# APPLICATION OF DIMENSIONALITY REDUCTION ON CLASSIFICATION OF COLON CANCER USING ICA AND K-NN ALGORITHM

**Rasheed G. Jimoh [1], Ridwan M. Yusuf [1], Yusuf O. Olatunde [1], Saheed Yakub Kayode [2]**

**[1] Computer Science Department, University of Ilorin, Ilorin, Kwara State, Nigeria**
**[2] Department of Physical Sciences, Al-Hikmah University, Ilorin, Nigeria**

Corresponding author: Yusuf O. Olatunde, hizyyusuf@gmail.com

**ABSTRACT:** Several sectors including engineering, health, academics and so on deals with very large number of information and few specimens. This highlights the need of a technique to improve data accuracy in order to enable professionals such as biologists, clinicians and so on to comprehend the structure of a complex microarray dataset and the gene expression in cells when reduced. This study employs Independent Component Analysis for feature extraction before using *k*-Nearest Neighbor algorithm to classify colon cancer dataset which contains DNA microarray gene expression data with 2000 features and 62 samples. The experiment was performed using MATLAB 2015a. The result shows that the dimensionality reduction applied improve the classification performance in terms of accuracy, sensitivity, specificity and precision by 11.3%, 25.2%, 36.3% and 12.8% respectively.

**KEYWORDS**: Microarray, k-NN, Dimensionality Reduction, ICA, Colon Cancer.

## 1. INTRODUCTION

Dimension reduction is an effective and essential tool used to analyze microarray datasets ([APM11]). A lot of algorithms and feature extraction techniques have been put forward in literature the reduction of dimensionality ([RT15]). Principal Component Analysis (PCA) is one of the most widely used and common dimensionality reduction techniques, it is seen as an unsupervised technique and relatively effective tool, but it's not considered as efficient for dataset that are complex and of high dimension ([APM11]). Therefore, there is need to address the inability of PCA to precisely retrieve the genuine latent features of complex datasets ([APM11]). Data in a very high dimensional space often exists in a lower dimensional space and unsupervised feature extraction technique such as PCA may not be totally efficient.

Independent Component Analysis (ICA) is a feature extraction method that utilizes the presence of free factors in multivariate information and breaks down an input dataset into statistically independent components ([X+05]). Independent Component Analysis can lessen the impacts of noise and is ideal for separating diverse signals. Recently, ICA have been used in extracting expression modes of genes in microarray analysis, Moreover, a recent report showed that ICA could enhance the biological legitimacy of the features measured with other methods ([LSC11]).

In this study, Independent Component Analysis (ICA) is use for feature extraction. It is a supervised feature extraction dimensionality reduction method. It choice result from the efficiency rate offered by supervised techniques compare to unsupervised techniques, because the connections among the latent features are much more effective ([LSC11]).

### 1.1 Independent Component Analysis (ICA)

Independent component analysis (ICA) is named as a computational strategy or system utilized for separating grouped flag into its diminished subcomponents. The "mixed drink party issue" can be portrayed as basic routine with regards to ICA, in which the central discourse signals are split from a specimen information including people speaking together inside a room. This situation is generally translated by considering the non-presence of time delays or echoes. A fundamental factor to be considered is that if N sources are available, at that point in any event N estimations (for example, Amplifiers) are required to mine the primal signs. ICA calculation has the ability to utilize higher request insights which may contain generous corresponding information. This is basically esteemed when the dispatch of information differentiates strikingly from its transcendent parameters, which for no less than two reasons is commonly the circumstance when managed microarray articulation information. Most ICA calculations require "brightened" information, by methods for a personality covariance network. This is more than a factual essential, wherein the calculation goes for crumbling the information past its initial two minutes which is a definitive objective of PCA and its related strategies ([C+14]).

