

A comprehensive simulation study on classification of RNA-Seq data

Gokmen Zararsiz ^{Corresp.} ¹, Dinçer Göksülük ², Selçuk Korkmaz ², Vahap Eldem ³, Gözde Ertürk Zararsız ¹, İzzet Parug Duru ⁴, Turgay Unver ⁵, Ahmet Öztürk ¹

¹ Biostatistics, Erciyes University, Faculty of Medicine, Kayseri, TURKEY

² Department of Biostatistics, Hacettepe University, Ankara, Turkey

³ Department of Biology, Istanbul University, Istanbul, Turkey

⁴ Department of Physics, Marmara University Istanbul, Istanbul, Turkey

⁵ Genomics, Izmir International Biomedicine and Genome Institute, Izmir, Turkey

Corresponding Author: Gokmen Zararsiz

Email address: gokmenzararsiz@hotmail.com

Background RNA sequencing (RNA-Seq) is a powerful technique for transcriptome profiling of the organisms that uses the capabilities of next-generation sequencing (NGS) technologies. Recent advances in NGS let to measure the expression levels of tens to thousands of transcripts simultaneously. Using such information, developing expression-based classification algorithms is an emerging powerful method for diagnosis, disease classification and monitoring at molecular level, as well as providing potential markers of disease. Microarray based classifiers cannot be directly applied due to the discrete nature of RNA-Seq data. One way is to develop count-based classifiers, such as poisson linear discriminant analysis (PLDA) and negative binomial linear discriminant analysis (NBLDA). Other way is to transform the data hierarchically closer to microarrays and apply microarray-based classifiers. In most of the studies, the data overdispersion seems to be an another challenge in modeling RNA-Seq data. In this study, we aimed to examine the effect of dispersion parameter and classification algorithms on RNA-Seq classification. We also considered the effect of other parameters (i) sample size, (ii) number of genes, (iii) number of class, (iv) DE (differential expression) rate, (v) transformation method on classification performance.

Methods We designed a comprehensive simulation study, also used two miRNA and two mRNA experimental datasets. Simulated datasets are generated from negative binomial distribution under different scenarios and real datasets are obtained from publicly available resources. We compared the results of several classifiers including PLDA with and without power transformation, NBLDA, single SVM, bagging SVM (bagSVM), classification and regression trees (CART), and random forests (RF).

Results Results from the simulated and real datasets revealed that increasing the sample size, differential expression rate, number of genes and decreasing the dispersion parameter and number of groups lead to an increase in the classification accuracy. To make an overall assessment, power transformed PLDA, RF and SVM classifiers performed the highest classification accuracies.

Discussion Overdispersion seems to be an important challenge in RNA-Seq classification studies. Similar with differential expression studies, classification of RNA-Seq data requires careful attention on handling data overdispersion. We conclude that, as a count-based classifier, power transformed PLDA; as a microarray based classifier vst or rlog transformed RF and SVM (bagSVM for high sample sized data) classifiers may be a good choice for classification. However, there is still a need to develop novel classifiers or transformation approaches for classification of RNA-Seq data. An **R/BIOCONDUCTOR** package MLSeq with a vignette is freely available at

<http://www.bioconductor.org/packages/2.14/bioc/html/MLSeq.html> .

1 **A Comprehensive Simulation Study on Classification of**

2 **RNA-Seq Data**

3 Gokmen Zararsiz¹, Dincer Goksuluk², Selcuk Korkmaz², Vahap Eldem³,

4 Gozde Erturk Zararsiz¹, Izzet Parug Duru⁴, Turgay Unver⁵, Ahmet Ozturk¹

5

6 ¹Department of Biostatistics, Erciyes University, Kayseri, Turkey

7 ²Department of Biostatistics, Hacettepe University, Ankara, Turkey

8 ³Department of Biology, Istanbul University, Istanbul, Turkey

9 ⁴Department of Physics, Marmara University, Istanbul, Turkey

10 ⁵Izmir International Biomedicine and Genome Institute, İzmir, Turkey

11

12 Corresponding Author:

13 Gokmen Zararsiz

14 Erciyes University Biostatistics, Kayseri, 38039, Turkey

15 Email address: gokmenzararsiz@erciyes.edu.tr, gokmenzararsiz@hotmail.com

16 Abstract**17 Background**

18 RNA sequencing (RNA-Seq) is a powerful technique for transcriptome profiling of the
19 organisms that uses the capabilities of next-generation sequencing (NGS) technologies. Recent
20 advances in NGS let to measure the expression levels of tens to thousands of transcripts
21 simultaneously. Using such information, developing expression-based classification algorithms
22 is an emerging powerful method for diagnosis, disease classification and monitoring at molecular
23 level, as well as providing potential markers of disease. Microarray based classifiers cannot be
24 directly applied due to the discrete nature of RNA-Seq data. One way is to develop count-based
25 classifiers, such as poisson linear discriminant analysis (PLDA) and negative binomial linear
26 discriminant analysis (NBLDA). Other way is to transform the data hierarchically closer to
27 microarrays and apply microarray-based classifiers. In most of the studies, the data
28 overdispersion seems to be an another challenge in modeling RNA-Seq data. In this study, we
29 aimed to examine the effect of dispersion parameter and classification algorithms on RNA-Seq
30 classification. We also considered the effect of other parameters (i) sample size, (ii) number of
31 genes, (iii) number of class, (iv) DE (differential expression) rate, (v) transformation method on
32 classification performance.

33 Methods

34 We designed a comprehensive simulation study, also used two miRNA and two mRNA
35 experimental datasets. Simulated datasets are generated from negative binomial distribution
36 under different scenarios and real datasets are obtained from publicly available resources. Data
37 normalization is applied using deseq median ratio approach. A variance stabilizing
38 transformation (vst) and regularized logarithmic transformation (rlog) methods are used before

39 applying microarray-based classifiers. We compared the results of several classifiers including
40 PLDA with and without power transformation, NBLDA, single SVM, bagging SVM (bagSVM),
41 classification and regression trees (CART), and random forests (RF).

42 **Results**

43 Results from the simulated and real datasets revealed that increasing the sample size, differential
44 expression rate, number of genes and decreasing the dispersion parameter and number of groups
45 lead to an increase in the classification accuracy. To make an overall assessment, power
46 transformed PLDA, RF and SVM classifiers performed the highest classification accuracies.

47 **Discussion**

48 Overdispersion seems to be an important challenge in RNA-Seq classification studies. Similar
49 with differential expression studies, classification of RNA-Seq data requires careful attention on
50 handling data overdispersion. We conclude that, as a count-based classifier, power transformed
51 PLDA; as a microarray based classifier vst or rlog transformed RF and SVM (bagSVM for high
52 sample sized data) classifiers may be a good choice for classification. However, there is still a
53 need to develop novel classifiers or transformation approaches for classification of RNA-Seq
54 data. An **R/BIOCONDUCTOR** package MLSeq with a vignette is freely available at
55 <http://www.bioconductor.org/packages/2.14/bioc/html/MLSeq.html>.

56 Introduction

57 With the advent of high-throughput NGS technologies, transcriptome sequencing (RNA-Seq) has
58 become one of the central experimental approaches for generating a comprehensive catalog of
59 protein-coding genes and non-coding RNAs and examining the transcriptional activity of
60 genomes. Furthermore, RNA-Seq has already proved itself to be a promising tool with a
61 remarkably diverse range of applications; (i) discovering novel transcripts, (ii) detection and
62 quantification of spliced isoforms, (iii) fusion detection, (iv) reveal sequence variations (e.g.,
63 SNPs, indels) (Wang, Gerstein & Snyder, 2009). Additionally, beyond these general applications,
64 RNA-Seq holds great promise for gene expression-based classification to identify the significant
65 transcripts, distinguish biological samples and predict clinical or other outcomes due to large
66 amounts of data, which can be generated in a single run. This classification is widely used in
67 medicine for diagnostic purpose and refers to the detection of small subset of genes that achieves
68 the maximal predictive performance. These genes are used afterwards for classification of new
69 observations into the disease classes (or tumor classes, cancer subtypes, cancer stage, etc.).

70 Although microarray-based gene expression classification have become very popular during last
71 decades, more recently, RNA-Seq replaced microarrays as the technology of choice in
72 quantifying gene expression due to some advantages as providing less noisy data, detecting novel
73 transcripts and isoforms, and unnecessary of prearranged transcripts of interest (Furey et al.,
74 2000; Zhu & Hastie, 2004; Uriarte & de Andres, 2006; Rapaport et al., 2007). However, to
75 measure gene expression, microarray technology provides continuous data, while RNA-Seq
76 technology generates discrete count data, which corresponds to the abundance of mRNA
77 transcripts (Witten, 2011). Another issue is the overdispersion problem, where the variance
78 exceeds the mean (Nagalakshmi et al., 2008). Various studies have been employed to deal with

79 the overdispersion problem for differential expression (DE) analysis of RNA-Seq data (Anders &
80 Huber, 2010; Robinson, McCarthy & Smyth, 2010; Di et al., 2011; Sonesson & Delorenzi, 2013;
81 Love, Huber & Anders, 2014).

82 One way is to use discrete probability distributions (e.g. poisson, negative binomial) to deal with
83 huge amount of RNA-Seq data for expression-based classification purpose. Witten et al. (Witten,
84 2011) proposed the sparse Poisson linear discriminant analysis (PLDA) classifier by extending
85 the popular microarray classifier, nearest shrunken centroids algorithm, to discrete RNA-Seq
86 data. The authors also suggested applying a power transformation, since Poisson distribution
87 underestimates the variation observed from the data. Dong et al. (Dong et al., 2016) proposed
88 negative binomial distribution by extending PLDA with the use of negative binomial
89 distribution. Another choice may be to use some transformation approaches (e.g. vst–variance
90 stabilizing transformation- or rlog–regularized logarithmic transformation-) to bring RNA-Seq
91 samples hierarchically closer to microarrays and apply known algorithms for classification
92 applications (Nagalakshmi et al., 2008; Anders & Huber, 2010; Robinson, McCarthy & Smyth,
93 2010).

