

DIFFERENTIAL EXPRESSION ANALYSIS FOR TUMOR AND NORMAL COLON TISSUES PROBED BY OLIGONUCLEOTIDE ARRAYS

ERALDA GJIKI (DHAMO), LULE BASHA (HALLAÇI)

Department of Applied Mathematics, Faculty of Natural Science, University of Tirana

e-mail: eralda.dhamo@fshn.edu.al

Abstract

Differential expression analysis used for detection of genes differentially expressed between cancer and healthy patients will be done in this study. Distributional assumption and appropriate multiple testing correction will be used to not miss out the signal. Combination of statistical methods such as heat map will be used to summarize the findings of up and down differentially expressed genes. Also PCA will be used to understand the number of components we need to explain more than 80% of the variation in the data. Cluster techniques will be an option to analyze samples using all differentially expressed genes. Advantages and limitations of our proposal will be discussed associated with the report of the classification accuracy of our chosen approach. We also propose a methodology which may be used to increase the clustering accuracy through dimension reduction, variable selection or a combination of both. At the end we implement a dimension reduction technique, variable selection, or a combination of both on the differentially expressed genes and report our clustering performance.

Key words: DEG, FDR, PCA, cluster, gene expression.

Përmbledhje:

Në këtë studim do të zhvillohet një analizë diferenciale e përdorur për zbulimin e gjeneve të rëndësishëm në pacientët e prekur nga tumori dhe pacientëve të shëndetshëm. Me qëllim ruajtjen e informacionit do të zhvillohet një analizë mbi shpërndarjen dhe testet e duhura për të korrigjuar të dhënat. Kombinimi i metodave statistikore, si hartat e nxehtësisë, do të përdoren për të përmbledhur gjetjet e gjeneve të rëndësishëm të klasifikuara si të rregulluara lart dhe poshtë. Gjithashtu metoda PCA do të përdoret për të kuptuar numrin e komponentëve që na nevojiten për të shpjeguar më shumë se 80% të variacionit në të dhëna. Teknikat e grupimit do të jenë një opsion për të analizuar mostrat duke përdorur të gjitha gjenet e shprehura në mënyrë diferenciale. Përparësitë dhe kufizimet e propozimit tonë do të diskutohen lidhur me raportin e saktësisë së klasifikimit të qasjes sonë. Ne propozojmë gjithashtu një metodologji që mund të përdoret për të rritur saktësinë e grupimit përmes reduktimit të përmasave, përzgjedhjes së variablave ose një kombinimi të të dyjave. Në fund do të zbatohet një teknikë të reduktimit të përmasave, përzgjedhjes së variablave, një kombinim të të dyjave në gjenet e shprehura në mënyrë diferenciale dhe do të raportojmë performancën tonë të grupimit.

Fjalë kyçe: DEG, FDR, PCA, grupe, shprehje e gjeneve.

1.1. Introduction and data

Many researches are done especially for differentially expressed genes pointing out the importance of this pre-processing step for a better process and evaluation of the genes. This study use data from Alon et al., (1999). Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays.

The dimension of our data are 62 rows and 2000 columns. Which means we have 2000 genes and 62 samples. From which 22 samples of class 1 which corresponds to normal tissues and 40 samples of class 2 which corresponds to tumor tissues. The dataset consists of three main variables:

X is a (62 x 2000) matrix giving the expression levels of 2000 genes for the 62 Colon tissue samples. Each row corresponds to a patient, each column to a gene.

Y is a numeric vector of length 62 giving the type of tissue sample (tumor or normal).

gene.names is a vector containing the names of the 2000 genes for the gene expression matrix X.

	1	2	3	4	5	6	7	8	9	10	11	12	13
1	8589.416	5468.241	4263.408	4064.936	1997.893	5282.325	2169.720	2773.421	7526.386	4607.6762	2598.0600	1522.6462	1300
2	9164.254	6719.529	4883.449	3718.159	2015.221	5569.907	3849.059	2793.387	7017.734	4802.2524	1672.9750	1792.1769	3792
3	3825.705	6970.361	5369.969	4705.650	1166.554	1572.168	1325.402	1472.259	3296.951	2786.5821	2441.4188	1487.6712	1315
4	6246.449	7823.534	5955.835	3975.564	2002.613	2130.543	1531.142	1714.631	3869.785	4989.4071	1723.5800	1298.7250	1305
5	3230.329	3694.450	3400.740	3463.586	2181.420	2922.782	2069.246	2948.575	3303.371	3109.4131	2724.2662	2557.7846	3164
6	2510.325	1960.655	1566.315	3072.816	1810.205	1673.564	1290.421	2465.846	1675.544	1312.8083	2289.0350	2162.2865	2070
7	7126.599	3779.068	3705.554	6594.514	2460.905	3775.682	2621.419	2047.281	6411.267	3857.1190	1496.4063	1604.2615	1485
8	4028.710	3156.159	2870.255	4417.591	1854.106	2828.304	1427.526	3390.706	4373.044	3080.4512	5784.1212	2222.2077	2502
9	9330.679	7017.230	4723.783	9491.534	5346.542	1557.143	1969.080	2295.403	6880.346	6162.8929	7927.6525	2022.6731	1427
10	5271.517	4740.768	3318.514	6792.348	2632.889	5449.207	4623.212	3277.404	4488.060	3343.8107	3830.4250	3046.9462	8878
11	14876.407	3201.905	2327.626	11248.680	5893.279	5319.857	9939.246	4058.644	14144.835	3282.6202	6707.6762	3870.9096	7921

Figure 1: Screenshot of the data ColonX

1.2. Pre-Processing the data Colon X

At this step we will try to have a better understanding of the data. Using histograms, Boxplot, Q-Q plot and other graphical test we will observe the distribution of our data and decide if any transformation (normalization) should be used to improve further analysis. Figure 2 shows the histogram for the first 4 genes. We do not observe a clear normal distribution but we may use other graphical visualization to decide.

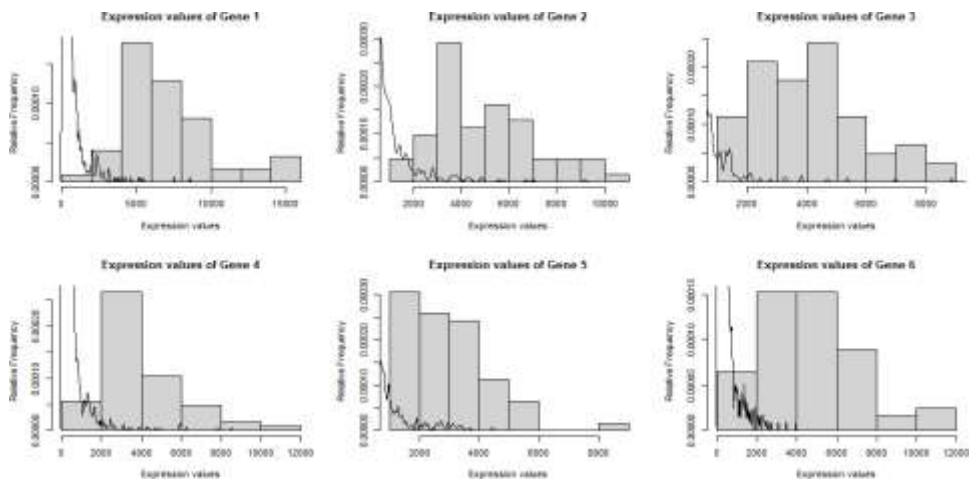


Figure 2: Histogram for the first 6 genes

Observing the boxplot of the first 70 genes we obtain Figure 3. We clearly observe a distribution which is not normal for many of the genes. We may suggest a normalization of the data based on the distributions showed from boxplot graphs.

