

1 **DOI: <https://doi.org/10.47391/JPMA.201>**

2
3 **Regulatory genes identification within functional genomics**
4 **experiments for tissue classification into binary classes via machine**
5 **learning techniques**

6
7 **Bushra Wazir¹, Dost Muhammad Khan², Umair Khalil³, Muhammad**
8 **Hamraz⁴, Naz Gul⁵, Zardad Khan⁶**

9 **1-5** Department of Statistics, Abdul Wali Khan University, Mardan, Pakistan; **6** Department of
10 Mathematical Sciences, University of Essex, United Kingdom

11 **Correspondence:** Zardad Khan **Email:** zkhan@essex.ac.uk

12
13 **Abstract**

14 **Objectives:** The aim of this study is to filter out the most informative genes that
15 mainly regulate the target tissue class, increase classification accuracy, reduce the
16 curse of dimensionality, and discard redundant and irrelevant genes.

17 **Methods:** This paper presented the idea of gene selection using bagging sub-forest
18 (BSF). The proposed method provided genes importance grounded on the idea
19 specified in the standard random forest algorithm. The new method is compared
20 with three state-of-the-art methods, i.e., Wilcoxon, masked painter and proportional
21 overlapped score (POS). These methods were applied on 5 data sets, i.e. Colon,
22 Lymph node breast cancer, Leukemia, Serrated colorectal carcinomas, and Breast
23 Cancer. Comparison was done by selecting top 20 genes by applying the gene
24 selection methods and applying random forest (RF) and support vector machine
25 (SVM) classifiers to assess their predictive performance on the datasets with

26 selected genes. Classification accuracy, Brier score, and sensitivity have been used
27 as performance measures.

28 **Results:** The proposed method gave better results than the other methods using
29 both random forest and SVM classifiers on all the datasets among all the feature
30 selection methods.

31 **Conclusion:** The proposed method showed improved performance in terms of
32 classification accuracy, Brier score and sensitivity, and hence, could be used as a
33 novel method for gene selection to classify tissue samples into their correct classes.

34 **Key Words:** Gene selection, classification, random forest, cancer, microarray gene
35 expression

37 **Introduction**

38 Cancer is a genetic disease due to changes in some of the genes that control the
39 way how our body cells function (for example, growth and division to make new
40 cells). Therefore, identification of such genes is important for the diagnosis of the
41 disease. High dimensional technologies produce a huge amount of data in many
42 research fields, such as biomedical science¹. These datasets, such as microarray
43 gene expression data, are known as high dimensional and contain a huge number
44 of irrelevant genes to the corresponding classification problem^{1,2}. High
45 dimensional gene expression data pose a number of challenges to the conventional
46 statistical tools, such as logistic regression and chi-square methods, used for their
47 analysis³. Analyzing high dimensional data, statistical models lose generalization
48 power and interpretability when applied to unseen data; in that very few genes
49 regulate the target class and rest are redundant (genes with similar expression
50 values) or non-informative^{4,5}. Moreover, these analyses also require considerable
51 computational resources³. Genes selection techniques are used to overcome these
52 problems with the main theme of selecting a subset of the most informative genes

53 to be used in model construction and prediction^{3,6,7}. Gene selection procedures
54 helps in discovering discriminative genes that regulate the target class and
55 eliminating irrelevant and useless genes that do not play any role in the regulation
56 of the response class³. Gene/feature selection methods are divided into three
57 categories, i.e Wrapper, Filter and Embedded. A brief discussion on these methods
58 is as follows:

59 **Filter methods:** These methods identify discriminative genes by calculating the
60 relevant score for each gene. Genes that have high relevance scores are selected for
61 the purpose of classification by the help of different classifiers. These methods do
62 not require classification algorithm for the selection of important genes. Moreover
63 filter methods deals with big data easily and are computationally fast and simple.
64 Examples of gene selection methods based on filtering approach could be found in
65 various studies¹.

66 **Wrapper methods:** In wrapper methods, gene subsets are evaluated by
67 partitioning the gene expression data into training and testing parts and running a
68 predictive model on the training parts corresponding to each gene subset. The
69 predictive model is then assessed by applying it on the testing part for each gene
70 subset and classification accuracy is calculated. Classification accuracy is used as
71 the corresponding score for each gene subset. Gene subset having the highest score
72 is selected as the final gene set to be used in tissue classification¹.

73 **Embedded methods:** Embedded methods select informative genes as part of
74 model construction. The model favours those genes that mainly regulate the target
75 class during the training phase of the model. Classification and regression tree is
76 one of the examples of embedded methods⁸.

