

# Customer Segmentation using Tableau Method

Ms.S.Sivaselvi.AP/CSE, M.Abi , G.Iswarya, M.Karthika  
Department of Computer Science and Engineering, Erode  
Sengunthar Engineering College,  
Erode, Tamilnadu

## Abstract

There is lot of competitions among businesses to attract new customers and hold on to existing ones as a result of the creation of numerous competitors and entrepreneurs. Because of the aforementioned, outstanding customer service is required, regardless of the size of the company. Any company will also benefit more from targeted customer services and the creation of unique customer care strategies if it can comprehend the wants of both of its clients. Structured customer service makes it feasible to have this understanding. Customers from each segment are similar in conditions of market characteristics. In contrast to traditional market analytics, which frequently fail, especially when the customer base is too broad, big data concepts and machine learning have increased the acceptance of the automated customer segmentation strategy. The Tableau clustering approach is employed in this paper to achieve this. The software is trained using a two-factor dataset of 100 patterns gathered from the retail industry and was created using the sklearn package for the Tableau approach (available in the appendix). Features of the average monthly customer visits and average monthly customer purchases.

Keywords: Customer Segmentation, k-Means algorithm, Mean shift algorithm, Agglomerative algorithm, Machine learning, Python.

---

## INTRODUCTION:

Over the years, the growing struggle between businesses and the availability of large-scale historical data has resulted in the extensive use of information pulling out techniques to discover important and strategic information that is hidden in the information of organizations. Data mining is the development of extract logical information from a dataset and presenting it in a human-accessible way for decision support. Data mining techniques distinguish areas such as statistics, artificial intelligence, machine learning and information systems. Information pulling out application contain but are not limited to bioinformatics, weather prediction, fraud detection, financial analysis and customer segmentation. Customer segmentation is the progression of grouping a company's customer base into subsets known as customer segments, each of which comprises of clients with comparable market criteria. These distinctions are created on variables such as product preferences or expectations, regions, behaviour, and other variables that might comprise a direct or indirect impact on the market or business. The ability of a company to tailor marketing strategies that will be suitable designed for each segment of its customers; support for business decisions base under a risky environment, such as debt relations with their customers; identification of products related to individual components and how to manage demand and supply power; and reveals the interdependence and interaction between consumers, between products, or between customers be immediately a little examples of the position of customer segmentation. Integrated was successful in locating hidden correlations or designs in a database of unencrypted data. [10][11]. This method of education falls under the category of supervised learning. The Tableau approach, Sorting Map (SOM), and other integration techniques are available. By comparing input patterns repeatedly until stable qualifiers in the exercise instances are formed, this method can discover clusters in data without having any prior knowledge of the information contingent on the topic or the procedure. The customer segment was subjected to the Tableau clustering algorithm. The Tableau method's sklearn library (Appendix) was created, and training began using a typical Silhouette -score and two attribute set of 100 preparation patterns identified in the retail industry.

### **PROPOSED SYSTEM:**

I employed the impression of transfer learning for picture categorization in this research. The key benefit of applying transfer learning is that the representation starts the learning process using patterns that have been learnt when tackling a different problem that is similar in nature to the one being solved rather than starting commencing score. In this manner, the model makes use of prior knowledge rather than initial from scratch. Transmission learning is typically expressed in image categorization through the usage of trained models. A prototypical that has already been proficient on a sizable standard dataset to address a problem that is related to the one we are trying to solve is referred to as pre-trained. A system's

structure, behaviour, and other aspects are all defined by its system architecture, a conceptual model. A proper description and illustration of a system that facilitates inferences about its behaviour and structure is called an architecture description. A system representation that incorporates the functional mapping of hardware and software mechanisms, the software architecture plotting to the hardware architecture, and the interaction of humans with these mechanisms.

### **SYSTEM DESIGN AND DEVELOPMENT:**

The process of creating a new system to either supplement or completely replace the current system is known as system design. The design phase's goal is to initiate the transition from the problem domain to the solution domain.

