# A Comparative Analysis of Feature Extraction Methods for Classifying Colon Cancer Microarray Data

M.O Arowolo[1,*], R.M. Isiaka[1], S.O. Abdulsalam[1], Y.K. Saheed[2] and K.A. Gbolagade[1]

[1,] Department of Computer Science, College of Information and Communication technology, Kwara State University, Malete, Nigeria.
[2,] Department of Physical Sciences, Al-Hikmah University, Ilorin, Kwara State, Nigeria.

## Abstract

Feature extraction is a proficient method for reducing dimensions in the analysis and prediction of cancer classification. Microarray procedure has shown great importance in fetching informative genes that needs enhancement in diagnosis. Microarray data is a challenging task due to high dimensional-low sample dataset with a lot of noisy or irrelevant genes and missing data. In this paper, a comparative study to demonstrate the effectiveness of feature extraction as a dimensionality reduction process is proposed, and concludes by investigating the most efficient approach that can be used to enhance classification of microarray. Principal Component Analysis (PCA) as an unsupervised technique and Partial Least Square (PLS) as a supervised technique are considered, Support Vector Machine (SVM) classifier were applied on the dataset. The overall result shows that PLS algorithm provides an improved performance of about 95.2% accuracy compared to PCA algorithms.

## 1. Introduction

Dimensionality reduction is a very helpful, important and necessary tool in the expression of microarray datasets. It endeavours to trim down, recognize and illustrates the collection of unified datasets by transforming a high-dimensional dataset into a lower dimensional dataset which signifies the most significant variables that triggers the distinctive data. This significant and essential tool attracts numerous researchers working in the aspect of bioinformatics and deals with gene expression datasets to work on the dimensionality reduction [1], [2]. Several methods for dimension reduction exist, but none is confirmed the best method for all circumstances due to the reason that at the time of the processing some information is lost [3]. To improve the performance of dataset, feature extraction as a universal method of dimension reduction is considered. In Feature extraction

method, the original high-dimensional feature space is estimated on to low-dimensional feature space, for typical microarray data analysis the training sample size is always limited. Due to classification algorithms may be short of efficiency or even fail in high dimensional microarray data analysis, dimension reduction is a good choice to variable selection in order to overcome the dimensionality problem; it uses a little quantity of features to substitute a feature subset containing well-built correlations in the original data [4]. Quite a lot of feature extraction algorithms and techniques have been proposed in literature, one of the most popular and widely used techniques is PCA [1]. PCA is an unsupervised method and an effective tool but it is not efficient for high dimensional and complex dataset, due to the fact that it cannot retrieve precisely the true latent variables of complex and supervised datasets [5], data in a very high dimensional space often exists in a lower dimension. With this kind of data, the intrinsic supervised structure could

---

[*]Corresponding author. Email:olliray2002@yahoo.com

not be found through an unsupervised feature extraction technique. Another drawback of PCA is that the size of the covariance matrix is proportional to the dimensionality of the data-points. In order to overcome the drawback of unsupervised feature extraction in a very high dimensional dataset, several supervised feature extraction methods have been developed. An improvement of supervised algorithm is Local Linear Embedding (LLE) [1], it is efficient and powerful for dimensionality reduction among the other algorithms [5], [6], [7]. Local Tangent Space Analysis (LTSA) is another nonlinear dimensionality reduction technique that describes local properties of the high-dimensional data using the local tangent space of each data point [8]. These techniques have been successfully applied on microarray data.

In this paper, a supervised feature extraction algorithm for dimensionality reduction is proposed to handle the curse of dimensionality of microarray data. PLS is a proposed algorithm for supervised feature extraction. The experiments show PLS outperforms PCA in reducing the dimension of supervised structures and visualization performance. This paper is organized as follows. Section 2 deals with the related work of dimensionality reduction for classification of microarray data. Section 3 describes the dataset used, methodology and algorithms. Section 4 deals with the discussion and results. Section 5 concludes the work.

## 2. Related Works

This Author Jian, Linh, and David [9] presented a dimension reduction models using PLS, SIR and PCA, the comparative performance of their classification procedures were similar to PCA and PLS, the complexity of microarray data analysis was reduced. PLS and SIR were both expensive in dimension reduction but extra effective than PCA, and the results are reliable with the scrutiny of the method.