77 In association with classification problems, gene selection concentrate on the
78 selection of the most informative and regulatory genes. In this connection many
79 gene selection methods have been employed. Among them is random forest (RF)

80 to solve the two issues of variable selection, i.e., to choose the most important
81 attributes and try to construct the best parsimonious predictive model⁹. Xu et al,
82 suggested an improved RF method, which exploits a new feature weighting tool for
83 selecting a gene subspace and hence boost classification accuracy on microarray
84 data¹⁰. Vladimir et al, used random forest as a classification and regression
85 techniques for compound classification and quantitative structure-activity
86 relationship (QSAR) modeling where they considered the technique for categorical
87 biological activities¹¹. They run models for six data sets and presented three
88 additional features of random forest. They claimed that random forest is one of the
89 most precise and dominant tool of delivering best performance. Diaz-Uriarte et al,
90 used random forest for classification of microarray data and proposed a new
91 method of feature selection based on random forest¹². Tran et al, proposed the idea
92 of combining bagging and feature selection¹³. In the selection of relevant gene
93 subset for bagging, a wrapper based feature selection method is employed.
94 Besides bagging and random forest, several other methods have also been
95 employed by various researchers for gene/feature selection. Apiletti et al, proposed
96 a method based on filtering gene selection approach where at first stage they detect
97 outliers in gene expression data for each gene¹⁴. Several gene selection techniques
98 evaluate the importance of genes in distinguishing the tissue samples in a given
99 target class by deciding a cut point or by fitting a statistical model to microarray
100 gene expression data¹⁵.

101 This idea was exploited to an expression range to build a gene mask¹⁴. The idea of
102 set covering approach was used to minimize gene subset and select the most
103 informative genes in the analysis of tumor and normal colon tissues probed by
104 oligonucleotide arrays¹⁶. The minimum gene subset was selected by replacing the
105 set covering approach with the greedy search approach¹⁴. To cope with the issue of
106 dimensionality and outliers in genes selection, the greedy search approach was

107 considered together with proportional overlapping analysis for classifying tissue
108 sample into binary classes¹⁷.

109 Most of the tree based feature selection methods discussed above use the idea given
110 in random forest to select genes, i.e., a random sample of genes is chosen to select
111 the splitting gene at each node of the tree⁹. In situations where the number of
112 samples is small compared to the number of genes, the random forest idea might
113 not give satisfactory results, as some of the genes might not get a chance to be used
114 as a splitting variable. Therefore, this study used bagged tree forest for feature
115 selection where all the features are assessed to decide on the best possible split.
116 This ensures that every gene plays its role in the construction of the tree model,
117 thus reducing the chance of missing out important genes.

118 **Datasets and Methods:**

119 A total of 5 data sets have been used to compare the methods by calculating
120 classification accuracies, Brier scores and sensitivities. A brief description of each
121 of these datasets is given as follows:

122 **Colon:** Colon is a type of cancer which is also known as colorectal cancer. This
123 cancer commonly initiates when strong tissues in the line of rectum change and
124 develop abnormally, resulting in a mass called tumor. The dataset comprised of
125 2000 genes with 62 colon tissues, out of which 22 tissues were normal and 42 were
126 cancerous. Colon dataset given by Alon et al, and Ben-Dar et al, was utilized for
127 binary tissue classification^{16,18}. Since then, this dataset is progressively utilized and
128 analyzed by many researchers.

129 **Breast Cancer:** Breast cancer is a type of cancer occurring in the breast cells. This
130 cancer is common in many parts of the world¹⁹. This dataset comprised of
131 observations on 4869 genes from 77 tissue samples, of which 33 were
132 noncancerous and 44 cancerous. The data matrix is a 77×4869 binary class
133 problem. This data was taken from the study conducted by Michiels et al.²⁰

134 **Lymph node breast cancer:** This data consisted of 144 lymph node breast cancer
 135 patients. Gene expression measurements of 70 genes were signals for metastasis-
 136 free survival. A class variable was represented by “event”. This data was used by
 137 Marc et al, to discover the most reliable means of signals in breast cancer that help
 138 in selection of patients for systemic treatment²¹.

139 **Leukemia:** Leukemia is a group of cancerous blood forming cells. It usually starts
 140 in the bone marrow of human body. This cancer occurs due to a mass production
 141 consisting of abnormal white blood cells, which fight against infection and
 142 toxicities²². The leukemia data was published by Golub et al,²³ consisting of
 143 observations on 72 tissue samples and 7129 genes.

144 **Serrated colorectal carcinomas:** Serrated colorectal carcinomas (CRCs), that are
 145 morphologically different from usual CRCs, have been proposed to follow a
 146 unique pathway of CRC formation. The dataset has been taken from a study to
 147 examine the gene expression profiling of 37 serrated CRCs against conventional
 148 CRCs, and to identify differentially expressed genes representing potential
 149 biomarkers for serrated CRC²⁴. Observations on a total of 22215 genes from the
 150 tissue samples have been made.

151 Table 1 gives a brief summary of the datasets showing the number of tissue
 152 samples, number of genes and classwise distribution of tissue samples as
 153 noncancerous and cancerous in each dataset.