The mainly imperative factor influencing the software's quality is the system's design. Top level design is another name for software design. The system's logical components are transformed into its physical components during the design process. The software engineer can define an analytical model using the modelling notation that is governed by a position of syntactic-semantic and pragmatic principles using the Unified Modelling Language (UML). Systems that require a lot of software are frequently designed and visualised using the Unified Modelling Language. The involvement of software has improved knowingly in current years, making it more difficult for developers to create sophisticated programmes quickly. These software programmes frequently have defects, even when they do, and it preserve take programmers weeks to uncover and fix them. This is time that was wasted since there was a method that could have been applied to cut down on issues before the programme was finished. Since UML is not a methodology, no official work products are necessary. However, it does offer a variability of diagrams that, when used with a certain technique, make it simpler to understand an application that is still being developed. Although there is more to UML than these diagram, for my needs they serve as a solid introduction to both the speech and the guiding principles for using it. By including typical UML diagrams in the work products of your approach, you make it simpler for UML-savvy individuals to join your project and start contributing right away.

Five views that each explain the UML system from a distinctively different angle are cast off to portray the scheme. A sequence of diagrams that characterize each view be because follow.

### **FILE DESIGN:**

File design is a difficult task that involves creating both input and output forms. The goal of form design is to ensure that data is gathered, analysed, and stored quickly and accurately. A code is an organised group of symbols used to uniquely identify an entry or feature. Sometimes, all of the item's physical or performance qualities or operating in instructions are stated in place of the item's name. Additionally, they can foster relationships and occasionally exist worn to preserve confidentiality or secrecy.

### **INPUT DESIGN:**

Errors in data processing are most frequently caused by inaccurate input data. Data entry operators' errors can be managed through input design. The practice of translating user-created data to a computer-based format is known as input design. The expanded data flow diagram identifies logical data stores, sources, and destinations throughout the system design process. The master file (database), transaction files, and computer applications are describe in a system flow chart. Once the proper input media have been determined, the involvement data are gathered and grouped together into groups of similar data for processing. This project makes use of text, list, and combo boxes to collect keyboard input from users. The text box and list box be worn to attain or provide data for the project, respectively. The project's input design is created using a variability of online forms and techniques. For instance, the Students username and password are not permitted in the Admin form.

If the username already exists in the database, the input is rejected as invalid. The student's name, Student ID, age, address, salary, and designation will all be recorded during the "Student Creation" procedure.

### OUTPUT DESIGN:

The most significant and immediate source of information for the user is computer output. The systems should be improved by effective, clear output design. It clarifies interactions with users and aids in making decisions. A hardcopy printed by the printer is a significant output format. Printouts should be created with the user's output needs in mind. The output expedient and system, the necessary numeral of copy, the expected print quality, and the response time requirements

This method employs a list box and combo box to display output. For output, we can typically use the Visual DisplayUnit (VDU). The bills are often employed as a grid structure for databases. The results are often shown as reports on the screens.

### SYSTEM DEVELOPMENT:

The procedure of defining a system's structural design, workings, module, interface, and information in order to meet predetermined requirements is known as system development. One could think of arrangement plan as the submission of system theory to the creation of products. The fields of system analysis, system architecture, and systems engineering have certain areas of overlap. The system's realinput and output operations are associated to the physical design. This is outlined in conditions of how information is entered into a system, validated or authorised, processed, and shown as output. The following system requirements are chosen during physical design.

