# SMART FINANCIAL FRAUD DETECTION AND CUSTOMER RISK PROFILLING IN FINALCIAL INSTITUTIONS TO IDENTIFY POTENTIAL CRIMINALS USING GENETIC MARKOV ALGORITHM.

**IDRIS MUHAMMAD SANI & DR. ISMAILA IDRIS**
*Computer Science Department. Federal University of Technology, Minna, Minna, Nigeria*

**Abstract**

As computing power

**INTRODUCTION**

63

*keep growing with tremendous growth of electronic transaction, so will the rate of electronic transaction fraud since people will rely more and more on computerized process for their daily activities. Hence, there are needs for more accurate and reliable approach for electronic transaction fraud detection which will help to reduce the illegal activity to the lowest minimum. The use of electronic transaction has increased to a great extent and it caused an explosion in the electronic fraud. Fraud has become one of the major ethical issues in the financial industry. Fraud associated with electronic transaction are also rising today as it is the major mode of payment for both online as well as regular purchase. In order to detect frauds from the mix of genuine as well as fraudulent transactions, efficient fraud detection techniques to detect them accurately are vital rather than simple pattern matching techniques. Here an approach is done to detect the electronic transaction fraud and classify the fraud as either low risk, medium risk or high risk transaction to the financial institution using a fusion approach of genetic and hidden markov algorithm which involve stages of pre-processing in which anonymous transactions were used, genetic algorithm was modelled for feature selection and hidden markov model for classification of fraud as low, medium and high risk transaction. The proposed model is done on existing electronic transaction dataset (anonymous and imbalanced). This research work propose the use of hidden markov model and genetic algorithm to build a model that is able to detect fraud and categorize the customer transaction into three risk levels as low risk, medium risk or high risk transaction to the financial institution to serve as a mechanism which can effectively detect and prevent fraud with great accuracy.*

### Introduction

In today's increasingly internet-dependent society, the use of electronic transactions has become convenient and necessary. electronic

transactions have become the de facto standard for business transactions and Internet ecommerce. The volume of electronic transactions continues to grow leading to higher risks of stolen account numbers and results into great losses to financial institutions (Nilson, 2013). With extensive use of electronic transactions, fraud appears as a major issue in the electronic transactions system. It is hard to have some figures on the impact of fraud, since companies and banks do not like to disclose the amount of losses due to frauds. Another problem in electronic transactions fraud loss estimation is that we can measure the loss of only those frauds that have been detected, and it is not possible to assess the size of unreported/undetected frauds. Other frauds are reported long after the criminal has completed the crime (Richard et al).

According to the Association for Payment Clearing Services (APACS) has estimated that total losses through electronic fraud in the United Kingdom have been growing rapidly from £122 million in 1997 to £440.3 million in 2010 (Linda et al).

According to The Nilson Report August 2013, Global Credit, Debit, and Prepaid Card Fraud losses reached $11.27 Billion in 2012 - Up 14.6% Over

2011. Gross fraud losses accounted for 5.22% total volume, up from 5.07% in 2011. In 2012, only in the USA fraud losses reached $5.33 billion. According to the Lexis Nexis 2014, fraudulent card transactions worldwide have reached around $11 billion a year, and the USA may account for about half of that.

Consumers' demand for electronic transactions due to its convenience and ease of use, and the rise in e-commerce has opened up new opportunities for criminals to steal electronic transaction numbers and consequently commit fraud (Tavan 2010).

Financial services institutions are well aware of the negative impact of fraud. Even at industry average levels, fraud hurts an institution's reputation, customer loyalty, and shareholder's confidence, but even most multi-national financial institutions also have major challenges in this area (Joyner, 2011). Fraud has evolved from being committed by casual fraudsters to being committed by organized crime and fraud rings that use sophisticated methods to hijack a customer's bank account and carryout fraudulent activities on it (Joyner, 2011).

Researchers have defined fraud as a deliberate act contrary to law, rule or policy with the intent of obtaining unauthorized financial benefit. The word fraud can also be used in financial institution as an intentional misstatements or omission of amount to deceive the users of that financial statement (Wang et al, 2006).

Financial fraud is a wide-ranging term for theft and fraud committed using or involving a payment, such as a electronic transaction or debit card, as a fraudulent source of funds in a transaction. The purpose may be to obtain goods without paying, or to obtain unauthorized funds from an account. electronic fraud is also an adjunct to identity theft. According to the United States Federal Trade Commission, while the rate of identity theft had been holding steady during the mid-2000s, it increased by 21 percent in 2008. In 1999, out of 12 billion transactions made annually, approximately 10 million or one out of every 1200 transactions turned out to be fraudulent. According to the Basel Committee on Banking Supervision, fraud can be divided into 2 types: internal fraud and external fraud (Basel Committee on Banking Supervision 2006). Businesses are always susceptible to internal fraud or corruption from its management or employees. While

external fraud is mainly about using the stolen, fake or counterfeit electronic transaction to consume or obtain cash in disguised forms. This thesis is focused on the investigation of the external card fraud, which accounts for the majority of electronic transaction frauds in Nigeria.

Electronic transaction fraud can be either an offline fraud or online fraud. The offline fraud is a stolen physical card at a storefront or call center. The institution issuing the card can lock the account before it is used in a fraudulent manner. Online fraud is committed via web, phone shopping or cardholder-not-present situations. The main objective in fraud detection is to identify fraud as quickly as possible once it is committed (Bolton and Hand 2012).

