# A chi-square-SVM based pedagogical rule extraction method for microarray data analysis

**Mukhtar Damola Salawu[1], Micheal Olaolu Arowolo[2], Sulaiman Olaniyi Abdulsalam[3],**
**Rafiu Mope Isiaka[4], Bilkisu Jimada-Ojuolape[5], Mudashiru Lateef Olumide[6], Kazeem A. Gbolagade[7]**
[1,2,3,4,7]Department of Computer Science, College of Information and Communication Science,
Kwara State University, Nigeria
[5]Department of Electrical Engineering and Computer Engineering, Kwara State University, Nigeria
[6]School of Computer Science, Universiti Sains Malaysia, Malaysia

## Article Info

## ABSTRACT

Support Vector Machine (SVM) is currently an efficient classification technique due to its ability to capture nonlinearities in diagnostic systems, but it does not reveal the knowledge learnt during training. It is important to understand of how a decision is reached in the machine learning technology, such as bioinformatics. On the other hand, a decision tree has good comprehensibility; the process of converting such incomprehensible models into an understandable model is often regarded as rule extraction. In this paper we proposed an approach for extracting rules from SVM for microarray dataset by combining the merits of both the SVM and decision tree. The proposed approach consists of three steps; the SVM-CHI-SQUARE is employed to reduce the feature set. Dataset with reduced features is used to obtain SVM model and synthetic data is generated. Classification and Regression Tree (CART) is used to generate Rules as the Last phase. We use breast masses dataset from UCI repository where comprehensibility is a key requirement. From the result of the experiment as the reduced feature dataset is used, the proposed approach extracts smaller length rules, thereby improving the comprehensibility of the system. We obtained accuracy of 93.53%, sensitivity of 89.58%, specificity of 96.70%, and training time of 3.195 seconds. A comparative analysis is carried out done with other algorithms.

*Corresponding Author:*

Micheal Olaolu Arowolo,
Department of Computer Science,
Kwara State University,
Malete, Nigeria.
Email: arowolo.olaolu@gmail.com

## 1. INTRODUCTION

Over the past two decades, the world has witnessed a true explosion of data, which has mainly been driven by innovative storage technology and the increasing popularity of the Internet. Today, a huge amount of data is generated in the medical domain. A popular source is microarray data. Microarray is a biological platform for gathering gene expressions [1]. The technology helps in the identification of new genes. They help to know about the functioning and expression levels under different conditions. It also helps to learn more about different diseases such as heart disease, mental illness, infectious disease and most importantly the study of cancer. Different types of cancer have been classified on the premise of the organs in which the tumors develop. With the help of microarray technology, it will be possible for the researchers to further classify the types of cancer on the basis of the patterns of gene activity in the tumor cells [2, 3]. Microarray is a popular

source of data. The intrinsic problem of a typical data set produced by microarrays is the sample size and the high dimensionality of the data set thus making analysis difficult. In addition the complicated relations among the different genes make analysis more difficult and removing excess features can improve the quality of the results.

In this paper, Feature selection using SVM-CHI-SQUARE algorithm is employed in the first phase to reduce dimensionality of the data by yielding the key attribute in the data. Thus, fewer number and smaller rules are obtained resulting in the improvement of the comprehensibility of the system. Also an approach for pedagogical rule-extraction was proposed which treats the classifier as a "black box" and directly extract rules which relate the input and the output of the SVM. Decision trees (CART) which have a significant advantage of producing interpretable rules were used. Rules were generated by integrating merits of both SVMs and decision tree. Rule set performance is then estimated using the measured rates of true positives (TPs) and false positives (FPs). Using this approach, we will show both an improved classification performance and comprehensibility compared to the previously proposed techniques.

## 2.    METHODS

In this section, we provide a brief introduction of Feature Selection, SVM, Rule extraction and CART

### 2.1.  Feature selection

A number of methods have been proposed for rule extraction from SVMs. Broadly speaking, these methods can be categorized into three main families which are: pedagogical, decomposition, and eclectic [4]. Some of these methods to date still produce relatively large rule sets, which limits their explanation capability [5]. Rule sets can only offer explanation if the number of rules in the rule set is relatively small and its classification accuracy is high. Simpler rules also offer better understanding and explanation [6]. To extract more comprehensible rules, irrelevant features which do not contribute to the classification decision should not be in the rule antecedents. This highlights a requirement to consider feature selection as an integral part of rule extraction. In feature selection, one selects only those input dimensions that contain the relevant information for solving the particular problem. There are three categories of feature selection which are: filters, wrappers, and embedded techniques. This work focuses on filter-based approach. Chi-square to be specific. The difference between Chi-square and other methods and the reason it will be used is that it is very robust with respect to distribution of the data, its simplicity of computation, the detailed information that can be derived from the test, and its flexibility in managing data from both two group and multiple group studies.

