

# Optimizing e-commerce sales through advanced customer segmentation and recommendation techniques

Shivani Shedole  
*Master's in Data Analytics*  
*National College of Ireland*  
Dublin, Ireland  
x22218688@student.ncirl.ie

**Abstract**—In this project, an actual retail data set is used to build a customer segmentation and recommendation system. Its main purpose is to divide customers according to their buying habits and to provide them with tailored advice aimed at improving both satisfaction levels and sales. This process includes pre-processing the data, creating different variables from it, and applying machine learning methods like clustering algorithms (for example, K-Means), which create groups of distinct customers. In addition, the implementation of a suggestion mechanism based on collaborative filtering feeds relevant products to each segment's clients. The approach taken to segment these clients has been seen to be effective, as it generates insights to target marketing strategies, as well as personalized customer interactions.

**Index Terms**—Data Pre processing, Retail Analytics, K-Means Clustering, Customer Segmentation, Predictive Analytics

## I. INTRODUCTION

Understanding customer behavior is the key to developing effective marketing strategies and improving customer retention in a highly competitive retail industry. Customer segmentation helps business categorize their customers into distinct groups, hence allowing more personalized marketing effort. This Project Responsible makes use of an extensive retail dataset to perform customer segmentation and builds a recommendation system that can suggest products for customers based on their purchasing patterns.

The project's main aim is to identify significant customer segments, and subsequently provide product recommendations which fit in line with every segment's preferences. Using clustering algorithms such as K-Means, we want to group clients according to various characteristics of theirs, e.g. frequency of buying from them, amount spent on shopping, and the choice of product made by them. A recommendation engine built through these segmentation results will take advantage of collaborative filtering techniques so as to propose tailored product suggestions, thus driving customer engagements and hence increasing sales.

This project not only adds on the understanding about clients' behavior but also shows how data-based processes can be used for real life problem solutions in retail industry context. The findings could be utilized for target campaigns marketing for the customer base high on business optimize.

## II. RELATED WORK: ANALYSIS OF RELEVANT ACADEMIC WORK

There has been much academic work in the area of customer segmentation and recommendation systems aimed at achieving a greater understanding of customer behaviors in order to improve marketing strategies. The relevant literature is filled with various methodologies, including collaborative filtering, clustering algorithms, and machine learning approaches to improve the performance of e-commerce platforms. Some key academic contributions that relate to this project are analyzed below.

1. Using Clustering Techniques for Customer Segmentation  
Clustering algorithms such as K-means have been frequently utilized in customer segmentation tasks. The Pioneering work by Wedel and Kamakura (2000) in "Market Segmentation: Conceptual and Methodological Foundations" provides a comprehensive framework for applying clustering techniques to market segmentation. Their efforts focus on the importance of recognizing similar customer categories in order to personalize their marketing campaigns.

Moreover, more recent works exemplified by Chen et al (2012) "Customer segmentation using RFM model and K-means algorithm" demonstrate how Recency Frequency and Monetary (RFM) model can be combined with K-means clustering to differentiate good customers from bad ones. This study emphasizes how RFM enables one to come up with targeted marketing plans.

2. Collaborative Filtering and Recommendation Systems  
Recommendation systems have become an essential component of e-Commerce platforms, yielding personalized product suggestions. The primary groundwork put forth by Schafer,

Konstan, and Riedl (2001), “E-commerce Recommendation Applications,” investigates collaborative filtering techniques based on user behavior data to recommend products. Their findings indicate that through appropriate suggestions of related products, collaborative filtering can minimize customer discontent and consequently increase sales.

Another significant study is by Koren, Bell, and Volinsky (2009) in “Matrix Factorization Techniques for Recommender Systems” which deals with advanced matrix factorization methods for collaborative filtering. This study is especially relevant, as it touches on concerns about scalability and accuracy with regard to large-scale recommendation systems that are paramount in e-Commerce applications.

3. Integration of Customer Segmentation with Recommendation Systems In various studies, the integration of customer segmentation with recommendation systems has been explored. Ngai et al. (2009), in their paper ‘The Application of Data Mining Techniques in Customer Relationship Management: A Classification Framework and an Academic Review of Literature,’ provide a comprehensive review of data mining techniques used in customer relationship management (CRM). Their work highlights the importance of combining customer segmentation with recommendation systems to achieve a more personalized customer experience.

