# CLIENT PARTITIONING IN ML TECHNIQUES USING K-MEANS CLUSTERING

## Calabe P S[1], Dr. Prabha[2]

[1]Student, Department of Computer Science and Engineering, Dr. Ambedkar Institute of Technology, Bangalore, India.

[2]Professor, Department of Computer Science and Engineering, Dr. Ambedkar Institute of Technology,

Bangalore, India.

**Abstract:** In sporting out a hit E-Commerce , the maximum critical matters are innovation and information what client wants. Now-a-days the benefit of the usage of ecommerce encourages the clients to shop for the usage of ecommerce. It runs on the idea of innovation having the capacity to enthral the clients with the merchandise, however with any such big raft of merchandise go away the clients pressured of what to shop for and what now no longer to. According to enterprise , a organization may also create 3 segments like High ( Group who buys often , spends greater and visited the platform lately ) , Medium ( Group which spends much less than excessive organization and isn't always that lots common to go to the platform) and Low (Group that's at the verge of churning out ). This is wherein Machine Learning presents a critical answer, numerous algorithms are implemented for revealing the hidden styles in statistics for higher selection making. In this paper we proposed a client segmentation idea wherein the consumer bases of an established order is split into segments primarily based totally at the clients' traits and attributes. This concept may be utilized by the B2C businesses to outperform the opposition through growing uniquely attractive services and products and make it attain to cappotential clients. This method is carried out the usage of "K-Means", an unmanaged clustering device mastering set of rules.

**Keywords:** Innovation, B2C, Machine Learning, E-Commerce, K-means Clustering, Client segmentation, Innovation, RFM Analysis, Loyalty Level**,** Cluster Creation, Business segments, Market Basket Analysis.

## I. INTRODUCTION

Client Segmentation is a way of grouping the clients primarily based totally on advertising corporations which stocks the similarity amongst clients. To be extra exact, it manner segmenting clients sharing the not unusual place traits that is the excellent manner of advertising. Client segmentation is accumulating records approximately every clients and analysing it to perceive the one of a kind styles for growing the segments. Some of the exceptional strategies for collecting records are face-to-face interviews, Telephonic interviews, thru surveys or via studies the usage of data which might be posted associated with marketplace categories. The primary records which incorporates billing statistics, delivery records, merchandise bought, promo codes, price approach etc., Beyond those a few groups additionally acquire data like cause for the buy, commercial channel which makes them to buy, age, gender etc., In B2B (Business to Business) advertising clients are grouped consistent with several elements like Industries, quantity of employers, Products bought from the organization in advance instances and location. On different-hand, in B2C (Business to Consumer) advertising businesses section the clients primarily based totally on Age, Gender of the clients, their marital status, lifestyles level of the clients like single, married, divorced, retired etc., On of the principle element of B2C is Location of clients (Rural, suburban, urban). Client segmentation may be practiced for all of the corporations though of length or industry. Common

segmentation kinds consist of Demographic, RFM (Recency, Frequency, Monetary) evaluation, HCS (High-fee client), patron status, Behavioral, psychographic etc., Some of the foremost advantages of client segmentation consist of Marketing strategy, promoting strategy, Budget efficiency, product improvement etc., In this text we implemented the fundamental analytics capability to offer the choice makers(in our case the enterprise investors) with the desired facts to make the proper choice. In this newsletter we outline an answer for decreasing threat elements and additionally make contributions to the choice making for brand spanking new enterprise investments. We proposed to apply K-method method for purchaser segmentation. Our answer is segmenting the clients primarily based totally on statistics analytics. Consumers may be divided into organizations in terms of not unusual place behaviors they share. Such behaviors

hyperlink to their information of, mindset toward, use of, or spending rating or reaction to a product. We used system getting to know Clustering set of rules K-Means for this client segmentation.