### 2.2.  SVMs

The SVM algorithm [7] is a classification algorithm that produces state-of-the-art performance in a vast variety of application domains, including bioinformatics. There are two key reasons for using the SVM in bioinformatics [8]. First, many biological issues involve high-dimensional, noisy data. The SVM is known to behave very well with these data compared to other statistical or machine learning methods. Second, contrary to most machine learning technique, kernel methods like the SVM can easily handle non vector inputs, such as variable length sequences or graphs [9].

For classification problems SVM finds a maximal margin hyperplane that divides two classes. The main intent of SVM is to find an optimal separating hyperplane that correctly classifies data points as much as possible by reducing the risk of misclassifying the training samples and unseen test samples. To address with non-linear issues, SVM first projects data into higher dimensional feature space and tries to find the linear margin in the new feature space [10].

Assuming $\{(x_1, y_1), \dots, (x_n, y_n)\}$ be a training set with $x_{1i} \in R^d$ and $y_i$ is the corresponding target class. SVM can be reformulated as (1) and (2):

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

$$\text{Subject to; } \sum_{i=1}^n \alpha_i y_i = 0 \ \ and \ \ \alpha_i \geq 0, i = 1,2, \dots, n \tag{2}$$

The kernel function is used to solve the problem. The Kernel function analyses the relationship among the data and it creates a complex division in the space [11].

### 2.2.1.  Rule extraction from svm

The support vector machine (SVM) method is a promising classification and regression technique proposed by Vapid and his coworkers [12]. The SVM has been successfully applied to a wide variety of

application domains [13] including bioinformatics [14]. It is especially important for the field of computational biology because it is used for pattern recognition problems including protein remote homology detection, microarray gene expression analysis, recognition of translation start sites, protein structure prediction, functional classification of promoter regions, prediction of protein–protein interactions, and peptide identification from mass spectrometry data [15].

The SVM has shown a superior performance than most traditional machine learning methods such as neural networks (NNs) in many applications. However, these technique still produces black box models with little or no explanation capability. In application areas such as medical diagnosis, there is an evident need for an explanation component to be associated with classification decisions in order to aid the acceptance of these methods by users [16]. To be able to use the extra accuracy of the SVM, which can lead to lives saved or money gained, as well as to obtain a usable, readable model, rules can be extracted from the complex, incomprehensible SVM models. These rules are interpretable by humans and keep as much of the accuracy of the black box as possible [17].

### 2.3. Rule extraction technique

Comprehensibility can be added to SVMs by extracting symbolic rules from the trained model. Rule extraction techniques attempt to open up the SVM black box and generate symbolic, comprehensible descriptions with approximately the same predictive power as the model itself. An advantage of using SVMs as a starting point for rule extraction is that the SVM considers the contribution of the inputs towards classification as a group, while decision tree algorithms like [18] CART measure the individual contribution of the inputs one at a time as the tree is grown.

In general, rule extraction techniques are divided into two major groups i.e. decomposition and pedagogical. Decomposition techniques view the model at its minimum (or finest) level of granularity (at the level of hidden and output units in case of ANN). Rules are first extracted at individual unit level, these subsets of rules are then aggregated to form global relationship [19]. On the other hand, a pedagogical algorithm considers the trained model as a black box. Instead of looking at the internal structure, these algorithms do not make use of the support vectors or SVM decision boundary, but directly extract rules using the input–output mapping defined by the SVM model. These techniques typically use the trained SVM model as an oracle to label or classify (artificially generated) training examples which are then used by a symbolic learning algorithm. The idea behind these techniques is the assumption that the trained model can better represent the data than the original dataset. That is, the data is cleaner, free of apparent conflicts. Since the model is viewed as a black box, most pedagogical algorithms lend themselves very easily to rule extraction from other machine learning algorithms [20].