94 In this study, we designed a comprehensive simulation study, also used four real datasets to
95 examine the effect of dispersion parameter and classification algorithms on RNA-Seq
96 classification. We also considered the effect of other parameters (i) sample size, (ii) number of
97 genes, (iii) number of class, (iv) DE rate, (v) transformation method on classification
98 performance. For each scenario, we performed PLDA and NBLDA as well as other machine
99 learning algorithms i.e. support vector machines (SVM), bagging support vector machines
100 (bagSVM), random forests (RF) and classification and regression trees (CART) algorithms.

101 **Materials and Methods**

102 **A workflow for RNA-Seq classification**

103 Providing a pipeline for classification algorithm of RNA-Seq data gives us a quick snapshot view
104 of how to handle the large-scale transcriptome data and establish a robust inference by using
105 computer-assisted learning algorithms. Therefore, we outlined the count-based classification
106 pipeline for RNA-Seq data in Fig. 1. NGS platforms produce millions of raw sequence reads
107 with quality scores corresponding to each base-call. The first step in RNA-Seq data analysis is to
108 assess the quality of the raw sequencing data for meaningful downstream analysis. The
109 conversion of raw sequence data into ready-to-use clean sequence reads needs a number of
110 processes such as removing the poor-quality sequences, low-quality reads with more than five
111 unknown bases, and trimming the sequencing adaptors and primers. In quality assessment and
112 filtering, the current popular tools are FASTQC
113 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>), HTSeq (Anders, Pyl & Huber,
114 2015), R ShortRead package (Morgan et al., 2009), PRINSEQ ([http://edwards.sdsu.edu/cgi-](http://edwards.sdsu.edu/cgi-bin/prinseq/prinseq.cgi)
115 [bin/prinseq/prinseq.cgi](http://edwards.sdsu.edu/cgi-bin/prinseq/prinseq.cgi)), FASTX Toolkit (http://hannonlab.cshl.edu/fastx_toolkit/) and QTrim
116 (Shrestha et al., 2014). Following these procedures, next step is to align the high-quality reads to
117 a reference genome or transcriptome. It has been reported that the number of reads mapped to the
118 reference genome is linearly related to the transcript abundance. Thus, transcript quantification
119 (calculated from the total number of mapped reads) is a prerequisite for further analysis. Splice-
120 aware short read aligners such as Tophat2 (Kim et al., 2013), MapSplice (Wang et al., 2010) or
121 Star (Dobin et al., 2012) can be preferred instead of unspliced aligners (BWA, Bowtie, etc.). After
122 obtaining the mapped reads, next step is counting how many reads mapped to each transcript. In
123 this way, gene expression levels can be inferred for each sample for downstream analysis. This

124 step can be accomplished with HTSeq (Anders, Pyl & Huber, 2015), bedtools (Quinlan & Hall,
125 2010) and FeatureCounts (Liao, Smyth & Shi, 2014) softwares. However, these counts cannot be
126 directly used for further analysis and should be normalized to adjust between-sample differences.
127 There is no standard tool for normalization, but the popular ones include deseq median ratio
128 (Anders & Huber, 2010), trimmed mean of M values (TMM) (Robinson & Oshlack, 2010), reads
129 per kilobase per million mapped reads (RPKM) (Mortazavi et al., 2008) and quantile (Bullard,
130 2010). For transformation, vst (Anders & Huber, 2010), rlog (Love, Huber & Anders, 2014) and
131 voom (Law et al., 2014) methods can be a method of choice. Apart from these approaches,
132 Witten considered power transformation to decrease the dispersion of data, before applying
133 PLDA classifier (Witten, 2011). Once all mapped reads per transcripts are counted and
134 normalized, we obtain gene-expression levels for each sample.

135 First way is to apply the count based classifiers, e.g. PLDA (Witten, 2011) and NBLDA (Dong et
136 al., 2016) directly to the count data or to the power transformed data. Second way is to use the
137 same workflow of microarray classification after transforming the data hierarchically to
138 microarrays. The crucial steps of classification can be written as feature selection, building
139 classification model and model validation. In feature selection step, we aim to work with an
140 optimal subset of data. This process is crucial to reduce the computational cost, decrease of noise
141 and improve the accuracy for classification of phenotypes, also to work with more interpretable
142 features to better understand the domain (Ding & Peng, 2005). Various feature selection methods
143 have been reviewed in detail and compared in (Xing, Jordan & Karp, 2001). Next step is model
144 building, which refers to the application of a machine-learning algorithm and to learn the
145 parameters of classifiers from training data. Thus, the built model can be used to predict class

146 memberships of new biological samples. The commonly used classifiers include SVM, RF and
147 other tree-based classifiers, artificial neural networks and k-nearest neighbors.

148 In many real life problems, it is possible to experience that a classification algorithm may
149 perform well and perfectly classify training samples, however perform poorly when classifying
150 new samples. This problem is called as overfitting and independent test samples should be used
151 to avoid overfitting and to generalize classification results. Holdout, k-fold cross-validation,
152 leave-one-out cross-validation and bootstrapping are among the recommended approaches for
153 model validation.

154 **Implementation of classifiers**

155 **Simulation study**

156 *Simulation setup*

157 A comprehensive simulation is conducted to investigate the effect of several parameters.
158 Simulated datasets are generated under 864 different scenarios using a negative binomial model
159 as follows:

$$160 \quad X_{ij}|y_i = k \sim NB(s_i g_j d_{kj}, \phi) \quad (1)$$

161 where, s_i is the number of counts per sample, g_j is the number of counts per gene, d_{kj} is the
162 differential expression probability of j^{th} gene between classes k and ϕ is the dispersion parameter.

163 The datasets contain all possible combination of:

- 164 • different dispersion parameters as $\phi=0.01$ (very slightly overdispersed), $\phi=0.1$
165 (substantially overdispersed), $\phi=1$ (highly overdispersed);
- 166 • number of biological samples (n) changing as 40, 60, 80, 100;
- 167 • number of differentially expressed genes (p') as 25, 50, 75, 100;
- 168 • differential expression probability (d_{kj}) as 1%, 5% and 10%;

- 169 • number of classes (k) as 2, 3, 4;
- 170 • method of transformation as rlog and vst.

171 In simulation setup, s_i and g_j are distributed identically and independently as s_i and g_j
172 respectively. Simulated datasets are generated using the CountDataSet function of PoiClaClu
173 package of R software (Witten, 2013) and manipulated based on the details given above. Seed
174 number is set to '10072013' in all analysis steps.

175 *Evaluation process*

176 All datasets are initially simulated for $p=10,000$ genes. Next, the data are split into training
177 (70%) and test sets (30%). All model building processes are applied in training datasets, model
178 performances are evaluated in test sets. We applied near-zero filtering to training data to filter the
179 genes with low counts to eliminate the effect of this genes for further analysis (Kuhn,
180 2008). Genes are filtered based on two criteria: (i) the frequency ratio of the most frequent value
181 to the second most frequent value is higher than 19 (95/5), (ii) the ratio of the number of unique
182 values to the sample size is less than 10%. Filtered genes are also excluded from the test datasets.
183 Next, DESeq2 method is applied to detect the most DE 25, 50, 75 and 100 genes (Love, Huber &
184 Anders, 2014). Selected genes are also selected in test datasets.

185 After selecting the DE genes, training data is normalized using median ratio approach to adjust
186 sample specific differences (Love, Huber & Anders, 2014). After normalization, datasets are
187 transformed using either rlog or vst transformation for SVM, bagSVM, RF and CART
188 algorithms. Classical logarithmic transformation approach transforms the data into a less skewed
189 distribution with less extreme values as well, however the genewise variances are still
190 unstabilized (Love, Huber & Anders, 2014). Normalized count datasets are directly used for
191 PLDA and NBLDA algorithms, since both algorithms use discrete probability distributions to fit

192 the models. In another scenario, a power transformation is applied to minimize the effect of
193 overdispersion and PLDA algorithms is applied to this transformed data. This approach is
194 defined as PLDA₂ in Results section. Note that, test datasets are normalized and transformed
195 using the same parameters as training datasets. Since, training and test datasets should be in same
196 scale and homoscedastic to each other. For instance, to normalize the test datasets, size factors of
197 test datasets are calculated based on the geometric means of training data. Dispersion estimations
198 are applied based on the training models as well.

199 After normalization and transformation processes, the parameters of each classifier are optimized
200 to avoid overfitting and underfitting. A five-fold cross-validation is applied to training data and
201 the parameters that achieves the highest accuracy rate are selected as optimal parameters. Same
202 folds are used for each classifier to make the results comparable. Each classifier is fit with the
203 optimal parameters. Fitted models are used in test datasets for prediction and performance
204 evaluation.

205 The sample sizes are very low relative to the number of genes, since we mimic the real datasets.
206 Thus, the model performances may vary depending on the split of training and test sets. To
207 overcome this limitation, we repeated the entire process 50 times and summarized the results in
208 single statistics, i.e. accuracy rates.

209 **Application to real datasets**

210 In addition to the simulated data, four real datasets, including both miRNA and mRNA datasets,
211 were also used as real life examples (Table 1).

212 *Experimental datasets*

213 *Cervical dataset:* Cervical dataset is a miRNA sequencing dataset obtained from (Witten et al.,
214 2010). miRNAs are non-coding small RNA molecules with average 21-23 bp length and take

215 role in the regulation of gene expression. The objective of this study was to both identify the
216 novel miRNAs and to detect the differentially expressed ones between normal and tumor
217 cervical tissue samples. For this purpose, the authors constructed 58 small RNA libraries,
218 prepared from 29 cervical cancer and 29 matched control tissues. After deep sequencing with
219 Solexa/Illumina sequencing platform, they obtained a total of 25 Mb and 17 Mb RNA sequences
220 from the normal and cancer libraries respectively. Of these 29 tumor samples, 21 of them had a
221 diagnosis of squamous cell carcinomas, 6 of them had adenocarcinomas and 2 were unclassified.
222 In our analysis, we used the data that contains the sequence read counts of 714 miRNAs
223 belonging to 58 human cervical tissue samples, where 29 tumor and 29 non-tumor samples are
224 treated as two distinct classes for prediction.

