hyperparameter configuration for scikit-learn.
896 In: *Proceedings of the 13th Python in Science Conference (SCIPY 2014)*. 33–39.
897 Kourou Konstantina, Exarcos T P, Exarcos K P, Karmouzis M V, Fotaidis D. 2015. Machine learning applications
898 in cancer prognosis and prediction. *Computational and Structural Biotechnology Journal*. 13: 8–17.
899 Kozlov Igor, Beason Kevin, Yu Cheng, Hughson Michael. 2005, CD79a expression in acute myeloid leukemia
900 t(8;21) and the importance of cytogenetics in the diagnosis of leukemias with immunophenotypic ambiguity. *Cancer*
901 *Genetics and Cytogenetics* 163(1):62-7.
902 Kozovska. Z, Patsalias A., Bajzik V, Durinikova E, Demkova L, Jargasova S, Smolkova B, Plava J, Kucerova L,
903 Matuskova M. 2018. ALDH1A inhibition sensitizes colon cancer cells to chemotherapy. *BMC Cancer*. 18:656.
904 Lagunas-Rangel Francisco Alejandro, Chávez-Valencia Venice, Gómez-Guijosa Miguel Ángel, Cortes-Penagos
905 Carlos. 2017. Acute Myeloid Leukemia—Genetic Alterations and Their Clinical Prognosis. *Int J Hematol Oncol*
906 *Stem Cell Res*. 11(4): 328–339.
907 Lai C, Reinders MJT, Wessels L. 2006. Random subspace method for multivariate feature selection. *Pattern*
908 *Recognition Letters*., 27: 1067-1076.
909 Lance Amanda, Druhan Lawrence J, Vestal C Greer, Steuerwald Nury M, Hamilton Alicia, Smith Mathew, Price
910 Andrea, Tjaden Elise, Fox Andee N, Avalos Belinda R. 2020. Altered expression of CSF3R splice variants impacts
911 signal response and is associated with SRSF2 mutations. *Leukemia*. 34(2):369-379.
912 Leung Yukyee and Hung Yeungsam. 2010. A Multiple-Filter-Multiple-Wrapper Approach to Gene Selection and
913 Microarray Data Classification”, *IEEE/ACM Transactions on Computational Biology and Bioinformatics*. 7(1).
914 Léveillard Thierry, Aït-Ali Najatel. 2017. Cell Signaling with Extracellular Thioredoxin and Thioredoxin-Like
915 Proteins: Insight into Their Mechanisms of Action. *Oxidative Medicine and Cellular Longevity*.
916 Li Z, Xie W, Liu T. 2018. Efficient feature selection and classification for microarray data. *Plos one*.
917 Liao B, Jiang Y, Liang W, Zhu W, Cai L, Cao Z. 2014. Gene selection using locality sensitive Laplacian score.
918 *IEEE/ACM Transactions on Computational Biology and Bioinformatics*. DOI: 10.1109/TCBB.2014.2328334.
919 Liu Jian, Cheng Yuhu, Wang Xuesong, Zhang Lin, Wang Z. Jane. 2018. Cancer Characteristic Gene Selection via
920 Sample Learning Based on Deep Sparse Filtering. *Scientific Reports, Nature*. 8: 8270.

- 921 Liu Q, Sung A H, Chen Z, Liu J, Chen L, Qiao M, Wang Z, Huang X, Deng Y. 2011. Gene selection and
922 classification for cancer microarray data based on machine learning and similarity measures. BMC Genomics.
923 12(5):S1.
- 924 Liu Sheng, Dissanayake Shamitha, Patel Sanjay V, Dang Xin. 2014. Learning accurate and interpretable models
925 based on regularized random forests regression. BMC Systems Biology. 8 Suppl 3(Suppl 3):S5.
- 926 M Lizet. Waals van der, M. H. Inne. Rinkes Borel, Kranenburg Onno. 2018. ALDH1A1 expression is associated
927 with poor differentiation, 'right-sidedness' and poor survival in human colorectal cancer. PLOS ONE.
- 928 Maji P, Das C. 2012. Relevant and Significant Supervised Gene Clusters for Microarray Cancer Classification. IEEE
929 Transactions on Nanobioscience. 11(2).
- 930 Nabendu Bhui, Pintu Kumar Ram, Pratyay Kuila. 2020. Feature Selection from Microarray Data based on Deep
931 Learning Approach. ICCCNT2020.
- 932 Nada A, Alshamlan H. 2019. A survey on hybrid feature selection methods in microarray gene expression data for
933 cancer classification", IEEE Access.
- 934 Nagi Sajid, Bhattacharyya Kr. Dhruva. 2013. Classification of microarray cancer data using ensemble approach.
935 Netw Model Anal Health Inform Bioinformatics. 2:159–173.
- 936 Osareh Alireza, Bitu Shadgar. 2013. An Efficient Ensemble Learning Method for Gene Microarray Classification.
937 BioMed Research International, Hindawi Publishing Corporation.
- 938 Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R,
939 Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E. 2011. Scikit-learn:
940 machine learning in python. Journal of Machine Learning Research 12:2825–2830.
- 941 Pilling M J, Henderson A, Gardner P. 2017. Quantum Cascade Laser Spectral Histopathology: Breast Cancer
942 Diagnostics Using High Throughput Chemical Imaging. Analytical Chemistry. 89 (14), 7348-7355.
- 943 Pouw Kraan T.C.T.M. Van der, Gaalen F.A. van, Kasperkovitz P.V., Verbeet N.L., Smeets T.J.M., Kraan M.C.,
944 Fero M., Tak P.-P., Huizinga T.W.J. , Pieterman E., Breedveld F.C., Alizadeh A.A. , and Verweij C.L. 2003.
945 Rheumatoid Arthritis is a Heterogeneous Disease:Evidence for Differences in the Activation of the STAT-1
946 Pathway between Rheumatoid Tissues. Arthritis and Rheumatism. 48(8):2132-2145.
- 947 Pouw Kraan T.C.T.M. van der, Wijbrandts C.A., Baarsen L.G.M. van, Voskuyl A.E., Rustenburg F., Baggen J.M.,
948 Ibrahim S.M., Fero M., Dijkmans B.A.C., Tak P.P., and Verweij C.L. 2007. Rheumatoid Arthritis Subtypes
949 Identified by Genomic Profiling of Peripheral Blood Cells: Assignment of a Type I Interferon Signature in a
950 Subpopulation of Patients. Annals of the Rheumatic Diseases. 66: 1008-1014.
- 951 Ritter Malte, Klimiankou Maksim, Klimenkova Olga, Schambach Axel, Hoffmann Dirk, Schmidt Amy, Kanz
952 Lothar, Link Daniel C., Welte Karl, Skokowa Julia. 2020. Cooperating, congenital neutropenia–associated Csf3r and
953 Runx1 mutations activate pro-inflammatory signaling and inhibit myeloid differentiation of mouse HSPCs. Annals
954 of Hematology. 99: 2329–2338.
- 955 Ruiz R., Riquelme J.C., Aguilar-Ruiz J.S. 2006. Incremental wrapper-based gene selection from microarray data for
956 cancer classification. Journal of Pattern Recognition. 39(12): 2383–2392.
- 957 Simmen Frank A, Su Ying, Xiao Rijin, Zeng Zhaoyang, Simmen CM Rosalia. 2008. The Krüppel-like factor 9
958 (KLF9) network in HEC-1-A endometrial carcinoma cells suggests the carcinogenic potential of dys-regulated
959 KLF9 expression. Reproductive Biology and Endocrinology.
- 960 Singh D., Febbo P.G., Ross K., Jackson D. G., Manola J., Ladd C., Tamayo P., Renshaw A. A., Amico A. V. D,
961 Richie J. P. , Lander E. S., Loda M., Kantoff P.W., Golub T.R., and Sellers W.R..2002. Gene expression correlates
962 of clinical prostate cancer behavior. Cancer Research. 62: 203–209.
- 963 Singh P, Dai B, Dhruva B, Widen S G 1994. Episomal Expression of Sense and Antisense Insulin-like Growth
964 Factor (IGF) binding Protein-4 Complementary DNA Alters the Mitogenic Response of a Human Colon Cancer Cell
965 Line (HT-29) by Mechanisms That Are Independent of and Dependent upon IGF-11. Cancer Research. 54: 6563-
966 6570.
- 967 Su AI, Welsh JB, Sapinoso LM, Kern SG, Dimitrov P, Lapp H, Schultz PG, Powell SM, Moskaluk CA, Frierson
968 HF, Hampton GM. 2001. Molecular classification of human carcinomas by use of gene expression signatures.
969 Cancer Research. 61:7388–7393.