## 1.2 *K*-Nearest Neighbor

Nearest neighbor search is among the predominant learning and classification methods as presented by ([FH51]), it has been turned out to be a straightforward yet capable acknowledgment calculation. In ([CH67]), it was demonstrated that the choice decide performs well considering that no unequivocal information of the information is accessible. *k*-NN govern is utilized to indicate the basic speculation of this technique, this is a circumstance whereby another example is characterized into the class with the individuals that are most present inside the K Nearest Neighbors, can be utilized to obtain great assessments of the Bayes blunder and its likelihood of mistake asymptotically approaches the Bayes blunder ([DH73]). There are three (3) confinements for the conventional *k*-NN content arrangement ([ST10]):

1. High calculation complexity: keeping in mind the end goal to distinguish the k closest neighbor tests, all similitudes between the preparation tests must be figured. The KNN classifier is not any more ideal when the quantity of preparing test is less. In any case, should the preparation set contain an enormous number of tests, at that point the KNN classifier needs more opportunity to ascertain the similitudes. There are three (3) approaches to take care of this issue which are diminishing the measurements of the element space, utilizing littler datasets and utilizing enhanced calculation.

2. Dependency on the training set: The classifier can be produced just with the preparation tests and it doesn't make utilization of any extra information. This influences the calculation to rely upon the preparation to set too much and recalculation is required paying little mind to whether a little change on preparing set exists

3. No weight difference between samples: All the preparation tests are dealt with similarly; no distinction exists between tests of information with modest number and tests of information with enormous number. So, it doesn't coordinate the genuine wonder where the examples have uneven appropriation ordinarily.

## 2. RELATED WORK

In 2014, Trstenjak, Mikac and Donko ([TMD14]) introduced the system of *k*-NN with TF-IDF for content order. Quality and speed of the arrangement was absolutely the reason for this system. It discovers protests that are comparable in view of the Euclidean separation and TF-IDF computes the weight for each term in each record. Both *k*-NN and TF-IDF implanted together, demonstrate great together, gave great outcomes and affirmed the

underlying desires. Structure is performed on various classifications of reports and the testing is performed. Amid testing, grouping gives exact outcomes because of *k*-NN calculation. This blend gives better outcomes and need to redesign and need to enhance the system for better and high exactness comes about.

Singh and Prakash ([SP14]) reported that consideration of scientists has been gotten by content classification in last ten (10) years with the ascent in electronic substance of the reports. Finding a specific report from the web or any vast storehouse, content or record arrangement is the most reasonable assignment. Some better framework and upgraded machine learning classifiers are requested to do the undertaking of report order. The study composed a multi-operator based framework which contains some product cross breed specialists that recovers classification of an archive and relates with each other in order to take ultimate conclusion about the classification and information is then bolstered to the machine learning classifier in order to upgrade the execution.

In 2014, Lin ([Lin14]) assessed the vitality cost of various classifiers and lessened the vitality cost by parallelization, attempting to distinguish the classifier that performs best on the two parts of adequacy and productivity. The study proposes a novel lexical way to deal with content arrangement in the bio-medicinal space. LK-NN (Lexical *k*-NN) calculation has been proposed, in which lexemes (or tokens) are utilized to speak to the therapeutic archives. These tokens are utilized to characterize abstracts by coordinating them with the standard rundown of watchwords indicated as Work (Medicinal Subject Headings). This naturally arranges diary articles of therapeutic space into particular classes. The study utilized gathering of therapeutic reports, called Ohsumed, as the test information for assessing the proposed approach. The research demonstrates that LK-NN beats the customary *k*-NN calculation as far as standard Fmeasure.

Ganiz, Tutkan and Akyokus ([GTA15]) presented another classifier for printed information, called Administered Meaning Classifier (SMC) was proposed. The new SMC classifier utilizes an importance measure, which depends on Helmholtz standard from Gestalt Hypothesis. In SMC, seriousness of the terms with regards to classes are figured and utilized for arrangement of the archive.

Perez and Marwala ([PM12]) reported that there is a need to isolate numerous accumulations of archives into comparative ones through bunching. Determining number of bunches is required in existing parceling calculations and the yield is absolutely reliant on given information. The study

presents an Enhanced Report Grouping (ERG) calculation by creating a number of bunches for any content records and uses cosine comparability measures to put comparative archives in appropriate bunches. Trial comes about demonstrated that precision of proposed calculation is high contrasted with existing calculation as far as F-Measure and time multifaceted nature.

## 3. METHODOLOGY

In this study, the method used consists of two different phase that is, Feature extraction and Classification. The dataset is loaded as an input to the feature extraction module (shown is figure 1) where ICA is used to remove redundant and irrelevant features. The output of the Feature extraction module (shown in figure 2) serves as input to the classification module where $k$-NN algorithm is used as a classifier and the final output is displayed.