225 *Alzheimer dataset:* This dataset is another miRNA dataset provided by Leidinger et al.
226 (Leidinger et al., 2013). The authors aimed to discover potential miRNAs from blood in
227 diagnosing alzheimer and related neurological diseases. In this purpose, the authors obtained
228 blood samples from 48 alzheimer patients that were evaluated after undergoing some tests
229 including Alzheimer Disease Assessment Scale-cognitive subscale (ADAS-Cog), Wechsler
230 Memory Scale (WMS), and Mini-Mental State Exam (MMSE) and Clinical Dementia Rating
231 (CDR). A total of 22 age-matched control samples were obtained and all sample libraries were
232 sequenced using 53 Illumina HiSeq2000 platform. After obtaining the raw read counts, the
233 authors filtered the miRNAs with less than 50 counts in each group. We used the data including
234 416 miRNA read counts of 70 samples, where 48 alzheimer and 22 control samples are
235 considered as two separate classes for classification.

236 *Renal cell cancer dataset:* Renal cell cancer (RCC) dataset is an RNA-Seq dataset that is
237 obtained from The Cancer Genome Atlas (TCGA) (Saleem et al., 2013). TCGA is a

238 comprehensive community resource platform for researchers to explore, download, and analyze
239 datasets. We downloaded this dataset (with options level 3, RNASeqV2 data) from this database
240 and obtained the raw 20,531 known human RNA transcript counts belonging to 1,020 RCC
241 samples. This RNA-Seq data has 606, 323 and 91 specimens from kidney renal papillary cell
242 (KIRP), kidney renal clear cell (KIRC) and kidney chromophobe carcinomas (KICH),
243 respectively. These three classes are referred as the most common subtypes of RCC (account for
244 nearly 90%-95% of the total malignant kidney tumors in adults) and treated as three separate
245 classes in our analysis (Goyal et al., 2013).

246 *Lung cancer dataset:* Lung cancer is another RNA-Seq dataset provided from TCGA platform.
247 Same options were used in the download process. The resulting count file contains the read
248 counts of 20,531 transcripts of 1,128 samples. The dataset has two distinct classes including lung
249 adenocarcinoma (LUAD) and lung squamous cell with carcinoma (LUSC) with 576 and 552
250 class sizes, respectively. These two classes are used as class labels in our analysis.

251 *Evaluation process*

252 A similar procedure is applied with the simulation study. Model building is applied in training
253 (70%) and tested in the test (30%) sets. Near-zero filtering is applied to the training set. Filtered
254 genes are also removed from the test set. For renal cell cancer and lung cancer datasets, 5,000
255 genes with highest variances are selected to eliminate the effect of non-informative mRNAs. All
256 miRNA's are used in model building process for cervical and alzheimer datasets. Differential
257 expression was performed to training data using DESeq2 method and genes are ranked from the
258 most significant to the less with increasing number of genes in steps of 25 up to 250 genes.
259 Selected differentially expressed genes in the training data are also selected in the test datasets.
260 Differentially expressed genes in training data are normalized using median ratio approach and

261 transformed using either vst or rlog approaches. Similar to simulation experiments, test datasets
262 are normalized based on the parameters obtained from the training data to make them in same
263 scale and homoscedastic to each other. Since, the sample size of cervical and alzheimer miRNA
264 datasets are relatively small, entire process is applied 50 times. Seed numbers in data selections
265 are set between 1 to 50 and results are summarized based on these 50 repeats. Other model
266 building process are applied as same as the simulation study.

267 **Implementation of classifiers**

268 Seven different algorithms are applied to both simulated and real datasets. In this section, we
269 summarize the background and use of each method.

270 **SVM:**SVM is a classification method based on statistical learning theory, which is developed by
271 Vapnik and his colleges, and has taken great attention because of its strong mathematical
272 background, learning capability and good generalization ability (Vapnik, 2000). Moreover, SVM
273 is capable of nonlinear classification and deal with high-dimensional data. Thus, it has been
274 applied in many fields such as computational biology, text classification, image segmentation
275 and cancer classification (Vapnik, 2000; Korkmaz, Zararsiz & Goksuluk, 2015).

276 In linearly separable cases, the decision function that correctly classifies the data points by their
277 true class labels represented by:

$$278 \mathbf{f}_{\mathbf{w},\mathbf{b}} = \mathbf{sign}(\mathbf{w} \cdot \mathbf{x}_i + \mathbf{b})(2)$$

$$279 \mathbf{i} = \mathbf{1,2},\dots,\mathbf{n}$$

280 In binary classification, SVM finds an optimal separating hyperplane in the feature space, which
281 maximizes the margin and minimizes the probability of misclassification by choosing \mathbf{w} and \mathbf{b} in
282 equation (2).For the linearly non-separable cases, slack variables $\{\xi_1, \dots, \xi_n\}$, which is a penalty
283 introduced by Cortes and Vapnik, can be used to allow misclassified data points, where

284 $\xi_i > 0$ (Cortes & Vapnik, 1995). In many classification problems, the separation surface is
285 nonlinear. In this case, SVM uses an implicit mapping Φ of the input vectors to a high-
286 dimensional space defined by a kernel function ($K(x,y) = \Phi(x_i)\Phi(x_j)$) and the linear classification
287 then takes place in this high-dimensional space. The most widely used kernel functions are linear
288 : $K(x,y) = x_i x_j$, polynomial: $K(x,y) = (x_i x_j + 1)^d$, radial basis function: $K(x,y) = \exp(-\gamma \|x_i - x_j\|^2)$ and
289 sigmoidal: $K(x,y) = \tanh(k(x_i x_j) - c)$, where d is the degree, $\gamma > 0$ sometimes parametrized as
290 $\gamma = 1/2\sigma^2$, and c is a constant. Normalized and transformed (either using vst or rlog) datasets are
291 used as input to SVM classifier. Radial basis kernel function is used in the analysis.

292 **BagSVM:** BagSVM is a bootstrap ensemble method, which creates individuals for its ensemble
293 by training each SVM classifier (learning algorithm) on a random subset of the training set. For a
294 given data set, multiple SVM classifiers are trained independently through a bootstrap method
295 and they are aggregated via an aggregation technique. To construct the SVM ensemble, k
296 replicated training sets are generated by randomly re-sampling, but with replacement, from the
297 given training set repeatedly. Each sample, x_i , in the given training set, may appear repeated
298 times, or not at all, in any particular replicate training set. Each replicate training set will be used
299 to train a specific SVM classifier. Normalized and transformed (either using vst or rlog) datasets
300 are used as input to BagSVM classifier. Number of bootstrap samples were set to 101, since
301 small changes were observed over this number.

302 **CART:** CART, which is introduced by Breiman et al., is one of the most popular tree classifiers
303 and applied in many fields (Breiman et al., 1986). It uses Gini index to choose the split which
304 maximizes the decrease in impurity at each node. If $p(i|j)$ is the probability of class i at node j ,
305 then the Gini index is $1 - \sum_i p^2(i|j)$. When CART grows a maximal tree, this tree is pruned upward
306 to get a decreasing sequence of subtrees. Then, a cross-validation is used to identify the subtree

307 that having the lowest estimated misclassification rate. Finally, the assignment of each terminal
308 node to a class is performed by choosing the class that minimizes the resubstitution estimate of
309 the misclassification probability (Breiman et al., 1984; Dudoit & Fridlyand, 2003). Normalized
310 and transformed (either using vst or rlog) datasets are used as input to CART classifier.

311 **RF:** A random forest is a collection of many CART trees combined by averaging the predictions
312 of individual trees in the forest (Breiman, 2001). The idea behind the RF is to combine many
313 weak classifiers to produce a significantly better strong classifier. For each tree, a training set is
314 generated by bootstrap sample from the original data. This bootstrap sample includes 2/3 of the
315 original data. The remaining of the cases are used as a test set to predict out-of-bag error of
316 classification. If there are m features, m_{try} out of m features are randomly selected at each node
317 and the best split is used to split the node. Different splitting criteria can be used such as Gini
318 index, information gain and node impurity. The value of m_{try} is chosen to be approximately either
319 $\frac{\sqrt{m}}{2}$ or \sqrt{m} or $2\sqrt{m}$ and constant during the forest growing. An unpruned tree is grown for each of
320 the bootstrap sample, unlike CART. Finally, new data is predicted by aggregating, i.e. majority
321 votes, the predictions of all trees (Liaw & Wiener, 2002; Okun & Priisalu, 2007). Normalized
322 and transformed (either using vst or rlog) datasets are used as input to RF classifier. Number of
323 trees was set to 500 in the analysis.

324 **PLDA₁ and PLDA₂:** Let X be an $n \times p$ matrix of sequencing data, where n is number of
325 observations and p is number of features. For sequencing data, X_{ij} indicates the total number of
326 reads mapping to gene j in observation i . Therefore, Poisson log-linear model can be used for
327 sequencing data,

$$328 \quad X_{ij} \sim \text{Poisson}(N_{ij}), \quad N_{ij} = s_i g_j(3)$$

329 where s_i is total number of reads per sample and g_j is total number of reads per region of interest.

330 For RNA-Seq data, equation (3) can be extended as follows,

$$331 \quad X_{ij}|y_i = k \sim \text{Poisson}(N_{ij}d_{jk}), \quad N_{ij} = s_i g_j \quad (4)$$

332 where $y_i \in \{1, \dots, K\}$ is the class of the i^{th} observation, and d_{1j}, \dots, d_{Kj} terms allow the j^{th}

333 feature to be differentially expressed between classes.