Furthermore, Li and Karahanna (2015) in their research “Online Recommendation Systems in B2C E-Commerce: A Review and Future Directions,” discuss the impact of integrating segmentation with recommendation systems. They argue that segmentation improves the relevance of recommendations by ensuring that they are tailored to the specific needs and preferences of distinct customer groups.

4. Machine learning has become a tool for improving customer segmentation and recommendation systems (CSRS). In their paper ‘A Machine Learning Framework for Customer Segmentation in E-Commerce’, Xia, Wang, and Zheng (2019) propose a framework for customer segmentation using machine learning that employs clustering algorithms and predictive modeling. Their findings show how advanced machine learning approaches can be useful in identifying customer segments more accurately.

In addition, Aggarwal (2016), in “Recommender Systems: The Textbook”, provides an overview of the various types of machine learning techniques used by contemporary recommender systems. This text is important for understanding how recommendation systems have developed from simple collaborative filtering to complex hybrid models integrating machine learning algorithms.

### A. Exploratory Data Analysis (EDA)

The Exploratory Data Analysis (EDA) process in this project, focused on customer segmentation and recommendation systems, was quite comprehensive. Here is a detailed breakdown of the EDA process based on the information provided.

#### 1. Exploring data sets initially:

df.head(10)								
	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	536365	85123A	WHITE HANGING HEART T-LIGHT HOLDER	6	12/01/10 8:26	2.55	17850.0	United Kingdom
1	536365	71053	WHITE METAL LANTERN	6	12/01/10 8:26	3.39	17850.0	United Kingdom
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8	12/01/10 8:26	2.75	17850.0	United Kingdom
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	6	12/01/10 8:26	3.39	17850.0	United Kingdom
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	6	12/01/10 8:26	3.39	17850.0	United Kingdom
5	536365	22752	SET 7 BABUSHKA NESTING BOXES	2	12/01/10 8:26	7.65	17850.0	United Kingdom
6	536365	21730	GLASS STAR FROSTED T-LIGHT HOLDER	6	12/01/10 8:26	4.25	17850.0	United Kingdom
7	536366	22633	HAND WARMER UNION JACK	6	12/01/10 8:28	1.85	17850.0	United Kingdom
8	536366	22632	HAND WARMER RED POLKA DOT	6	12/01/10 8:28	1.85	17850.0	United Kingdom
9	536367	84879	ASSORTED COLOUR BIRD ORNAMENT	32	12/01/10 8:34	1.69	130470	United Kingdom

Fig. 1. Data Summary

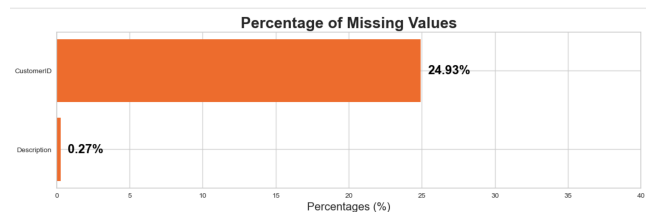


Fig. 2. Missing Vale

Summary of the data set: First off was loading and examining the dataset consisting of transactional data from a British retailer. The dataset had some columns like InvoiceNo, StockCode, Description, Quantity, InvoiceDate, UnitPrice, CustomerID and Country. The data set has 541,909 rows and 8 columns with different data types such as integers, floating objects, or objects. Summary Statistics: Initial statistical summaries were generated for numerical columns that indicate issues such as negative quantities (potentially returns) and also indicating the presence of missing values. 2. Handling Missing Values:

Let me mention that in the case of data preprocessing, particularly in sections explaining customer behavior or features related to products, specific measures were taken toward missing values. First, there were problems with crucial columns as CustomerID, which is necessary to distinguish separate buyers and analyze their activity. Since, CustomerID is essential for proper customer segmentation for clustering and for the analysis of their behavior, the rows with missing CustomerID were deleted. This action although as a result led to the deletion of roughly twenty five percent of the entire set, was important to ensure that customer based analyses were not skewed. In addition, there were missing values in the description feature, which is very significant in identifying the types of products that are being bought; the percentage of missing values was only 0.27. Due to the fact that there could be some discrepancies between the research findings that involve products and other analyzes, the remaining rows with missing product descriptions were also excluded. Even this small measure used was relevant to guaranteeing the quality of the analysis carried out on the whole.