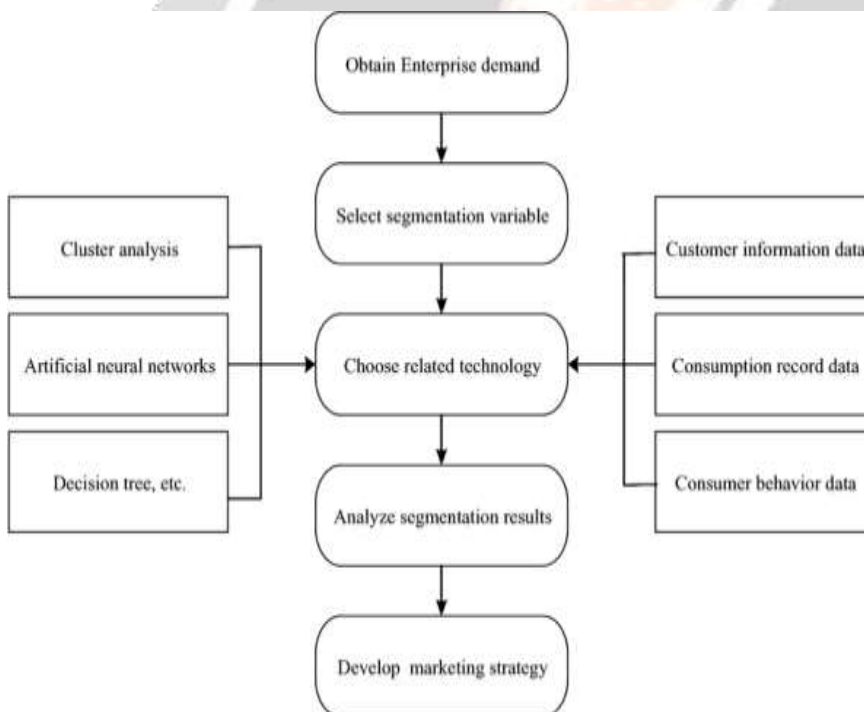


Fig.1.Process of system development

**Data Pre-Processing:** The data set wants to exist prepared via data pre-processing consequently to facilitate it may be breast-fed to the chosen model. The dataset must first be split into the three categories of Train Dataset, Validation Dataset, and Test Dataset in command to attain more accurate results. The Train dataset will receive the largest percentage of the photos because it is used to train the chosen model, even if each of these category has its own significance. The validation set is use to test the trained model's correctness, and we may find out the accuracy as we train for each epoch. As a outcome, the validation set is crucial. The built-in model's accuracy indicates how powerful and effective it is. The ideal range is between 85 and 98%.[7] The Test usual is innovative to verify the finalmodel that was created. The image set from the Test usual is breast-fed into the reproduction to

see if the output is accurate or not, indicating how effectively the typical can function with real-time images. To prevent the model from becoming overfitted, it is best practise to use distinct photos designed for every of these sets.

**a) Data Argumentation:** Data argumentation is a procedure for increasing the quantity of information by adding copies of existing photos with minor modifications. When the typical is being trained, this feature helps to minimise overfitting. This procedure is used in the current model since the dataset that was obtained from the organisation isn't big enough, which prevents the model from properly extracting the characteristics of the feed and, as a end result, from effectively training the model, which could lead to an incorrect outcome in the end. Data transformations including zooming, cropping, rotation, noise injection, and horizontal and vertical flipping are the applied to the existing dataset as measurement of data argumentation. The trained data set is used in the current system to apply horizontal flipping, image rescaling, zooming in on photos, etc.

**b) Model Building:** The most crucial stage of the entire system is thought to be this Phase. The idea of transfer learning is being worn in this situation. In added words, we use the predefined models that are already in use and modify them by replacing the top-most layer with the output classes we need.

**c)** The model weights used to connect the nodes are left unchanged. Only the output layer, which is the topmost layer, is fitted.

The precision of the system is significantly influenced by the typical choice. Every model has a unique set of parameters, including image size. Resnet, the Inception series, MobileNet, VGG, DenseNet, and others be a little example of these Keras Application models. Based on their weights, these models' build times, accuracy, and classification times vary. within arrange to build the system, I chose the Inception V3 Model from among these models based on its accuracy. It supports 3 channels with an image size of 224x224. The truthfulness of the form, which was trained over 20 epochs, is about 96.1% Training and Testing: The chosen model is assembled, trained, and authorized using the effort data after a predetermined number of epochs. The trained model's accuracy can then be determined. Additionally, we can build using many models so so as to the accuracy differences are known and the optimum algorithm for the system is selected.

