



**JOMO KENYATTA UNIVERSITY OF AGRICULTURE AND  
TECHNOLOGY**

**SCHOOL OF COMPUTING AND INFORMATION  
TECHNOLOGY**

**BSC. INFORMATION TECHNOLOGY**

**PROJECT TITLE: CUSTOMER SEGMENTATION IN E-  
COMMERCE**

**STUDENT NAME: TEDMACK MUTUMA KIRIMI**

**REGISTRATION NUMBER: SCT221-C004-0305/2020**

**SUPERVISOR: MS. JUDY GATERI**

This project has been submitted in partial fulfillment of the requirements for the award of the degree of Bachelor of Science in Information Technology in the year 2023.

**DECLARATION**

I hereby declare that this project report is based on my original work apart from the citations and quotations included which have been duly acknowledged. The project has also not been previously and concurrently submitted for any other degree or award at Jomo Kenyatta of Agriculture and Technology.

TEDMACK MUTUMA KIRIMI

REG NO: SCT221-C004-0305/2020

Signature.....Date .....

**SUPERVISOR**

I the undersigned do hereby certify that this is a report for the project undertaken by the above-named student under my supervision and that it has been submitted to Jomo Kenyatta University of Agriculture and Technology with my approval.

Ms. JUDY GATERI

LECTURER - DEPARTMENT OF INFORMATION TECHNOLOGY

.

Signature ..... Date .....

## **ACKNOWLEDGEMENT**

I am profoundly grateful to The Almighty God, whose unwavering support, strength, and wisdom have been the driving force behind the successful completion of this project. Without His divine intervention, this endeavor would not have been possible, and I humbly acknowledge His grace that has guided me throughout.

My heartfelt appreciation extends to my family, the bedrock of my support system. Their encouragement, both emotionally and financially, has been pivotal in sustaining me on this challenging journey. Their unwavering belief in my abilities has provided me with the motivation to persevere until the very end. I am truly blessed to have such a dedicated and loving family by my side.

A special note of gratitude goes to my esteemed supervisor, Ms. Judy Gateri. Her invaluable guidance, understanding, and sagacious direction have played a pivotal role in shaping the trajectory of this project. I am thankful for her mentorship, which has enriched my knowledge and skills, contributing significantly to the project's success.

In addition, I extend my thanks to my classmates, a remarkable cohort whose collective support has been instrumental. Their collaborative spirit, willingness to share insights, and assistance on various aspects of the project have been invaluable. In the spirit of camaraderie, they have created an environment conducive to growth and learning, for which I am sincerely thankful.

## **DEDICATION**

This thesis is devoted to my family's unfailing love and belief in my skills, which has been my biggest motivator throughout this academic journey. To my parents, who instilled in me a love of study and the fortitude to overcome obstacles; and to my siblings, who provided consistent encouragement and understanding throughout the difficult phases of this research.

This work is also dedicated to my mentors and advisors, whose advice and knowledge have molded my intellectual evolution. Their excellent thoughts and constructive feedback were critical in the development of this study. Finally, I'd want to thank the many people who have helped me comprehend the issue by their contributions and encouragement.

## Table of Contents

DECLARATION.....	ii
ACKNOWLEDGEMENT.....	iii
DEDICATION .....	iv
ABBREVIATIONS AND ACRONYMS .....	x
ABSTRACT.....	xi
CHAPTER ONE .....	1
INTRODUCTION.....	1
1.1 BACKGROUND OF THE STUDY .....	1
1.2 PROBLEM STATEMENT .....	2
1.3 OBJECTIVES .....	3
1.3.1 GENERAL OBJECTIVES .....	3
1.3.2 SPECIFIC OBJECTIVES.....	3
1.4 RESEARCH QUESTIONS.....	3
1.5 JUSTIFICATION.....	3
1.6 SCOPE AND LIMITATIONS .....	4
CHAPTER TWO .....	5
LITERATURE REVIEW .....	5
2.0 INTRODUCTION.....	5
2.1 HISTORY OF CUSTOMER SEGMENTATION .....	5
2.1.1 Early application in Traditional Retail (19 <sup>th</sup> century) .....	5
2.1.2 Phase 2: Mass marketing (Early 20 <sup>th</sup> century).....	5
2.1.3 Phase 3: Post-war segmentation (Mid-20 <sup>th</sup> century) .....	6
2.1.4 Phase 4: Market Research Era (Mid to Late 20th Century).....	7
2.1.5 Phase 5: Database and Technology Revolution Era (Late 20th Century) .....	7
2.1.6 Phase 6: Digitalization and Online Era .....	8
2.2 Theoretical Framework .....	8
2.2.1 Supervised learning approach .....	8
2.2.1.1 Decision Trees .....	8
2.2.1.2 Random Forests .....	9
2.2.1.3 Support Vector machines.....	10
2.2.1.4 Logistic Regression.....	11
2.2.1.5 Neural Networks (Deep learning) .....	13
2.2.1.6 K-Nearest Neighbors (K-NN) .....	17
2.2.1.7 Gradient Boosting Machines.....	17
2.2.2 Unsupervised machine learning .....	18

2.2.2.1 K-Means Clustering .....	18
2.2.2.2 Hierarchical Clustering .....	19
2.2.2.3 DBSCAN .....	19
2.2.2.4 Principal Component Analysis (PCA) for Dimensionality Reduction .....	20
2.2.2.5 Latent Dirichlet Allocation (LDA) for Topic Modeling .....	21
2.2.2.6 Gaussian Mixture Model (GMM) .....	21
2.2.2.7 Self-Organizing Maps (SOM) for Customer Segmentation .....	22
2.2.2.8 Associate rule mining .....	23
2.3 TECHNIQUES IN CUSTOMER SEGMENTATION .....	24
2.4 Datasets and Data Sources in Customer segmentation .....	26
2.4.1 Available datasets for Customer Segmentation .....	27
2.5 Case Studies on Customer Segmentation in e-commerce using Machine learning .....	29
2.6 Synthesis and Analysis .....	31
2.7 Gaps in the Literature .....	33
2.7.1 Data Quality and Quantity .....	33
2.7.2 Algorithm Selection and Validation .....	33
2.7.3 Interpretable Models .....	33
2.7.4 Privacy and Ethical Concerns .....	33
2.7.5 Dynamic Segmentation .....	34
2.7.6 Scalability and Resource Constraints .....	34
2.7.7 Evaluation Metrics .....	34
2.8 FUTURE TRENDS IN CUSTOMER SEGMENTATION .....	34
2.9 CONCLUSION .....	36
CHAPTER 3 .....	38
RESEARCH METHODOLOGY .....	38
3.0 INTRODUCTION .....	38
3.1 RESEARCH DESIGN .....	38
3.1.1 Description of research design .....	38
3.1.2 Justification .....	39
3.2 SYSTEM DEVELOPMENT METHODOLOGY .....	39
3.2.1 PHASES OF AGILE DEVELOPMENT .....	39
3.2.2 AGILE DEVELOPMENT PRONCIPLES .....	40
3.3 PARTICIPANTS OR SAMPLE .....	41
3.3.1 SAMPLING METHODOLOGY .....	41
3.3.2 SAMPLE CHARAECTERISTICS .....	41
3.4 MODEL DEVELOPMENT .....	41

3.4.1 DATA COLLECTION .....	41
3.4.2 DATA PREPROCESSING .....	42
3.4.3 TRAINING THE MODEL .....	42
3.4.4 TESTING THE MODEL .....	43
3.5 ETHICAL CONSIDERATIONS .....	44
3.5.1 ETHICAL FRAMEWORK .....	44
3.6 RESEARCH QUALITY ASSURANCE .....	45
3.7 CONCLUSION .....	45
CHAPTER 4.....	47
SYSTEM ANALYSIS.....	47
4.0 INTRODUCTION .....	47
4.1 DATA COLLECTION METHODS .....	47
4.1.1 PREVIOUS TRANSACTION ANALYSIS.....	47
4.2 GATHERING USER REQUIREMENTS.....	48
4.2.1 STAKEHOLDERS AND THEIR NEEDS .....	48
4.3 FEASIBILITY STUDY .....	49
4.3.1 TECHNICAL FEASIBILITY.....	49
4.3.2 OPERATIONAL FEASIBILITY .....	49
4.3.3 ECONOMIC FEASIBILITY.....	49
4.4 DETERMINING PROJECT VIABILITY .....	49
4.5 SOFTWARE REQUIREMENT SPECIFICATION .....	50
4.5.1 FUNCTIONAL REQUIREMENTS .....	50
4.5.2 NON-FUNCTIONAL REQUIREMENTS .....	51
4.6 CONCLUSION .....	52
CHAPTER 5.....	53
SYSTEM DESIGN .....	53
5.0 INTRODUCTION .....	53
5.1 SYSTEM ACHITECTURE .....	54
5.2 SYSTEM FLOW CHART .....	54
5.3 SYSTEM WORKFLOW .....	55
5.4 K-CLUSTERING WORKFLOW IN CUSTOMER SEGMENTATION.....	55
5.5 SEQUENCE DIAGRAM .....	56
5.6 CONCLUSION .....	57
CHAPTER 6.....	58
SYSTEM IMPLEMENTATION AND TESTING .....	58
6.0 INTRODUCTION .....	58

6.1 TOOLS FOR IMPLEMENTATION.....	58
6.1.1 PROGRAMMING LANGUAGE.....	58
6.1.2 PACKAGES.....	59
6.2 TESTING.....	60
6.2.1 UNIT TESTING.....	60
6.2.2 INTEGRATION TESTING.....	61
6.2.3 SYSTEM TESTING .....	61
6.2.4 User Acceptance Testing (UAT) .....	61
6.3 SCREENSHOTS OF THE SYSTEM .....	61
6.4 CONCLUSION .....	62
CHAPTER 7.....	64
CONCLUSION AND RECCOMENDATION .....	64
7.1 CONCLUSION .....	64
7.2 FUTURE RECCOMENDATIONS .....	64
REFERENCES .....	66
APPENDICES .....	69
APPENDICE A: SAMPLE CODE FOR CLUSTER PREDICTION .....	69



## TABLE OF FIGURES

Figure 1. Desicion Tree.....	9
Figure 2. Random Forests.....	10
Figure 3. Support Vector Machines .....	11
Figure 4. Logistic Regression .....	12
Figure 5. Neural Networks.....	13
Figure 6. Imaged based Customer Segmentation .....	14
Figure 7. Dimensionality Reduction.....	15
Figure 8. GAN Image Generation.....	16
Figure 9. K.N Nearest .....	17
Figure 10. Gradient Boosting Machine.....	18
Figure 11. Unsupervised Learning.....	19
Figure 12. DBSCAN.....	20
Figure 13. Latent Dirichlet Allocation (LDA) .....	21
Figure 14. Gaussian Mixture Model (GMM) .....	22
Figure 15. Self-Organizing Maps (SOM) .....	23
Figure 16. Associate rule mining .....	24
Figure 17. Variables Table.....	27
Figure 18. Variables Table 2.....	28
Figure 19. Variables Table Accuracy .....	28
Figure 20. Variables Table Precision.....	29
Figure 21. K- Mean Architecture .....	31
Figure 22. Synthesis and Analysis .....	31
Figure 23. Random Forest learning curve.....	32
Figure 24. Training Examples.....	32
Figure 25. Training Result .....	33
Figure 26. Agile Methodology.....	39
Figure 27. System Architecture.....	54
Figure 28. System Flowchart .....	54
Figure 29. System Workflow .....	55
Figure 30. K-Clustering Workflow .....	55
Figure 31. Sequece Diagram.....	56
Figure 32. Screenshot 1.....	61
Figure 33. Screenshot 2.....	62
Figure 34. Screenshot 3.....	62

## **ABBREVIATIONS AND ACRONYMS**

**E-commerce:** Electronic Commerce

**RFM:** Recency, Frequency, Monetary

**CLV:** Customer Lifetime Value

**AI:** Artificial Intelligence

**ML:** Machine Learning

**DFD:** Data Flow Diagram

**GPU:** Graphics Processing Unit

**IDE:** Integrated Development Environment

**ROI:** Return on Investment

**SVM:** Support Vector Machines

**COHN-KANADE:** Cohn-Kanade Database

**LBP:** Local Binary Patterns

**CDF:** Cumulative Distribution Function

**PCA:** Principal Component Analysis

**KPIs:** Key Performance Indicators

**API:** Application Programming Interface

**CRM:** Customer Relationship Management

**B2B:** Business-to-Business

**B2C:** Business-to-Consumer

**ROI:** Return on Investment

**CSV:** Comma-Separated Values

**GDPR:** General Data Protection Regulation

**UX:** User Experience

**A/B Testing:** Split Testing or A/B Split Testing

**IoT:** Internet of Things

## **ABSTRACT**

Understanding client behavior is critical in the dynamic landscape of e-commerce for targeted marketing tactics and increased user experiences. This study dives into customer segmentation, using advanced analytics approaches to categorize consumers based on their preferences, purchasing habits, and interactions with an online platform. The study uses a diversified dataset from an e-commerce platform to investigate the various factors that lead to client segmentation.

The research finds various client segments, each with its own set of qualities and interests, using machine learning algorithms and statistical approaches. E-commerce enterprises can adapt marketing efforts, enhance product recommendations, and refine user interfaces to better resonate with the individual requirements and desires of each client category by gathering insights about these categories.

This study's findings not only add to the academic understanding of client segmentation in e-commerce, but they also have practical consequences for firms looking to improve their marketing tactics in an increasingly competitive online industry. The findings of this study lay the groundwork for designing focused tactics that improve consumer satisfaction, create brand loyalty, and ultimately drive long-term success in the dynamic and ever-changing world of e-commerce.

# **CHAPTER ONE**

## **INTRODUCTION**

### **1.1 BACKGROUND OF THE STUDY**

In the digital age, e-commerce has transformed the way we shop and engage with businesses. Online retailers and wholesalers offer a diverse array of items and services, attracting customers from diverse background inclinations. With such diversity among clients, it becomes inherent for e-commerce platforms to cater to their individual needs and preferences. This is where customer segmentation comes into play.

Customer segmentation is the art of breaking a company's clientele into discrete segments based on qualities and their preferences. It allows companies to learn more about their customers, better their marketing strategies so as to better the overall experience. Through identifying unique segments in their client base, e-commerce platforms deliver personalized content by conveying appropriate product recommendations which boosts customer satisfaction and improves company sales.

The validity of this problem arises from the enormous potential for e-commerce enterprises to obtain a gain competitive advantage in the market. Furthermore, customer segmentation allows companies to allocate their resources efficiently by targeting high-value segments, bettering marketing campaigns and maximizing the (ROI) return on investment.

Customer segmentation is widely used in multiple e-commerce industries and business strategies. This can range from fashion, groceries, electronics and vacation travels. E-commerce platforms whether for a large multinational company or a small niche online company, customer segmentation can provide meaningful insights into consumer behavior, preferences, and spending patterns.

Implementing an effective segmentation plan comes with its own obstacles. A key obstacle involves gathering and analyzing appropriate data so as to identify significant segments. E-commerce businesses must gather extensive consumer data including demographic information such as gender and age, browsing behavior and purchase history. Integrating this data from various sources can be intricate which necessitates the need for robust analytical tools.

To overcome these challenges, advanced machine learning and data analysis techniques are going to be employed in the project. Businesses will identify patterns, categorize their clients based on their similarities, and predict future behavior through employing clustering algorithms, classification models and decision trees.

This chapter will therefore explore the importance of customer segmentation and its significance so as to drive business success.

## **1.2 PROBLEM STATEMENT**

The issue at hand is a lack of efficient client segmentation in e-commerce, which prevents businesses from providing tailored experiences and maximizing marketing efforts. To overcome this issue, it is critical to draw on current research that has addressed comparable difficulties and produced answers in the field of client segmentation in e-commerce. Several scholars have investigated various approaches and techniques for circumventing the constraints of classic segmentation methods.

Here are a few studies that have helped to solve the client segmentation problem:

Chen and Popovich (2003) concentrated on improving consumer segmentation through the use of web mining and data analytics approaches. They used text mining and association rule mining to derive relevant data from client browsing activity and website interactions. Their research demonstrated how examining unstructured data from online platforms might yield useful information for client segmentation and customized marketing strategies.

Verhoef et al. (2014) addressed the issue of scalability in consumer segmentation for large-scale e-commerce enterprises in their study. They created a scalable segmentation system that analyzed and segmented millions of clients in real-time using machine learning techniques and parallel processing. Their research proved how sophisticated technology can overcome standard segmentation methodologies' scaling constraints and boost the efficiency of tailored marketing operations.

These studies are useful resources for understanding the problem, exploring potential solutions, and adapting best practices to manage consumer segmentation challenges and provide more tailored experiences to customers.

## **1.3 OBJECTIVES**

### **1.3.1 GENERAL OBJECTIVES**

The project intends to identify different client segments based on shared characteristics, habits, and preferences by leveraging data analytic tools and machine learning algorithms so as to better customer engagement.

### **1.3.2 SPECIFIC OBJECTIVES**

The deliverables of this project will be:

1. Analyze consumer data to discover important segmentation variables.
2. Create a strong segmentation structure and customer segmentation criteria.
3. Use machine learning techniques to divide your clients into various segments.
4. Using important metrics, assess the performance of the consumer segmentation plan.

## **1.4 RESEARCH QUESTIONS**

The research questions are as follows:

1. What are the precise traits, behaviors, and preferences of each client segment determined by the e-commerce business's segmentation strategy?
2. How can personalized marketing strategies be established and targeted to each consumer category in the e-commerce firm to increase customer engagement, generate conversions, and promote customer loyalty?
3. Which machine learning algorithms and data analysis approaches are best suited for segmenting clients based on their similarities and behaviors in the context of an e-commerce business?
4. What are the primary factors and attributes in available customer data that can be used to effectively segment customers in the e-commerce business?

## **1.5 JUSTIFICATION**

To begin, understanding customers and providing personalized experiences is critical for business success in the competitive e-commerce landscape. Customer segmentation enables

firms to obtain a more in-depth understanding of their customer base, adjust marketing strategies, and build long-term consumer loyalty.

Additionally, with the exponential expansion of online shopping and the availability of enormous amounts of client data, effective data utilization is critical. Customer segmentation enables organizations to find trends and group clients based on similarities using data analysis tools and machine learning algorithms. This enables firms to make data-driven decisions, optimize marketing efforts, and efficiently manage resources.

In summary, the need for the project lies in the need for personalized experiences, the abundance of customer data, and the potential to enhance marketing strategies in e-commerce. By addressing challenges, providing insights and offering practical solutions. This project holds great potential to benefit e-commerce businesses and contribute to the advancement of customer segmentation practices in the industry.

## **1.6 SCOPE AND LIMITATIONS**

The project's goals are to identify essential variables, create a comprehensive segmentation framework, deploy machine learning algorithms for segmentation, assess the strategy's efficacy and provide practical recommendations for individualized marketing tactics. The study will therefore be undertaken in the context of e-commerce enterprises, taking into account various industry sectors and business types. The scope includes both small and large-scale e-commerce firms, with an emphasis on increasing marketing effectiveness, customer engagement, and sales performance.