#### 3. Handling duplicates:

To increase the quality of the dataset and the lack of homogeneity, a common issue is what method should be used to erase twin values in the data set. Duple represented

```

#Displaying the number of duplicate rows
print(f"The dataset contains {df.duplicated().sum()} duplicate rows that need to be removed.")

# Removing duplicate rows
df.drop_duplicates(inplace=True)

The dataset contains 5225 duplicate rows that need to be removed.

```

Fig. 3. Handling Duplicates

records may inflate sales figure, give wrong impressions about customer behavior, or result in wrong clusters which lead to wrong conclusions, therefore offering the right business decisions. In order to do this, the check was done in a rowwise manner such that a complete check was performed on the columns of the data set. This approach ensured that the matching rows were selected solely, and only the rows that were similar in every aspect. It was done so as to remove an entry similar in some way to a row, that could contain a different or crucial value in one or more fields. Having seen the duplicity of records, one is able to notice that the provided dataset has 5,225 records. Among these, some were dismissed in a bid to erase any likelihood of bias in the analysis that ensued; The data set was purified prior to the sales forecast and customer segmentation and behavior imitation with the purpose of removing such duplicates that could skew the results of the computations. The removal of replicated records meant that the quality of the dataset was increased and, hence, the ability to derive superior insights and make sound decisions was made possible.

#### 4. Treating canceled transactions

To this end, in order to identify the canceled transactions and the kind of customer behaviour for example where the customer is unhappy or no decision has been made, it was crucial to filter out the canceled transactions. These transactions were identified using the 'InvoiceNo' field where if the first character of the number was 'C' it signified a cancellation. This means that all cancellations were preserved, while the appropriate tag was set in front of each of them to indicate their status; this allowed for the exclusion of identified records, but it was not done under no conditions, kinds which could pose to the loss of important information. This approach was important because cancellation data reveal about customer dissatisfaction or reluctance or any other behavioural traits relevant to clustering/segmentation study. Where such transactions are effected, it will enable businesses to discover why certain clients are always canceling their orders, an aspect that will form the basis for solving such problems specifically as opposed to offering generic solutions such as improved client relations, change in product mix, or altering the marketing strategy. However, managing canceled transactions to enhance the recycling advantage impacts both the data and the corresponding customer interactions, leading to improved decision-making outcomes.

#### 5. Identification of Outliers:

To ensure accurate and reliable analysis, addressing outliers was a key focus. Outliers, which are data points significantly different from others, can distort analyses, leading to

	Quantity	UnitPrice
<b>count</b>	8872.000000	8872.000000
<b>mean</b>	-30.774910	18.899512
<b>std</b>	1172.249902	445.190864
<b>min</b>	-80995.000000	0.010000
<b>25%</b>	-6.000000	1.450000
<b>50%</b>	-2.000000	2.950000
<b>75%</b>	-1.000000	4.950000
<b>max</b>	-1.000000	38970.000000

Fig. 4. Treating Cancelled Transactions

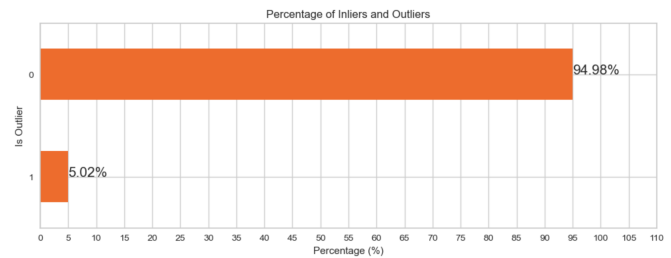


Fig. 5. Outliers

misleading conclusions and inaccuracies in sales forecasts, customer segmentation, and other critical evaluations. Outliers were identified by analyzing key variables such as 'Quantity' and 'UnitPrice' using statistical methods like the interquartile range (IQR) and z-scores to detect values that fell far outside the normal range. This approach considered both extreme high and low values. Instead of simply removing outliers that could result in the loss of meaningful data, such as legitimate but rare events such as bulk purchases, each outlier was carefully reviewed. Data points identified as likely errors (e.g., negative quantities where not applicable) were corrected or removed, while legitimate outliers were retained with careful consideration of their impact on the analysis. In addition, transformations, such as logarithmic transformations, were applied where appropriate to minimize the impact of these outliers on statistical analyzes. This balanced approach preserved the integrity of the dataset, ensuring that the analysis remained accurate and reflective of genuine trends, providing a realistic understanding of customer behavior, sales patterns, and other key metrics without the distortions that outliers can introduce.

