# International Journal of Research Publication and Reviews

Journal homepage: www.ijrpr.com  ISSN 2582-7421

# Customer Segmentation for Enhancing Business Strategy

## [1]Suryakanthan. M, [2]Vimal. K, [3]Sanjay Raj. R, [4]Mr. Thiruselvan. P M. E

[1]Dept. of Computer Science, P.S.R Engineering College, Sivakasi, India E-mail: 20cs106@psr.edu.in
[2]Dept. of Computer Science, P.S.R Engineering College , Sivakasi, India, e-mail: 20cs119@psr.edu.in
[3]Dept. of Computer Science, P.S.R Engineering College, Sivakasi, India, E-mail: 20cs090@psr.edu.in
[4]Associate Professor, Dept. of Computer Science, P.S.R Engineering College, Sivakasi, India, e-mail: thiruselvan@psr.edu.in

## ABSTRACT

Understanding consumer behavior and preferences is essential for developing successful marketing strategies and ensuring customer satisfaction in today's cutthroat corporate environment. The technique of grouping a company's consumers according to shared traits is known as customer segmentation, which enables businesses to successfully and correctly market to each group. The objective of this project is to segment clients according to their demographic data and purchase patterns using the K-Means clustering technique. Through the examination of a heterogeneous dataset that encompasses a range of consumer traits, our goal is to identify discrete customer categories that can inform tailored marketing campaigns, product suggestions, and enhanced customer experiences. This project has a lot of potential advantages. Companies can increase revenue and customer happiness by customizing their marketing strategies to each customer segment's specific needs. Furthermore, product development and inventory management can be optimized to cater to specific customer preferences.

Keywords: Clustering, Elbow Method, K-Means, Algorithm, Customer Segmentation, Visualization.

## 1. Introduction

In business the most important component is data. With the help of grouped or ungrouped data, we can perform some operations to find customer interests. Data mining helpful to extract data from the database in a human-readable format. However we may not know the actual beneficiaries in the whole dataset. Customer segmentation is a strategic approach that involves dividing a customer base into distinct groups based on specific characteristics, behaviors, or preferences. This segmentation strategy is crucial for businesses aiming to personalize their marketing efforts, improve customer satisfaction, and enhance overall business performance. In this project, the primary objective is to analyze and segment customers within a particular business or industry context.

By leveraging data-driven techniques, statistical analyses, and machine learning algorithms, the project aims to identify meaningful patterns and clusters within the customer data. The rationale behind customer segmentation lies in the recognition that not all customers are identical. Each customer group may exhibit unique purchasing behaviors, preferences, and responses to marketing strategies. By understanding and categorizing these differences, businesses can tailor their products, services, and marketing campaigns to better meet the distinct needs of each segment.

Segmentation allows for more targeted and personalized approaches, enabling businesses to allocate resources efficiently, design specific marketing strategies, and provide customized offerings. Ultimately, this approach contributes to improved customer satisfaction, increased customer retention, and more effective utilization of resources. This project's findings will empower businesses with insights to make data-driven decisions, optimize marketing strategies, and enhance their competitive edge in the market. Through the segmentation of customers, businesses can build stronger relationships, drive loyalty, and achieve sustainable growth by meeting the diverse needs of their customer base.

### 1.1 Problem Statement

Customer segmentation is helpful in breaking out the big dataset of data into multiple groups according to factors like gender, age, spending patterns, income, and demographics. These groups areal so known as clusters. This allows us to learn things like which product sold a lot and what age group is buying it, among other things. We can also supply that goods much more readily for improved income creation. We will start with the historical data. Since the old is gold, we will apply the K-means clustering technique on the old data and first determine the number of clusters. Finally, the data must be visualized. Looking at the representation makes it easy to identify the possible set of data. This paper's objective is to discover consumer subgroups through the use of the K-means clustering algorithm, a partitioning technique used in data mining.

The best clusters are identified using the PCA approach. Upon applying the K-Means clustering method to this dataset, consumer groups exhibiting comparable purchasing behaviors and demographic traits will be identified. The intrinsic structure of the data and, if required, iterative testing will be

used to determine the number of clusters. Businesses can classify their clientele into discrete categories by gaining a deep understanding of them once the clusters are established.