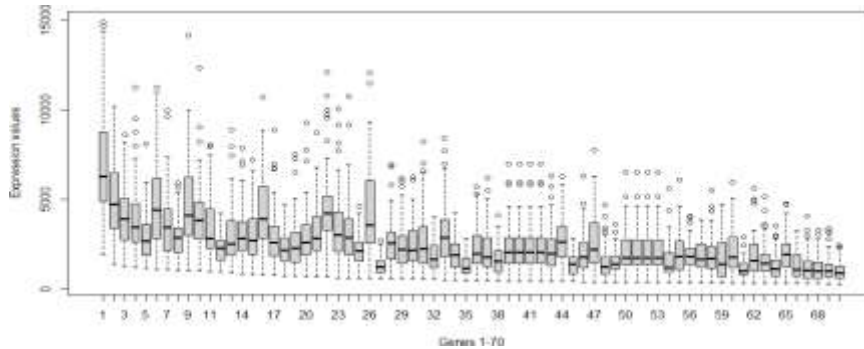


Figure 3: Boxplot for the first 70 genes

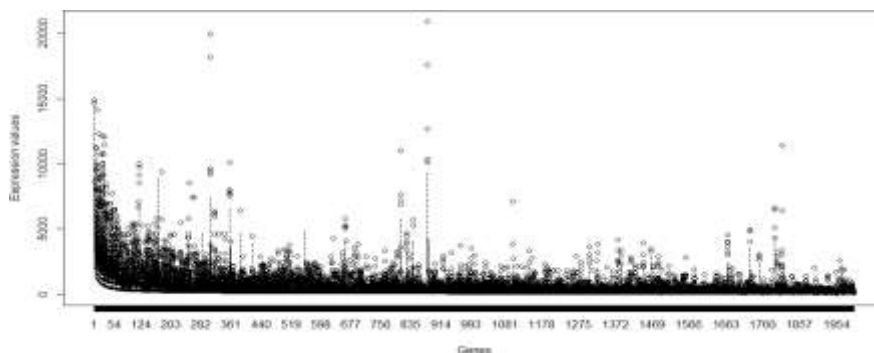


Figure 4: Boxplot for all genes

When we increase the number of genes from 20 up to 2000 we observe a change in the distribution. Which seem to be right skewed distributed (fat

tail or heavy tail distributions are a special probability distribution that exhibits a large skewness or kurtosis). The first genes seem to be different from the other genes (box plot and also density plot shows it). So, again here we confirm the necessity of transforming the data.

Also from the QQplot (Figure 5) we observe some deviations at the tail of the graph which suggest outliers and so not normally distributed data. For example, gene 1 and Gene 4 have visible deviations from the theoretical quantiles and so they may not be normally distributed. (Figure 5)

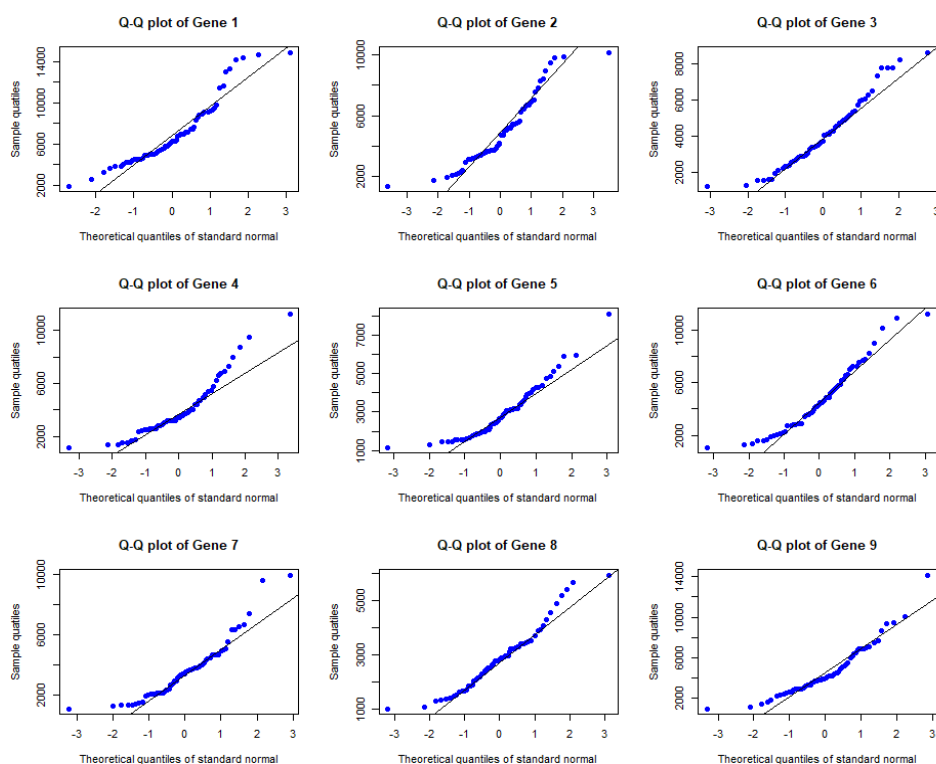


Figure 5: QQ plot for the first 9 genes

1.3. Transformation and distributional assumption

First approach: `getNormMatrix()`. In the first approach we tried to normalize the data by using the reference package `getNormMatrix()`. Based from the review of this function in R it is mentioned that normalization without evaluation, TU is recommended for the normalization of gene expression data, as it has been already ranked as the best method for both scRNA-seq and bulk RNA-seq data, Zhenfeng et al., (2019). The results for this normalization procedure in our data does not show significant improvements in the distribution of the data.

Second approach: $\log_2()$ In the second approach we used logarithm of base 2 to transform the data. And based on the visualization it looks like the second approach performed better than the first one.

Third approach: ZScore. We have also tried the Z – score transformation on the data ($Z = (X - \text{mean})/\text{sd}$). But again this process was not successful and did not overpass \log_2 transformation. Below are a set of graphs showing how these proposed transformations performed and how the data look like after the transformation.

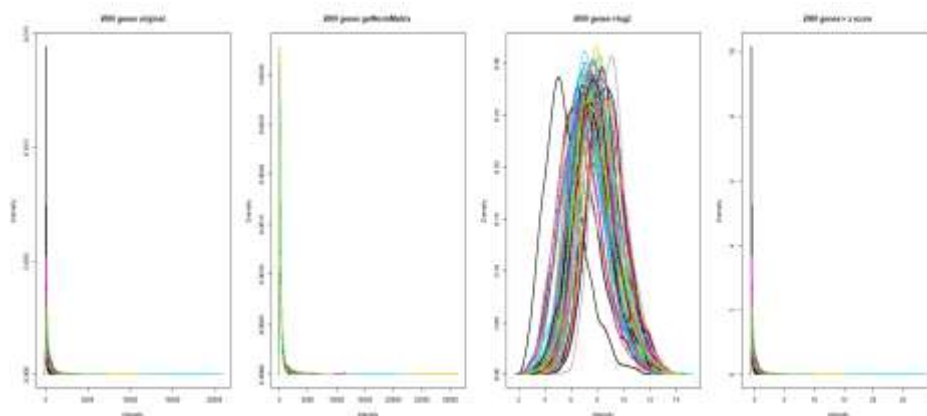


Figure 6: Density plot comparison of transformation

Figure 6 gives (from left to right) the density plot for the original data, getNormMatrix, \log_2 and Zscore transformation. Again here, we observe that \log_2 transformation offers a clearer distribution close to normal distribution of the data.