In 2012 [10] carried out normalization to regulate all the features in the dataset and dimensionality reduction to carry out clustering. The diabetic dataset which contains 768 instances and 8 attributes has been taken and PCA algorithm is used to reduce the dimensions. Out of 8 features, 4 features are selected without the loss of information.WEKA3.7 tool is used to investigate the diabetes data. After performing dimensionality reduction density based clustering algorithm is used to find the maximal set of density. Dimensionality reduction is used to increase the accuracy of the clustered data.

In 2008 [11] has compared two different feature extraction algorithms. The features of the products based on the review of the customer, is considered as the dataset. In the first algorithm the candidate features are identified and they are pruned. In the second approach association rule mining is used to find the frequent pattern. Here, the dataset is based on the customer review which are collected from the social website such as

amazon,cnet and it is based on five different products(two digital cameras, a DVD player, an MP3 player and a cell phone). Likelihood Ratio Test is the method used to extract the features of the product.

In 2016 [12], has presented the methods for visual data mining in order to mine the data and to make cognitive. Here, the author has performed the attribute selection method i.e. wrapper method and filter method. Here, seven types of dataset (lung cancer, promoter, sonar, Arrhythmia, Colon Tumour, and Central Nervous System) have been used and the accuracy has been calculated before reduction and after reduction. In this reduction framework, the numbers of attributes have been reduced. The data visualization is represented in order to determine the relationship between the data. The algorithms like (LDA, QDA, and KNN) have been used in this work and found that LDA have performed efficiently and reduces the attributes effectively.

In 2014 [16] reviewed numerous development application to help users implement feature extraction of gene expression data, the paper presented review of software for feature extraction methods such as PCA, ICA, PLA and LLE. The software applications have limitations in terms of computational performance and there is need for development of classification methods to improve performances of these feature extraction methods.

2015 [16] compared dimension reduction based on logic regression models for the case-control genome-wide association by employing PCA and PLS, there were limitations in the interaction of the genes of dataset used affecting the goodness of fit and accuracy of the parameter estimation of PLS and needed further investigations.

## 3. Datasets Used

Colon cancer dataset was used for this experiment, it contains an expression of 2000 genes with highest minimal intensity across 62 tissues, derived from 40 tumour and 22 normal colon tissue samples [13].The gene expression was analyzed with an Affymetrix oligonucleotide array complementary to more than 6,500 human genes. The gene intensity has been derived from about 20 feature pairs that correspond to the gene on the DNA microarray chip by using a filtering process. Details for data collection methods and procedures are described in [13], and the data set is available from the website http://microarray.princeton.edu/oncology/.

## 3.1. Methods and Algorithms

Dimensionality reduction is an important factor used in reducing original data features without the loss of information. The main objective of this analysis work is to compare two different feature extraction algorithms namely PCA and PLS. The architecture is as follows:
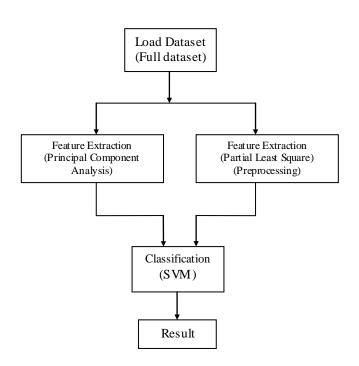
**Figure 1.** Technique Workflow

## 3.1.1. Principal Component Analysis (PCA)

In dimension reduction, PCA is one of the well-known techniques; its conception is to lessen the high dimensionality of a given dataset, while keeping enough of the variation existing in the initial predictor variables.

This is attained by transforming the $p$ initial variables $X=[x_1, x_2, \ldots, x_p]$, [16] to a latest set of $q$ predictor variables.

PCA is a widely used unsupervised feature extraction technique; it works by replacing the original variables in a data with numerical variables called principal component by capturing the most descriptive features with respect to the most relevant ones [15]. PCA mathematically transforms data by referring them to a different coordinate system in order to obtain the greatest variance. A number of correlated variables into a smaller number of uncorrelated variables called principal components [10].