- 970 Su Chun-Wen, Lin Chiao-Wen, Yang Wei-En, Yang Shun-Fa. 2019. TIMP-3 as a therapeutic target for cancer.
971 Therapeutic Advances in Medical Oncology.
- 972 Su Li, Luo Yongli, Yang Zhi, Yang Jing, Yao Chao, Cheng Feifei, Shan Juanjuan, Chen Jun, Li Fangfang, Liu
973 Limei, Liu Chungang, Xu Yanmin, Jiang Lupin, Deyu Guo, Prieto Jesus, Ávila A Matías, Shen Junjie, Qian Cheng.
974 2016. MEF2D Transduces Microenvironment Stimuli to ZEB1 to Promote Epithelial–Mesenchymal Transition and
975 Metastasis in Colorectal Cancer. *Molecular and Cellular Pathobiology*. 76(17):5054-5067.
- 976 Swan A L, Mobasheri Ali, Allaway David, Liddell Susan, Bacardit Jaume. 2013. Application of Machine Learning
977 to Proteomics Data: Classification and Biomarker Identification in Postgenomics Biology. *OMICS*. 17(12): 595–
978 610.
- 979 Szuber Natasha, Tefferi Ayalew. 2018. Chronic neutrophilic leukemia: new science and new diagnostic criteria.
980 *Blood Cancer Journal*. 19.
- 981 Tabares-Soto Reinel, Orozco-Arias Simon, Romero-Cano Victor, Segovia Bucheli Vanesa, Rodríguez-Sotelo José
982 Luis , Jiménez-Varón Cristian Felipe. 2020. A comparative study of machine learning and deep learning algorithms
983 to classify cancer types based on microarray gene expression data. *PeerJ Computer Science*.
- 984 Tabakhi S, Moradi P, Akhlaghian F. 2014. An unsupervised feature selection algorithm based on ant colony
985 optimization. *Engineering Applications of Artificial Intelligence*. 32: 112-123.
- 986 Tabakhi S, Najafi, A., Ranjbar, R., &Moradi, P. 2015. Gene selection for microarray data classification using a
987 novel ant colony optimization. *Neurocomputing*, 168: 1024-1036.
- 988 Tang Zhenjie, Yuan Shuqiang, Hu Yumin, Zhang Hui, Wu Wenjing, Zeng Zhaolei, Yang Jing, Yun Jingping, Xu
989 Ruihua, Huang Peng. 2012. Over-expression of GAPDH in human colorectal carcinoma as a preferred target of 3-
990 Bromopyruvate Propyl Ester” Published in final edited form as: *J Bioenerg Biomembr*. 44(1): 117–125.
- 991 Theodoridis S, Koutroumbas K. 2008. *Pattern Recognition*, fourth ed., Elsevier Science.
- 992 Thorsen Kasper, Sørensen D Karina, Brems-Eskildsen Sofie Anne, Modin Charlotte, Gaustadnes Mette, Hein K
993 Anne-Mette, Kruhøffer Mogens, Laurberg Soren, Borre Michael, Wang Kai, Brunak Søren, Krainer R Adrian,
994 Tørring Niels , Dyrskjøl Lars, Andersen Claus L, Orntoftet Torben F.2008. Alternative Splicing in Colon, Bladder,
995 and Prostate Cancer Identified by Exon Array Analysis. *Molecular & Cellular Proteomics*. 7: 1214-1224.
- 996 Tong Dong Ling, Ball Graham R. 2014. Exploration of Leukemia Gene Regulatory Networks Using A Systems
997 Biology Approach. 2014 IEEE International Conference on Bioinformatics and Biomedicine.
- 998 Vapnik V. 1995. *The Nature of Statistical Learning Theory*. New York: Springer-Verlag. Veer L. J. V., Dai H.,
999 Vijver M. J. V. D., He Y. D., Hart A. A.M., Mao M., Peterse H. L., Kooy K. v. d., Marton M. J., Witteveen A. T.,
1000 Schreiber G.J. , Kerkhoven R.M., Roberts C. , Linsley P. S., Bernards R., and Friend S. H..2002. Gene expression
1001 profiling predicts clinical outcome of breast cancer. *Nature*. 415(6871): 530–536.
- 1002 Wang Ching Wei. 2006. New Ensemble Machine Learning Method for Classification and Prediction on Gene
1003 Expression Data. *Proceedings of the 28th IEEE EMBS Annual International Conference New York City, USA*.
- 1004 Wang Shu-Lin, Li Xue-Ling, Fanget Jianwen. 2012. Finding minimum gene subsets with heuristic breadth-first
1005 search algorithm for robust tumor classification. *BMC Bioinformatics*.13:178.
- 1006 West M., Blanchette C., Dressman H., Huang E., Ishida S., Spang R., Zuzan H. , Olson J. A., Marks J. R., and
1007 Nevins J. R.. 2001. Predicting the clinical status of human breast cancer by using gene expression profiles. *Proc.*
1008 *Natl. Acad. Sci. USA*. 98(20):11 462–11 467.
- 1009 Yang P, Zhou B Bing, Zhang Zili, Zomaya Y, Albert. 2010. A multi-filter enhanced genetic ensemble system for
1010 gene selection and sample classification of microarray data. *BMC Bioinformatics*. 11: S5.
- 1011 Yan Zhi, Li Jiangeng, Xiong Yimin, Xu Weitian, Zheng Guorong. 2012. Identification of candidate colon cancer
1012 biomarkers by applying a random forest approach on microarray data. *Oncology Reports*, SPANDIDOS
1013 Publications. 28(3):1036-42.
- 1014 Yunsong Q, Yang Xibei. 2013. Interval-valued analysis for discriminative gene selection and tissue sample
1015 classification using microarray data. *Genomics*. 101(1):38-48.
- 1016 Ying Mingyao, Tilghman Jessica, Wei Yingying, Guerrero-Cazares Hugo, Alfredo Quinones-Hinojosa, Ji Hongkai,
1017 Lateral John. 2014. KLF9 Inhibits Glioblastoma Stemness through Global Transcription Repression and Integrin-α6
1018 Inhibition. *Journal for Biochemistry and Molecular Biology*. 289(47): 32742–32756.

1019 Yu Herbert, Rohan Thomas. 2000. Role of the Insulin-Like Growth Factor Family in Cancer Development and
 1020 Progression. *Journal of the National Cancer Institute*. 92(18): 1472–1489.
 1021 Zhang Jin-Ying , Zhang Fan, Hong Chao-Qun, Giuliano Armando E., Cui Xiao-Jiang, Zhou Guang-Ji , Zhang Guo-
 1022 Jun, Cui Yu-Kun. 2015. Critical protein GAPDH and its regulatory mechanisms in cancer cells. *Cancer Biol Med*
 1023 12(1):10-22.
 1024 Zhang Yang, Wang Fang, Chen Xue, Zhang Yu, Wang Mingyu, Liu Hong, Cao Panxiang, Ma Xiaoli, Wang Tong,
 1025 Zhang Jianping, Zhang Xian, Lu Peihua, Liu Hongxing 2018. CSF3R Mutations are frequently associated with
 1026 abnormalities of RUNX1, CBFB, CEBPA, and NPM1 genes in acute myeloid leukemia. *Cancer*. 124(16):3329-
 1027 3338.
 1028 Zhou Nina and Wang Lipo. 2007. A Modified T-test Feature Selection Method and Its Application on the HapMap
 1029 Genotype Data. *Genomics Proteomics Bioinformatics*. 5(3-4): 242-9.
 1030 Zhu Kongxi, Wang Yunxia, Liu Lan, Li Shuai, Yu Weihua. 2019. Long non-coding RNA MBNL1-AS1 regulates
 1031 proliferation, migration, and invasion of cancer stem cells in colon cancer by interacting with MYL9 via sponging
 1032 microRNA-412-3p. *Clinics and Research in Hepatology and Gastroenterology*. 44(1):101-114.