## II.    METHODOLOGY

There are 3 most important sections:
- Pre-processing of Data.
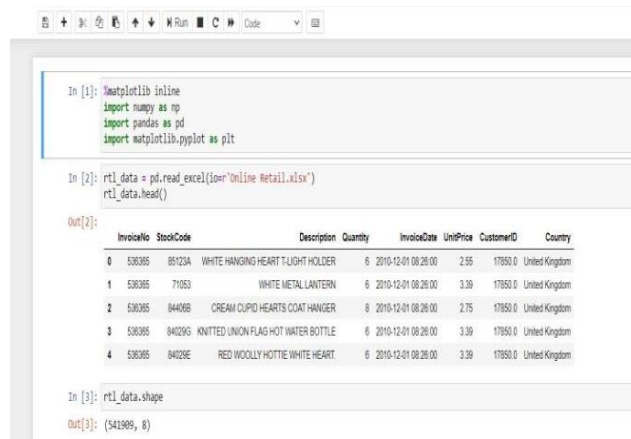- Creating a Cluster.
- Calculating RFM rating.

A python application has been evolved and this system is performed in Jupyter Notebook through uploading the subsequent important packages:
- **Pandas** – Loading & Data Pre-processing.
- **Numpy** – numeric calculation of data entered.
- **matplotlib, seaborn** – Visualization.
- **scikit** – Learning of M L library.

### 2.1 Pre-processing of Data:

The dataset that's an excel record is loaded with the use of Pandas.
Figure.1: Dataset uploaded using Pandas



We dropped all of the replica entries in Client ID and Country column. Count of clients for every Country is calculated and looked after the rely in descending order to peer from which us of a the most variety of clients purchases.



| Segment | % Conversion ▼ | % Improvement over Baseline ⇕ | % Chance to Outperform ⇕ |
|---|---|---|---|
| ● United States | 8.20% | 58.0% | 100% |
| ● Australia | 8.05% | 55.2% | 89.8% |
| ● Netherlands | 6.12% | 17.9% | 64.9% |
| ● Mexico | 5.81% | 12.0% | 59.7% |
| ● China | 5.56% | 7.03% | 55.3% |
| ● Baseline | 5.19% | - | - |
| ● Brazil | 3.90% | -24.9% | 28.1% |
| ● India | 2.29% | -55.9% | 0.00328% |

Figure.2: Count of Clients based on Country in Descending Order

From the primary few row's we found that most quantity of clients are from UK in Country. So we grouped the clients primarily based totally on nations. So we filtered out different nations the use of "query" technique. In the consolidated dataset, best the Description and Client ID columns had null values. Those entries are dropped off, also the bad entries in Quantity because it can't be bad. Invoice

Data that is in string layout and this is transformed to this point and time layout because it is vital even as calculating Recency. A new column Total Amount is introduced as a made of Quantity and Unit Price for every consumer. This is beneficial in Monetary calculation.

| Country/Region Code | Description | No Of Customer Entries | Percent % | METHOD USED |
|---|---|---|---|---|
| AT | Country/Region Code | 3 | 4.41 | Legacy |
| BE | Country/Region Code | 3 | 4.41 | Legacy |
| CA | Country/Region Code | 3 | 4.41 | Legacy |
| CH | Country/Region Code | 3 | 4.41 | Legacy |
| CZ | Country/Region Code | 3 | 4.41 | Legacy |
| DE | Country/Region Code | 4 | 5.88 | Legacy |
| DK | Country/Region Code | 5 | 7.35 | Legacy |
| ES | Country/Region Code | 3 | 4.41 | Legacy |
| FR | Country/Region Code | 3 | 4.41 | Legacy |
| GB | Country/Region Code | 3 | 4.41 | Legacy |
| IS | Country/Region Code | 3 | 4.41 | Legacy |
| MA | Country/Region Code | 3 | 4.41 | Legacy |
| MY | Country/Region Code | 3 | 4.41 | Legacy |
| NL | Country/Region Code | 11 | 16.18 | Legacy |
| NO | Country/Region Code | 3 | 4.41 | Legacy |
| SE | Country/Region Code | 3 | 4.41 | Legacy |
| SI | Country/Region Code | 3 | 4.41 | Legacy |
| US | Country/Region Code | 3 | 4.41 | Legacy |
| ZA | Country/Region Code | 3 | 4.41 | Legacy |

Figure.3: Customer Grouping primarily based on Country

## 2.2 Cluster Creation

K-Means is an unmanaged mastering set of rules and used for clustering responsibilities which goes truly properly with complicated dataset. It is an iterative set of rules that walls the dataset into "okay" pre-described non- overlapping subgroups (clusters) wherein every records factor belongs to most effective one organization.
The set of rules works as follows:

**Step-1:** Specifying the range of clusters – okay fee.

**Step-2:** Centroids are initialized via way of means of shuffling the dataset after which randomly choosing okay facts factors for the centroids with out replacement.

