# Predicting Customers Behavior on an Online Retail System Using Association and Clustering Machine Learning Algorithms (Comparative Analysis using SAS and R Languages)

IBEKWE ABUNDANCE EMERIE[1], OGUOMAIKECHUKWU STANLEY[2], EZURIKE ONYEWUCHI[3], VICTORY CHIBUIKE ONUMAKU[4], EBELE PRECIOUS OKEMBA[5]

[1] Carribean and African Health Network (CAHN), United Kingdom
[2] Department of Computer Science, University of Agriculture and Environmental Science (UAES) Umuagwo, Nigeria
[3] Department of Computer Science, Fed Polytechnic Nekede Owerri, Nigeria
[4] Edge Hill University, Ormskirk, United Kingdom
[5] University of Kent, United Kingdom

**Abstract-** *The aim of this research is to adopt Machine Learning and Data Mining tools in predicting Customers behaviour on an Online Retail System Using Association and Clustering Machine learning Algorithms (Comparative Analysis using SAS and R Languages). The objective of the research includes analyzing a dataset using Association and Clustering mining tool to predict customers behaviors' with likely products to purchase based on customers purchase records and to facilitate customer's product decision on an online retail systems. The research was motivated due to the high demand of online services and products, inadequate customer's product satisfaction and poor product response time and product decision. Two different data mining methodologies were applied differently in the study, association and clustering algorithm while apriori and k-mean modeling tools where adopted. The data was analyzed with R and SAS Enterprise Miner. The experiments are done on UK-based and registered non-store online retail dataset, sourced from UCI machine learning repository. The result after the experiment was able to provide a model that could facilitate customer product decision making and hence provide fast online retail activities to customers and after a comparative analysis on the two different approaches, it was proven that R is one of the best data modeling tools as it is easier to clean or explore a dataset accurately than in SAS.*

*Indexed Terms- Artificial Intelligence, Machine Learning, Online Retail, Association, Clustering Algorithm, K-Means Modeling, Apriori modeling and Data Mining tools*

## I. INTRODUCTION

Data science has set the way for better analysis of data, looking at the amount of data generated on daily basis in most online transactional systems, more especially marketing on ecommerce platforms. Online Retail services are taking over the internet as number of online marketing website has increased rapidly over the globe with advanced technological tools and devices. Ability for customers to search and take decision based on a particular product online is a very huge task that requires speed and accuracy. Providing an easy means of assistance towards these issues is one of the objectives of applying association and clustering rule in this paper. For example, a customer might visits an online shop in purchase of a product say "A" which maybe has been purchased by the customer earlier, therefore the system will provide a similar type of same product say "B" base on the previous purchase made by the customer. Here, the association data mining rule has been applied to the marketing activity in this process. Now

let's look at what data mining is all about, data mining also sometimes referred as knowledge discovery, is a process of analyzing data from different perspectives and summarizing it into useful information [1]. More so [1] further stated that data mining deals with huge amount of data kept in the database, so as to locate required information and hence provides reliable facts, the birth of data mining tool towards data discovery, prediction and extraction has helped researcher and scientist in numerous ways to uncover hidden facts both in physical and biological sciences. Data mining (DM) is one of the hidden extractions of predictive tool in progression from a very big database, and is now a new great technology by means of large potential to provide huge support to companies looking into the direction of large data discovery from knowledge warehouses.[2].It is now clear that the introduction of DM tool in data processing has help to predict the future trends and behaviours of events, enabling various business owners create practical knowledge-driven selections in such that it could help to answer business queries that factually remained too time irresistible to resolve. On the other hand, Cluster analysis according to [1] stated that it is a procedure which learns the substructure of a data set by distributing it into several clusters. Clustering is defined by [3]as methods for grouping of unlabeled data. Nevertheless, the clustering algorithm was adopted so that it will be easy to group customer transaction on the online retail system based on their purchase behavior and history. Customer's records are being tracked based on purchasing behavior and designed strategic initiatives. The objectives of this research is to provide a more unique means of providing and predicting customers behaviors' with likely products to purchase based on customers purchase records and hence facilitate customers product decision on an online retail systems. This research is organized as follows: Introduction: presents general introduction of machine learning (ML) and its importance to online business, it gave meaning and definitions of association rule and clustering algorithms, it also looked at data mining (DM) and how it is an important tool for effective prediction and making discovery for future use, it also highlighted the objectives of the study and gave significant facts why adoption of machine learning and data mining tools are good for decision making

more especially for an online businesses. Literature Review: looks at generally the literature review on related works, machine learning modeling tools and technique, association and clustering algorithm, types of association rule approaches, and different techniques of the clustering, Methodology: methodology adoption for the study, proposed system diagram, system algorithm while Results: present results, conclusion and recommendation of the study.