#### 6. Cleaning Description Column

Since there is much more variation in the entries in the

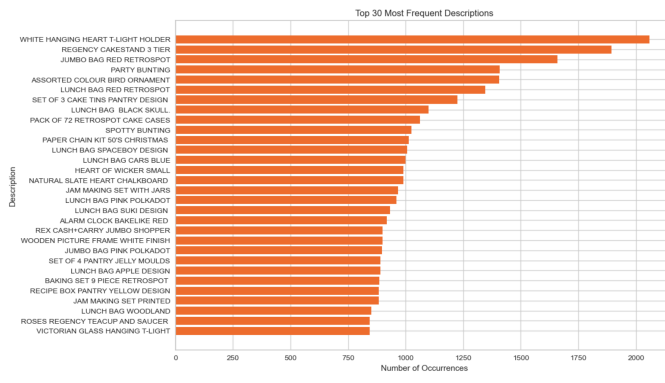


Fig. 6. Cleaning Description Column

Description column than in the other numerical columns, few major operations were done to pre-process the data set in the manner that was done in the Cleaning Step 2. This column is very important when it comes to identification of products being purchased and if analysis is to be made on them. Editing the data suppressed this finding: some entries referring to certain products were 'Next-Day Carriage' or 'High Resolution Image', meaning they were services or extra charges that could be viewed as noise; others were important only from the point of view of transactional data. For this reason, the authors chose a systematic approach to such non-product descriptions to ensure that the final data set was purely product. Additionally, and for the purpose of enhancing fairness; all words in all the product descriptions were capitalized except for the first word, which was already altered. Standardization is beneficial where perhaps you do not desire issues arising out of case differences, for instance, between 'Gift Box' and 'GIFT BOX'. All the descriptions are made in uppercase and this also enhances the standardization of field and can also remove other possible sources of confusion between different types of products. This logical and methodology approach of cleaning and standarization ensure that the data set is both meaningful and relevant, excluding every incidence that may be irrelevant in the process of data analysis or modeling and absolutely focus on only real products delivered, thus improving the quality of the data analysis or modeling process.

### III. METHODOLOGY

Clustering is a fundamental technique in unsupervised machine learning, used to group a set of objects in such a way that objects in the same group (or cluster) are more similar to each other than to those in other groups. Various clustering methods exist, each with its own strengths, assumptions, and applications. Here is an overview of some common clustering types:

#### A. Partitioning Clustering(K-Mean)

Partitioning methods are a class of clustering algorithms with the goal of developing the clustering of resources in such a way that there are several clusters of the resources and each resource belongs to a single cluster. Clustering is

done in such a way that objects in any given cluster contain similarity compared to objects in other clusters and a method of determining the number of clusters. These methods are preferred because they prove to be efficient in many ways and not so complicated to perform.

Among all the partitioning methods, the most common is the K-Means clustering. The data is then divided into K clusters in case of K-means clustering where the number of clusters K is specified in advance. The algorithm in simple terms involves continually assigning each of the data points to the cluster whose centroid – a statistic that can be thought of as the average position of all the points in that cluster – is nearest to the point in question. The centroids are then calculated and the process of classifications is further carried out in order to get new values of the centroids for each of the classes and the process continues until the centroids no longer change much and this will be a hint that the clusters have settled. K-means are fast and simple and thus are a frequently used approach to solve many clustering problems. However, k-means decided on the number of clusters a priori, which may be an inconvenience if the number of clusters has not been determined. Also, the assumption of equal sizes and spherical forms of the clusters also poses a principle limitation to the use of K-means.

Another notable partitioning method is K-Medoids Clustering, which is similar to K-means but with a key difference: As opposed to centroids, K-medoids employs medoids as the cluster representatives – actual data samples, belonging to a given dataset. This makes K-medoids more resistant to noise since; while calculating the centroid weighted of all the points in the cluster it is easily affected by the extreme values where as medoids are not. On the downside, because K-medoids have to find the medoids, which involves the calculation of distances between each pair of points, this is computationally more costly than K-means, particularly for large databases.

