

MA321-7-SP Assignment

University of Essex

Department of Mathematics, Statistics and Actuarial Science

Abstract:

This paper provides a comprehensive overview of Applied Statistics, and also explores both unsupervised and supervised dimension reduction techniques applied to a dataset containing 2000 observed gene expression values. Unsupervised learning approaches, such as Principal Component Analysis (PCA), k-means clustering, and hierarchical clustering, are used to identify groupings in the gene expression dataset.

Furthermore, supervised learning models are then applied, such as Logistic Regression, Linear Discriminant Analysis (LDA), k-nearest Neighbours (k-NN), Random Forest, Naïve Bayes, QDA and Support Vector Machines (SVM) are employed. The specific hyperparameters for each supervised learning model are described, considering model complexity and interpretability. Resampling strategies are used to compare the performance of various machine learning models and determine the best approach.

The study looks into whether clusters formed through unsupervised learning improve the prediction performance of the best machine learning model. This detailed analysis contributes to the knowledge of gene expression patterns and helps to construct predictive models for patient classification based on gene expression data.

Table of Contents

Abstract 2

Introduction 2

Preliminary Analysis..... 3

Analysis..... 4

Discussions..... 12

Conclusion..... 12

Contribution 13

Ratings..... 14

References 13

Appendix 15

Word Count: 3,134

Introduction:

Gene expression data analysis plays a crucial role in unravelling biological mechanisms and identifying potential biomarkers for various disorders, notably cancer. In this project, our focus is to utilize both unsupervised and supervised machine learning algorithms to delve into a dataset comprising 2000 observed gene expression levels.

The gene expression dataset typically consists of measurements of gene expression levels across samples. Each row in the dataset represents a gene, and each column represents a sample (e.g., patient or experimental condition). The values in the dataset denote the expression levels of each gene in each sample, often represented as normalized intensities or counts. Our ultimate objective is to predict whether forthcoming patients are predisposed to invasive or non-invasive forms of cancer using supervised and unsupervised learning methods.

Preliminary Analysis:

Handling Missing Values:

The given sample dataset have 79 missing values. Among them, 76 were present in one row and the other row has 3 missing values. We have plotted box plots as shown in fig(1) to see the data distribution and it is observed that outliers are present in the data. Imputing values with the mean is not the correct way. So, the experiment was conducted to decide on which method to use to fill the missing values. Missing values imputed using k-NN Imputation method has given least misclassification error compared to values imputed using median method. The rows having missing values were filled with k-NN Imputation method to give the better results of the models.

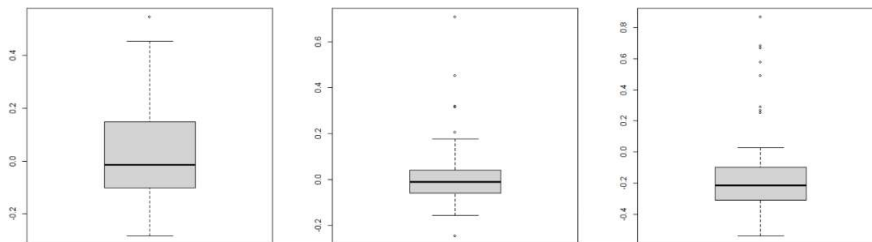


Fig.1

Analysis

We gathered the initial dataset from the provided csv file and we have set the seed value of 2312122 and we have considered top 2000 variables of the given data. We have replaced the class value having 1's with the value of zero's and 2's with the value of ones.

We have checked for the null values and we replaced the values using knn Imputation method. On the other hand we have replaced the null values using median.

Part 1

By consider supervised and unsupervised learning, dimensionality reduction of the 2000 observed gene expression data set is as follows:

Supervised Learning

Two-Sample T-test method: A t-test is a statistical analysis that compares data from two groups to determine the likelihood that the results obtained are significantly different from what is usually expected. We have used the Two-Sample T-test method, and the dataset has been reduced from 2000 to 314 columns. These maintained traits most likely record crucial information that helps distinguish between the classes. The decrease may lead to better classification performance and improved model interpretability.

Unsupervised Learning

Variance Method: A group of data points' dispersion or spread around their mean (average) is measured statistically via a concept called variance. The variance method aims to reduce a dataset's dimensions by selecting its most significant features that capture the most variance. The variance method reduced the dataset from 2000 to 56 columns having the most informative features and eliminated features with minimal variance.