The entire experiment performed is explained in stages.

Stage 1: All the data in colon cancer dataset are classified using $k$-NN algorithm.

Stage 2: ICA is use for feature extraction on the colon cancer dataset, it reduces the data as a sub-optimal of the whole dataset, and classification was then carried out using $k$-NN.

Stage 3: The result of stage 1 and stage 2 are compared based on accuracy, sensitivity, specificity and precision.
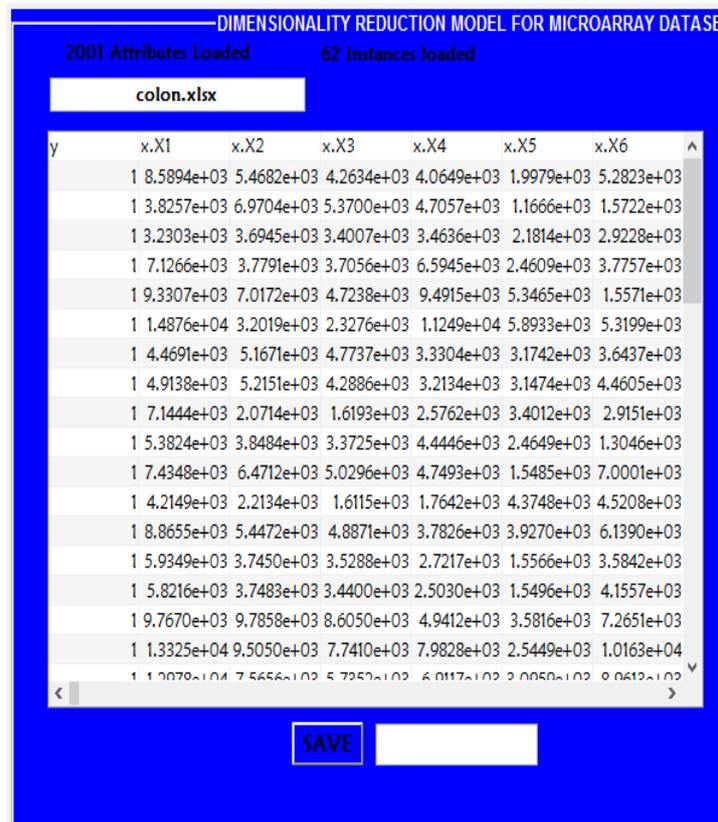
### 3.1 Description of Dataset

This study uses Colon cancer dataset by Alon. It contains DNA microarray gene expression data with 2000 features and 62 samples. MATLAB 2015a was used as implementing tools.

## 4. RESULT AND DISCUSSION

Based on the evaluation metrics used, ICA method performs better compare to the KNN method for the microarray analysis of colon cancer as shown in table 1. The result is also presented graphically in figure 3.

**Table 1: Comparative Analyses of the Performance Metrics Used**

| Performance Metrics | ICA based method | $k$-NN based method |
|---|---|---|
| Accuracy (%) | 88.7 | 77.4 |
| Sensitivity (%) | *100* | 74.1 |
| Specificity (%) | *72.7* | 36.4 |
| Precision (%) | 86.9 | 74.1 |



**Figure 1: Colon cancer dataset loaded for feature extraction process**

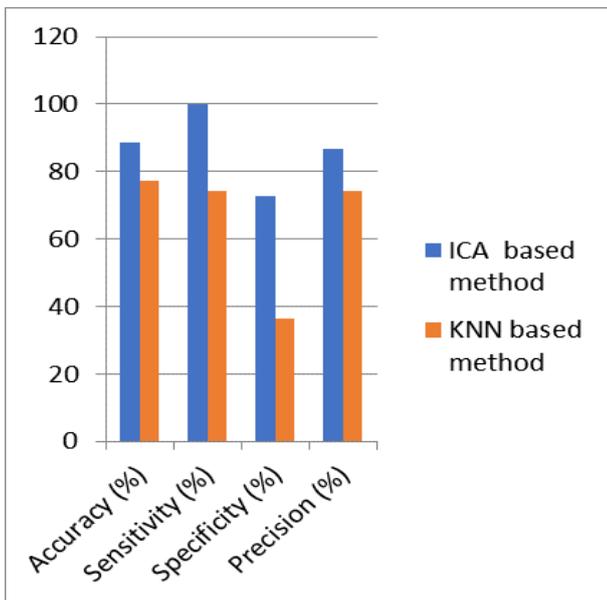**Figure 2: Result of extracted features using ICA**