334 Let $(x_i, y_i), i = 1, \dots, n$, be a training set and $x^* = (X_1^*, \dots, X_p^*)^T$ be a test set. Using the Bayes' rule

335 as follows,

$$336 \quad P(y^* = k|x^*) \propto f_k(x^*)\pi_k \quad (5)$$

337 where y^* denotes the unknown class label, f_k is the density of an observation in class k and π_k is

338 the prior probability that an observation belongs to class k . If f_k is a normal density with a class-

339 specific mean and common variance then a standard LDA is used for assigning a new

340 observation to the class (Hastie, Tibshirani & Friedman, 2009). In case of the observations are

341 normally distributed with a class-specific mean and a common diagonal matrix, then diagonal

342 LDA methodology is used for the classification (Dudoit, Fridlyand & Speed, 2001). However,

343 neither normality nor common covariance matrix assumptions are not appropriate for sequencing

344 data. Instead, Witten (Witten, 2011) assumes that the data arise from following: Poisson model,

$$345 \quad X_{ij}|y_i = k \sim \text{Poisson}(N_{ij}d_{kj}), \quad N_{ij} = s_i g_j \quad (6)$$

346 where y_i represents the class of the i^{th} observation and the features are independent. The equation

347 (4) specifies that $X_j^*|y^* = k \sim \text{Poisson}(s^* g_j d_{kj})$. First, the size factors for the training data,

348 s_1, \dots, s_n , is estimated. Then s^* , g_j , d_{kj} and π_k are estimated as described in (Witten, 2011).

349 Substituting these estimations into equation (4) and recalling independent features assumption,

350 equation (5) produces,

$$\begin{aligned}
 351 \quad \log P(\widehat{y^*} = k | x^*) &= \log \hat{f}_k(x^*) + \log \hat{\pi}_k + c \\
 352 \quad &= \sum_{j=1}^p X_j^* \log \hat{d}_{kj} - s^* \sum_{j=1}^p \hat{g}_j \log \hat{d}_{kj} + \log \hat{\pi}_k + c', \quad (7)
 \end{aligned}$$

353 where c and c' are constants and do not depend on the class label. The classification rule that
 354 assigns a new observation to the one of the classes for which equation (7) is the largest and it is
 355 linear in x^* (Witten, 2011).

356 Normalized count data is used as input to PLDA₁ classifier. After normalization, a power
 357 transformation ($X'_{ij} = \sqrt{X_{ij} + 3/8}$) is applied to reduce the overdispersion effect and make genes
 358 have constant variance. These normalized and power transformed datasets are used as input to
 359 PLDA₂ classifier. To optimize the tuning parameter, a grid search (30 searches) is applied and
 360 the sparsest model with the highest accuracy rates are selected for classification.

361 **NBLDA:** Dong et al. generalized that PLDA using an extra dispersion parameter (ϕ) of negative
 362 binomial distribution and named the method as negative binomial linear discriminant analysis
 363 (NBLDA) (Dong et al., 2016). This extra dispersion parameter is estimated using a shrinkage
 364 approach detailed in (Yu, Huber & Vitek, 2013). A new test observation will be assigned to its
 365 class based on the following NBLDA discriminating function:

$$\begin{aligned}
 366 \quad \log P(\widehat{y^*} = k | x^*) &= \sum_{j=1}^p X_j^* [\log \hat{d}_{kj} - \log (1 + s^* \hat{g}_j d_{kj} \phi_j)] - \\
 367 \quad &\sum_{j=1}^p \phi_j^{-1} \log (1 + s^* \hat{g}_j d_{kj} \phi_j) + \log \hat{\pi}_k + c', \quad (8)
 \end{aligned}$$

368 Decreasing the dispersion parameter will approximate the data distribution from negative
 369 binomial to poisson, thus will approximate NBLDA to PLDA. More details about this method
 370 can be found in (Dong et al., 2016).

371 Evaluation criteria

372 To validate each classifier model, 5-fold cross-validation was used, repeated 10 times and
373 accuracy rates were calculated to evaluate the performance of each model. Same folds are used
374 for all classifiers to make the results comparable to each other. Accuracy rates are calculated as
375 $(TP + TN)/n$ based on the confusion matrices of test set class labels and test set predictions. For
376 multiclass scenarios, these measures are calculated via one-versus-all approach. Since, class sizes
377 are unbalanced in alzheimer and renal cell cancer datasets, accuracies are balanced using the
378 formula: $(Sensitivity + Specificity) / 2$.

379 **MLSeq R/BIOCONDUCTOR Package**

380 We presented an R package in BIOCONDUCTOR network to make RNA-Seq classification less
381 complicated for researchers and allow users to fit classifiers using single functions. MLSeq
382 package requires from users to upload their raw count data in which can be obtained from feature
383 counting tools (e.g. HTSeq (Anders, Pyl & Huber, 2014), bedtools (Quinlan & Hall, 2010) and
384 FeatureCounts (Liao, Smyth & Shi, 2014) etc.) and allow them to normalize, transform and build
385 classifiers including SVM, bagSVM, RF and CART. Users can access MLSeq package from
386 <https://www.bioconductor.org/packages/release/bioc/html/MLSeq.html> .

387 Results and Discussion

388 Datasets and Classifiers

389 A comprehensive simulation study is designed under 864 different scenarios. Negative binomial
390 distribution is used in all simulation settings. Simulated datasets contain possible combinations
391 of different dispersion parameters, number of biological samples, number of differentially
392 expressed genes, differential expression rate, number of class and transformation method.
393 Moreover, four real mRNA (lung and renal cell cancer) and miRNA (alzheimer and cervical
394 cancer) datasets were used alongside the simulated datasets (Table 1). Support vector machines
395 (SVM), bagging support vector machines (bagSVM), random forests (RF), classification and
396 regression trees (CART), Poisson linear discriminant analysis without power transformation
397 (PLDA1), Poisson linear discriminant analysis with power transformation (PLDA2) and negative
398 binomial linear discriminant analysis (NBLDA) classifiers were applied to each simulated and
399 real datasets. More detailed information about the datasets, classifiers and analysis settings can
400 be found in Methods section.

401 Experimental Results and Discussion

402 Genewise dispersion parameters are estimated for each classifier with method of moments
403 approach and given in Fig. 2. It is seen from the figure that cervical and alzheimer miRNA
404 datasets are very highly overdispersed, while lung and renal cell cancer datasets are substantially
405 overdispersed. Simulation results for $k=2$, $d_{kj}=10\%$ for vst and rlog transformations are given in
406 Fig. 3 and Fig. 4. All other simulation results are given in
407 <http://www.biosoft.hacettepe.edu.tr/MLSeqSupplementary/> and in Supp. file-1. More detailed
408 results are given in Supp. file-2. Results for real datasets are given in Fig. 5.

409 *Effect of simulation parameters*

410 Since combining each significant gene on class conditions is equivalent to combining their
411 predictive abilities, increased number of differentially expressed genes leads to an increase in the
412 classification accuracy (Fig. 4-5). Similarly, in most scenarios, working with more samples and
413 genes has a positive impact on the overall model accuracies. This relationship between number
414 of genes and accuracy is mostly available in $d_{kj}=10\%$ scenarios. Likewise, slight increases is
415 observed in real dataset classification accuracies, since this leads to an increase in the probability
416 of a differentially expressed gene to be included into classification model. For PLDA classifier,
417 high number of selected genes provides alternative options for the lasso shrinkage method to test
418 more genes in classification models. On the other hand, RF builds trees with bagging approach,
419 thus using more genes, and enhances its probability to specify the optimal tree. Increasing
420 sample size improves the discrimination power, as well as the classification accuracy.
421 Conversely, overall accuracies decrease as the number of classes increases. This is due to the fact
422 that the misclassification probability of an observation may be arised depending on the increase
423 in class number.

424 *Dispersion effect on classification accuracies*

425 The performance of each method was increasing depending on the decrease in dispersion
426 parameter. In fact, only decreasing the dispersion parameter makes a significant contribution to
427 classification accuracy, even for the same data and the same scenario. This is mostly clear in $k=2$
428 and $d_{kj}=10\%$ scenarios. When the data is overdispersed, the variance increases; thus we need
429 more sample sizes to achieve the same discrimination power. When we stabilize the sample size
430 and increase the dispersion parameter, this will decrease the discrimination power and lead to a
431 decrease in the classification accuracies. Nagalakshmi et al. mentioned that using biological
432 replicates instead of technical replicates leads to an increase in the dispersion of the data

433 (Nagalakshmi et al., 2008). Based on this idea, increasing the biological variance of the
434 observations will lead to an increase in the data dispersion, thus the classification of observations
435 will be much harder. In differential expression studies of RNA-Seq data, overdispersion is one of
436 the major problems in analysis settings. Many studies are made to overcome this problem
437 (Robinson, McCarthy & Smyth, 2010; Robinson & Oshlack, 2010; Love, Huber & Anders,
438 2014; Anders & Huber, 2012; Law et al., 2014). When we look at the classification accuracy
439 results, overdispersion seems to be a major challenge in classification studies as well. Unless we
440 work with technical replicates, RNA-Seq data is overdispersed and that leads for same gene,
441 counts from different biological replicates have variance exceeding the mean (Nagalakshmi et
442 al., 2008). This overdispersion can be seen in other studies (Robinson & Smyth, 2007, Bloom et
443 al., 2009; Robinson, McCarthy & Smyth, 2010; Zhou, Xia & Wright, 2011; Auer & Doerge,
444 2011). Results of our study revealed that overdispersion has a significant and negative effect on
445 classification accuracies and should be taken into account before model building.