### 2.4. Decision tree

The comprehensibility of decision trees is one their most useful characteristics, since domain experts can easily understand the principle of the tree, and why a certain object is classified to belong to a specific class. Moreover, decision trees are probably the most extensively researched machine learning method, can deal with any kind of input data (discrete, continuous, binary, attributes). They can also cope with missing values, since the information that attribute values are missing for specific objects can be processed by most decision tree algorithms. The learning process of decision trees is usually quite fast compared to other methods like support vector machines or neural networks, and since most trees are pruned, their classification process is usually also very fast. Several studies [21] have shown that the classification accuracy is generally comparable to the quality of kNN and rule-based learners, but cannot reach the quality of support vector machines or embedded methods, which, on the contrary, are hardly comprehensible (difficult to understand for domain experts) and are not good in handling missing data (since missing data has to be replaced with alternative values such as mean or zero values before classification).

### 2.5. Combined svm and decision tree

The inspiration of combining the SVM and decision tree is to merge the strong generalization ability of the SVM and the strong comprehensibility of rule induction. Specifically, our algorithm employs the SVM as a preprocess of decision tree and consists of three major steps. First, a labelled data set is used for SVM learning purposes, i.e. to build a model with acceptable accuracy. A second data set is generated with the same attributes but different values to explore the generalization behavior of the SVM. That is, the SVM is used to get the class labels for this data set. Hence a synthetic data set is obtained. Finally, the synthetic set is then used to train a machine learning technique with explanation capability. Thereby, rules are generated that represent the generalization behavior of the SVM [22].

## 3.　PROPOSED RULE EXTRACTION APPROACH

The proposed system is a pedagogical/Learning-Based procedure for extracting rules from SVM and is composed of three phases. Firstly, feature selection using Chi-square-SVM is first employed and the actual target values of training instances are replaced by the predictions of SVM models and Case-P (i.e. training instances with corresponding predicted target values) datasets are generated. For the high dimensionality data set produced by microarray, Wisconsin breast cancer dataset is analysed. It is observed that reduced features reduce the complexity of the system and increases the comprehensibility of the rules Secondly, the reduced dataset is used to train the SVM, a second dataset is generated i.e SVM is used to get the class labels for the dataset. Hence a synthetic data is obtained. Finally, Rules are generated using NBTree. Other cases will be discussed is subsequent work. The architecture of the approach proposed is shown in Figure 1.

## 4.　EXPERIMENTAL SETUP
### 4.1.　Dataset description

For building and testing the effectiveness of our algorithm, we performed experiment on cancer dataset. The breast cancer database was obtained from the University of Wisconsin Hospitals, Madison from Dr. William H. Wolberg. This dataset is chosen because of its public accessibility and has previously been used for several Machine Learning studies. The problem is to classify breast masses as either benign or malignant, using nine attributes, all with integer values between 1 and 10. For the class label (2 for benign, 4 for malignant).
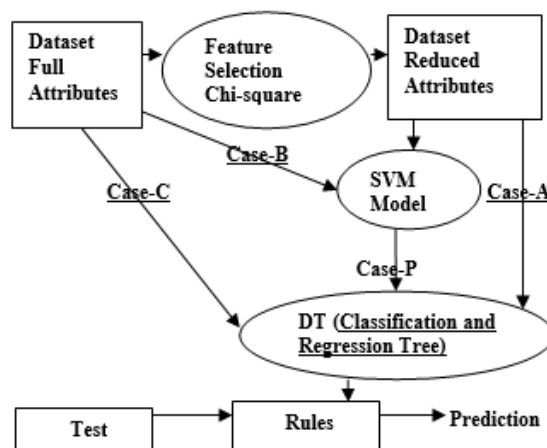


Figure 1. Architecture of the proposed chi-square-SVM pedagogical rule extraction approach

### 4.2.　Experimental setup

This data set has 699 samples. In this experiment they were divided into two parts of 80:20 ratios. 80% and 20% training and test/validation sets with 560 and 139 samples respectively. The 20% of the data is kept aside for later use. Using validation set, efficiency of the rules generated during the experiment is evaluated. The experimental setup were carried out and developed with the Matlab programming (MATLAB 2015A) various functions were developed and linked to a graphic user interface for user interactivity and responsiveness. During the training of the SVM model the following are the parameter settings that was used Cost=0.762942779291553, Kernel function=RBf Kernel the developed systems made use of various component environments in Matlab to develop and output result of the data mining task.