Fraud detection and fraud prevention are two major effective strategy that have been used in tackling fraud (Bart et al, 2015). Fraud detection refers to the ability to recognize or discover fraudulent actives, while as fraud prevention refers to measures that can be taken to avoid or reduce fraud. The different between the two is that the former is an ex post approach while as the latter is an ex ant approach.

Fraud detection is a subject applicable to many industries ranging from banking and financial sectors, insurance, government agencies and law enforcement, and more (Reurink, 2016). Cases of Fraud have drastic increased in recent years, making fraud detection more crucial than ever (Yifu & Isabel, 2017). Despite all prevention mechanism on the part of the affected institutions of electronic transaction fraud, huge amount of money are lost to fraud yearly because finding fraud is still tricky since relatively few cases show fraud in a large population (John, 2017).

The purpose of this work is to apply machine learning strategies to a unique Nigeria Financial Institutions transaction dataset (electronic transaction filtered dataset), and to investigate whether a meta-learning strategy (a combination methodology of GA and HMM) has the potential to save money and improve fraud detection.

This work primarily aims to evaluate the performance of GA, HMM and Genetic Markov algorithms on the task of customer risk classification with a view to determine which algorithm performs better to improve current fraud detection processes by improving the prediction of fraudulent accounts.

## RELATED WORK

Many techniques have been applied to the field of fraud detection ranging from supervised learning to unsupervised learning. Fraud detection has been usually seen as a data mining problem where the objective is to correctly classify the transactions as legitimate or fraudulent. For classification problems many performance measures are defined most of which are related with correct number of cases classified correctly.

A more appropriate measure is needed due to the inherent structure of electronic transaction transactions. When a card is copied or stolen or lost and captured by fraudsters it is usually used until its available limit is depleted. Thus, rather than the number of correctly classified transactions, a solution which minimizes the total available limit on cards subject to fraud is more prominent.

Since the fraud detection problem has mostly been defined as a classification problem, in addition to some statistical approaches many data mining algorithms have been proposed to solve it. Among these, decision trees and artificial neural networks are the most popular ones. The study of Bolton and Hand provides a good summary of literature on fraud detection problems.

Dal and Bontempi (2015) investigated how machine learning algorithms could be used to address the issues of electronic transaction fraud. The study focused on a framework that is able to report the transactions with the highest risk to investigators by means of algorithms that can deal with unbalanced and evolving data streams.

MohdAvesh et al. (2014) proposed a model using Hidden Markov Model (HMM) and K-clustering to detect electronic transaction fraud. In this model HMM categorizes card holder's profile as low, medium and high spending base on their spending behaviour in terms of amount. A set of probabilities for amount of transaction is being assigned to each cardholder and amount of each incoming transaction  is then matched with card owner's category, if it justifies a predefined threshold value then the transaction is decided to be legitimate else declared as fraudulent.

Zareapoor et al. (2012) reviewed various techniques that have been used for electronic transaction fraud detection. This review analyzed the

68

working principles and the performance of a total of nine(9) machine learning algorithm which are Neural network , Bayesian Network , Support Vector Machine , K-Nearest Neighbor algorithm , Decision tree , Fuzzy logic based system , Hidden Markov Model , Artificial Immune System and Genetic Algorithm (GA).

Delamaire el al. (2009), identified the different types of electronic transaction frauds and various alternative techniques that have been used for fraud detection was reviewed. In their work they outlined common terms in the electronic transaction fraud and key statistics and figure in the field of electronic transaction fraud was also highlighted. The types of fraud highlighted in their work are bankruptcy fraud, theft fraud/ counterfeit fraud, application fraud and behavioral fraud. According to their research pair-wise matching, decision tree, genetic algorithm, clustering and neural network were the various technique used to detect electronic transaction fraud.

However, there are other forms of electronic transaction fraud which is as follows:

### a.  Application fraud

Application fraud takes place when a person uses stolen or fake documents to open an account in another person's name. Criminals may steal documents such as utility bills and bank statements to build up useful personal information. Alternatively, they may create fake documents with this information, they could open a electronic transaction account or loan account in the victim's name, and then obtain monetary benefit from the account.

### b.  Account takeover

An account takeover occurs when a criminal poses as a genuine customer, gains control of an account and then makes unauthorized transactions. The most common method of account takeover is a hacker gaining access to a list of user names and passwords. Other methods include dumpster diving to find personal information in discarded mail, and outright buying lists of 'Fullz,' a slang term for full packages of identifying information sold on the black market.

### c.  Skimming

Skimming is the crime of getting private information about somebody else's electronic transaction used in an otherwise normal transaction. The thief can procure a victim's card number using basic methods such as photocopying receipts or more advanced methods such as using a small electronic device (skimmer) to swipe and store hundreds of victims' card numbers.

## HIDDEN MARKOV MODEL
A Hidden Markov Model is a finite set of states; each state is linked with a probability distribution.
Transitions among these states are governed by a set of probabilities called transition probabilities. In a particular state a possible outcome or observation can be generated which is associated symbol of observation of probability distribution. It is only the outcome, not the state that is visible to an external observer and therefore states are ``hidden'' to the outside; hence the name Hidden Markov Model. Hidden Markov Model will be helpful to find out the fraudulent transaction by using spending profiles of user. It works on the user spending profiles which can be divided into major three types such as:
   i.    Lower profile
  ii.    Middle profile
 iii.    Higher profile.