#### **DESCRIPTION OF MODULES:**

In attendance is rejection need for a separate module in this project. The software recognises consumer segmentation using machine learning techniques from Tableau.

#### **ARCHITECTURE:**

A system's structure, behaviour, and other aspects are all defined by its system architecture, a theoretical model. A proper description and illustration of a scheme that facilitates inferences about its behaviour and structure is called an architecture description. A system representation that incorporates the functional mapping of hardware and software components, the software architecture mapping to the hardware architecture, and the interaction of humans with these components. In order to provide an objective analysis, the system gives each consumer a numerical score based on these variables. The marketing proverb "80% of your business originates from 20% of your consumers" serves as the foundation for RFM analysis.



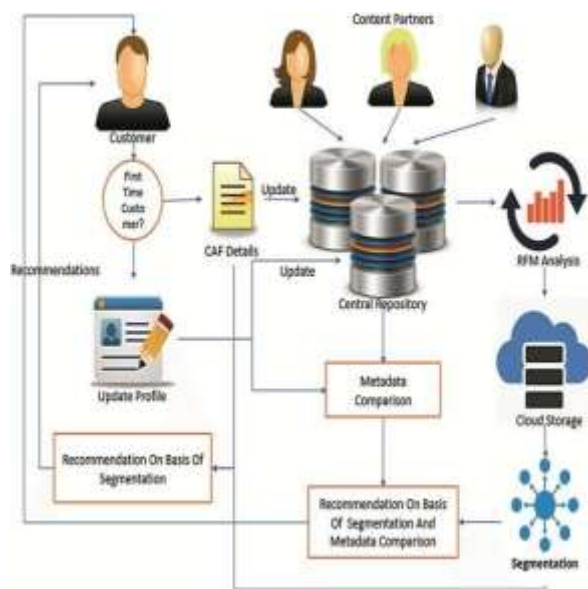
**ARCHITECTURE DESIGN:**

Fig.2. Architecture diagram

**TESTING AND IMPLEMENTATION: SYSTEM TESTING:**

Testing is decisive for error correction. Otherwise, the project or programmer is deemed incomplete. Software testing is a process that determines whether the created system is operating in accordance with the initial goals and specifications. In the direction of create confident the system meets the needed specifications, it should be experimentally tested using test data. Software testing serves as the final assessment of the specification, design, and coding and is a crucial component of software quality assurance. Following the coding stage, it is possible to run computer programmers for testing. This suggests that testing must find problems that were introduced not only during coding but also in the earlier phases.

Testing objectives: objective of the test

The next three steps are a summary of the testing goals.

- Testing is the process of running a software in search of errors.
- A trial case with a high likelihood of detecting an error that hasn't been found yet is a good trial case.
- A effective trial is one that identifies an error that hasn't been found yet.
- Testing guidelines
- Each test should be able to be linked back to a client requirement.
- Tests should be scheduled well in advance of the start of testing, ideally as soon as the requirement model is finished.
- There is a decent chance that a test will identify a mistake.
- An effective test is non-redundant.
- A decent test would not be overly straightforward nor overly difficult.

**Testing methodologies:** Before the scheme is put into actual operation, accuracy and efficiency are checked during the system testing phase of implementation. The system's success depends on testing. Scheme testing makes the logical premise that the objective will be effectively attained if every component of the scheme is correct.

**The testing steps are:**

- Unit Testing
- Integration Testing
- Validation Testing
- Output Testing
- User acceptance Testing

**Unit testing:** Unit testing concentrates verification efforts on the modules, which are the minimum component of software design. It is also referred to as "Module Testing". Each module is tested independently. This testing is complete right together with the code. To achieve complete coverage and maximum error detection, unit testing specifies tracks in the module's control structure. This test focuses on each module separately to brand definite that they all work together properly.