## II. LITERATURE REVIEW

This is one of the most recommended and researched data mining technique, these technique is aimed at extracting correlations, relationships or associations on various objects (sets of items or transactions) in either relational or distributed databases or from other related repositories within a system. According [3] stated that association is an issue surrounding the process of finding relationships in different attributes in very large customers data center of database. The attributes here may be in zeros or once ($0_s$ and $1_s$) literals, or in a quantitative state. The problem in an association data mining rule is to identify the nature of the causalities between values of different attributes [4]. The researcher used an example to illustrate an association rule stating that a supermarket that is being maintained, containing customers information from different transactions are sets of items bought by each of the customers, hence finding the purchase behavior of any item would affect another customer's purchase behavior respectively. Therefore applying association rule in this regards would provide an easy process of finding the relationship between such transactions accurately without affecting other information in the database. The identified customers data would be used to make accurate decision in respect to the purchase behavior of the customer [4].ASSOCIATION RULE APPROACH: Let A = {A1, A2, A3, A4………..An}be a set of binary literals known as items. Each of the transactions (T) is a set of items, such that T$\subseteq$ A. These correspond to the set of items which a customer may buy in a market transaction. Now the association rule here is a condition of the form z $\Rightarrow$ $Y$ $where$ $Z \subseteq$ A and Z$\subseteq$A are two sets of items. Looking at the above equations, the idea of a data mining association rule is to provide a systematic means by which a user can figure out how

to infer the presence. Furthermore, there are different types of association algorithms as stated by [5], such as: Apriori algorithm, Eclat algorithm, FP growth algorithm, Node set based algorithm, GUHA procedure ASSOC, OPUS search and Context based association rule mining algorithm. On the Other hand, Clustering Rule:is a type of data mining technique that deals with the grouping of similar records together in a very large database of multidimensional records, these clustering rule creates segments of data with the same similarity within a given group of points, based on the type of application use for the analysis, each of these segments may be treated differently [4]. Figure 1 below illustrates the various phases of clustering rule algorithm.



Figure 1: Diagram Showing Clustering Algorithm (Source: [1])

## III. ADOPTED METHODOLOGY

- Association and Clustering Rule Algorithm

This study adopted two different data mining algorithm for an effective and more accurate result provision. The study implemented the both methodology on the topic differently (An online Retail Services) while each methodology is being adopted, a given model from the methodology was selected and apply on the online retail services dataset. The methodologies used are: Association Rule Algorithm and Clustering Rule Algorithm. The adopted algorithm for the study includes Apriori Rule Algorithm (Bottom up Approach) on dataset which is a known approach by any data scientist as it goes together with market basket analysis (MBA) in helping retailers boost business by proving accurate prediction of items to customers by using R language.

The variables in the dataset are:
1. InvoiceNo – number uniquely assign to each transaction in which some starts with "C" meaning cancelled transaction (6-digit integer)

2. StockCode– SKU number that uniquely differentiate a product from another product (5-digit integer)
3. Description—Product (item) name. Nominal
4. Quantity—The quantities of each product (item) per transaction. Numeric
5. InvoiceDate- Invoice Date and time. Numeric, the day and time when each transaction was generated.
6. UnitPrice- Unit price. Numeric, Product price per unit in pound sterling.
7. CustomerID—Customer identification number number. Nominal, a 5-digit integral number uniquely assigned to each customer
8. Country—Country name. Nominal, the name of the country where each customer resides.

## IV. THE PROPOSED SYSTEM DIAGRAM



Figure 2: The Proposed System Diagram (Source: Fieldwork 2023)

## V. ASSOCIATION RULE EXPERIMENTS ON THE DATASET USING R LANGUAGE

After proper packages installation, the Arule was used to initiate the association rules where the files are read by using the readxl script which was used to read or load the file named online_retail into the RStudio environment. Figure 3 shows the file after reading with 541909 observations of 8 variables. In trying to clean the data set so as to remove some values with the same invoiceN0, cancel transactions with c alphabet witnessed in the data set. Some of the cause of the double invoiceN0 might be because of some items were bought at the same time therefore cleaning the data set will help in correcting and arranging them as one transaction.