**Step-3:** Repeat the new release till there may be no extrade to the centroids. i.e, project of facts factors to the clusters does now no longer alternate.

Frequency, Monetary and Recency, are delivered to the identical scale and the records is normalized earlier than clustering process. It is crucial to decide the most advantageous wide variety of clusters i.e, "okay price". For this we used "Elbow technique".
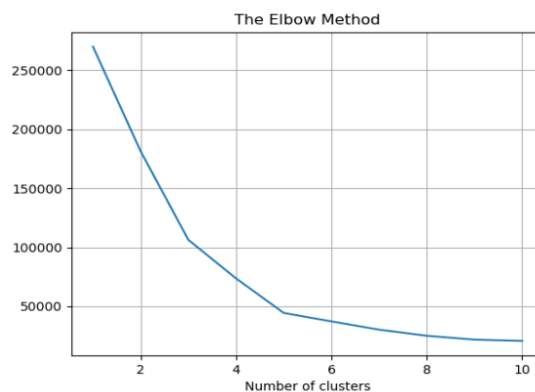


Chart-1: K value with the use of Elbow technique

It entails strolling the set of rules more than one instances over a loop with more and more cluster preference after which plotting a rating as a feature of the quantity of clusters. When "okay" increases, the centroids are toward cluster centroids. The development will decline sooner or later unexpectedly developing an elbow-like form in graph
and this is the entire motive this technique is referred to as elbow. We take the depend of cluster, okay-fee on the factor in which this elbow is bending. When we done this metric, the end result turned into now no longer apparent and the bend isn't clean as there has been fast decline at 3 values – three, five and seven. Silhouette technique that is taken into

consideration as a higher metric than elbow is used to decide the superior variety of clusters. Silhouette rating for every pattern is calculated the use of the formula:

Calculate the silhouette s(I) as follows, ratio of the distinction among cluster brotherly love and separation to the more of the two :

$$S(i) = b(i) - a(i) / [max\{a(i), b(i)\}]$$

Silhouette co-green degrees from [-1,1]. The calculated common rating is 0.3. Any cost better than this indicates it's far nicely matched to its personal cluster.



Figure.4: Silhouette evaluation for K-means clustering on pattern information with n clusters=3

From the outputs, we will infer that cluster three has a rating 0.35 and the fluctuation length is similar. To word that the thickness of the silhouette plot representing every cluster Additionally contributes to the choice. Therefore the most fulfilling fee of okay is three. Okay-method is imported from sklearn library and the arguments are set and the operation is implemented at the scaled records. In RFM we were given four businesses while in okay-manner the silver and bronze companies are merged as follows:



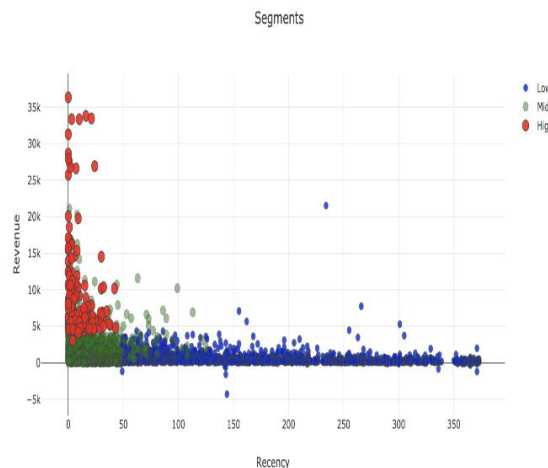Figure.5: Bronze and Silver merged to Silver organization.



Figure.6: Graph illustration of the three Clusters

### 2.3 Calculating RFM rating

Our dataset is restricted to income record, we are able to use a RFM primarily based totally version for locating segments wherein R is Recency (how currently a buy happened),
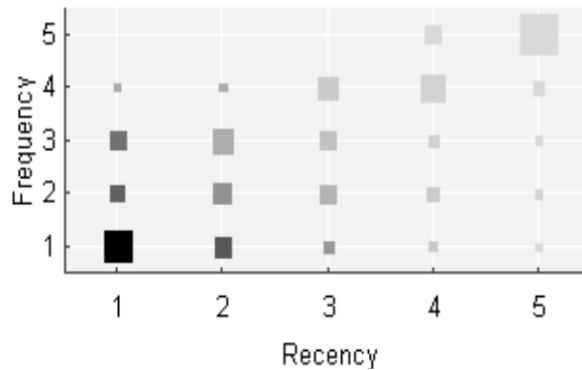


