# *RFM Segmentation and K-Means Clustering: Optimizing Customer Insights Using Elbow Method*

## I. Background

RetailCo is a medium-sized retail company specializing in consumer goods, operating through multiple channels, including in-store purchases, online sales, and direct marketing campaigns. Offering a wide range of products, such as groceries, electronics, fashion, and household goods, RetailCo aims to improve customer satisfaction and increase profitability by understanding customer behavior and optimizing its marketing strategies. The primary objective of RetailCo is to increase customer retention and maximize customer lifetime value by leveraging data analytics to identify high-value customers who are more likely to make repeat purchases, allowing the company to focus on those who contribute most to revenue. Additionally, RetailCo aims to optimize marketing efforts by targeting customers based on their purchasing behavior, thereby ensuring more relevant and effective marketing campaigns. Improved inventory management is another key focus, as understanding product demand patterns will help align stock levels with customer needs, reducing excessive inventory costs.

To thrive in a competitive market, RetailCo recognizes that it cannot satisfy all potential buyers of its products or services. Instead, it must strategically choose groups of customers it can serve better than its competitors. This necessitates segmenting customers and creating detailed profiles that enable targeted marketing strategies aimed at the most profitable segments. Customer segmentation is a tool to communicate effectively with customers. By means of the partitioning process, the characteristics of the hidden customer groups are defined in the data. The partitioning process breaks down customers according to similar features. Partitioning is a preliminary step that was created to classify the identified customer groups. With segmentation, marketers can better guide resources, and are more effective in exploring opportunities (Sabuncu, 2020).

The segmentation process serves as a foundational step that helps classify customers into groups with shared features, making it easier for RetailCo to allocate resources efficiently and identify new business opportunities. Following segmentation, profiling provides deeper insights into each segment by collecting and analyzing demographic and personal data to create comprehensive customer profiles. These profiles not only enable RetailCo to enhance service and communication with existing customers but also assist in identifying potential new customers using external data sources.

RetailCo plans to achieve its strategic objectives through key analytical goals, including conducting RFM (Recency, Frequency, Monetary) analysis along with tenure metrics to gain deeper insights into customer purchasing behavior. The RFM analytic model is proposed by Hughes (1994), and it is a model that differentiates important customers from large data by three variables (attributes), i.e., interval of customer consumption, frequency and money amount. The detail definitions of RFM model are described as follows: Recency of the last purchase (R), R represents recency, which refers to the interval between the time that the latest consuming behavior happens and present. The shorter the interval is, the bigger R is. Frequency of the purchases (F), F represents frequency, which refers to the number of transactions in a particular period, for example, two times of one year, two times of one quarter or two times of one month. The many the frequency is, the bigger F is. Monetary value of the purchases (M), M represents monetary, which refers to consumption money amount in a particular period. The much the monetary is, the bigger M is (Cheng, 2007).

RFM model has been widely applied in many practical areas in a long history, particularly in direct marketing. By adopting RFM model, decision makers can effectively identify valuable customers (Roshan, 2017). As Kahan (1998) notes, RFM is easy to use and can generally be implemented very quickly. Furthermore, it is a method that managers and decision makers can understand (McCarty, 2006). To further enhance the effectiveness of RFM, clustering techniques can be employed to group customers based on similar purchasing behaviors. Clustering is the process of grouping a set of physical or abstract objects into groups of similar objects. A cluster is a collection of data objects that are similar to one another within the same cluster and are dissimilar to the objects in other clusters (Han & Kamber,2001). K-means is one of the well-known algorithms for clustering, originally known as Forgy's method (Forgy, 1965), and it has been used extensively in various fields including data mining, statistical data analysis and other business applications. Thus, this study proposes the K-means algorithm to build clusters by attributes (i.e. R–F–M attributes). The K-means algorithm for partitioning is base on the mean value of the objects in the cluster (Cheng, 2009).

RetailCo aims to strategically identify and serve customer groups that it can cater to better than its competitors by creating detailed profiles of these segments. This approach enables the company to develop targeted marketing strategies focused on the most profitable customer segments, thereby enhancing its competitive position in the market. To achieve this, the study employs RFM analysis, which differentiates valuable customers based on recency, frequency, and monetary value of their purchases, providing a data-driven foundation for understanding and segmenting the customer base. Additionally, the study proposes the use of K-means clustering to effectively group customers with similar purchasing behaviors according to their RFM attributes. This method allows RetailCo to tailor strategies for each segment, thereby optimizing resource allocation and improving overall marketing effectiveness.