However, it should be noted that the strategy's efficacy may vary based on external factors like as market dynamics, client preferences e.g., brand preference, political factors and the competitive environment. The project does not include a thorough examination of external influences and instead concentrates on internal segmentation approach.

## **CHAPTER TWO**

### **LITERATURE REVIEW**

#### **2.0 INTRODUCTION**

Understanding clients and providing personalized experiences has become critical for businesses looking to flourish in a competitive field in the changing world of e-commerce. With this there is a need to harness of customer data created from numerous online interactions in order to obtain deeper insights into customer behavior and preferences. Customer segmentation, which divides the client base into discrete groups based on similar features and behaviors, has emerged as a significant technique for addressing this difficulty.

In this literature review we will look at a thorough overview of client segmentation in e-commerce. We will investigate the theoretical underpinnings of segmentation practices, investigate and understand various methodologies used in customer segmentation, challenges faced in implementing effective customer segmentation, evaluate the outcomes and impacts of segmentation on business performance as well as analyze future directions and trends in this industry.

#### **2.1 HISTORY OF CUSTOMER SEGMENTATION**

Customer segmentation has been in evolution since the pre 19<sup>th</sup> Century till now, the 21<sup>st</sup> Century and it has gone through seven critical phases. A detailed analysis of the phases is as follows:

##### **2.1.1 Early application in Traditional Retail (19<sup>th</sup> century)**

Prior to the industrial revolution, retailers and craftspeople attempted rudimentary consumer segmentation by getting to know their clients on a personalized level. (John Maynard Keynes, 1938) They customized products and services to meet the interests and demands of each individual. As their trade grew, they realized the value of geographic segmentation and adapted their inventories and marketing techniques to different regions' interests and preferences.

##### **2.1.2 Phase 2: Mass marketing (Early 20<sup>th</sup> century)**

(Wedel and Kamakura, 2001) stated that, as production efficiency improved and product variation rose in the early twentieth century, the concept of market segmentation became a formal component of marketing practice. "Industrial development in various sectors of the economy induced mass production and marketing strategies." Those tactics were



manufacturing-oriented, with an emphasis on lowering production costs rather than customer happiness.

However, as manufacturing processes became more flexible and consumer affluence resulted in demand diversification, enterprises that understood the specific demands of groups of customers were able to design the proper offer for one or more sub-markets and therefore gained a competitive advantage. (Chamberlin, 1933).

The rise of mass production in the early twentieth century enabled corporations to make standardized products on a large scale. This was done in order to reach the greatest number of potential clients and to appeal to the broadest segment of the market while ignoring niche demographic variations. This strategy prioritized huge sales volumes at affordable prices. An example was with (Henry Ford, 1913) whose mass marketing method was centered on reaching large, undifferentiated markets.

(Joan Robinson, 1938) developed the economic theory of imperfect competition by expanding on this premise. It stated that firms have some pricing power in this paradigm due to product differentiation, allowing them to establish prices, participate in price discrimination, and potentially generate short-run profits or losses. (John Maynard Keynes, 1938) also added that in the long run, these businesses may have excess capacity and produce below their technical efficiency level. Robinson's research calls into question the standard perfect competition paradigm, providing a more nuanced understanding of market dynamics and economic wellbeing in situations when enterprises have some degree of market dominance.

### **2.1.3 Phase 3: Post-war segmentation (Mid-20<sup>th</sup> century)**

In the mid-twentieth century, demographic segmentation became popular. Advertisers started categorizing consumers based on demographics such as age, gender, income, occupation, and education. It was influenced by the significant shift in marketing in response to the social and economic changes after World War 2.

(J. Smith, 1956) acknowledged "the existence of heterogeneity in the demand for goods and services, based on Robinson's economic theory of imperfect competition" (Wedel and Kamakura, 2000). "Market segmentation involves viewing a heterogeneous market as a number of smaller homogeneous markets in response to differing preferences, attributable to consumers' desires for more precise satisfaction of their varying wants," Smith claimed.

This was a crucial stage as it laid the foundation in the modern practices for customer segmentation through improved advertising efficiency by targeting certain demographic groups.

#### **2.1.4 Phase 4: Market Research Era (Mid to Late 20th Century)**

Market research methodologies developed (Anna-Lena, 2001), Stated that segments should be based on consumer/user desires, and a company should be better able to serve these needs when some segments within a larger market have been defined" allowing businesses to gain a better understanding of their customers' views, values, and behaviors. Psychographic segmentation took into account aspects such as lifestyle, values, and personality. Customer activities, such as purchase history and brand loyalty, were utilized to categorize customers. This gave rise to retail loyalty programs in late 20<sup>th</sup> century where which firms collected client data via cards and offered benefits in exchange.

#### **2.1.5 Phase 5: Database and Technology Revolution Era (Late 20th Century)**

As technology advanced, businesses began to use database marketing to collect, store, and analyze customer data more effectively. CRM (client Relationship Management) systems were developed, allowing businesses to track client interactions and customize marketing activities accordingly. (Bergeron, B, 2002)

With the help of CRM, it gave rise to the factor that consumer segmentation refers to the practice of identifying the most receptive groups within a population and targeting them with the most relevant messaging (Frank, Massy, Wind 1972, McDonald & Dunbar 2004, Anna-Lena 2001, Jiang and Tuzhilin, 2006). (Wind and Cardozo,1974) recognized that consumer segmentation also known as market segmentation at the time "involves appropriate grouping of individual customers into a manageable and efficient in a cost/benefit sense number of market segments, for each of which a different marketing strategy is feasible and likely profitable" (Wind and Cardozo, 1974). Since then, several marketers have argued for the creation of consumer segments in order to better align products and services with certain groups.

Direct Marketing Emerged based on segmentation and it became a popular marketing method. Businesses employed methods such as direct mail, telemarketing, and, email marketing so as to reach specific client segments.

### **2.1.6 Phase 6: Digitalization and Online Era**

As the internet and e-commerce grew in popularity, customer segmentation expanded into the digital arena. (Kotler and Armstrong, 2007) Customers' online behavior, clickstream data, and social media interactions have all become useful sources of information. Businesses began segmenting online clients based on their digital footprint using algorithms and machine learning thus creating a precedent for e-commerce personalization.

Six major criteria are recognized in contemporary segmentation methodologies for evaluating segmentation effectiveness (Wedel and Kamakura, 2000) These criteria evaluate segment development and profitability to determine effectiveness and say that customer segments should be: identifiable/measurable, sizable, accessible, stable actionable, and differentiable. The 21<sup>st</sup> Century saw a substantial shift towards customization. Customer segmentation is used by businesses to provide individualized product suggestions, targeted advertising, and customized user experiences. The big data age has provided a lot of information that can be used for segmentation. (Anna-Lena, 2001). Businesses may identify detailed client segments and forecast future behavior with unparalleled accuracy using advanced analytics, such as predictive modeling and AI algorithms. This can be seen in highly tailored services include Netflix's content recommendations and Spotify's music recommendations.

## **2.2 Theoretical Framework**

### **2.2.1 Supervised learning approach**

#### **2.2.1.1 Decision Trees**

Decision trees are a popular method for customer segmentation. They construct a tree-like structure, with each internal node representing a feature or attribute and each leaf node representing a client segment. (J. Han, 2006). The model makes decisions based on feature values at each node, eventually assigning customers to certain segments. They are simple to interpret and explain and can deal with both categorical and numerical data. (M. Kamber, 2006).

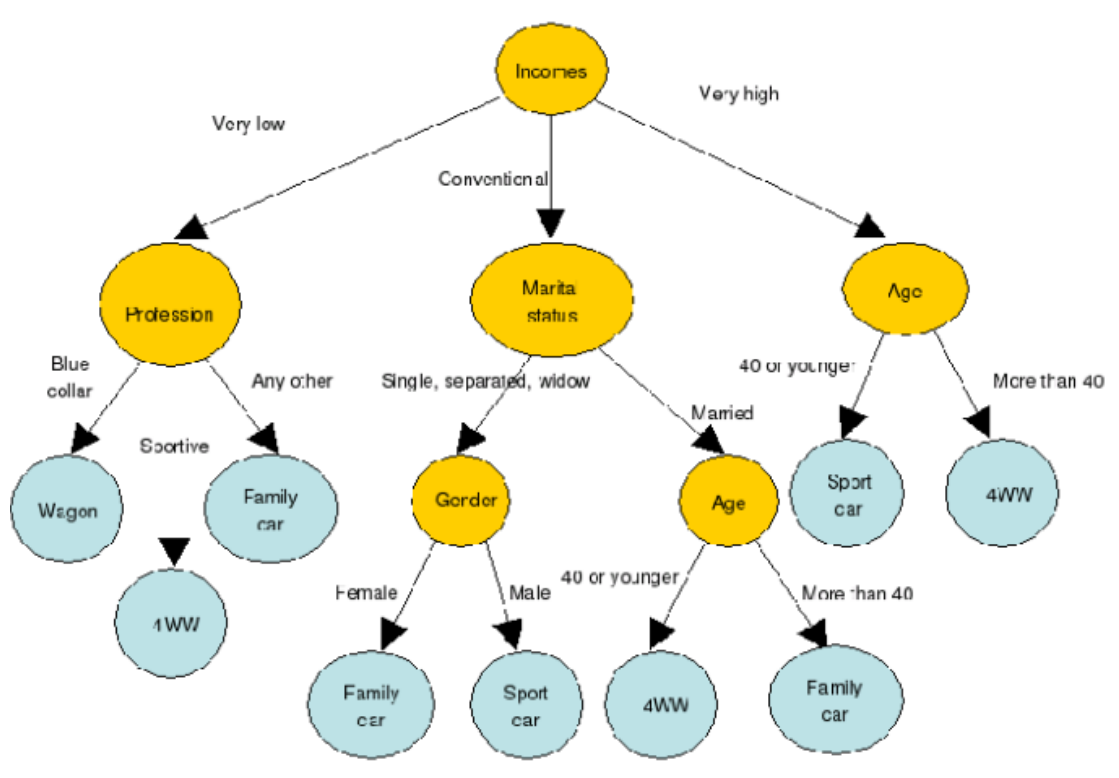


Figure 1. Decision Tree

Both classification and regression problems make use of decision trees. They can be used to classify clients into segments in the context of customer segmentation.

The program recursively divides the data into subsets depending on the most important features, with the goal of maximizing information gain or Gini impurity.

The process is repeated until a stopping requirement, such as a maximum depth or a minimum number of samples in a node, is reached.

Decision trees are useful for analyzing client segments since they are easily viewed and analyzed.

### 2.2.1.2 Random Forests

(Leo Breiman, 2001), stated that a random forest is a collection of decision trees. They build numerous decision trees based on random subsets of data and features, and the final segmentation choice is based on a majority vote or an average of all trees. They outperformed a single decision tree in terms of accuracy and resilience. We deal with overfitting.

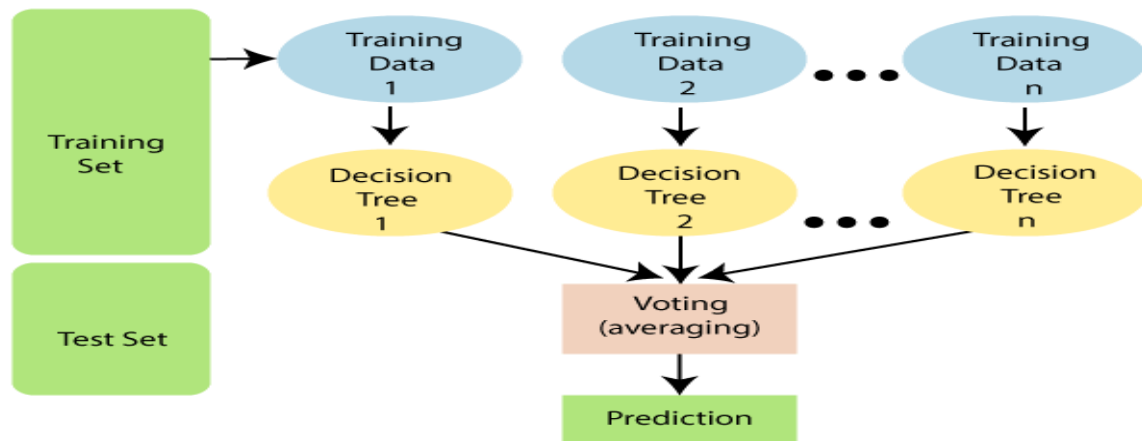


Figure 2. Random Forests

(Leo Breiman, 2001) is well-known for his work in decision trees, but he was also a forerunner in the development of the random forest model. In consumer segmentation, random forests have been widely used. This ensemble learning technique combines numerous decision trees to increase accuracy and decrease overfitting, making it a useful tool for identifying client segments.

### 2.2.1.3 Support Vector machines

Vladimir Vapnik's pioneering work on Support Vector Machines (SVMs) focuses on finding a hyperplane that best separates different classes in data. SVMs maximize the margin between classes while minimizing classification errors. (Vladimir Vapnik, 1960). SVMs are effective in high-dimensional spaces and have been applied to customer segmentation as part of their broader application in classification tasks.

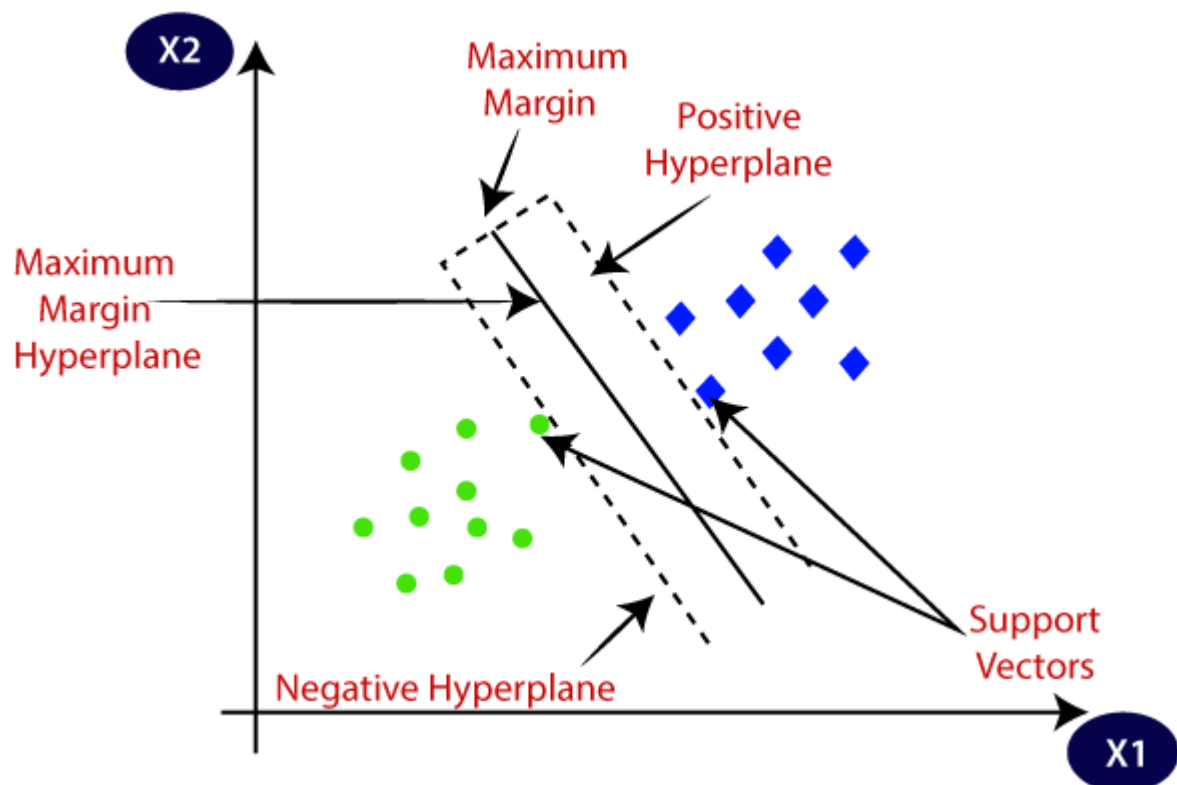


Figure 3. Support Vector Machines

SVM is a binary classification technique that can be extended to multi-class issues. It can classify clients into segments in customer segmentation.

The program searches for a hyperplane that maximizes the margin between data points of distinct classes.

SVM determines support vectors (data points nearest to the decision border) during training.

New data points are mapped into the feature space during prediction, and their position relative to the decision boundary defines the class.

#### 2.2.1.4 Logistic Regresssion

(Josph Berkson, 1974) stated that based on input features, logistic regression models the likelihood that a given observation belongs to a specific class or category. (David Cox, mid-20<sup>th</sup> century). It is used for binary or multiclass classification in the context of customer segmentation.

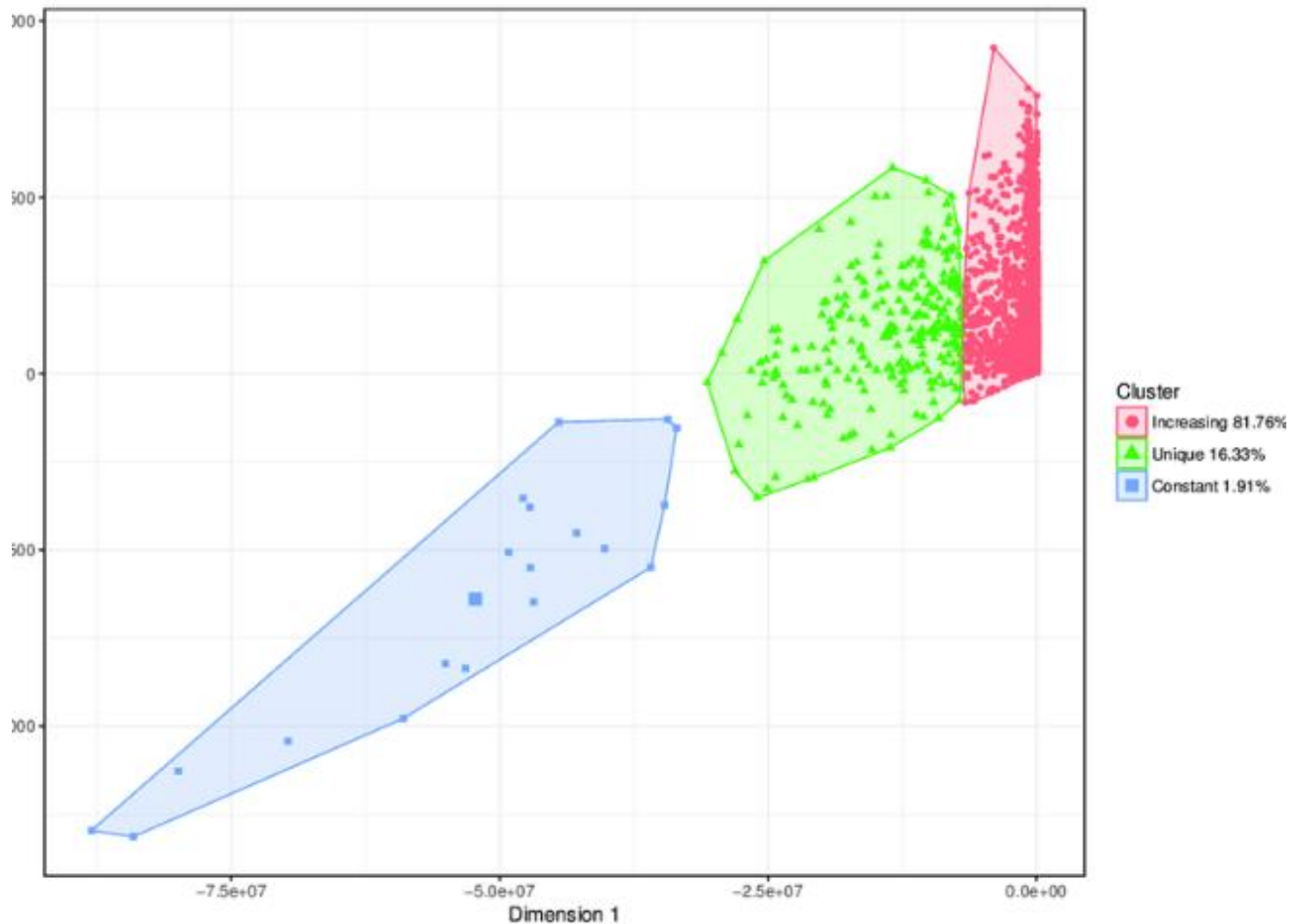


Figure 4. Logistic Regression

Logistic regression is most commonly used for binary classification problems, where the goal is to estimate the likelihood that an input belongs to a specific class (customer segment).