Chart-2: Recency plot

F is Frequency (how common transactions are made), M is Monetary fee (Value of all transactions). Recency, Frequency and Monetary rating for every patron is calculated. The contemporary date is assigned as a placeholder to calculate current purchases.
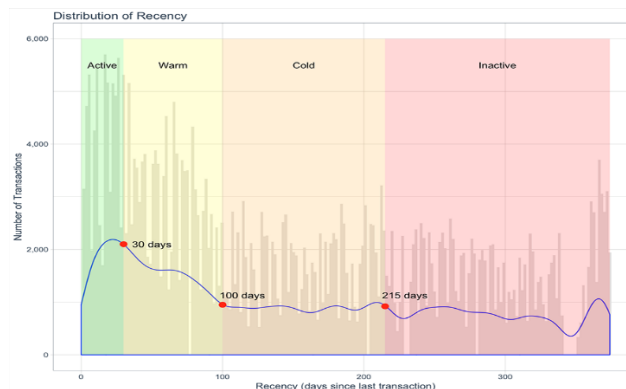


Figure.7: Graph illustration of RFM

All the transactions are grouped the usage of Client ID after which combination lambda operations are performed. As an end result of this operation numbers may be received which depicts the recency., frequency and what kind of a particular consumer spent until date. All those are saved in
a brand-new data frame RFM rate. To notice, the distribution for recency is proper skewed. Using quantiles information is split into 4 businesses (we will select quantile values of our personal). After RFM calculation the newly introduced columns are as follows:

| RFM_Score | Recency mean | Frequency mean | Monetary mean | count |
|---|---|---|---|---|
| 3.0 | 258.3 | 2.9 | -57.5 | 409 |
| 4.0 | 172.1 | 4.5 | 60.0 | 372 |
| 5.0 | 143.5 | 6.5 | 106.5 | 495 |
| 6.0 | 104.0 | 9.3 | 191.8 | 471 |
| 7.0 | 81.3 | 13.0 | 237.5 | 418 |

Figure.8: Calculating RFM

The person recency, frequency and financial values are concatenated and transformed to thread the usage of map characteristic. This is performed to effortlessly test which organization the client belongs to. This RFM score column suggests the loyalty of engagement of the client. In our case, the decrease the price of RFM score, extra dependable the purchaser may be in addition to greater engaged he/she might be. Based in this withinside the subsequent step Loyalty Level like Platinum, Gold, Silver and Bronze ranges are assigned to every client.

From this we should derive a end that if the consumer is in platinum institution we are able to say that they may be the nice clients while in bronze institution, the patron haven't bought for an extended time. With this a corporation can determine like supplying unique attention, gives and precedence get right of entry to to newly released merchandise to their platinum clients. On the opposite hand, if the purchaser falls into the bronze institution, the agency can deliver a few rewards or coupons to inspire the spending rating of them.
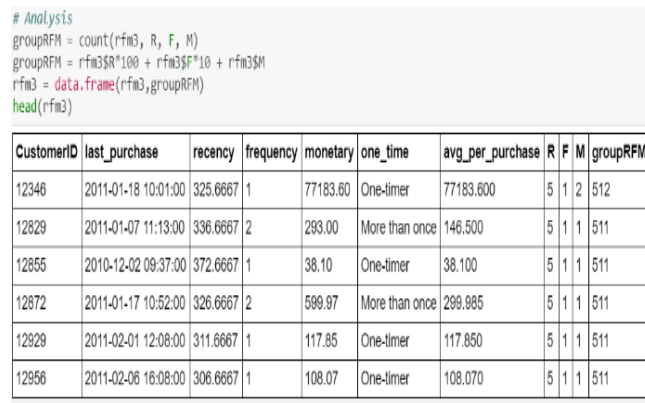