## II. Related Work

Christy et al. conducted a study about RFM Ranking-An Effective Approach to Customer Segmentation. The study builds on the Pareto principle, highlighting that a small proportion of customers significantly contribute to the company's revenue. Various segmentation methods were explored, including traditional RFM analysis, K-Means clustering, Fuzzy C-Means, and a novel RM K-Means algorithm. The RM K-Means approach modified the conventional K-Means by calculating centroids based on effective medians, resulting in faster execution and fewer iterations. The results of this study show that segmentation based on RFM values enables companies to tailor marketing strategies according to customer behavior, thereby improving communication and engagement with key customer segments. Future directions involve analyzing product preferences in each segment to optimize promotional efforts and further refine customer targeting strategies (Christy, 2018).

Dursun and Caber (2016) conducted a study in Turkey with a sample of 369 hotel customers from a population of 5,939, aiming to segment hotel customers using RFM analysis combined with K-means clustering. Their research identified eight distinct customer clusters based on RFM scores: loyal customers, loyal summer season customers, collective buying customers, winter season customers, lost customers, high potential customers, new customers, and winter season high potential customers. This segmentation allowed for a detailed understanding of customer types and their purchasing behaviors, and the researchers also compared the customers' card types with the newly established segmentation, providing valuable insights into customer loyalty and preferences within the hotel industry (Dursun, 2016).

Ina Maryani et al. conducted a study on customer segmentation based on the RFM model and clustering techniques using the K-means algorithm, focusing on transactions from Nine Reload Credit between January and December 2017. This study aimed to leverage data mining processes to perform customer segmentation, utilizing RFM (Recency, Frequency, Monetary) attributes to categorize customers. Using the K-means algorithm, implemented through Rapidminer 5.2, the study processed 82,648 transactions, resulting in the identification of 102 customers. These customers were further segmented into two clusters: Cluster 1 with 63 customers and Cluster 2 with 39 customers. The findings provide valuable insights for the company, enabling it to understand customer categories and develop targeted strategies for customer retention and engagement based on their specific RFM profiles (Maryani, 2017).

Dogan et al. conducted a study on Customer Segmentation by Using Rfm Model and Clustering Methods: A Case Study In Retail Industry. This study highlights the importance of understanding customer behavior and preferences for effective engagement and strategy development. The research demonstrates that traditional customer segmentation based solely on expenditure is insufficient. Instead, it proposes two advanced clustering models using RFM values to provide a more nuanced understanding of customer segments. By applying these models, distinct clusters are identified that offer more actionable insights compared to previous methods. The findings suggest that incorporating a broader range of data points, beyond just spending, allows for more accurate customer segmentation. This, in turn, enables businesses to tailor their marketing strategies more effectively, leading to improved customer retention and more targeted promotions. The study emphasizes the value of detailed customer profiling in developing strategic initiatives and enhancing customer relationship management (Dogan, 2018).

Abirami and Pattabiraman (2016) conducted a study in the retailing sector in India, employing RFM analysis combined with K-means clustering and association rules to enhance customer segmentation. Their approach involved classifying customers using the RFM model to analyze and estimate customer behavior. This methodology provided a systematic way to group customers based on their recency, frequency, and monetary values, offering valuable insights into customer segments and their purchasing patterns. The integration of association rules further enriched the analysis, enabling a more comprehensive understanding of customer relationships and preferences within each segment. This study emphasizes the potential of combining RFM analysis with clustering techniques and association rules to refine customer classification and develop targeted marketing strategies (Abirami, 2016).

## III.   Dataset & Features

The dataset comes from Retail Store Sales Transactions (Scanner Data) (kaggle.com). This dataset includes 64.682 transactions of 5.242 SKU's sold to 22.625 customers during one year. The dataset contains the following key attributes:
- Date: The data when the transaction occured.
- Customer_ID : A unique identifier for each customer.
- Transaction_ID : A unique identifier for each transaction.
- SKU_Category : The category of the product purchased (e.g., electronics, groceries)
- SKU : Stock Keeping Unit, a unique identifier for each product.
- Quantity : The quantity of each product purchased in a transaction.
- Sales_Amount : The total sales amount of the transaction.