R-tsne Method: t-SNE is a dimensionality reduction technique which points high-dimensional data to a lower-dimensional space using a measure like Euclidean distance to identify patterns between pairs of data. It reduces the data to two variables having the most specific features. We have used R-tsne method to perform QDA model as it requires very less number of variables for the analysis.

Part 2

We are examining groups of genes and groups of patients clustering through unsupervised learning techniques such as Principal Component Analysis (PCA), k-means clustering, and hierarchical clustering. The aim is to detect clusters/groups of genes and patients by analysing gene expression data. The process includes data preprocessing, dimensionality reduction via PCA, and clustering using k-means and hierarchical approaches.

We have performed all the techniques on both patients and genes.

In PCA analysis for patients, we use the dataset we got using the variance method. For groups of genes, we use the dataset we got from the two-sample t-test method. For others, like k clustering and hierarchical clustering for patients and genes, we use the dataset we got using the two-sample t-test method.

We are performing Principal Component Analysis (PCA) on gene/patient expression data, we summarised the result for groups of patients having 26 components showing cumulative proportion, which can explain the data variability to 90% as shown in fig(b) and for groups of genes having 37 components showing cumulative proportion which can explain the data variability to 90% as shown in fig(a), generating a bar plot to illustrate variance explained by each principal component. The analysis includes employing PCA to reduce data dimensionality and visualize gene clusters.

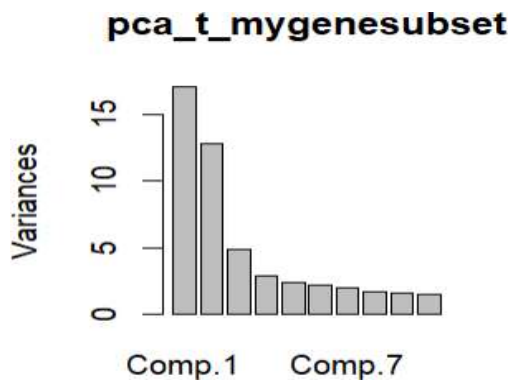


Fig.a

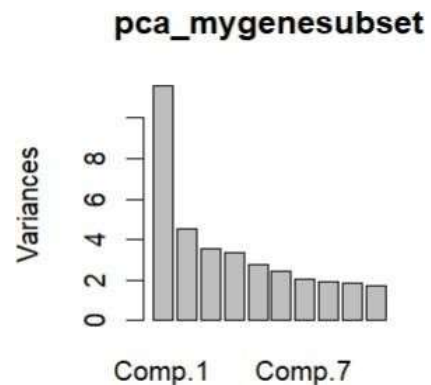


Fig.b

K-means clustering is a popular unsupervised machine learning algorithm used for partitioning a dataset into groups or clusters. This seeks to find the best arrangement of clusters based on the distribution of the data.

We are using the reduced data which we got the value from Two sample T-test method on to perform clustering on both, groups of genes and groups of patients. We scaled the data, Visualise the cluster and found the optimal value for the number of clusters to be used by using elbow method.

We got the value as 4 as shown in fig(c), which is the optimal value, and we plot the different clusters using the K-means function based on the data points. We used different colours to differentiate the clusters and plotted the graph shown in fig(d).

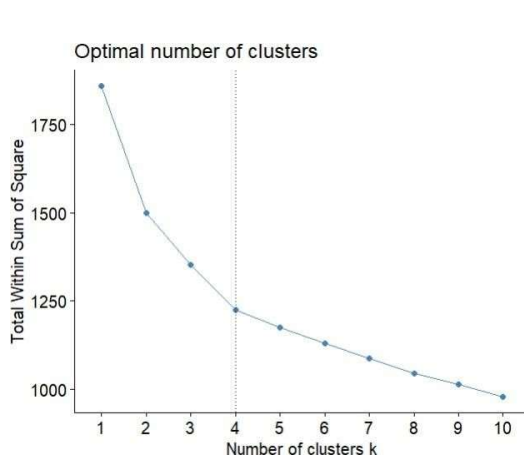


Fig.c

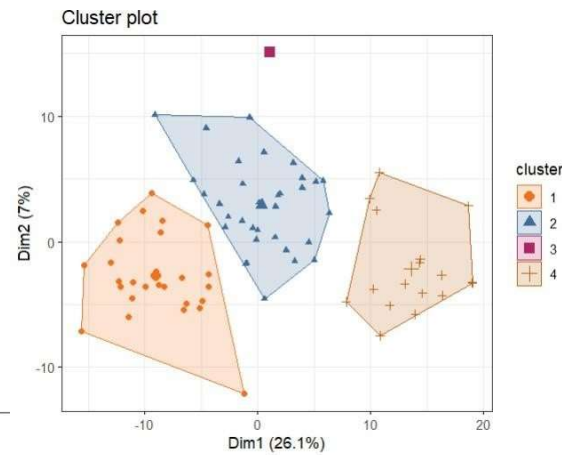


Fig.d

Hierarchical clustering is a clustering method that constructs a hierarchy of clusters by treating each data point as a separate cluster and then progressively merging the closest pairs of clusters until all data points are linked into a single cluster or a specified stopping condition is met. This technique is valuable for identifying clusters within a dataset and revealing the relationships between data points in terms of their similarities or differences.