Using the logistic function, the method models the relationship between the independent variables (consumer qualities) and the binary dependent variable (segment).

The logistic curve is fitted to the data during training using methods such as maximum likelihood estimation.

During prediction, the model computes the likelihood of a client belonging to a segment and applies a threshold to classify them into one of two groups.

Aside from previous citations of the logistic regression model, we can add the contribution of academic David W. Hosmer, Jr., (1989), whose work in biostatistics and epidemiology has stressed the use of logistic regression in evaluating customer behavior. Logistic regression models are commonly used to forecast customer outcomes and segment them based on past data and attributes.

### 2.2.1.5 Neural Networks (Deep learning)

Deep learning, particularly neural networks, has evolved with contributions from various scholars such as (Geoffrey Hinton, 1980) who stated that Neural networks consist of layers of interconnected nodes that learn complex patterns and relationships between features and customer segments. Yann (LeCun, Yoshua Bengio, 1986) stated that they have been widely used for customer segmentation and are capable of capturing intricate patterns in data.

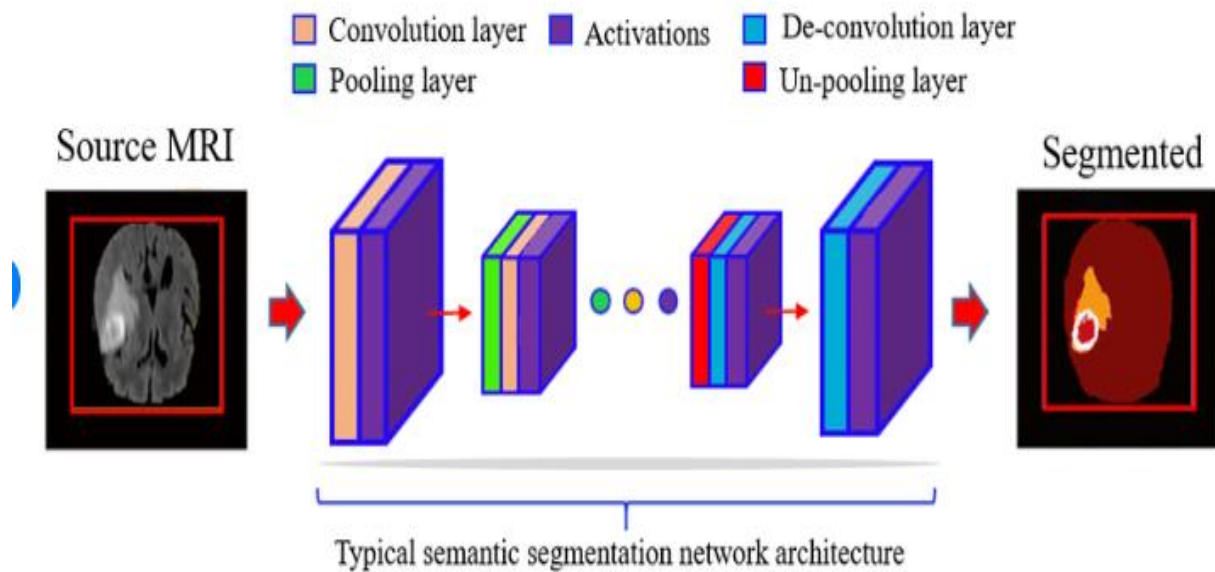


Figure 5. Neural Networks

#### 2.2.1.5.1 Image-Based Customer Segmentation with CNNs

Yann LeCun is well-known for his contributions to convolutional neural networks (CNNs) and related image recognition applications. While his research is more broadly focused on computer vision, CNNs play an important role in image-based customer segmentation.



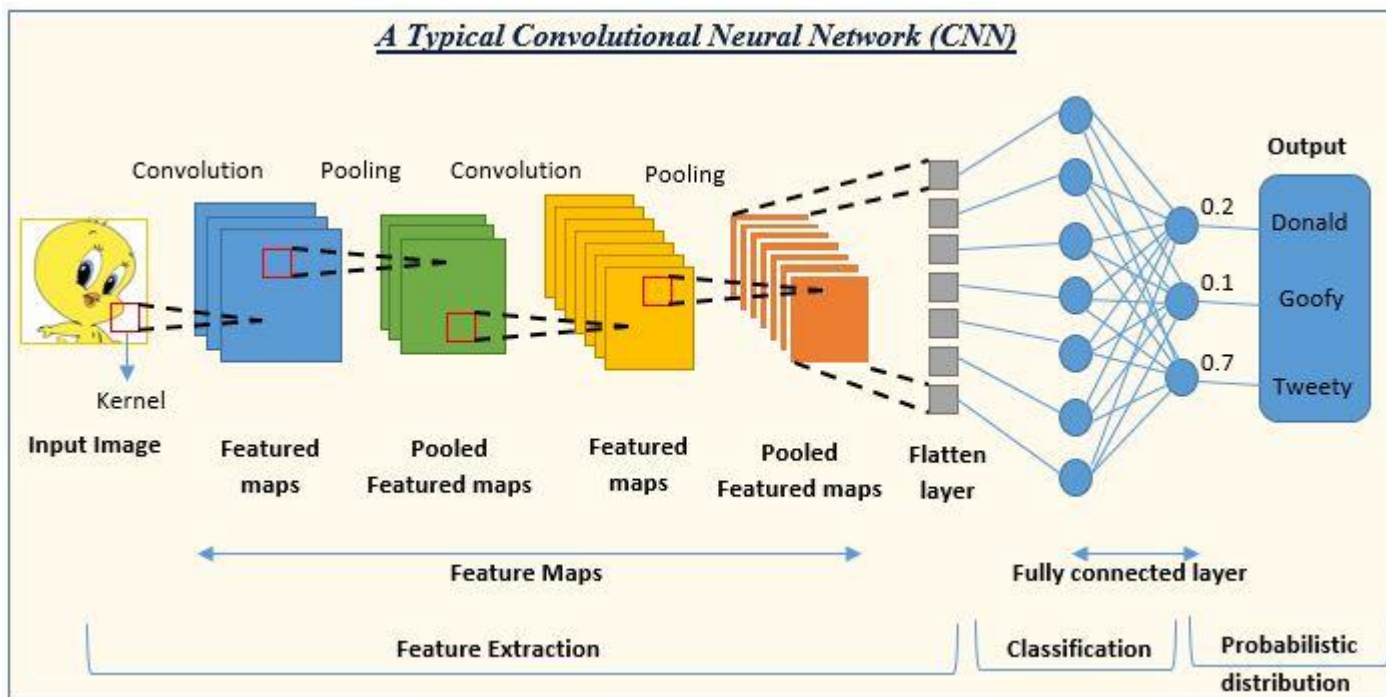


Figure 6. Imaged based Customer Segmentation

CNNs are deep learning algorithms that are used to analyze and segment images. CNNs use client photos or image-related data as input in image-based consumer segmentation. Convolutional layers are used to extract features, followed by fully linked layers for segmentation. To produce pixel-wise segmentations or categorize complete images into segments, training entails minimizing a segmentation loss.

#### 2.2.1.5.2 Sequential Data Analysis with RNNs and LSTMs

(Juergen Schmidhuber,1990) is a well-known expert in recurrent neural networks (RNNs). His work has expanded our understanding and development of RNNs, as well as its applications in sequential data analysis for customer segmentation.

RNNs and LSTMs are neural network architectures that are used to analyze sequential data.

They're ideal for consumer data with temporal relationships, including time series data.

RNNs and LSTMs may recognize sequential patterns by recurrently processing input and taking past time steps into account.

Backpropagation through time (BPTT) is used in training to update network weights and optimize predictions.

### 2.2.1.5.3 Dimensionality Reduction with Autoencoders

(Geoffrey Hinton's, 1980) contributions to neural networks, notably autoencoders, have been significant since then and continue to this day. His contributions to deep learning have included autoencoders. DBSCAN is a clustering technique based on density.

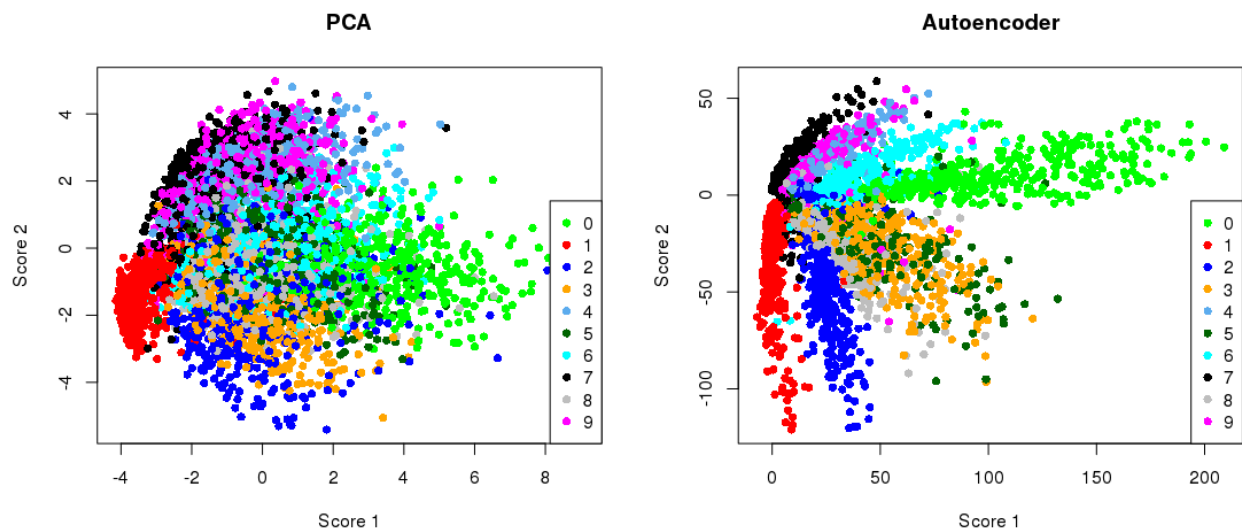


Figure 7. Dimensionality Reduction

It clusters data points depending on density, essentially detecting clusters of various forms.

The algorithm specifies core points, border points, and noise points, and clusters develop with a minimum number of neighbors around core locations.

DBSCAN is useful for locating clusters of various forms and sizes.

### 2.2.1.5.4 Customer Image Generation with GANs

(Ian Goodfellow, 2014) is well-known for his research on generative adversarial networks (GANs). GANs are frequently used for picture production, and they can be utilized to generate images indicating consumer preferences for segmentation.

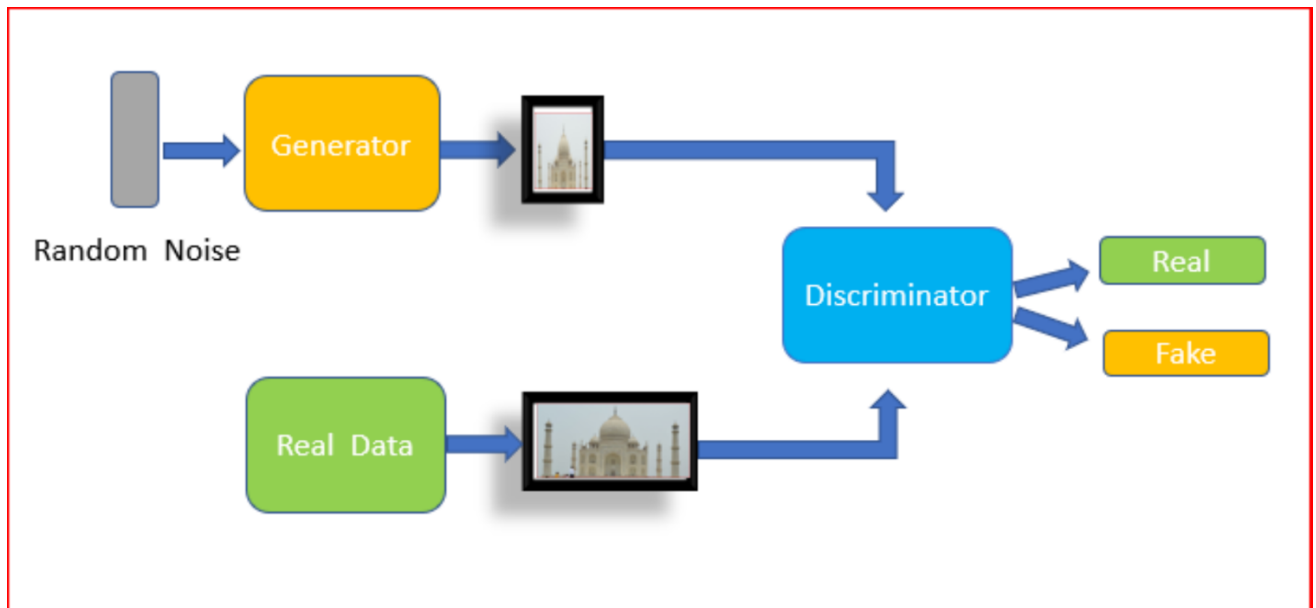


Figure 8. GAN Image Generation

#### 2.2.1.5.5 Natural Language Processing with Transformers

(Thomas Wolf, Lysandre Debut, Victor Sanh, 2019) The researchers who created the Hugging Face "Transformers" library have helped to popularize and apply transformer models like BERT in natural language processing. These models are essential for assessing textual data in the context of consumer segmentation.

These scholars (Julien Chaumond, Clement Delangue, 2019) Customer feedback and complaints can be classified and prioritized with the use of transformers. This can assist in more effectively identifying and addressing issues, as well as improving customer satisfaction.

NLP approaches can be used to evaluate conversational data from customer care chatbots in order to understand common client inquiries and difficulties. (Anthony Moi, 2019) This information can be utilized to segment consumers based on their assistance requirements.

### 2.2.1.6 K-Nearest Neighbors (K-NN)

Based on feature similarity, K-Nearest Neighbors allocates a client to the segment that is most frequent among its k-nearest neighbors. (D. L. Johnson and S. Kotz, 1970). It measures similarity using distance measurements. K-NN is non-parametric and easy to construct, making it helpful for segmentation of local patterns and nonlinear decision boundaries. (Evelyn Fix and Joseph L. Hodges, 1951).

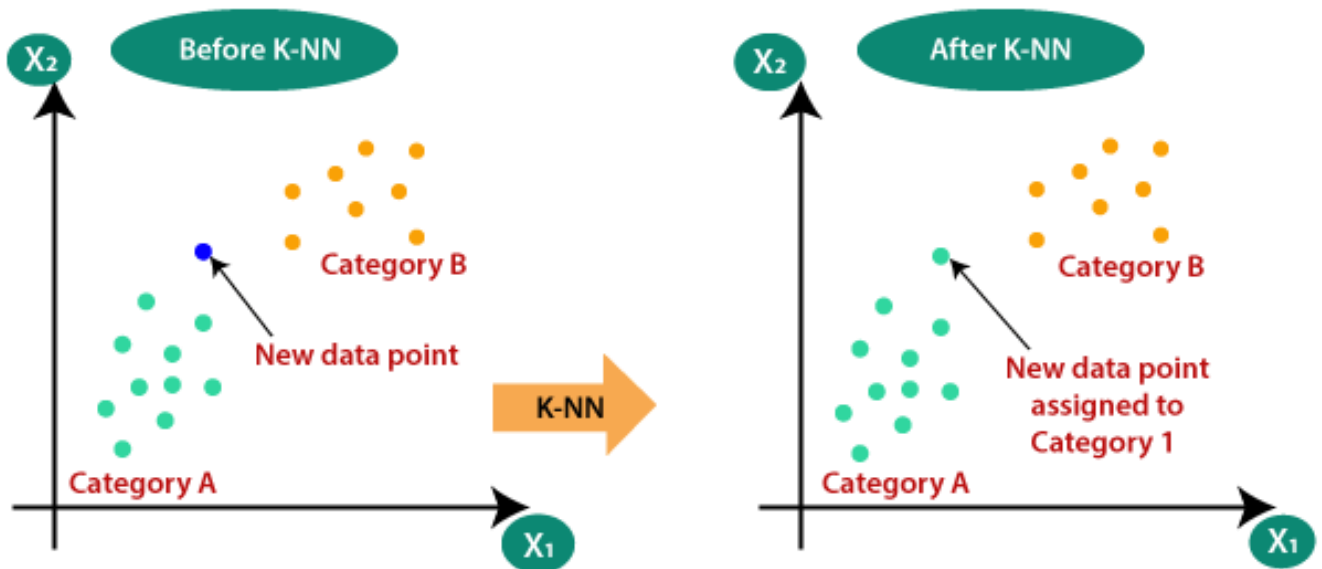


Figure 9. K.N Nearest

The K-NN algorithm computes the distance between the new data point and all of the training instances first. The distance metric employed is determined on the type of data being classified. The Euclidean distance is a typical distance metric in the case of photographs.

After calculating the distances, the K-NN method selects the K most similar training instances to the new data point. The K-NN algorithm selected the three most comparable training samples

The K-NN method then guesses the new data point's class label based on the class labels of the K most comparable training samples.

### 2.2.1.7 Gradient Boosting Machines

Gradient Boosting Machines, which include popular libraries such as XGBoost and LightGBM, generate a set of decision trees. (Tianqi Chen, Mu Li, 2010). To improve segmentation accuracy, each tree corrects the mistakes of the previous ones. These strategies

are well-known for their great accuracy and efficiency in dealing with complex consumer data linkages.