Here are some of the values in the dataset :

| | Date | Customer_ID | Transaction_ID | SKU_Category | SKU | Quantity | Sales_Amount |
|---|---|---|---|---|---|---|---|
| 0 | 02/01/2016 | 2547 | 1 | X52 | 0EM7L | 1.0 | 3.13 |
| 1 | 02/01/2016 | 822 | 2 | 2ML | 68BRQ | 1.0 | 5.46 |
| 2 | 02/01/2016 | 3686 | 3 | 0H2 | CZUZX | 1.0 | 6.35 |
| 3 | 02/01/2016 | 3719 | 4 | 0H2 | 549KK | 1.0 | 5.59 |
| 4 | 02/01/2016 | 9200 | 5 | 0H2 | K8EHH | 1.0 | 6.88 |

## IV. Methods

1. RFM (Recency, Frequency and Monetary)

RFM is very valuable in predicting response and can boost a company's profits in a short term, and it is a long-familiar method to measure the strength of customer relationship as RFM can effectively identify valuable customers. In this model, recency measures the interval between the most recent transaction time and the analyzing time. Frequency measures the purchase frequency within a specified period. Monetary measures the total monetary expenditure within a specified period. This section introduces a designated RFM model to analyze the relative profitability for each customer cluster from the segmentation result after purchase-based segmentation algorithm. With this model, an enterprise can quickly find the target clusters and adjust its marketing programs and business initiatives to provide the right products, services, and resources to the target clusters. Since the RFM model measures the customer value based on Recency (R), Frequency (F), and Monetary (M) criteria (Roshan, 2017).

2. Clustering

Clustering is the process of grouping a set of physical or abstract objects into groups of similar objects.Acluster is a collection of data objects that are similar to one another within the same cluster and are dissimilar to the objects in other clusters (Cheng, 2009).

3. K-Means

K-means is the mostly popular and widely used algorithm for grouping data into groups to get right number of clusters. K-means is an iterative Algorithm which try to partition the data into k distinct groups (Shirole, 2021). Clustering Steps to follow while using K -means Algorithm:
 • Predetermine Number of clusters K.
 • Initialize Centroid by randomly selecting K data points.
 • Compute the distance of the next data points with all centroids.
 • Assign the data point to the nearest cluster
 • Repeat this step until all data points converges to a cluster.

Formula for Centroid Determination:

$$C_i = {}^1\!/_M \sum_{j=1}^{m} X_j$$

Formula for Euclidean Distance:

$$d\,(p,\,q) = \sqrt{(p1\text{-}q1)^2 + (p2\text{-}q2)^2}$$

4. Elbow Method

Elbow method is used to determine the optimal number of clusters based on the dataset. The idea is simple behind it, i.e., plotting the SSE (Sum squared Error) against suitable no of cluster value. Then we will select the value at which there is maximum curve in the graph (Shirole, 2021).
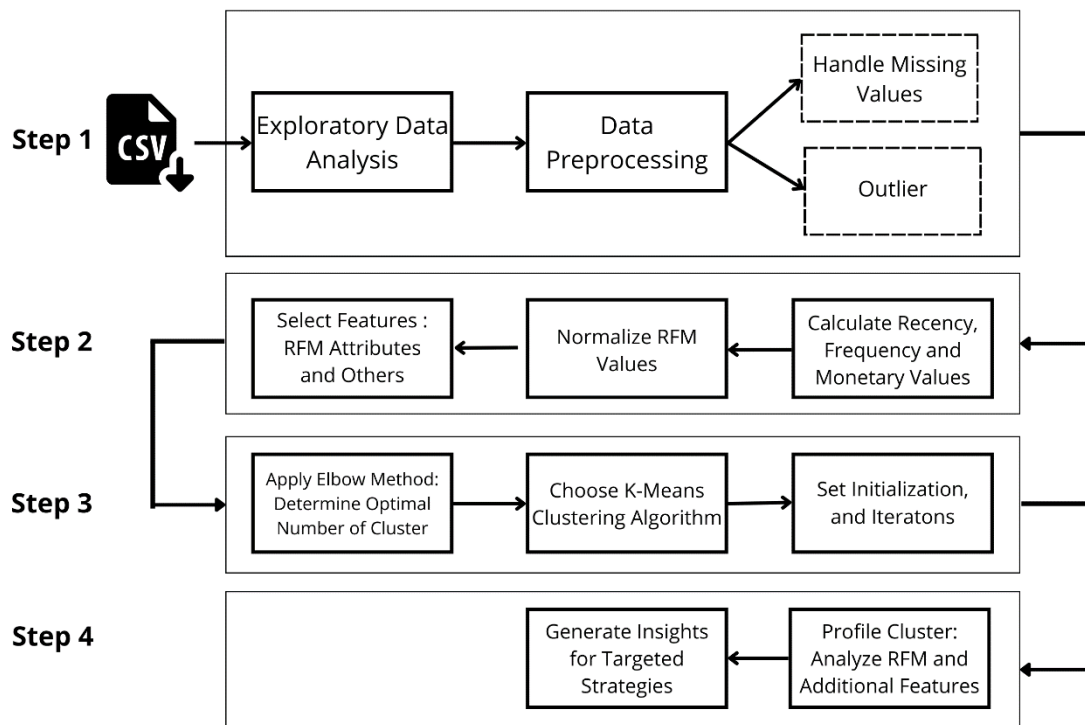
SSE (Sum Square Error) is one of the statistical methods used to measure the total difference from the actual value of the value achieved.

$$SSE = \sum_{i=1}^{n} (d)^2$$

Where, d is the distance between the data and the Cluster center. Sum of Square Error (SSE) is a formula used to measure the difference between the data obtained by the prediction model that has been done previously. SSE is often used as a research reference in determining optimal clusters (Nainggolan, 2019).