PCA identifies patterns of similarities and differences in a data, these patterns are determined and can be compressed by reducing the numbers of dimensions without much loss of information.

In order to conduct the PCA analysis for the input data, the following steps are performed by adopting [16]:

- Create N x d data matrix with one row vector $x_n$ per data input.
- Subtract mean from each row vector $x_n$ in **X**
- Calculate the covariance matrix of X
- Find Eigen Vectors and Eigen values of $\sum$
- Fetch the Eigen vector with the largest Eigen values.

The PCA are uncorrelated and the components explain the largest percentage in the dimensional dataset with results

in extracting 10 components which are considered relevant in the colon cancer dataset used.

## 3.1.2. Partial Least Square (PLS)

Partial Least Square (PLS) is a supervised feature extraction technique, which is widely used as a procedure in modeling associations linking blocks of experimental variables by means of latent variable, it tries finding uncorrelated linear transformations (latent components) of the original predictor variables which have high covariance with the response variables [16]. The goal of PLS is to find the linear relationship between the response and explanatory variables y and X:

$$X = TP^T + E_x \qquad (4)$$
$$y = TC^T + E_y \qquad (5)$$

Where T represents the scores (latent variables) P and C are loadings, and E$x$ and E$y$ are the residual matrices obtained the original X and $y$ variables.

Feature extraction using PCA ignores the response variable and its equivalence. PLS integrates the response variable during the dimensionality reduction procedure. PLS outperforms PCA in the case of microarray gene expression, PLS only consists of indicating the amount of gene components whereas PCA necessitates choosing the essential gene components [17].

## 3.1.3 Support Vector Machine (SVM)

In this step, the results for classification are computed using SVM for classification. SVM is a recently developed technique used for classification suggested by Vapnik, which was consecutively applied to several domains. SVM is applied to microarray cancer data which comprises of several gene expressions. SVM is applied after many steps after analysis to finally classify cancer tissues as part of an integrated algorithm. SVM is a constructive learning procedure based on statistical knowledge theory [18], it is used for classification tasks, and it uses linear models in implementing non-linear class boundaries by transforming input space using a non-linear mapping into a new space. SVM produces an accurate classifier with less over fitting and it is robust to noise. Assuming $\{(x_1, y_1), \ldots, (x_n, y_n)\}$ be a training set with $x_{1i} \in R^d$ and $y_i$ is the corresponding target class. SVM can be reformulated as:

Maximize:

$$J = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j (x_i^T, x_j) \qquad (6)$$

Subject to;

$$\sum_{i=1}^{n} \alpha_i y_i = 0 \quad and \quad \alpha_i \geq 0, i = 1, 2, \ldots, n \qquad (7)$$

This is the weighted average of the training features. Here, αi is a Lagrange multiplier of the optimization task and αi is a rank label. Values of α′si are non zero for all the points lying inside the margin and on the correct side

of the classifier. The kernel function is used to solve the problem. The Kernel function analyses the relationship among the data and it creates a complex divisions in the space [19].

## 3.1.4  MATLAB

PLS, PCA, clustering, dimension reduction, factor analysis, visualization, and others. In the statistical toolbox of MATLAB, several PLS and PCA functions are provided for multivariate analysis. Most of these functions are used for dimensional reduction. All of these functions are implemented in MATLAB.

## 4. Results and Discussions

The colon cancer dataset extracted were classified, the classification results obtained show the features capability for classifying the colon's status. The average classification accuracy, which is using features with PCA and PLS are recorded in tabular form below.

The proposed methodology was applied to the publicly available colon cancer database.



**Figure 1:** The Loaded Colon Cancer Microarray Dataset.

**Table 1:** Result Evaluations

| Dataset | Feature Extracted (PCA) | Feature Extracted (PLS) |
|---|---|---|
| Colon Cancer (2001x62) | 10 Components | 20 Components |



**Figure 2:** PCA Features Extracted 10 Components



**Figure 3:** PLS Feature Extracted.

In this experiment, PCA as a feature extraction method is used to reduce the high-dimension and SVM is used as the classifier. PCA is used to de-correlate the data and 10 components was achieved, in Fig. 2, the overall accuracy on all the datasets obtained using PCA as feature extraction to transform and extract the dataset is reported in a confusion matrix.