Figure 3: first 36 rows view of the dataset
(Source: Fieldwork 2023)



By further exploring the dataset using the head script which shows the first 6 values with tail of first 6 set structure of the data set. Applying skinr enable the output of the statistical analysis of the data set shown in figure 4 above. Therefore trying to remove the canceled transactions, the script was used "length(which(substr(str_to_UPPER(Online_Retail& InvoiceN0)1,1)== "C"))" which shows a value of

9288 observations with cancelation, then the script to remove all observations with alphabet C and save it in a new variable name called onlineRetail_non_Cancel <Online_Retail( which(substr(str_to_UPPER(Online_Retail &Invoic eN0),1,1)=="C")). Now to totally roved all variable containing missing values, the script was used "Online_Retail <- Online_Retail_nonCance(Complete.case(Online Retail_nonCancel))" which gave values of observation of 397924 of 8 variables. Adding rows or editing some of the columns in the data set, a function was called "mutate" because there is need to carry out some conversions like the description variable to factor variable, the same goes to conversion of country variable to a factor and editing orseparation of date/time and set independently then lastly, the invoiceN0 variable was converted to numeric then bind the new columns together. Then after the various conversions, the data set was further converted into a transaction format

where a script was written to check every transaction to find out if there is any transaction with the same invoiceN0, Country, date and description if there is, it will join them together and separate with a comma which shows 18536 observations of 4 variables, because to enable the creation of the data set to basket format requires removing of all variables and leave one, therefore there is need to rename the newly created variable to "Product" by using the script "colnames(transactionData<-c("Product"))," then save it with a new name called clrd_online_retail_mkt_bskt.csv in the directory. Then loading of the new dataset into an object of the transaction class and assign the name "tr" to it, this is done by using the R function read.transactions of the arules package which shows transactions in sparse format with 18537 transactions in rows and 7896 items in columns. Figure 7 shows the summary is of the most frequent items with (White hanging heart T-Light holder of 1729 times, Regency cakestand of 1611 times, Jumbo bag red retrospot 1363 times, Party bunting 1285 and Assorted color bird ornament 1180 times and others with 319810) while Figure 6: absolute frequency plot which shows how frequent they appeared independently while figure 6 below shows relative item frequency bar plot of the top 20 products/items after dataset cleaning and color with "RColorBrewer" package. The relative frequency plot shows how each item is related while absolute frequency bar plot shows the absolute state of the items.

## MODEL DEVELOPMENT USING APRIORI ALGORITHM

Themodel was achieved usingapriori algorithmwhich is a function in the arules package. After using the default parameters of rule which are minimum support of 0.1, minimum confidence of 0.8, maximum of 10 items (maxlen), and a maximal time for subsetchecking of 5 seconds (maxtime) and after a successive run, there was no result generated shown figure 8a.



Figure 8a: The association rule using apriori algorithm



Figure 8b: The association rule using apriori algorithm

Therefore, there is need to reduce the parameters to a smaller state so as to generate a small rules. Because the researcher is interested in stronger rules in order to make the report and the visualization simple. Hence the parameters were set as follows (supp=0.01, conf=0.8, maxlen=10), and 18 rules were generated shown in figure 8b while Figure 9 below shows the summary of the apriori algorithm rules on the data set.



Figure 9: Summary of the apriori algorithm rules on the data set

In other to inspect the generated rules, figure 10 below inspection, shows that: 100% of the customer that purchased "SUGAR" also purchased "SET 3 RETROSPOT TEA", 100% of the customer that purchased "COFFEE and SET 3 RETROSPOT TEA" also purchased "SUGAR" Also from bar plot; "WHITE HANGING LIGHT T-HEART HOLDER" was discovered to be the most frequent item. Some exploration was made on it to see rule where customers are purchasing more of the item, and this

was carried out by reducing the support to 0.001, set the item at right hand side and maximum number of 3 item so as not to generate too many rules. Hence 36 rules were generated shown figure 10. Then lastly, carrying out data visualization, so many plot and graph were used which are made available by the arules package in R, the first graph was the scatter plot follows by the two key plot which shows a more way of using colors to add more visibility to items interactive plot which provides access to hover around the each rules and view all possible qualities, measures, support, confidence, lift of each item and attributes. Other visualization of the model plots carried out is shown from figure 11 to figure 16 respectively.



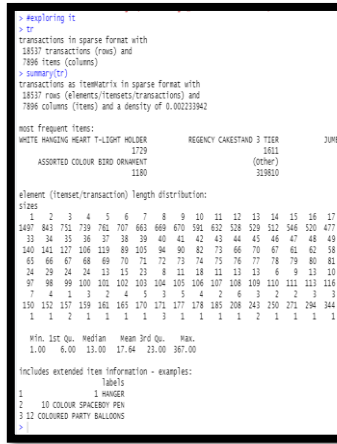Figure 6: absolute frequency plot which shows how frequent they appeared independently



Figure 7:Summary showing the most frequent items
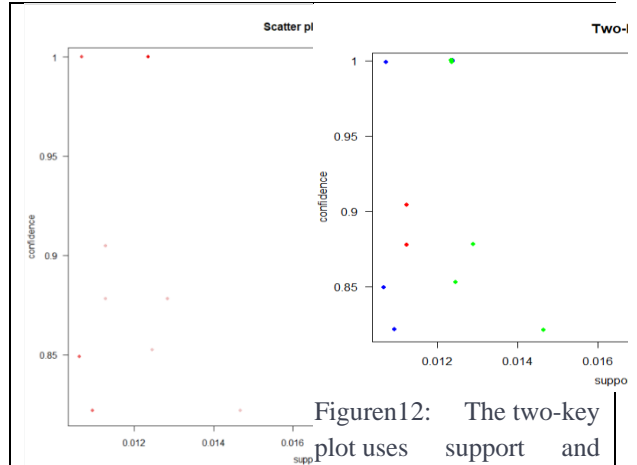