## V. Experiments/Result/Discussion

In this project, 5 stages are carried out which can be seen in the figure below:
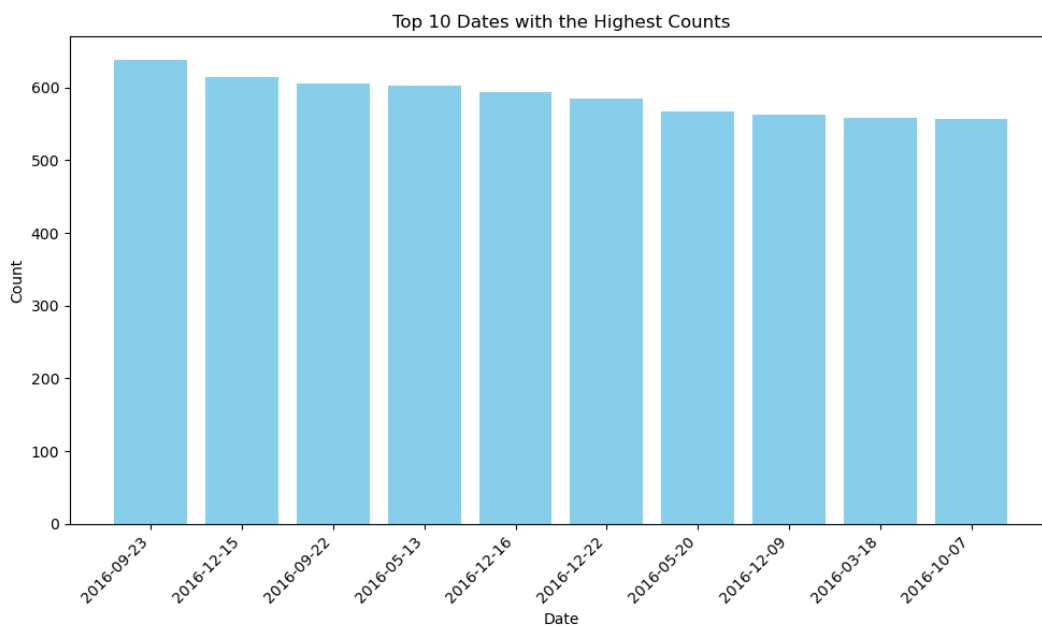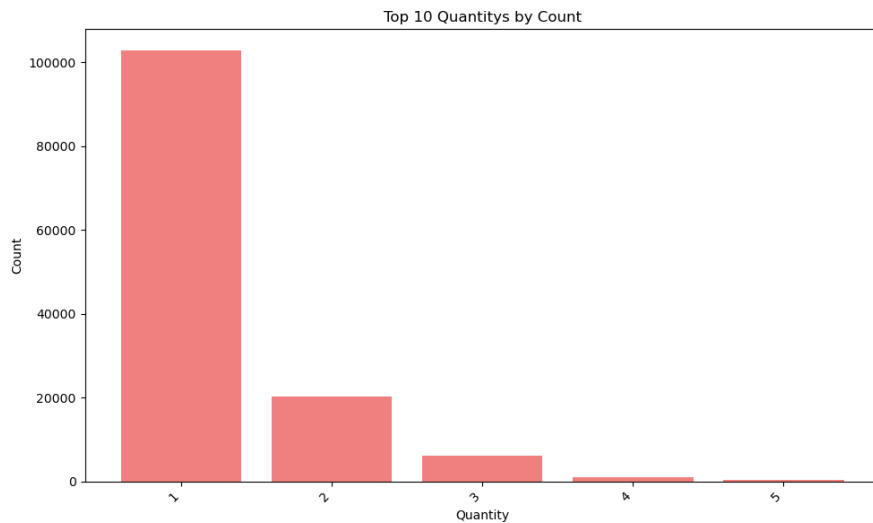
1. Data Collection and Preprocessing

In the data collection and preprocessing stage, the first step is to collect data from relevant sources. Once the data is available, the next step is to conduct Exploratory Data Analysis (EDA). EDA helps understand the overall distribution of the data, including the scatter, central tendency, and skewness of key features such as Recency, Frequency, and Monetary values. This is important to identify patterns and trends in customer behavior that can guide further analysis.