**Figure 4.** Confusion Matrix of Proposed Classification, using PCA-Based for Classification

True Positive Rate 68.2% and False Negative Rate yields 90.0% .
TP=36 FP=7 FN=4 TN=15
ACCURACY: (TP + TN) / (TP+TN+FP+FN) = 82.3%
SENSITIVITY:   TP/ (TP+FN)   =   90.00
SPECIFICITY:   TN/ (FP+TN)   =   68.18
PRECISION:   TP/(TP+FP)   =   83.72

In this experiment, to evaluate the performance of colon cancer dataset PLS is used as a feature extraction method with SVM as a classifier. It obtained a better overall accuracy compared to the first and second experiment as shown in the confusion matrix of Figure. 3, it de-correlated the features of colon cancer dataset into 20-components which are considered relevant.
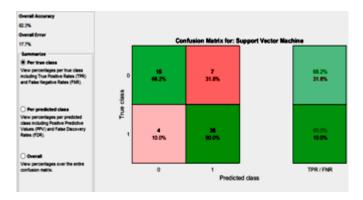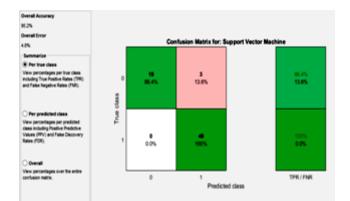
**Figure 5.** Confusion Matrix of Proposed Classification, using PLS for Classification

True Positive Rate 86.4% and False Negative Rate yields 100.0%.

TP=40 FP=3 FN=0 TN=19
ACCURACY: (TP + TN) / (TP+TN+FP+FN) = 95.16%
SENSITIVITY: TP/ (TP+FN)     =     100
SPECIFICITY:   TN/ (FP+TN)     =     86.36
PRECISION:     TP/ (TP+FP)     =     93.02

**Table II**: Performance Evaluation of Proposed PCA and PLS Methods.

| S/No | Performance Metrics | PCA Method | PLS Method |
|---|---|---|---|
| 1 | Training Time | 47.2686 | 5.0088 |
| 2 | Accuracy (%) | 82.30 | 95.16 |
| 3 | Sensitivity (%) | 90.00 | 100 |
| 4 | Specificity (%) | 68.18 | 86.36 |
| 5 | Precision (%) | 83.72 | 93.02 |
| 6 | Area Under Curve | 0.848864 | 0.994318 |
| 7 | Error (%) | 17.7 | 4.8 |

Table II illustrates a comparative chart between the three methods used in terms of several performance measures such as accuracy, sensitivity, specificity, precision, error, time and area under curve. This comparison shows the integrity of the proposed approach with respect to the state of the art. The colon cancer dataset used to generate our result achieved its best on PLS for feature extraction; it makes this method suitable for practitioners.

## Conclusion

In this paper, a widely used colon cancer datasets was used for the evaluation of the algorithms used. The dimension reduction algorithms used to eliminate high dimensional data were PCA and PLS, it uses SVM as its classifier, and it was successfully implemented on MATLAB. For the purpose of finding the smallest gene subsets for accurate cancer classification, PLS method is highly effective compared to PCA.

PLS Based method showed a better performance than PCA-based method with 95.16% to 82.30% accuracy. Hence it can be stated that the PLS based dimensionality reduction scheme is suitable for microarray gene classification as it extracts relevant and a reduced amount of information from the feature selection based technique. In future studies PLS can be compared with another feature extraction method with the aforementioned criteria. Another dataset will be a good avenue for further research of dimensionality reduction.

## References

[1] Anaissi Ali, Paul Kennedy, and Goyal Mahdu: A. Dimension Reduction of Microarray Data Based on Local Principal Component, World Academy of Science, Engineering and Technology. IJCEACIE, 5(5): 529-534.