Figure 10: Generated Rules



Figure 11: showing scatter plot for 18 rules, showing rules that have minimum support have higher confidence



Figuren12: The two-key plot uses support and confidence on x and y-axis, respectively. It uses *order* for colouring. The order is the number of items in the rule.
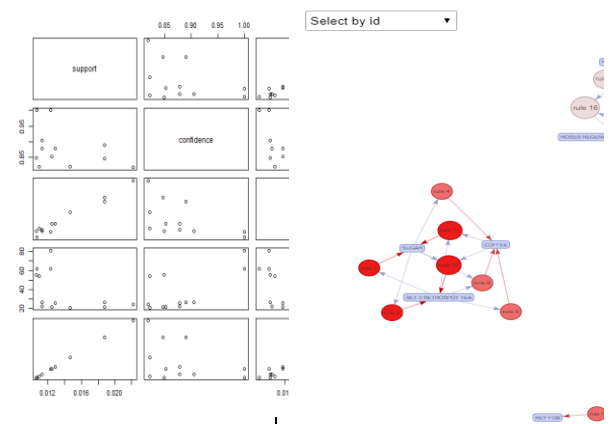


Figure 13: above is a scatterplot matrix to compare the support, confidence, and lift
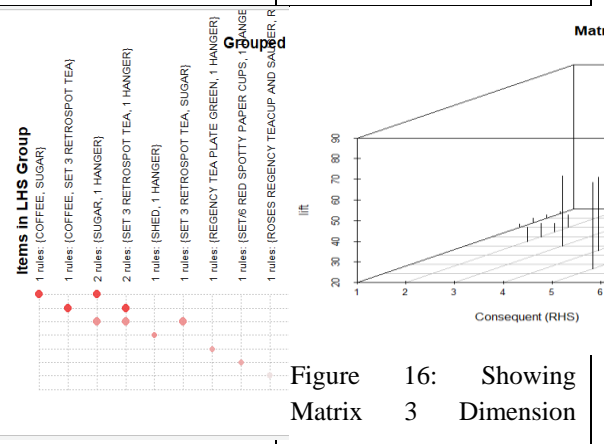


Figure 14: Showing the Graph-Based Visualizations



Figure 16: Showing Matrix 3 Dimension

| Figure 15: Showing Grouped Matrix | Plot/Graph |
| --- | --- |

ASSOCIATION RULE EXPERIMENTS ON THE DATASET USING SAS ENTERPRISE MINER

This is summary on how the experiment was done usingSaS enterprise miner software. The enterprise miner software was launched after which a project was created. Then the cleaned data set (Online_Retail) was imported by dragging and dropping from the import node from sample menu bar of the software then exploration of the dataset was done which shows the first 55 rows shown in figure 17 below whilethe propertiesand Sample statistics of the data set was shown below in figure 18 and 19 respectively.
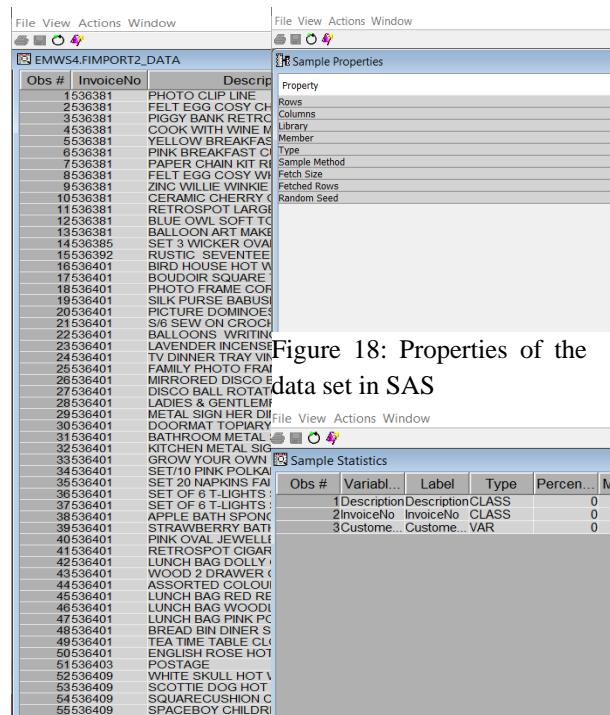


Figure 17: view of dataset showing the first 55 rows.