Why K- Mean CLustering ? K-Means clustering is a widely used algorithm known for its simplicity and efficiency. It partitions data into a predetermined number of clusters (K) by iteratively adjusting centroids to minimize the distance between data points and their assigned centroid. This method scales well with large datasets, making it ideal for practical applications like market segmentation and image compression. K-Means is particularly effective when data forms spherical clusters of similar size. However, it requires the number of clusters to be specified in advance and is sensitive to outliers, which can distort the clustering results. Despite these limitations, K-Means remains a popular choice due to its straightforward implementation and versatility.

Example:-K-Medoids clustering like K-Means is also used for clustering of data into clusters but the measure 'center' of a cluster is substituted by a called medoid which is a data point. This makes K-Medoids more resistance to noise and outliers; it will now be easier to determine an object that is different from the group mean. Here are some examples of where K-Medoids clustering can be effectively applied:Here are some examples of where K-Medoids clustering can be effectively applied:



1. Customer Segmentation in Retail: To cluster them based on the purchasing behaviour, retailers can employ K-Medoids. While using the K-Means method, it is susceptible to the external influences such as high and low spending customers (outliers), but K-Medoids will give the cluster median or the most representative customers (medoids).

2. Bioinformatics: For clustering the genes which have similar expression patterns, K-Medoids can be used in gene expression data analysis. Due to the noisy and large fluctuation in biological data sets, K-Medoids provide a better clustering as it works on actual gene data points not on centroids.

3. Document Clustering: In case of textual data like clustering of the articles, or research papers, K-Medoids can be used to partition documents according to the resemblance of their content. In text clustering, K-Medoids is useful since text data often contains many outliers such as very short or very long documents and here, the medoid of each cluster is always a real document and hence makes more sense.

These examples illustrate real benefits of K-Medoids especially where outliers are to be found and where it may be preferred to have the actual data points as initial approximations of the means.

1) *Elbow Method*: Hence, the Elbow Method is one of the common techniques through which the right number of clusters to be formed when conducting the algorithm is established. The method consists in employing the decision-making algorithm for more iterations with various amounts of clusters (K) and calculating the WCSS for each K. The WCSS is equal to the sum of squared distance vectors of every summative point to the centroid of the assigned cluster, and the lower the WCSS, the more compact the clusters are. Plotting the number of clusters with WCSS gives a graph whose features normally have the following phenomena: there is an acute reduction in WCSS when the number of clusters rises in a sharp manner but then a moment comes where the reduction in WCSS trends to be sluggish. It is the point that is typically like 'elbow' in the graph and represents the right number of clusters. The rationale, however, is that incorporating additional clusters at some point does not change the model creditably much, because the WCSS reduction ceases. Hence, the "elbow" point is the tradeoff between a minimal WCSS and non overfitting by the use of many clusters.

Advantages :- 1. Simplicity and Intuitiveness: The Elbow Method is straightforward to recognize and put in force. It involves plotting the within-cluster sum of squares (WCSS) in opposition to the number of clusters and seeking out an "elbow" point where the price of lower sharply slows down. This point usually shows a very good stability between cluster compactness and range of clusters.

2. Helps Avoid Overfitting: By identifying the most useful range of clusters, the Elbow Method enables save you overfitting. Too many clusters can lead to overfitting, wherein the model captures noise in place of meaningful styles. The Elbow Method helps discover a affordable number of clusters that generalizes well to unseen facts.

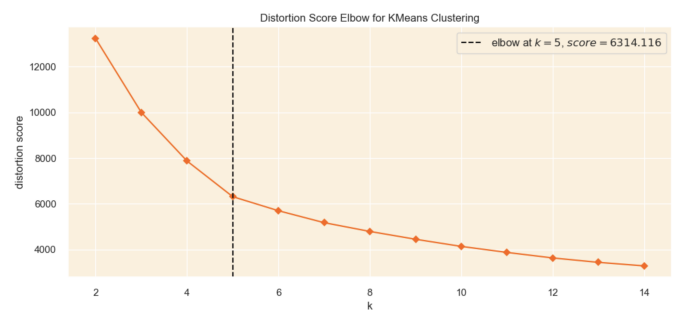


Fig. 7. Elbow Method

3. Improves Interpretability: Selecting the precise variety of clusters using the Elbow Method can result in more significant and interpretable clusters. When the wide variety of clusters is chosen carefully, the ensuing clusters are probable to symbolize true underlying patterns inside the data, making the evaluation extra useful.