We are executing hierarchical clustering, a method for grouping similar data points into clusters based on their distance from each other. Specifically, the code calculates and illustrates the hierarchical clustering of scaled gene data by calculating the distance matrix

using the Euclidean distance. The resulting dendrogram shown in fig(e) offers insights into the hierarchical organization of the data and how data points are grouped into clusters according to their similarities.

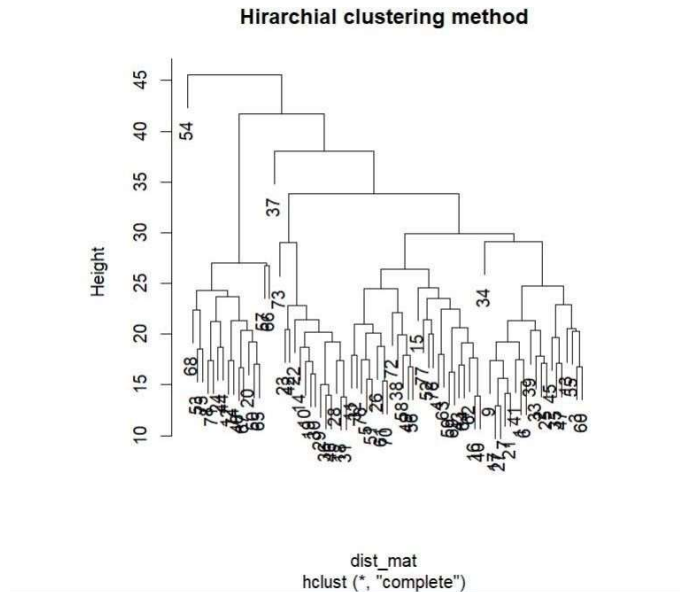


Fig.e

Part 3:

We filled the null values of the raw data using KNN Imputation and the median values, and we are using three sampling methods to predict the outcome variable.

1. We took the 2000 samples of data directly to predict the outcome variable through different models. We split the data into 70% train data and 30% test data and performed different models to predict the outcome variable.
2. We reduced the raw data using the t-test sampling method, by which we received 314 columns and 78 rows of data and data and performed different models to predict the outcome variable.
3. We used the re-sampling method for the t-test reduced data and performed different models to predict the outcome variable.

Logistic Regression:

Logistic regression is a statistical method to analyze data, predict the value of one variable based on another, and has limited outcomes. We are considering a binary logistic regression

where there is a single binary dependent variable where there are two values labelled '0' and '1', and various independent variables which are having real values (continuous values).

We used t-test sampling data to train a logistic model and generate predictions. We converted the predictions to numeric data, assigned class labels, and constructed a confusion matrix. The misclassification error rate was 0.65, which was consistent for both KNN imputed and median data.

We used K-fold cross-validation with 5 re-sampling, dividing data into 4 training sets and 1 test set for validation. We stored misclassification errors and calculated their mean. KNN Imputed data had a mean error value of 0.3, while median data had a mean error value of 0.36.

Finally, we considered the data having 2000 columns and 78 rows and performed the model to calculate the misclassification. 0.61 is the value that we received for misclassification error in sampling and 0.44 for re-sampling and we are getting the same error for both KNN Imputed data as well as for median data.

Linear Discriminant Analysis (LDA):

LDA is a statistical technique that identifies a linear combination of features to separate classes of objects or events. It helps to find the most important factors that predict the outcome of interest. We used t-test sampling method to reduce the data. Then, we trained the LDA model with the train data, made predictions using the test data and created confusion matrix. The misclassification error was 0.42, which was the same for both KNN imputed data and median data.

We performed K-fold cross-validation with $K=5$ to train LDA models. The dataset was divided into four train sets and one test set for resampling. We stored misclassification errors and calculated the mean value. Mean error for KNN imputed data was 0.26, and for median data, it was 0.28.