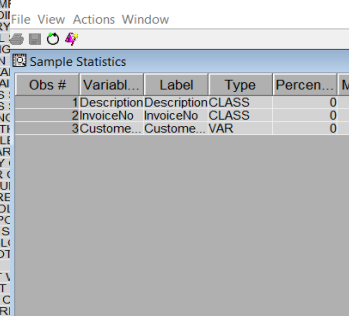Figure 18: Properties of the data set in SAS



Figure 19: Sample statistics of the data set in SAS

The data exploration further shows thefrequency histogram with the scale of 0 to 1500 of customerID, where customerID was segmented in range, figure 21 shows frequency Bar plot with the scale of 0 to 60 of description, showing "Rex cash+carry jumbo shopper" and "Paper chain kit 50's Christmas" as the two most frequent items and lastly figure 20 showing frequency Bar plot with the scale of 0 to 150 of

Invoice number, showing invoice "537781", "536796" as the two most invoice that contains the highest number of items in a transaction respectively. MODEL BUILDING USING SAS: Now to set target before application of the association rule algorithm, an editing of the variable to select the variables needed for specification of role (description = target, invoiceNO = ID, and customer id rejected) is required which will enable applying association rule algorithm on the data set by dragging association node from the sample menu bar on the enterprise miner and then connect the association node to the file import node, then proceed to generation of the rule, because using the entire data set to generate rules will give huge one, therefore there is need to set parameters as follows (supp(count)=100, conf=80, maxlen=2), and 49 rules was generatedwith the rule description while figure 22 shown below shows the rule table, with the confidence of 90.08%, support of 1.09%, and analysis can be made such as: 100% of the customer that purchased "REGENCY TEA PLATE PINK" also purchased "REGENCY TEA PLATE GREEN ", also with the confidence of 82.81%, support of 1.27%, analysis can be made such as: 100% of the customer that purchased "SET/6 RED SPOTTY PAPER CUPS " also purchased "SET/6 RED SPOTTY PAPER PLATES". By applying visualization, various graph plots were achieved such as Figure 23: statistics plot showing confidence against support, figure 24 showing rule matrix, RHS against LHS, figure 25 showing link graph visualizations and figure 26 showing the statistic line-plot of the rule. The achieved design model was shown in figure 27 below.
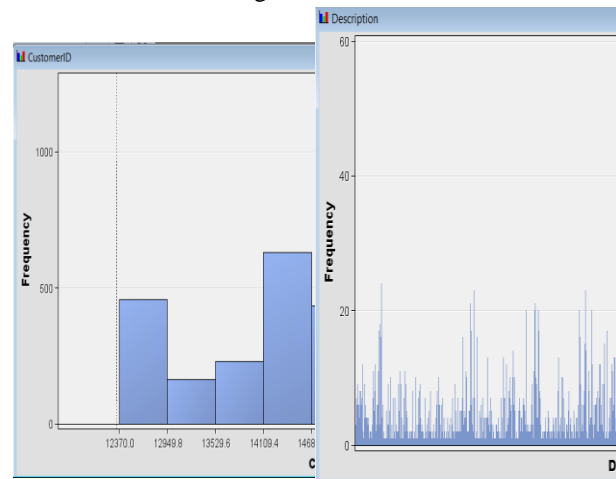


Figure 20: frequency



Figure 21: frequency Bar

histogram with the scale of 0 to 1500 of customerID, where customerID was segmented in range

plot with the scale of 0 to 60 of Description, showing "Rex cash+carry jumbo shopper" and "Paper chain kit 50's Christmas" as the two most frequent items respectively
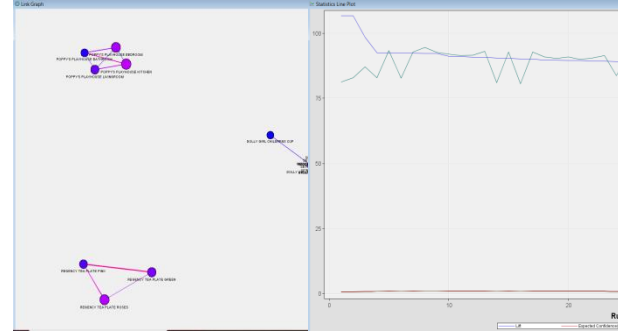


Figure 25: Showing Link Graph Visualizations



Figure 26: Showing the Statistic Line-Plot of the Rule



Figure 22: diagram showing the rules table;



Figure 23: statistics plot showing confidence against support.