These client categories have numerous practical ramifications. Companies can create specialized marketing plans by identifying the particular requirements and preferences of every target niche. This personalization may appear as customized Personalized product recommendations and marketing messaging are two ways that this personalization can deepen the bond between the consumer and the company. Additionally, companies can streamline their product creation procedures by matching them to the tastes of particular clientele groups. Additionally, inventory management can be improved to minimize stockouts and overstock situations by guaranteeing that products that appeal to each segment are easily accessible. Simply said, consumer segmentation involves assembling groups of clients based on shared traits.

Geographical, demographic, behavioral, purchasing power, situational, personality, lifestyle, psychographic, and other aspects are some examples of these traits. Customer acquisition, customer retention, profitability growth, customer happiness, resource allocation through program or marketing measure design, and enhanced target marketing are the objectives of customer segmentation. An effective method for client segmentation is clustering. Clustering arranges similar data points. These groupings are referred to as clusters. Each cluster's objects differ from those of other clusters even though they are similar among themselves.

Unsupervised learning is a subset of machine learning that includes clustering as one of its data mining techniques. This is a result of its capacity to extract information and patterns from unlabeled data. It is widely applied to pattern recognition, classification, and machine learning. The K-means method, hierarchical clustering, and DBSCAN are examples of clustering algorithms. The k-means clustering method has been used in this project to segment customers. The partitioning principle is the foundation of the K-means clustering technique. The number of clusters selected is denoted by the letter k. It is the most widely used method based on centroids. The split of a market into distinct customer groups with comparable traits is known as customer segmentation. Finding unmet customer demands can be effectively accomplished through the use of customer segmentation. Companies can then surpass the competition by creating distinctively appealing goods and services using the aforementioned data.

## 2. Existing System

Recency, frequency, and monetary value, or RFM, models did not take into account the interactions that occur between users and items; instead, they concentrated on a single form of user behavior data: purchase behavior.

Based on the Self-Organizing Map (SOM) algorithm, it deconstructs the various behaviors of customers within a certain time frame and carries out customer segmentation in an application promotion system known as multi-behavior RFM (MB-RFM). We used the superiority chart and entropy value approaches to examine the weight relationship between various user actions and items, taking into account the R, F, and M values. An enhanced SOM neural network was used to classify customers based on the values of the MB-RFM model, which were thought to include every behavior that a customer had stated.

## 3. Proposed System

K-Means Clustering is an unsupervised learning algorithm that groups the unlabeled dataset into different clusters and is used to solve the clustering conflicts by detecting the centroid of each cluster. Principal Component Analysis is used for dimensionality reduction It is a statistical process that converts the observations of correlated features into a set of linearly uncorrelated features with the help of orthogonal transformation. By implementing the K-means and PCA algorithm we will sort the targeted customers so that the sense can be given to the marketing team and plan the strategy accordingly. To deal with segmenting the shopping mall dataset that contains attributes like Customer ID, age, gender, annual income, and spending score by applying the k means algorithm and principal component analysis.

### 3.1 System Architecture

Data preprocessing refers to the technique of cleaning, transforming, and organizing raw data into a format suitable for analysis or modeling in the context of machine learning or data analysis. It's a crucial initial step in the data analysis process that significantly influences the accuracy and effectiveness of any subsequent data-driven modeling or analysis. Gather relevant data from various sources such as CRM systems, transactional databases, social media, etc. This data might include demographics, purchase history, browsing behavior, etc. Missing values, outliers, and inconsistencies in the dataset should be focused.

Impute missing values using techniques like mean, median, or sophisticated imputation methods. Outliers might be treated by removing them or transforming them based on the context. Identify which features are relevant for your segmentation. Eliminate irrelevant or redundant features that do not contribute to the segmentation goal.