By visualizing and summarizing data, EDA makes it possible to detect outliers and anomalies that may distort the results of RFM calculations or clustering. In addition, EDA helps in exploring the relationships between various variables, which is useful in selecting the most relevant features for segmentation. EDA also helps assess data quality by highlighting issues such as missing values, incompatible data types, or inconsistent data entries.
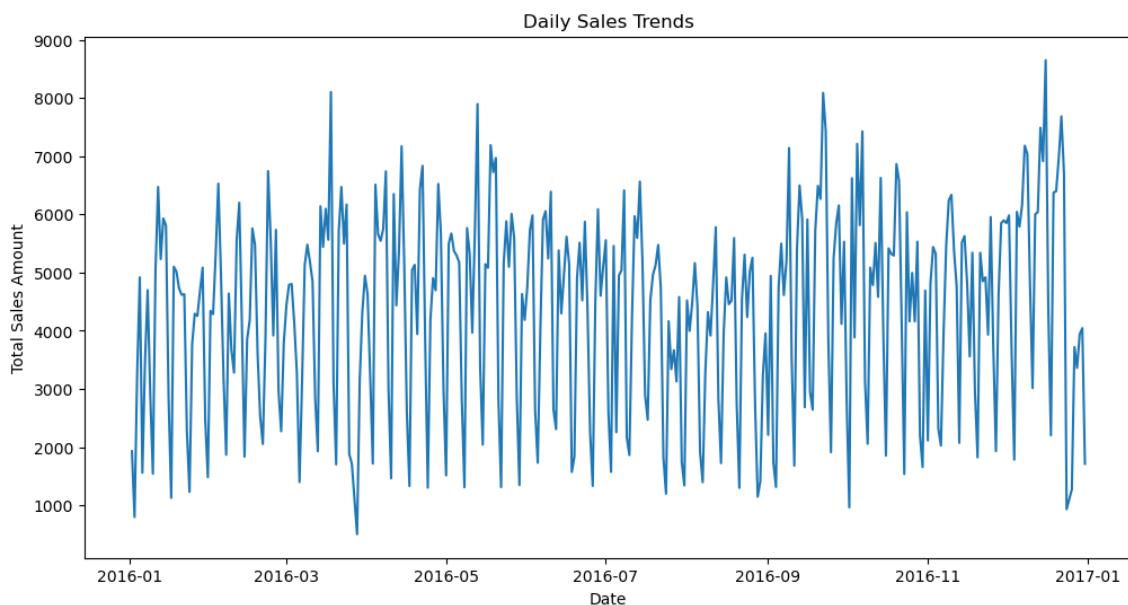
The following are some of the results of Exploratory Data Analysis (EDA):



From the image above, it can be seen that the dates with the highest transaction counts are at the top of the list, such as September 23, 2016, with 638 transactions, December 15, 2016, with 614 transactions, and September 22, 2016, with 606 transactions. These dates indicate periods of peak activity, which could coincide with sales events, holidays, or promotional campaigns that drove a large number of transactions. In contrast, dates like March 28, 2016, had only 73 transactions, suggesting lower customer engagement or normal, non-promotional periods. The transactions are spread across 363 unique dates, indicating that the data covers nearly a full year's worth of daily transactions, although some dates have significantly more activity than others.

Top 10 Quantitys by Count

From the figure above, it can be seen that 1 is the most commonly found quantity, with 102,741 occurrences. This shows that single-unit purchases are very prevalent in the data set. This may indicate that many transactions are for individual items rather than bulk purchases.



Daily Sales Trends

The line chart above illustrates daily sales trends throughout 2016, showing fluctuations in sales amounts over time. The x-axis represents the dates from January 2016 to early January 2017, while the y-axis shows the total sales amount. The sales amounts vary noticeably, with several peaks and troughs throughout the year, suggesting periods of heightened activity followed by dips. There are significant spikes in sales during certain periods, such as around September and December, which could indicate seasonal patterns, promotional events, or holidays that drive increased sales. Despite the frequent fluctuations, the overall trend appears relatively stable, without a clear long-term increase or decrease, indicating a consistent sales pattern with periodic surges. High activity periods are particularly visible in mid-2016 and towards the end of the year, especially around December, possibly due to holiday shopping or end-of-year promotions.