446 *Microarray based classifiers and transformation effect on classification accuracies*

447 Hundreds of microarray based classifiers are developed and able to work in large p and small n
448 settings. However, the technological improvements makes RNA-Seq state-of-the-art approach
449 for quantified transcriptomics. Currently, much of these microarray based classifiers are no
450 longer to be applied to RNA-Seq data, because of the different data types of microarrays and
451 RNA-Seq. Microarray data consists the continuous log-intensities of probes, while RNA-Seq
452 data consists the discrete and overdispersed mapped read counts of sequencing technologies.
453 Results of this study revealed that, transforming the data hierarchically to microarrays (e.g.
454 through $rlog$ and vst) will be a proper approach to recover these classifiers for RNA-Seq
455 classification.

456 Witten et al. stated that normalization strategy has little impact on the classification performance
457 but may be important in differential expression analysis (Witten, 2011). However, data
458 transformation has a direct effect on classification results, by changing the distribution of data. In
459 this study, we used *deseq* normalization with *vst* and *rlog* transformations and had satisfactory
460 classification performances. Love et al. discussed that *vst* transformation does not consider the
461 size factors during the transformation (Love, Huber & Anders, 2014). However, there were no
462 substantial differences between *rlog* and *vst* transformation approaches on classification
463 accuracies. Both transformations can be applied with microarray based classifiers.

464 *Power transformed PLDA and other count based classifiers*

465 Without transformation, PLDA seemed to perform well in very slightly overdispersed datasets.
466 This can be seen in both simulated and real datasets (Fig. 5). For instance, in renal cell carcinoma
467 dataset, the dispersion parameter is very low and the data seem to follow a Poisson distribution.
468 In this dataset, PLDA₁ and PLDA₂ shows similar performances (Fig. 5). However, the
469 performance of this method decreases, when the data becomes more overdispersed. The reason is
470 that PLDA classifies the data using a model based on Poisson distribution. It minimizes the
471 dispersion parameter and makes a significant improvement on classification accuracy using a
472 power transformation (Witten, 2011). Therefore, we suggest that this transformation is very
473 useful and should be applied to be used with PLDA classifier, even in very slightly overdispersed
474 datasets. NBLDA extends this classifier using a negative binomial model. However,
475 classification accuracies of this method is not as higher as PLDA with power transformation.
476 Hence, we believe that this may be due to the dispersion parameter estimation or the unparsed
477 property of the classifier. We conclude that, novel count-based classifiers are still needed for
478 accurate and robust classification of RNA-Seq data.

479 *Overall performances of classifiers*

480 In simulated datasets, power transformed PLDA performed to be the best classifier. RF and
481 NBLDA performed moderately similar. On the other hand, SVM and bagSVM performed the
482 highest classification accuracies in real datasets. PLDA₂, RF and NBLDA have still comparable
483 and high classification accuracies, but lower than SVM and bagSVM. This slight differences
484 may arise from the differences between negative binomial distribution which is used in
485 simulation settings and exact distributions of real RNA-Seq data. In real datasets, SVM and
486 bagSVM classifiers put forward their classification abilities. Moreover, it can be seen from the
487 simulated and real datasets that, the performance of bagSVM classifier increases as the sample
488 size increases. A possible explanation for such observation is that bagSVM uses bootstrap
489 technique and trains better models in datasets with high number of samples. The performance of
490 CART and PLDA₁ were seemed to be lower than the other classifiers.

491 All assessments in this study are made based on the classification accuracies. Another important
492 measure may be the sparsity of classifiers. Since we included mostly the unsparse classifiers to
493 this study, we leave the effect of dispersion parameter on sparsity as a topic for further research.

494 Conclusions

495 A considerable amount of evidence collected from genome-wide gene expression studies
496 suggests that the identification and comparison of differentially expressed genes have been a
497 promising approach of cancer classification for diagnosis and prognosis purposes. Although
498 microarray-based gene expression studies through a combination of classification algorithms
499 such as SVM and feature selection techniques have recently been widely used for new
500 biomarkers for cancer diagnosis (Lee, 2008; Statnikov, Wang & Aliferis, 2008; Anand &
501 Suganthan, 2009; George & Raj, 2011), it has its own limitations in terms of novel transcript
502 discovery and abundance estimation with large dynamic range. Thus, one choice is to utilize the
503 power of RNA-Seq techniques in the analysis of transcriptome for diagnostic classification to
504 surpass the limitations of microarray-based experiment. As mentioned in earlier sections,
505 working with less noisy data can enhance the predictive performance of classifiers, and the novel
506 transcripts may be a biomarker in interested disease or phenotypes.

507 Hundreds of studies are published for microarray based classification. The goal of these studies
508 were to develop or adapt novel approaches to identify a small subset of genes and predict the
509 class labels of a new observation. This has a particular importance in biomedical studies for
510 molecular diagnosis of diseases. In this study, we demonstrated how researchers can classify the
511 RNA-Seq data, which is the state-of-the-art technique for quantification of gene expression. We
512 designed a comprehensive simulation study and also used four real experimental miRNA and
513 mRNA datasets.

514 Besides its technological advantages of RNA-Seq as compared to microarrays, the data obtained
515 from this method is overdispersed due to the inherent variability. This overdispersion seemed to
516 be a drawback for differential expression studies of RNA-Seq data. In this study, we showed that

517 this overdispersion is also a drawback for classification studies, since an increase in the variance
518 will lead to a decrease in the discrimination power. We reach a conclusion that three solutions
519 are available to handle classification of overdispersed RNA-Seq data: (i) increasing the sample
520 size, (ii) transforming the data hierarchically closer to microarrays with variance stabilizers, e.g.
521 vst and rlog transformations, (iii) using count based classifiers, e.g. PLDA₂ and NBLDA. Our
522 simulation study revealed that both microarray based classifiers after an rlog/vst transformations
523 and count based classifiers (that are dealing with the overdispersion) can be efficiently used for
524 classification of RNA-Seq data.

525 To make an overall assessment for the performances of classifiers, PLDA after a power
526 transformation may be a good choice as a count based classifier. Furthermore, its sparsity seems
527 to be an advantage for researchers, however further researches are needed. Surprisingly, the
528 performance of the NBLDA was not satisfactory enough as a count based classifier. Dong et al.
529 mentioned that NBLDA has a better performance than PLDA in moderate and highly
530 overdispersed data (Dong et al., 2016). However, these comparisons are made with same number
531 of genes. Our analysis are performed based on the sparse PLDA classifiers, where the best subset
532 of genes are used in classification. Sparse PLDA classifier after a power transformation
533 performed more accurately in all dispersion settings. We believe that extending NBLDA
534 algorithm into a sparse classifier may improve its classification performance by selecting the
535 most significant genomic features.

536 Moreover, an alternative option may be to transform the data hierarchically closer to microarrays
537 and perform microarray based classifiers. Our results revealed that RF, SVM and bagSVM may
538 perform accurate results after an rlog or vst transformation. Moreover, the efficiency of the
539 bagSVM is improved observably with the increasing sample size.

540 We conclude that, the data with less overdispersion, highly differentially expressed genes, lower
541 number of groups and large sample size may improve the accuracy of the classifiers. Finally, we
542 developed an R/BIOCONDUCTOR package, MLSeq, to make the computation less complicated
543 for researchers and allow them to learn a classification model using various classifiers with one
544 single function. This package can be accessed and downloaded through
545 <https://www.bioconductor.org/packages/release/bioc/html/MLSeq.html>.

546 **Supplemental Information**

547 **Supp-1.** All figures for simulation results

548 **Supp-2.** MLSeq package source

549 **Supp-3.** Simulation R Codes

550 **Supp-4.** Computational Infrastructure

551 **Supp-5.** Computational costs of classifiers

552 **Competing Interests**

553 The authors declare that they have no competing interests.

554 **Author Contributions**

555 GZ developed the method's framework, DG and SK contributed to algorithm design and
556 implementation. GZ, VE and IPD surveyed the literature for other available methods and
557 collected performance data for the other methods used in the study for comparison. GZ, VE,
558 IPD, DG carried out the simulation studies and data analysis. GZ, DG and SK developed MLSeq
559 package. GZ, DG, SK and VE wrote the paper, TU and AO supervised the research process,
560 revised and directed the manuscript and contributed statistical concepts and ideas. All authors
561 read and approved the final manuscript.

562 **Funding**

563 This research was by the Research Fund of Erciyes University [TDK-2015-5468] and Istanbul
564 University [29506].The funders had no role in study design, data collection and analysis,
565 decision to publish, or preparation of the manuscript.

566 **Acknowledgements**

567 We would like to thank A. Keller for sharing the alzheimer data, also thank B. Klaus, S.Anders
568 and M.I. Love for insightful discussions on the simulation settings of this paper.

569 **References**

- 570 Anand A, and Suganthan PN. 2009. Multiclass cancer classification by support vector machines
571 with class-wise optimized genes and probability estimates. *J Theor Biol* 259:533-540.
572 10.1016/j.jtbi.2009.04.013
- 573 Anders S, and Huber W. 2010. Differential expression analysis for sequence count data. *Genome*
574 *Biol* 11:R106. 10.1186/gb-2010-11-10-r106
- 575 Anders S, and Huber W. 2012. Differential expression of RNA-Seq data at the gene level—the
576 DESeq package. Heidelberg, Germany: European Molecular Biology Laboratory (EMBL).
- 577 Anders S, Pyl PT, and Huber W. 2015. HTSeq—a Python framework to work with high-
578 throughput sequencing data. *Bioinformatics* 31:166-169. 10.1093/bioinformatics/btu638
- 579 Auer PL, and Doerge RW. 2011. A two-stage Poisson model for testing RNA-seq data.
580 *Statistical Applications in Genetics and Molecular Biology* 10:1-26.
- 581 Bloom JS, Khan Z, Kruglyak L, Singh M, and Caudy AA. 2009. Measuring differential gene
582 expression by short read sequencing: quantitative comparison to 2-channel gene expression
583 microarrays. *BMC Genomics* 10:221. 10.1186/1471-2164-10-221
- 584 Breiman L. 2001. Random forests. *Machine learning* 45:5-32.
- 585 Breiman L, Friedman J, Stone CJ, and Olshen RA. 1984. Classification and regression trees:
586 CRC press.
- 587 Bullard JH, Purdom E, Hansen KD, and Dudoit S. 2010. Evaluation of statistical methods for
588 normalization and differential expression in mRNA-Seq experiments. *BMC Bioinformatics* 11:94.
589 10.1186/1471-2105-11-94
- 590 Cortes C, and Vapnik V. 1995. Support-vector networks. *Machine learning* 20:273-297.