Table 1 (on next page)

Description of Cancer Gene Expression Datasets

Table 1. Description of Cancer Gene Expression Datasets

Dataset	Data Dimension Gene \times Sample (Original)	Data Dimension Gene \times Sample (Used)	Sample Class Labels	Dataset	Data Dimension Gene \times Sample (Original)	Data Dimension Gene \times Sample (Used)	Sample Class Labels
Leukemia	7129×72	7070×72	2	Breast	7129×49	7129×49	2
Colon	2000×62	2000×62	2	MLL	12582×72	12582×72	3
Prostate	12600×136	12600×136	2	SRBCT	2308×63	2308×63	4
Lung	12533×181	12533×181	2	RAHC	41057×50	41057×50	2
Rbreast	24481×97	24188×97	2	RAOA	18433×30	18433×30	2

Table 2 (on next page)

Classification Accuracy of MFSAC-EC depending on varying number of genes selected by each Filter

This table shows the impact of parameter P with respect to sample classification accuracy(%) in terms of both LOOCV and 10-Fold Cross Validation approach. P defines the number of top ranked genes selected by each filter method.

1

Table 2: Classification Accuracy of MFSAC-EC depending on varying number of genes selected by each Filter

Dataset	Evaluation Metric	MFSAC-EC																							
		P=100				P=200				P=500				P=1000				P=1200				P=1500			
		NB	KNN	DT	SVM	NB	KNN	DT	SVM	NB	KNN	DT	SVM	NB	KNN	DT	SVM	NB	KNN	DT	SVM	NB	KNN	DT	SVM
Leukemia	LOOCV	98.6	98.6	98.6	98.6	98.6	98.6	98.6	98.6	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
	10 Fold	98.6	98.6	98.6	98.6	98.6	98.6	98.6	98.6	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
RAHC	LOOCV	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
	10 Fold	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
MLL	LOOCV	98.6	100	100	97.2	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
	10 Fold	97.2	100	100	97.2	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
RAOA	LOOCV	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
	10 Fold	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
SRBCT	LOOCV	98.4	98.4	100	98.4	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
	10 Fold	100	98.4	98.4	98.4	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
Breast	LOOCV	98	95.9	93.9	95.9	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
	10 Fold	100	95.9	95.9	98	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
Lung	LOOCV	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
	10 Fold	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
Rbreast	LOOCV	92.6	93.7	93.7	95.8	99	97.9	100	99	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
	10 Fold	91.6	96.8	95.8	97.9	97.9	99	99	99	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
COLON	LOOCV	91.9	91.9	93.6	91.9	91.9	91.9	96.8	91.9	98.4	98.4	96.8	98.4	98.4	98.4	98.4	98.4	98.4	98.4	100	100	100	98.4	100	100
	10 Fold	91.9	91.9	93.6	91.9	91.9	91.9	95.2	91.9	98.4	96.8	96.8	96.8	98.4	100	98.4	98.4	100	98.4	98.4	100	100	100	98.4	100
Prostrate	LOOCV	83.8	90.4	92.7	86.8	88.2	92.7	92.7	88.2	91.2	95.6	97.1	91.9	94.9	97.1	97.8	94.1	98.5	98.5	97.8	98.5	98.5	99.3	98.5	99.3
	10 Fold	85.3	88.2	91.9	86.8	88.2	91.9	93.4	89.7	91.9	95.6	97.8	91.2	95.6	97.8	97.8	94.1	99.3	98.5	97.8	98.5	99.3	99.3	97.8	99.3

2

Table 3 (on next page)

Total Execution Time in a single run of MFSAC-EC on Different Datasets

Total execution time in a single run of MFSAC-EC including Bootstrapped dataset creation, Feature Selection by filter methods and supervised attribute clustering approach, Training, Testing using LOOCV, 5-Fold, 10-Fold, and Random Splitting is given in the first row. While execution time using only 10-Fold Cross Validation is given in the 2nd row. Here the time for the best P value is shown here.

1 **Table 3: Total Execution Time (Bootstrapping, Feature Selection, Training, Testing for LOOCV, 5-Fold, 10-Fold, Random Splitting) in a single run of MFSAC-EC on Different Datasets**

	Leukemia	RAHC	MLL	RAOA	SRBCT	Breast	Lung	Rbreast	COLON	Prostrate
No. of Feature selected for best result	500	100	200	100	200	200	100	500	1200	3000
Total Time Taken	8mins 23secs	7mins 32secs	7mins 54secs	4mins 43secs	5mins 17secs	4mins 2secs	11mins 14secs	10mins 22secs	17mins 40secs	1hr 18mins 41secs
Time Taken for only 10 fold	35secs	30secs	41secs	36secs	30secs	32secs	36secs	33secs	30secs	36secs

2

Table 4(on next page)

Classification Accuracy of the proposed MFSAC-EC model with respect to LOOCV

Classification accuracy (%) of MFSAC-EC model has been shown in terms of LOOCV with respect to four ensemble classifiers MFSAC-EC + NB, MFSAC-EC+KNN, MFSAC-EC+DT, and MFSAC-EC+SVM. Every ensemble classifier is run 50 times using LOOCV for every dataset and the accuracy is shown which is obtained maximum number of times.

Table 4. Classification Accuracy of the proposed MFSAC-EC model with respect to LOOCV

Dataset	Proposed Model		Cluster Representatives			Dataset	Proposed Model		Cluster Representatives		
			1	2	3				1	2	3
COLON	MFSAC-EC	NB	100	98.39	98.39	MLL	MFSAC-EC	NB	100	100	100
		KNN	98.39	100	100			KNN	100	100	100
		DT	98.39	98.39	98.39			DT	100	100	100
		SVM	100	98.4	98.4			SVM	100	100	100
Prostate		NB	97.06	97.79	98.53	SRBCT		NB	96.83	100	100
		KNN	97.79	97.79	98.53			KNN	96.83	100	100
		DT	97.79	98.53	97.79			DT	96.83	98.41	100
		SVM	98.53	99.26	99.26			SVM	82.54	98.41	100
Leukemia		NB	100	100	100	Lung		NB	100	100	100
		KNN	100	100	100			KNN	100	100	100
		DT	100	100	100			DT	100	100	100
		SVM	100	100	100			SVM	100	100	100
RAOA		NB	100	100	100	RAHC		NB	100	100	100
		KNN	100	100	100			KNN	100	100	100
		DT	100	100	100			DT	100	100	100
		SVM	100	100	100			SVM	100	100	100
Breast		NB	100	100	100	RBreast		NB	100	100	100
		KNN	100	100	100			KNN	100	100	100
		DT	100	100	100			DT	100	100	100
		SVM	100	100	100			SVM	100	100	100

Table 5 (on next page)

Classification Accuracy of the proposed MFSAC-EC model with respect to 5-Fold Cross Validation

Classification accuracy (%) of MFSAC-EC model has been shown in terms of 5-Fold Cross Validation with respect to four ensemble classifiers MFSAC-EC + NB, MFSAC-EC+KNN, MFSAC-EC+DT, and MFSAC-EC+SVM. Every ensemble classifier is run 50 times using 5-Fold Cross Validation for every dataset and the accuracy is shown which is obtained maximum number of times.

Table 5. Classification Accuracy of the proposed MFSAC-EC model with respect to 5-Fold Cross Validation

Dataset	Proposed Model		Cluster Representatives			Dataset	Proposed Model		Cluster Representatives		
			1	2	3				1	2	3
COLON	MFSAC-EC	NB	96.77	96.77	96.77	MLL	MFSAC-EC	NB	100	100	100
		KNN	98.39	96.77	96.77			KNN	98.61	100	100
		DT	98.39	96.77	98.39			DT	98.61	100	100
		SVM	98.39	96.77	96.77			SVM	100	100	100
Prostate		NB	97.06	97.79	98.53	SRBCT		NB	98.41	100	100
		KNN	97.79	97.79	99.26			KNN	96.83	100	100
		DT	97.06	97.79	94.85			DT	96.83	98.41	100
		SVM	97.79	98.53	99.26			SVM	96.83	100	100
Leukemia		NB	100	100	100	Lung		NB	100	100	100
		KNN	100	100	100			KNN	100	100	100
		DT	100	100	100			DT	100	99.44	100
		SVM	100	100	100			SVM	100	100	100
RAOA		NB	100	100	100	RAHC		NB	100	100	100
		KNN	100	100	100			KNN	100	100	100
		DT	100	100	100			DT	100	100	100
		SVM	100	100	100			SVM	100	100	100
Breast		NB	100	100	100	RBreast		NB	100	100	100
		KNN	100	100	100			KNN	100	100	100
		DT	100	100	100			DT	100	100	100
		SVM	100	100	100			SVM	100	100	100

Table 6 (on next page)

Classification Accuracy of the proposed MFSAC-EC model with respect to 10-Fold Cross Validation

Classification accuracy (%) of MFSAC-EC model has been shown in terms of 10-Fold Cross Validation with respect to four ensemble classifiers MFSAC-EC + NB, MFSAC-EC+KNN, MFSAC-EC+DT, and MFSAC-EC+SVM. Every ensemble classifier is run 50 times using 10-Fold Cross Validation for every dataset and the accuracy is shown which is obtained maximum number of times.