4. Effective for Various Data Types: The Elbow Method may be implemented to exceptional varieties of facts and clustering algorithms. Whether you're the usage of K-approach, hierarchical clustering, or different strategies, the Elbow Method provides a honest manner to assess the quantity of clusters.

5. Visualization: The Elbow Method affords a visual representation of the WCSS or some other metric towards the variety of clusters. This visual resource helps in making informed decisions by means of truly showing where adding greater clusters yields diminishing returns.

The figure shows the distortion score, also known as the within-cluster sum of squares, on the Y-axis using the Elbow Method, a method for determining the optimal number of clusters (k) for a K-Means clustering plot. WCSS), which measures the smallness of the cluster As the number of clusters (k) increases, the distortion generally decreases because data points are divided into more clusters, resulting in significantly smaller clusters. The main point of this plot is "corner", where decreasing distribution decreases sharply, implying diminishing returns to a higher number of clusters. In this case k = 5 holds ahead, implying that five clusters are optimal. Beyond this point, adding more clusters does not significantly improve the accuracy of the clusters, as reflected by the flat circle and thus the fingerprint method helps to identify them the right number of clusters by determining where the benefits of adding more clusters begin to disappear.

2) *Silhouette method*?: The Silhouette Method is some other method used to determine the optimum range of clusters in a dataset, frequently used along side clustering algorithms like K-Means. This method evaluates the pleasant of clustering by measuring how comparable every statistics point is to its personal cluster as compared to other clusters.

The Silhouette value for each point is calculated using the formula:

$$\text{Silhouette Score} = \frac{b-a}{\max(a,b)}$$

wherein: - (a) is the average distance from the factor to all

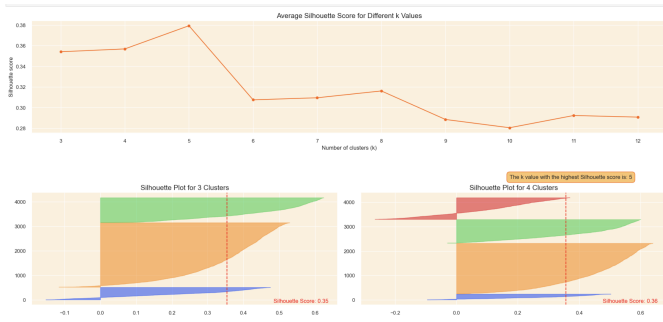


Fig. 8. Silhouette Method

other factors within the equal cluster (intra-cluster distance). - (b) is the common distance from the factor to all points within the nearest neighboring cluster (inter-cluster distance).

The Silhouette Score degrees from -1 to one: - A rating near 1 indicates that the statistics factor is well clustered, as it's far a good deal closer to factors in its personal cluster than to factors in other clusters. - A rating around zero shows that the point is on or very close to the selection boundary among two clusters. - A score close to -1 indicates that the point might be assigned to the wrong cluster, as it is in the direction of points in a neighboring cluster than to points in its very own cluster.

The photograph gives a comprehensive analysis of clustering effectiveness the use of the Silhouette Score, a metric that evaluates how well records factors healthy inside their assigned clusters relative to others. The top plot suggests the average Silhouette Score as the number of clusters ( $k$ ) varies from three to twelve, with the rating peaking at around 0.38 for  $k =$  five, indicating that 5 clusters can be the maximum appropriate desire for this dataset. This top suggests an top of the line stability between concord within clusters and separation among them. The backside plots illustrate the distribution of Silhouette Scores for  $k = 3$  and okay = four, displaying that even as both configurations provide a few degree of clustering satisfactory, they nevertheless showcase overlap between clusters, specifically for okay = 3. The moderate development at ok = four, with a median score of zero.36, nonetheless would not surpass the readability located at ok = 5. This analysis indicates that ok = 5 gives the most wonderful and meaningful clustering, making sure that the agencies fashioned appropriately constitute the underlying statistics structure without needless complexity or over-segmentation.