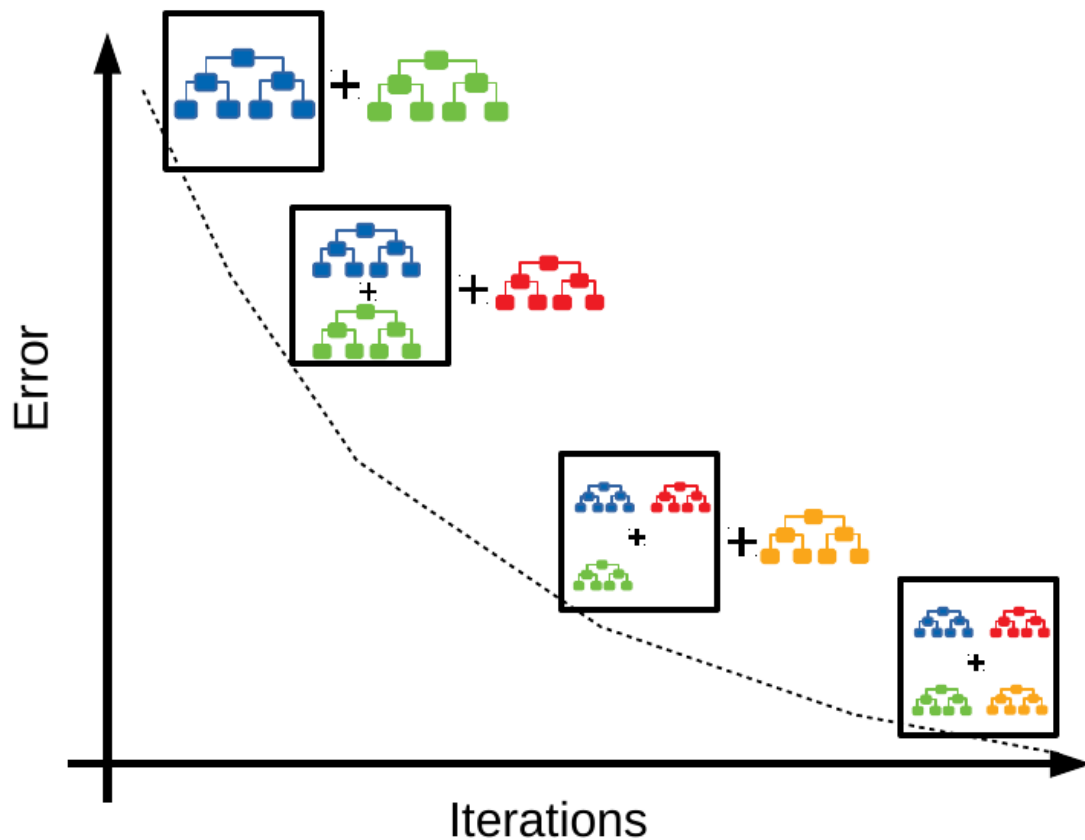


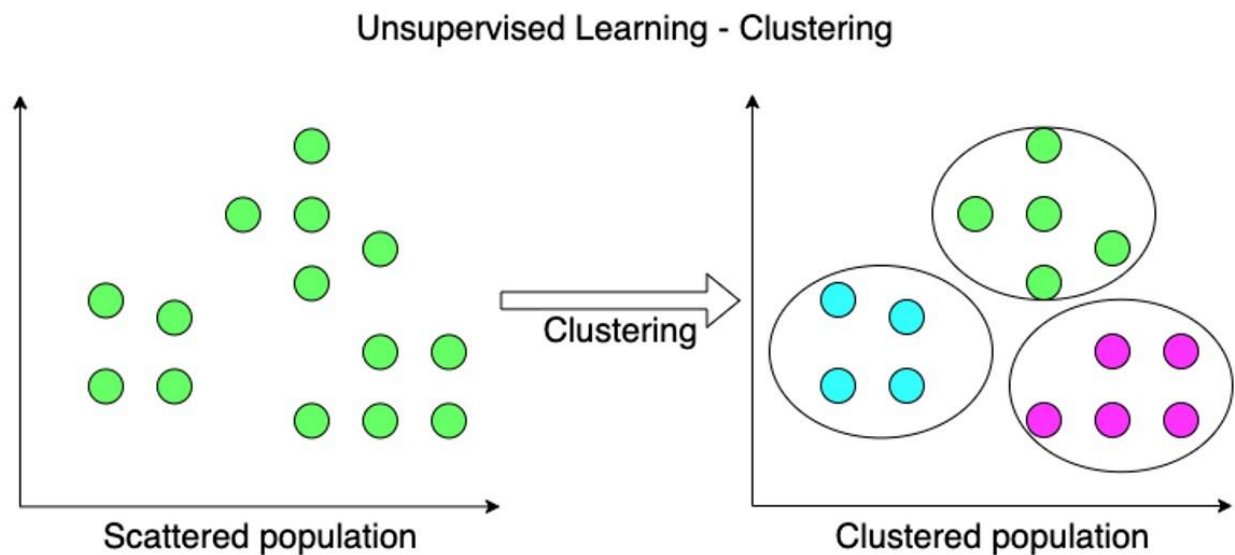
Figure 10. Gradient Boosting Machine

GBM is a classification and regression ensemble learning algorithm. To produce a powerful predictive model, it integrates numerous weak learners (usually decision trees). The algorithm trains consecutive trees, with each tree focusing on rectifying previous trees' mistakes. The results from all trees are merged to generate a final prediction during prediction.

## 2.2.2 Unsupervised machine learning

### 2.2.2.1 K-Means Clustering

K-Means is a popular unsupervised learning method used in consumer segmentation. (Lloyd, 1982) It intends to divide clients into K clusters based on their similarities. It reduces volatility within each cluster while increasing variance between clusters. (Hastie, T., Tibshirani, R., & Friedman, J. 2009) Each client is assigned to the cluster with the nearest centroid, and the process is repeated until convergence is reached.



*Figure 11. Unsupervised Learning*

According to (Anil K. Jain, 1988) a pioneer in the field of clustering and pattern recognition, K-means clustering is a frequently used technique in consumer segmentation, putting customers into groups based on similar attributes.

#### **2.2.2.2 Hierarchical Clustering**

Hierarchical clustering creates a nested tree-like structure. All customers are in a single cluster at the top, and when the tree is visited, clusters divide into smaller ones. (Rokach, L., & Maimon, O. 2005) This method does not require a predetermined number of clusters, making it adaptable for client segmentation.

Hierarchical clustering generates a hierarchy of clusters that can be agglomerative (bottom-up) or divisive (top-down).

It starts with a single cluster for each data point and then merges or divides clusters based on a defined linking mechanism (e.g., single, complete, or average linkage).

The ultimate result is a hierarchical tree (dendrogram) that represents the interactions between clusters and allows for segmentation at multiple levels.

#### **2.2.2.3 DBSCAN**

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) may detect clusters of any shape. (Ester, M., Kriegel, H. P., Sander, J., & Xu, X. 1996) Customers who are close to each other in terms of density are grouped together. Customers in low-density areas are seen as noise or outliers. (Jain, A. K., & Dubes, R. C. 1988).

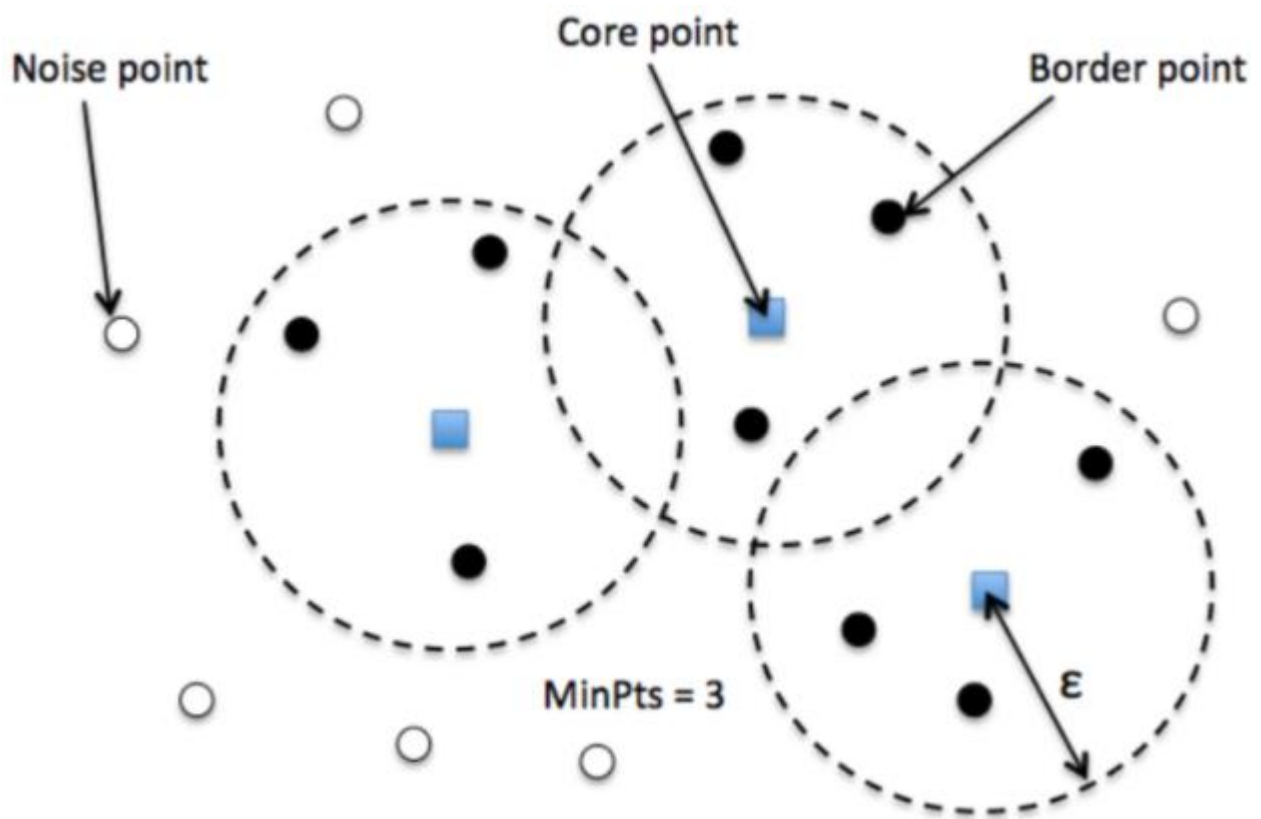


Figure 12. DBSCAN

DBSCAN is a clustering technique based on density.

It clusters data points depending on density, essentially detecting clusters of various forms.

The algorithm specifies core points, border points, and noise points, and clusters develop with a minimum number of neighbors around core locations.

DBSCAN is useful for locating clusters of various forms and sizes.

#### 2.2.2.4 Principal Component Analysis (PCA) for Dimensionality Reduction

PCA is a dimensionality reduction approach that can be used before clustering, not a clustering method. (Jolliffe, I. T. 2002) It maps high-dimensional consumer data to a lower-dimensional space while retaining as much variance as possible. Reduced data can thus be more effectively grouped.

PCA is a technique for extracting the most essential characteristics from data while retaining variation.

It converts the data into a new coordinate system (principal components) in which the first component captures the most variation, the second the second, and so on.

PCA can minimize the dimensionality of data in consumer segmentation, making it easier to display and evaluate.

#### 2.2.2.5 Latent Dirichlet Allocation (LDA) for Topic Modeling

LDA is often used in text mining, (Blei, D. M., Ng, A. Y., & Jordan, M. I. 2003) but it can also be used to segment customers when they submit textual comments or reviews. It finds latent themes in data and groups clients into categories depending on their topic preferences.

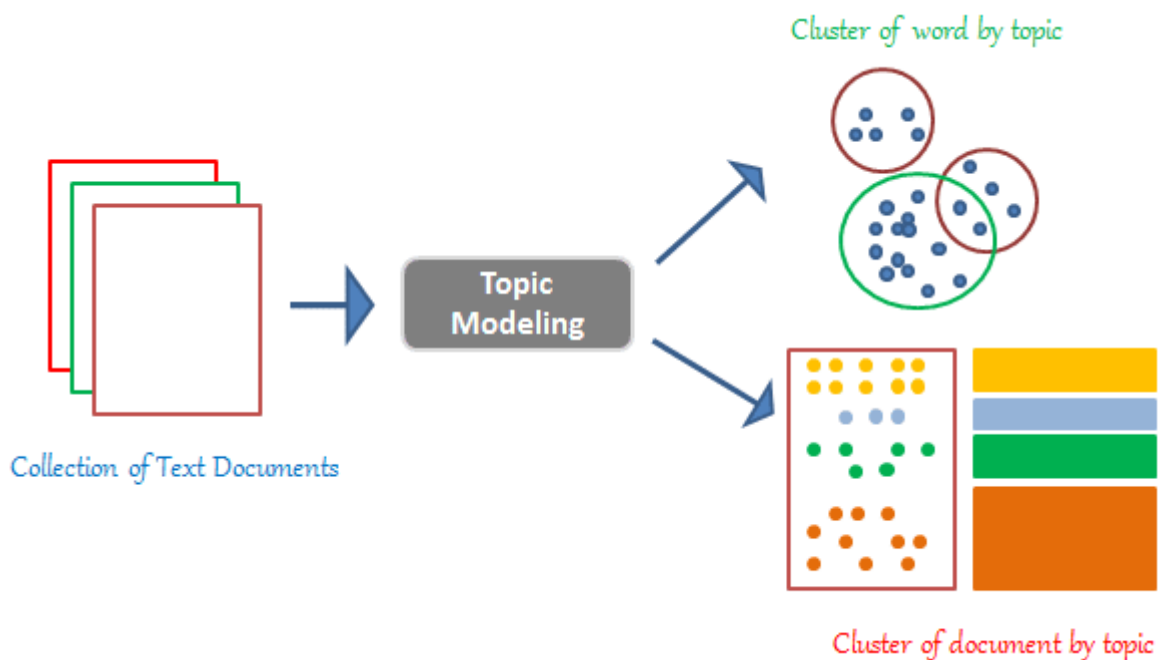


Figure 13. Latent Dirichlet Allocation (LDA)

LDA is a subject modeling generative probabilistic model. It is assumed that documents are made up of themes, and topics are made up of words. During training, LDA recognizes and assigns words to subjects in a set of documents. The result includes a list of subjects as well as the word distribution associated with each topic.

#### 2.2.2.6 Gaussian Mixture Model (GMM)

A Gaussian Mixture Model (GMM) is a probabilistic model that describes data point distributions as a mixture of numerous Gaussian distributions. (Fraley, C., & Raftery, A. E. 2002) It is especially effective for modeling data that may be divided into segments or clusters. GMM can uncover underlying patterns in consumer behavior and attributes in the



context of customer segmentation, allowing clients to be assigned to segments based on the Gaussian distributions to which they belong.

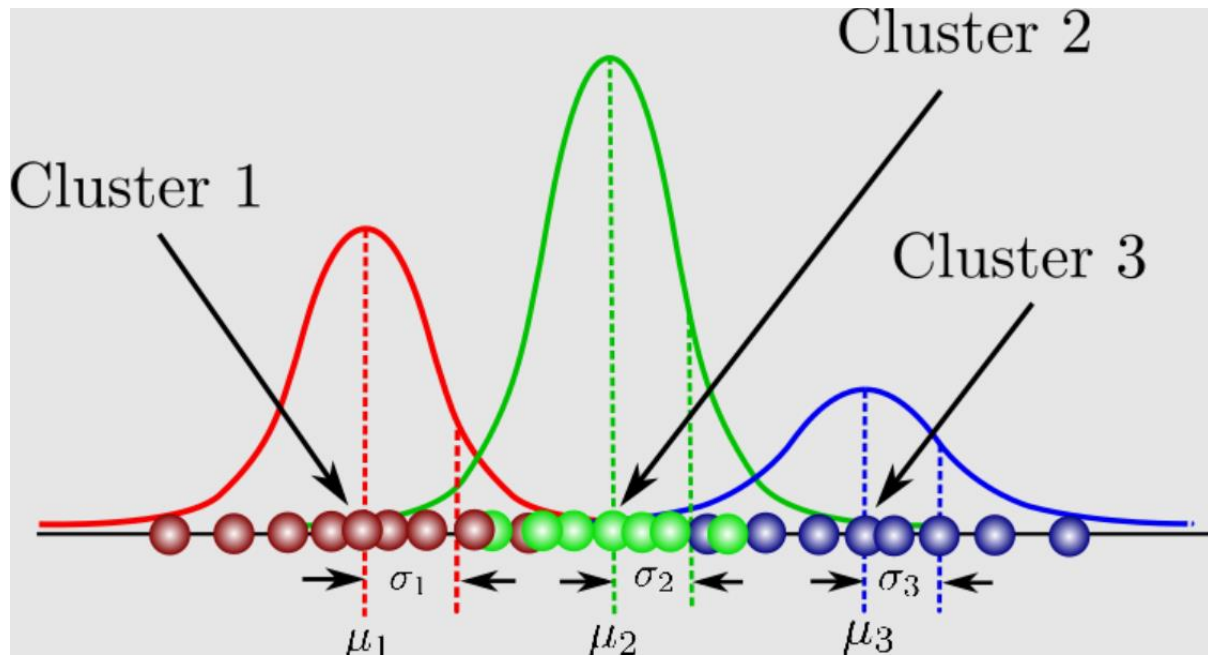


Figure 14. Gaussian Mixture Model (GMM)

GMM is a probabilistic model that is used for clustering and estimating density. It is assumed that data points are created by a Gaussian mixture. During training, the method estimates the Gaussian components' parameters (means, variances, and weights). GMM can be used for soft clustering, in which each data point is assigned to many groups, each with its own probability.

#### 2.2.2.7 Self-Organizing Maps (SOM) for Customer Segmentation

Self-Organizing Maps (SOM) are a type of artificial neural network that organizes high-dimensional data into a lower-dimensional grid. (Kohonen, T. 1982) SOMs are useful for visualizing and clustering complex data. In customer segmentation, SOMs can reveal the topological relationships between customer groups and help identify clusters of similar customers based on their preferences or behaviors.

SOM is an unsupervised grouping and dimensionality reduction artificial neural network. The algorithm maps data onto a grid of nodes, each of which represents a cluster prototype.