Figure 24: showing rule matrix, RHS against LHS



Figure 27: Showing the designed model on the SaS enterprise miner Software

RESULT COMPARISM OF ASSOCIATION RULE BETWEEN R AND SAS

CRITICAL FINDINGS

1. The both tools predicted (PINK REGENCY TEACUP AND SAUCER) => (GREEN REGENCY TEACUP AND SAUCER) as their highest support with

   R= 0.02217187 (2.22%)

   SAS= 2.48%

2. In R, 7rules with confidence of 1 =100% was generated while

   SAS has the highest confidence of 94.62%

3. In R, highest lift was 80.94760 with 4 rules while

SAS was 106.81 with 2 rules

4. Comparing parameter in both tools.

R has parameter set with maxlen=10, Min_supp=, MIN_conf=, where R generated only (8 rules with maxlen=3) and 2 rules were generated with maxlen=4. While,

SAS generated 49 rule parameter = maxlen=2, min_support_count = 100, min_confidence=80%

5. R is a very interesting language which has the ability for data manipulation. Hence it is easy to clean a data set in R while in SaS enterprise miner is not which is one of the disadvantages of using SaS in data prediction.

6. Results generated with R is more accurate than in SaS because R provided more ability to explore data to any point in other to achieve a perfect and more accurate results than in SaS.

7. Parameters can be sets for more clarifications of results in R than in SaS

APPLYING CLUSTERING RULE ALGORITHM USING R AND SAS

Clustering algorithm uses an approach where customer's records in the database are grouped and hence creates segments of the data within the group based on data similarity, then after which each segment can be treated differently depending. Figure 1 above illustrates the clustering rule algorithm which is used to analyze customer's records in the database.

EXPERIMENTS ON THE DATASET USING R

The main grail of this in R is to find the clusters of customers with the amount likely to be spent, how frequent and recent they are likely to buy based on customers purchase records, so phasing the clients based on RFM (Recency, Frequent, Monetary) which have 4339 rows and 4 column in order that the business

enterprise can target its customers efficaciously. The experiments where done and a model was developed. The approach towards the experiment includes, loading of the required packages, loading the dataset to R dataframe, exploring the data which when the first six rows) and the last six rows are viewed shows the data-types of the variables, the dimension of the dataset, and also used data-explorer package to see the frequency bar plot of the "country" variable shown in figure 28 below, data cleaning is another stage, data preparationvisual exploration of prepared dataset. After a successive exploration of the dataset,

figure 28, 29, 30, 31, 32, 33 and 34below shows the result after the data exploration.
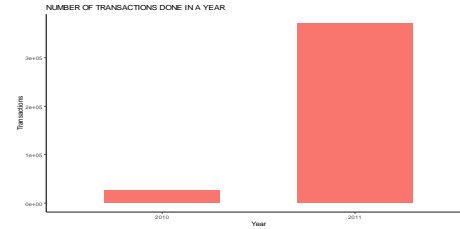


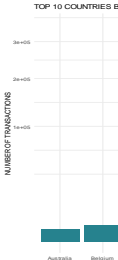Figure 28: Bar graph showing plot of number of Transactions done in a year



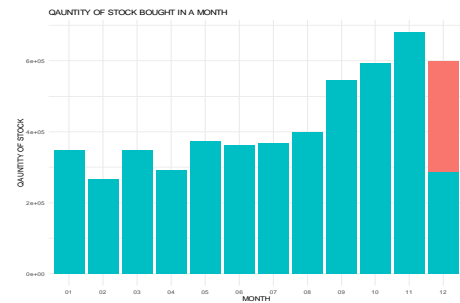Figure 29:



Figure 30: highest num



Figure 31: Plot showing the quantity of stock bought in a month within a year
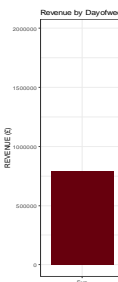


Figure 32: basic in a v



Figure 34:

Figure 33: Plot showing revenue generated in a month of a year

CALCULATING RFM

To get the RFM, where R is the closest time to the last day of transaction or the minimum difference of days calculated in STEP 14, F is the frequency of observation/row of customerID or how many times the customer had a transaction within, and M is the summation of the totalsales of a customer or adding up all the money generated by the customer and conversion of the recency to numeric by replacing all n/a observations with zero (0) shown in figure 35 and 36 respectively.



Figure 35: Showing the summarized group by customer id to get the RFM



Figure 36: output after conversion of recency to numeric

APPLYING THE CLUSTERING ALGORITHM

In other to determine optimal number of cluster using three (3) methods namely, silhouette method, gap statistic method and elbow method shown in figure 37, 38 and 39 below.



Figure 37: (SILHOUETTE METHOD) From above ofgraph it evidently shows that k=2 is the Optimal number of Cluster

Figure 38: showing OPTIMAL number of clusters with gap statistics method (which indicate k=5, hence we can determine that K=5 is the Optimal Cluster.)