[2] Tenenbaum Joshua, Vin de Silva, and Langford John. A Global Geometric framework For Nonlinear Dimensionality reduction. Science. 290: 2319-2323 (2009).

[3] Araujo Daniel, Neto Adriao, Martins Alan, and Melo Jorge: Comparative Study on Dimension Reduction Techniques for Cluster Analysis of Microarray Data, Proceedings of International Joint Conference on Neural Networks, San Jose, California, USA. 1(4): 1835-1842 (2011)

[4] Hua Jianping, Chao Sima, and Doughertya, Edward: Performance of feature selection methods in the classification of high-dimension data, Pattern Recognition, Volume. 42(3):409– 424 (2009).

[5] Lee John, and Verleysen Michel: Nonlinear Dimensionality Reduction, Springer Publishing Company. 8226: 617-622 (2007).

[6] Quansheng Jiang, Jiayun Lv, and Minping Jia: New approach of intelligent fault diagnosis based on LLE algorithm. Control and Decision Conference. 10: 522-526 (2008).

[7] Varini C, Nattkemper T: Breast MRI Data Analysis by LLE. Neural Networks. Proceedings. IEEE International Joint Conference. (10)3: 2449-2453 (2004).

[8] Zhang Zhenyue and Zha Hongyuan: Principal Manifolds and Nonlinear Dimensionality Reduction Via Local tangent Space Alignment. SIAM Journal of Scientific Computing. 26 (1): 313-338 (2006)

[9] Dai J, Lieu L, and Rocke D: Dimension reduction for classification with gene expression microarray data, Statistic. Appl. Genet. Mol. Biol. 5(6): 1-19 (2006)

[10] Rekha Awasthi, Anil Kumar, and Seema Pathak: An Analysis of Density Based Clustering Technique With Dimensionality Reduction for Diabetic Patient, International Journal of Computer Engineering and Applications. 9(4): 165-171 (2012)

[11] Liliana Ferreira, Niklas Jakob, and Iryna Gurevych. Comparative Study of Feature Extraction Algorithms in Customer. Proceeding of the 2nd IEEE International Conference on Semantics Computing (ICSC). 10(40): 4-7 (2008)

[12] Vijayarani S, Maria Sylviaa. Comparative Analysis of Dimensionality Reduction Techniques. International Journal of Innovative Research in Computer and Communication Engineering. 4(1): 23-29

[13] Alon U, Barkai N, Notterman D, Gish A, Ybarra S, Mack, D, and Levine, A: Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. Proc Natl myAcad Sci USA 8; 96(12): 6745-6750 (1999)

[14] Smitarani Satpathy, and Pratikshya Mahapatra: Microarray Classification Using Intelligent Techniques, International Journal of Scientific & Engineering Research. 5(7): 1663-1670 (2014).

[15] Nebu Varghese, Vinay Verghese, Prof Gayathri. and Dr. Jaisankar: A Survey of Dimensionality Reduction and Classification Methods, International Journal of Computer Science & Engineering Survey. 3(3): 45-54 (2012).

[16] Ching Siang, Wai Soon, Mohd Saberi, Weng Howe, Safaai Deris, and Zuraini Ali: A Review of Feature Extraction Software for Microarray Gene Expression Data. BioMed Research Int. 10(1155): 15: 1-15 (2014)

[17] Honggang Yi, Hongmei Wo, Yang Zhao, Ruyang Zhang, Junchen Dai, Guangfu Jin, Hongxia Ma, Tanchun Wu, Zhibin Hu, Dongxin Lin, Hongbing Shen, and Feng Chen: Comparison of Dimension Reduction-based Logistic Regression Models for Case-control Genome-wide Association Study. The Journal of Biomedical Research. 29(4):289-307 (2015).

[18] Haozhe Xie, Jie Li, Qiaosheng Zhang, and Yadong Wang: Comparison among Dimensionality Reduction Techniques Based on Random Projection for Cancer Classification. airXiv. 2: 1-11 (2017).

[19] Seung Jun, Yichao Wu, Hao Helen, and Yufeng Liu: principal Weighted Support Vector Machine for Sufficient Dimension Reduction in Classification. Biometrika. 104(1): 67-81. (2017).