154 Gene expression data are commonly given in the form of an expression matrix,
 155 $X = [x_{ji}]$, such that $X \in \mathbb{R}^{n \times d}$ and $[x_{ji}]$ is the expression value of gene j for the
 156 i th observation where $j = 1, \dots, d$ and $i = 1, \dots, n$. Each tissue sample also has a
 157 response class label, y_i , that showed the phenotype of the observation (tissue
 158 sample) being observed. Let $Y \in \mathbb{R}^n$ be the set of class labels given that its
 159 element, y_i , has a unique value c which is either 1 or 0.

160 The method exploited the idea of bagging and variable importance using random
161 forest mechanism to select the relevant genes. Bagged classification trees used to
162 ensure that all genes are given the chance to contribute in tree construction. At
163 each node all the “d” features were considered for choosing the best split. On the
164 other hand, trees in random forest were grown on bootstrap samples by considering
165 a random set of $p < d$ features for node splitting. A large number of bagged
166 forests, each consisting of a small number of trees, were grown on bootstrap
167 samples from the training part and the most accurate forests were selected based on
168 out-of-bag error estimates. The selected forests are combined together and used to
169 rank genes in the random forest style. Gene scoring was done as follow:

170 In each tree of the forest, a binary split was made on the gene that gives the best
171 possible partition of the tissue samples. This was done by computing an impurity
172 measure (Gini index here) of class distribution in the original sample set and the
173 sets formed due to the binary split. Gene that gives the least value of the Gini index
174 is chosen for splitting the tissue samples. This process was iterated recursively on
175 each consequent partition until there remains a single tissue sample in the resulting
176 node. Each time a partition of a node is made on a gene, the Gini impurity measure
177 for the two resulting nodes is less than the parent node. Gini decreases for each
178 individual gene were added over all the trees in the forest. Genes with the highest
179 added decrease were selected as the final set of genes.

180 Let M be the total number of forests each consisting of B trees. After estimating
181 the out-of-bag errors of all the M forests, top L forests with the smallest errors were
182 selected and combined to form a bagged tree ensemble. The final ensemble thus
183 consist of $T = L \times B$ trees. To calculate the importance of a gene x_i , i.e., $Imp(x_i)$,
184 for predicting Y , weighted impurity decreases $p(t)\Delta i(s_t, t)$ for all nodes t where
185 x_i is used as the split variable, were added and averaged over all the T trees in the
186 final forest, i.e.,

$$187 \quad Imp(x_i) = \frac{1}{T} \sum_T \sum_{\forall t: v(s_t=x_i)} p(t) \Delta i(s_t, t),$$

188 where $p(t)$ is the proportion T_t/n of tissues reaching t and $v(s_t)$ is the gene used
 189 in split s_t . $\Delta i(s_t, t)$ is the maximum decrease in impurity at node t for the split s_t
 190 that divides the n_t node samples into t_L and t_R and is given by

$$191 \quad \Delta i(s_t, t) = i(t) - p_L i(t_L) - p_R i(t_R),$$

192 where $i(t)$ is the Gini impurity measure, $p_L = n_{t_L}/n_t$ and $p_R = n_{t_R}/n_t$.

193 The proposed BSF method took the following steps for genes selection.

- 194 1. Grow a large number M of bagged tree forests each consisting of a small
 195 number of trees, say B and rank them with respect to their out-of-bag error.
- 196 2. Select a certain number of the top ranked forests, say L .
- 197 3. Combine the top ranked L forests to form a final ensemble of $L \times B$ bagged
 198 trees.
- 199 4. Rank genes using the selected forest by calculating the Gini score of each
 200 gene.

201 To evaluate different gene selection approaches, it is necessary to check the
 202 accuracy of a classifier applied after the gene selection procedure, where
 203 classification is done only on the selected genes. By this evaluation, one could
 204 check the ability of genes that regulate the tissue target class. Authors have used
 205 different gene selection methods and have shown that gene selection methods have
 206 significant effect on a classifier accuracy^{25,26}. The same approach has been utilized
 207 in another study¹⁴.

208 Random forest and support vector machine approaches were utilized to assess the
 209 predictive power based on the selected genes in comparison with three other state-
 210 of-the-art methods, i.e., Wilcoxon rank sum test¹, masked painter and proportional
 211 overlapping analysis. A brief description of the classifiers is given as follows:

212 **Support vector machine (SVM)**

213 Support vector machine (SVM) is one of the most commonly used classifiers²⁷.
214 The simplest kind of SVM is the linear SVM classifier. In linear classification, a
215 separating plane is said to be the best if it has the maximum margin from both the
216 classes. The margin line is the space between two equivalent hyper planes, each of
217 which goes into the support vectors of one class.

218 **Random Forest (RF)**

219 Random forest is an ensemble learning procedure for classification and regression
220 problems that consisted of many decision trees, where each tree is grown on
221 bootstrap sample from the training data⁹. A new tissue sample is classified on the
222 basis of majority voting from the decision trees in the forest.