**Integration-testing:** Data loss across the interface is a possibility, and one module may negatively impact others. Systematic testing for building programmer structure is called integration testing. simultaneously running tests to find interface-related problems The concerns associated to the concurrent issues of confirmation and programmer construction are addressed by integration testing. A sequence of high order sets are undertaken after the software has been fully integrated. The objective is to merge components that have been unit tested and test the organization as a whole. Therefore, all faults found during the integration testing stage are fixed in preparation for the following testing phases.

**Validation Testing:** The inputs that enter the system produce the outputs that come out of the system. So, in demand for the system to produce the desired and expected outputs, the inputs must be appropriate and correct. Therefore, this testing is carried out to confirm that the inputs are genuine and correct before they are entered into the system for processing.

**Output Testing:** The proposed system's output testing comes after the validation testing because no system can be helpful if it cannot deliver the necessary output in the needed format. The outputs produced or displayed by the system under consideration are tested by asking the users about the format they desire. As a effect, there are two ways to think about the output format: on-screen and printed.

**User acceptance Testing:** A system test examines the integration of each system module slightly than the software itself. Additionally, it checks for inconsistencies between the system's current specs, initial aim, and system documentation.

**SYSTEM IMPLEMENTATION:**

The project stage known as system implementation is where the conceptual design is transformed into a functional system. The implementation stage can result in error if it is not carefully planned and managed. As a result, it canister be said to be the most significant step in ensuring the success of a new system and in giving users confidence that the scheme will operate as intended. Typically, at this stage, a coordinating committee is formed to serve as a forum for ideas, grievances, and problems. Implementation planning, or choosing the procedures and timeline to be used, is the initial task. The emphasis during the implementation stage must be on training in new skillsto give workers confidence they can utilise the system. Education takes place far sooner in the project than planning, which is the other primary duty of preparing for implementation. The system can be verified after staff training. Evaluation and maintenance are needed to bring the new system up to standards once the implementation phase is over and the user staff has had time to adjust to the changes brought about by the entrant system. The implementation phase's actions can be summed up as follows:

- Implementation planning
- Education planning
- System planning

**Datasets:** A dataset is a grouping of different kind of evidence to has been digitally stored. Any project consuming machine learning desires data as its primary input. Datasets typically include text, audio, video, photos, and quantitative data points. This data set was designed specifically to teach the market basket analysis and customer segmentation principles. You own a grocery mall, and you have some basic in order regarding your clients, such

as Customer ID, age, gender, annual income, and spending score, thanks to membership cards. I'll show you how this works by utilising the most basic version of the unsupervised ML approach (KMeans Clustering Algorithm).

Fig.3. Comparison of Customer Datasets

The following dataset's data quality was evaluated before moving on near the information cleaning process. Following a data quality evaluation, the following concerns with data quality were identified, and the appropriate steps to address them were taken:

**Customer Demo graphics.xlsx:** There was one irrelevant column, and those columns were removed from the dataset.

There were missing values in 5 of the columns. According to the number of missing values for such columns, either records were discarded or appropriate values were imputed at the locations where missing values were present.

There was no standardisation of evidence for the gender column. The column data was standardised to eliminate data inconsistencies built on the morals that were provided. In directive to look for differences in the age distribution, the Date of Birth column was changed to generate the new feature columns "Age" and "Age Group." The record was altered after an outlier was discovered.

checked toward observe if the dataset contains any duplicate records. There were no duplicate records in this collection.

**New Customer List.xlsx:** The dataset contained an irrelevant column, which was removed. There were missing values in 4 of the columns. Contingent on the volume of missing values for such columns, either records were dropped or appropriate values were imputed at the locations where missing values were present. In order to look for differences in the age distribution, the Date of Birth column was changed to generate the new feature columns "Age" and "Age Group." There were no inconsistent data points checked to see if the dataset contains any duplicate records. There were no duplicate records in this collection.

