# GROUPING CUSTOMERS USING CLASSIFICATION ALGORITHMS

*Ravi teja Yamsani[1]*

*Student, Department of Computer Science and Engineering, Kalasalingam Academy of Research and Education, Anand Nagar, Krishnankoil, India*

*Dr. A. Saravanan*

*Associate Professor, Department of Computer Science and Engineering, Kalasalingam Academy of Research And Education, Anand Nagar, Krishnankoil, India*

*Abstract* - Meeting of many competitors at once entrepreneurs have created a major conflict between competitors businesses to find new customers and keep old ones. As a result formerly, the need for unique customer service becomes relevant regardless of the size of the business. In addition, the ability of any business to understand the needs of each customer will supply a larger customer support for the provision of targeted customer services and development customized customer service plans. This understanding possible through systematic customer service. Each part has customers who share similar market features. Big data ideas and machine learning have encouraged greater acceptance of automatic customer segregation methods in favor traditional market statistics that often do not work there the customer base is very large. Customer segregation is important in customer relationship management software. The most common way to differentiate one customer from another is to develop a customer group according to their interests and factors like demographic factors, Psychological factors, Behavioral factors, and geographical factors . In this work, the company's hand-separated customer data is analyzed. As machine learning methods are useful for solving problems with data management and data processing, the solution is searched within machine learning methods. The different classification methods are used to differentiate the people into different groups and help with suitable solutions. For increasing the efficiency and precision we used all the classification algorithms as a ensemble algorithm and created a model.

## I. INTRODUCTION

From many years, Due to the rise in competitive nature among the businesses and the use of the past data has motivated in implementation of data mining techniques for finding the critical and useful insights from the information. Data mining procedure is helpful for extraction of the logical information from the dataset of the company and project in a human understanding manner such that human can take required pre requisite to solve the problem.



**Figure:** Gartner company proves that segregating of customers brings 2nd highest sales

## II. LITERATURE REVIEW

In literature survey we discussed various machine learning techniques have been proposed by the scientists for identification and grouping the customers into different clusters. This research study shows some existing machine learning based techniques in order to explain importance of the implemented work. The market separation or grouping is based on only few factors for determining the customer group. Over the years, the commercial world has become more competitive, as organizations such as these have to meet the needs and desires of their customers, attract new customers, and thus improve their businesses. For amazon, flip kart, and Netflix the customer grouping is done by item to item and only built by the data received by their website. All their systems algorithms are confidential and not shared with the all but only the basis of their fundamental are shared for privacy concern of the customers. Few data scientists did the project using KNN algorithm for separation but the data considered for the project is likely small and the accuracy is about 78.3%.. Grouping also includes following factors like Demographic, geographic, technographic and behavioural segmentation. Data collection is part of the research in all fields of study including natural and social sciences, people and businesses. The purpose of all data collection to obtain quality evidence leading to analysis creates straightforward and misleading answers to questions presented. Integration is the process of combining information into a Database based on similarities. There are several algorithms, which can be used in data sets based on given condition. [7] However, there is no general

consolidation. An algorithm exists, so it is important to choose appropriate methods of integration. Customer segmentation used to be done manually and wasn't too precise. They would manually create and populating different data tables, and analyze the data.

## III.EXISTING SYSTEM

In the existing system only KNN algorithm is used for the project which cannot define all the essential factors or features for efficient precision. Recently, Big Data research has gained momentum. Where big data describes a large number of formal and informal data, which cannot be analysed using traditional methods and algorithms. No proper data visualization is used in the existing system such that only tech geek can understand the results from the implementation.

## IV.PROBLEM DEFINATION

The marketing methods able to get results by using the traditional methods by the assumption that every customer buy whatever they need without going to the other shops and such that only such type of customers are attracted to the business and the products availability will be limited. For instance if we make groups and assign customers into those categories such that where one size fits method. This can be obtained through only the successful segmentation. The segmented customers are grouped into clusters where all of them almost likes the same things and think similarly, This helps to drive dynamic content and personalization tactics for increase in the business. However, in order for the division to be used effectively, it needs to take into account that different customers are buying for different reasons, and retailers need to use wisely a lot of considerations that could affect their buying decisions. A Harvard Business School professor even said that of the 30,000 new products introduced each year, 95% fail because of poor market performance. There are numerous reasons why the company fail in selecting their target group. It is a wild guess that segmentation is only done on the one factor like demographics or the cluster is too broad to be under the one roof. There may be no strategic goal, such as lead discovery or customer retention. Or, it may be that a lack of One Customer View and aggregated data has led to a misunderstanding or misinterpretation of your customer base. So for segmentation in this project we use almost all factors like demographic, geographic, psychographic and Behavioural.

## V.PROPOSED SYSTEM

In Proposed system we Classify customers into segments and Anticipate the purchases that will be made by a customer, during the following year.For classification, I took data from Kaggle. After pre processing the data we will visualize the data and its factors in a graph model using matploit and seaborn libraries. To handle the big data problem I used all the classification algorithms and selected best classifiers in each algorithm and made the prediction by combining all the parameters

For better understanding even for the non tech geek I used word cloud and made the visualizations of each group as samples.In this project large data set is used such that the results which are produced are not biased and the system which is proposed can be used for real time purpose. Ensembling the algorithms and producing the results make the prediction more accurate. The geographical parameters are also took into consideration before classification..

## VI. SYSTEM DESIGN

**Proposed Model for Customer Grouping:** The proposed model helps us in grouping only after providing the suitable data for retrieving the information. In this model the company provides the data to the model and then perform the ensemble model on the data for segregating the customers for increasing their sales. The database used here for machine learning comes from sales of the UK retailer business. The goal of segregation is to increase the sales of the company and attract new customers by showing the their needs in the websites. The data set consists of above 5lakh rows and 8 columns.
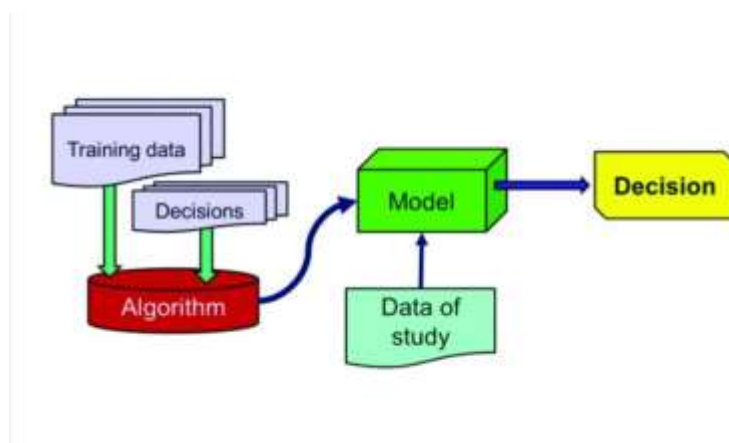


**Figure:** Training and Decision process of ML Model.

# VII. METHODOLOGY

**Classification Algorithm**
The attributes mentioned in Table 1 are presented as inputs to various ML algorithms such as Random Forest, Decision Tree, Logistic Retreat, Naive Bayes, XG Boost and Light GBM split strategies. The set of input data is divided into 80% of the training database and the remaining 20% of the test database. Database training is a database used to train a model. The test database is used to test the performance of a trained model. In each algorithms the performance is calculated and analyzed based on the various metrics used such as accuracy, precision, recall and F scoring scores as described further. The various algorithms tested in this paper are listed below.