Average Sales by Cohort and Duration in Months

From the heatmap, it is evident that some cohorts exhibit strong sales performance in the initial months, such as the January 2016 cohort, which shows consistently high average sales values in the early months and again later around month 11. The March 2016 cohort also stands out with notable increases in sales, reaching its peak around the 8th and 9th months, suggesting a surge in customer engagement or purchasing activity during this period.

The overall pattern reveals that many cohorts start with moderate average sales in the first month, with some increasing or maintaining strong sales in the following months. However, a general decline in average sales can be observed as the duration increases, which is a typical trend in cohort analyses where customer engagement tends to decrease over time. Notably, the March 2016 cohort demonstrates a sustained and even rising trend in average sales over a longer period, indicating successful retention or re-engagement strategies that kept these customers active.
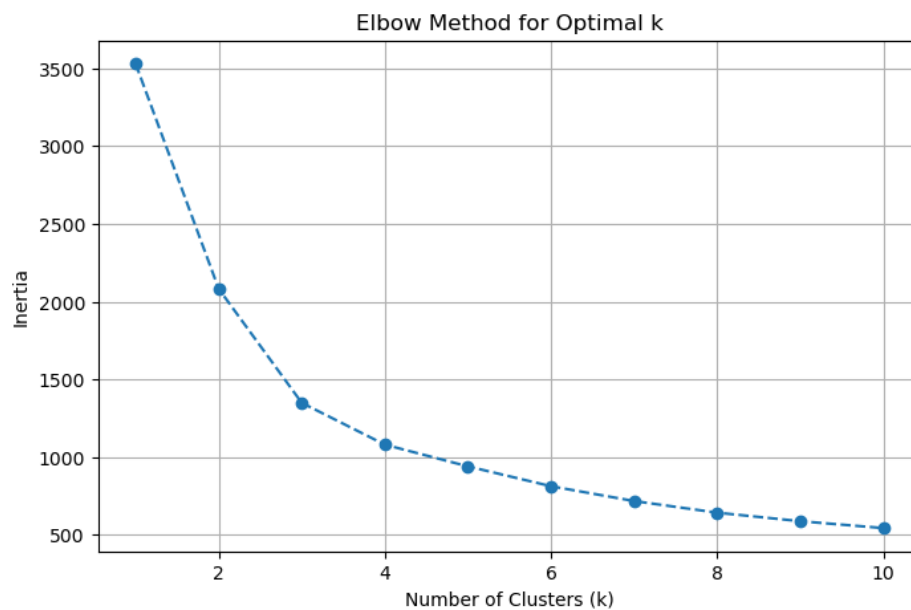
2.  RFM Analysis
    At the RFM Analysis stage, recency, frequency and monetary values are calculated until the results are obtained like this.

| Customer_ID | Recency | Tenure | Frequency | Monetary |
|---|---|---|---|---|
| 1 | 345 | 0 | 1 | 16.29 |
| 2 | 196 | 87 | 2 | 22.77 |
| 3 | 365 | 0 | 1 | 10.92 |
| 4 | 53 | 121 | 2 | 33.29 |
| 5 | 180 | 147 | 5 | 78.82 |

3. Determine Optimal Clusters and K-Means Clustering

At this stage use the Elbow Method to plot the sum of squared errors (SSE) against the number of clusters. The results of the elbow method calculation can be seen in this picture.



Elbow Method for Optimal k

Based on the elbow plot, the inertia drops significantly when moving from 1 to 2 clusters and again from 2 to 3 clusters, indicating that clustering improves substantially by increasing the number of clusters from 1 to 3. From k = 3 onwards, the decrease in inertia becomes less steep, suggesting diminishing returns for additional clusters. Although there is still a reduction in inertia when increasing from 3 to 4, 5, and so on, it is not as pronounced as in the earlier steps. K = 3 emerges as the strongest candidate since it balances a significant reduction in inertia with a reasonable number of clusters, capturing the major structure in the data without over-complicating the model. Adding more clusters beyond k = 3, such as k = 4 or 5, shows diminishing returns in terms of inertia reduction, with only marginal improvements in cluster quality, which suggests that additional clusters do not significantly enhance differentiation among data points. Based on this analysis, the K-means clustering algorithm is chosen with the optimal number of clusters set to 3. Next Choose the K-means clustering algorithm based on the optimal number of clusters identified, which is 3. Then Configure the algorithm settings (e.g., distance metric, initialization). The centroid results of each cluster can be seen in the following image.

|         | Recency  | Frequency | Monetary |
|---------|----------|-----------|----------|
| Cluster |          |           |          |
| 0       | 0.331689 | 0.495232  | 0.600857 |
| 1       | 0.210824 | 0.065928  | 0.175488 |
| 2       | 0.766271 | 0.059773  | 0.184537 |