To follow the Silhouette Method, you calculate the Silhouette Score for each records factor after which common these scores to get an ordinary Silhouette Score for each price of  $K$  (the wide variety of clusters). By plotting those common Silhouette Scores in opposition to the quantity of clusters, you may perceive the optimal variety of clusters as the one with the best common Silhouette Score, indicating that clusters are well-separated and accurately assigned. The Silhouette Method is particularly beneficial because it now not

most effective considers how well the clusters are fashioned but also how wonderful they're from every different, supplying a extra nuanced asses

3) *MapReduce*: MapReduce is a programming model and processing approach mainly designed to address and technique massive datasets in a dispensed computing surroundings. The system entails important stages: Map and Reduce. In the Map phase, input records is damaged down into smaller key-fee pairs, which might be then processed independently and in parallel across a couple of nodes in a cluster. This phase produces intermediate key-cost pairs. Next, during the Shuffle and Sort phase, those intermediate pairs are grouped by means of key and organized to put together for the Reduce phase. In the Reduce phase, the grouped facts is aggregated, filtered, or in any other case processed to produce a final output of key-cost pairs. This model is noticeably efficient for big-scale facts processing tasks as it leverages parallelism, allowing responsibilities to be dispensed across more than one machines, thereby enhancing overall performance and scalability. Additionally, MapReduce is fault-tolerant; if a node fails, tasks may be re-executed on another node, ensuring reliability. This framework is broadly used in various programs, together with search engines, records mining, log processing, and advice systems, making it a cornerstone of huge information processing in allotted environments.

#### IV. DIMENSIONALITY REDUCTION

Dimensionality reduction offers numerous benefits in facts evaluation, especially when coping with complicated datasets containing multicollinear features. By decreasing the dimensionality, we can dispose of redundant records and address multicollinearity, thereby improving the quality of our models. For clustering algorithms like K-manner, that are touchy to the quantity of functions, dimensionality reduction enables in figuring out more wonderful and well-separated clusters by focusing on the most sizeable variables. Additionally, this process can lessen noise inside the statistics, main to greater accurate and strong clusters. Beyond these technical blessings, dimensionality reduction enhances the interpretability of the results with the aid of permitting us to visualize customer segments in two or three dimensions, providing clearer insights. Finally, it additionally improves computational performance through reducing the complexity of the dataset, dashing up the modeling process and making it more green.

The picture provides a visualization of defined variance and cumulative variance as a part of a Principal analysis (PCA), a technique commonly used to reduce the dimensionality of datasets. It features a bar chart that shows the defined variance for each important thing, showing how plenty facts each thing captures. The first factor captures the very best amount of variance, around 58 , with next additives capturing step by step less. Alongside this, a line plot illustrates the cumulative explained variance, which accumulates the variance captured via every aspect as more are added. For example, by using consisting of the second factor, approximately 78 of the overall variance is captured, and with the aid of the 0.33, ninety

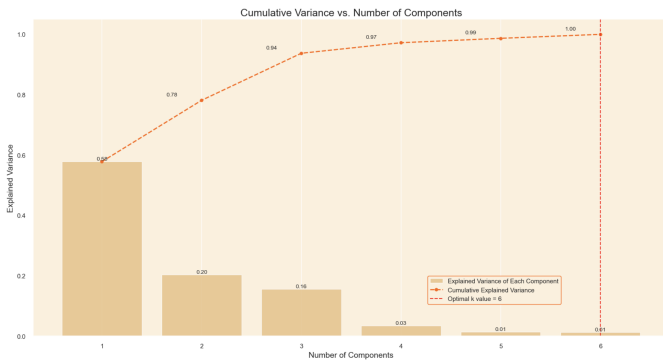


Fig. 9. PCA Explained and Cumulative Variance Plot

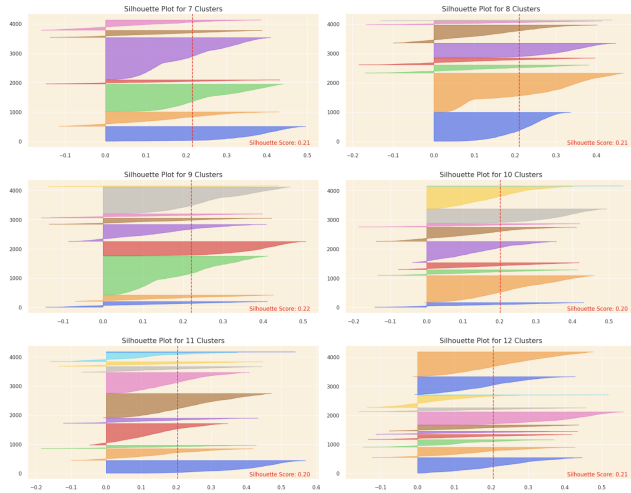


Fig. 10. Silhouette plots

four is captured. The plot additionally consists of a vertical line marking the gold standard variety of components (okay), which is about 6 in this situation, indicating that these six additives collectively capture almost all of the variance inside the statistics. Text annotations similarly make clear the exact variance values, and the legend explains the factors of the plot. This visualization is important for figuring out the gold standard quantity of components to maintain in PCA, making sure that dimensionality is decreased with out sizable lack of information, for that reason it helps in greater effective records analysis and version construction.

