# Customer Segmentation using Machine Learning in R

## Gurram Venkata Sai Lakshmi Tejaswani[1], Gurram Satya Sravani[2], Kesavarthini M[3], Deepika L P[4]

*[1,2,3,4]Student(B.E), CSBS Department, R.M.D Engineering College, Kavaraipettai, Chennai-601 206, Tamil Nadu, India.*

---------------------------------------------------------------------***---------------------------------------------------------------------

**ABSTRACT -** Customers Segmentation has been a subject of revenue for a ton of industry, scholastics, and showcasing pioneers. The expected worth of a client to an organization can be a center fixing in dynamic. One of the huge difficulties in client based associations is client cognizance, understanding the distinction among them, and scoring them. However, presently with all capacities we have, utilizing new innovations like AI calculation and information treatment we can make an exceptionally amazing structure that permit us to best comprehend clients needs and practices, and act suitably to fulfill their requirements. In the current paper, we propose another model dependent on RFM model Recency, Frequency, and Monetary and k-mean calculation to determine those difficulties. This model will permit us to utilize grouping, scoring, and conveyance to have an unmistakable thought regarding what move we should make to improve consumer loyalty.

**Keywords**: Data mining; machine learning; customer segment; k-Mean algorithm; sklearn; extrapolation

## INTRODUCTION

Customer Segmentation is the cycle of division of client base into a few gatherings of people that share a comparability in various manners that are pertinent to showcasing like sexual orientation, age, interests, and different ways of managing money.

Organizations that convey client division are under the thought that each client has various necessities and require a particular advertising exertion to address them suitably. Organizations mean to acquire a more profound methodology of the client they are focusing on. In this manner, their point must be explicit and ought to be customized to address the necessities of every single individual client. Besides, through the information gathered, organizations can acquire a more profound comprehension of client inclinations just as the prerequisites for finding important portions that would harvest them greatest benefit. Along these lines, they can plan their showcasing strategies all the more productively and limit the chance of hazard to their venture.

The procedure of client division is subject to a few key differentiators that partition clients into gatherings to be focused on. Information identified with socioeconomics, geology, financial status just as personal conduct standards assume a significant part in deciding the organization heading towards tending to the different sections.
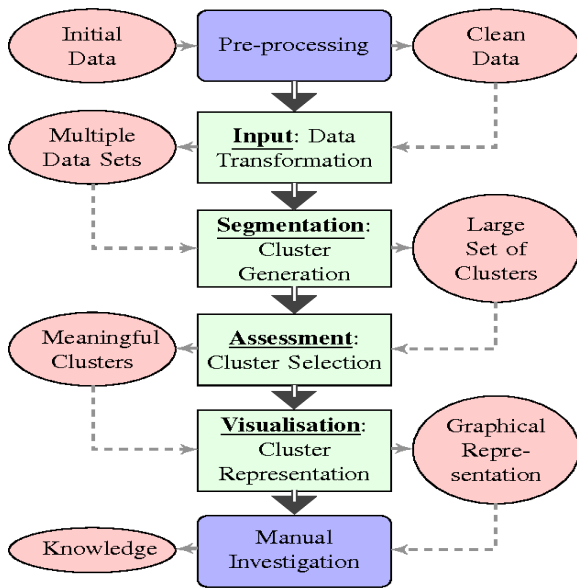
Covered in an information base of incorporated information end up being powerful for identifying unpretentious yet unobtrusive examples or connections. This method of learning is ordered under regulated learning. Reconciliation calculations incorporate the KMeans calculation, K-closest calculation, arranging map (SOM), and more.[4] These calculations, without earlier information on the information, can distinguish bunches in them by more than once contrasting info designs, as long as static inclination in preparing models is accomplished dependent on topic or interaction. Each set has information focuses that have close similitudes however vary incredibly from the information points of different gatherings. Coordination has incredible applications in design acknowledgment, picture examination, and bioinformatics thus on.[15] In this paper the k-implies bunching calculation was executed in the client portion. The scalar library (Appendix) of the K-Means calculation was created, and preparing was begun utilizing a standard outline - score with two capabilities of 100 preparing designs found in the retail exchange. After a few signs, four stable stretches or client sections were recognized.

## Proposed Model:

To start, we import the fundamental bundles to do our examination and afterward the xlsx (Excel bookkeeping page) information record. In the event that you need to circle back to a similar information, you need to download it from UCI. For this model, I place the xlsx document in the envelope (catalog) where I present Jupiter's scratch pad. Subsequent to bringing in the bundle and information, we will see that the information isn't actually that supportive, so we need to clean and sort out this information such that we can make more noteworthy bits of knowledge.

The K-implies region unit is touchy to the size of the data utilized, like grouping calculations, so we might want to sum up the information.A screen capture of the StackExchange answer beneath talks about why normalization or standardization is important for information utilized in K-implies bunching. The screen capture is connected to the StackExchange question, so you can tap on it and read the sum of the conversation in the event that you need more data.