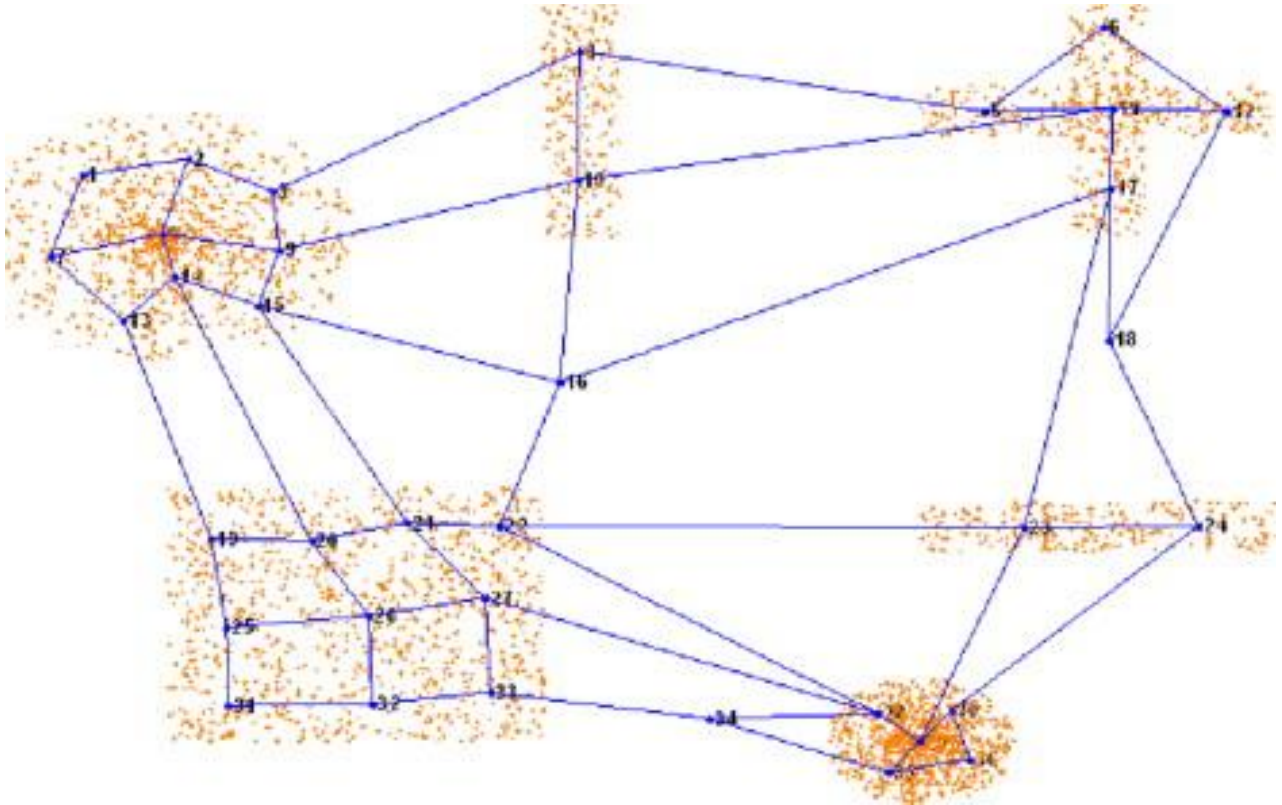


Figure 15. Self-Organizing Maps (SOM)

SOM alters prototypes iteratively to best reflect data and produces clusters. It is very effective for viewing and detecting high-dimensional data.

#### 2.2.2.8 Associate rule mining

The seminal work of (Agrawal, Imielinski, and Swami's, 1993) in the subject of data mining, particularly association rule mining, is important. Association rule mining identifies patterns in consumer behavior by revealing relationships between various products or services that customers typically buy together. This strategy is useful for analyzing market baskets and segmenting customers based on their purchasing habits.

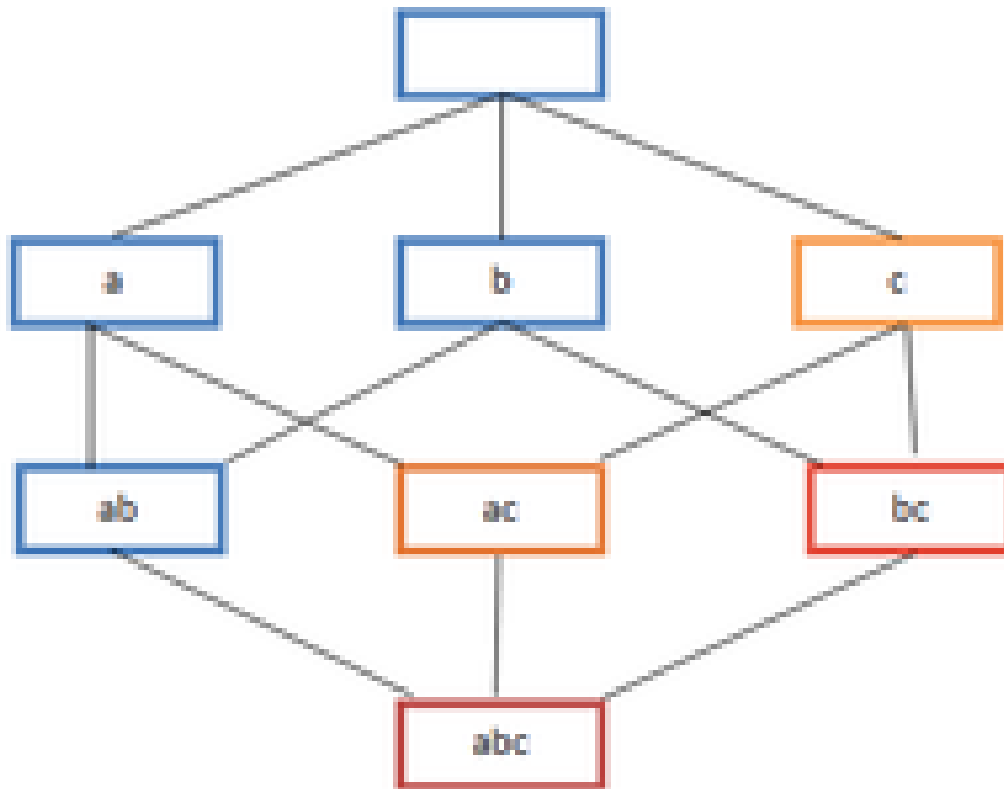


Figure 16. Associate rule mining

The mining of association rules identifies associations between items in transactions.

It can be utilized in the context of e-commerce to uncover correlations between products that buyers regularly purchase together.

To determine meaningful relationships, common variables employed include support, confidence, and lift.

Product suggestions and client segmentation can be guided by rules with high support and confidence.

## 2.3 TECHNIQUES IN CUSTOMER SEGMENTATION

### Demographic Segmentation:

Demographic segmentation divides clients into groups based on demographic criteria such as age, gender, income, and education. (Philip Kotler, 1960s), He stated that it helped in tailoring products and marketing messaging to various age groups, genders, or economic levels.

### Geographic Segmentation:

Geographic segmentation divides clients into groups depending on their geography, such as country, region, city, or address. (Kotler, Philip, 1972) stated that it helped in optimizing business locations, localizing marketing campaigns, providing location-based advertising and Urban vs. Rural: Distinguishing between customers in urban and rural areas, considering different lifestyles and needs.

### **Psychographic Segmentation:**

Customers are classified into psychographic segments based on their lifestyle, values, interests, attitudes, and personality attributes. (William Wells, 1966). Creating material and advertising that caters to specific lifestyle groups, as well as personalizing items to match hobbies and values

### **Behavioral Segmentation:**

Customers are segmented depending on their behavior, which includes purchase history, internet activity, loyalty, and response to marketing initiatives. (William Lazer ,1963). Personalizing product recommendations, targeting based on browsing history, and re-engaging dormant consumers are all examples of use cases.

### **Customer Lifecycle Segmentation:**

Customer Lifecycle Segmentation is forecasting customer behavior or preferences using algorithms and statistical models, (Trevor Hastie and Robert Tibshirani, 2009), Based on previous data, predict customer attrition (Identifying at-risk customers who may be likely to leave and targeting them with retention strategies), forecast revenues, and recommend items. Also, it includes New Customers vs. Returning Customers, distinguishing between first-time buyers and repeat customers. Advocates: Recognizing and nurturing customers who are loyal and serve as brand advocates. (Sung Ho Ha & Sung Min Bae, 2006)

### **Needs-Based Segmentation:**

Customers are classified using Needs-Based Segmentation based on their individual needs, goals, and preferences, with an emphasis on their motivations and intended solutions. (Malcolm McDonald and Ian Dunbar, 1998). Tailoring product development, marketing communications, and customer assistance to specific customer demands and effectively addressing problem points.

### **Value-Based Segmentation:**

Value-Based Segmentation categorizes clients based on their economic value to the firm, taking into account aspects such as customer relationship profitability and long-term revenue generation potential. (Sunil Gupta, 2005) Targeting high-value consumers for marketing and retention, recognizing opportunities for increased profitability, and personalizing offers to enhance client lifetime value.

### **Surveys and Feedback Analysis**

To learn client preferences and sentiments, surveys and feedback analysis are used. (Valarie Zeithaml, 1988). Customer input is used to improve products and services, alleviate pain points, and personalize customer assistance.

### **Channel and Interaction-Based Segmentation:**

Customers are segmented based on their preferred communication channels and interaction history, taking into account multiple touchpoints such as email, social media, website, and in-store visits. (Stephen J. Shaw, 2012). Tailoring communication and marketing tactics for customers depending on their chosen channels, optimizing multichannel marketing efforts, and providing personalized experiences.

### **RFM Analysis (Recency, Frequency, Monetary)**

Customers are evaluated using RFM analysis based on their recent purchases, frequency of transactions, and monetary value. (Robert C. Blattberg and John Deighton, 1996).

Recognizing high-value customers, implementing customized loyalty programs, and re-engaging customers who haven't purchased in a while.

## **2.4 Datasets and Data Sources in Customer segmentation**

Datasets are essential in a machine learning context for customer segmentation since they serve as the foundation for building and refining segmentation models. These datasets provide raw consumer data that is required for training machine learning models to uncover patterns and relationships in the data. These datasets are used in feature engineering, a vital phase in segmentation, to select, manipulate, or create relevant variables. Datasets are critical for model evaluation, with a subset of the dataset being saved for validation and testing to

examine the model's performance and ability to generalize to previously unseen customer data. They support hyperparameter adjustment, allowing the model's design to be fine-tuned for optimal results. Datasets assist data preprocessing since they require cleansing, normalization, and preparation, which affects data quality of the output.

### 2.4.1 Available datasets for Customer Segmentation

Examples include:

The two datasets linked to red and white "Vinho Verde" wines from Portugal. (Cortez et al., 2009) which can be thought of as classification or regression tasks. The classes are organized but not balanced, for example, there are far more ordinary wines than exceptional or terrible wines. (Telmo Matos, J. Reis. 2009) Outlier identification algorithms could be used to identify the few good or bad wines.

Another example is the transnational data set containing all transactions for a UK-based and registered non-store internet retailer that occurred between 01/12/2010 and 09/12/2011. The company primarily sells one-of-a-kind all-occasion presents. Many of the company's clients are wholesalers.

Variable Name	Role	Type	Demographic	Description	Units	Missing Values
InvoiceNo	ID	Categorical		a 6-digit integral number uniquely assigned to each transaction. If this code starts with letter 'c', it indicates a cancellation		no
StockCode	ID	Categorical		a 5-digit integral number uniquely assigned to each distinct product		no
Description	Feature	Categorical		product name		no
Quantity	Feature	Integer		the quantities of each product (item) per transaction		no
InvoiceDate	Feature	Date		the day and time when each transaction was generated		no
UnitPrice	Feature	Continuous		product price per unit	sterling	no
CustomerID	Feature	Categorical		a 5-digit integral number uniquely assigned to each customer		no

Figure 17. Variables Table

Another example is the iris data set which is one of the oldest datasets used in the classification literature, and it is widely used in statistics and machine learning. The data set

has three classes of 50 instances each, with each class referring to a different species of iris plant.

**Variables Table**

Variable Name	Role	Type	Demographic	Description	Units	Missing Values
sepal length	Feature	Continuous			cm	no
sepal width	Feature	Continuous			cm	no
petal length	Feature	Continuous			cm	no
petal width	Feature	Continuous			cm	no
class	Target	Categorical		class of iris plant: Iris Setosa, Iris Versicolour, or Iris Virginica		no

Figure 18. Variables Table 2

The performance of this dataset as per accuracy of different models is as follows:

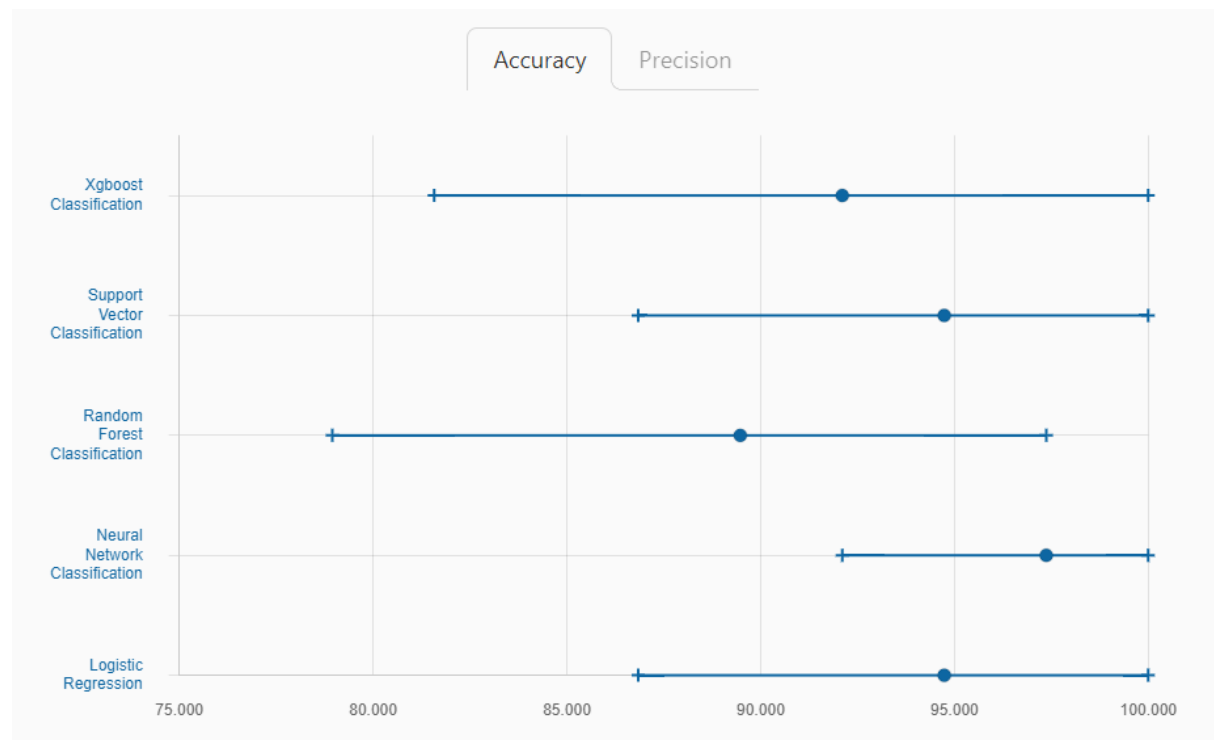


Figure 19. Variables Table Accuracy

The performance of this dataset as per precision of different models is as follows:

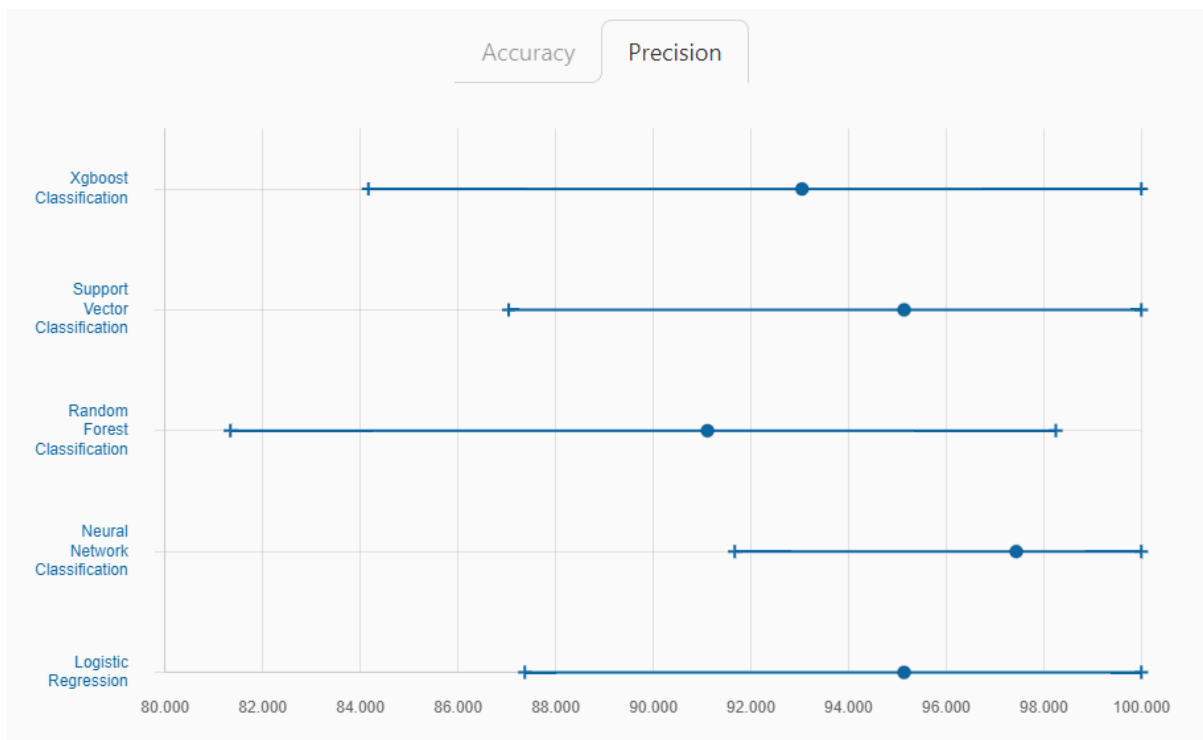


Figure 20. Variables Table Precision

## 2.5 Case Studies on Customer Segmentation in e-commerce using Machine learning

First (Li and Zeying, 2011) proposed a method in which a retail supermarket was used as the research object, and data mining methods were used to segment retail enterprise customer segments, and then association rules obtained using the appropriate algorithm were applied to different groups of customers to obtain rules about customer characteristics in order to efficiently analyze customer characteristics. Finally, the author provided some references to the supermarket's marketing and management efforts, which aided in comprehension. Data mining was utilized effectively to deal with a vast amount of history and current data from the database in order to identify some viable, helpful, and important information for retail establishments that will help us target customers.

Second (Wang, Zhenyu, Yi Zuo, Tieshan Li, 2019) investigated consumer segmentation based on a broad learning system, providing an alternate view of learning in a deep structure. To begin, RFID (Radio Frequency Identification) data was included in addition to customer purchasing behavior, which can correctly represent consumers' in-store behavior. Second, this



article examined consumer segmentation using the Broad Learning System (BLS). (CL Philip Chen, and Katsutoshi Yada,2019) BLS is a superb machine learning algorithm that is both efficient and effective for classification jobs. Third, the customer behavior data used in this work came from a real-world Japanese supermarket. Based on both POS and RFID data, customer segmentation was viewed as a multi-label classification challenge.

Another example will be from (Kansal, Tushar, Suraj Bahuguna, Vishal Singh, and Tanupriya Choudhury, 2018) used K-means clustering to segment customers. A Python software was created, and the program was trained using a standard scaler on a dataset of 200 training samples from a local retail shop. Both features are the annual average of the quantity of shopping done by customers and the annual average of the customer's visits to the shop. Using clustering, 5 cluster segments were generated, labeled as Careless, Careful, Standard, Target, and Sensible customers. However, using mean shift clustering, the authors discovered two new clusters, High purchasers and regular visits and High buyers and occasional visitors.