Figure 39: With the Elbow method, the graph starts to bend cluster 2, therefore we can ascertain that k=2 is the optimal cluster

Then the visualization of the clusters with the optimal clusters of the 3 methods was achieved using k-means clusters on a k=2, k=5 for better understanding shown in figure 40 and figure 41 below respectively.



Figure 40: showing K-means cluster, where k=2

Figure 41: showing K-means, where k=5

EXPERIMENTS ON THE DATASET USING SAS

The data set was first cleaned in and then prepared using R by segmenting customers to determine clusters of customers with the amount likely to be spent, how frequent and recent they are likely to buy (RFM) and that is the objective. Next is to normalize and save it in the working directory. This is done because SAS cannot clean the data set before it could be used to carry out another experiment in SAS. After a successive data exploration, the clustering rule algorithm was done by dragging clustering node from the explore menu bar on the enterprise miner software. Then connect the clustering node to the file import node which shows the victuals on figure 42 while figure 43 shows the mean statistics of the dataset using SAS enterprise miner and figure 45 shows index screenshot of the generated model on the enterprise miner softwarebelow:



Figure 42: showing result generated at specification method set at automatic



Figure 43: showing the mean statistics

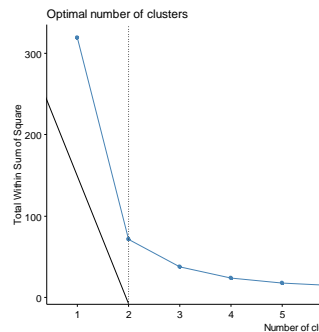| Clustering Criterion | Maximum Relative Change in Cluster Seeds | Improvement in Clustering Criterion | Segment Id | Frequency of Cluster | Root-Mean-Square Standard Deviation | Maximum Distance from Cluster Seed | Nearest Cluster |
|---|---|---|---|---|---|---|---|
| 0.899239 | 0 | 0 | 1 | 4 | 4.87603 | 10.13493 | 2 |
| 0.899239 | 0 | 0 | 2 | 4335 | 0.890565 | 34.32612 | 1 |



Figure 44: showing the result after the model was applied and summary of the variable

INDEX



Figure 45: Index screenshot of the generated model on the enterprise miner software

CONCLUSION

Every business wants to satisfy its customers and also provide a free and easy access to her customers' needs remotely, therefore understanding your customers behaviors is one of the best approach in pro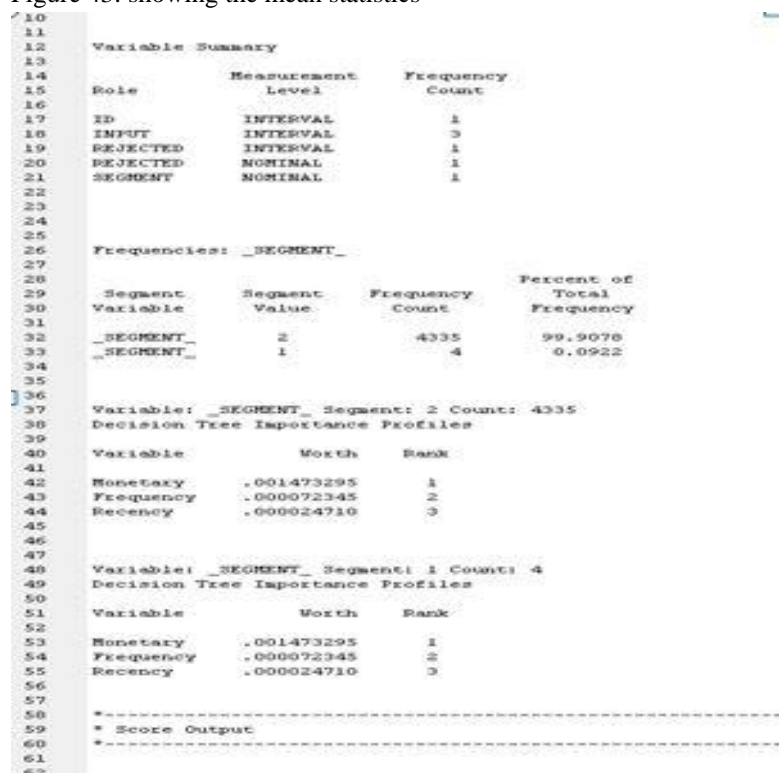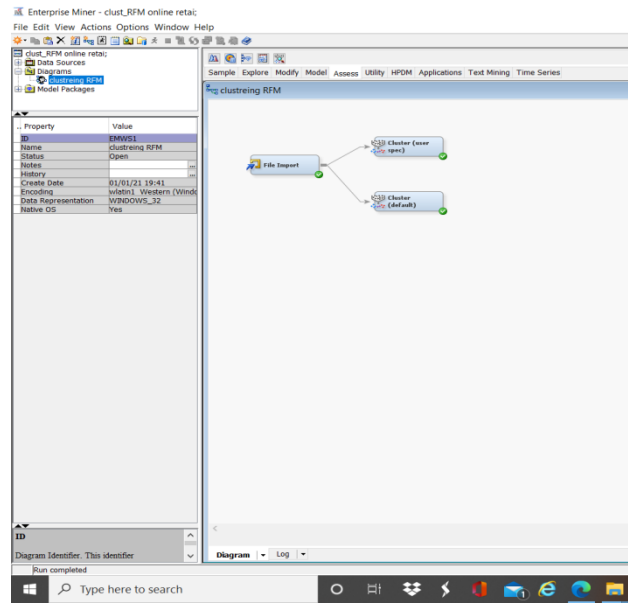viding the goods and services they require and it is through their patronize the growth of the organization could be guaranteed. This study tends to predict customers behavior on an online retail system by using association and clustering machine learning algorithms a (comparative analysis using SAS and R languages). The analysis was first conducted using R language which showed values of observation of 397924 and 8 variables. The model was built by using apriori algorithm and it shows in figure 10 that: 100% of customer that purchased "SUGAR" also purchased "SET 3 RETROSPOT TEA",    100% of the customer that purchased "COFFEE and SET 3 RETROSPOT TEA" also purchased "SUGAR" and also "WHITE HANGING LIGHT T-HEART HOLDER" was discovered to be the most frequent item purchase by customers meaning that the organization can go ahead to produce more of the identified products than others. When the SaS enterprise miner software was used to perform the same association by using the same dataset, the result