1

Table 6. Classification Accuracy of the proposed MFSAC-EC model with respect to 10-Fold Cross Validation

Dataset	Proposed Model		Cluster Representatives			Dataset	Proposed Model		Cluster Representatives		
			1	2	3				1	2	3
COLON	MFSAC-EC	NB	98.39	98.39	98.39	MLL	MFSAC-EC	NB	100	100	100
		KNN	98.39	98.39	100			KNN	100	100	100
		DT	98.39	98.39	98.39			DT	100	100	100
		SVM	98.39	98.39	98.39			SVM	100	100	100
Prostate		NB	97.06	97.79	98.53	SRBCT		NB	96.83	96.83	100
		KNN	97.79	97.79	99.26			KNN	92.06	100	100
		DT	97.06	97.79	94.85			DT	95.24	96.83	100
		SVM	97.79	98.53	99.26			SVM	80.95	92.06	100
Leukemia		NB	100	100	100	Lung		NB	100	100	100
		KNN	100	100	100			KNN	100	100	100
		DT	100	100	100			DT	100	100	100
		SVM	100	100	100			SVM	100	100	100
Breast		NB	100	100	100	RBreast		NB	100	100	100
		KNN	100	100	100			KNN	100	100	100
		DT	100	100	100			DT	100	100	100
		SVM	100	100	100			SVM	100	100	100
RAOA		NB	100	100	100	RAHC		NB	100	100	100
		KNN	100	100	100			KNN	100	100	100
		DT	100	100	100			DT	100	100	100
		SVM	100	100	100			SVM	100	100	100

2

Table 7 (on next page)

Classification Accuracy of the proposed MFSAC-EC model with respect to Random Splitting of the Datasets

Classification accuracy (%) of MFSAC-EC model has been shown in terms of random splitting with respect to four ensemble classifiers MFSAC-EC + NB, MFSAC-EC+KNN, MFSAC-EC+DT, and MFSAC-EC+SVM. Every ensemble classifier is run 50 times using random splitting for every dataset and the accuracy is shown which is obtained maximum number of times. For random splitting the dataset is divided into training (2/3) and testing (1/3) part 50 times randomly.

Table 7. Classification Accuracy of the proposed MFSAC-EC model with respect to Random Splitting of the Datasets

Dataset	Proposed Model		Cluster Representatives			Dataset	Proposed Model		Cluster Representatives		
			1	2	3				1	2	3
COLON	MFSAC-EC	NB	98.39	98.39	98.39	MLL	MFSAC-EC	NB	100	100	100
		KNN	98.39	98.39	98.39			KNN	100	100	100
		DT	98.39	98.39	98.39			DT	98.61	100	98.61
		SVM	98.39	100	98.39			SVM	100	100	100
Prostate		NB	94.68	95.74	93.62	SRBCT		NB	95	85	95
		KNN	97.87	96.81	92.55			KNN	95	100	90
		DT	94.68	94.68	94.68			DT	80	90	95
		SVM	94.68	96.81	94.68			SVM	65	75	95
Leukemia		NB	100	100	100	Lung		NB	100	100	100
		KNN	100	100	100			KNN	100	100	100
		DT	100	100	100			DT	100	100	100
		SVM	100	100	100			SVM	100	100	100
RAOA		NB	100	100	100	RAHC		NB	100	100	100
		KNN	100	100	100			KNN	100	100	100
		DT	100	100	100			DT	100	100	81.25
		SVM	100	100	100			SVM	100	100	81.25
Breast		NB	100	100	100	RBreast		NB	91.94	91.94	91.94
		KNN	100	100	100			KNN	85.48	87.10	83.87
		DT	100	100	100			DT	83.87	79.03	80.65
		SVM	100	100	100			SVM	93.55	91.94	91.94

Table 8(on next page)

Evaluation of MFSAC-EC classifier based on SN, SP, PPV, NPV, FPR for two class data sets with respect to LOOCV

The performance of the MFSAC-EC model for two class datasets is represented using Receiver Operator Characteristic (ROC) analysis. SN represents Sensitivity, SP represents Specificity, PPV represents Positive Predicted Value, NPV represents Negative Predicted Value, and FPR represents False Positive Rate.

Table 8. Evaluation of MFSAC-EC classifier based on SN, SP, PPV, NPV, FPR for two class data sets with respect to LOOCV

Dataset	Proposed Model	SN	SP	PPV	NPV	FPR	Dataset	Proposed Model	SN	SP	PPV	NPV	FPR		
Leukemia	MFSGC-EC	NB	100	100	100	0	Breast	MFSGC-EC	NB	100	100	100	0		
		KNN	100	100	100	0			KNN	100	100	100	0		
		DT	100	100	100	0			DT	100	100	100	0		
		SVM	100	100	100	0			SVM	100	100	100	0		
Prostate		NB	98.7	98.3	98.7	98.3	1.7		Rbreast	NB	100	100	100	100	0
		KNN	98.7	98.3	98.7	98.3	1.7			KNN	100	100	100	100	0
		DT	100	96.61	97.46	100	3.4			DT	100	100	100	100	0
		SVM	100	98.3	98.7	100	1.7			SVM	100	100	100	100	0
Colon		NB	100	100	100	100	0		Lung	NB	100	100	100	100	0
		KNN	100	100	100	100	0			KNN	100	100	100	100	0
		DT	100	100	100	100	0			DT	100	100	100	100	0
		SVM	100	100	100	100	0			SVM	100	100	100	100	0
RAHC		NB	100	100	100	100	0		RAOA	NB	100	100	100	100	0
		KNN	100	100	100	100	0			KNN	100	100	100	100	0
		DT	100	100	100	100	0			DT	100	100	100	100	0
		SVM	100	100	100	100	0			SVM	100	100	100	100	0

Table 9(on next page)

Comparison of MFSAC-EC + DT with different existing Ensemble Classifiers using DT in terms of 10-Fold Cross Validation

Here MFSAC-EC + DT model is compared with existing ensemble classifiers where DT is used as base classifier. C4.5 algorithm is used as DT.

Table 9. Comparison of MFSAC-EC using DT with different existing Ensemble Classifiers using DT in terms of 10-Fold Cross Validation

	MFSAC-EC	PCA-based RotBoost	ICA-based RotBoost	AdaBoost	Bagging	Arcing	Rotation Forest	EN-NEW1	EN-NEW2
Colon	98.39	95.48	96.1	94.97	94.92	69.35	95.21	79.03	83.87
Leukemia	100	98.75	98.77	98.22	97.47	Not Found	97.97	Not Found	Not Found
Breast	100	94.39	97.88	98.89	92.74	80.41	98.6	94.85	95.88
Lung	100	98.11	99.54	96.3	97.08	97.24	97.56	98.34	99.45
Prostate	97.79	Not Found	Not Found	90.44	94.12	87.5	Not Found	94.85	97.06
MLL	100	98.86	99.31	97.63	97.11	91.67	97.61	93.06	98.61
SRBCT	100	99.5	99.59	98.16	96.46	Not Found	97.44	Not Found	Not Found

Table 10(on next page)

Comparison of MFSAC-EC using DT, KNN, NB, SVM with different existing Ensemble Classifiers using DT, KNN, NB, SVM in terms of 10-Fold Cross Validation

Here classification accuracy (%) of four ensemble classifiers MFSAC-EC + NB, MFSAC-EC + KNN, MFSAC-EC+DT, and MFSAC-EC+SVM are shown with respect to results of other existing ensemble classifiers with the same base learners. The best accuracy (%) for every dataset is shown in bold.