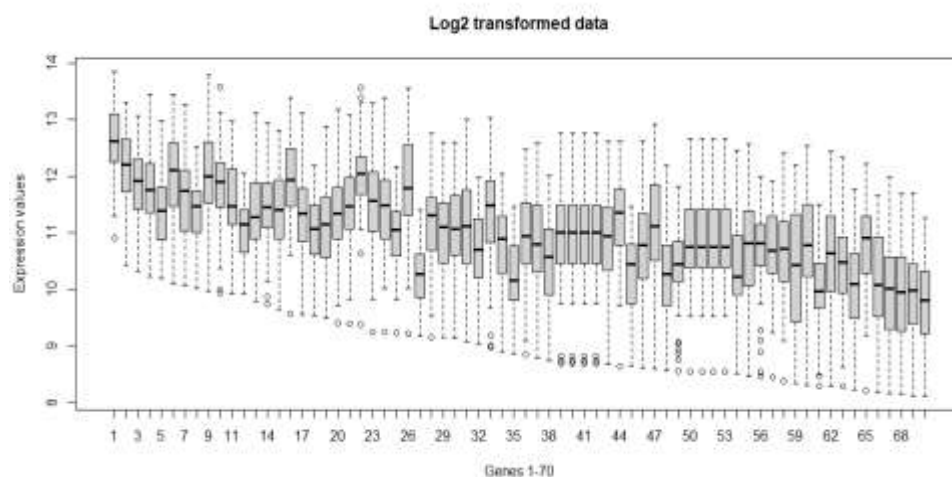


Figure 7: Boxplot of 70 genes for \log_2 transformation

Figure 7 visualize the boxplot for \log_2 transformation approach. Observing all genes, we clearly see different behavior of the first genes (numbered as in Colon\$X) and the last genes. This difference is clear and will be present when we will discuss in the other steps of gene analysis. Here we may ask: do we have clear clusters of genes which will help us identify and explain this behavior among genes?

Materials and methods

1. Multiple testing correction methods

Getting a list of differentially expressed genes means that we need to choose an absolute threshold for the \log_2 fold change (column $\log_2\text{FoldChange}$ in the output of the functions from R packages) and the adjusted p – value (column padj). Therefore, any one can make different list of differential genes based on the selected thresholds. It is common to choose a \log_2 fold change threshold of $|1|$ or $|2|$ and an adjusted p – value of 0.01 for instance. Below we are analyzing the distribution of the p -value by histogram visualization. Based on the distribution of the p -values (Figure 10) we may notice that: For the adjusted p – value we have a peak at 0, but we also have a peak close to 1. What do we do in this situation? What is this behavior telling us? Is this something that was proceeded from boxplot and density plot above?

A p -value close to 1 may indicate increase of gene in response to a drug or they belong to a pathological case! recommendations for this situation are to filter these genes out beforehand (it's not like we are losing any information!), (Crow et al., 2019; Barbey et al., 2020; Barbey et al., 2021; Rodriguez-Esteban and Jiang, 2017)

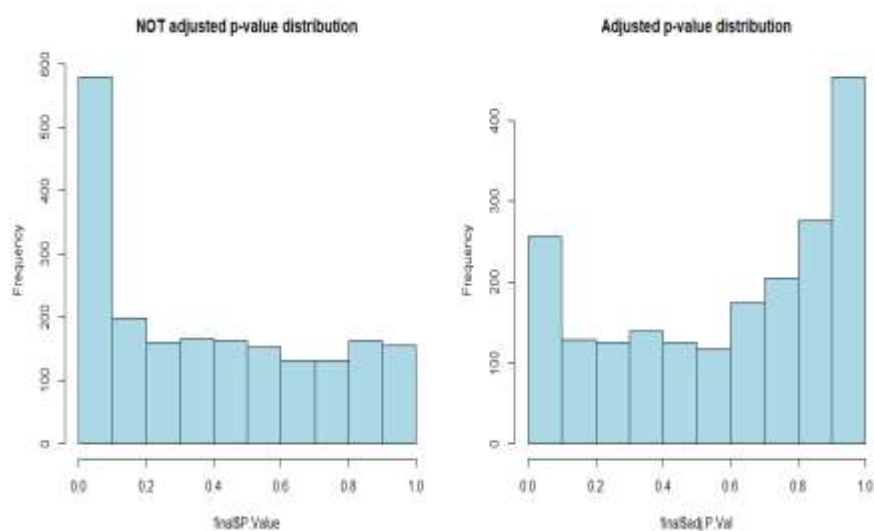


Figure 8: Histogram of p -value and adj. p -value of all genes for \log_2 transformation

Going through commonly used methods for controlling FWER: Bonferroni's method, Holm's method and multiple testing error measure known as false discovery rate FDR. We obtain the following results.

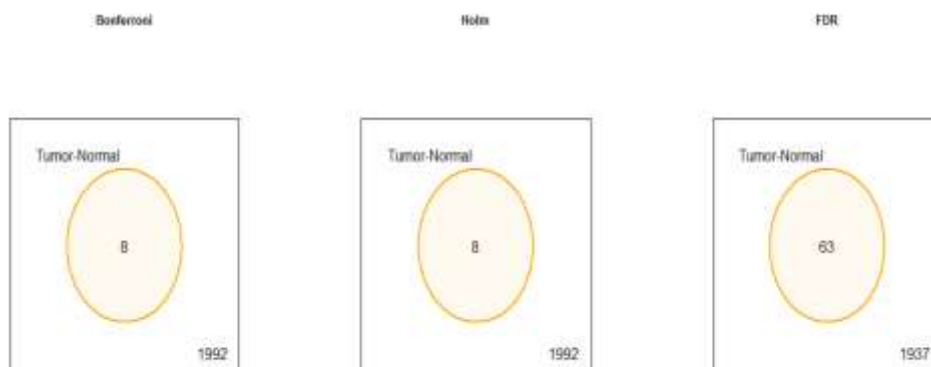


Figure 9: Bonferroni Holm FDR on the original data

Observe above the set of genes when directly applied to the original data which is less than the number of differentially expressed genes obtained after we use log2 transformation. We have from FDR: 114 Up and 29 Down and 1857 Not-significant genes. (Figure 10). The ID of these genes are those in rows: Now our subset has 143 genes and 62 samples.

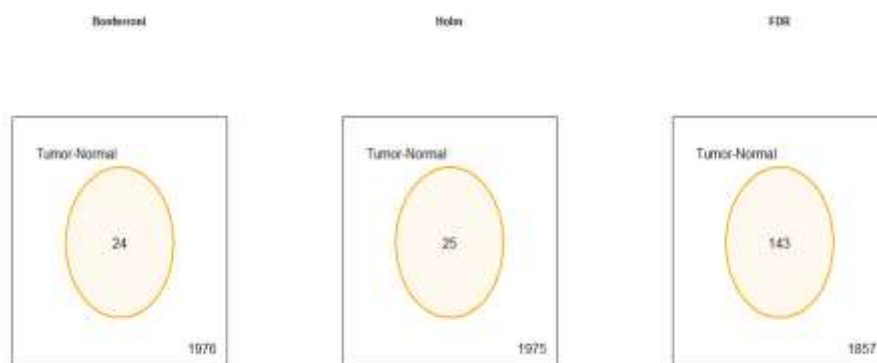


Figure 10: Bonferroni Holm FDR after log2 transformation

So, again here we observe that FDR is less conservative and because the other methods give just 24, 25 genes we will decide to go with FDR decision of considering as differentially expressed genes 143 from total of 2000. The output from FDR for the log2 transformed data is shown in Figure 11. By transforming the data, we went from 3.15% (63 genes) of considered genes to 7.15% (143 genes from which 29 Down and 114 UP). Volcano plot (Figure 13) shows how our genes are represented from each method.