223 Several other classification methods could be considered in addition to random
224 forest and support vector machine to evaluate the performance of the genes
225 selection methods^{28,29}.

226 The given datasets were divided into 70% training part (for feature selection and
227 model fitting) and 30% testing (for performance evaluation) part. In the first phase,
228 feature selection methods were applied on the training part. Top 20 genes were
229 selected by all the gene selection methods. In the second phase, the two classifiers
230 were applied on the training part of each dataset with selected set of genes and then
231 the required metrics were calculated on the testing parts. This process was iterated
232 500 times and final results were the average from all the combination of the runs.
233 For all the analysis, R programming language has been used. For the SVM and RF
234 classifiers, '*kernlab*' and '*randomForest*', R packages were used, respectively. The
235 default values of the parameters as given in the corresponding packages have
236 been used. For random forest, number of trees have been fixed at 500, node size
237 was set to 1 and the number of genes selected randomly at each node of the tree
238 was the square root of the total number of genes. In the case of SVM, the linear

239 kernel has been used along with the default automatic selection for the alpha
240 parameter.

241 **Performance Evaluation**

242 To assess the predictive performance of the classifiers based on the selected set of
243 genes identified by the gene selection methods, classification accuracy, sensitivity
244 and Brier score have be used as performance measures. These measures are
245 explained as follows:

246 **Classification accuracy:** Classification accuracy was obtained by dividing the
247 number of correct classifications on the total number of tissue samples in the test
248 data. This measure could be easily obtained from a matrix formed by cross
249 tabulating true tissue status vs prediction made by a model. This matrix was known
250 as confusion matrix and is given below.

251

Prediction	True Status	
	Postive	Negative
Postive	TP (n_{11})	FP (n_{12})
Negative	FN (n_{21})	TN (n_{22})
	<i>Sensitivity</i> $= \frac{n_{11}}{n_{11} + n_{21}}$	<i>Specificity</i> $= \frac{n_{22}}{n_{12} + n_{22}}$

252

253 In the above table, TP indicated the number of positive cases classified as positive
254 in the test data; FP was the count of negative cases labelled as positive; FN was the
255 count of positive cases labelled as negative and TN was the count of negative cases
256 labelled as negative. Based on this, classification accuracy was given as

257

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

258 *Sensitivity*: Sensitivity, also called true positive rate (TPR) is the proportion of
 259 positive instances that are accurately classified as positive. It is also called recall
 260 and hit rate and is given by

$$261 \quad \text{Sensitivity} = \frac{TP}{TP + FN}$$

262
 263 *Brier Score*: Brier score is a scoring criteria that measures the accuracy of
 264 probabilistic prediction³⁰. A probabilistic prediction indicated an exact
 265 outcome/event. This score is commonly used for binary class problems. For true
 266 predictions, Brier score would be closed to 0 and close to 1 otherwise.

267 Brier score could be computed as follows

$$268 \quad BS = \frac{\sum_{i=1}^{\# \text{ of test point}} (y_i - \hat{p}_i)^2}{\# \text{ of tissue samples in test data}}$$

269 where,

270 y_i = A particular tissue class value in 0, 1 form;

271 \hat{p}_i = Estimated class probability

272

273 **Results**

274 Tables 2-6 gives the results of all the gene selection methods i.e. the proposed BSF
 275 method, Wilcoxon, proportional overlapping score (POS) method and masked
 276 painter (MP), on the 5 datasets via random forest (RF) and support vector machine
 277 (SVM). In the given tables, the first column showed the classifiers, the second was
 278 the performance measure used against each of the classifiers and the subsequent
 279 columns were the values obtained for the measures for each gene selection method.
 280 The result of the best performing method is shown in bold.

281 Table 2 gives the results of the methods for colon data, Table 3 demonstrates the
 282 results for breast cancer data and Table 4 gives the results of all the methods for

283 nki dataset. Table 5 and Table 6 give the results of the gene selection methods on
284 GSE4045 and leukemia datasets, respectively.

285 Top 20 genes were selected by each method and the rest were discarded.
286 Classification accuracy, sensitivity and Brier score (to know the degree of belief in
287 classification) were calculated on tissue samples in the test set. This process has
288 been repeated 500 times on different random partitions of the datasets. The BSF
289 method gave better result than the other competitors. The proposed BSF method
290 achieved the smallest Brier score values in most of the cases. Further discussion on
291 individual dataset is as follows:

292 For colon dataset, the random forest and SVM classifiers gave the highest accuracy
293 on genes selected via the proposed BSF method as compared to the Wilcoxon,
294 POS and Masked Painter methods. Accuracy of the proposed method by using
295 random forest classifier was 0.8508, whereas accuracies of the competitors, i.e.,
296 Wilcoxon, POS and Masked Painter methods were 0.8495, 0.8391 and 0.7371,
297 respectively. Likewise, random forest attained the highest sensitivity for the
298 proposed method, which is 0.6698. The sensitivity for Wilcoxon, POS and masked
299 painter by using random forest classifier were 0.6479, 0.6345 and 0.5341,
300 respectively. Similar conclusion was drawn from the SVM classifier results. In
301 terms of Brier score, the proposed method has also outperformed the other three
302 methods via both the classifiers.