Lastly, we will look at (Vani Ashok, Rahul R Kamath, Adithya RK, Supreeth Singh, Ajay Bhati, 2021) who used a dataset acquired from the University of California's machine learning repository. The dataset initially had 5,41,909 rows and 8 columns (Invoice No, Stock Code, Description, Quantity, Invoice Date, Unit Price, Customer ID, Country).

The data frame had 4,06,829 rows and 8 columns after deleting all NULL values. It was discovered that as the number of major components increases, so does the loss. Each customer's mean, sum, count, min, max, and percentage of the customer's allocation to each group were calculated.

The diagram below depicts the system architecture used for consumer segmentation. The resulting data set was sent to the pre-processing stage, where duplicate and null rows were removed. A dataset may contain cancelled orders, which are also eliminated. Column datatypes were transformed to the appropriate datatype for pre-processing. The improved data set was then utilized to classify clients into multiple categories using the K-Mean technique for unsupervised learning. (Adithya RK, Supreeth Singh, Ajay Bhati, 2021) The resulting dataset or matrix, which included customers and their categories, was used to train and test supervised training models such Random Forest, K-Nearest Neighbors, and Gradient Boost.

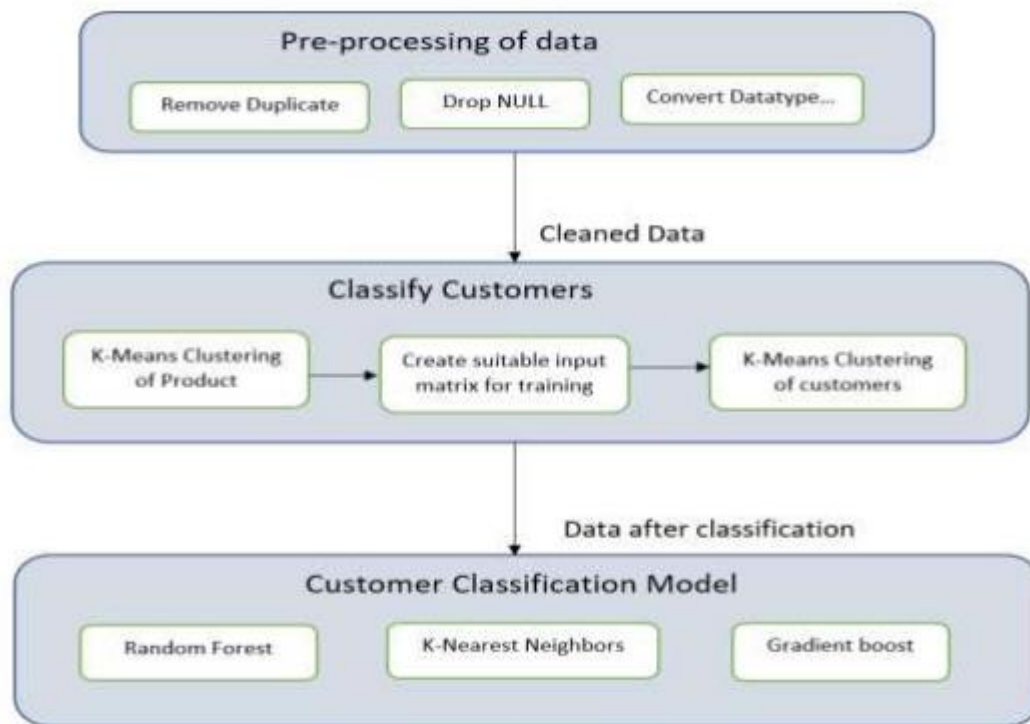


Figure 21. K- Mean Architecture

## 2.6 Synthesis and Analysis

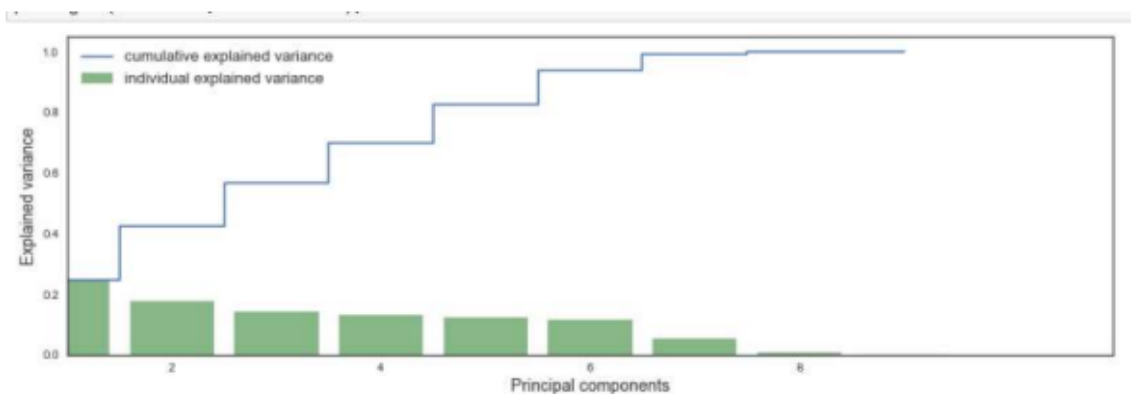


Figure 22. Synthesis and Analysis

From the figure below it was concluded that any clusters are well isolated, while some overlap. At this stage, client clusters were formed using the previously determined standardized matrix, and a silhouette score of 0.218 was obtained using the K-means algorithm.

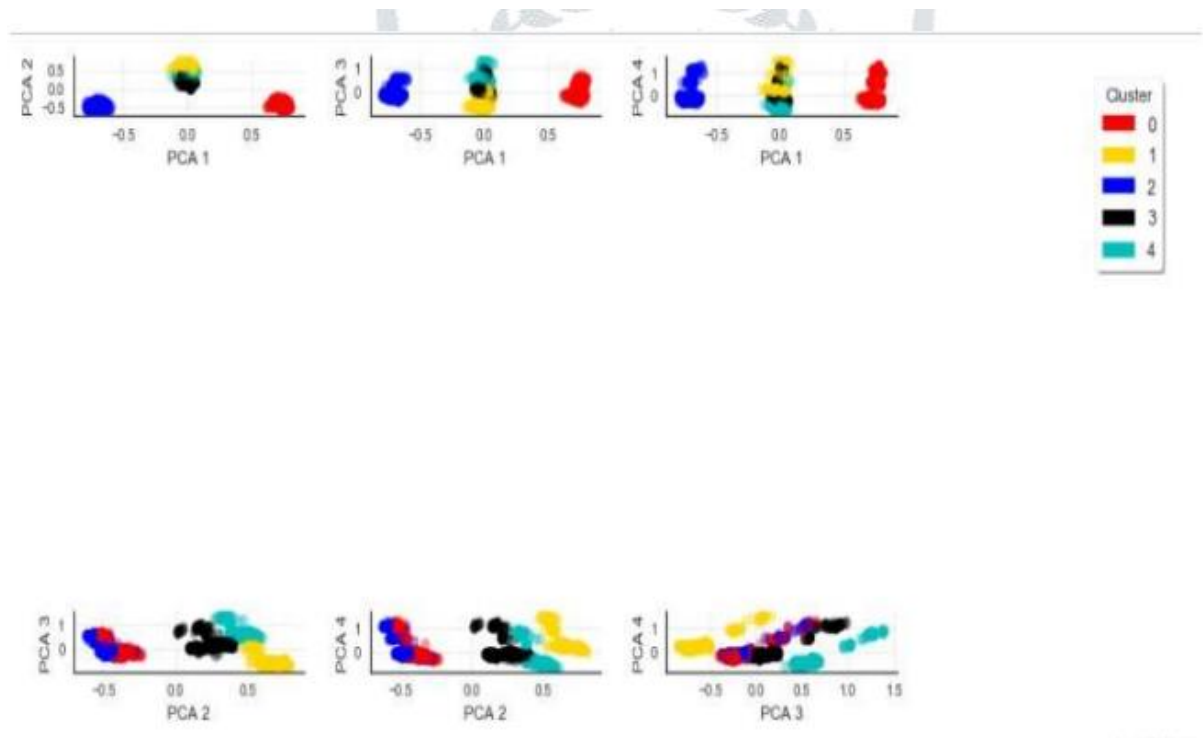


Figure 23. Random Forest learning curve

The figure below shows the results of the learning curve which shows that Random Forest's learning curve used. The learning curve shows that as the number of samples rises, the training score declines but the cross-validation score increases. Gradient Boost training and testing accuracy was 90.17 and 75.77, respectively.

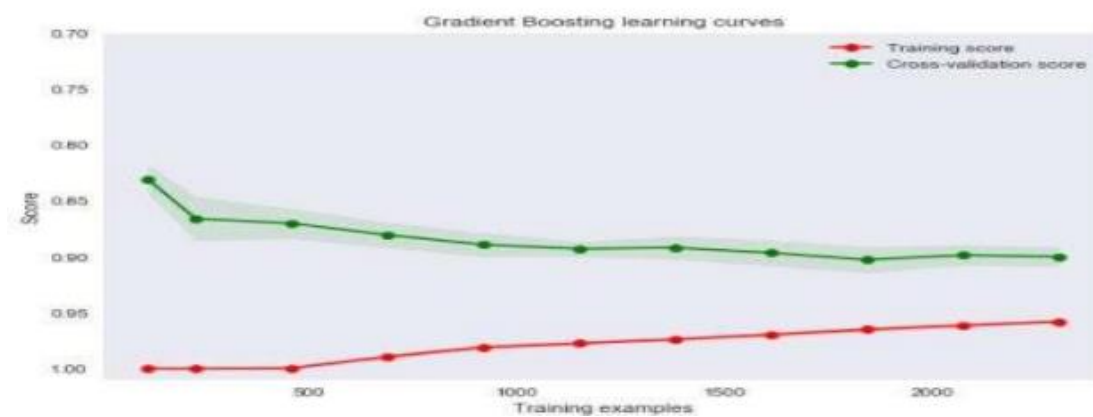


Figure 24. Training Examples

The accuracy of training and testing for the three distinct classifiers. Random Forest and Gradient Boost achieved good accuracy of 75.50 and 75.77, respectively, however KNN produced a relatively low accuracy of 68.29.

**Table 1: Training and testing accuracy**

CLASSIFIER	TRAINING ACCURACY	TESTING ACCURACY
K-Nearest Neighbors (k=5)	81.44	68.29
Random Forest	90.30	75.50
Gradient boost	90.17	75.77

*Figure 25. Training Result*

## **2.7 Gaps in the Literature**

### **2.7.1 Data Quality and Quantity**

Many e-commerce businesses struggle to collect high-quality and adequate data for effective machine learning-based client segmentation. Data that is incomplete or noisy can produce inferior findings. (Richard S. Baskerville 1997). Baskerville has talked about data quality and how it affects many areas of data-driven decision-making, such as consumer segmentation.

### **2.7.2 Algorithm Selection and Validation**

It is crucial to choose the correct machine learning algorithms for client segmentation. However, it is not always obvious which algorithms are best suited for a specific dataset and business scenario. Validating the effectiveness of selected algorithms is also difficult. (Trevor Hastie and Robert Tibshirani, 2009).

### **2.7.3 Interpretable Models**

Machine learning models' interpretability is critical for understanding why particular clients are assigned to specific segments. Many complicated models are difficult to interpret, making it difficult to explain segmentation results to stakeholders. (Christoph Molnar, 2021).

### **2.7.4 Privacy and Ethical Concerns**

Concerns about data privacy and ethical use have arisen as e-commerce corporations acquire massive volumes of user data. It is becoming increasingly difficult to ensure that segmentation approaches respect client privacy and comply to ethical requirements. (Solon Barocas, Woodrow Hartzog, and Edward W. Felten, 2019).

#### **2.7.5 Dynamic Segmentation**

Many e-commerce businesses require dynamic and real-time customer segmentation to adapt to changing customer behaviors. Traditional batch processing methods may not be suitable for such requirements. (Gordon Linoff and Michael J. A. Berry, 2011).

#### **2.7.6 Scalability and Resource Constraints**

Implementing machine learning models for customer segmentation at scale can be time-consuming and resource-intensive, necessitating tremendous processing capacity. Scalability issues may arise for smaller e-commerce enterprises with limited resources. (Andreas C. Müller and Sarah Guido 2016).

#### **2.7.7 Evaluation Metrics**

It is critical to define proper evaluation measures for analyzing the quality of consumer segmentation. The measurements chosen can have an impact on the results and alignment with corporate goals. (Ron Kohavi, 1995).

### **2.8 FUTURE TRENDS IN CUSTOMER SEGMENTATION**

#### **Ephemeral Segmentation:**

Ephemeral segmentation is the temporary classification of customers based on fleeting or short-lived attributes or actions. This enables companies to seize specific, time-sensitive opportunities.

#### **Genetic and Biometric Segmentation:**

Genetic and biometric data will play an increasingly important role in customer segmentation, (Roderic Guigó, 2000) allowing firms to personalize products and services to individuals based on their genetic or biometric profiles.

#### **Dynamic and Agile Segmentation:**

Segmentation models will become more dynamic, continuously adapting to changes in customer behavior. Client segments will become more adaptive to fast changing market

conditions, allowing businesses to quickly pivot in response to client wants. Scott Brinker, 2010)

### **Ethical and Privacy-Centric Segmentation:**

As concerns about data privacy and ethics grow, businesses will need to prioritize responsible data usage. Ethical segmentation practices will focus on obtaining customer consent, ensuring data security, and providing transparent data usage policies. With growing worries about data privacy and ethics, firms will implement segmentation strategies that protect client privacy. Regulations like as (GDPR and the Data Protection act, 2019) for Kenya, will continue to have an impact on these activities.

### **Emotional and Experiential Segmentation:**

Beyond transactional data, businesses will analyze emotional and experiential data to understand how customers feel about their interactions. Emotional and experiential segmentation will concentrate on comprehending and categorizing customers based on their emotional responses and overall brand experiences. Its goal is to create emotionally charged customer interactions. Sentiment analysis and emotional segmentation will become essential for delivering emotionally resonant experiences.

### **Behavioral Micro-Segmentation:**

Behavioral micro-segmentation is the process of creating extremely fine-grained consumer categories based on highly specific behaviors, preferences, and interactions. (William Lazer, 1963). This level of specificity enables highly targeted marketing and personalization.

### **Generational Segmentation**

Businesses will continue to segment customers by generation (Michael Solomon, 2002) (e.g., Gen Z, Millennials, and Baby Boomers) in order to cater to the varied preferences and behaviours of these groups. (Gordon Linoff, 2011).

### **Hyper-Local Segmentation**

Customer segmentation at the hyper-local level will become more important, allowing firms to target clients based on their precise geographic locations and neighbourhoods' features.

### **Cross-Channel Integration**

Businesses will place a greater emphasis on integrating data from numerous consumer touchpoints, both online and offline, to develop a single view of the customer journey. This will improve segmentation accuracy.

### **Personalization at Scale**

Personalization will become more complex and automated, (Paul Marsden, 2006) allowing businesses to offer highly personalized experiences to customers on a broad scale. (Eric Siegel, Pedro Domingos (2000). Machine learning models will be critical in this trend.

### **Real-Time Segmentation**

Real-time segmentation will become increasingly crucial as it allows firms to adapt to changes in client behaviour promptly. The need for quick customisation and relevance is driving this trend.

## **2.9 CONCLUSION**

We have dived into a plethora of variables that define this subject while examining the varied domain of customer segmentation in the context of e-commerce and the deployment of machine learning models. The evolution of client segmentation has been nothing short of astounding, from its historical roots to its present uses.

We started by learning about the fundamentals of consumer segmentation, which were presented to us in the mid-20th century by scholars such as Wendell Smith and Joan Robinson. Robinson's theory of imperfect competition, as well as the recognition of customer non-uniformity, established the framework for today's segmentation techniques. Although modest by modern standards, these early concepts spurred a paradigm shift that has since altered marketing and the e-commerce sector.

Our discussions focused on the incorporation of machine learning models into consumer segmentation. In this field, supervised learning models such as logistic regression, decision trees, and k-nearest neighbors have played critical roles. Scholars such as Jerome Friedman, Leo Breiman, and Ross Quinlan have made lasting contributions to decision trees, while David Aha and Ian Witten have advanced our understanding of k-nearest neighbors. With their different capabilities, these models ushered in a new era of data-driven segmentation, allowing organizations to make data-backed decisions with better precision and efficiency.

We looked into the future of consumer segmentation in our investigation, noting that this subject is constantly growing. Personalization at scale, real-time segmentation, AI-powered segmentation, and ethical and privacy-conscious techniques are altering the environment. While particular scholars for these trends may not be as easily available, their importance cannot be overstated. Researchers, industry executives, and practitioners from a variety of disciplines are actively participating in these revolutionary talks.

Customer segmentation in e-commerce is a complex, ever-changing profession that owes its advancement to early scholars' pioneering insights and the relentless creativity of modern researchers and practitioners. Traditional notions combined with cutting-edge machine learning algorithms have enabled organizations to better understand and communicate with their client base. Looking ahead, consumer segmentation trends offer a more personalized, dynamic, and ethical approach to responding to customers' different requirements. With each innovation, this industry reveals new chances to provide long-term value to organizations and their customers. It exemplifies the never-ending evolution of marketing techniques in the digital age, where the client stays at the center of all decisions.



## **CHAPTER 3**

### **RESEARCH METHODOLOGY**

#### **3.0 INTRODUCTION**

The methodology used in this study is critical in the aim of uncovering the subtle dynamics of client segmentation within the vast domain of e-commerce. This chapter acts as a strategic guide, outlining the investigation's methodical methodology, research design, and ethical considerations. The methodology chapter maintains the precision and rigor of the study by offering a full overview, but it also describes the paths through which the research objectives will be properly handled.

As such this chapter attempts to decipher the varied nature of consumer segmentation through a comprehensive assessment. Understanding these complexities is more than just an intellectual exercise; it has far-reaching ramifications for practical implementations in the e-commerce sector. This study's findings are expected to make a significant contribution to the improvement of marketing tactics and the delivery of individualized services. As a result, the methodology serves as a compass guiding this inquiry, enabling a thorough and rigorous study of client segmentation in the dynamic terrain of e-commerce. The next sections will go into the methodology's specific components, covering the research design, participant selection, data collection methods, analysis tactics, and ethical concerns that contribute to the investigation's robustness and credibility.

#### **3.1 RESEARCH DESIGN**

The research design acts as the architectural blueprint, guiding the systematic and scientific investigation of customer segmentation in the e-commerce area. A deliberate decision was made to use a quantitative strategy in the dynamic world of online customer behavior. This methodological synthesis utilizes quantitative precision for analysis.

##### **3.1.1 Description of research design**

The nature of the study design is at the heart of it. A more extensive and nuanced picture of client segmentation results from the seamless use of quantitative technique. The quantitative component is meticulously analyzing transactional data to identify patterns, trends, and statistical linkages. (Bhandari, 2020).

### 3.1.2 Justification

This quantitative strategy was purposely designed to capitalize on the capabilities of the statistical robustness and generalizability and the conclusions by adding richness and context. In the complex world of e-commerce, where client behavior is influenced by a variety of circumstances, this method ensures a more comprehensive and thorough examination.