## Architecture Diagram



**DATASET:** informational collection is known as the methodical course of action of the patient details.in any informational index the individual subtleties are covered up like name,address,phone number etc.a informational collection is the fundamental thing for any AI project.generally, informational indexes are available as CSV(Comma Separated Values).

### PREPROCESSING DATA:

This is a significant stage prior to going to the following phase. In this stage we eliminate all the invalid values, invalid objects, symbols, strings, etc.

### CLASSIFICATION TECHINQUES:

This method is to distinguish the classification or class to which another information will fall under.for model if u got a mail it will consequently order on the off chance that it is spam or not.

### ALGORITHM:
This is the primary stage or the core of the system. After the preprocessing and characterization of information we will pass clean informational collection into the calculation which is otherwise called preparing of the model. In this undertaking we utilized various calculations like K closest neighbor algorithm, support vector algorithm, linear regression, decision trees, random woodland and so on

### DATA:

The data set contains of columns and 200 rows. Data set contains integer and float data type values after the processing of data

| CustomerID | Gender | Age | Annual.Income..k.. | Spending.Score..1.100. |
|---|---|---|---|---|
| 1 | Male | 19 | 15 | 39 |
| 2 | Male | 21 | 15 | 81 |
| 3 | Female | 20 | 16 | 6 |
| 4 | Female | 23 | 16 | 77 |
| 5 | Female | 31 | 17 | 40 |
| 6 | Female | 22 | 17 | 76 |

### ATTRIBUTE DETAILS:

| Attribute | Description | Type |
|---|---|---|
| CustomerID | Customer has a unique id by default | Numeric |
| Gender | Customers gender(male represented as Male and female as Female) | String |
| Age | Customer's age in Completed year | Numeric |
| Income | Total income achiving the customer in an year | Numeric |
| Spending Score | It derives the Spending score of the annual amount | Numeric |

### TRAINING AND TESTING METHOS:

The below graphs will show the accuracy percentage while we using different machine learning algorithms as follows

### Pie Chart:

A pie outline (or a circle diagram) is a round measurable realistic, which is isolated into cuts to show mathematical extent The underneath pie graph show the Depicting Ratio of Female and male.

**FIGURE - 1**

### Histogram:

A **histogram** is a bar graph-like representation of data that buckets a range of outcomes into columns along the x-axis. The y-axis represents the number count or percentage of occurrences in the data for each column and can be used to visualize data distributions. Figure 2 represents about the count of an age class and Figure 3 represents about the Annual income of an overall customers.



FIGURE - 2



FIGURE – 3

Box Plot:

In descriptive statistics, a box plot or boxplot is a method for graphically depicting groups of numerical data through their quartiles. Figure 4 represents about the Desciptive Analysis of Age.



**FIGURE – 4**

### OUTPUT:

In the output of our kmeans operation, we observe a list with several key information. From this, we conclude the useful information being –

- o **cluster –** This is a vector of several integers that denote the cluster which has an allocation of each point.
- o **totss –** This represents the total sum of squares.
- o **centers –** Matrix comprising of several cluster centers
- o **withinss –** This is a vector representing the intra-cluster sum of squares having one component per cluster.
- o **tot.withinss –** This denotes the total intra-cluster sum of squares.
- o **betweens –** This is the sum of between-cluster squares.
- o **size –** The total number of points that each cluster holds.

```
# compute gap statistic
k6<-kmeans(customer_data[,3:5],6,iter.max=100,nstart=50,algorithm="Lloyd")
k6

## K-means clustering with 6 clusters of sizes 45, 22, 21, 38, 35, 39
##
## Cluster means:
##        Age Annual.Income..k.. Spending.Score..1.100.
## 1 56.15556           53.37778               49.08889
## 2 25.27273           25.72727               79.36364
## 3 44.14286           25.14286               19.52381
## 4 27.00000           56.65789               49.13158
## 5 41.68571           88.22857               17.28571
## 6 32.69231           86.53846               82.12821
##
## Clustering vector:
##  [1] 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3
## [36] 2 3 2 3 2 1 2 1 4 3 2 1 4 4 4 1 4 4 1 1 1 1 1 4 1 1 4 1 1 1 4 1 1 4 4
## [71] 1 1 1 1 1 4 1 4 4 1 1 4 1 1 4 1 1 4 4 1 1 4 1 4 4 4 1 4 1 4 4 1 1 4 1
```

**Conclusion**: In this information science project, we went through the client division model. We fostered this utilizing a class of AI known as solo learning. In particular, we utilized a bunching calculation called K-implies grouping. We examined and imagined the information and afterward continued to execute our calculation. Expectation you delighted in this client division venture of AI utilizing

**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