Feature selection techniques like correlation analysis or feature importance can be utilized. Normalize or standardize numerical features to bring them to a similar scale. This step is crucial for algorithms sensitive to feature scales, like K-means or neural networks. Convert categorical variables into numerical format. Techniques like one-hot encoding or label encoding can be used based on the nature of the data and the algorithm to be employed. If dealing with high-dimensional data, consider techniques like Principal Component Analysis (PCA) or t-SNE to reduce the dimensionality while preserving the most
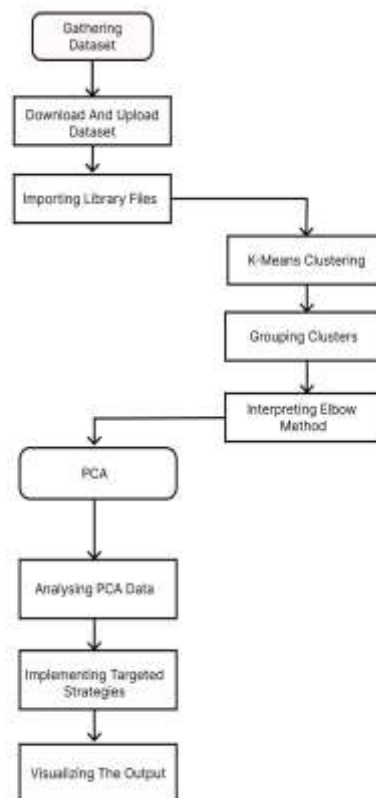
critical information. If the dataset is imbalanced, meaning one segment greatly outnumbers the others, techniques like oversampling, under-sampling, or using algorithms that handle class imbalance can be applied. Partition the dataset into training and testing sets. The training set is used to build the model, while the testing set is used to evaluate its performance. Ensure the data is in a suitable format for the chosen algorithms. Some algorithms perform better with normalized data. For example, scaling data to a range between 0 and 1 can be beneficial.

Creating new features from existing ones if necessary. For example, you might derive new features from the existing dataset that could be more informative for segmentation purposes. Some of the algorithms might perform better with transformed data. For instance, logarithmic transformations or Box-Cox transformations can make the data more normally distributed, which might benefit certain models.

Validating the preprocessing steps and iterate as necessary. Ensure that the preprocessing steps applied do not introduce biases or information leaks into the model. Each and every step should be adapted based on the nature of your dataset, the segmentation goals, and the machine learning algorithms you intend to use. Always keep in mind that data preprocessing significantly impacts the performance and accuracy of your machine learning model.

### 3.1.2 K-means Clustering

unsupervised learning technique called K-means clustering is used in data science and machine learning to address clustering issues. We will study the definition of the K-means clustering method, its operation, and its Python implementation in this topic. Each data point is assigned to a group iteratively, and over time, the data points are grouped according to common qualities.
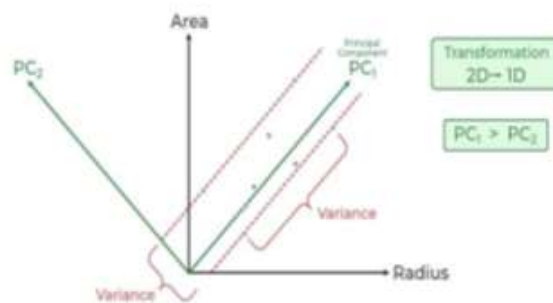


### 3.1.2 Principal Component Analysis:

In machine learning, dimensionality reduction is achieved by the use of PCA, an unsupervised learning technique. It is a statistical procedure that uses orthogonal transformation to turn a set of correlated feature observations into a set of linearly uncorrelated features. The Principal Components are those newly altered features. It is a widely used tool for both predictive modeling and exploratory data analysis. It is a method for minimizing variances in order to extract meaningful patterns from the provided dataset. looks for a lower-dimensional surface to project the high-dimensional data onto in most cases. It decreases dimensionality by taking into account each characteristic's variance because a high attribute indicates a good split between the classes. Real-world uses include image processing, movie recommendation systems, and power allocation optimization across several channels of communication. Since it uses a feature extraction technique, the significant variables are retained while the less significant ones are removed.

It functions under the requirement that the variance of the data in the lower dimensional space should be maximal even when data in a higher dimensional space is mapped to data in a lower dimensional space. An orthogonal transformation is used in the statistical process of principal component analysis

(PCA) to turn a set of correlated variables into a set of uncorrelated variables. It is the most used tool for machine learning predictive models and exploratory data analysis.