Finally, we considered the data having 2000 columns and 78 rows and performed the model to calculate the misclassification. 0.5 is the value that we received for misclassification error in sampling and 0.33 for re-sampling and we are getting the same error for both KNN Imputed data as well as for median data.

Random Forest:

Random forest is a statistical method that clusters data into functional groups and calculates the probability of a data point belonging to a group based on the number of trees grown and variables used at each split. We used t-test sampling method on data and trained random forest model with train data. Predictions were made on test data which was converted to numeric data. We assigned class labels '0' and '1' and calculated 0.26 misclassification error for a tree of 3. Interestingly, we found the same error for KNN imputed and median data.

We used K-fold cross-validation with $K=5$ to resample the data and perform Random Forest. We calculated the mean misclassification errors, which were 0.29 for KNN Imputed data and 0.26 for median data, both for $\text{tree}=50$.

Finally, we considered the data having 2000 columns and 78 rows and performed the model to calculate the misclassification. 0.34 for $\text{tree} = 8$ is the value that we received for misclassification error in sampling and 0.36 for $\text{tree} = 50$ for re-sampling and we are getting the same error for both KNN Imputed data as well as for median data.

Support Vector Machine (SVM):

SVM is a supervised learning algorithm used for classification and regression tasks. Its main objective is to find the best line or decision boundary to separate data points belonging to different classes. We used t-test sampling to reduce the data, then trained an SVM model with the training data and used it to predict the test data. The predictions were converted to numeric data and assigned class labels. We created a confusion matrix to calculate the misclassification error, which was 0.42 for both KNN imputed data and median data.

We used a K-fold cross-validation with $K=5$ to perform SVM on the dataset. The mean misclassification error value was calculated and found to be 0.25, which was the same for both KNN imputed and median data.

Finally, we considered the data having 2000 columns and 78 rows and performed the model to calculate the misclassification. 0.57 is the value that we received for misclassification error in sampling and 0.33 for re-sampling and we are getting the same error for both KNN Imputed data as well as for median data.

Naïve Bayes:

Naive Bayes classifier is a machine learning model that calculates the probability of an input belonging to a particular class, assuming independence between features. We used t-test

reduced data and Naïve Bayes model to make predictions on test data. Numeric functions were used to convert predictions into numeric data, and class labels '0' and '1' were assigned based on threshold values. The misclassification error calculated from the confusion matrix was 0.42, which was the same for both KNN imputed and median data.

We used K-fold cross-validation with $K=5$ to perform Naïve Bayes model, storing misclassification errors and considering the mean value. The error mean value is 0.21 for KNN Imputed data and 0.29 for median data.

Finally, we considered the data having 2000 columns and 78 rows and performed the model to calculate the misclassification. 0.42 is the value that we received for misclassification error in sampling and 0.35 for re-sampling and we are getting the same error for both KNN Imputed data as well as for median data.

K-Nearest Neighbors (k-NN):

KNN is a non-parametric, supervised learning classifier that uses proximity to predict the grouping of a data point. Initially, we used the data which was reduced by t-test sampling method. We loaded the necessary library for performing the KNN model and trained the model using train data and we created confusion matrix using the numeric data and calculated the misclassification error which is 0.19 for $k=6$ in KNN imputed data and 0.26 for $k=5$ in median data.

We used K-fold cross-validation ($K=5$) to resample the dataset and perform KNN model. We performed the model and created the list for storing misclassification errors. The mean misclassification errors for KNN Imputed data were 0.28 ($k=6$) and for median data was 0.21 ($k=6$).

Finally, we considered the data having 2000 columns and 78 rows and performed the model to calculate the misclassification. 0.42 at k value 2 is the value that we received for misclassification error in sampling and 0.34 at k value 24 for re-sampling and we are getting the same error for both KNN Imputed data as well as for median data.

By considering the table 1, the misclassification error that we received is from KNN model in which 0.19 at $k=6$ for KNN Imputed data having two sample T-test and 0.21 at $k=6$ for median data having K-fold resampling of two sample T-test.