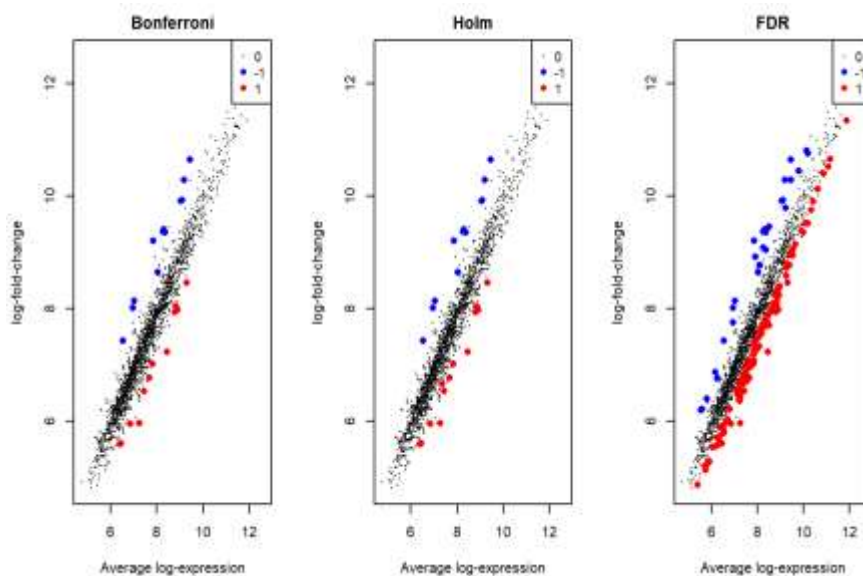


Figure 11: Volcano plot for Bonferroni Holm FDR after log2 transformation

	logFC	AveExpr	t	P.Value	adj.P.Val	B
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
493	-1.533126	8.367340	-6.292594	1.517618e-08	3.035236e-05	9.315900
1671	-2.005973	7.260691	5.939360	6.891831e-08	5.596730e-05	7.906059
249	-1.845727	9.451096	-5.893216	8.380096e-08	5.596730e-05	7.723967
1423	-2.081832	7.860515	-5.733383	1.642684e-07	8.213418e-05	7.097333
625	1.452040	8.910696	5.458614	5.139644e-07	2.055858e-04	6.036233
1042	1.392027	7.668984	5.389082	6.835260e-07	2.278420e-04	5.771224

Figure 12: FDR output after log2 transformation

Figure 12 gives the ID of the genes corresponding to differentially expressed genes among Tumor and Normal (sample 1 and sample 2) which were classified by FDR (in total 143 genes).

In R markdown file we have subtracted from the original (2000 genes, remember: log2 transformed dataset) the genes which were classified from FDR as the most differentially expressed genes and created a subset which will be used for further analysis. Figure 13 show a heat map of the genes before and after we created the subset of the genes selected by FDR.

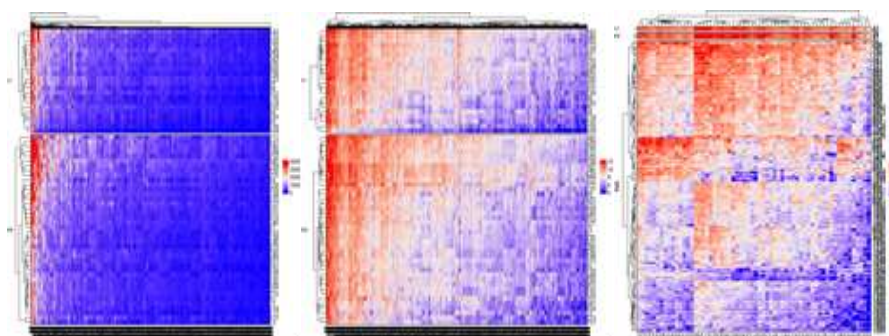


Figure 13: Heat map of original, log2 and after FDR subset created (from left to right)

2. PCA with prcomp()

To check which genes differ the most across the two groups (normal, tumor), we perform the PCA on the data using the `prcomp()` function. By default, the `prcomp()` function centers the variables to have mean zero. By using the option `scale = TRUE`, we scale the variables to have standard deviation one. The rotation matrix provides the principal component loadings; each column of `pr.out$rotation` contains the corresponding principal component loading vector. In this case, the loading can be considered as the weight of each gene in both the groups. We may notice (Figure 14) that the first component explains approximately 67% of the variance and together with component 2 they go up to 80% of variance explained. Associated with the output is the Scree plot (Figure 14) which shows how each Dim is explaining the variation in the data.

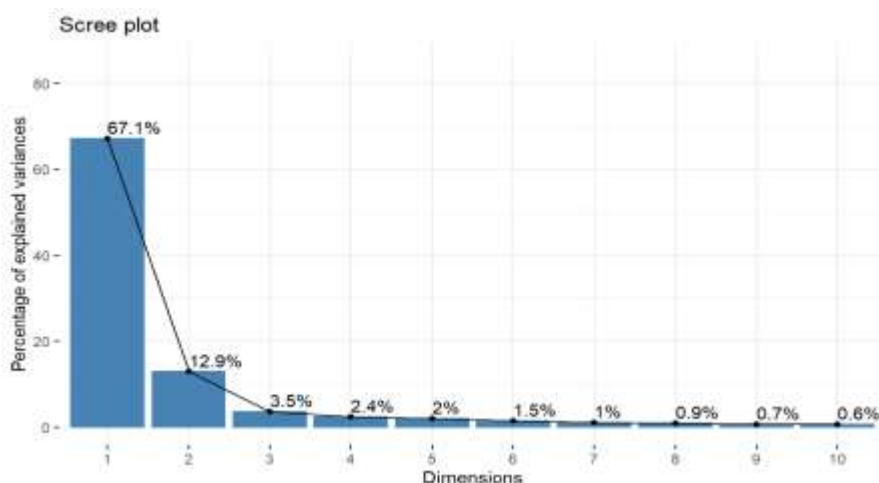


Figure 14: PCA scree plot

The correlation between a variable and a principal component (PC) is used as the coordinates of the variable on the PC. The representation of variables

differs from the plot of the observations: The observations are represented by their projections, but the variables are represented by their correlations.

The plot in Figure 15 (left) is also known as variable correlation plots. It shows the relationships between all variables. It can be interpreted as follow: Positively correlated variables are grouped together. Negatively correlated variables are positioned on opposite sides of the plot origin (opposed quadrants). The distance between variables and the origin measures the quality of the variables on the factor map. Variables that are away from the origin are well represented on the factor map. The plot to the right shows the individual contributions to dimensions.

A high \cos^2 value indicates a good representation of the variable on the principal component. In this case the variable is positioned close to the circumference of the correlation circle. A low \cos^2 indicates that the variable is not perfectly represented by the PCs. In this case the variable is close to the center of the circle. It's possible to color variables by their \cos^2 values using the argument `col.var = "cos2"`. This produces a gradient colors. In this case, the argument `gradient.cols` can be used to provide a custom color. For instance, `gradient.cols = c("white", "blue", "red")` means that: variables with low \cos^2 values will be colored in "white" variables with mid \cos^2 values will be colored in "blue" variables with high \cos^2 values will be colored in red. We may observe a clear division in two groups for the individuals (genes). From Figure 2 when clustering in 2 clusters we observe a clear division of genes but when increasing to 3 clusters the clusters are not clearly divided (more graphs are in R Markdown file. See Appendix).