## 5.　RESULT AND DISCUSSION

We evaluate and discuss the performance of our approach SVM+NBTree using Case-P with respect to specificity, sensitivity, and accuracy. During first first phase of the proposed approach, SVM-Chi-square algorithm is employed for feature selection and six attributes are then selected, those are, Clump thickness,

Uniformity of cell size, Marginal adhesion, Single epithelial cell size, Bare nuclei, and Bland chromatin. After selection of optimal subset the selected subset was hence divided into training set and test set in the ratio of 80%: 20%. Then, the training set was passed to SVM model for the prediction of new class label which forms the synthetic class label. A training time of 3.19512 seconds was obtained. The synthetic set is now used to train the decision tree. The CART decision tree was used to generalize the obtained results. Thereby, a total of nine rules were generated that represents the generalization behavior of the SVM. Figure 2 shows the CART tree generated.
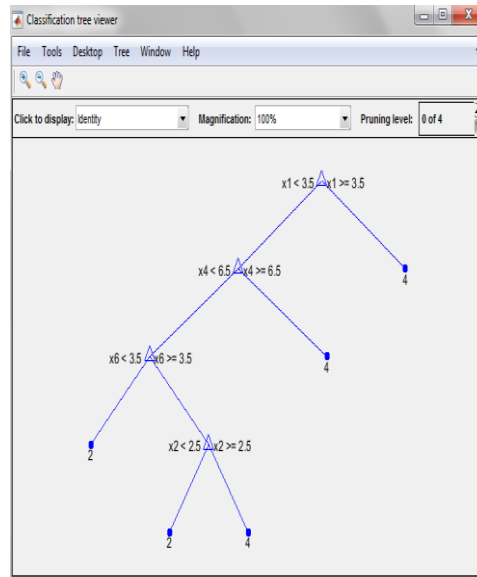


Figure 2. CART tree generated

## 5.1. Confusion matrix for our proposed model

The effectiveness of our approach is evaluated using Accuracy, Sensitivity, and Specificity. The measures used to evaluate the performance of our model is defined as follows:
Sensitivity = TP/ (TP+FN), Specificity = TN/ (FP+TN,), Accuracy = (TP + TN) / (TP+TN+FP+FN), TP = True Positive, TN = True Negative, FP = False Positive, FN = False Negative. Figure 3 shows the confusion matrix for our proposed model i.e Case P



Figure 3. Confusion matrix for our proposed model

From the Figure 1 above our values can be obtained:

TP = 47, TN = 83, FP = 8, FN = 1
Sensitivity = 47/47 + 1 = 0.9791 = 97.9166%
Specificity = 83/3 + 88 = 0.967033 = 96.7033%
Accuracy = 47 + 83/ 47 + 83 + 8 + 1 = 0.9352 = 93.5251%

### 5.2. Comparison of our technique to other svm rule extraction technique

Table 1 shows comparative results of SVM+CART+Chi-square against a number of other previously published SVM rule extraction methods.

Table 1. Comparison with other rule extraction method

| Technique | Accuracy (%) | Sensitivity (%) | Specificity (%) | Number of Rules | Time(sec) |
|---|---|---|---|---|---|
| SVM+NBTree-Farquad (2009) | 77.07 | 68.52 | 78 | N/A | N/A |
| SVM+C4.5- Martens (2007) | 94.6 | N/A | N/A | 9 | N/A |
| SVM+CART+ Chi-square (2016) | 93.5251 | 89.5833 | 96.7033 | 9 | 3.195 |

The churn prediction dataset was analysed by Farquad, while wiscosin breast cancaer data set was analyzed by Martens and us. In these tables, it is evident that SVM+CART+Chi-square has, on the dataset, significantly small rule sets and, therefore, potentially improved comprehensibility. In addition, SVM+CART+Chi-square has good performance, as measured by the overall accuracy. It will be noted that C4.5 outperformed our algorithm in terms of accuracy alone but sensitivity, specificity, and time were not reported for a comparison. It is concluded that there is no one algorithm that can be favored in general. There is currently no one method that can fulfill all criteria simultaneously. There is always tradeoffs.

## 6. CONCLUSION

In this work, we treated Feature Selection as an integral part of rule extraction. Employing Feature selection lead to the removal of irrelevant features which do not contribute to classification decision and extraction of more comprehensible rules. SVMs have proven to be a classification technique with excellent predictive performance. As for many applications, the opaqueness of the trained nonlinear model is an unbreachable barrier; more comprehensible solutions need to be found. The most comprehensible classification models being rule sets, SVM rule extraction tries to combine the predictive accuracy of the trained SVM model with the comprehensibility of the rule set format

Rule extraction techniques generate classification models that have clear advantages. First of all, they are comprehensible and therefore easy to incorporate in real-life applications where clarity of the classifications made is needed. Secondly, the extracted rules only lose a small percentage in accuracy of the black box model from which they are generated. Since support vector machines are among the best performing classifiers, rules extracted from SVMs achieve an accuracy that often surpasses that of the classical methods, such as CART and C4.5. Using the SVM model instead of the original data points eliminates the apparent conflicts and creates a cleaner dataset. In our experiments, the rules generated by CART on the data with labels predicted by the SVM even outperform the CART rules that result from the dataset with the actual class labels. These advantages make it appropriate to consider SVMs and their extracted rules for applications where both accuracy and comprehensibility are required. One no longer needs to settle for the traditional comprehensible, yet less accurate classification methods.