For every electronic transaction, the spending profile is different, so it can figure out an inconsistency of user
profile and try to find fraudulent transaction. It keeps record of spending profile of the card holder by both ways, either offline or online. Thus, analysis of purchased commodities of cardholder will be a useful tool in fraud detection system and it is assuring way to check fraudulent transaction, although fraud detection system does not keep records of number of purchased goods and categories. Every user represented by specific patterns of set which containing information about last 10 transaction using electronic transaction (Chiu and Tsai, 2014). The set of information contains spending profile of card holder, money spent in every

transaction, the last purchase time, category of purchase etc. The potential threat for fraud detection will be a deviation from set of patterns.
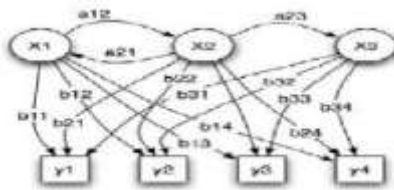


**Figure 1** Architecture of HMM

## Genetic Algorithm

Genetic algorithms, inspired from natural evolution were first introduced by Holland (1975). Genetic algorithms are evolutionary algorithms which aim at obtaining better solutions as time progresses.

GA is used in data mining mainly for variable selection and is mostly coupled with other data mining algorithms and their combination with other techniques has a very good performance (Ekrem et al, 2011). GA has been used in electronic transaction fraud detection for minimizing the wrongly classified number of transactions and is easy accessible for computer programming language implementation, thus, make it strong in electronic transaction fraud detection (Ekrem et al, 2011)..

Zahira Benkhellat and Ali Belmehdi (2012), proposed Genetic Algorithms in Speech Recognition Systems using training model for a speech pattern recognition, this does not only enhances the speed of recognition tremendously, but also improves the quality of the overall performance in recognizing the speech utterance. In general, there are two classic approaches for this development, namely Dynamic Time Warping (DTW) and Hidden Markov Model (HMM). In this article, Genetic Algorithm (GA) is applied to solve involved nonlinear, discrete and constrained problems for DTW .Because of the intrinsic properties of GA, the associated non trival K-best paths of DTW can be identified without extra computational cost. The obtained results show the important contribution of the genetic algorithms in temporal alignment through the increasingly small factor of distortion.

## General Concept of Genetic Algorithm

Being motivated by the principles of natural genetics and natural selection, GA is a stochastic global search method that works on populations of individuals instead of single solutions. It starts with the initial population that has no knowledge of the correct solution. Then GA searches in parallel and depends completely on responses from its evolution operators, i.e. crossover, mutation, and reproduction, to arrive at the best solution. Figure 2 shows the basic structure of a genetic algorithm.
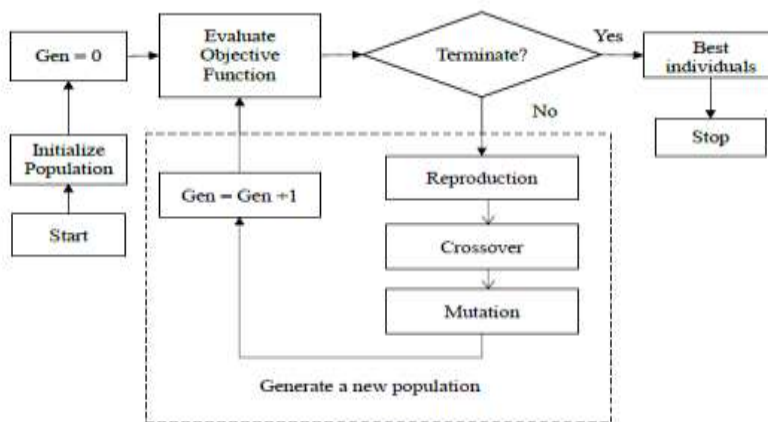


Figure 2: Block diagram of Genetic Algorithm

The candidate solution to a problem is called a chromosome. The first step on applying a GA optimization is to generate an initial population (a collection of chromosomes). The standard way of generating an initial population is random selection. Initial population should be a uniform representation of the entire search space. Otherwise, there are regions of the search space that are not covered by the initial population. Consequently, they can be neglected by the search process. The size of the initial population depends on the computational complexity and exploration abilities. Then GA evolves the population through multiple generations by using the genetic operators, i.e. reproduction, crossover, and mutation, in the search for a good solution.

Genetic Algorithm:

Initiate the strategy parameters
Create and initialize the initial gains

For each gain
Evaluate the objective function J
End

While stopping condition(s) not true do
For i =1, n do

Choose i $\geq$ 2, new gains at random
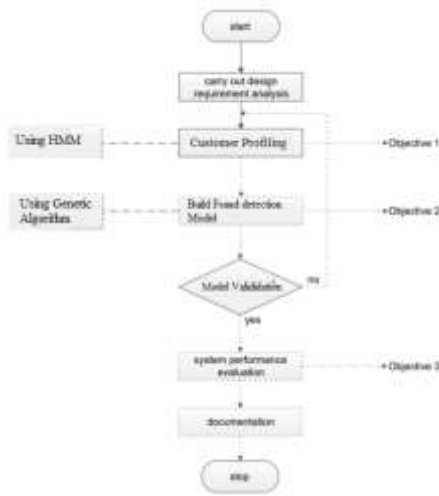Create offspring through application of crossover operator

Mutate offspring strategy parameters
Evaluate the objective function of new gains
If value of objective function is less than epsilon
Best gains
End

Select the new population
t = t +1
End

### Hybrid approach of using hidden markov model and genetic algorithm:

In the process of HMM each incoming transaction is submitted to the Fraud Detection System (FDS) for verification. FDS receives the card details and the value of purchase to verify whether the transaction is genuine or not. If the FDS confirms the transaction to be malicious, it categorize the transaction as either low risk, medium or high risk transaction.