shown in figure 22 showsthe rule table, with the confidence of 90.08% and support of 1.09%, and analysis was made such as: 100% of the customer that purchased "REGENCY TEA PLATE PINK" also purchased "REGENCY TEA PLATE GREEN ", also with the confidence of 82.81%, support of 1.27%, analysis was made such as: 100% of the customer that purchased "SET/6 RED SPOTTY PAPER CUPS " also purchased "SET/6 RED SPOTTY PAPER PLATES".While figure 27: shows the designed model on the SaS enterprise miner Software whit the developed dashboard design.  On the other hand, by application of clustering algorithm on both R and SaS revealed another predictive model that can help various business organization progress and understand their customers behavior and hence improve on their service delivery. To achieve this, R was first used to discover the clusters of customers with the amount likely to be spent on a product, how frequent and recent they are likely to purchased based on customers purchase records by calculating the customers RFM (Recency, Frequent, Monetary) status as the predictive formula and when the K-means model was applied on the dataset figure 35 shows the summarized group by customer id to get the RFM and figure 36 shows output after conversion of recency to numeric with number of customers recency, frequent and monetary on a particular product while application of SAS enterprise miner on the dataset shows a result on figure 44 after the model was applied and summary of the variable.

RECOMMENDATION

Machine learning models are part of artificial intelligence that helps developers predict or forecast an event and make intelligent decision base of knowledge discovery from data. This study after a comparative analysis on an online retail System as a dataset to predict customers market behavior therefore make the following recommendations:

1. Business owners should explore the intelligent information on this study for an effective understanding of their data to day business activities and also on how to handle their customers' needs always.

2. In other not to loss in business activities especially those involved with online transactions, customers comments/feedback is highly needed for effective and smooth running of the organization and hence

provide those required products and series needed by their customers.

3. It is very good for academic scholars to explore other data mining tools for more accurate prediction and development of intelligent models as this study has shown us the importance and huge difference between R language and SAS enterprise miner for data scientist and data exploration and visualization.

REFERENCE

1. MeenuSharma(2014)Clustering In Data Mining : A Brief Review;
International Journal Of Core Engineering & Management (IJCEM) Volume 1, Issue 5, accessed from
file:///C:/Users/IPROSO~1/AppData/Local/Temp/Data_Clustering_Using_Data_Mining_Techni.pdf

2. Pravarti Jain AndSantosh Kr Vishwakarma (2017)A Case Study on Car
Evaluation and Prediction: Comparative Analysis using Data Mining Models, International Journal of Computer Applications (0975 – 8887) Volume 172 – No.9

3. Zhong and L. Zhou (1999): PAKDD'99, LNAI 1574, pp. 13–23, 1999. Springer-
Verlag Berlin Heidelberg.

4.Charu C. Aggarwal and Philip S. Yu (2012) Data Mining Techniques for
Associations, Clustering and Classification, IBM T. J. Watson Research Center, Yorktown Heights, NY 10598                accessed                from
https://link.springer.com/article/10.1057/dbm.2012.17

5.Aayushi Maheshwari, Garima Kharbanda and Harsh Patel(2015)Association Rules
in      Data      Mining,      accessed      from
Data_mining_for_the_online_retail_industry_A_case_study_of_RFM_model-
based_customer_segmentation_using_data_mining