## REFERENCES

[1] Zena, M. H., and Gillies, D. F., "A review of feature selection and feature extraction methods applied on microarray data", *Department of Computing*, Imperial College, 2015.
[2] Selvaraja, S., "Microarray Data Analysis Tool (MAT)",. Akron: Graduate Faculty of The University of Akron, 2008.
[3] Arowolo M.O., Abdulsalam S.O., Saheed Y.K., and Salawu M.D., "A feature selection based on one-way-ANOVA for microarray data classification," *Al-Hikmah Journal of Pure and Applied Sciences*, vol. 3, no. 1, pp. 30-35, 2016.
[4] Andrews, R., Diederich, J., and Tickle, A. B., "A Survey and critique of techniques for extracting rules from trained Artificial Neural Networks", *Knowledge Based Systems*, 1995.
[5] Barakat, N. H., & Bradley, A. P., "Rule extraction from Support Vector Machines: A Sequential Covering Approach IEEE transactions on knowledge and data engineering", 2007.
[6] Duch, W., Setiono, R., & Zurada, J., "Computational intelligence methods for rule-based data understanding", 2004.

[7]    Vapnik, V., "Statistical Learning Theory" New York: Wiley, 1998.
[8]    Schoelkopf, B., Tsuda, K., and Vert, J. P., "Kernel Methods in Computational Biology", *MA: MIT Press*, pp. 71–92, 2004.
[9]    Arowolo M.O., Isiaka R.M., Abdulsalam S.O., Saheed Y.K., and Gbolagade K.A., "A comparative analysis of feature extraction methods for classifying colon cancer microarray data," *EAI Endorsed Transaction on Scalable Information System*, vol. 4, pp. 1-14, 2017.
[10]   Farquad, M. A. H., Ravi, V., Sriramjee, and Praveen G., "Credit Scoring Using PCA-SVM Hybrid Model", *Springer-Verlag Berlin Heidelberg*, 2011.
[11]   Isabelle, G., Jason, W., Stephen B., and Vapnik, V., "Gene selection for cancer classification using support vector machines", *Mach. Learn*, pp. 389-422, 2002).
[12]   Cortes, C., and Vapnik, V., "Support-vector networks", *Mach. Learn*., vol. 20, pp. 237–297, 1995 .
[13]   Cristianini, N., and Shawe-Taylor, J., "An Introduction to Support Vector Machines and Other Kernel-based Learning Methods", *Cambridge Univ. Press*, 2000.
[14]   Schoelkopf, B., Tsuda, K., and Vert, J. P., "Kernel Methods in Computational Biology", *MA: MIT Press*, pp. 71–92, 2004.
[15]   Noble, W. S., In B. Schoelkopf, K. Tsuda, & J.-P. Vert (Eds.), "Kernel methods in computational biology", *MA: MIT Press*, ", pp. 71–92, 2004.
[16]   Fung, G., Sandilya, S., and Rao, R., "Rule Extraction from Linear Support Vector Machines", *Proc. 11th Int'l Conf. Knowledge Discovery and Data Mining*, 2005.
[17]   Martens, D., Baesens, B., Van-gestel, T., and Vanthienen, J., "Comprehensible credit scoring models using rule extraction from support vector machines", *European Journal of Operational Research*, vol. 18, pp. 1466–1476, 2007.
[18]   Martens, D., Baesens, B., Van-gestel, T., and Vanthienen, J., "Comprehensible credit scoring models using rule extraction from support vector machines", *European Journal of Operational Research*, vol. 18, pp. 1466–1476, 2007.
[19]   Farquad, M.A.H., Ravi, D., and Raju, S.V., "Churn prediction using comprehensible support vector machine: An analytical CRM application", 2014.
[20]   Martens, D., Baesens, B., Van-gestel, T., and Vanthienen, J., "Comprehensible credit scoring models using rule extraction from support vector machines", *European Journal of Operational Research*, vol. 18, pp. 1466–1476, 2007.
[21]   Janecek G.K., Gansterer W. N, Demel M., and Ecker G. F., "On the Relationship Between Feature Selection and Classification Accuracy", *JMLR: Workshop and Conference Proceedings*, 2008.
[22]   Arowolo M.O., Adebiyi M.O., and Adebiyi A.A., "Dimensional reduced model for the classification of RNA-Seq Anopheles gambiae data," *Journal of Theoretical and Applied Information Technology*, vol. 97, no. 23, pp.1-10, 2019.