HMM never check the original user as it maintains a log. The log which is maintained will also be a proof for the bank for the transaction made. HMM reduces the tedious work of an employee in bank since it maintains a log. HMM produces high false alarm as well as high false positive. That is overcome by using fusion of this with genetic algorithm which identify fraud accurately and prevent them to classify a genuine transaction as fraud provided that only the relevant fields from the database are extracted into a simple text file by applying appropriate SQL queries which reduce the accessing time and help to identify the fraud easily. The FDS raises an alarm and the issuing bank declines the transaction. The concerned cardholder may then be contacted and alerted about the possibility that the card is misused.



## METHODOLOGY
The research methodology to be employed is depicted in Figure 3.

**Figure 3: Overview of the Research Methodology**
**Design Requirement Analysis (Data Collection Mode and Source)**
The data collected for the study includes, anonymous Customers Transactions from the Financial Intelligence Unit (FIU) of the Economic and Financial Crime Commission (EFCC), The dataset was collected in Microsoft Excel form which include Customers Details with their Transaction Information.

TABLE I. ATTRIBUTES OF TRAINING SAMPLE DATA SET

| Attribute number | Attribute |
| --- | --- |
| 1 | Customer Id |

| | |
|---|---|
| 2 | Authentication type |
| 3 | Current balance |
| 4 | Average bank balance |
| 5 | Times of Overdraft |
| 6 | Electronic transaction age |
| 7 | deducted amount |
| 8 | location of CC used |
| 9 | Time of the CC used with respect to location |
| 10 | Average daily Over draft |
| 11 | Amount of transaction |
| 12 | Electronic transaction type |
| 13 | The Time of using electronic transaction |
| 14 | Card holder income |
| 15 | Card holder age |
| 16 | Card holder position |
| 17 | Card holder profession |
| 18 | Card holder marital status |
| 19 | Average daily spending |
| 20 | Card frequency |

## Customer Profiling

To achieve the first objective, a comprehensive customer profiling by analyzing the spending profiles of electronic transaction holder. Spending profiles of the user can be calculated according to user's past history of transaction in terms of attributes like transaction amount, IP address, shipping address & location of last transaction, etc. This research categorizes the spending profiles of the users into 3 different categories such as high, medium and low based on the level of risk it can cause the user. The customer profiling is carried out in two steps, in first step, the model is trained on the basis of past transaction history and in second step, the model takes the input and check whether transaction details are accepted by trained model or not, otherwise it raises an alarm.
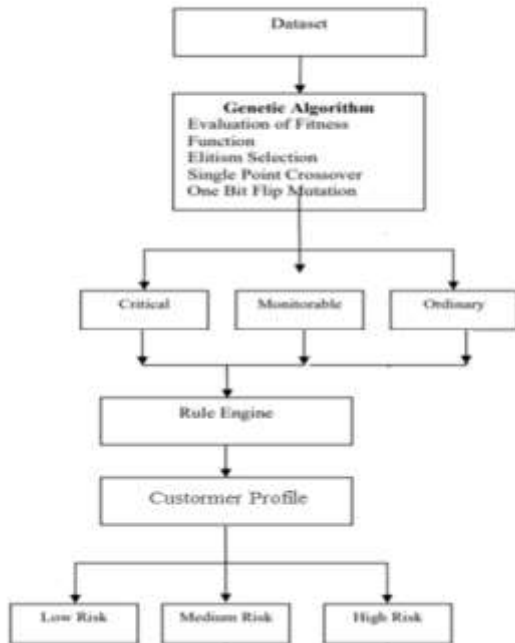
Fig 4: Flow Chart for Customer Risk Profiling

The **Hidden Markov Model** is a finite set of states; every state is associated with a probability distribution. Transitions among these states are administered by a set of probabilities called transition probability. In hidden Markov model only, output is clearly visible but states are not visible which means all states are hidden. In this model firstly, all the transaction sequence need for deciding the category. Into three clusters namely,

*1)* Low Risk Category
*2)* Medium Risk Category
*3)* High Risk Category

According to the user electronic transaction limit. After deciding the categories, the fraud detection of incoming transaction will be verified by last 10 transactions and finding electronic transaction fraud in a system. A

Hidden Markov Model is checking the normal behavior of a card user. The HMM can be well defined with the following elements:

N number of states that are hidden denoted by a set $S = \{S1, S2, S3, .... SN\}$, where $i = 1, 2, ......... N$, is count of state and $Si$, is an individual state.

M denotes the total number of observation symbols. When observations are continuous then M is infinite. We denote the set of symbols

$V = \{V1, V2,......VM\}$ where $Vi$, is an individual symbol.