KNN Imputed Values				
Misclassification Errors				
Models	Two Sample T-test	K-fold Re-Samples of t-test	2000 Gene Samples	K-fold Re-Samples of 2000 Gene
Logistic Regression	0.65	0.3	0.61	0.44
LDA	0.42	0.26	0.5	0.33
Random Forest	0.26, tree = 3	0.29, tree = 50	0.34, tree = 8	0.36, tree=50
SVM	0.42	0.25	0.57	0.33
Naïve Bayes	0.42	0.21	0.42	0.35
KNN	0.19, k=6	0.28, k=6	0.42, k=2	0.34, k=24
Median Values				
Misclassification Errors				
Models	Two sample T-test	K-fold Re-Samples of t-test	2000 Gene Samples	K-fold Re-Samples of 2000 Gene
Logistic Regression	0.65	0.36	0.61	0.44
LDA	0.42	0.28	0.5	0.33
Random Forest	0.26, tree = 3	0.26, tree = 50	0.34, tree = 8	0.36, tree = 50
SVM	0.42	0.25	0.57	0.33
Naïve Bayes	0.42	0.29	0.42	0.35
KNN	0.26, k=5	0.21, k=6	0.42, k=2	0.34, k=24

Table 1

Quadratic Discriminant Analysis (QDA):

Quadratic Discriminant Analysis (QDA) assumes that each class follows a Gaussian distribution. It is a generative model, and the class-specific prior is simply the proportion of data points that belong to the class. We performed R-tsne dimensionality reduction method and we used R-tsne reduced data to operate the QDA model and make predictions on test data. The misclassification error calculated from the confusion matrix was 0.53, which was the same for both KNN imputed and median data.

We used K-fold cross-validation with K=5 to perform QDA model, storing misclassification errors and considering the mean value. The error mean value is 0.28 for KNN Imputed data and 0.33 for median data.

Misclassification Errors				
Model	KNN Imputed Values		Median	
	R-tsne Method	K-fold Re-Samples of R-tsne	R-tsne Method	K-fold Re-Samples of R-tsne
QDA	0.53	0.28, k=5	0.53	0.33, k=5

Table 2

Part 4:

Apart from all the models we used, we considered k-NN to be the best model from part 3, which gives the least misclassification error. Based on the results obtained from both K-means and hierarchical clustering of genes and patients, we used these clusters in the k-NN model. We observed that at k=7, the misclassification error was the same for both clusters,

which was 0.22. By fitting the clusters for improvising the model, there is not much difference in the misclassification error compared to the k-NN model done in part 3.

Discussions

We have used the original data Initial dataset which contains 4949 variables. We have taken subset of 2000 variables and assigned 0 for all values of 1 and 1 for all vales of 2. Later we used k-NN Imputation and Median to fill the null values. We performed Two sample t-test which is supervised dimensionality reduction method which resulted having 314 variables. Also, we have used variance method which is unsupervised dimensionality reduction method which reduced the dataset to 56 columns. In addition to this, we used R-tsne method which is also unsupervised learning for only QDA analysis which reduced the data to only 2 columns, i.e. x1 and x2.

We have used unsupervised learning methods like PCA, k-means clustering, and hierarchical clustering to a gene expression dataset to identify clusters of genes and patients. To identify groups of patients we have used the dataset generated by variance method to perform PCA analysis while to identify groups of genes we have used the two-sample t-test generated dataset. In the next part, we have used supervised learning methods such as Logistic regression, random forest, k-NN, LDA, QDA, Naïve bayes, SVM models to find the misclassification errors. Compared to the sampling method resampling methods were giving better results by producing less misclassification errors as we were using k fold cross validation techniques. Among all the models, k-NN model gives the best result having least misclassification errors. k-NN imputed values performs better compared to the imputation of median values.

Conclusion

Considering all the results after the analysis of all the models, we have come to the conclusion that k-NN provides the best result having least misclassification error from the k-NN Imputed values. So, we fitted cluster to k-NN model to check whether we can improvise the model. But, in our case, there was not much difference in the misclassification error. Hence, we came to the conclusion that k-NN model gives the best outcome with or without fitting the clusters.

References

- [1]. <https://www.datacamp.com/tutorial/introduction-t-sne>
- [2]. https://www.researchgate.net/publication/228657549_Dimensionality_Reduction_A_Comparative_Review
- [3]. J. C. Avendano, L. D. Otero and C. Otero, "Application of Statistical Machine Learning Algorithms for Classification of Bridge Deformation Data Sets," 2021 IEEE International Systems Conference (SysCon), Vancouver, BC, Canada, 2021.
- [4]. J. Oyelade et al., "Data Clustering: Algorithms and Its Applications," 2019 19th International Conference on Computational Science and Its Applications (ICCSA), St. Petersburg, Russia, 2019.
- [5]. K. G. Nisha and K. Sreekumar, "A review and analysis of machine learning and statistical approaches for prediction," 2017 International Conference on Inventive Communication and Computational Technologies (ICICCT), Coimbatore, India, 2017.