3. Variable contribution

Variables that are correlated with PC1 (i.e., Dim.1) and PC2 (i.e., Dim.2) are the most important in explaining the variability in the data set. Variables that do not correlated with any PC or correlated with the last dimensions are variables with low contribution and might be removed to simplify the overall analysis.

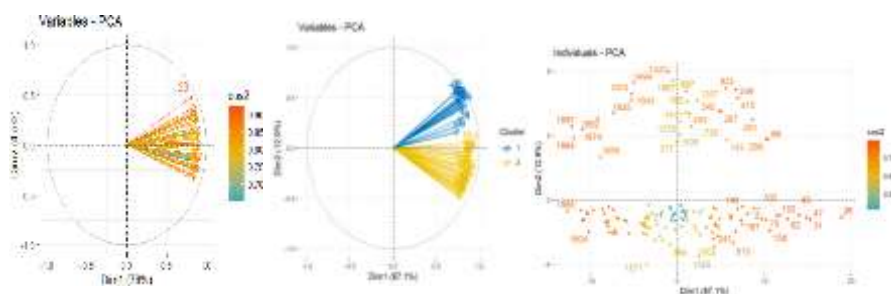


Figure 15: PCA variable and individual contribution to dimensions

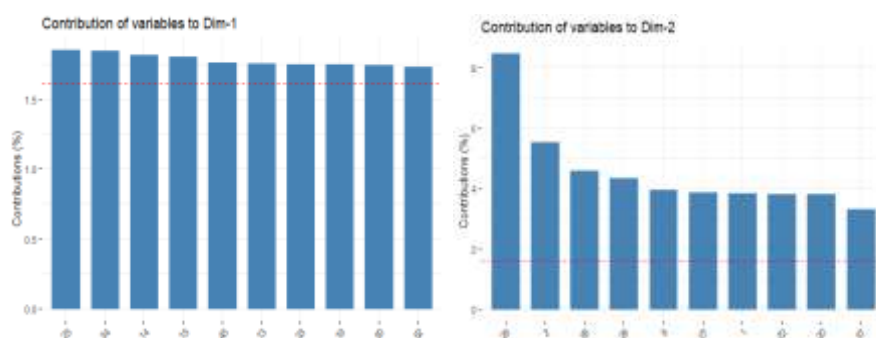


Figure 16: Top 10 contribution to Dim1 and Dim2

The larger the value of the contribution, the more the variable (sample) contributes to the component.

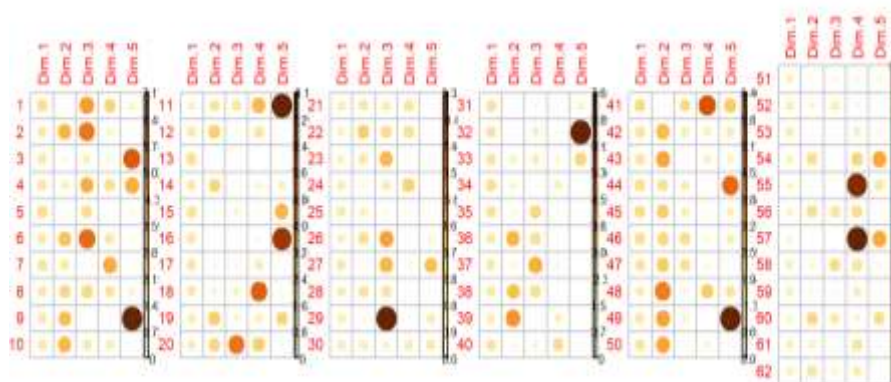


Figure 17: PCA variable contribution and first 5 dimensions

As we may observe from the contribution plots (Figure 17) there are different samples which contribute to Dim1 and Dim 2. What is clearly observed is the fact that for Dim 1 and Dim2 the number of genes contributing to is larger compared to other dimensions. This contribution may be also seen by a barplot of contributions only to Dim 1, only to Dim 2 and Dim1 and Dim2 together (For more see Appendix) To help understanding the above graph below is shown the bar graph for the contribution of samples to each of the two dimensions (Dim1 and Dim2). Figure 17 shows the contribution of samples (variables) to each of the first 5 dimensions of PCA.

This result will be used in next section to create a subset of the data and then use it further to increase accuracy in cluster analysis and prediction. The color and the size of the circle shows the contribution of the variable (row) to the dimension (column). The dark color shows a high contribution and the light color shows a lower contribution of the variable to the dimension. Note also that, the function `dimdesc()` [in FactoMineR], for dimension description, can be used to identify the

most significantly associated variables with a given principal component. So, the first 2 dimensions will explain up to 80% of the variance of the data.

4. Cluster Analysis

The distance matrix is a key element when working with clusters since it helps to identify the observations which are close to each other and group them into one cluster. In the previous parts we have used the *heatmap* which is also based in one cluster methodology which is the hierarchical cluster. The distance matrix shows some clear patterns of genes which are close and it also shows a division of genes/samples starting from 2 up to 4 clusters (Figure 18).

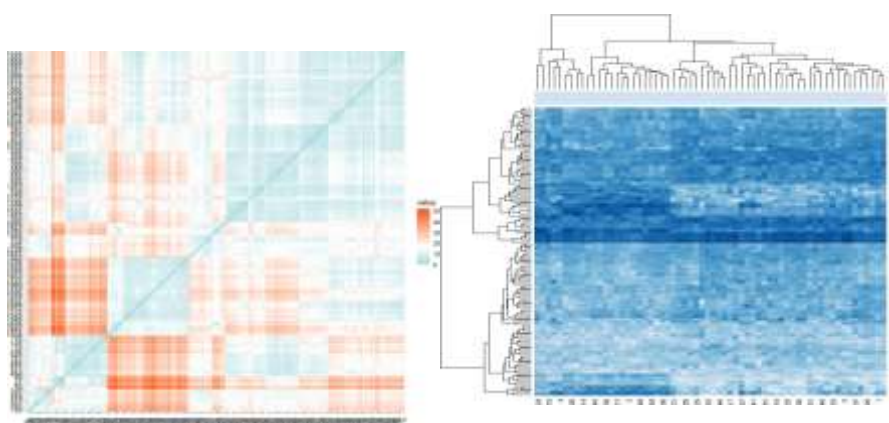


Figure 18: Distance matrix and clusters (log2 transformed data)

At this part we have considered two main cluster techniques: *Hierarchical* and *kmeans*.

4.1 Hierarchical clustering

If we don't know how to reason and decide on the number of clusters, we may use an alternative methodology which is Hierarchical Clustering. Hierarchical Clustering has an added advantage that it produces a tree based representation of the observations, called a Dendrogram. To choose the number of clusters we just draw horizontal lines across the dendrogram. We can form any number of clusters depending on where we draw the break point. Implementing hierarchical clustering involves one obvious issue. How do we define the dissimilarity, or linkage? Some of the options are:

Complete Linkage : Largest distance between observations

Single Linkage : Smallest distance between observations

Average Linkage : Average distance between observations

Centroid: distance between centroids of the observations.