303 Likewise Table 3 reflected the highest accuracy and sensitivity for breast cancer of
304 the proposed BSF method via random forest classifier i.e. 0.7931 and 0.8372
305 respectively, with minimum Brier score of 0.1453. Moreover on SVM classifier
306 the proposed BSF method outperformed all the competitors in terms of Brier score
307 and accuracy, while, Wilcoxon excelled other methods, with regard to sensitivity.

308 For lymph node breast cancer dataset the proposed method has achieved the best
309 results in terms of classification accuracy, sensitivity and Brier score among all the
310 other methods, on both the classifiers, as can be seen in Table 4.

311 The results on the cerated colerectal carcinoma dataset were tabulated in Table 5.
312 Substantial difference could be noted in terms of accuracy, sensitivity and Brier
313 score of the proposed method via both the classifiers. The Accuracy and
314 sensitivity of the proposed method via random forest classifier are 0.9177 and
315 0.6842, respectively, which is much higher than the other methods. The Accuracy
316 of Wilcoxon, POS and masked painter were 0.7847, 0.8025 and 0.7331,
317 respectively, via random forest classifier. On random forest classifier the Brier
318 score of the proposed method was less than all the other methods. The SVM
319 classifier too yielded the highest accuracy, sensitivity and the smallest Brier score
320 for the proposed method.

321 Likewise, it could be noted from Table 6, that the accuracy, sensitivity, and Brier
322 score of the proposed method via the random forest and SVM classifiers were
323 comparatively better than all the other methods on the leukemia dataset. Using
324 random forest classifier the accuracy and sensitivity of the proposed BSF method
325 was 0.9214 and 0.8753, surpassing all the other methods. The accuracy of the
326 competitors, i.e., Wilcoxon, POS and masked painter via random forest is 0.829,
327 0.88871 and 0.7536, respectively. The proposed method has the smallest Brier
328 score as compared to all the other methods via random forest classifier. Similarly
329 SVM classifier has given the highest accuracy, sensitivity and the smallest Brier
330 score for the proposed method, which is 0.7552, 0.7661 and 0.0332, respectively.

331

332 **Discussion**

333 This work has proposed a gene selection method using gene ranking via Gini score
334 method in conjunction with bagged tree ensemble. The proposed method is based

335 on a greedy search approach that selects the best bagged tree forests from a large
336 pool of forests based on their out-of-bag error. The proposed method achieves
337 improved gene selection in two folds, i.e. selecting the best tree forests while
338 discarding those that do not perform well in the training phase, and considering all
339 the genes for finding the best splitting variable at each node of the trees in the
340 forest. On the other hand, the standard random forest⁹ considers node splitting on a
341 randomly selected subset of genes/features and might have the likelihood of
342 ignoring important features/genes in the model building process. The proposed
343 method has outperformed all the methods on almost all the datasets considered in
344 this paper. In addition to classification accuracy and sensitivity, Brier scores³⁰ has
345 also been used as performance measure to know the degree of belief in classifying
346 observations to their correct target classes based on the selected set of genes. The
347 proposed method has outperformed all the state-of-the-art methods on all the
348 datasets given in this paper in terms of Brier score. This suggests that the proposed
349 method could be used for gene selection to classify observations into their correct
350 target classes with higher degree of belief as compared to the other methods, i.e.
351 POS¹⁷, Wilcoxon rank sum test¹ and MP¹⁴. Genes selected via the proposed
352 method also give higher values of sensitivity for majority of the datasets considered
353 using random forest⁹ and support vector machine²⁷ classification algorithms. This
354 means that the proposed method could be more effective in avoiding false positives
355 as compared to the rest of the methods. Moreover, as the proposed method selects
356 the best performing sub-forests and discards those with poor performance, this
357 reduces the size of the final ensemble which in turn saves computational costs in
358 terms of storage resources.

359

360

361

362 **Conclusion and Future Work**

363 A novel gene selection method has been proposed in this paper that is based on an
364 ensemble of the most accurate forests chosen from a large pool of small size
365 forests grown by the method of bagging. The method has been compared with
366 other state-of-the-art methods used in literature for gene selection on a total of 5
367 gene expression datasets. The analyses have revealed that the proposed BSF
368 method has outperformed the other methods in almost all the cases. This means
369 that the proposed method could effectively identify those genes that have the
370 highest discriminative ability to classify a given tissue sample to its correct target
371 class. The main limitation of the method is computational complexity. This could
372 be avoided by using parallel computing.