- 591 Di Y, Schafer DW, Cumbie JS, and Chang JH. 2011. The NBP negative binomial model for
592 assessing differential gene expression from RNA-Seq. *Statistical Applications in Genetics and*
593 *Molecular Biology* 10.
- 594 Díaz-Uriarte R, and De Andres SA. 2006. Gene selection and classification of microarray data
595 using random forest. *BMC Bioinformatics* 7:3.
- 596 Ding C, and Peng H. 2005. Minimum redundancy feature selection from microarray gene
597 expression data. *J Bioinform Comput Biol* 3:185-205.
- 598 Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, and
599 Gingeras TR. 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 29:15-21.
600 10.1093/bioinformatics/bts635
- 601 Dong K, Zhao H, Tong T, and Wan X. 2016. NBLDA: negative binomial linear discriminant
602 analysis for RNA-Seq data. *BMC Bioinformatics* 17:369. 10.1186/s12859-016-1208-1
- 603 Dudoit S, Fridlyand J, and Speed TP. 2001. Comparison of discrimination methods for the
604 classification of tumors using gene expression data. *Journal of the American statistical*
605 *association* 97:77-87.
- 606 Dudoit S, and Fridlyand J. 2003. Classification in microarray experiments. *Statistical analysis of*
607 *gene expression microarray data* 1:93-158.
- 608 Furey TS, Cristianini N, Duffy N, Bednarski DW, Schummer M, and Haussler D. 2000. Support
609 vector machine classification and validation of cancer tissue samples using microarray
610 expression data. *Bioinformatics* 16:906-914.
- 611 George G, and Raj VC. 2011. Review on feature selection techniques and the impact of SVM for
612 cancer classification using gene expression profile. *arXiv preprint arXiv:11091062*.

613 Goyal R, Gersbach E, Yang XJ, and Rohan SM. 2013. Differential diagnosis of renal tumors
614 with clear cytoplasm: clinical relevance of renal tumor subclassification in the era of targeted
615 therapies and personalized medicine. *Arch Pathol Lab Med* 137:467-480. 10.5858/arpa.2012-
616 0085-RA

617 Hastie T, Tibshirani R, and Friedman J. 2009. The elements of statistical learning 2:1. NY
618 Springer.

619 Kuhn M. 2008. Building Predictive Models in R Using the caret Package. *J Stat Softw*:28-5

620 Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, and Salzberg SL. 2013. TopHat2: accurate
621 alignment of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome*
622 *Biol* 14:R36. 10.1186/gb-2013-14-4-r36

623 Korkmaz S, Zararsiz G, and Goksuluk D. 2015. MLViS: A Web Tool for Machine Learning-
624 Based Virtual Screening in Early-Phase of Drug Discovery and Development. *PLoS One*
625 10:e0124600. 10.1371/journal.pone.0124600

626 Law CW, Chen Y, Shi W, and Smyth GK. 2014. voom: Precision weights unlock linear model
627 analysis tools for RNA-seq read counts. *Genome Biol* 15:R29. 10.1186/gb-2014-15-2-r29

628 Lee ZJ. 2008. An integrated algorithm for gene selection and classification applied to microarray
629 data of ovarian cancer. *Artif Intell Med* 42:81-93. 10.1016/j.artmed.2007.09.004

630 Leidinger P, Backes C, Deutscher S, Schmitt K, Mueller SC, Frese K, Haas J, Ruprecht K, Paul
631 F, Stahler C, Lang CJ, Meder B, Bartfai T, Meese E, and Keller A. 2013. A blood based 12-
632 miRNA signature of Alzheimer disease patients. *Genome Biol* 14:R78. 10.1186/gb-2013-14-7-
633 r78

- 634 Liao Y, Smyth GK, and Shi W. 2014.featureCounts: an efficient general purpose program for
635 assigning sequence reads to genomic features. *Bioinformatics* 30:923-930.
636 10.1093/bioinformatics/btt656
- 637 Liaw A, and Wiener M. 2002.Classification and regression by randomForest.*R news* 2:18-22.
- 638 Love MI, Huber W, and Anders S. 2014.Moderated estimation of fold change and dispersion for
639 RNA-seq data with DESeq2.*Genome Biol* 15:550. 10.1186/s13059-014-0550-8
- 640 Morgan M, Anders S, Lawrence M, Aboyoun P, Pages H, and Gentleman R. 2009. ShortRead: a
641 bioconductor package for input, quality assessment and exploration of high-throughput sequence
642 data. *Bioinformatics* 25:2607-2608. 10.1093/bioinformatics/btp450
- 643 Mortazavi A, Williams BA, McCue K, Schaeffer L, and Wold B. 2008. Mapping and quantifying
644 mammalian transcriptomes by RNA-Seq. *Nat Methods* 5:621-628. 10.1038/nmeth.1226
- 645 Nagalakshmi U, Wang Z, Waern K, Shou C, Raha D, Gerstein M, and Snyder M. 2008. The
646 transcriptional landscape of the yeast genome defined by RNA sequencing. *Science* 320:1344-
647 1349. 10.1126/science.1158441
- 648 Okun O, and Priisalu H. 2007. Random forest for gene expression based cancer classification:
649 overlooked issues. Iberian Conference on Pattern Recognition and Image Analysis: *Springer*. p
650 483-490.
- 651 Quinlan AR, and Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
652 features. *Bioinformatics* 26:841-842. 10.1093/bioinformatics/btq033
- 653 Rapaport F, Zinovyev A, Dutreix M, Barillot E, and Vert JP. 2007. Classification of microarray
654 data using gene networks. *BMC Bioinformatics* 8:35. 10.1186/1471-2105-8-35

- 655 Robinson MD, McCarthy DJ, and Smyth GK. 2010. edgeR: a Bioconductor package for
656 differential expression analysis of digital gene expression data. *Bioinformatics* 26:139-140.
657 10.1093/bioinformatics/btp616
- 658 Robinson MD, and Oshlack A. 2010. A scaling normalization method for differential expression
659 analysis of RNA-seq data. *Genome Biol* 11:R25. 10.1186/gb-2010-11-3-r25
- 660 Robinson MD, and Smyth GK. 2007. Moderated statistical tests for assessing differences in tag
661 abundance. *Bioinformatics* 23:2881-2887. 10.1093/bioinformatics/btm453
- 662 Saleem M, Padmanabhuni SS, Ngomo A-CN, Almeida JS, Decker S, and Deus HF. 2013.
663 Linked cancer genome atlas database. *Proceedings of the 9th International Conference on*
664 *Semantic Systems*: ACM. p 129-134.
- 665 Shrestha RK, Lubinsky B, Bansode VB, Moinz MB, McCormack GP, and Travers SA. 2014.
666 QTrim: a novel tool for the quality trimming of sequence reads generated using the Roche/454
667 sequencing platform. *BMC Bioinformatics* 15:33.
- 668 Sonesson C, and Delorenzi M. 2013. A comparison of methods for differential expression analysis
669 of RNA-seq data. *BMC Bioinformatics* 14:91. 10.1186/1471-2105-14-91
- 670 Statnikov A, Wang L, and Aliferis CF. 2008. A comprehensive comparison of random forests
671 and support vector machines for microarray-based cancer classification. *BMC Bioinformatics*
672 9:319. 10.1186/1471-2105-9-319
- 673 Vapnik VN. 2000. The nature of statistical learning theory, ser. Statistics for engineering and
674 information science. New York: *Springer* 21:1003-1008.
- 675 Wang Z, Gerstein M, and Snyder M. 2009. RNA-Seq: a revolutionary tool for transcriptomics.
676 *Nat Rev Genet* 10:57-63. 10.1038/nrg2484

677 Wang K, Singh D, Zeng Z, Coleman SJ, Huang Y, Savich GL, He X, Mieczkowski P, Grimm
678 SA, Perou CM, MacLeod JN, Chiang DY, Prins JF, and Liu J. 2010. MapSplice: accurate
679 mapping of RNA-seq reads for splice junction discovery. *Nucleic Acids Res* 38:e178.
680 10.1093/nar/gkq622

681 Witten D. 2013. PoiClaClu: Classification and clustering of sequencing data based on a Poisson
682 model. R package version 1.0.2., <https://CRAN.R-project.org/package=PoiClaClu>

683 Witten DM. 2011. Classification and clustering of sequencing data using a Poisson model. *The*
684 *Annals of Applied Statistics*:2493-2518.

685 Witten D, Tibshirani R, Gu SG, Fire A, and Lui WO. 2010. Ultra-high throughput sequencing-
686 based small RNA discovery and discrete statistical biomarker analysis in a collection of cervical
687 tumours and matched controls. *BMC Biol* 8:58. 10.1186/1741-7007-8-58

688 Xing EP, Jordan MI, and Karp RM. 2001. Feature selection for high-dimensional genomic
689 microarray data. *ICML: Citeseer*. p 601-608.