1 **Table 10. Comparison of MFSAC-EC using DT, KNN, NB, SVM with different existing Ensemble Classifiers using DT, KNN, NB, SVM**
2 **in terms of 10-Fold Cross Validation**

Dataset	MFSAC-EC				Bagging			Boosting			Stacking			HBSA		SD_Ens	Meta_Ens
	DT	NB	KNN	SVM	DT	NB	KNN	DT	NB	KNN	DT	NB	KNN	KNN	SVM		
Leukemia	100	100	100	100	94.12	88.23	73.53	91.18	88.24	75.53	91.18	91.18	91.18	88.46	88.46	92.45	94.12
Colon	98.39	98.39	100	98.39	95.16	66.13	90.32	98.39	87.1	91.94	98.39	93.59	93.59	75	85	94.4	99.21
Prostate	97.79	99.26	99.26	99.26	26.47	26.47	38.24	26.47	26.47	52.94	26.47	26.47	52.94	85.29	97.06	52.94	52.94
Lung	100	100	100	100	91.28	96.64	97.32	81.88	95.3	97.99	97.99	97.99	96.64	Not Found	Not Found	81.88	97.99
Breast	100	100	100	100	78.95	36.84	68.42	68.42	36.84	68.42	68.42	68.42	68.42	Not Found	Not Found	73.49	79.87

3
4

Table 11(on next page)

Comparison of MFSAC-EC using SVM and KNN with respect to different existing deep learning Classifiers using random splitting

Here classification accuracy (%) of two ensemble classifiers MFSAC-EC + KNN, and MFSAC-EC+SVM are shown with respect to results of other existing ensemble classifiers with the same base learners. The best accuracy (%) for every dataset is shown in bold.

Table 11. Comparison of MFSAC-EC using SVM and KNN with respect to different existing deep learning Classifiers using random splitting

Dataset	SVM			KNN		
	MFSAC-EC	Folded Autoencoder	Autoencoder	MFSAC-EC	Folded Autoencoder	Autoencoder
Colon	100	90.15	73.11	98.39	81.09	56.97
Prostate	96.81	84.16	64.3	97.87	76.48	52.1
Leukemia	100	93.62	84.12	100	85.24	77.13

Table 12 (on next page)

List of genes selected by MFSAC-EC model for Colon and Leukemia cancer Datasets

Here second column represents the gene names while third column indicate the gene accession number. The fourth column indicates the description of the gene while the fifth column indicates the literature where it has been referred as cancer biomarker.

Table 12. List of genes selected by MFSAC-EC model for Colon and Leukemia cancer Datasets

Dataset	Gene Name	Accession Number	Description	Validation of Genes
Colon	TPM1	Hsa.1130	Human tropomyosin isoform mRNA, complete cds.	Gardina 2006, Thorsen 2008, Botchkina 2009
	IGFBP4	Hsa.1532	Human insulin-like growth factor binding protein-4 (IGFBP4) gene, promoter and complete cds.	Durai 2007, Singh 1994, Yu 2000
	MYL9	Hsa.1832	Myosin Regulatory Light Chain 2, Smooth Muscle Isoform (Human); contains element TAR1 repetitive element	Yan 2012, Zhu 2019
	ALDH1L1	Hsa.10224	Aldehyde Dehydrogenase, Mitochondrial X Precursor (Homo sapiens)	Feng 2018, Waals 2018, Kozovska 2018
	KLF9	Hsa.41338	Human mRNA for GC box binding protein/ Kruppel Like Factor 9, complete cds	Brown 2015, Ying 2014, Simmen 2008
	MEF2C	Hsa.5226	Myocyte-Specific Enhancer Factor 2, Isoform MEF2 (Homo sapiens)	Chen 2017, Giorgio 2018, Su2016
	GADPH	Hsa.1447	Glyceraldehyde 3-Phosphate Dehydrogenase	Zhang 2015, Tang 2012
	TIMP3	Hsa.11582	Metalloproteinase Inhibitor 3 Precursor	Su 2019, Bai 2007
Leukemia	TXN	X77584_at	TXN Thioredoxin	Kamal 2016, Léveillard 2017, Karlenius 2010
	CSF3R	M59820_at	CSF3R Colony stimulating factor 3 receptor (granulocyte)	Zhang 2018, , Ritter 2020, Klimiankou 2019, Lance 2020
	MPO	M19508_xpt3_s_at	MPO from Human myeloperoxidase gene	Szuber 2018, Kim 2012, Lagunas-Rangel 2017, Handschuh 2019
	LYZ	M21119_s_at	LYZ Lysozyme	Wang 2013, Liu 2018, Tong 2014
	CST3	M27891_at	CST3 Cystatin C (amyloid angiopathy and cerebral hemorrhage)	Austin 2010
	ZYX	X95735_at	Zyxin	Austin 2010, Yunsong 2013
	CTSD	M63138_at	CTSD Cathepsin D (lysosomal aspartyl protease)	Wang2013
	CD79A/ MB-1 gene	U05259_rna1_at	MB-1 membrane glycoprotein	Wang 2013, Kozlov 2005

Figure 1

Cluster Representative Refinement Procedure

Each row of the table represents the gene with its class relevance value in terms of Pearson correlation coefficient with respect to sample class row. TR+ and TR- represent the augmented gene with their class relevance score in terms of Pearson correlation coefficient with respect to sample class row.

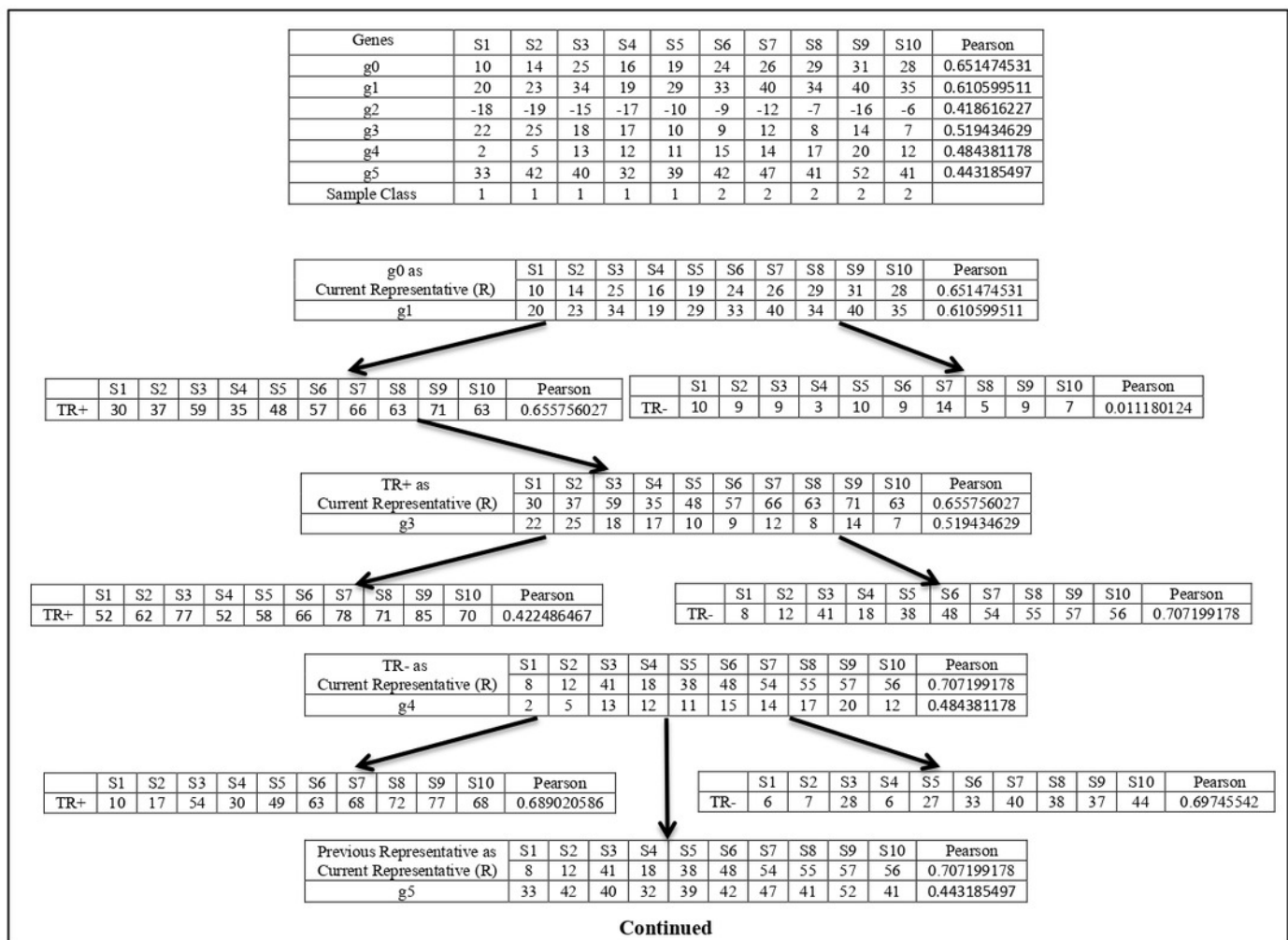


Figure 2

Block Diagram of the Proposed MFSAC-EC Model

Here BK_1, BK_2, \dots, BK_D are D number of bootstrapped datasets. $RSD_{11}, \dots, RSD_{17}$ represent different reduced sub datasets of BK_1 bootstrapped datasets after applying MFSAC method. IC_{11} to IC_{17} represent individual classifiers applied on $RSD_{11} \dots RSD_{17}$ respectively.