373

374 **Disclaimer:** None

375 **Conflict of Interest:** None

376 **Funding Sources:** None

377

378 **References**

- 379 1. Saeys Y, Inza I, Larrañaga P. A review of feature selection techniques in
380 bioinformatics. *Bioinformatics*. 2007;23:2507-17.
- 381 2. Mahmood MS, Kureshi N, Frossard PM. Gene markers and complex
382 disorders: a review. *J Pak Med Assoc*. 2004;54:584-9.
- 383 3. Khan Z, Naeem M, Khalil U, Khan DM, Aldahmani S, Hamraz M. Feature
384 Selection for Binary Classification Within Functional Genomics
385 Experiments via Interquartile Range and Clustering. *IEEE Access*. 2019 Jun
386 12;7:78159-69.
- 387 4. Witten DM, Tibshirani R. A framework for feature selection in clustering. *J*
388 *Am Stat Assoc*. 2010;105:713-26.

- 389 5. Mahmood MS, Kureshi N, Frossard PM. Gene markers and complex
390 disorders: a review. *J Pak Med Assoc.* 2004;54:584-9.
- 391 6. Khalid M, Khan S, Ahmad J, Shaheryar M. Multivariate Covariance using
392 Principal Component Analysis for Reconstruction of Bidirected Gene
393 Regulatory Networks. In 2017 International Conference on Frontiers of
394 Information Technology (FIT) 2017 Dec 18 (pp. 229-234). IEEE.
- 395 7. Khalid M, Khan S, Ahmad J, Shaheryar M. Identification of self-regulatory
396 network motifs in reverse engineering gene regulatory networks using
397 microarray gene expression data. *IET Systems Biology.* 2018 Dec
398 13;13(2):55-68.
- 399 8. Breiman L. *Classification and regression trees.* Routledge; 2017.
- 400 9. Breiman L. Random forests. *Mach Learn.* 2001;45:5-32.
- 401 10. Xu B, Huang JZ, Williams G, Wang Q, Ye Y. Classifying very high-
402 dimensional data with random forests built from small subspaces. *Int J Data*
403 *Warehous.* 2012;8:44-63.
- 404 11. Svetnik V, Liaw A, Tong C, Culberson JC, Sheridan RP, Feuston BP.
405 Random forest: a classification and regression tool for compound
406 classification and QSAR modeling. *J Chem Inf Comput Sci.* 2003;43:1947-
407 58.
- 408 12. Díaz-Uriarte R, De Andres SA. Gene selection and classification of
409 microarray data using random forest. *BMC Bioinformatics.* 2006;3.
- 410 13. Tran CT, Zhang M, Andrae P, Xue B. Bagging and Feature Selection for
411 Classification with Incomplete Data. In *Europ Conf Ap Evol Comp 2017;*
412 471-486.
- 413 14. Apiletti D, Baralis E, Bruno G, Fiori A. Maskedpainter: feature selection for
414 microarray data analysis. *Intel Data Anal.* 2012;16:717-37.

- 415 15. Marczyk M, Jaksik R, Polanski A, Polanska J. Adaptive filtering of
416 microarray gene expression data based on gaussian mixture decomposition.
417 BMC Bioinformatics. 2013;14:101.
- 418 16. Alon U, Barkai N, Notterman DA, Gish K, Ybarra S, Mack D, Levine AJ.
419 Broad patterns of gene expression revealed by clustering analysis of tumor
420 and normal colon tissues probed by oligonucleotide arrays. Proc Natl Acad
421 Sci. 1999;96:6745-50.
- 422 17. Mahmoud O, Harrison A, Perperoglou A, Gul A, Khan Z, Metodiev MV,
423 Lausen B. A feature selection method for classification within functional
424 genomics experiments based on the proportional overlapping score. BMC
425 Bioinformatics. 2014;15:274.
- 426 18. Ben-Dor A, Bruhn L, Friedman N, Nachman I, Schummer M, Yakhini Z.
427 Tissue classification with gene expression profiles. J Comput Biol.
428 2000;7:559-83.
- 429 19. Mujtaba S, Haroon S, Faridi N, Lodhi FR. Correlation of human epidermal
430 growth factor receptor 2 (HER-2/neu) receptor status with hormone
431 receptors Oestrogen Receptor, Progesterone Receptor status and other
432 prognostic markers in breast cancer: an experience at tertiary care hospital in
433 Karachi. J Pak Med Assoc. 2013;63:854-8.
- 434 20. Michiels S, Koscielny S, Hill C. Prediction of cancer outcome with
435 microarrays: a multiple random validation strategy. Lancet. 2005;365:488-
436 92.
- 437 21. Reich M, Liefeld T, Gould J, Lerner J, Tamayo P, Mesirov JP. Gene Pattern
438 2.0. Na Genet. 2006;38:500.
- 439 22. Khan A, Shafiq I, Shah MH, Khan S, Shahid G, Arabdin M. Chronic
440 myeloid leukaemia presenting as priapism: A case report from Khyber
441 Pakhtunkhwa. J Pak Med Assoc. 2018;68:942-4.