**Figure 3: Graphical representation of ICA based method and KNN performance.**

## 5. SUMMARY AND CONCLUSION

This study presents an extraction approach to identify informative features in the bioinformatics field. The ICA principle is use to extract features from high dimensional data and imbalanced data. The approach using feature extraction facilities exploits the features of ICA as a feature extraction and works efficiently on colon cancer data in the field. Using the universal computation of $k$-NN and the evolutionary ability of ICA to extract informative genes from a specific microarray dataset. The result of this study reveal that the ICA method combined with KNN algorithm performs better compare to when only $k$-NN is used by 11.3%, 25.2%, 36.3% and 12.8% in terms of accuracy, sensitivity, specificity and precision respectively. In future, there is need to improve the ICA technique in order to handle data with multiclass more efficiently, reproduce identical sets of extracted features and reuse the selection results.

## REFERENCES

[APM11]   **Ali A., Paul J., Madhu G.** - *Dimension reduction of microarray Data Based on Local Principal Component*. World Academy of Science, Engineering and technology, (IJCEACIE), Vol. 5, No. 5, pp. 65-71, 2011.

[CH67]   **Cover T. M., Hart P. E.** - *Nearest neighbor pattern classification*, IEEE Transactions on Information Theory, 13, pp. 21–27, 1967.

[C+14]   **Ching S. T., Wai S. T., Mohd S. M., Weng H. C., Safaai D., Zuraini A. S.** - *A Review of Feature Extraction Software for Microarray Gene Expression Data*. BioMed Research International. Pp. 1-12, 2014.

[DH73]    **Duda R. O., Hart P. E. -** *Pattern classification and scene analysis*, New York: Wiley, 1973.

[FH51]    **Fix E., Hodges J. -** *Discriminatory analysis – Nonparametric discrimination: Consistency properties*. Technical Report 4, USAF School of Aviation Medicine, Randolph Field, Texas, 1951.

[GTA15]   **Ganiz M. C., Tutkan T., Akyokus S. -** *A Novel Classifier Based on Meaning for Text Classification ICA*, European Journal of Human Genetics, Vol. 13, 2015.

[Lin14]   **Lin H. -** *Research on energy-efficient text classification*. 2nd International Conference on Information Technology and Electronic Commerce (ICITEC 2014)

[LSC11]   **Liu Q., Sung A.H., Chen Z.** *- Gene selection and classification for cancer microarray data based on machine learning and similarity measures*. BioMed Genomics, vol. 12, No. 5, 2011.

[PM12]    **Perez M., Marwala T. -** *Microarray data feature selection using hybrid genetic algorithm simulated annealing*, in Proceedings of the IEEE 27th Convention of Electrical and Electronics Engineers in Israel (IEEEI '12), pp. 1–5, 2012.

[RT15]    **Jindal R., Taneja S. -** *A Lexical Approach for Text Categorization of Medical Documents*. Proceeding National Academic. Science, USAVol.96, pp.6745–6750, 2015.

[SP14]    **Singh S., Prakash C. -** *Document Categorization in Multi-Agent Environment with Enhanced Machine Learning Classifier,* 2014 Seventh International Conference on Contemporary Computing (IC3), 2014.

[ST10]    **Suguna N. N., Thanushkodi K. -** *An Improved k-NN Classification Using GA*. International Journal of Computer Science, Vol. 7, No. 4, pp. 18-21, 2010.

[TMD14]   **Trstenjak B., Mikac S., Donko D. -** *KNN with TF-IDF Based Framework for Text Categorization*. Massive Computing Series, Springer, Heidelberg, Germany, pp. 227, 2014

[X+05]    **Xue W. Z., Yee L. Y., Dong W., Feng C., Antoine D. -** *Molecular Diagnosis of Human Cancer Type by Gene Expression Profiles and ICA*, European Journal of Human Genetics, Vol. 13, pp. 1301-1311, 2005.