The ML process starts in the pre-processing phase of data followed by feature selection and feature rating and we will use it with the following algorithms such as KNN, SVM, Decision trees, XGBOOST , Support vector Machine  and Random Forest algorithms.

**Logistic Regression :** Regression usually refers to the analysis of segmentation problems, despite your regression name in the segmentation algorithm. It can be used for both binary classification and multiple classification. The predictable rate in depreciation is by category. The sigmoidal function is used to reduce the value from a range of 0 to 1.

**XGBOOST:** The gradient boosted trees technique is implemented in XGBOOST, a popular and efficient open-source implementation. Gradient boosting is a supervised learning approach that combines the estimates of a collection of smaller, weaker models to attempt to accurately predict a target variable. The data is built in a sequential manner, with each successive data aiming to minimise the entire data's faults. Extreme Gradient Boosting is abbreviated as XGBoost. The word xgboost, on the other hand, alludes to the technical objective of pushing the computational resources for boosted tree algorithms to their limits. XGBoost is a software library that you may download and install on your computer, then use from a number of different interfaces.

**Support vector machines:** SVM selects extreme points / vectors that help create hyperplane. These extreme cases are called supporting vectors, which is why the algorithm is called Vector Support Machine. Consider the diagram below in which there are two distinct categories divided by resolution or hyperplane.

**Random Forest:** Random forest are built from decision trees. Generally, tree works well with the  data they have or familiar but fails to classify new samples. Random forest adds up the simplicity of the decision tree with flexibility resulting in the large improvement and in the accuracy. At each time we take the bootstrapped data and consider only a subset of the variables at each step which results in a wide variety of trees by applying the decision tree approach. The final result depends on the forests (multiple trees) that we created with the variety of trees. The remaining data after bootstrapping at each step is called out-of-bag data. The correct classification of the out-of-bag data is solemnly responsible for the accuracy of the random forest algorithm. It is an example of the bagging technique.

**KNN:** K-means the algorithm is one of the most popular segmentation algorithms. This integration algorithm depends on it centro, where each data point is placed in one of scattered, pre-filtered in the K-algorithm. Creating collections corresponding to hidden patterns in data that provides the information needed to help determine execution. process.
Decision Tree: In Decision tree algorithm the data is split by using the conditions and create the classes for each branch. These trees can continuously take the real number values it can also called as regression trees. In this algorithm the identification of the root node parameter is difficult. The process of identifying is called attribute selection. For this we need to take 2 factors under consideration they are information gain and Gini index.



**Figure** : Market Grouping Factors.

3

## IX. EXPERIMENTAL RESULTS AND DISCUSSION



**Figure 8.1:** Importing the essential dependencies



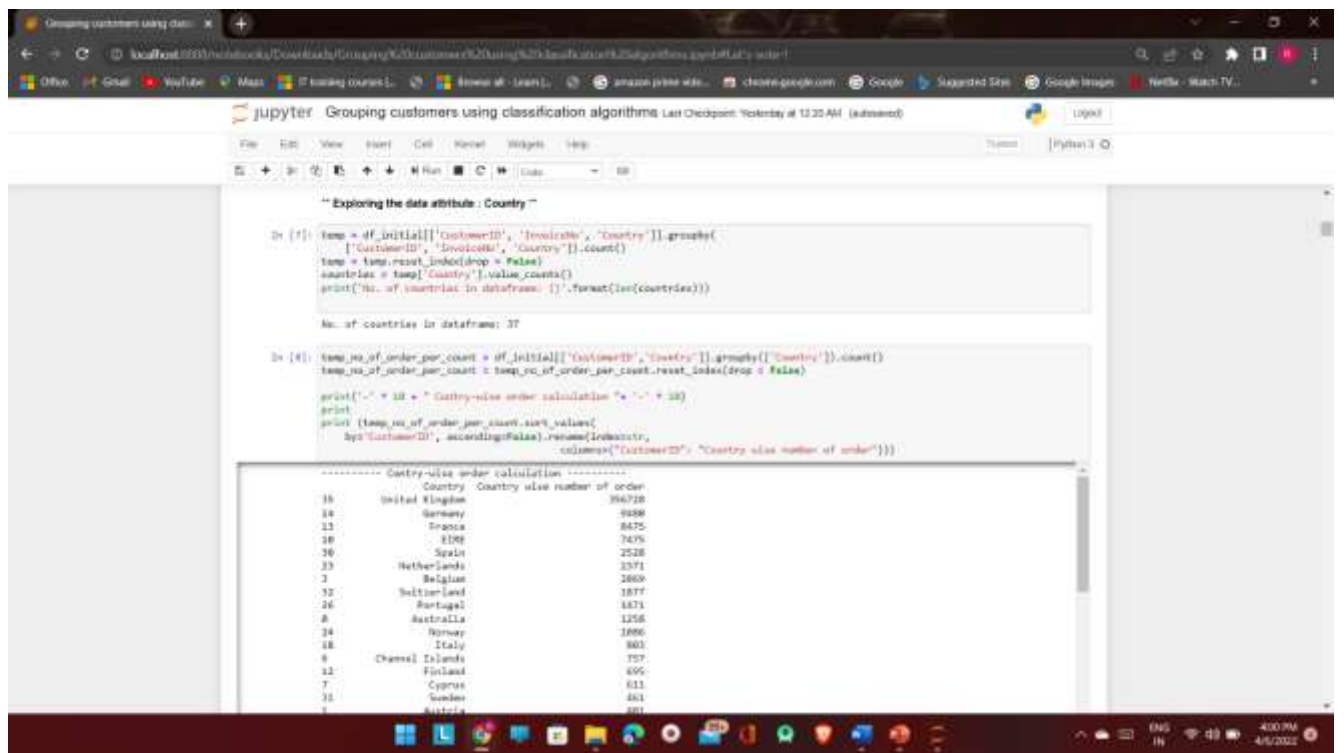**Figure 8.2:** Data preprocessing and data cleaning