- 442 23. Golub TR, Slonim DK, Tamayo P, Huard C, Gaasenbeek M, Mesirov JP,
443 Coller H, Loh ML, Downing JR, Caligiuri MA, Bloomfield CD. Molecular
444 classification of cancer: class discovery and class prediction by gene
445 expression monitoring. *Science*. 1999;286:531-7.
- 446 24. Laiho P, Kokko A, Vanharanta S, Salovaara R, Sammalkorpi H, Järvinen H,
447 Mecklin JP, Karttunen TJ, Tuppurainen K, Davalos V, Schwartz Jr S.
448 Serrated carcinomas form a subclass of colorectal cancer with distinct
449 molecular basis. *Oncogene*. 2007;26:312.
- 450 25. Jirapech-Umpai T, Aitken S. Feature selection and classification for
451 microarray data analysis: Evolutionary methods for identifying predictive
452 genes. *BMC Bioinformatics*. 2005;6:148.
- 453 26. Mahmoud O, Harrison A, Gul A, Khan Z, Metodiev MV, Lausen B.
454 Minimizing redundancy among genes selected based on the overlapping
455 analysis. In *Analysis of Large and Complex Data 2016* (pp. 275-285).
456 Springer, Cham.
- 457 27. Vapnik V, Golowich SE, Smola AJ. Support vector method for function
458 approximation, regression estimation and signal processing. In *Advances in
459 neural information processing systems 1997* (pp. 281-287).
- 460 28. Gul A, Perperoglou A, Khan Z, Mahmoud O, Miftahuddin M, Adler W,
461 Lausen B. Ensemble of a subset of kNN classifiers. *Advances in data
462 analysis and classification*. 2018 Dec 1;12(4):827-40.
- 463 29. Khan Z, Gul A, Perperoglou A, Miftahuddin M, Mahmoud O, Adler W,
464 Lausen B. Ensemble of optimal trees, random forest and random projection
465 ensemble classification. *Advances in Data Analysis and Classification*. 2019
466 Jun 2:1-20.
- 467 30. Brier GW. Verification of forecasts expressed in terms of probability. *Mon
468 Weather Rev*. 1950;78:1-3.

469

470

471

472 **Table 1: Summary of the datasets.**

473

474

475

476

477

478

479

480

481

482 **Table 2: Performance of the methods on colon dataset.**

Classifier		BSF	Methods		
			Wilcoxon	POS	MP
RF	Sensitivity	0.6698	0.6479	0.6345	0.5341
	BS	0.1431	0.1575	0.1457	0.1511
	Accuracy	0.8508	0.8495	0.8391	0.7371
SVM	Sensitivity	0.4211	0.2396	0.4474	0.4131
	BS	0.1691	0.1924	0.1694	0.1832
	Accuracy	0.7554	0.7109	0.7542	0.6632

483

484

485

486 **Table 3: Performance of the methods on breast cancer dataset.**

Classifier		BSF	Methods		
			Wilcoxon	POS	MP
RF	Sensitivity	0.8372	0.8298	0.7746	0.6351
	BS	0.1453	0.1828	0.2069	0.2113
	Accuracy	0.7931	0.7482	0.6309	0.6310
SVM	Sensitivity	0.7798	0.8214	0.7876	0.6115
	BS	0.1521	0.1592	0.2377	0.2641
	Accuracy	0.7807	0.7759	0.6414	0.6192

487
488
489
490

Table 4: Performance of the methods on lymph node breast cancer dataset.

Classifier		Methods			
		BSF	Wilcoxon	POS	MP
RF	Sensitivity	0.8124	0.8030	0.7599	0.7991
	BS	0.0996	0.1095	0.1149	0.2013
	Accuracy	0.8754	0.8630	0.862	0.7539
SVM	Sensitivity	0.7697	0.7475	0.6227	0.6422
	BS	0.1194	0.1329	0.15671	0.2110
	Accuracy	0.8536	0.8515	0.8113	0.7566

491
492
493
494
495

Table 5: Performance of the methods on serrated colorectal carcinoma dataset.