Furthermore, an unsupervised learning algorithm technique called Principal Component Analysis is utilized to look at how a set of variables relate to one another. Regression is used to find the line of greatest fit, and this method is sometimes referred to as generic factor analysis. Without any prior knowledge of the target variables, Principal Component Analysis's primary objective is to decrease a dataset's dimensionality while maintaining the most significant patterns or correlations between the variables. By identifying a new set of variables that are smaller than the original set of variables, preserve the majority of the sample's information, and are helpful for data regression and classification, PCA seeks to minimize the dimensionality of a given data set.

## 4. Conclusion

By maximizing marketing tactics, organizations can establish a connection with their core consumer base through effective customer segmentation. Businesses can use segmentation analysis to identify the highest-value customer segments in their markets and target them with the appropriate product at the right moment. Organizations that grasp this strategy can boost revenue and establish connections with a larger portion of their target market. Finding the segment that generates the highest return on investment for your marketing efforts is one of the main advantages of customer segmentation study. Sales teams may broaden their customer base, boost income, and close more leads by concentrating their efforts on their most lucrative market segments. Once you understand the company's customers and identify the right strategies for connecting with them, you can personalize aspects of its offerings and marketing materials. In conclusion, the implementation of K-means clustering in conjunction with Principal Component Analysis (PCA) in the customer segmentation project has proven to be a powerful and effective methodology for identifying and categorizing distinct customer groups based on their characteristics, behaviors, and preferences. The utilization of K-means clustering has allowed for the creation of homogeneous customer segments by grouping individuals with similar traits, enabling businesses to gain valuable insights into their customer base. By iteratively assigning customers to the nearest cluster center based on their features, K-means has efficiently partitioned the data, facilitating the identification of meaningful and well-defined customer groups. Furthermore, integrating PCA into the segmentation process has been instrumental in reducing dimensionality and enhancing the quality of clustering. PCA enabled the reduction of the feature space, maintaining relevant information and enhancing the clustering algorithm's efficiency by capturing the variance in the data. Through this project, we have successfully unveiled clusters that represent different customer personas, each exhibiting distinct preferences, purchase behaviors, and responses to marketing stimuli. These segments provide a framework for developing targeted marketing strategies, personalized offerings, and tailored communication approaches to cater to the specific needs of each group. The outcomes of this project offer significant potential for businesses to optimize resource allocation, refine marketing campaigns, and enhance customer experiences. By leveraging the insights derived from K-means clustering and PCA, companies can better tailor their products, services, and marketing efforts to meet the diverse needs of their customer segments, ultimately leading to improved customer satisfaction, heightened customer engagement, and potentially increased business performance. While K-means and PCA have proven to be robust techniques for customer segmentation, it's essential to note that continuous refinement and validation of these segments are necessary for ongoing relevance and accuracy. Incorporating additional data sources, employing different algorithms, and periodically reassessing the segments' validity will ensure that businesses maintain a comprehensive understanding of their customer base. Overall, the integration of K-means clustering and PCA within the customer segmentation process has provided a solid foundation for businesses to better understand, categorize, and effectively cater to their diverse customer base, ultimately fostering stronger relationships and driving success in today's competitive market landscape.

### 4.1 Future Enhancement

Implementing dynamic or adaptive segmentation models that can continuously evolve based on real-time data updates. This allows for the adaptation of segments as customer behaviors and preferences change. Integrating diverse data sources (social media, geospatial, etc.) to enrich the customer segmentation process. This would enable a more comprehensive understanding of customer behaviors and preferences, thereby refining the segmentation accuracy. Utilizing advanced feature engineering and selection techniques to enhance the quality of input variables for clustering. This involves identifying and incorporating the most impactful features and removing noise from the dataset, thereby improving the clustering performance. Developing real-time segmentation capabilities to enable immediate, personalized marketing actions. Integrating these segments into marketing automation systems