A set containing probability of moving from one state to another, defined as transition probability.

$aij = P\{qt+1 = Sj \mid qt = Si\}$, $1 \leq i, j \leq N$

where $qt$ denotes the present state. Transition probabilities should satisfy two constraints

$aij \geq 0$, $1 \leq i, j \leq N$

and

summation $= 1$, $1 \leq i \leq N$

Matrix B, indicating observation symbol probability

$B = \{bj(k)\}$

It was found that existing dataset of electronic transaction transaction is highly imbalance, that is for a given dataset the number of fraudulent transactions will be very small relatively to non-fraudulent transaction. This nature of electronic transaction transaction dataset ultimately affects the choice of using a given learning algorithm. Hence, this research adopts the use of Genetic Algorithm and normal multivariate distribution which is a statistical-based technique so as to tackle the imbalance nature of the dataset and F-score was adopted as a means of evaluation since it is ideal method of evaluation when working with imbalance dataset.

### Detection phase

Electronic transaction Frauds are detected based on the source of the fraud, either offline or online detection. offline detection is based on the offline metrics, namely- Card usage frequency, CC Usage Location, rate of unsatisfied transactions and amount of money used in transaction at offline interface.

While online detection is based on the following metrics:

1. Card Usage Frequency.
2. CC Usage Location, Proxy Port check , IP Address Check .
3. Wrong Password Attempt Check, Authentication Type Check.
4. CC Balance, CC Overdraft, Execution Time, Average Daily Spending.

## Detection rule

As shown in Table. 1, first the dataset is loaded into the system. In second step, the detection rules will be applied on each dataset from rule engine module. The rule engine contains the following rules: Average daily spending, CC Usage Frequency, CC Usage Location, Proxy Port check, IP Address Check, Wrong Password Attempt Check, Authentication Type Check, CC Balance, CC Overdraft.

The various parameters involved in the data set are:

CCfreq= number of times card used

CCloc = location at which CCs in the hands of fraudsters

CCoverdraft = the rate of overdraft time

CCbank balance = the balance available at bank of CC

CCdailyspending = the average daily spending amount

## Critical Value identification:

### i. Based on CC usage Frequency

Total number of card used  CCfreq = (CU) / CC age

If CCfreq is less than 0.2, it means this property is not applicable for fraud and critical value = CCfreq Otherwise, it check for condition of fraud (i.e)

Fraud condition = number of time Card used Today (CUT) > ( 5 * CCfreq)

If true, there may chance for fraud using this property and its critical value is CUT*CCfreq

If false, no fraud occurrence and critical value =CCfreq

### ii. Based on CC usage Location

Number of locations CC used so far (loc) obtained from dataset (loc)

If loc is less than 5, it means this property is not applicable for fraud and critical value =0.01

Otherwise, it checks for condition of fraud (i.e)

Fraud condition = number of locations Card used Today (CUT) > (5 * loc)

78

If true, there may chance for fraud using this property and its critical value is loc/CUT

If false, no fraud occurrence and critical value =0.01

### iii.    Based on CC OverDraft

Number of times CC overdraft with respect to CU occurred so far Consider the (OD) can be found as, OD with respect to CU = OD/CU

If OD with respect to CU is less than 0.02, it means this property is not applicable for fraud and critical value = Od with respect to CU Otherwise, it checks for condition of fraud (i.e)

Fraud condition = check whether overdraft condition occurred today from (ODT dataset)

If true, there may chance for fraud using this property and its critical value is ODT * OD with respect to CU

If false, no fraud occurrence and critical value = Od with respect to CU

### iv. Based on CC Book Balance

 Standard Book balance can be found as,

Bb = current BB / Avg. BB

If bb is less or equals than 0.25, it means this property is not applicable for fraud and critical value = BB Otherwise, it check for condition of fraud (i.e)

If true, there may chance for fraud using this property and its critical value is currBB * BB

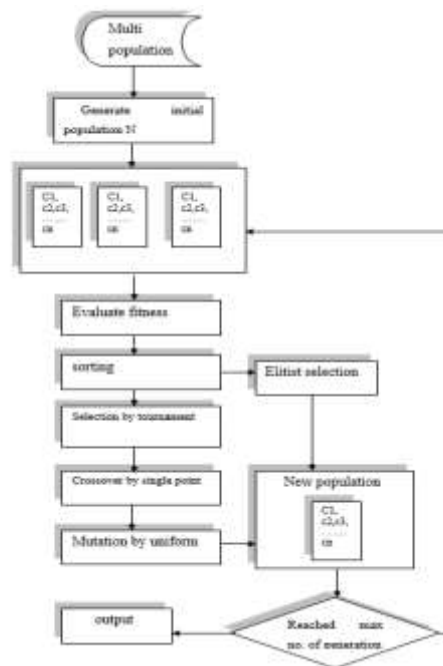If false, no fraud occurrence and critical value = BB



**Figure 5** Flow of Genetic algorithm

For every transaction summation of all critical values by each rule is computed and then k-means clustering algorithm applied on summation of

79

critical values for each transaction. In order to overcome execution time of electronic transaction fraud risk assessment model k-means algorithm is used which will form three clusters low risk, medium risk, and high risk. Genetic algorithm select fittest individuals from medium risk and high risk cluster then perform single point crossover and one bit flip mutation to produce new population until best result found and give results in terms of critical, monitorable and ordinary records. Genetic algorithm is robust algorithm and gives optimization solution but due to more number of iteration to produce new population increases execution time. Hence in electronic transaction fraud risk assessment model we have used Genetic k-means algorithm to use features of simple GA and to overcome execution time.