**New vs Old Customers Age Distribution:** Additionally, the popular of both new and existing consumers are between the ages of 40 and 49. Clients who are under 20 and over 80 years old variety up the smallest percentage of both sorts of customers. The data preparation, quality evaluation, and data cleaning stage was the initial step in producing valuable insights from the data. Exploratory data analysis on the dataset and the identification of client purchase habits can be done after the cleaning phase to produce insights. The automaker enjoys a in height leveled of the popularity among original customers in the age data

brackets of 20–29 and 40–49. Customers in the 30-39 age bracket who are new customers show a sharp decline. Old Clientele Distribution by Age

#### Age Distribution of New Clients:

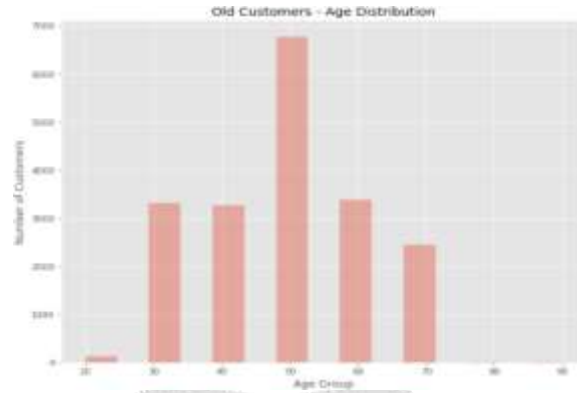


Fig.4.Age distribution of old customer

**Old Customers Wealth by Age Group:**

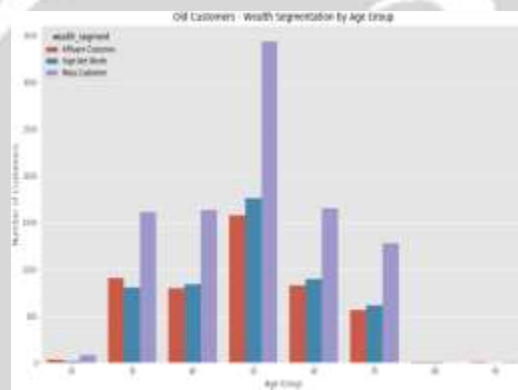


Fig.5.Wealth distribution of old customer

**New Customers Wealth by Age Group:**

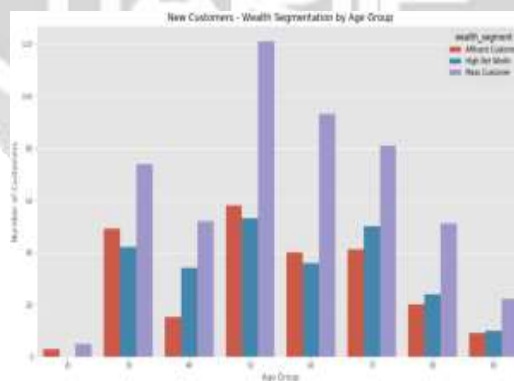


Fig.6.Wealth distribution of new customer

**Future Work:** The project's goals are achieved by the software during execution. With minor adjustments, additional extensions to this system can be made necessary. The inventions container be worn popular computer hardware, firmware, software, digital electronic circuitry, or in combinations of these. By processing raw data and producing output, inventions carry out their intended functions.



**Literature Review:**

Customer Segmentation:

Jayant et al. [1] states that in customer segmentation, the customers are divided into different groups where customers of the same group are similar to each other in conditions of marketing. Customers are divided into different clusters based on various attributes such as age, interests, age, spending habits etc. Sulekha et al.[7] provides the four popular bases for segmentation

- 1) **Geographic Segmentation:** segmentation on the basis geographic region, population density or climate.
- 2) **Demographic Segmentation:** market segment on the basis of age, size and family type, etc.
- 3) **Psychographic Segmentation:** segmentation based on customer's life style variables like interests, opinions, attitudes etc.
- 4) **Behavioural Segmentation:** This caring of segmentation bases its conclusions on real consumer behaviour toward a product, such as brand loyalty, user status, purchase readiness, etc.