690 Yu D, Huber W, and Vitek O. 2013. Shrinkage estimation of dispersion in Negative Binomial
691 models for RNA-seq experiments with small sample size. *Bioinformatics* 29:1275-1282.
692 10.1093/bioinformatics/btt143

693 Zhou YH, Xia K, and Wright FA. 2011. A powerful and flexible approach to the analysis of
694 RNA sequence count data. *Bioinformatics* 27:2672-2678. 10.1093/bioinformatics/btr449

695 Zhu J, and Hastie T. 2004. Classification of gene microarrays by penalized logistic regression.
696 *Biostatistics* 5:427-443. 10.1093/biostatistics/5.3.427

697 **Figure legends**

698 **Fig. 1.** RNA-Seq classification workflow

699 **Fig. 2.** Genewise dispersion estimations for real datasets

700 **Fig. 3.** Simulation results for $k=2, d_{kj}=10\%$, transformation: vst. Figure shows the performance
701 results of classifiers with changing parameters of sample size (n), number of genes (p) and type
702 of dispersion ($\varphi=0.01$: very slight, $\varphi=0.1$: substantial, $\varphi=1$: very high)

703 **Fig. 4.** Simulation results for $k=2, d_{kj}=10\%$, transformation: rlog. Figure shows the performance
704 results of classifiers with changing parameters of sample size (n), number of genes (p) and type
705 of dispersion ($\varphi=0.01$: very slight, $\varphi=0.1$: substantial, $\varphi=1$: very high)

706 **Fig. 5.** Results obtained from real datasets. Figure shows the performance results of classifiers
707 for datasets with changing number of most significant number of genes

Table 1 (on next page)

Description of real RNA-Seq datasets used in this study

Table 1 - Description of real RNA-Seq datasets used in this study

1 **Table**2 **Table 1 - Description of real RNA-Seq datasets used in this study**

Dataset	Number of groups	Sample size	Number of features
Cervical cancer (Witten et al., 2010)	2	58 (29 cervical cancer, 29 control)	714 miRNAs
Alzheimer (Leidinger et al., 2013)	2	70 (48 alzheimer, 22 control)	416 miRNAs
Renal cell cancer (Saleem et al., 2013)	3	1,020 (606 KIRP, 323 KIRC, 91 KICH)	20,531 mRNAs
Lung cancer (Saleem et al., 2013)	2	1,128 (576 LUAD, 552 LUSC)	20,531 mRNAs

3

Figure 1 (on next page)

RNA-Seq classification workflow

Fig 1 - RNA-Seq classification workflow

Raw sequence reads

```
graph TD; A[Raw sequence reads] --> B[Quality assessment (trimming, filtering, etc.)]; B --> C[Mapping to the reference genome]; C --> D[Feature counting]; D --> E[Data normalization and transformation]; E --> F[Feature selection]; F --> G[Building classification model]; G --> H[Model validation];
```

Quality assessment (trimming, filtering, etc.)

Mapping to the reference genome

Feature counting

Data normalization and transformation

Feature selection

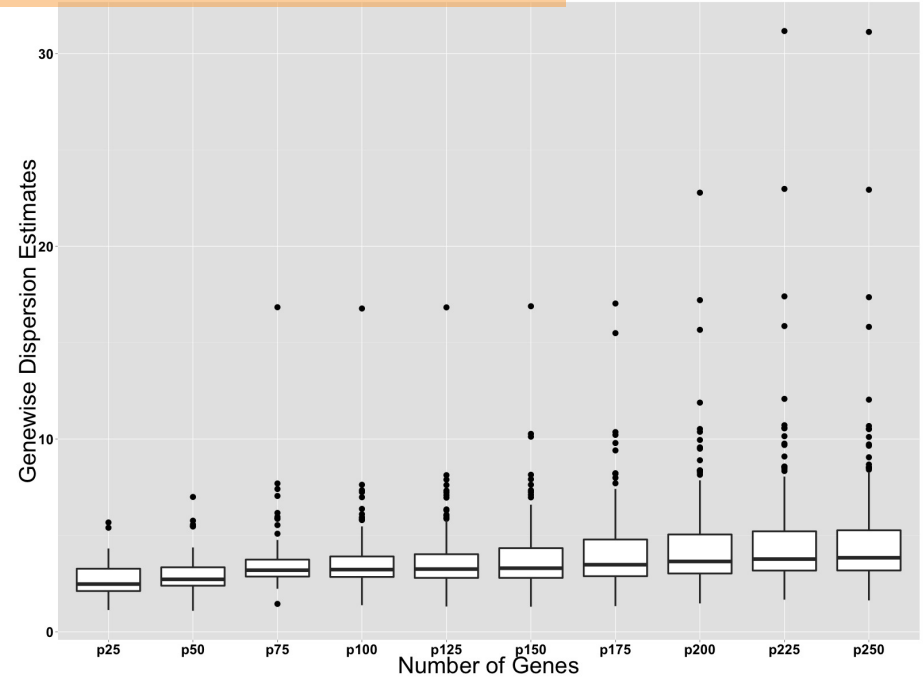
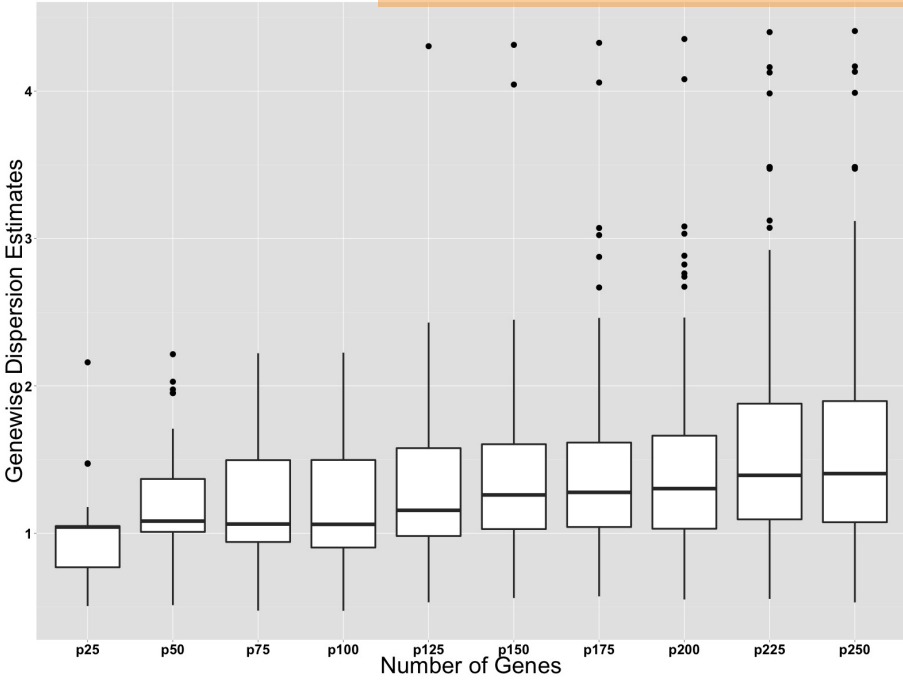
Building classification model

Model validation

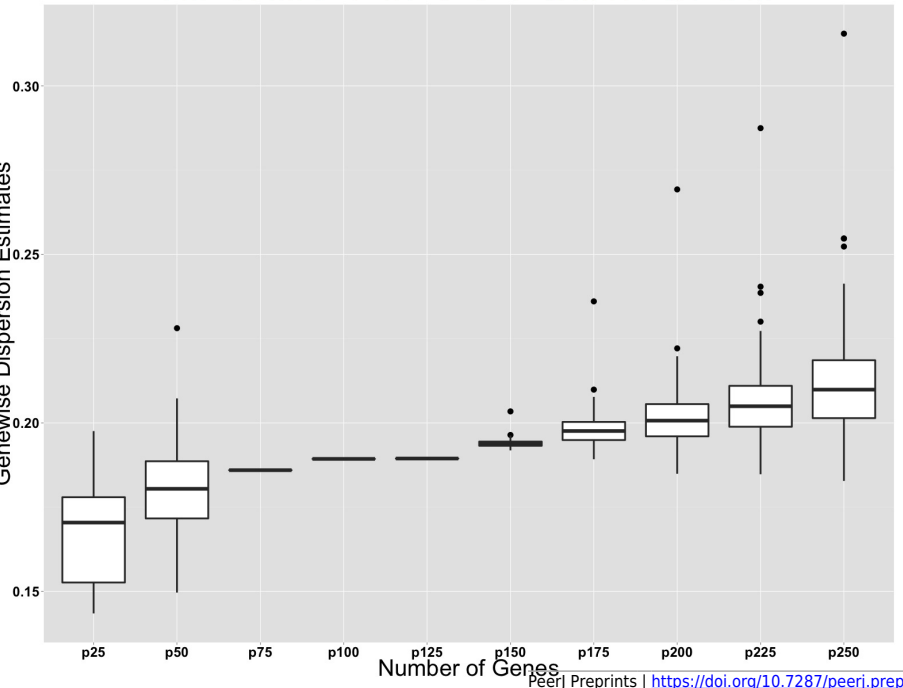
Figure 2 (on next page)

Genewise dispersion estimations for real datasets

Fig 2 - Genewise dispersion estimations for real datasets



RENAL CELL CANCER



LUNG

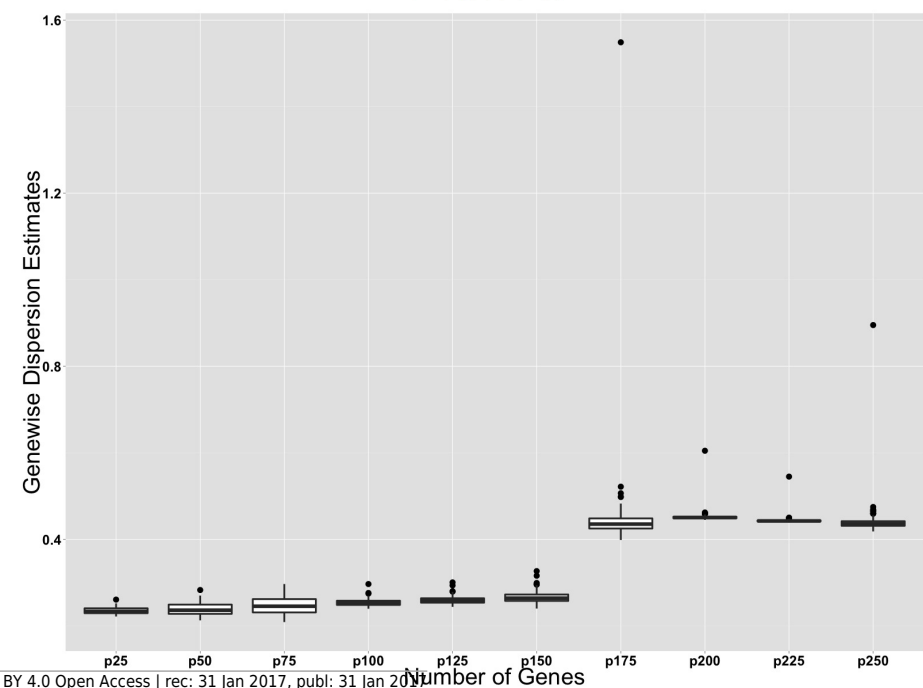


Figure 3(on next page)