## Electronic transaction Fraud Detection Model
## Mathematical model

This Mathematical model can be derived using in six (6) stages which are:

Step 1: The whole dataset which contains thirty thousand (30,000) observation was randomly divided into seventy percent (70%), fifteen percent (15%) and fifteen percent (15%) as training dataset, validating dataset and testing dataset respectively.

Step 2: The training dataset was used to compute the mean vector ($\mu$) which is a column vector with each entry corresponding to the mean a column in the training dataset. That is

$$\mu = \{\mu_1 \, \mu_2 \ldots \mu_{n-1} \, \mu_n\}$$

Where n is the number of features.

Also the covariant matrix ($\Sigma$) of training dataset was computed, which is an n by n matrix. That is

$$\Sigma = \begin{pmatrix} \Sigma_{1\,1} & \cdots & \Sigma_{1\,n} \\ \vdots & \ddots & \vdots \\ \Sigma_{n\,1} & \cdots & \Sigma_{n\,n} \end{pmatrix}$$

$\mu$ and $\Sigma$ were computed using MATLAB.

Step 3: Hidden Markov Model $N(X, \mu, \Sigma)$ were used to compute probability vector ($P$) of all the observation ($X$) in the validation dataset. That is,

$$P = \{P_1 \, P_2 \ldots P_{i-1} \, P_i\} \text{ and}$$

$P_i = N(X_i, \mu, \sum)$

Where i = the number of observations in the validation dataset.

Step 4: Genetic Algorithm (GA) was to select $\varepsilon$ (a real number) that minimizes the misclassification rate of the validation dataset such that

$$\begin{cases} P_i < \varepsilon \text{ fraudulent} \\ otherwise\ not\ fraudulent \end{cases}$$

Step 5: $N(X, \mu, \sum)$ and $\varepsilon$ are bundled to form a machine learning classifier model (called Highbred Electronic transaction Fraud Detection HCCFD model) such that for an observation $x$

$P = N(x, \mu, \sum)$ then

$$\begin{cases} P < \varepsilon \text{ fraudulent} \\ otherwise\ not\ fraudulent \end{cases}$$

Step 6: The performance of the model was evaluated using the testing dataset.


**Standard Performance Measures**

The intended approach and the selected metrics for carrying out the performance analysis is outlined as follows:

**Classification Accuracy**: this can be defined as the percentage of the test set that the model projected as correct. It is the ratio of number of correct predictions to the total number of input samples.

$$Accuracy = \frac{Number\ of\ correct\ predictions}{Total\ number\ of\ predictions\ made}$$

It works well only if there are equal number of samples belonging to each class. For example, consider that there are 98% samples of class A and 2% samples of class B in our training set. Then our model can easily get 98% training accuracy by simply predicting every training sample set with 60% samples of class A and 40% samples of class B, then the test accuracy would drop down to 60%. Classification Accuracy is great, but gives us the false sense of achieving high accuracy.

The real problem arises, when the cost of misclassification of the minor class samples are very high. If we deal with a rare but fatal disease, the cost

of failing to diagnose the disease of a sick person is much higher than the cost of sending a healthy person to more tests.

**Speed of Prediction**: it is the percentage of the test set that the model predicts correctly. It is denoted as:

**Misclassification Rate**: Also referred to as Logarithmic Loss or Log Loss, works by penalising the false classifications.

False Positive Rate: this can also be called the false alarm rate, it is the percentage of the test set that the model predicts falsely as positive when it was actually negative. It is denoted as:

False Negative Rate: this can be referred to the percentage of the test set that the model predicts falsely as negative when it is actually positive.It is denoted as:

It works well for multi-class classification. When working with Log Loss, the classifier must assign probability to each class for all the samples. Suppose, there are N samples belonging to M classes, then the Log Loss is calculated as below:

$$\frac{-1}{N}\sum_{i=1}^{N}\sum_{j=1}^{M} y_{ij} * \log(p_{ij})$$

where,

yij, indicates whether sample i belongs to class j or not

pij, indicates the probability of sample i belonging to class j

Log Loss has no upper bound and it exists on the range $[0, \infty)$. Log Loss nearer to 0 indicates higher accuracy, whereas if the Log Loss is away from 0 then it indicates lower accuracy.

In general, minimising Log Loss gives greater accuracy for the classifier.

## F1 Score

F1 Score is used to measure a test's accuracy. F1 Score is the Harmonic Mean between precision and recall. The range for F1 Score is [0, 1]. It tells you how precise your classifier is (how many instances it classifies correctly), as well as how robust it is (it does not miss a significant number of instances).

82

High precision but lower recall, gives you an extremely accurate, but it then misses a large number of instances that are difficult to classify. The greater the F1 Score, the better is the performance of our model. Mathematically, it can be expressed as:

$$F1\ Score = \ 2 * \frac{1}{\frac{1}{precision} + \frac{1}{recall}}$$

F1 Score tries to find the balance between precision and recall.