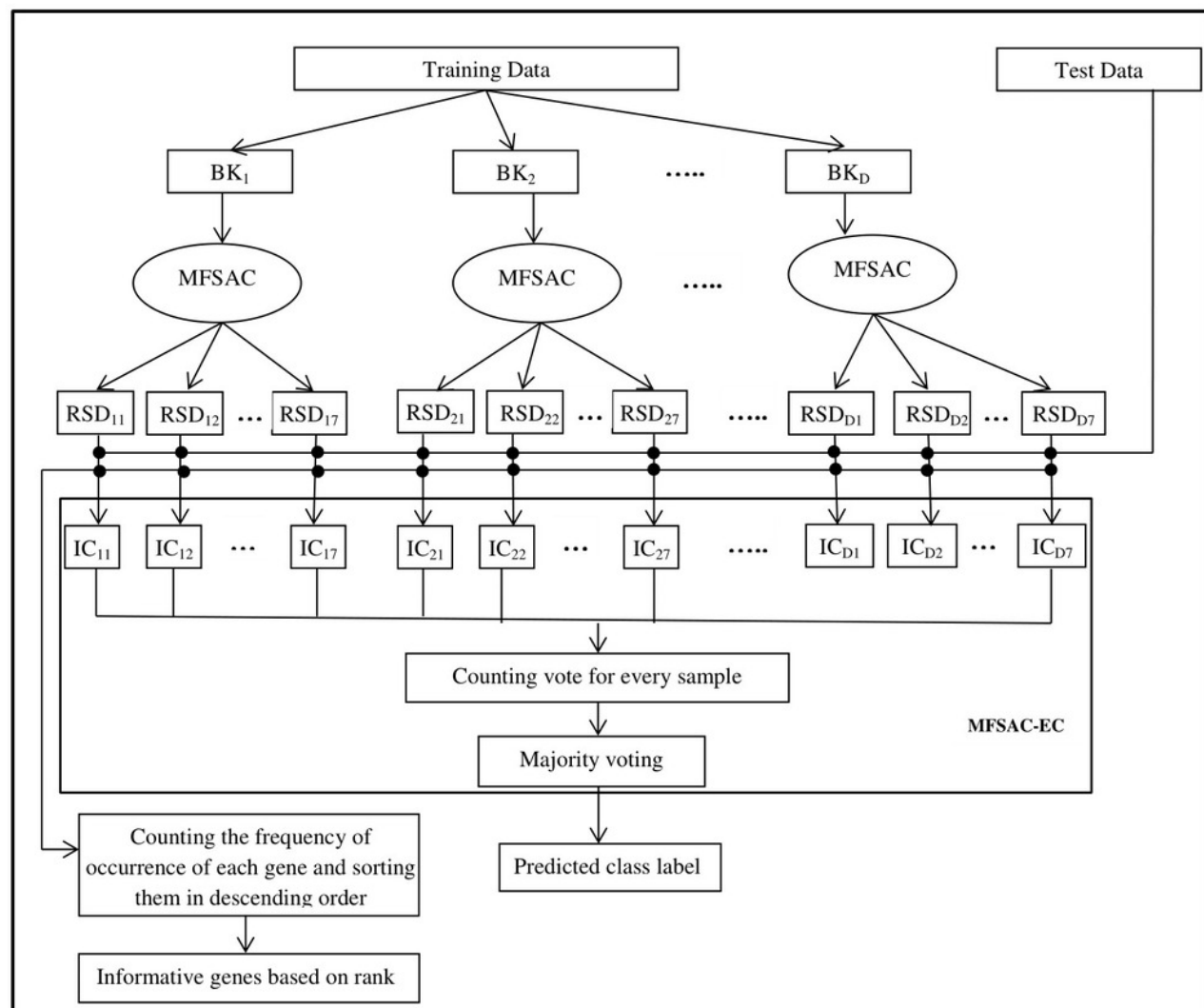


Figure 3

Block Diagram of MFSAC Method

BK_l is the l th bootstrapped dataset. FT_1 FT_7 are the seven filter score functions as Supplemental Table S1. SD_{11} SD_{17} are sub datasets created after applying filter score functions. SAC is the Supervised attribute clustering method applied to generate RSD_{11} RSD_{17} reduced sub datasets.

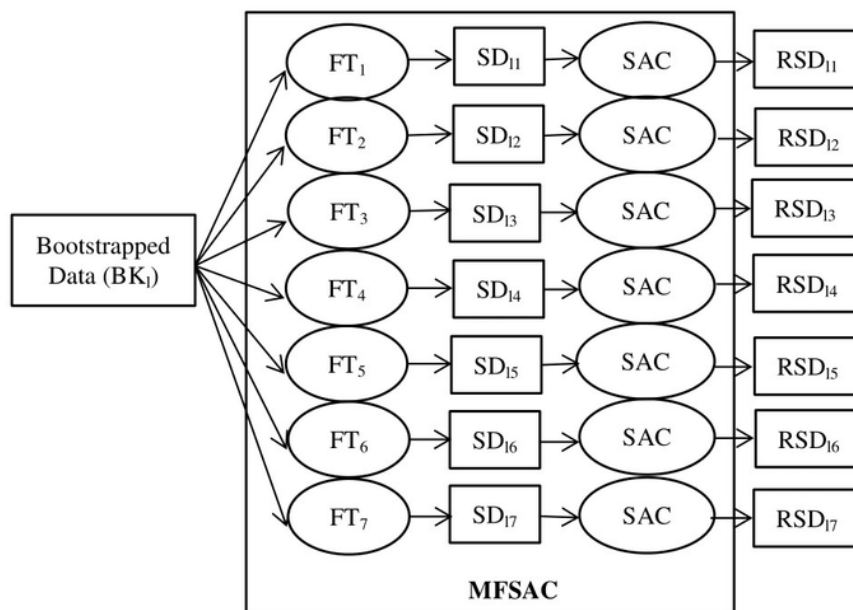


Figure 4

AUC for for three datasets using MFSAC-EC+ KNN, MFSAC-EC+NB, MFSAC-EC+ DT and MFSAC-EC+ SVM Classifiers.