Next display the first few rows of the DataFrame with the original scale data and cluster labels.

|             | Recency | Frequency | Monetary | Label |
|-------------|---------|-----------|----------|-------|
| Customer_ID |         |           |          |       |
| 1           | 345.0   | 1.0       | 16.29    | 2     |
| 2           | 196.0   | 2.0       | 22.77    | 2     |
| 3           | 365.0   | 1.0       | 10.92    | 2     |
| 4           | 53.0    | 2.0       | 33.29    | 1     |
| 5           | 180.0   | 5.0       | 78.82    | 0     |

4. Cluster Profiling and Insights
   At this stage, cluster profiling is done to understand the characteristics of each segment.
   - Cluster 0 = Top Tier Customers (low recency, high frequency and high monetary)
   - Cluster 1 = Regular Customers (low recency, medium frequency and low monetary)
   - Cluster 2 = Passive Customers (high recency, low frequency and low monetary

   The data can be seen in the following figure.

|             | Recency | Frequency | Monetary | Label | Customer Type     |
|-------------|---------|-----------|----------|-------|-------------------|
| Customer_ID |         |           |          |       |                   |
| 1           | 345.0   | 1.0       | 16.29    | 2     | Passive Customer  |
| 2           | 196.0   | 2.0       | 22.77    | 2     | Passive Customer  |
| 3           | 365.0   | 1.0       | 10.92    | 2     | Passive Customer  |
| 4           | 53.0    | 2.0       | 33.29    | 1     | Regular Customer  |
| 5           | 180.0   | 5.0       | 78.82    | 0     | Top Tier Customer |
| ...         | ...     | ...       | ...      | ...   | ...               |
| 22620       | 61.0    | 1.0       | 8.60     | 1     | Regular Customer  |
| 22621       | 22.0    | 1.0       | 9.69     | 1     | Regular Customer  |
| 22622       | 16.0    | 1.0       | 6.07     | 1     | Regular Customer  |
| 22624       | 30.0    | 1.0       | 19.60    | 1     | Regular Customer  |
| 22625       | 1.0     | 4.0       | 83.62    | 0     | Top Tier Customer |

Next, perform aggregation and statistical analysis based on the cluster profiling results. After identifying various customer types (such as Top Tier Customers, Regular Customers, and Passive Customers) using RFM (Recency, Frequency, Monetary) clustering, the data is then summarized to gain further insight into the performance of each segment. The data can be seen in the following figure.

| Customer Type | Recency | Frequency | Monetary | Count Customers | Total Monetary | % Revenue | Avg Revenue per Transaction | % of Total Customers | Median Recency | Monetary Std Dev |
|---|---|---|---|---|---|---|---|---|---|---|
| Passive Customer | 279.90 | 1.30 | 18.31 | 8290 | 151789.90 | 33.36 | 14.08 | 45.59 | 284.0 | 15.75 |
| Regular Customer | 77.75 | 1.33 | 17.43 | 6773 | 118053.39 | 25.94 | 13.11 | 37.25 | 74.0 | 12.70 |
| Top Tier Customer | 121.68 | 3.48 | 59.34 | 3121 | 185200.14 | 40.70 | 17.05 | 17.16 | 97.0 | 20.23 |

Based on the analysis and insights from the RFM segmentation, total monetary value highlights the absolute contribution of each customer segment to the total revenue. For example, if Top Tier Customers contribute the highest total monetary value, it underscores their critical role in driving revenue, even though they might represent a smaller portion of the customer base. These Top Tier Customers are expected to have a high average revenue per transaction, reflecting their substantial spending power per purchase, whereas lower values in Passive Customers suggest potential for upselling or cross-selling opportunities. The percentage of Total Customers helps to assess the size of each segment relative to the entire customer base, with a notably high percentage coming from Passive Customers.

The median recency metric offers a clearer understanding of the recency of interactions within each segment. Top Tier Customers typically display low recency values, indicating more recent transactions and ongoing engagement. However, a high standard deviation in monetary value among Top Tier Customers might point to a wide range of spending behaviors within this group. Despite being the most valuable in terms of revenue per transaction and total monetary contribution, Top Tier Customers may be fewer in number. Therefore, focusing on personalized offers and loyalty programs could be crucial for enhancing retention and maximizing value from this segment. Regular Customers, on the other hand, show consistent spending patterns but at a lower frequency, suggesting that increasing engagement through reminder communications, personalized offers, and loyalty incentives could help elevate them to Top Tier status.