**Precision** : It is the number of correct positive results divided by the number of positive results predicted by the classifier.

$$Precision = \frac{True\ Positives}{True\ Positives = False\ Positives}$$

**Recall** : It is the number of correct positive results divided by the number of all relevant samples (all samples that should have been identified as positive).

$$Precision = \frac{True\ Positive}{True\ Positives + Flase\ Negatives}$$

**RESULTS**

The model was trained on customers transactions labeled dataset with the results obtained shown in Table II. Genetic Algorithm fitted averagely on both the training data and test data. The model obtained an accuracy of 51.67% on the test data. However, it can be observed that the default Hidden Markov Model fitted fairly good on both the training data and test data, but had poor performance. This model could not separate low risk customers from high risk customers and those that do not contain any medium risk accurately therby leading to high level of misclassification rate. It also had the problem of separating high risk related transactions from medium risk transactions. This may be due to the data being imbalanced and the algorithm might be biased towards the majority classes because the loss function did not take the data distribution into consideration. The fusion of Genetic and Markov model obtained an

accuracy of 79.9% on the test set. From the result in Figure 7, it can be seen that the model fitted well on all the classes. It classified 25% of the transaction dataset  as medium risk, 10% as low risk and 65% as high risk transactions.
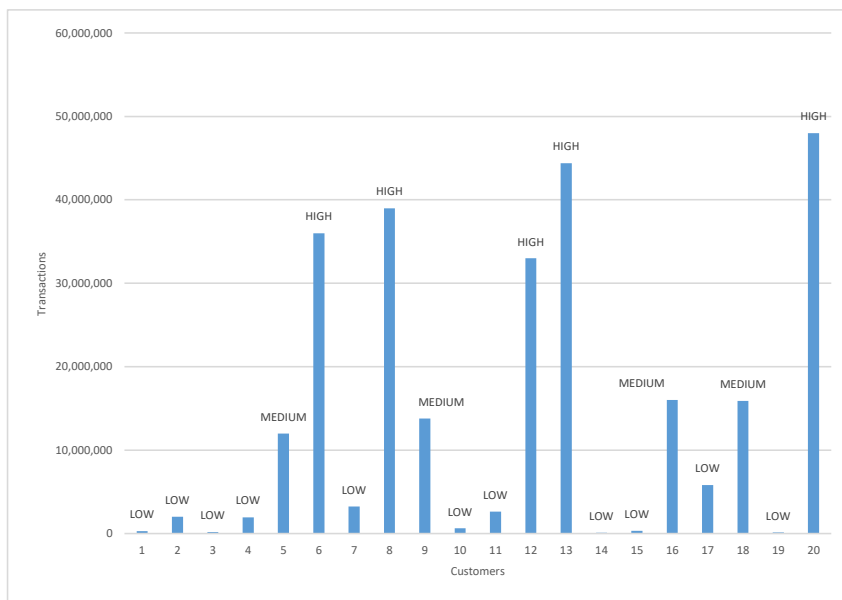


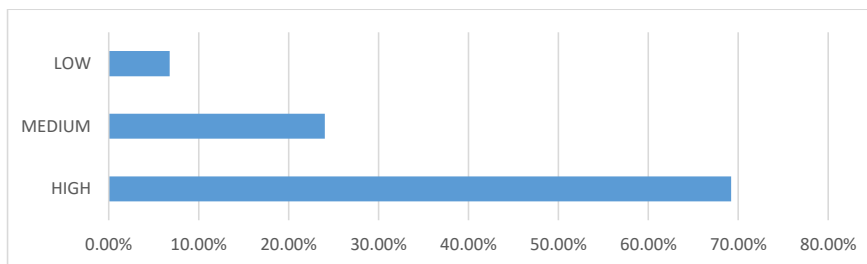Figure 6. Spending Pattern of Customers Transaction



Figure 7. Percentage of Spending Profile

84

| Approach | Accuracy (%) | Precision (%) | Sensitivity (%) | Specificity (%) |
|----------|--------------|---------------|-----------------|-----------------|
| GA | 51.67 | 67.5 | 32.5 | 25 |
| HMM | 43.23 | 62.8 | 30.23 | 23.25 |
| GA & HMM | 79.9 | 76.2 | 73.5 | 61.2 |

Table II. Results Of Genetic And Hidden Markov Models

The HMM model obtained the least accuracy of 43.23% on the transaction dataset meaning that it did not fit well on both the training data and test data. Also, as can be seen from the Table II above, it did not fit well on all the classes as it classified majority of all the transactions into the low risk class.

The GA model performed a little better than the HMM model as shown in Table II. However, the model still classified majority of all the transactions as low. This may be due to the data being imbalanced and the algorithm might be biased towards the majority class because the loss function did not take the data distribution into consideration.

From the results in Table II, the fusion of GA and HMM obtained an accuracy of 79.9% outperforming ordinary HMM which got an accuracy of 43.23% on the test dataset and it also shows that GA obtained an accuracy of 51.67% which is a little higher than the accuracy of HMM which obtained.

## DISCUSSION

The results from all the three experiment showed that Genetic Markov models generally performed well in classifying transactions into three different categories. HMM which was the worst performing model got the lowest accuracy of 43.23%. This may be due to imbalance in the dataset and the algorithms may be biased towards the majority classes because the loss function did not take the data distribution into consideration. To prove this assumption to be true, there is need to balance the dataset or to increase the number of transactions in the low and high-risk classes for training. The Genetic Markov model achieved higher accuracies than the individual algorithms.