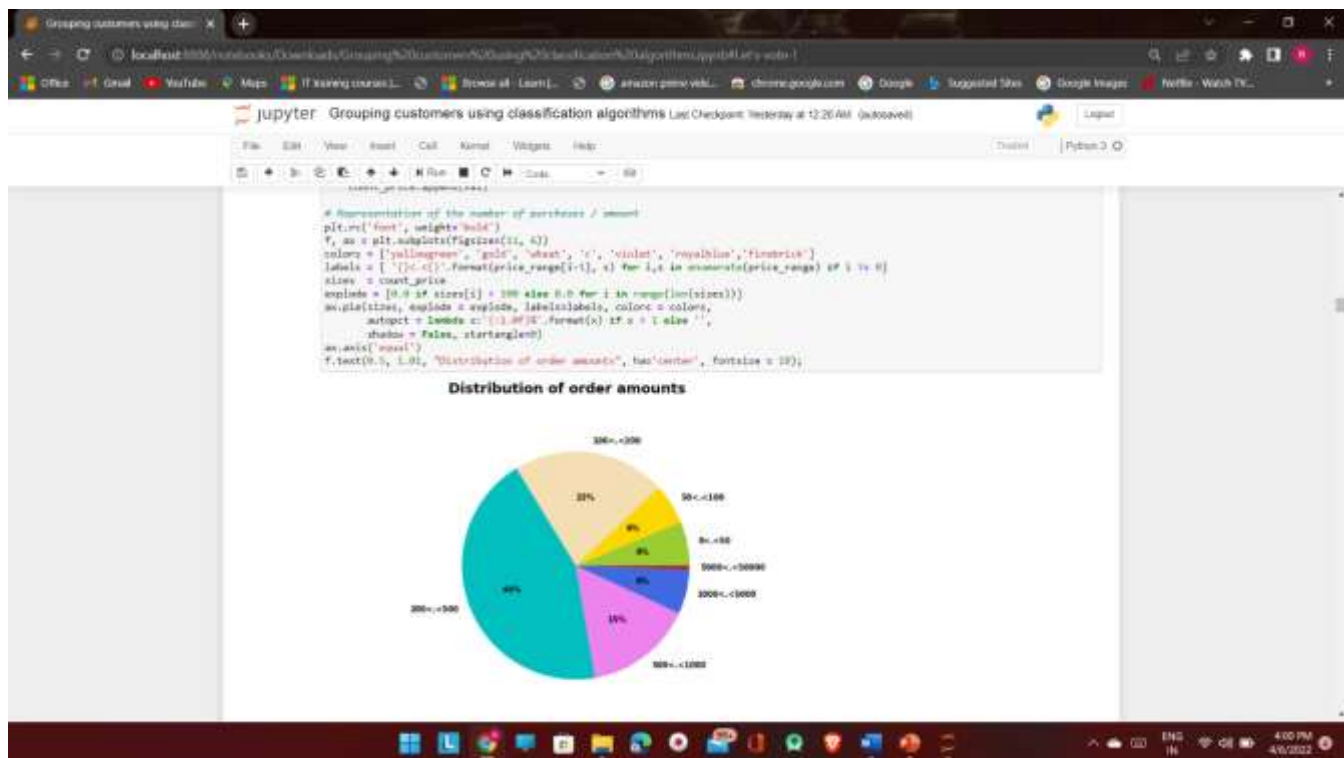**Figure 8.3:** Exploring the data and geographical purchases



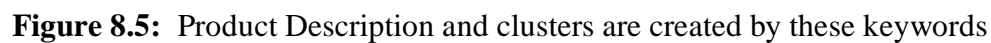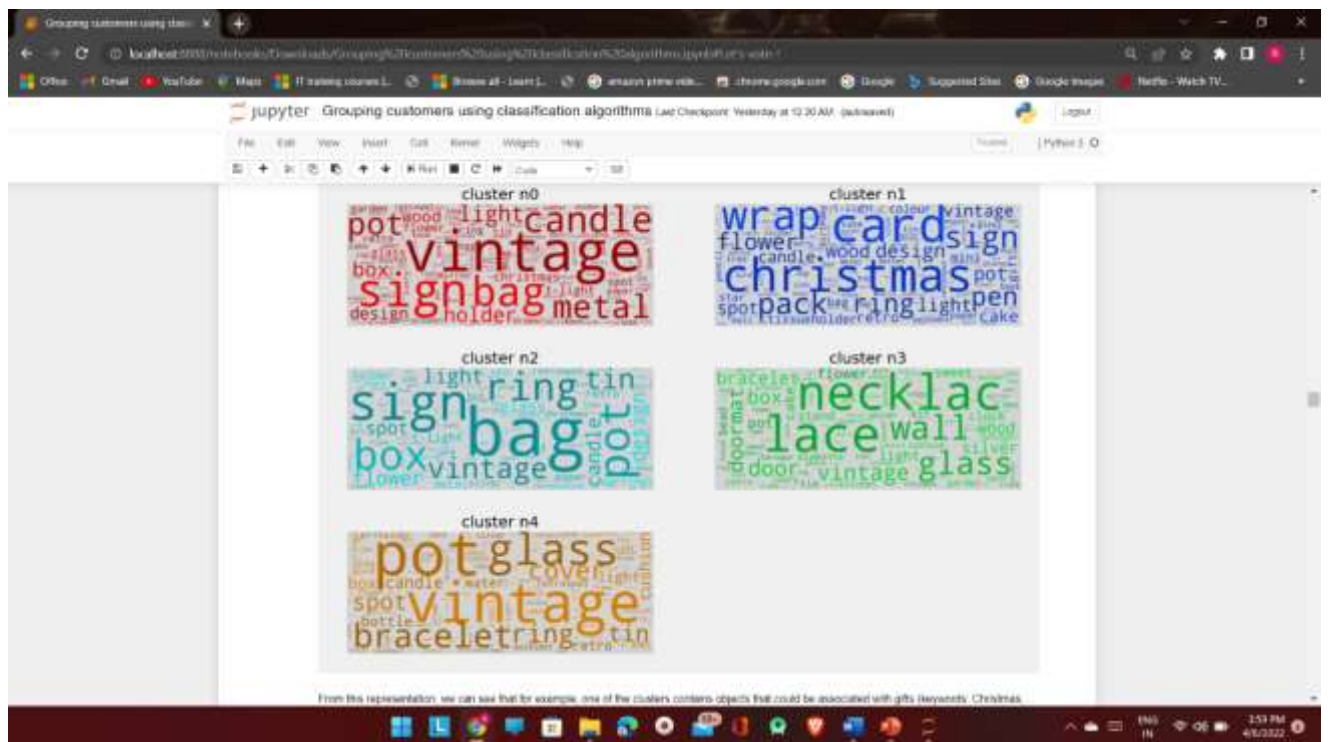**Figure 8.4:** Flow chart of the transactions and order of amounts.