We may create a dendrogram and then also visualize the clusters. Hierarchical clustering can be divided into two main types: agglomerative and divisive. Agglomerative clustering: It's also known as AGNES (Agglomerative Nesting). It works in a bottom-up way. That is, each object/observation is initially considered as a single element cluster (leaf). At each step of the algorithm, the two clusters which are the most similar are combined into a new bigger cluster (nodes). This procedure is iterated until all points are member of just one single big cluster (root). The result is a tree which can be plotted as a dendrogram.

Divisive hierarchical clustering: It's also known as DIANA (Divise Analysis) and it works in a top-down approach. The algorithm is an inverse order of AGNES. It begins with the root, in which all objects/observations are included in a single cluster. At each step of iteration, the most heterogeneous cluster is divided into two. The process is iterated until all objects are in their own cluster (see figure below). These two functions behave very similarly; however, with the *agnes* function you can also get the *agglomerative coefficient*, which measures the amount of clustering structure found (values closer to 1 suggest strong clustering structure).

As we can see from Figure 19 the *ac* value obtained from *agnes* is 0.76 which is a value close to 1 and also *dc* value obtained from *diana* is very close to 1 (0.88) thus we can say the clusters are accurately formed and are valid. And also the visualizations suggest: Agnes has 2 to 4 clusters where the differences are observed. Diana has 2 clear clusters and up to 6 very small clusters. No matter which approaches to take either DIANA or AGNES both will give you clusters with same meaning just the number of entries in the clusters might differ (Figure 20). As from the above results a number of cluster from 2 up to 4 will give accurate results in gene classification for the dataset considered. Further we will also observe suggestions from the kmeans methodology and advice how to proceed.

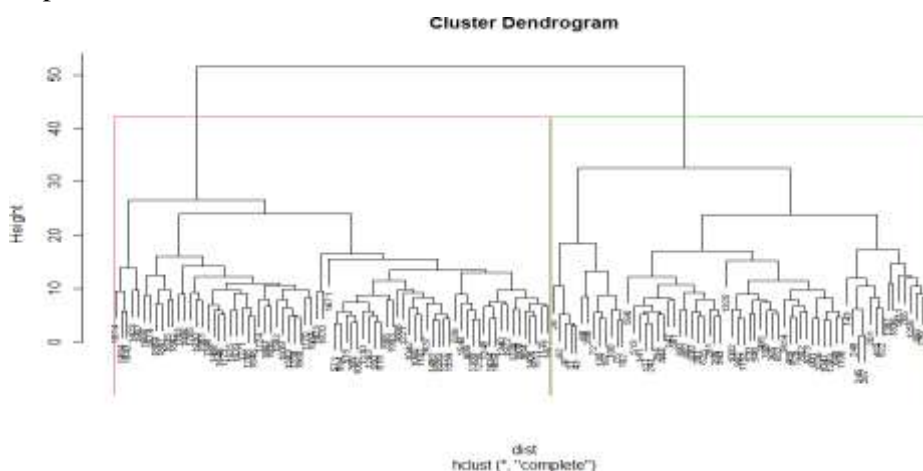


Figure 19: Hierarchical clustering dendrogram with 2 clusters

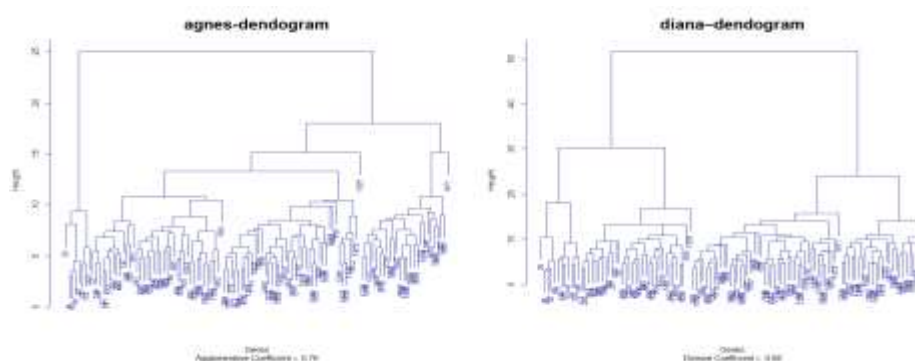


Figure 20: Distance matrix and clusters (log2 transformed data)

4.2 Kmeans

Determining Optimal Clusters is one of the challenges when using *kmeans*. To help taking a decision we may use three most popular methods for determining the optimal clusters, which includes: *Elbow method*, *Silhouette method* and *Gap statistic*. From a start visualization we may observe that our genes are clearly organized in different clusters. And using a number of clusters from 2 up to 5 the clusters remain clear and we do not have any cross location (overlapping) of the clusters. (Figure 21) But for a better decision we have also considered the below methods to decide on the optimal number of clusters. Based on Figure 22 results for the optimal number of clusters we may advise: The *WSS* suggest using 2 or 3 clusters, the *silhouette* suggest 2 clusters and the *Gapstatistic* method suggest 1 or 2 cluster. Analysing also the *kmeans* visualizations above we may suggest at this step that 2 clusters will perform good in classification of the genes.

End note cluster: Clustering is an unsupervised machine learning algorithm in which we compute analytics mostly without a pre-defined aim to understand the relationships between the data. Once we get the understanding and trends in the data we can accordingly take necessary actions and data-driven decisions. Both methodologies suggested here have their advantages and disadvantages. So, a step by step reasoning and further investigation of the data will help on deciding the optimal number of clusters used for a “best” classification of observations.

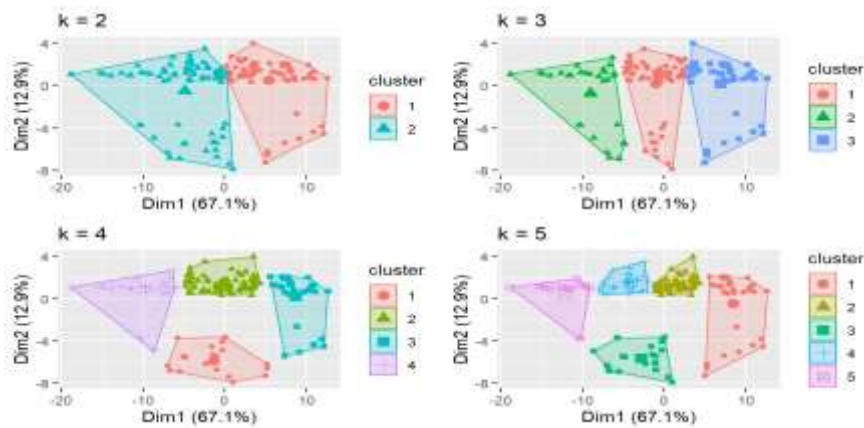


Figure 21: Kmeans for 2:5 clusters

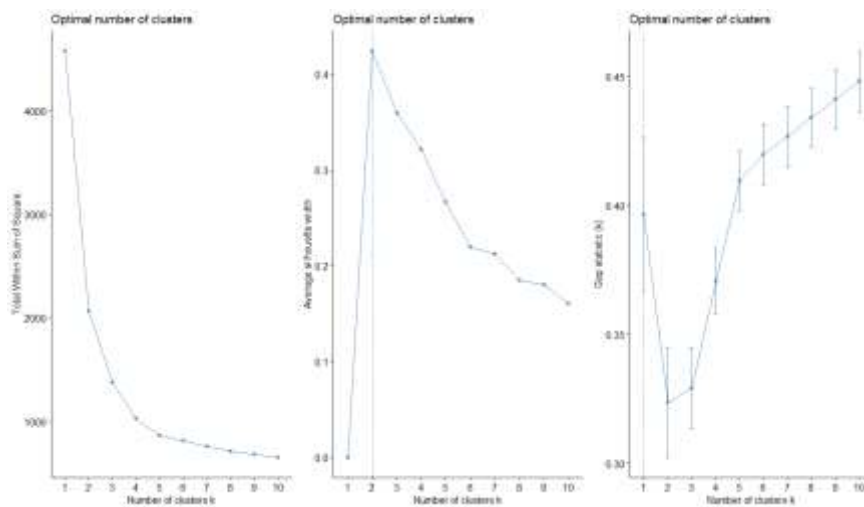


Figure 22: Optimal number of clusters from FDR subset

5. Different approach

The high amount of data in genomic dataset is extremely important especially in the pre- processing phase when we want to extract as much information as we can. But the high dimension of information is not always good. It happens that even when using machine learning methods at a certain point, increasing the number of features or dimensions can decrease the model accuracy.