(a) For Breast Cancer dataset using LOOCV. (b) For Colon Cancer dataset using 5-Fold Cross Validation. (c) For RAHC dataset using 10-Fold Cross-Validation

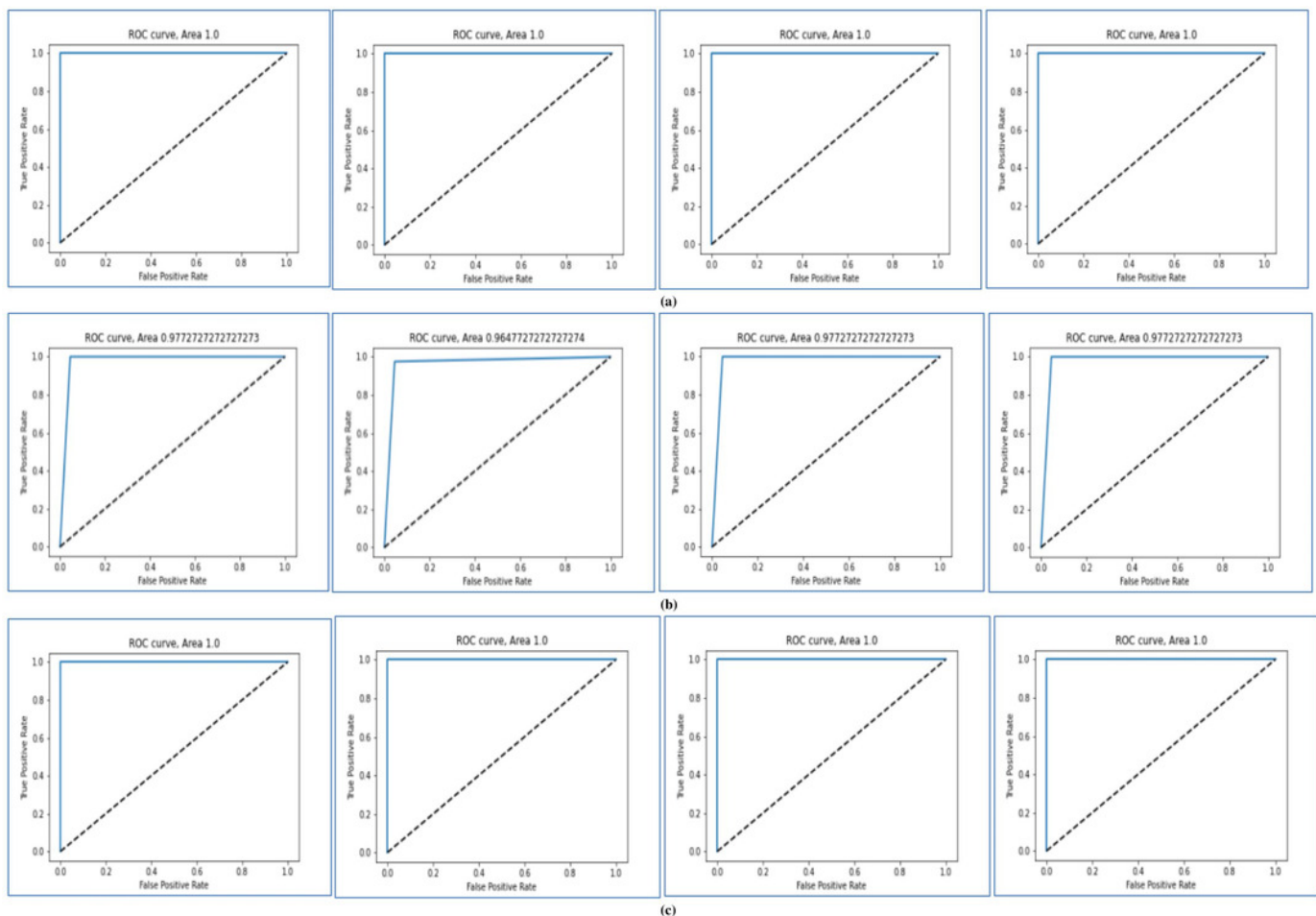


Figure 5

Heatmap of MFSAC-EC with base classifiers NB, KNN, DT and SVM respectively for multiclass datasets.

(a) For SRBCT dataset using 5-Fold Cross-validation. (b) For MLL dataset using 10-Fold Cross-Validation.

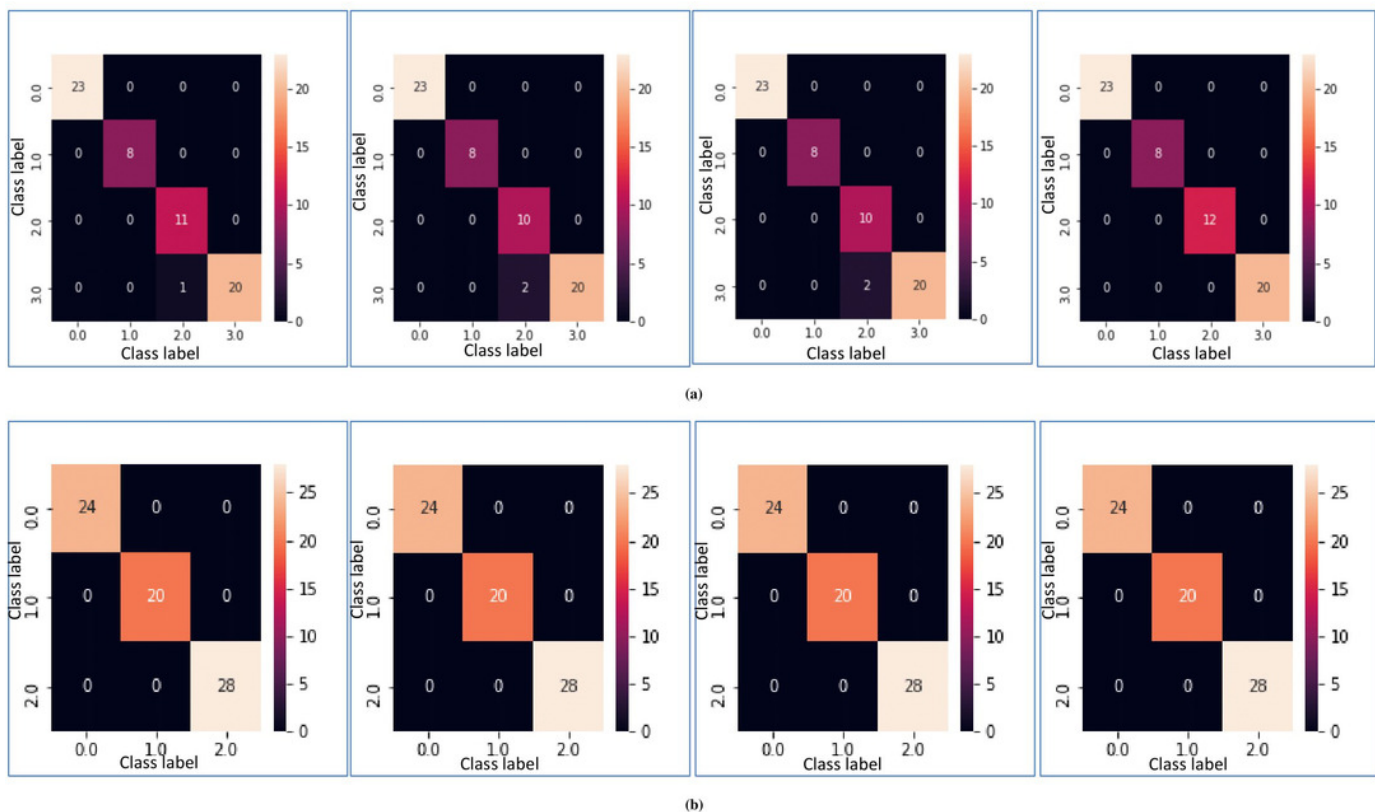


Figure 6

Comparison of MFSAC-EC with other well-known supervised gene selection methods and full gene set in terms of 10-Fold Cross-Validation for all Datasets.

In each figure classification accuracy(%) of MFSAC-EC model along with other supervised gene selection methods for all datasets are represented using different colored bars using (a) NB (b) KNN (c) DT and (d) SVM as base classifier.