**Figure 8.5:** Product Description and clusters are created by these keywords



**Figure 8.6:** Silhouette score of the cluster groups are determined.

6

**Figure 8.7:** All the customers are classified by words and visualized



**Figure 8.8:** After the customers are grouped now we should group the prodcuts

**Figure 8.9 :** Next I created the combinations of the customers transactions made\



**Figure 8.10 :** Data encoding turning words into numerical values for processing
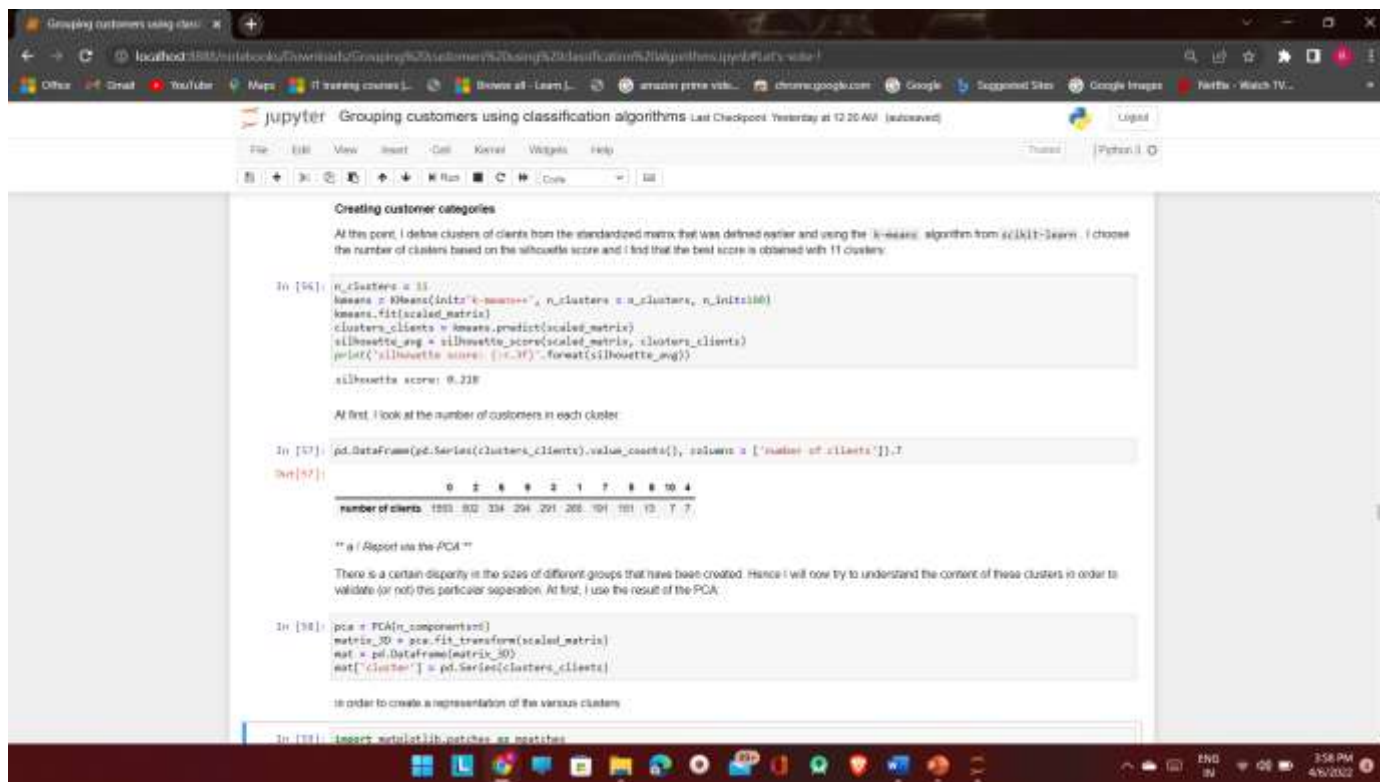
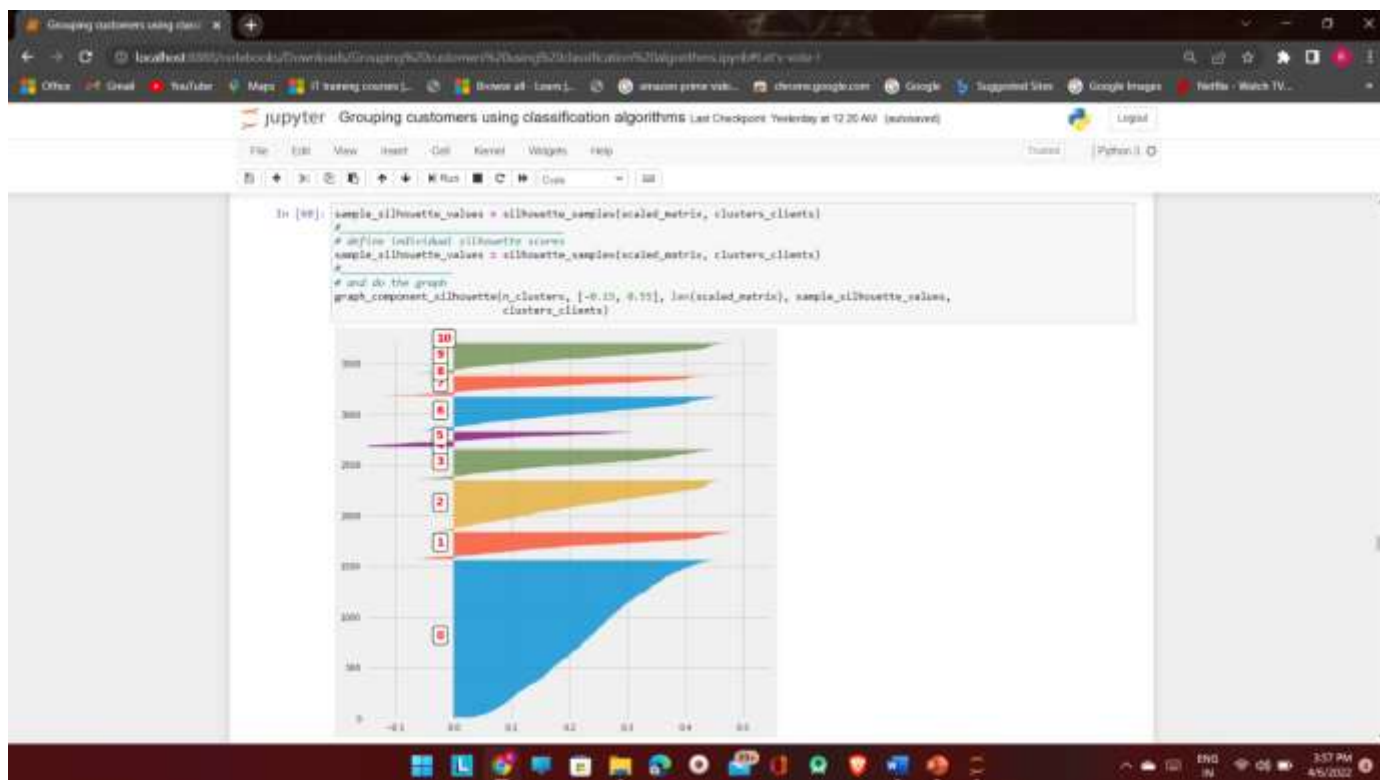**Figure 8.11 :** Customer categories created.



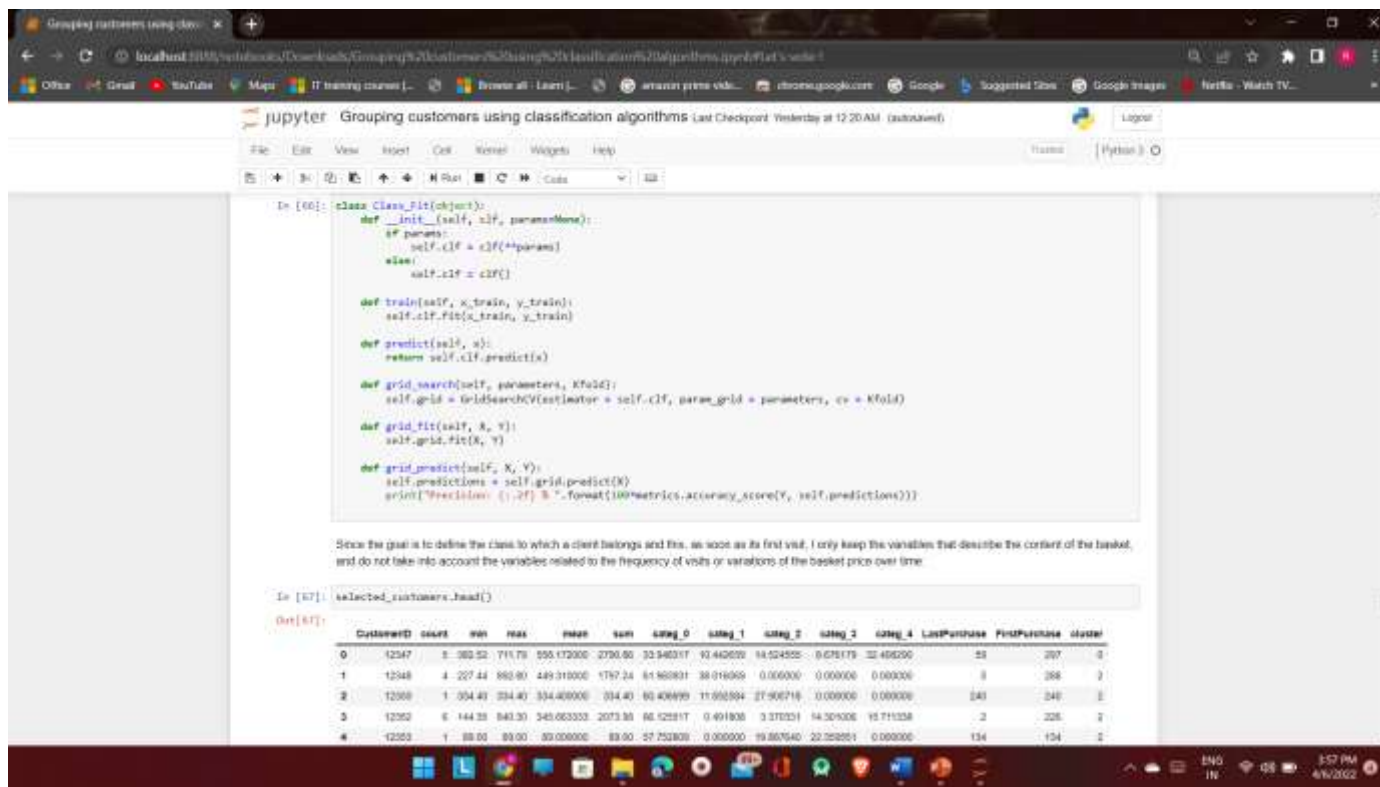**Figure 8.12 :** Customer categories are visualized