can allow for real_time, tailored customer interactions. Moving beyond descriptive segmentation to predictive models that forecast potential future behaviors or preferences of customers. Utilizing predictive analytics can help in proactively meeting customer needs. Enhancing the interpretability of the segmentation results by implementing visualization techniques or providing clear, actionable insights. This can help business stakeholders easily comprehend and act upon the segmentation findings. Incorporating contextual information, such as the customer's journey or their interactions with the business, to create more nuanced segments. Understanding the context of customer actions can lead to more precise segment definitions. Focusing on responsible data usage, ensuring compliance with privacy regulations, and adopting ethical practices in handling customer data.
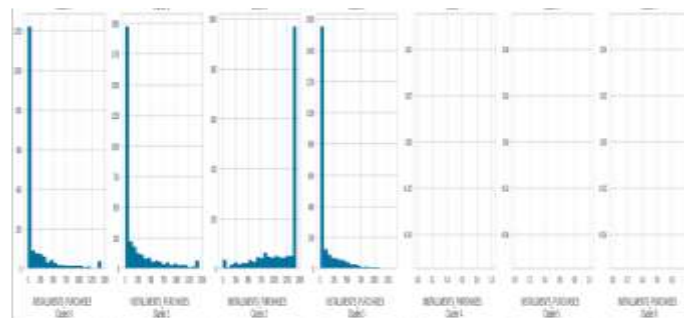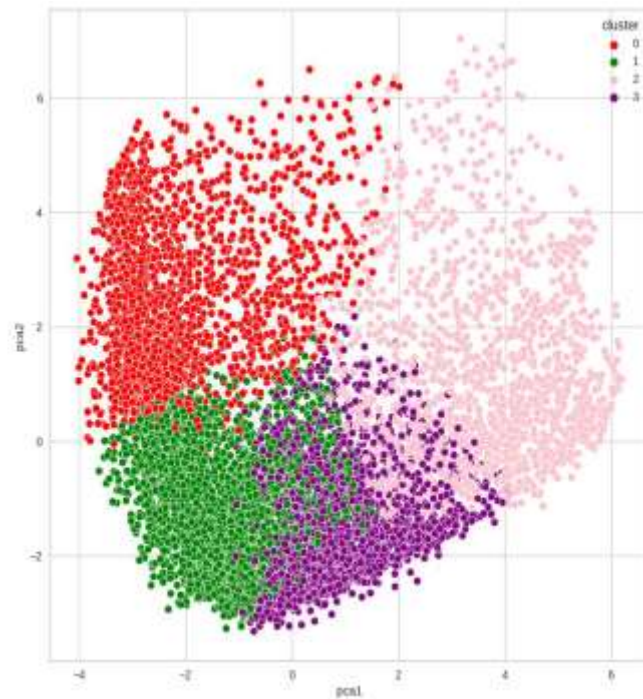
## 5. Sample Output

| | CUST_ID | BALANCE | BALANCE_FREQUENCY | PURCHASES | ONEOFF_PURCHASES | INSTALLMENTS_PURCHASES | CASH_ADVANCE | PURCHASES_FREQUENCY | ONEOFF_PURCHASES_FREQUENCY |
|---|---|---|---|---|---|---|---|---|---|
| 0 | C10001 | 40.900749 | 0.818182 | 95.40 | 0.00 | 95.40 | 0.000000 | 0.166667 | 0.000000 |
| 1 | C10002 | 3202.467416 | 0.909091 | 0.00 | 0.00 | 0.00 | 6442.945483 | 0.000000 | 0.000000 |
| 2 | C10003 | 2495.148862 | 1.000000 | 773.17 | 773.17 | 0.00 | 0.000000 | 1.000000 | 1.000000 |
| 3 | C10004 | 1666.670542 | 0.636364 | 1499.00 | 1499.00 | 0.00 | 205.788017 | 0.083333 | 0.083333 |
| 4 | C10005 | 817.714335 | 1.000000 | 16.00 | 16.00 | 0.00 | 0.000000 | 0.083333 | 0.083333 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 8945 | C19186 | 28.493517 | 1.000000 | 291.12 | 0.00 | 291.12 | 0.000000 | 1.000000 | 0.000000 |
| 8946 | C19187 | 19.183215 | 1.000000 | 300.00 | 0.00 | 300.00 | 0.000000 | 1.000000 | 0.000000 |
| 8947 | C19188 | 23.398673 | 0.833333 | 144.40 | 0.00 | 144.40 | 0.000000 | 0.833333 | 0.000000 |
| 8948 | C19189 | 13.457564 | 0.833333 | 0.00 | 0.00 | 0.00 | 36.558778 | 0.000000 | 0.000000 |
| 8949 | C19190 | 372.708075 | 0.666667 | 1093.25 | 1093.25 | 0.00 | 127.040008 | 0.666667 | 0.666667 |