## BIOGRAPHIES OF AUTHORS

Mukhtar Damola Salawu received his B. Sc from the department of Computer Science, Al-Hikmah University. He is a CCNA, CCNP and Microsoft Certified Expert, with his research interests in Data Science, Data Mining, Bio-informatics, Machine Learnine, Computer Arithmetics. He currently resides in the United State of America.

Micheal Olaolu Arowolo received his B. Sc and M. Sc from the department of Computer Science, Al-Hikmah University and Kwara State University, Nigeria, in 2012 and 2017 respectively. He is presently a PhD Student at Landmark University, Omu-Aran, Nigeria. He is a Lecturer at the Institute of Professional Studies, Kwara State University, Malete, he is an Oracle Certified Expert, a member of the IAENG and SDIWC, with his research interests in Data Science, Data Mining, Bio-informatics, Machine Learnine, Computer Arithmetics.

Sulaiman Olaniyi Abdulsalam Holds a Bachelor degree and Master of Science Degree in Computer Science, both from University of Ilorin, Nigeria. He is currently a lecturer at the Department of Computer, Library and Information Science, Kwara State University, Malete, Nigeria. He is a member of Nigeria Computer science. His research areas include Data Mining, Software Engineering and artificial Intelligence.

Rafiu M. Isiaka he has his Ph.D. in Computer science, he is a lecturer in the Department of Computer Science, Kwara State University Malete since 2009. His research interest includes soft computing, e-learning, data mining and information security.

Bilkisu Jimada-Ojuolape received her B. Sc degree in Electrical Engineering and Electronic Engineering from the Kwame Nkrumah University of Science and Technology, Kumasi, Ghana in 2011, and her M. Sc Degree in Systems Engineering, Loughboroh University, Loughborough, United Kingdom in 2013. She is currently pursuing her PhD degree with the School of Electrical and Electronic Engineering, University Sains Malaysia.
She worked with the Nigerian Electricity Regulatory Comission (NERC) for about a year, where she was involved in license application evaluation for Independent Power Producers and Inspection of substations within Abuja, Nigeria, amongst other duties. She was also a National System Engineer with the NTA-Star TV Network (StarTimes), where she managed the operations of ten cities and was also part of a think tank committee that was brainstorming new products for the company.
She joined Kwara State University (KWASU) in 2015. Her current research interests are in renewable energy, power systems, smart grid reliability, ICTs in smart grid, and small hydropower generation, Machine Learning. Mrs. Jimada-Ojuolape is a member of the Association of Practicing Women Engineers in Nigeria (APWEN) and the Nigerian Society of Engi- neers (NSE) and certified by the Council for the Regulation of Engineering in Nigeria (COREN).

Mudashiru Lateef Olumide B.sc holder and an intel Retail certified specialist, he is assistant Manager at National identity management commission (NIMC). He is presently an M. Sc Student of Universiti Sains Malaysia majoring in Informatics with a research interest in Business Intelligence and Data Warehousing.

Kazeem A. Gbolagade, A Professor and Provost at the College of Computer in Information Science, Kwara State University, Malete, Nigeria. was born in Iwo (Osun State), Nigeria, on the 27th of August, 1974. He received his B. Sc degree in 2000 in Computer Science from the University of Ilorin, Kwara State, Nigeria. In 2004, he obtained his Masters degree from the University of Ibadan, Nigeria. In April 2007, he joined the Computer Engineering Laboratory group at the Delft University of Technology (TU Delft), The Netherlands. In TU Delft, he pursued a PhD degree under the supervision of Prof. Sorin Cotofana. He is a member of the IEEE. His research interests include Digital Logic Design, Computer Arithmetic, Residue Number Systems, VLSI Design, and Numerical Computing. His research interests include Digital Logic Design, Computer Arithmetic, Residue Number Systems, VLSI Design, and Numerical Computing