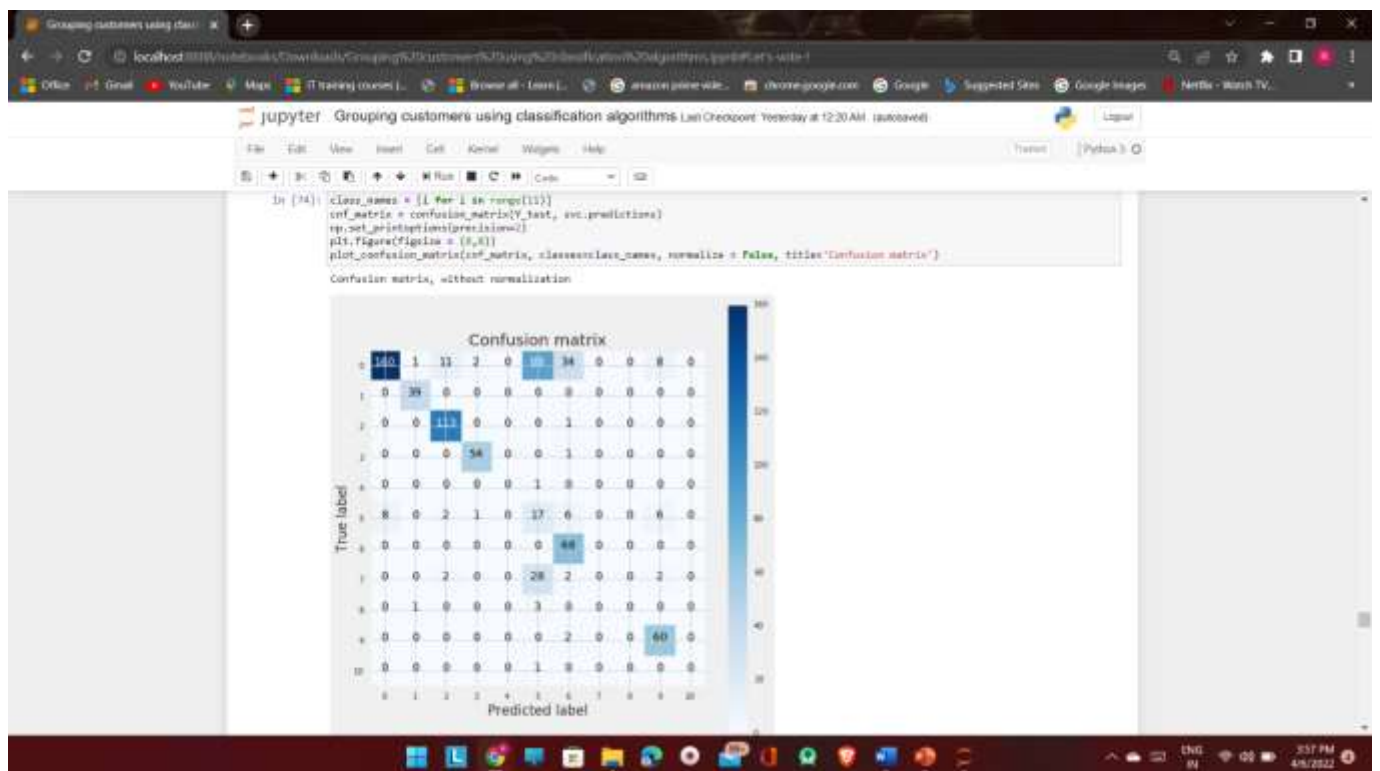**Figure 8.13 :** Program for creating the categories display and enhancement



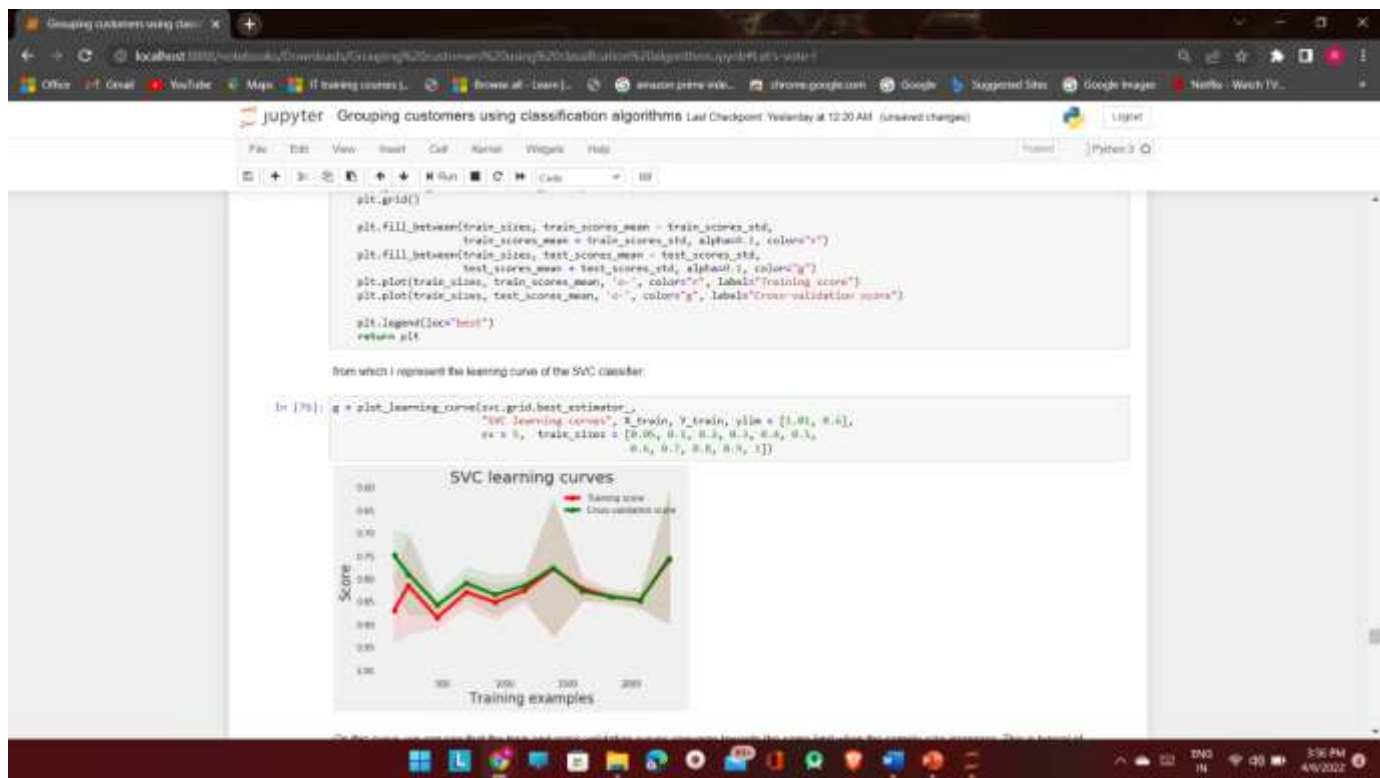**Figure 8.14 :** Displaying the confuse matrix of the MLmodel