8950 rows × 18 columns

| PURCHASES_INSTALLMENTS_FREQUENCY | CASH_ADVANCE_FREQUENCY | CASH_ADVANCE_TRX | PURCHASES_TRX | CREDIT_LIMIT | PAYMENTS | MINIMUM_PAYMENTS | PRC_FULL_PAYMENT | TENURE |
|---|---|---|---|---|---|---|---|---|
| 0.083333 | 0.000000 | 0 | 2 | 1000.0 | 201.802084 | 139.509787 | 0.000000 | 12 |
| 0.000000 | 0.250000 | 4 | 0 | 7000.0 | 4103.032597 | 1072.340217 | 0.222222 | 12 |
| 0.000000 | 0.000000 | 0 | 12 | 7500.0 | 622.066742 | 627.284787 | 0.000000 | 12 |
| 0.000000 | 0.083333 | 1 | 1 | 7500.0 | 0.000000 | NaN | 0.000000 | 12 |
| 0.000000 | 0.000000 | 0 | 1 | 1200.0 | 678.334763 | 244.791237 | 0.000000 | 12 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 0.833333 | 0.000000 | 0 | 6 | 1000.0 | 325.594462 | 48.886365 | 0.500000 | 6 |
| 0.833333 | 0.000000 | 0 | 6 | 1000.0 | 275.861322 | NaN | 0.000000 | 6 |
| 0.666667 | 0.000000 | 0 | 5 | 1000.0 | 81.270775 | 82.418369 | 0.250000 | 6 |
| 0.000000 | 0.166667 | 2 | 0 | 500.0 | 52.549959 | 55.755628 | 0.250000 | 6 |
| 0.000000 | 0.333333 | 2 | 23 | 1200.0 | 63.165404 | 88.288956 | 0.000000 | 6 |

**Viewing Datas from datasets**

## 6.1 SUMMARY

The red color represents the cluster number zero, and data regarding about maximum purchase rate. The green color represents cluster number one, which indicates frequent purchase data of customers. The purple color represents cluster number two, regarding data about negative purchase rates. The pink color represents the cluster number three, that represents data about neutral customers.

## 8. References

[1] Kayalvily Tabianan, Shubashini Velu, Vinayakumar Ravi "K-MEANS CLUSTERING APPROACH FOR INTELLIGENT CUSTOMER SEGMENTATION USING CUSTOMER PURCHASE BEHAVIOR DATA," IEEE, May 2022.

[2]Ranjan Pandey, Saikat Ruj "CUSTOMER SEGMENTATION USING K MEANS CLUSTERING," Towards Data Science, Apr. 2019.

[3] V. Vijilesh, "CUSTOMER SEGMENTATION USING MACHINE LEARNING," International Research. [4] Tushar Kansal, Suraj Bahuguna, Vishal Singh; Tanupriya Choudhury, "CUSTOMER SEGMENTATION USING K-MEANS CLUSTERING," IEEE, Jul. 2019.

[5] Shubashini "CUSTOMER SEGMENTATION USING K-MEANS CLUSTERING" International Conference on Computational Techniques, Electronics and Mechanical Systems (CTEMS), Jun. 2018.

 [6] Mamta Juneja; Parvinder Singh Sandhu "DESIGNING OF ROBUST IMAGE STEGANOGRAPHY TECHNIQUE BASED ON LSB". May 2021

[7] Bahuguna, Tanupriya, Vinayakumar Ravi "CUSTOMER SEGMENTATION OF RETAIL CHAIN CUSTOMERS USING CLUSTER ANALYSIS", credits stockholm, Sweden, May 2019

[8] Asif Iqbal, Rajeev Ranjan Pandey, Subhraneel Bagchi, Saikat Ruj, Sujata Dawn.June "CUSTOMER SEGMENTATION USING MACHINE LEARNING WITH A COUPON GENERATOR GUI", IEEE, June 2023

[9] Kilari, Sailesh Edara , Guna Ratna Sai Yarra , Dileep Varma Gadhiraju "CUSTOMER SEGMENTATION USING K-MEANS CLUSTERING HEMASHREE", IOP Conference, March-2022.