Lastly, Passive Customers represent a large segment with significant potential, yet their low frequency and high recency indicate they are less engaged. To unlock their potential, re-engagement strategies such as targeted marketing, personalized content, and reactivation campaigns are essential. By understanding these insights, businesses can tailor their strategies to effectively engage each segment, maximizing their revenue potential and strengthening customer relationships.

## VI. Conclusion/Future Work

The strategic plan for RetailCo, derived from an in-depth RFM segmentation of Top Tier Customers, Regular Customers, and Passive Customers, is designed to maximize customer lifetime value and drive overall revenue growth through tailored approaches for each customer group. Top Tier Customers, who make frequent purchases with high monetary value, are pivotal to revenue despite comprising a smaller proportion of the customer base.

Immediate strategies involve implementing a VIP rewards program, personalized promotions, and premium services to enhance engagement. Over the mid to long term, introducing subscription models, referral incentives, and a tiered loyalty program will aim to sustain their high-value engagement. Regular Customers, with moderate frequency and lower transaction values, contribute significantly to revenue but have potential for increased engagement. Short-term actions focus on personalized product suggestions and targeted campaigns, while mid-term strategies include escalating discounts and seasonal promotions to drive more frequent purchases. Long-term efforts involve continuous feedback to refine their shopping experience and personalized retargeting to maintain engagement. Passive Customers, the largest group with moderate recency but lower frequency and spending, are targeted with time-sensitive offers and reactivation campaigns in the short term to boost engagement. In the mid to long term, promoting product bundles and providing valuable content are key to retaining their interest.

To further enhance the effectiveness of these strategies, future work could involve leveraging advanced algorithms or methods such as predictive modeling to identify at-risk customers or potential high-value segments early on. Machine learning algorithms like clustering, decision trees, or even neural networks could provide deeper insights into customer behavior patterns, allowing for more dynamic and responsive marketing strategies. Additionally, incorporating A/B testing or multivariate testing of various engagement tactics could refine and optimize the strategies, ensuring they are both effective and aligned with evolving customer preferences.

## Reference

Abirami, M. & Pattabiraman, V. 2016. Data mining Approach for Intelligent Customer Behavior Analysis for A Retail Store. In: proceedings of the 3rd international symposium on big data and cloud computing challenges (ISBCC–16). 283-291.

Ching-Hsue Cheng, You-Shyang Chen C.-H. Cheng, Y.-S. Chen. 2008. Classifying the Segmentation of Customer Value Via RFM Model and RS Theory. Expert Systems with Applications 36 (2009) 4176–4184

Christy, A.J., Umamakeswari, A., Priyatharsini, L., Neyaa, A. 2018. RFM Ranking – An Effective Approach to Customer Segmentation, Journal of King Saud University - Computer and Information Sciences

Dursun, A. & Caber, M. (2016). Using Data Mining Techniques for Profiling Profitable Hotel Customers: An Application of RFM Analysis. Tourism Management Perspectives, 18, 153-160

Hadi Roshan, Masoumeh Afsharinezhad. 2017. The new Approach in Market Segmentation by Using RFM Model. Journal of Applied Research on Industrial Engineering

Ina Maryani, Dwiza Riana, Rachmawati Darma Astuti, Ahmad Ishaq, Sutrisno, Eva Argarini Pratama. 2018. Customer Segmentation based on RFM model and Clustering Techniques With K-Means Algorithm. Third International Conference on Informatics and Computing (ICIC)

John A. McCarty a, Manoj Hastak. 2006. Segmentation Approaches in Data-Mining: A Comparison of RFM, CHAID, and Logistic Regression Journal of Business Research 60 (2007) 656–662.

Onur Dogan, Ejder Aycin, Zeki Atil Bulut. 2018. Conducted A Study About Customer Segmentation By Using RFM Model And Clustering Methods: A Case Study In Retail Industry. International Journal of Contemporary Economics and Administrative Sciences ISSN: 1925– 4423 Volume :8, Issue: 1, Year:2018, pp. 1-19

Rahul Shirole, Laxmiputra Salokhe, Saraswati Jadhav. 2021. Customer Segmentation using RFM Model and K-Means Clustering. Department of Computer Engineering, Vishwakarma Institute of Technology Pune, Maharashtra, India

Rena Nainggolan, Resianta Perangin-angin, Emma Simarmata and Astuti Feriani Tarigan. 2019. Improved the Performance of the K-Means Cluster Using the Sum of Squared Error (SSE) optimized by using the Elbow Method. Journal of Physics: Conference Series

Sabuncu, İ., Türkan, E., & Polat, H. 2020. Customer Segmentation And Profiling With RFM Analysis, TUJOM, 5(1): 22-36