**Clustering:** The technique of arranging objects into based on some similarities across the dataset's metadata. Numerous algorithms exist. which, depending on the circumstances, can be chosen to be used on a dataset. However, there is no one clustering algorithm, therefore choosing the right strategies for clustering becomes crucial. Vaishali and others [8]. Using the Python Scikit- Learn module, we constructed three clustering methods in this paper [9].

**K-Means Clustering:** single of the the majority well- liked partitioning clustering algorithms is the K-Means algorithm. According to Chinedu et al., the clustering technique is dependent on the centroid, which is the location of each data point on one of the K non- overlapping clusters that are chosen prior toward the program's execution. [2]. The clusters that were created lined awake among the data's underlying pattern, providing the necessary details to aid indecision-making.

**Clustering by agglomeration:**

This clustering falls under the grouping of hierarchical clusters, which are created via a hierarchy. Using hierarchical clustering It is base on the idea that objects that are closer to one another are more linked to one another than persons that are farther apart. [3]. This challenge is lucrative because it results in lower calculation costs by letting go of a combinatorial variety of options. Yogita and others [4]. In hierarchical clustering, present are two different approaches: the top-down approach, besides known as divisive clustering, and the bottom-up approach, also known as agglomerative clustering. According to Omar et al. [5,] agglomerative clustering is often slower than divisive clustering but offers greater flexibility because it enables the consumer to input any arbitrary similarity function that specifies what counts as a similar cluster pair to merge together.

**CONCLUSION:**

This study demonstrates that dividing customers based on behavioural characteristics is a better solution for the current customer segmentation problem, and Tableau method is identified as a excellent choice for this approach. Customer segmentation is perform on the company's customers data and with the help of Tableau method, customers are divided using features like total spending and annual income.

**References:**

- [1] Jayant Tikmani, Sudhanshu Tiwari, Sujata Khedkar "Telecom customer segmentation based on cluster analysis An Approach to Customer Classification using k- means", IJIRCCE, Year: 2015.
- [2] Chinedu Pascal Ezenkwu, Simeon Ozuomba, Constance kalu Electrical/Electronics & Computer Engineering Department, University of Uyo, Uyo, Akwa Ibom State, Nigeria "Application of K-Means Algorithm for Efficient Customer Segmentation: A Strategy for Targeted Customer Services", IJARAI, Year: 2015.
- [3] T.Nelson Gnanaraj Dr.K.Ramesh Kumar onica "Survey on mining clusters using new k-mean algorithm from structured and unstructured data", IJACST, Year: 2014.

- [4] Yogita Rani and Dr. Harish Rohil “A Study of Hierarchical Clustering Algorithm”, IJICT, Year: 2013.
- [5] Omar Kettani, Faycal Ramdani, Benaissa Tadili “An Agglomerative Clustering Method for Large Data Sets”, IJCA, Year: 2014.
- [6] Snekha, Chetna Sachdeva, Rajesh Birok “Real Time Object Tracking Using Different Mean Shift Techniques–aReview”, IJSCE, Year: 2013.
- [7] Sulekha Goyat “The basis of market segmentation: acritical review of literature”, EJBM, Year: 2011.
- [8] Vaishali R. Patel and Rupa G. Mehta “Impact of Outlier Removal and Normalization Approach in Modified k- MeansClustering Algorithm”, IJCSI, Year: 2011.
- [9] Scikit-learn: <https://scikit-learn.org>
- [10] Tanupriya Choudhury, Vivek Kumar, Darshika Nigam, Intelligent Classification & Clustering Of Lung & Oral Cancer through Decision Tree & Genetic Algorithm, International Journal of Advanced Research in Computer Science and Software Engineering,2015.
- [11] Tanupriya Choudhury, Vivek Kumar, Darshika Nigam, An Innovative and Automatic Lung and Oral Cancer Classification Using Soft Computing Techniques, International Journal of Computer Science & Mobile Computing,2015 .