Recognizing that client segmentation is not a static phenomenon, the study design is adaptive. It evolves in response to market conditions, technology improvements, and altering consumer expectations. As a result, the chosen quantitative technique corresponds with the dynamic nature of e-commerce by giving a comprehensive lens through which to investigate the intricate interaction of variables driving client segmentation.

## 3.2 SYSTEM DEVELOPMENT METHODOLOGY

A system development methodology refers to the framework that is used to structure, plan and control the process of developing a system. Agile was chosen as the system development methodology for this project because it is a dynamic and iterative approach that can adapt to the changing nature of e-commerce and customer segmentation. In contrast to the structured and sequential Waterfall paradigm, Agile encourages flexibility, collaboration, and rapid iterations to enable continuous improvement throughout the development lifecycle.

### 3.2.1 PHASES OF AGILE DEVELOPMENT

The technique is made up of several fundamental elements, each of which contributes to the systematic creation of the consumer segmentation system:



Figure 26. Agile Methodology

1. **Requirements Analysis:** Perform a detailed analysis of system requirements, drawing insights from research objectives and e-commerce business demands.
2. **Design Phase:** Create a detailed system design, including data models and system architecture. Check to see if the design adheres to the principles of effective customer segmentation.
3. **Development Iterations:** Use of an iterative development approach to allow for continuous refinement and adaptation based on ongoing feedback from stakeholders.
4. **Testing and Quality Assurance:** Perform rigorous testing to confirm the system's dependability, accuracy, and security. To protect system integrity, address any discovered faults as soon as possible.
5. **Deployment and Training:** Deploy the system in stages, providing users and stakeholders with the necessary training. Monitor the system's performance and fix any issues that arise after deployment.
6. **Review:** Implement methods for continuous improvement that allow for updates and enhancements based on new trends, technical breakthroughs, and changing customer habits.

### **3.2.2 AGILE DEVELOPMENT PRINCIPLES**

1. **Iterative Improvement:** It accepts an iterative approach to development, reducing the process down into manageable iterations. Each iteration yields a functional system component, enabling for continuing modifications based on user feedback and evolving requirements.
2. **Communication and Collaboration:** It Encourage developers, data analysts, and end users to work together closely. Regular communication ensures that the system responds to changing market circumstances and matches with evolving business needs.
3. **Participation of Users:** Throughout the development phase, it prioritizes user involvement collection of customer feedback on a regular basis to improve features and functionality and ensure the system is in sync with the actual needs of e-commerce firms.

### **3.3 PARTICIPANTS OR SAMPLE**

In untangling the complex landscape of customer segmentation in the area of e-commerce, the participation or sample selection is a critical factor that determines the depth and breadth of insights gained. The careful selection of a representative sample helps to the study's findings' validity and application.

#### **3.3.1 SAMPLING METHODOLOGY**

Recognizing the heterogeneous nature of the e-commerce client base, a stratified random sample approach was used. This strategy assures that participants are drawn from different strata depending on factors such as age, gender, geographic area, and previous buying history. The research attempts to capture a broad cross-section of the various criteria that influence consumer segmentation by stratifying the sample.

#### **3.3.2 SAMPLE CHARACTERISTICS**

The targeted sample size is 200 people, which is judged sufficient for statistical significance and diversity. This sample will include a diverse range of e-commerce customers with varying demographics and purchasing habits. The deliberate diversity intends to promote a comprehensive assessment of the factors impacting client segmentation and to capture the breadth of perspectives available in the e-commerce industry.

The sample's diversity is critical because it reflects the varied nature of client segmentation in the e-commerce market. This study aims to find trends, preferences, and behaviors that may differ across different client segments using a large sample. The results of such a diverse sample are expected to provide substantial and practical insights that resonate with the dynamic reality of e-commerce markets.

The subsequent sections will elaborate on the data gathering methodology, offering light on the instruments and strategies utilized to collect both quantitative and qualitative data. Moving forward, the emphasis will be on ensuring the inclusivity and representativeness of the sample in order to deepen the investigation of client segmentation in e-commerce.

### **3.4 MODEL DEVELOPMENT**

#### **3.4.1 DATA COLLECTION**

The data collection phase is a critical component of this research endeavor in the aim of thoroughly comprehending client segmentation within the e-commerce market. The approaches and instruments utilized are carefully designed to capture quantitative measurements and ensuring a comprehensive examination of the varied nature of customer behavior.

Methods and Tools: The data collection procedure will be of a quantitative manner:

Quantitative Data Collection: The quantitative arm examines transactional records acquired from the e-commerce platform. This extensive dataset includes purchase history, browsing habits, and demographic information. The research tries to discover patterns and trends indicative of various client segments by exploiting this quantitative data.

Data Validity and Reliability: It is critical to ensure the validity and reliability of the obtained data. A test will be done for quantitative data to develop and validate the survey tools. Before the full-scale data collecting, this method allows for the identification of potential faults and corrections. Data correctness, completeness, and consistency will be closely monitored.

### **3.4.2 DATA PREPROCESSING**

Data preprocessing is a fundamental stage in refining obtained data to ensure accuracy, consistency, and preparation for analysis. Several key processes are used in this phase to improve the quality of the quantitative dataset and lay the framework for useful interpretations.

This involves:

Data Cleaning: Examine transactional records thoroughly to discover and correct errors or inaccuracies. Address any missing variables to ensure the quantitative dataset is full.

Data transformation: Normalize quantitative values to a standardized scale, avoiding any biases produced by different units of measurement. This ensures that different measurements within the quantitative collection are compared fairly.

Missing Data Management: Implementing strong strategies for dealing with missing quantitative data pieces. Depending on the nature of the missing values, employ appropriate imputation techniques to keep the dataset intact.

### **3.4.3 TRAINING THE MODEL**

This next phase in the methodology involves building an unsupervised machine learning model for customer segmentation using the preprocessed quantitative dataset. K-Means clustering, a popular and effective technique, is chosen for its ease of use and effectiveness in detecting unique groups within data.

The following are the steps:

Normalization: Ensure that quantitative features are standardized to a common scale to avoid any one variable dominating the clustering process owing to measurement unit variations.

Determining the Number of Clusters (K): To estimate the best number of clusters (K), use methods such as the elbow method or silhouette analysis. This phase is critical for determining the number of distinct customer segments contained within the dataset.

Feature Selection: Select quantitative features that are relevant for clustering analysis. Consider factors that are indicative of distinct customer habits, such as purchase frequency, average transaction value, and browsing behavior metrics.

Interpretation of Clusters: Analyze the produced clusters to comprehend the model's distinct consumer categories. Investigate the centroids of each cluster to learn about the average quantitative values that define each segment.

Validation and Refinement: Use metrics such as the silhouette score or within-cluster sum of squares to evaluate the quality of the clustering results. If necessary, fine-tune the model parameters to improve the accuracy and significance of the discovered clusters.

Segment Profiling: Create profiles for each client segment based on quantitative criteria. These profiles offer a comprehensive overview of the various behaviors and interests associated with each cluster.

Visualization: We will use visuals to depict the distribution of clients inside each cluster, such as scatter plots and charts. Visualization will aid in communicating clustering results in an understandable manner.

#### **3.4.4 TESTING THE MODEL**

Following the training of the unsupervised K-Means clustering model, the next critical stage in the process is to test and validate its performance. This assures that the consumer categories discovered are useful, dependable, and applicable to real-world settings.

This will entail:

Testing dataset: Preserve a subset of the preprocessed quantitative dataset as a testing dataset. This dataset was not used during the training phase and will be used as an independent sample to assess the model's generalization capacity.

**Trained Model Application:** Apply the trained K-Means clustering model to the testing dataset. Based on their quantitative characteristics, assign each data point in the testing set to one of the previously established clusters.

**Visual inspection:** Examine the distribution of data points in the testing dataset across the selected clusters visually. Visualization helps to see how well the model generalizes to new, previously unseen data.

**Cluster Stability:** Evaluate cluster stability by comparing the results of repeated runs or subsamples. A stable clustering approach produces consistent results over iterations, increasing trust in the dependability of recognized segments.

**Iterative Refinement:** Iteration on the model as needed by modifying parameters or refining features. Continuous testing-based refining increases the model's robustness and ensures its applicability to a wide range of datasets.

### **3.5 ETHICAL CONSIDERATIONS**

Ethical issues are crucial in any research endeavor, guaranteeing the appropriate and courteous treatment of participants as well as the responsible processing of data. Several ethical considerations govern our technique as we investigate client segmentation in the e-commerce arena.

#### **3.5.1 ETHICAL FRAMEWORK**

This organized set of rules or norms that offer a foundation for making ethical decisions and engaging in ethical behavior within a certain context or topic. It provides a foundation for individuals, organizations, and professions to traverse ethical quandaries, make moral decisions, and provide a framework for responsible behavior.

**Informed consent:** Before collecting any data, obtain informed consent from participants. Communicate the research's aim, the type of data collecting, and how the information will be used clearly. Participants should be free to give or withhold consent.

**Data Privacy and secrecy:** Make a strong case for data privacy and secrecy. Assure participants that their personal information and will be anonymized, and that safeguards will be put in place to protect their identity. To preserve sensitive data, use secure and encrypted storage mechanisms.

Optional Participation: Ensure that participation in the research is optional, and that participants can withdraw at any time without penalty. Throughout the study process, respect participants' decisions and autonomy.

Data Handling Transparency: Clearly explain the methods for data handling, storage, and disposal. Transparent data management techniques contribute to the research's integrity and reflect a commitment to ethical behavior.

### **3.6 RESEARCH QUALITY ASSURANCE**

Maintaining the integrity and credibility of client segmentation research in the e-commerce arena is critical. The research quality assurance procedures include a rigorous approach that begins with well stated research objectives that fit with the scope and purpose of the study. The study strategy, which relies mostly on quantitative methodologies like the K-Means clustering algorithm, was chosen for its potential to give valuable insights about client segmentation. To achieve a strong fit with the study objectives, methodological decisions are justified.

Participant selection and sample procedures prioritize representation through stratified random sampling, assuring diversity across demographics and improving the findings' external validity. A carefully justified sample size of 200 participants is chosen to attain statistical significance and adequately capture various customer behaviors.

The outlined strategies work together to maintain high research quality, ensuring that the investigation into customer segmentation in the e-commerce domain yields reliable, credible, and actionable insights that advance our understanding of customer behaviors in this dynamic landscape.

### **3.7 CONCLUSION**

Finally, the methodology part of this research on customer segmentation in the e-commerce area employs a thorough and rigorous quantitative technique to generate strong findings and useful insights. The well-defined research objectives inform the choice of a quantitative research methodology, which predominantly use the unsupervised K-Means clustering technique, which is well-known for its efficiency in identifying unique client segments. The careful evaluation of participant selection and sampling procedures, such as stratified random sampling with a justified sample size of 200 participants, ensures the diversity and statistical significance required for a thorough understanding of customer behaviors.



The data gathering procedure combines quantitative transactional data with structured dataset to provide a comprehensive picture of customer interactions with the e-commerce platform. To improve the quality and integrity of the quantitative dataset, rigorous data preprocessing procedures such as resolving missing data, outliers, and normalizing are used. Systematic procedures, such as feature selection and normalization, are used throughout the model training phase to ensure that the unsupervised K-Means clustering model captures meaningful patterns in consumer behavior.

Throughout the study process, ethical considerations are integrated, with a focus on informed permission, data privacy, and continual ethical monitoring. The testing and validation phase thoroughly evaluates the trained model's performance on an independent dataset using evaluation metrics and visual inspections. The emphasis on continual quality monitoring, including progress reviews and peer review processes, adds to the research's dependability and trustworthiness.

## **CHAPTER 4**

### **SYSTEM ANALYSIS**

#### **4.0 INTRODUCTION**

System analysis is a vital phase in the software development lifecycle, where the complexities of a proposed software solution are rigorously investigated. It entails conducting a thorough analysis of existing processes, identifying user requirements, and assessing business needs in order to provide a full grasp of the project's scope. The importance of system analysis stems from its capacity to create the framework for educated decision-making, ensuring that succeeding development phases are aligned with corporate goals and user expectations. The key goals of this phase are to define the system's boundaries, understand its functional needs, and lay a solid basis for the succeeding design and implementation stages.

System analysis, as a prelude to the design and development stages, sets the tone for the entire software development process. It acts as a link between the conceptualization of a software solution and its practical execution, guiding the development team in developing a system that not only fits user needs but also aligns with the organization's broader goals. System analysis offers the framework for a successful and well-informed software development journey by giving a disciplined way to comprehending the intricacies of the problem domain. The next sections dig into the methodologies used during system analysis, covering everything from data gathering approaches to the validation and verification of requirements acquired.

#### **4.1 DATA COLLECTION METHODS**

Effective data collection is essential for successful system analysis because it serves as the foundation for comprehending the complexities of present operations and requirements. To acquire information extensively, several procedures and techniques are used, ensuring that the analysis phase is built on a firm basis. The following methods of data collecting are essential to the system analysis process:

##### **4.1.1 PREVIOUS TRANSACTION ANALYSIS**

Examining previous transactions can help you comprehend the historical background and complexities of present procedures. The document analysis supplements other methods of data collecting by giving additional context and insights of what type of segmentation can be undertaken within the dataset provided.

## **4.2 GATHERING USER REQUIREMENTS**

Gathering user requirements is a critical component of system analysis, focused on understanding and documenting end-user needs, expectations, and preferences. This phase entails a comprehensive investigation of user interactions with the proposed system, with the goal of capturing a detailed and accurate depiction of their requirements. Several strategies are used in this procedure:

### **4.2.1 STAKEHOLDERS AND THEIR NEEDS**

Identifying and engaging with all key stakeholders is a critical step in acquiring user requirements. The stakeholders for this project will be:

- Managers
- Sales and Marketing team
- Supplier

The requirements for the stakeholders above are but not limited to:

#### **1. Insights into Customer Segmentation**

Managers may need a system that gives extensive information on customer segments. Tools for assessing client behaviors, preferences, and purchase habits are included, allowing managers to adjust marketing tactics based on segmented consumer profiles.

#### **2. Efficacy of Marketing Campaigns:**

It is critical to have a system in place that measures the performance of marketing initiatives. Managers may require capabilities that allow them to assess campaign performance, analyze client responses, and alter plans to maximize return on investment (ROI).

#### **3. Analytics in Real Time:**

Managers may require real-time information to track ongoing marketing campaigns. Real-time insights enable marketing strategies to be adjusted in real time, allowing managers to respond swiftly to changing market conditions.

#### **4. Customer Satisfaction and Feedback:**

Managers who want to improve customer satisfaction must collect and analyze consumer feedback. A system should contain mechanisms for gathering consumer feedback, performing surveys, and gauging overall customer satisfaction.

#### **5. Product Performance and Inventory:**

Managers may require elements relating to inventory management and product performance. Tools for tracking product popularity, maintaining inventory levels, and adjusting product offerings based on customer preferences are all included.

#### **4.3 FEASIBILITY STUDY**

The feasibility analysis phase of the system analysis process evaluates the practicality of the proposed software solution from a variety of angles. Conducting a thorough feasibility analysis is critical in the context of our e-commerce client segmentation project to assure the practicality and success of the suggested software solution. To address our project's particular requirements and constraints, the assessment will be undertaken across many dimensions:

##### **4.3.1 TECHNICAL FEASIBILITY**

We will examine the compatibility of the client segmentation software with the existing e-commerce technology stack in our technical feasibility analysis. It will determine whether the proposed solution interfaces seamlessly with the existing infrastructure, ensuring that existing technologies and resources are efficiently leveraged. Furthermore, the investigation will assess the technical capability to deploy advanced segmentation algorithms as well as data processing capabilities.

##### **4.3.2 OPERATIONAL FEASIBILITY**

In the e-commerce market, where seamless operations directly affect client experiences, operational feasibility is critical. To verify that the proposed system complies with existing workflows, the analysis will focus on user acceptance and flexibility. It will also take into account any delays throughout the implementation phase, emphasizing a phased approach to reduce operational effect and ensure a smooth transition.

##### **4.3.3 ECONOMIC FEASIBILITY**

Economic feasibility will entail a complete cost-benefit analysis specific to the e-commerce domain. The expenses of creating, installing, training, and maintaining the customer segmentation system will be estimated. Concurrently, the advantages will be quantified, taking into account prospective revenue growth, improved marketing effectiveness, and cost savings through focused initiatives. This economic analysis seeks to demonstrate a favorable return on investment in line with the goals of the e-commerce firm.

#### **4.4 DETERMINING PROJECT VIABILITY**

Determining the profitability of a consumer segmentation project in the e-commerce area entails a thorough examination of the feasibility analysis findings and the technical needs provided. Based on the insights gathered, this stage tries to make informed judgments on

whether the project should progress, be adjusted, or reconsidered. The following are important factors to consider while considering project viability:

#### Adaptability to Market Dynamics:

Determining the profitability of a consumer segmentation project in the e-commerce area entails a thorough examination of the feasibility analysis findings and the technical needs provided. Based on the insights gathered, this stage tries to make informed judgments on whether the project should progress, be adjusted, or reconsidered. The following are important factors to consider while considering project viability

#### Input and approval from stakeholders:

It is critical to solicit feedback from important stakeholders such as marketing teams, IT departments, and decision-makers. Stakeholder feedback is taken into account when determining whether the suggested consumer segmentation solution fits their expectations and needs.

#### Alignment with Business Goals:

It is critical to ensure that the project corresponds with the e-commerce platform's broader commercial objectives. The viability of a project is closely related to its capacity to contribute effectively to marketing initiatives, improve customer experiences, and align with the organization's overall goals.

#### Resource Availability:

Assessing the availability of resources, both personnel and financial, is critical for project viability. Ensure that the relevant talents are accessible and that the project can be completed within budget restrictions. This contributes to the project's overall feasibility.

## **4.5 SOFTWARE REQUIREMENT SPECIFICATION**

The Software Requirement Specification (SRS) defines the functional and non-functional requirements required for the effective development and implementation of e-commerce consumer segmentation software.

### **4.5.1 FUNCTIONAL REQUIREMENTS**

Functional requirements specify the features, capabilities, and behaviors the software system must have. They specify what the system should do in terms of input, processing, and output.

They include:

Input of Customer Data:

The system should be able to accept consumer data from a variety of sources, such as purchase history, browsing behavior, and demographic information.

Algorithm for Segmentation Implementation:

The program should use advanced machine learning algorithms, such as clustering K-Means, which will be used to categorize customers based on predefined criteria.

Generation of Segment Profiles:

Create detailed profiles for each client category, detailing demographic information, preferences, and habits.

Integration of Marketing Strategy:

Allow for the integration of client segments with focused marketing techniques, enabling individualized campaigns tailored to each segment.

#### **4.5.2 NON-FUNCTIONAL REQUIREMENTS**

Non-functional requirements specify the characteristics or features that characterize how well the system performs specific functions. They define the overall qualities and restrictions to which the system must conform.

They include:

Performance:

The system should be fast, with the ability to handle big datasets and provide real-time answers during segmentation and data processing.

Security

To protect consumer information, ensure strong security measures such as data encryption, access controls, and compliance with data protection requirements.

Usability:

The software should have an intuitive and user-friendly that makes it easy to use for both technical and non-technical users.