Classifier		Methods			
		BSF	Wilcoxon	POS	MP
RF	Sensitivity	0.6842	0.2997	0.2043	0.2110
	BS	0.0720	0.1603	0.1611	0.1562
	Accuracy	0.9177	0.7847	0.8025	0.7331
SVM	Sensitivity	0.3100	0.0069	0.0124	0.2741
	BS	0.0241	0.1730	0.1799	.01526
	Accuracy	0.8215	0.7828	0.7848	0.6317

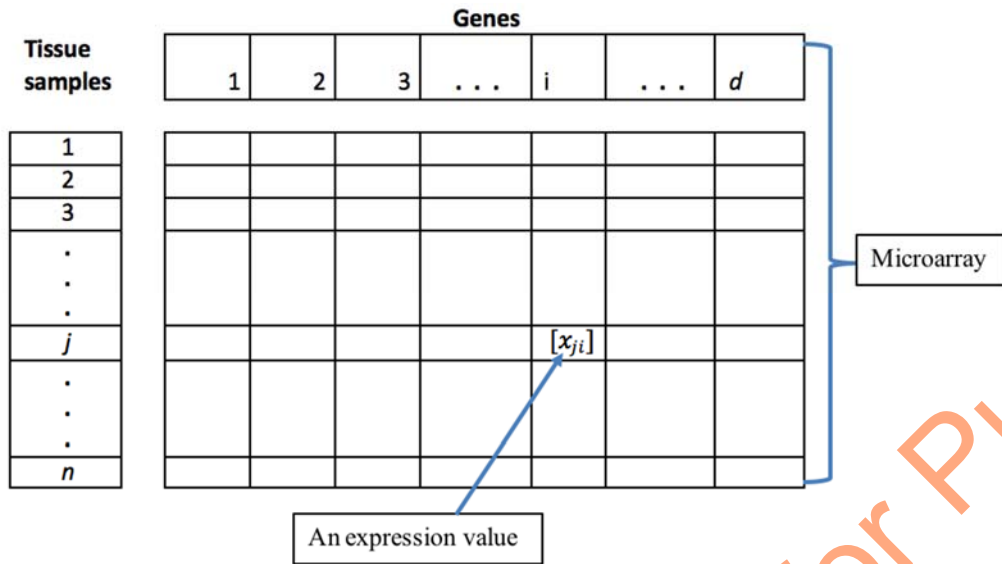
496
497
498
499

Table 6: Performance of the methods on leukaemia dataset.

Classifier		Methods			
		BSF	Wilcoxon	POS	MP
RF	Sensitivity	0.8753	0.7701	0.7831	0.6119
	BS	0.1103	0.2110	0.203	0.2142
	Accuracy	0.9214	0.8391	0.8871	0.7536
SVM	Sensitivity	0.7661	0.7335	0.6552	0.6112
	BS	0.0332	0.1130	0.2411	0.1142
	Accuracy	0.7552	0.7315	0.7112	0.6351

500
501

502



503

504 **Figure 1: Microarray gene expression dataset.**

505

506

507 **A flowchart of the method is given in Figure 1.**

508

509

510

511

512

513

514

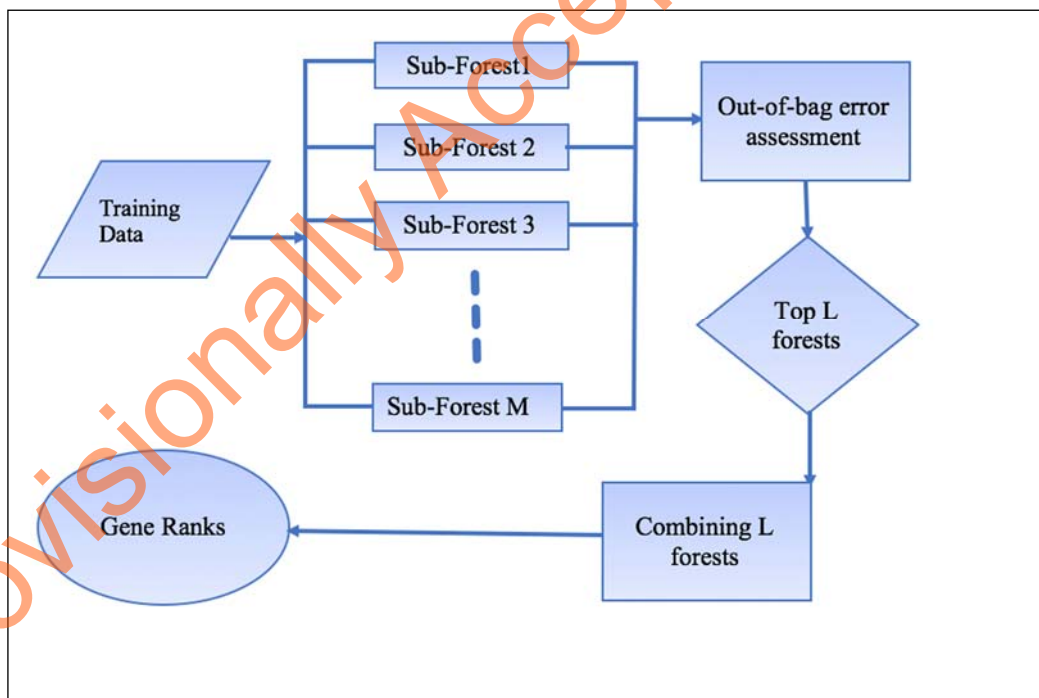
515

516

517

518

519

519 **Figure 2: Flow chart of the proposed method.**