Depending on the problem and the nature of the data it is possible to increase the clustering accuracy through dimension reduction or variable selection. Sometimes an ensemble approach may help to achieve a better clustering performance.

PCA- Principal Component Analysis

Dimensionality reduction is defined as way to reduce the complexity of a model and avoid over-fitting. In our case we may use a feature selection approach which is done by selecting a subset of the original features (in our case samples). Based on the above results of PCA we saw that it suffices two components to explain more than 80% of variation in data. So, we may select a subset of the features (samples) (observe correlation plot of Dim1 to Dim 5 and 62 samples). By selecting those samples which contribute mostly to Dim1 and Dim 2 we may then start analyzing again the clustering approaches. Let's bring at our attention Figure 22. In this visualization we may select those samples who are contributing more to the first and second dimensions (which together explain more than 80% of the variation). Returning to PCA results from *factoextra* and *FactoMineR* packages we may obtain the correlations of variable and dimensions.

Bring to your attention Figure 17, which shows the results of correlation and link between variables and dimensions. Based on these values we may proceed and construct a subset based on a given threshold of the correlation. In this case we have considered correlation coefficient greater than 0.85. And the sample ID where this correlation is ≥ 0.85 is compound of 18 samples. Which are almost 30% of the total number of samples (62 samples). So, with 30% of the samples we aim to achieve a higher accuracy for clustering genes in our dataset. (Figure 23)

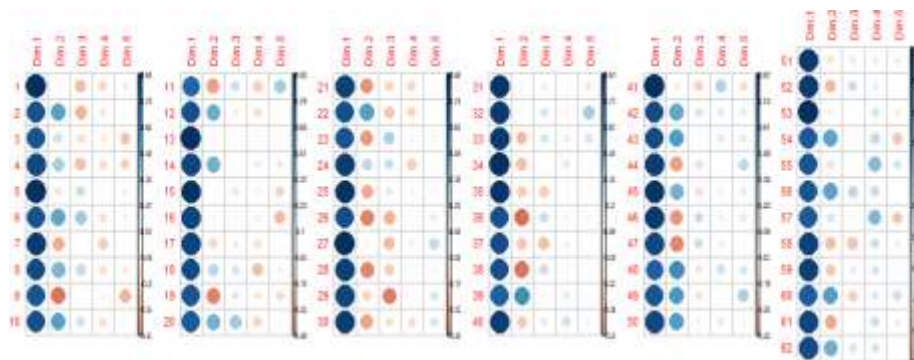


Figure 23: Variable correlation from PCA (first 5 Dim)

5.1 Basic test for normality assumption

Below is the graphical output for the normality assumption in our reduced subset.

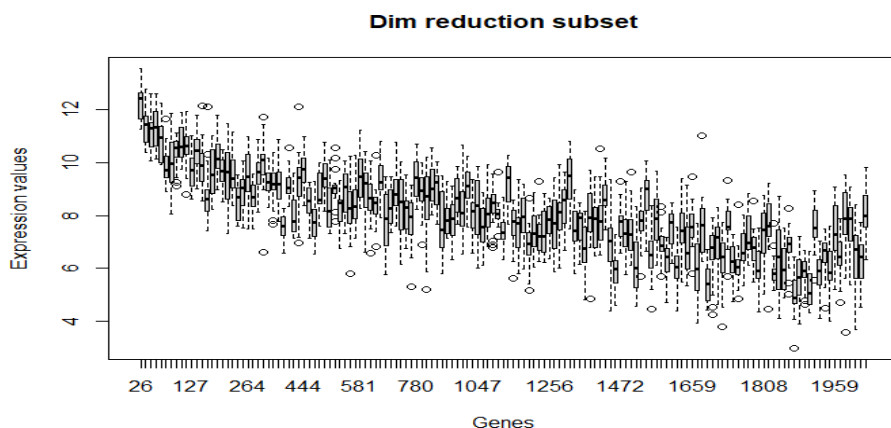


Figure 24: Boxplot after dimension reduction

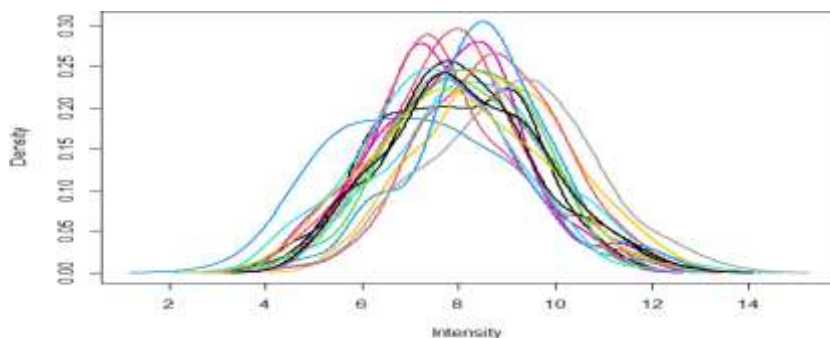


Figure 25: Density plot after dimension reduction

5.2 Hierarchical and Kmeans Clusters

The results after dimension reduction for diana and agnes approach are shown in Figure 26. The coefficients respectively for diana and agnes are increased (from 0.75 to 0.79 agnes and from 0.87 to 0.91 for diana). Observing the clustering process for *kmeans* were PCA is also visualized we also have an improvement of Dim 1 from 67% to 80.6%. Visually comparing our first try of clustering using *kmeans* (Figure 21) and after dimensional reduction (Figure 27) we do not have a significant change in the performance, since our individuals (genes) were correctly classified to a cluster in both steps.

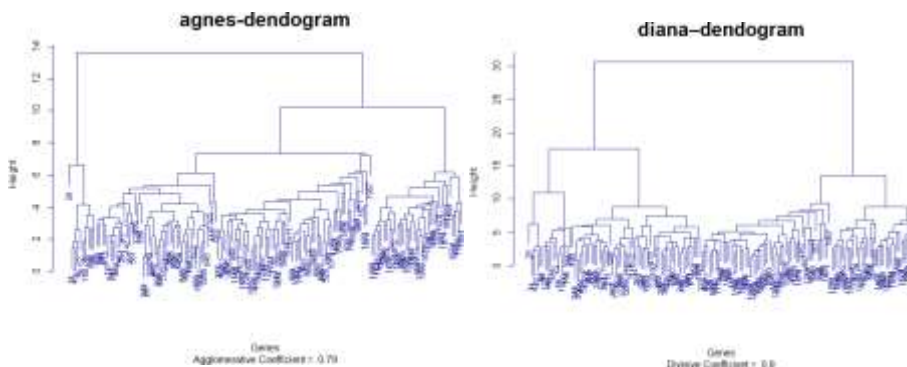


Figure 26: Diana and Agnes after dimension reduction

Dimension reduction is an important step in cluster analysis. Apart from making the high dimensional data addressable it also reduces the computational cost, and can also provide users with a clearer visualization of the data of interest. Below is the final visualization of the heatmap and clusters for both genes and samples. These results may be used further as a training set for all the data and classify each gene in the appropriate cluster. Even after dimension reduction we notice two clusters of individuals (genes) which contribute to the explanation of the variance (Figure 29). Size of the circles shows the contribution value and the color shows the value of \cos^2 . To the left and right we observe the creation of two clouds.