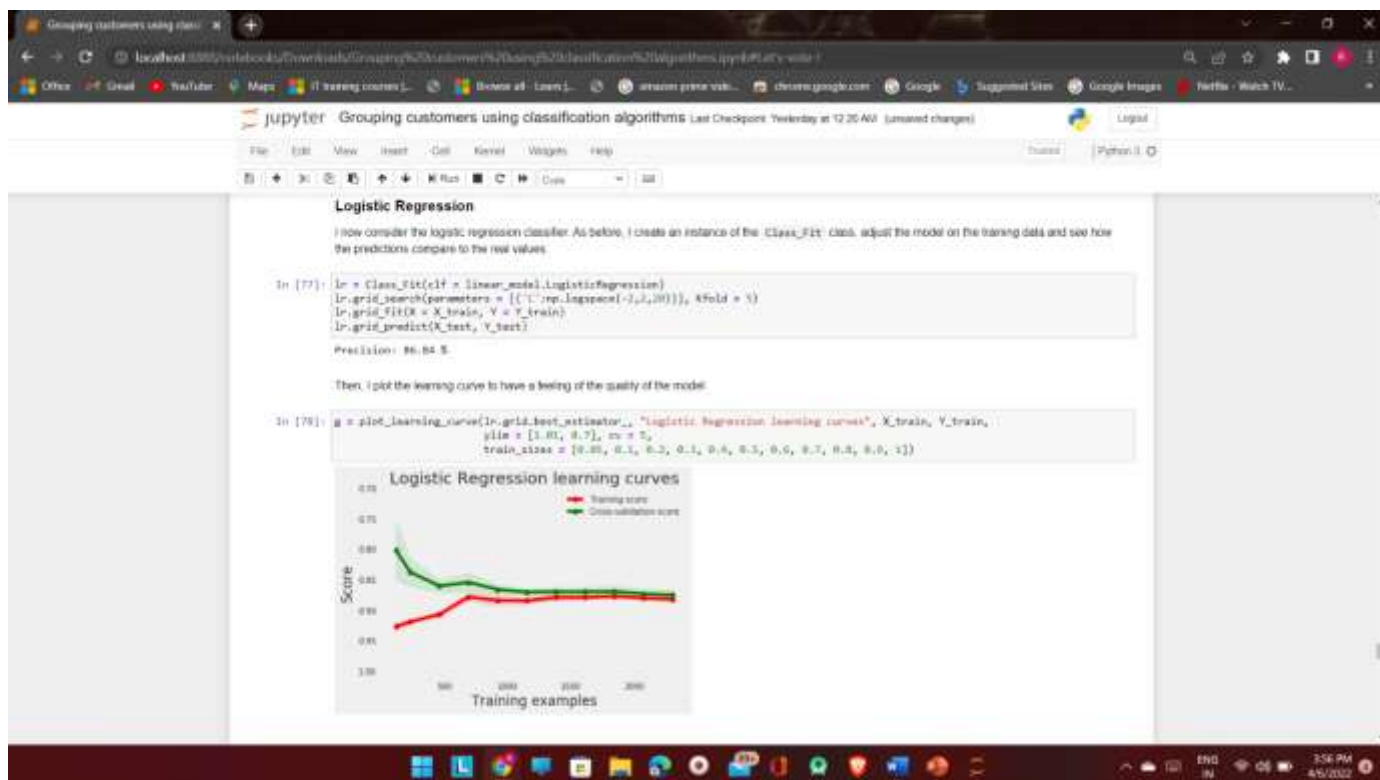**Figure 8.15 :** Accuracy of the SVC.