Simulation results for $k=2, d_{kj}=10\%$, transformation: vst. Figure shows the performance results of classifiers with changing parameters of sample size (n), number of genes (p) and type of dispersion ($\varphi=0.01$: very slight

Fig 3 - Simulation results for $k=2, d_{kj}=10\%$, transformation: vst. Figure shows the performance results of classifiers with changing parameters of sample size (n), number of genes (p) and type of dispersion ($\varphi=0.01$: very slight, $\varphi=0.1$: substantial, $\varphi=1$: very high)

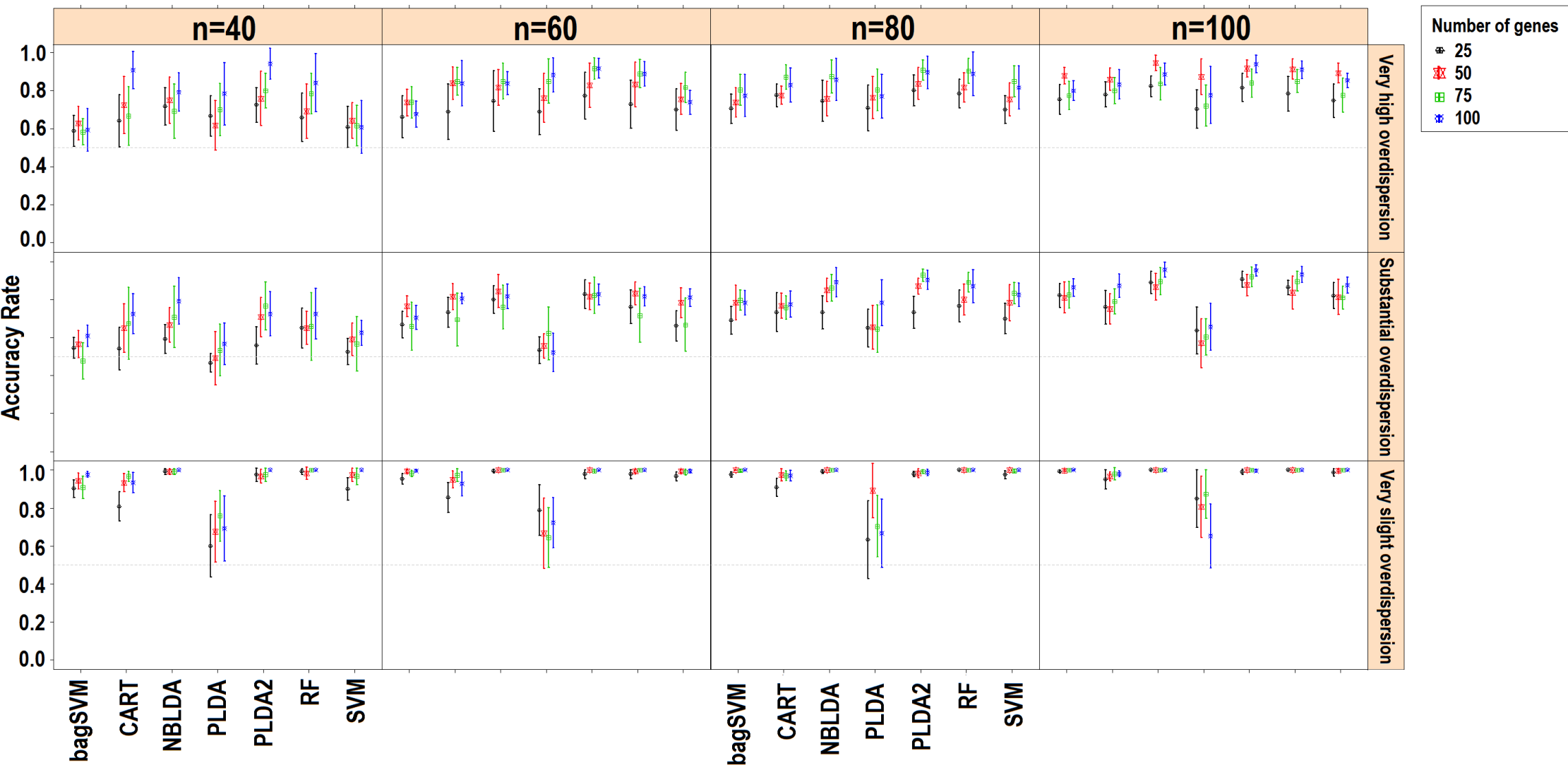


Figure 4(on next page)

Simulation results for $k=2, d_{kj}=10\%$, transformation: rlog. Figure shows the performance results of classifiers with changing parameters of sample size (n), number of genes (p) and type of dispersion ($\varphi=0.01$: very sli

Fig 4 - Simulation results for $k=2, d_{kj}=10\%$, transformation: rlog. Figure shows the performance results of classifiers with changing parameters of sample size (n), number of genes (p) and type of dispersion ($\varphi=0.01$: very slight, $\varphi=0.1$: substantial, $\varphi=1$: very high)

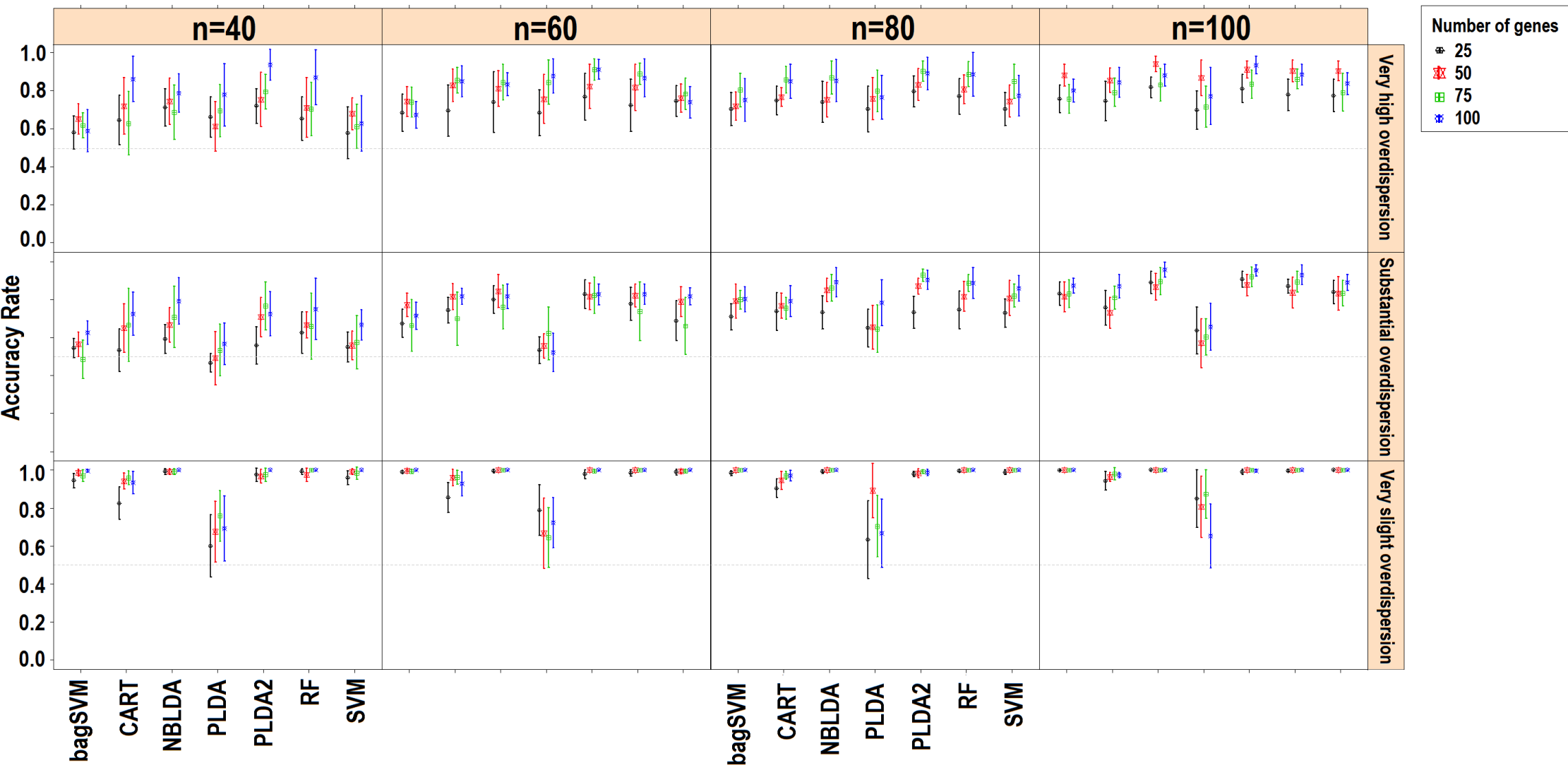


Figure 5(on next page)

Results obtained from real datasets. Figure shows the performance results of classifiers for datasets with changing number of most significant number of genes

Fig 5 - Results obtained from real datasets. Figure shows the performance results of classifiers for datasets with changing number of most significant number of genes