The photo provides a chain of Silhouette plots evaluating the best of clustering for extraordinary numbers of clusters (okay) starting from 7 to twelve. These plots are instrumental in assessing how nicely records factors are grouped within their clusters, with the Silhouette Score indicating the diploma of separation between clusters. In every subplot, clusters are represented through coloured regions, wherein the width of every location correlates with the number of records points, and the position on the X-axis displays the Silhouette Score. The plots display usually low Silhouette Scores, ranging between 0.20 and 0.22, suggesting bad clustering pleasant

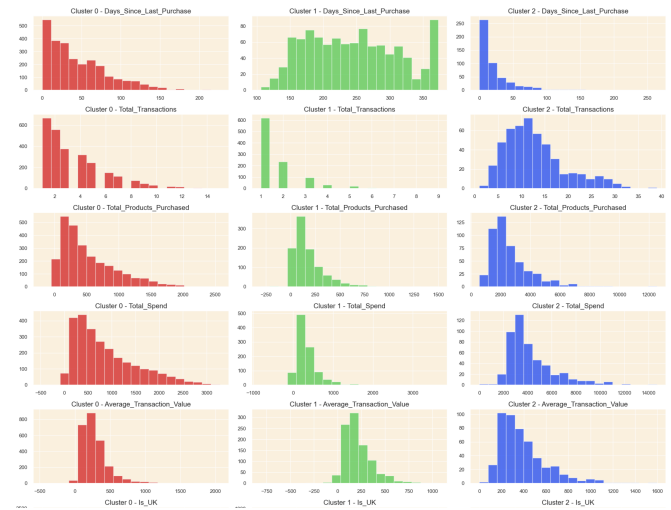


Fig. 11. Overall Purchase

throughout all configurations. The presence of negative ratings and overlapping clusters suggests that many information points are not properly-labeled within their assigned clusters, implying that they might belong to extraordinary clusters. Moreover, increasing the range of clusters past a sure threshold does now not seem to enhance the clustering satisfactory, rather main to over-segmentation, in which clusters grow to be less wonderful and significant. Overall, the evaluation underscores the importance of choosing the suitable quantity of clusters that balances the clarity and separation of organizations, in preference to merely growing okay, which can result in ambiguous and less useful clustering outcomes.

## REFERENCES

- [1] MacQueen, J. : Some methods for classification and analysis of multi-variate observations. In Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability. 1, no. 14, 281-297. Univ. of California Press, Berkeley and Los Angeles, 1967.
- [2] Lloyd, Stuart P. "Least squares quantization in PCM. " IEEE Transactions on Information Theory 28, no. 2 (1982): This association was evident in the study among almost all the subgroups of the research communities with slight variations in the number; a number ranging from 129 to 137.
- [3] Hartigan, John A. and Manchek A. Wong , "Algorithm AS 136: A k-means clustering algorithm. " Journal of the Royal Statistical Society: Series C , Applied Statistics, volume 28 , number 1, 1979;; 100-108.
- [4] Kodinariya, M. Trupti, Makwana, P. R. "Review on determining number of cluster in K-means clustering. " International Journal of Advance Research in Computer Science and Management Studies 1, no. 6 (2013): 90 – 95..
- [5] Kodinariya, Trupti M. and Prashant R. Makwana. "Review on determining number of cluster in K-means clustering. " International Journal of Advance Research in Computer Science and Management Studies 1, no. 6 (2013): 90-95.
- [6] Rousseeuw, Peter J. "Silhouettes: "A graphical aid to the interpretation and validation of cluster analysis". Journal of Computational and Applied Mathematics 20, (1987), 53-65.
- [7] Ketchen, David J. and Shook Christopher L. "The application of cluster analysis in strategic management research: argumentation of the article, which is an analysis and critique. Strategic Management Journal 17, no. 6 (1996): 441–458.