## CONCLUSION

The aim of this paper was to evaluate the performance of GA, HMM and Genetic Markov algorithms on the task of customer risk classification with a view to determine which algorithm performs better to improve current fraud detection processes by improving the prediction of fraudulent accounts.

The best performing model was identified by comparing the accuracies of the three machine learning models trained in this experiment. Out of all the models evaluated, the Genetic Markov models outperformed the individual machine learning algorithms trained in this experiment.

How well each model performs on fraud detection classification can be influenced by different factors such as the size of the dataset, how balanced the dataset is, the chosen parameters and how the preprocessing of the raw data is performed.

The results in this study showed that the individual machine learning models performed poorly in classifying transactions fraud risk with low level of accuracy, precision, sensitivity. This leaves the authors with the conclusion that the performance of the models was improved with the fusion of the two models.

## FUTURE WORK

The models in this paper were trained based on the amount contained in the transactions to profile customers with tendencies of performing fraudulent transactions. The multimodal analysis of transactions that includes customer behavior and politically exposed person (PEPS) is an important future work. Also, profiling customers based on their transactions to detect customers with tendency to perform money laundering, or financial terrorism is another future work.

## REFRENCES

Advance Fee Fraud and Other Fraud Related Offences Act (2006). *Laws of the Federation of Nigeria.* Aggarwal, C. (2016). Outlier Analysis Second Edition.

APWG (Anti-Phishing Working Group) (2008). Phishing Activity Trends Report. Retrieved November 15, 2017, from http://www.antiphishing.org

Atherton, M., (2010). Criminals switch attention from cheques and plastic to internet transactions. *The Sunday Times of March 10, 2010.*

Barnett, V., & Lewis, T. (1994). Outliers in Statistical Data. *John Wiley & Sons.* Hyperlink "https://www.sas.com/sas/books/authors/bart-baesens.html" \t "_blank"

Bart, B., Veronique V. & Wouter, V. (2015) . Fraud Analytics Using Descriptive, Predictive, and Social Network Techniques: A Guide to Data Science for Fraud Detection. John Wiley Sons Inc

Bolton,R.J., & Hand,D.J.(2012). "Statistical Fraud Detection: a review, Statistical Science" International Journal of Computer Application Volume 17 (3), (2012), page 235–255

Cameroon, D. (2011). Cyber crime costs UK 27 Billion pounds. *Reuters media briefs.*

Dal Pozzolo, A., & Bontempi, G. (2015). Adaptive machine learning for electronic transaction fraud detection.

Chiu,C., and Tsai, C., 2014. A Web Services-Based Collaborative Scheme for Electronic transaction Fraud Detection, Proceedings of IEEE International Conference e-Technology, e-Commerce and e-Service(2004),pp. 177-181

Delamaire, L, Abdou, HAH and Pointon, J (2009). Electronic transaction fraud and detection techniques: a review. University of Salford

Denning D., (1999). Information Warfare and Security. *ACM Press USA.* EFCC/ NBS/ (2010). Business Survey on Crime & Corruption and Awareness of EFCC in Nigeria. *Summary Report.*

Ellen, J. (2011). Enterprisewide Fraud Management. *SAS Institute Inc. Cary, NC, USA.* 029-2011

Ellison, L., & Akdeniz, Y.,(1998). Cyberstalking the Regulation of Harassment on the Internet. *1998 Criminal Law Review December Special Edition: Crime, Criminal Justice and the Internet.* pp 29-48.

Ekrem Duman, M. Hamdi Ozcelik "Detecting electronic transaction fraud by genetic algorithm and scatter search". Elsevier, Expert Systems with Applications, (2011). 38; (13057–13063).

Hawa D., Amrizah K., Zuraidah M., Khairun S.K., (2014), "Detecting Fraudulent Financial Reporting through Financial Statement Analysis", Journal of Advanced Management Science Vol. 2, No. 1, March 2014.

Joyner Ellen, (2011), "Detecting and Preventing Fraud in Financial Institutions", Enterprise wide Fraud Management, SAS Institute Inc., Page 9.

Masoumeh Zareapoor, Seeja.K.R,M. Afshar A. "Analysis of Electronic transaction Fraud Detection Techniques: based on Certain Design Criteria" International Journal of Computer Applications (0975–8887)Volume 52–No.3, August 2012.

Mikolov, T. Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space, CoRR, vol. abs/1301.3781. doi: https://arxiv.org/abs/1301.3781

Nilson Report. (February, 2016). "Purchase Volume at Merchants in the U.S." Issue 988, Retrieved from https://nilsonreport.com/upload/content_promo/The_Nilson_Report_10-17-2016.pdf

Pennington, J., Socher, R., & Manning, C.D. (2014). Glove: Global vectors for word representation in Empirical Methods in Natural Language Processing (EMNLP), pp. 1532–1543.

Richard, J Sullivan., 2013. "The U.S. adoption of computer-chip payment cards: implications for payment fraud," Economic Review, Federal Reserve Bank of Kansas City, issue Q I, pages 59-87.

Tayan, B. & Larcker, D. (2012), "Corporate Governance Matters: A Closer Look at Organizational Choices and Their Consequences", New Jersey: Pearson Education Inc., Drake Management Review, 2011, Volume 2, Issue 1, page 5.