**Figure 8.16 :** Accuracy of the Logistic regression

**Figure 8.17 :**  Accuracy of the KNN



**Figure 8.18 :**  Accuracy of the Decision tree.

12

**Figure 8.19:** Accuracy of the Random Forest



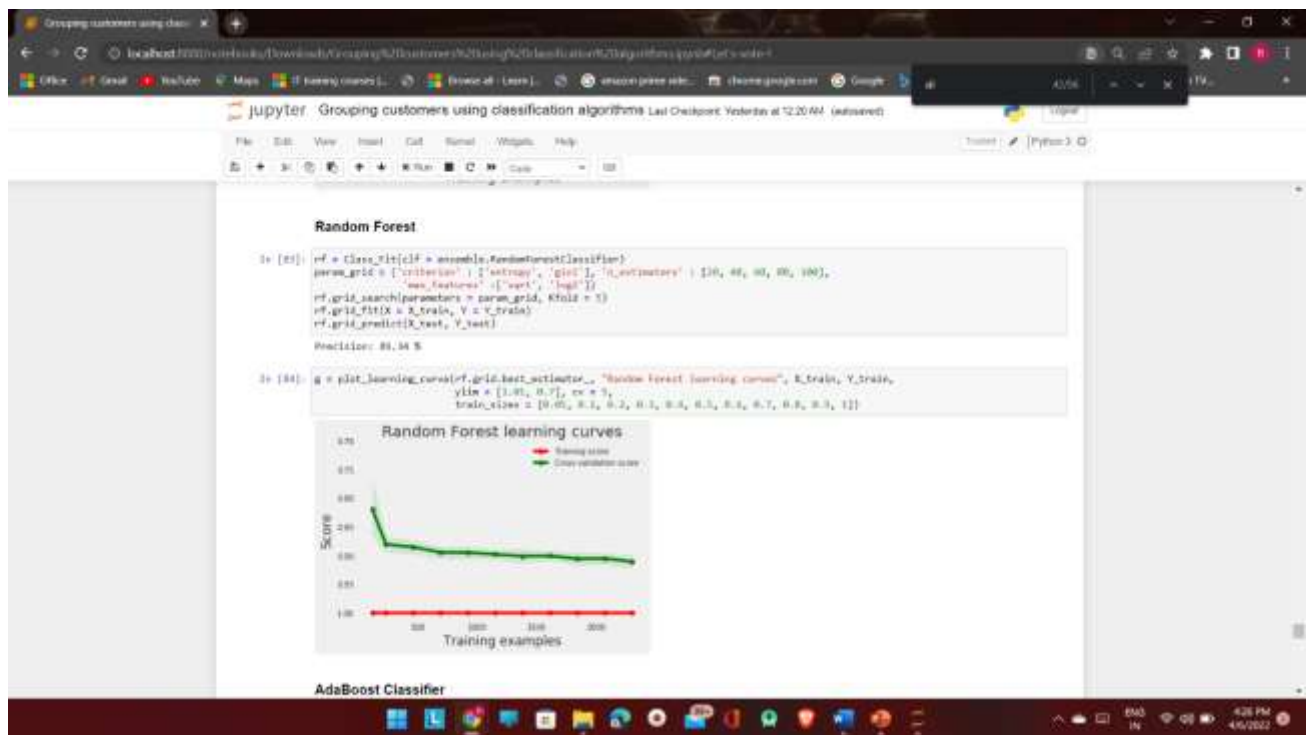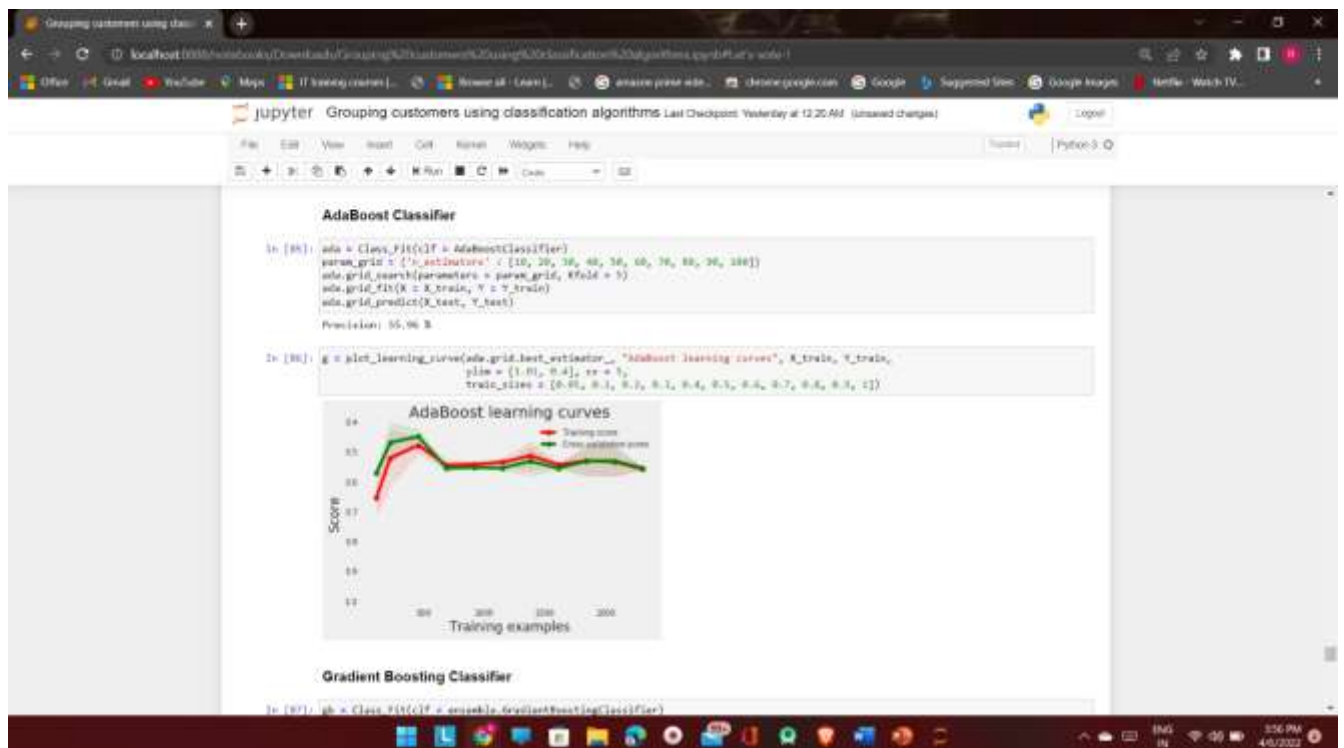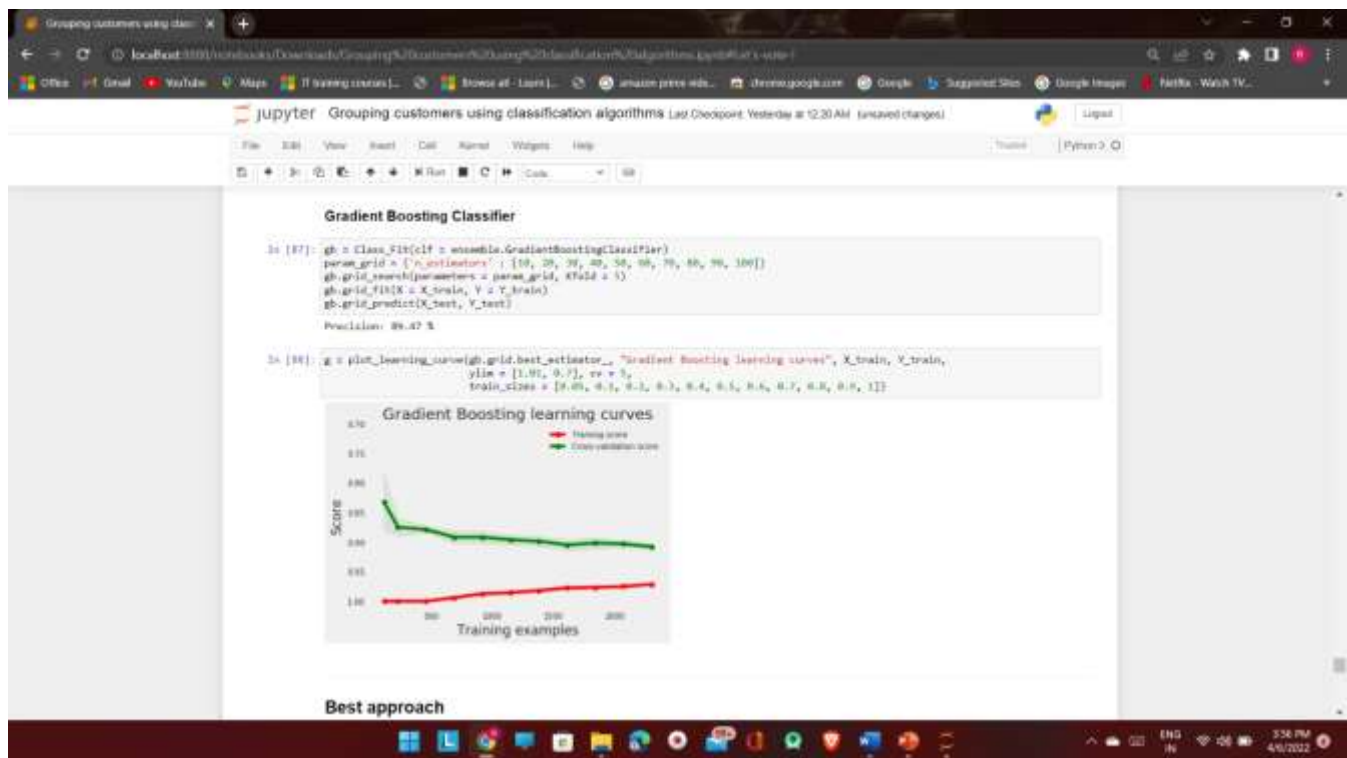**Figure 8.20 :** Accuracy of the Ada boost

**Figure 8.21 :** Accuracy of the XGBOOST



**Figure 8.22 :** Accuracy of the Ensemble model in the training phase

**Figure 8.23 :** Testing phase where the data is took from the last 2 months transactions of dataset



**Figure 8.24 :** Accuracy of the   all the algorithms and the ensemble model in the testing phase.

## X. CONCLUSION

Since our database was inconsistent as it is a real time data , we need to have a good computing power to run the billions of data. At last we provided the testing phase results where the ML model achieved about 90% accuracy. This is achieved due to the ensemble model learning where it chose the best parameters from each classification algorithm and used as the parameter. The visualization of the clusters is done by the word clouds in the boxes such all the customers who belong to one category has similar interests for the products. In the training phase we can observe that XG boost gave the more accuracy comparatively to other algorithms.