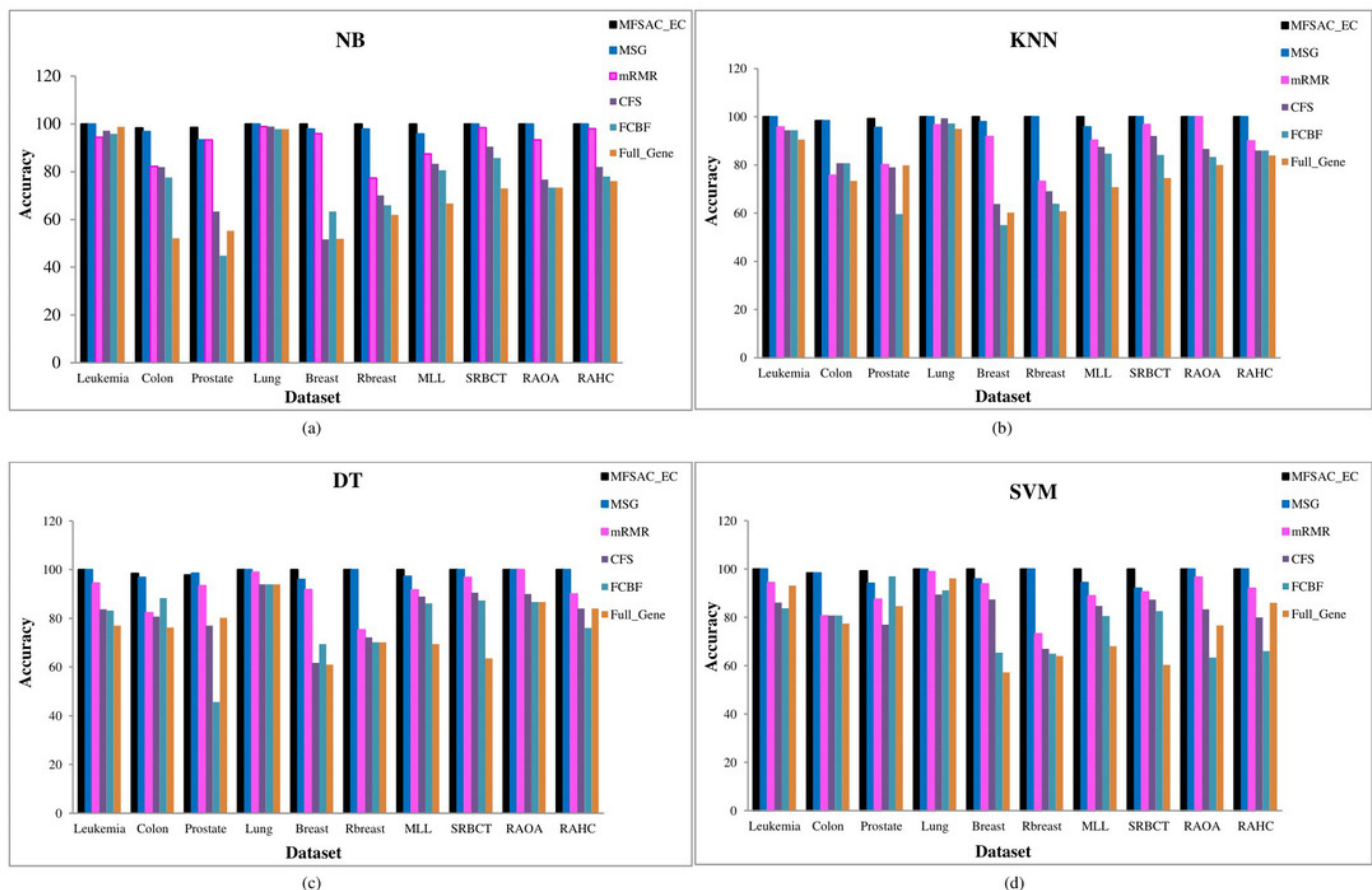
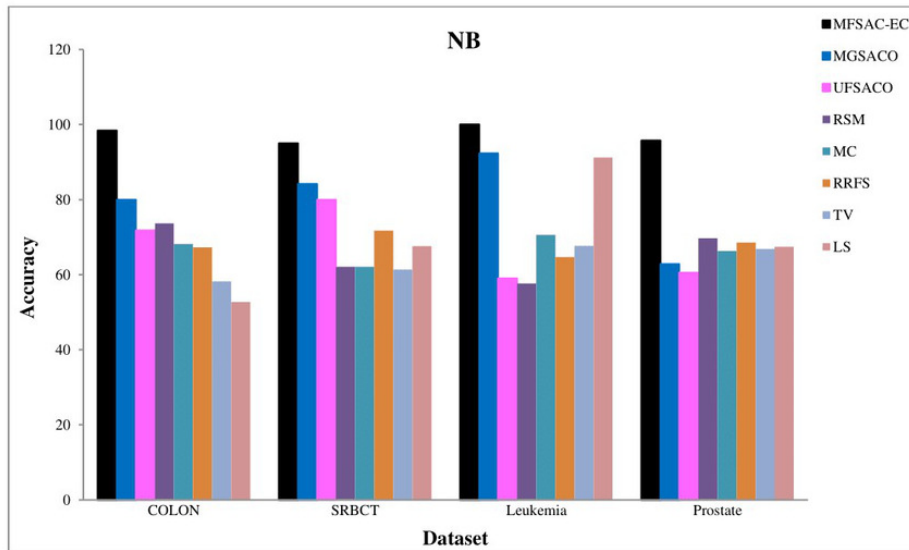


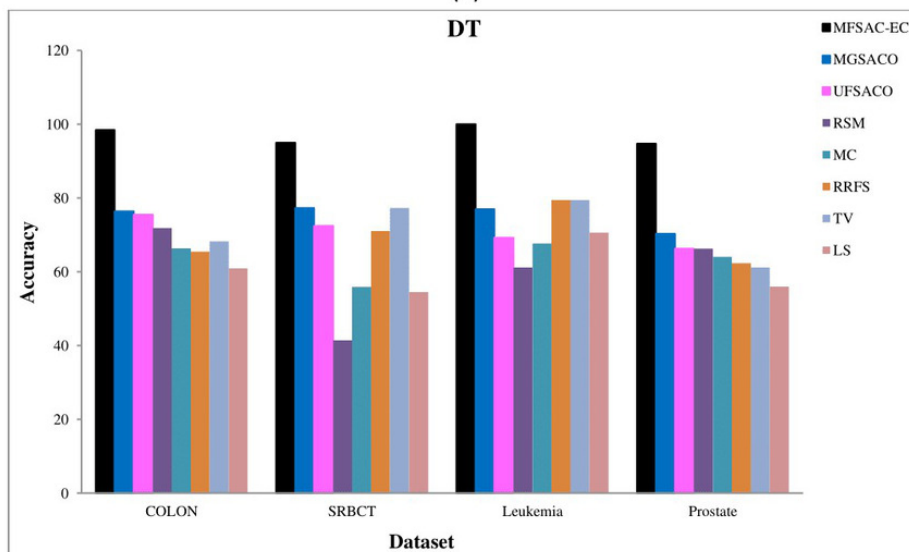
Figure 7

Comparison of MFSAC-EC with other well-known unsupervised gene selection methods in terms of random splitting for different datasets.

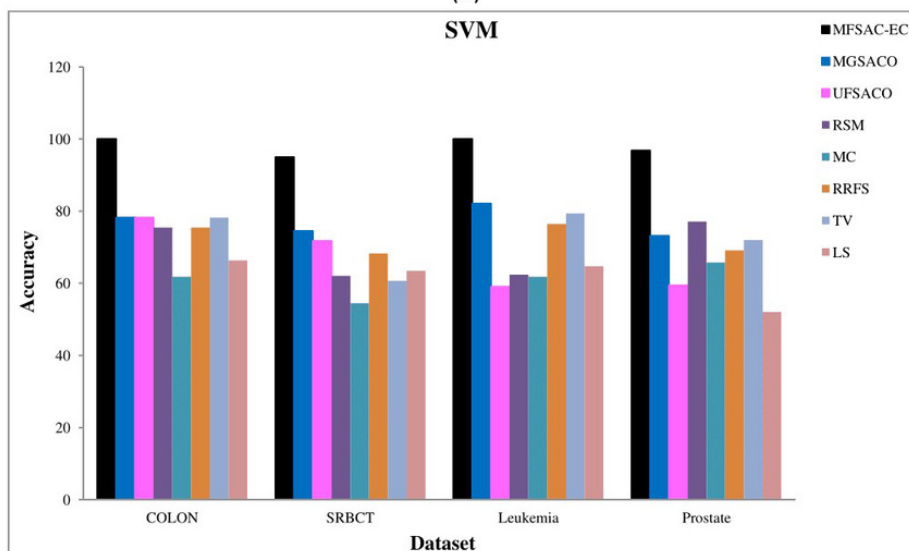
In each figure classification accuracy(%) of MFSAC-EC model along with other unsupervised gene selection methods for four datasets are represented with different colored bars using (a) NB (b) DT and (c) SVM as base classifier.



(a)



(b)



(c)

Figure 8

Original gene (different class label with different color) and corresponding Augmented gene with respect to different filter methods for Breast Cancer dataset

Seven figures for seven different filter score function are shown here. In each figure the original gene and augmented gene are plotted with respect to sample class label. X-axis represents class label while Y-axis represents expression value. Two different class labels are represented by Blue and Red color. The difference of expression values of two classes in the augmented gene shows class discrimination ability of that gene. Gene number is the column number in the original dataset.