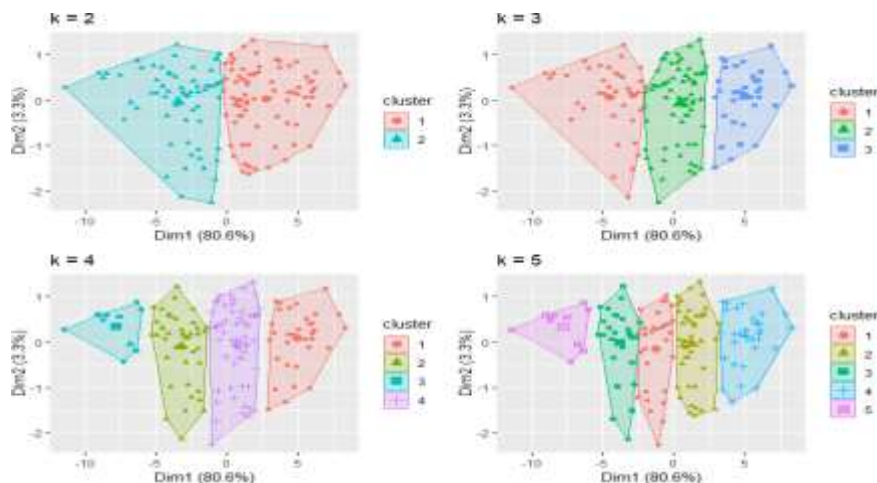


Figure 27: K-means after dimensionality reduction

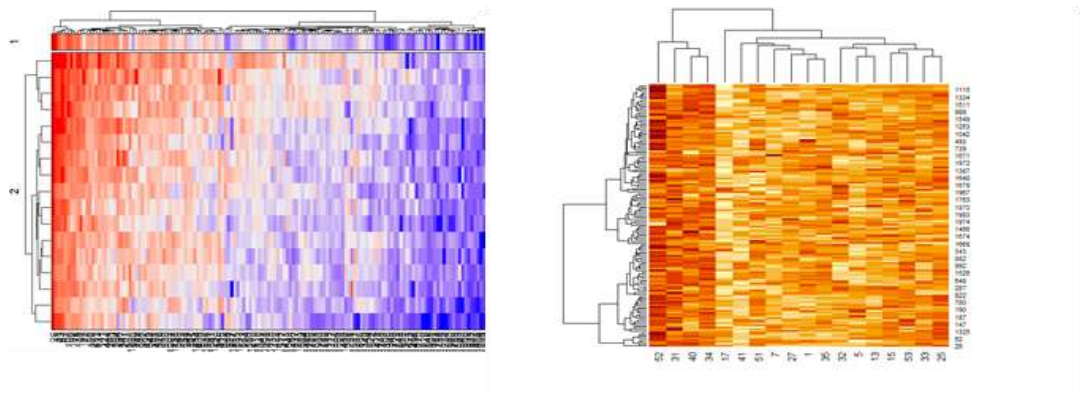


Figure 28: Heatmap after dimension reduction

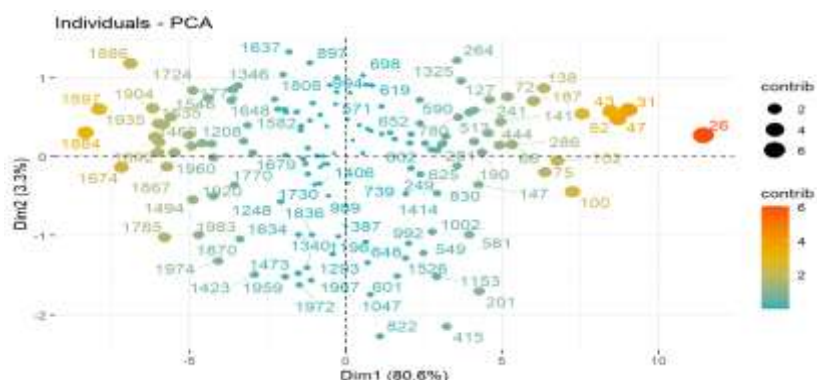


Figure 29: Individuals contribution after dimension reduction

Conclusion

Differential expression analysis has been used as a pre processing step in gene analysis data. Distributional assumption to normality and appropriate multiple testing correction are advised to use with the main aim of not missing out the information on our data. It is very important that we care about the assumption of statistical methods and then use them to summarize the findings. In this study we have used different transformation approach (standardization, log2, z-score) to achieve a set of up and down differentially expressed genes. T

o further analyse the data we have used PCA as a tool to better understand the number of components needed to explain the variation in the data. We found that more than 80% of the variation was explained with the first two components. We also used cluster techniques for a better understanding on the analysis of potential samples using from our differentially expressed genes. Advantages of our approach are related to the dimensionality reduction especially when the amount of information is large. We showed an increase in variance explanation and cluster accuracy. The approach used in this study may be useful to other studies on genomic data and genes.

Appendix

R codes with the analysis included in this material and more may be downloaded from: R Markdown Code Output-Github

<https://github.com/EGjika/Statistical-Method-for-Genomic-Data-PCA-Application/blob/main/Assignment-1-Genomic%20-%20Copy.nb.html>

References

Alon U., Barkai N., Notterman D. A., Gish K., Ybarra S., Mack D., and Levine A. J. (1999). Microarray Databases

<http://microarray.princeton.edu/oncology/affydata/index.html>

Barbey C., Hogshead M., Schwartz A. E., Mourad N., Verma S., Lee S., Whitaker V. M., Folta K. M. (2020). The Genetics of Differential Gene Expression Related to Fruit Traits in Strawberry (*Fragaria ×ananassa*). *Frontiers in Genetics*. <https://doi.org/10.3389/fgene.2019.01317>

Chung M., Bruno V.M., Rasko D.A. et al. (2021). Best practices on the differential expression analysis of multi-species RNA-seq. *Genome Biol* 22, 121. <https://doi.org/10.1186/s13059-021-02337-8>

Crow M., Lim N., Ballouz S., Pavlidis P., Gillis J. (2019). Predictability of human differential gene expression. *Proceedings of the National Academy of Sciences*. <https://www.pnas.org/doi/full/10.1073/pnas.180297311>

Rodriguez-Esteban R., Jiang X. (2017). Differential gene expression in disease: a comparison between high-throughput studies and the literature. *BMC Med Genomics* 10, 59

<https://doi.org/10.1186/s12920-017-0293-y>

Zhenfeng W., Weixiang L., Xiufeng J., Haishuo., Hua W., Gustavo G., Max R., Lin L., Jishou R., Shan G. (2019). NormExpression: An R Package to Normalize Gene Expression Data Using Evaluated Methods. *Frontiers in Genetics*, <https://doi.org/10.3389/fgene.2019.004>