**Figure:** Displaying all the categories customers.

# XI. REFERENCES

[1] Blanchard, Tommy. Bhatnagar, Pranshu. Behera, Trash. (2019). Marketing Analytics Scientific Data: Achieve your marketing objectives with Python's data analytics capabilities. S.l: Packt printing is limited [2] Griva, A., Bardaki, C., Pramatari, K., Papakiriakopoulos, D. (2018). Sales business analysis: Customer categories use market basket data. Systems Expert Systems, 100, 1-16.
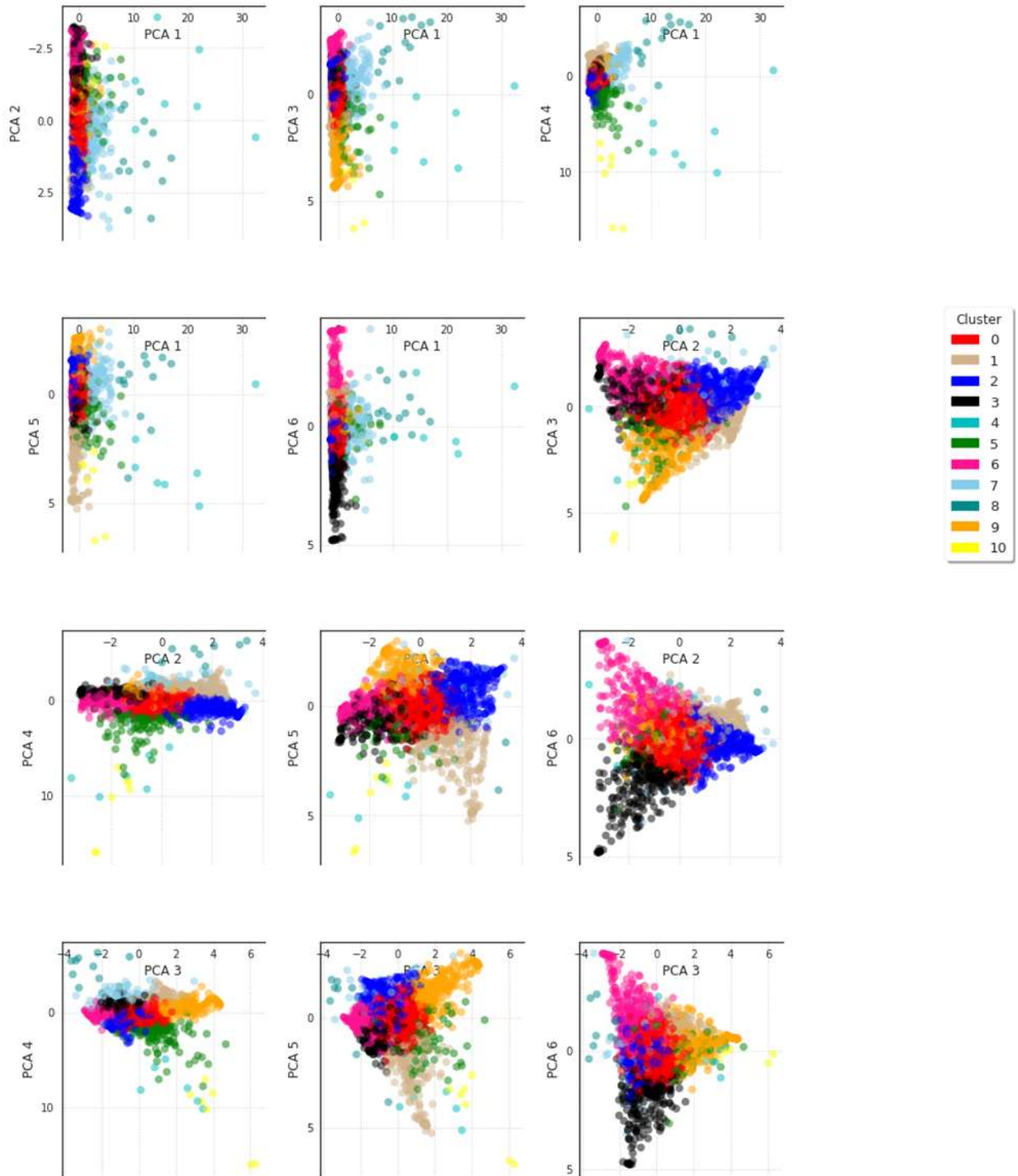
[3] Hong, T., Kim, E. (2011). It separates consumers from online stores based on factors that affect the customer's intention to purchase. Expert System Applications, 39 (2), 2127-2131.

[4] Hwang, Y. H. (2019). Hands-on Advertising Science Data: Develop your machine learning marketing strategies… using python and r. S.l: Packt printing is limited

[5] Puwanenthiren Premkanth, - Market Classification and Its Impact on Customer Satisfaction and Special Reference to the Commercial Bank of Ceylon PLC.‖ Global Journal of Management and Business Publisher Research: Global Magazenals Inc. (USA). 2012. Print ISSN: 0975-5853. Volume 12 Issue 1.

[6] Puwanenthiren Premkanth, - Market Classification and Its Impact on Customer Satisfaction and Special Reference to the Commercial Bank of Ceylon PLC.‖ Global Journal of Management and Business Publisher Research: Global Magazenals Inc. (USA). 2012. Print ISSN: 0975-5853. Volume 12 Issue 1.

[7] Sulekha Goyat. "The basis of market segmentation: a critical review of the literature. European Journal of Business and Management www.iiste.org. 2011. ISSN 2222-1905 (Paper) ISSN 2222-2839 (Online). Vol 3, No.9, 2011

[8] By Jerry W Thomas. 2007. Accessed at: www.decisionanalyst.com on July 12, 2015.

[9] T.Nelson Gnanaraj, Dr.K.Ramesh Kumar N.Monica. AnuManufactured cluster analysis using a new algorithm from structured and unstructured data. International Journal of Advances in Computer Science and Technology. 2007. Volume 3, No.2.

[10] McKinsey Global Institute. Big data. The next frontier is creativity, competition and productivity. 2011. Accessed at: www.mckinsey.com/mgi on July 14, 2015.

[11] Jean Yan. - Big Data, Big Opportunities- Domains of Data.gov: Promote, lead, contribute, and collaborate in the big data era. 2013. Retrieved from http://www.meritalk.com/pdfs/bdx/bdxwhitepaper-090413.pdf July 14, 2015.

[12] A.K. Jain, M.N. Murty and P.J. Flynn.‖ Data Integration: A Review‖. ACM Computer Research. 1999. Vol. 31, No. 3.

[13] Vishish R. Patel1 and Rupa G. Mehta. MpImpact for External Removal and Standard Procedures for JCSI International International Science Issues Issues, Vol. 8, Appeals 5, No 2, September 2011 ISSN (Online): 1694-0814  7

[14] Jayant Tikmani, Sudhanshu Tiwari, Sujata Khedkar "Telecom Customer Classification Based on Group Analysis of K-methods", JIRCCE, Year: 2015.


[15] Vaishali R. Patel and Rupa G. Mehta "Impact of Outlier Removal and Normalization Approach in Modified k-Mean