Figure.9: RFM Loyalty degree

To visualize, following is the scatter plot of Loyalty degree and RFM rating of recency in opposition to frequency. It is found that ranges are grouped together.
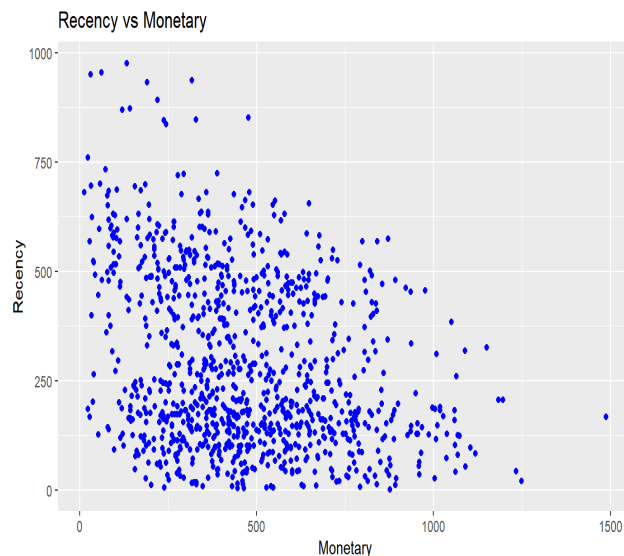


Figure.10: Scatter plot of Loyalty stage & RFM rating of recency towards Recency

## III. FUTURE WORK

The proposed fundamental cluster version given the shortage of information about the adjustments in client behaviors. Therefore distinct guidelines and techniques are essential to locate the hidden styles and buying traits of the clients. RFM and K-way helped to discover clusters of capacity clients. In addition to this Cross Selling and Market Basket Analysis strategies may be used to examine and provide extra merchandise to clients as a proposal withinside the wish that they might purchase reaping rewards the consumer and the retail status quo.

## IV. CONCLUSION

This paper offered an implementation of the okay-Means clustering set of rules for purchaser segmentation the use of statistics amassed from a web retail outfit. Our version has partitioned clients into jointly extraordinary organizations, 3 clusters in our case. This might be beneficial for making use of in addition information mining techniques and the derived insights are useful in choice making for the commercial enterprise wings.

## REFERENCES

[1]     T. Kansal, S. Bahuguna, V. Singh and T. Choudhury, "Customer Segmentation using K - means Clustering," International Conference on Computational Techniques,

Electronics and Mechanical Systems (CTEMS), Belgaum, India, 2018, pp. 135-139, doi: 10.1109/CTEMS.2018.8769171.

[2]     Aloise, A. Deshpande, P. Hansen, and P. Popat, "The Basis Of Market Segmentation"Euclidean sum-of-squares clustering," Machine Learning, vol. 75, pp.245-249, 2009.

[3]     Nurma Sari, Juni & Nugroho, Lukito & Ferdiana, Ridi & Santosa, Paulus. Review on Customer Segmentation Technique on Ecommerce. Advanced Science Letters. 2016.

[4]     Haiying Ma, "A study on Customer Segmentation for E-Commerce using the Generalized Association Rules and Decision Tree", American Journal of Industrial and Business Management , 2015,5,813-818.

[5]     Jyoti, Savita Bisnoi " A Predictive Analytics of Cluster using Associative Techniques Tool", IJRDET, Vo - 5,Issue 6, June 2016.

[6]     R.Kaur ,K.Kaur, "Data Mining on Customer Segmentation:A Review", International Journal of Advanced Research in Computer Science,Volume No-5,2017

[7]     Rachel Blasucci. Event triggered Customer Segmentation. DZone, July 23, 2018.

[8]     Kanungo, D. M. Mount, N. S. Netanyahu, C. D. Piatko, R.Silverman, and A. Y. Wu, "An efficient K-means clustering algorithm," IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 24, pp. 881-892, 2002.

[9]     Effective Cross Selling, Springer Science and Business Media LLC, 2018.

[10]     Chinedu Pascal, Simeon Ozuomba. "Application of K-Means Algorithm for Efficient Customer Segmentation: A Strategy for Targeted Customer Services", International Journal of Advanced Research in Artificial Intelligence, 2015.

[11]     Ling Luo, Bin Li et. al. "Tracking the Evolution of Customer Purchase Behaviour Segmentation via a Fragmentation-Coagulation Process", Twenty-Sixth International Joint Conference on Artificial Intelligence (IJCAI-17).

[12]     Mark K.Y.Mak,George T.S.Ho,S.L.Ting,"A Financial Data Mining Model for extracting Customer Behaviour", INTECH open access publisher, 23 July 2011.