User-Friendliness:

The application's user interface should be straightforward and in accordance with accepted web design principles, providing a seamless user experience. It should also be responsive, meaning it should adjust to multiple screen sizes and resolutions.

#### **4.6 CONCLUSION**

System Analysis is a detailed tour through the fundamental steps of the software development lifecycle. The chapter began with an excellent introduction that defined the essence and importance of system analysis. It stressed that system analysis is a strategic activity that sets the framework for the entire development process, not just a procedural step.

The following sections went over several strategies and techniques for gathering relevant information during the analysis phase. The armory of data collection methods, ranging from interviews and surveys to workshops and observation, emphasized the significance of a nuanced approach to gathering insights. This comprehensive toolbox ensures a comprehensive understanding of user requirements, laying the groundwork for succeeding stages of development.

The chapter closed with a summary of significant principles, emphasizing the need of precise and thorough analysis in the success of software development. It underlined that the rigorous groundwork created during the system analysis phase serves as the foundation for the entire software development architecture. The emphasis on precision, stakeholder participation, and alignment with organizational goals assures the success of following development phases. In essence, Chapter 4 is a strategic guidepost that directs the software development path toward efficacy and alignment with user and organizational goals.

## **CHAPTER 5**

### **SYSTEM DESIGN**

#### **5.0 INTRODUCTION**

The system design phase is a critical stage in the software development lifecycle, during which the hypothesized software solution begins to assume physical form. This section provides a thorough explanation of the relevance and goals of the system design process. System design, at its heart, is a methodical way to transforming user needs and system functionality into a well-defined and executable blueprint. It establishes the framework for following phases of development by ensuring that the envisioned software system not only fits user expectations but also adheres to key concepts of efficiency, scalability, and maintainability.

In the intricate dance between user needs and system capabilities, the system design phase plays a central role in bridging these realms. It serves as the nexus where abstract concepts are translated into concrete structures, offering a framework that not only meets immediate user requirements but also anticipates the evolving needs of the software throughout its lifecycle. Through the lens of system design, the chapter explores how the design process serves as a linchpin, aligning the envisioned software system with the dynamic interplay of user expectations and technical feasibility. As we delve into the subsequent sections, ranging from use case design to database implementation, the overarching theme remains the translation of abstract ideas into tangible, functional components. The narrative unfolds to elucidate the intricate methodologies and techniques employed during the system design phase, elucidating their role in ensuring a robust, scalable, and user-centric software solution.



## 5.1 SYSTEM ARCHITECTURE

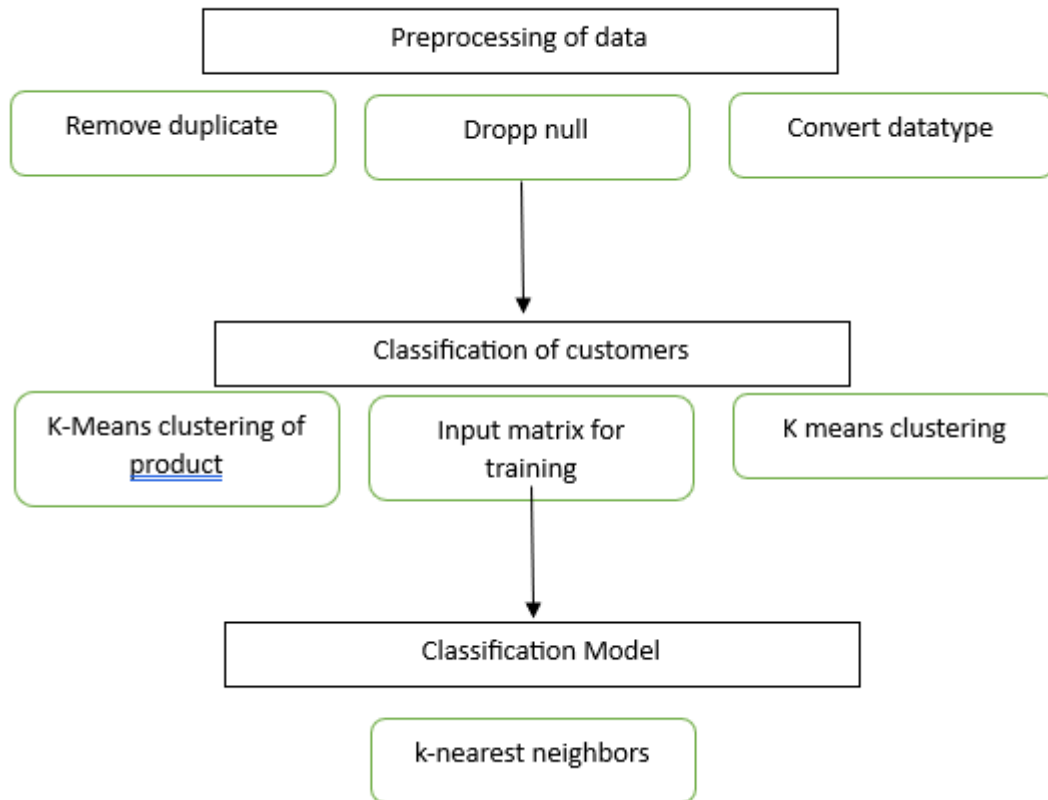


Figure 27. System Architecture

## 5.2 SYSTEM FLOW CHART

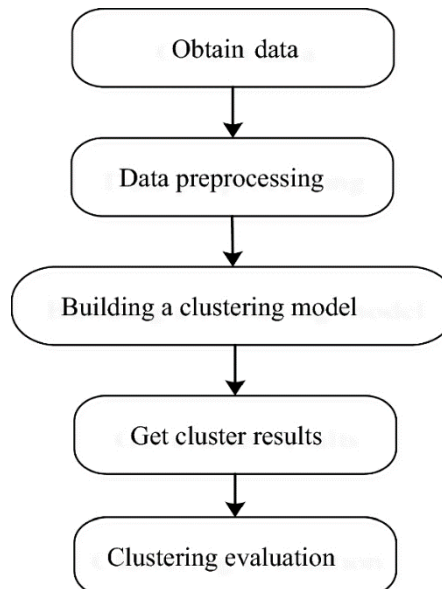


Figure 28. System Flowchart

### 5.3 SYSTEM WORKFLOW

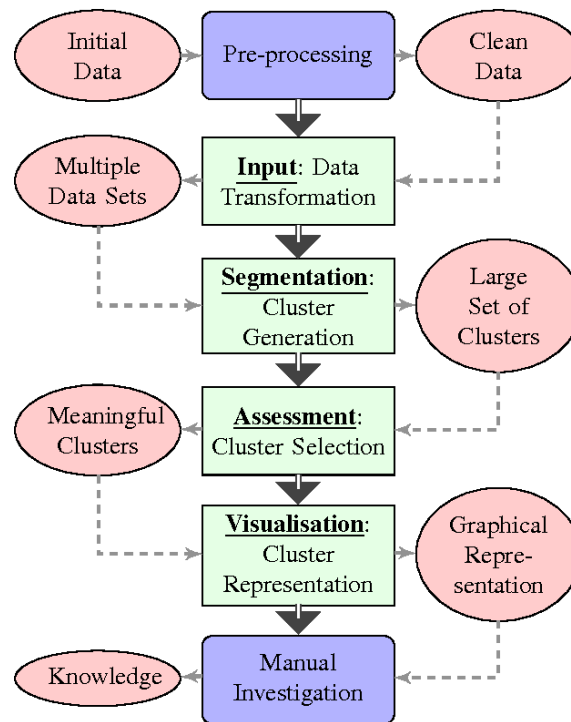


Figure 29. System Workflow

### 5.4 K-CLUSTERING WORKFLOW IN CUSTOMER SEGMENTATION

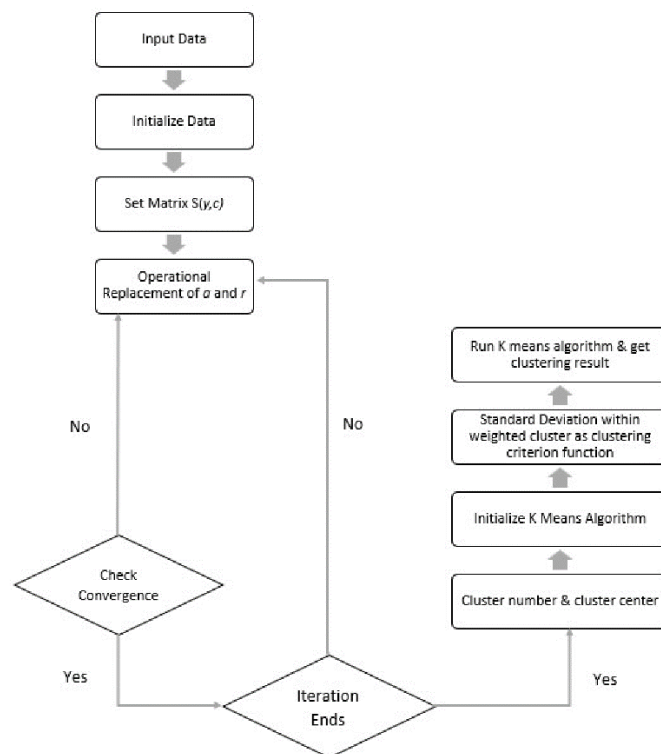


Figure 30. K-Clustering Workflow

## 5.5 SEQUENCE DIAGRAM

A sequence diagram is a form of interaction diagram that displays how different actors or objects in a system interact over time. It depicts the flow of messages or events between various entities. Vertical lines represent the lifelines of participants (actors or objects) in a sequence diagram, while horizontal arrows represent the messages transferred between them.

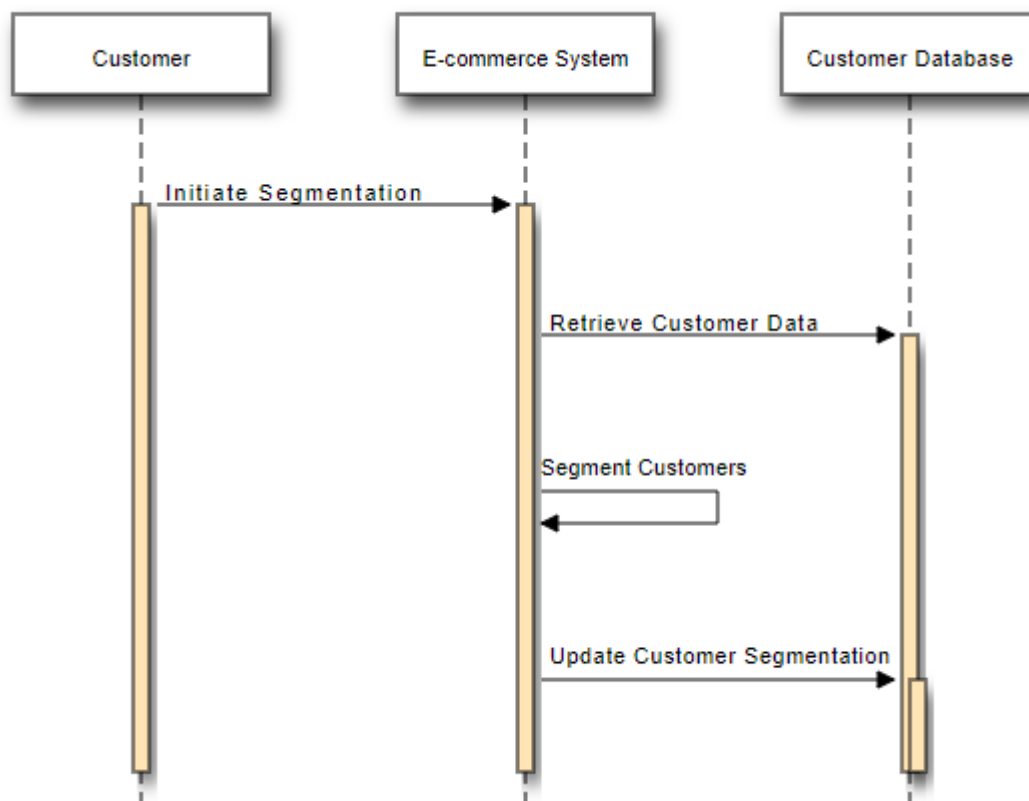


Figure 31. Sequence Diagram

The steps are:

Initiation by the customer:

The sequence starts with the customer starting the segmentation process by a purchase from a client's product line.

Data Extraction:

Data is then extracted from the dataset provided by the client e.g supermarket.

Segmentation of customers:

The system evaluates the received data and Segment Customers indicating the customer segment in which a customer belongs..

Results of Segmentation Update:

The result of segmentation is then updated adding where the client belongs

## **5.6 CONCLUSION**

This chapter describes the architectural design of an e-commerce customer segmentation system that uses k-means clustering techniques to precisely categorize various consumer behaviors. The components, interactions, and dynamics of the implemented system are described, modeled with Data Flow Diagrams (DFD), and illustrated with a use case diagram and a system sequence diagram.

The system architecture is made up of three main parts: segmentation, profiling, and targeting. The segmentation process is driven by K-means clustering techniques, which allow the identification of discrete client segments. Data Flow Diagrams depict the smooth flow of information inside the system, demonstrating how data goes through the processes of segmentation, profiling, and targeting. To clarify step-by-step operations, flow charts showing the system's training and testing/predicting phases are shown. The sequence diagram depicts the interactions between users and the segmentation system in detail.

## **CHAPTER 6**

### **SYSTEM IMPLEMENTATION AND TESTING**

#### **6.0 INTRODUCTION**

The essential phase of system implementation is introduced in this chapter, illustrating the pivotal transition from well written designs to the actualization of a functional system. The importance of this phase rests in the tangible transformation of theoretical blueprints into practical components, during which codes are created, databases are populated, and algorithms are brought to life. This section provides a complete description of the objectives during implementation, with an emphasis on translating design blueprints into workable code and validating the practicality of the proposed solution. A primary goal is the smooth integration of varied components, ensuring the implementation of a coherent system consistent with the envisioned architecture. The story emphasizes the precise approaches used as we progress through the following sections, which focus on coding practices, database population, and algorithm execution.

#### **6.1 TOOLS FOR IMPLEMENTATION**

##### **6.1.1 PROGRAMMING LANGUAGE**

###### **A) PYTHON**

Python's simplicity, readability, and accessibility make it ideal for developing a consumer segmentation system in e-commerce. Its vast libraries, such as NumPy and scikit-learn, provide sophisticated data analysis and machine learning tools, which are critical for effective segmentation. The versatility of the language in interacting with databases and other services provides easy data integration, and its rapid prototyping features allow for rapid testing and iteration in the dynamic e-commerce context. Overall, Python's distinct features make it an excellent choice for creating a comprehensive and adaptive customer segmentation system.

###### **B) ANACONDA IDE**

Anaconda IDE is ideal for developing an e-commerce consumer segmentation system because it provides a user-friendly environment, efficient package management, and interoperability with important data science libraries such as NumPy and scikit-learn. Jupyter Notebooks facilitate iterative algorithm development, and conda environments assure constant project dependencies. Anaconda's collaboration and documentation tools improve

teamwork and information sharing, making it an excellent choice for creating a complex customer segmentation system within the limits of e-commerce.

### C) JUPYTER NOTEBOOK

Jupyter Notebook, a strong open-source tool, is ideal for creating an e-commerce customer segmentation system. Its interactive computing environment enables seamless data and algorithm experimentation, giving a platform for rapid exploration and analysis. Jupyter Notebooks, which support a variety of computer languages, including Python, allow for the construction and iterative refining of segmentation algorithms. The visualization and documentation elements improve the interpretability of results and foster collaboration. Furthermore, Jupyter Notebooks interface smoothly with data science libraries, making it easier to include crucial tools like NumPy and scikit-learn. Overall, Jupyter Notebook's versatility and interactivity make it an excellent asset in the development and refinement of customer segmentation strategies in the ever-changing landscape of e-commerce.

#### 6.1.2 PACKAGES

In our project we used following packages:

- Pandas (version : 1.1.5)
- Numpy (version : 1.19.2)
- Matplotlib (version : 3.3.2)
- Scikit Learn (version : 0.23.2)
- Seaborn (version : 0.11.1)

Pandas:

Pandas is a software library written for the Python programming language for data manipulation and analysis. In particular, it offers data structures and operations for manipulating numerical tables and time series.

Numpy:

NumPy is the fundamental package for scientific computing in Python. It is a Python library that provides a multidimensional array object, various derived objects (such as masked arrays and matrices), and an assortment of routines for fast operations on arrays,

including mathematical, logical, shape manipulation, sorting, selecting, I/O, discrete Fourier transforms, basic linear algebra, basic statistical operations, random simulation and much more.

Matplotlib:

Matplotlib is a plotting library for the Python programming language and its numerical mathematics extension NumPy. It provides an object-oriented API for embedding plots into applications using general-purpose GUI toolkits.

Scikit Learn:

Scikit-learn (Sklearn) is the most useful and robust library for machine learning in Python. It provides a selection of efficient tools for machine learning and statistical modeling including classification, regression, clustering and dimensionality reduction via a consistency interface in Python

Seaborn:

Seaborn is a library for making statistical graphics in Python. It builds on top of matplotlib and integrates closely with pandas data structures. Seaborn helps you explore and understand your data. Its plotting functions operate on data frames and arrays containing whole datasets and internally perform the necessary semantic mapping and statistical aggregation to produce informative plots.

## **6.2 TESTING**

Robust testing and quality assurance techniques are required in the implementation of the consumer segmentation system for e-commerce to assure the system's reliability and operation prior to deployment. The testing methodologies used cover a wide range of levels in order to thoroughly analyze various aspects of the system.

### **6.2.1 UNIT TESTING**

Unit testing is performed at the fundamental level to assess the functionality of individual components or modules in isolation. This guarantees that each unit of the system functions as planned, finding and correcting any issues as they arise.

### 6.2.2 INTEGRATION TESTING

Integration testing is used to assess the interactions and interfaces of interconnected components. This ensures that when the pieces are integrated, they work flawlessly and meet the design and architectural standards.

### 6.2.3 SYSTEM TESTING

System testing examines the behavior of the system as a whole in the context of the full e-commerce environment. This phase validates the customer segmentation system's interaction with the current infrastructure and assures overall system functionality.

### 6.2.4 User Acceptance Testing (UAT)

User acceptance testing entails using real-world scenarios to determine whether the system meets the user's expectations. This phase allows stakeholders to assess whether the system satisfies their expectations and requirements, ensuring that it is in line with company goals.

## 6.3 SCREENSHOTS OF THE SYSTEM

now, we need to visualize the data which we are going to use for the clustering. This will give us a fair idea about the data we're working on.

```
[4]: fig, ax = plt.subplots(figsize=(15,7))
sns.set(font_scale=1.5)
ax = sns.scatterplot(y=clustering_data['Spending_Score'],x=clustering_data['Annual_Income'], s=70, color='#f73434', edgecolor='b')
ax.set_ylabel('Spending Scores')
ax.set_xlabel('Annual Income (in X 1000 KSH)')
plt.title('Spending Score per Annual Income', fontsize = 20)
plt.show()
```



Figure 32. Screenshot 1



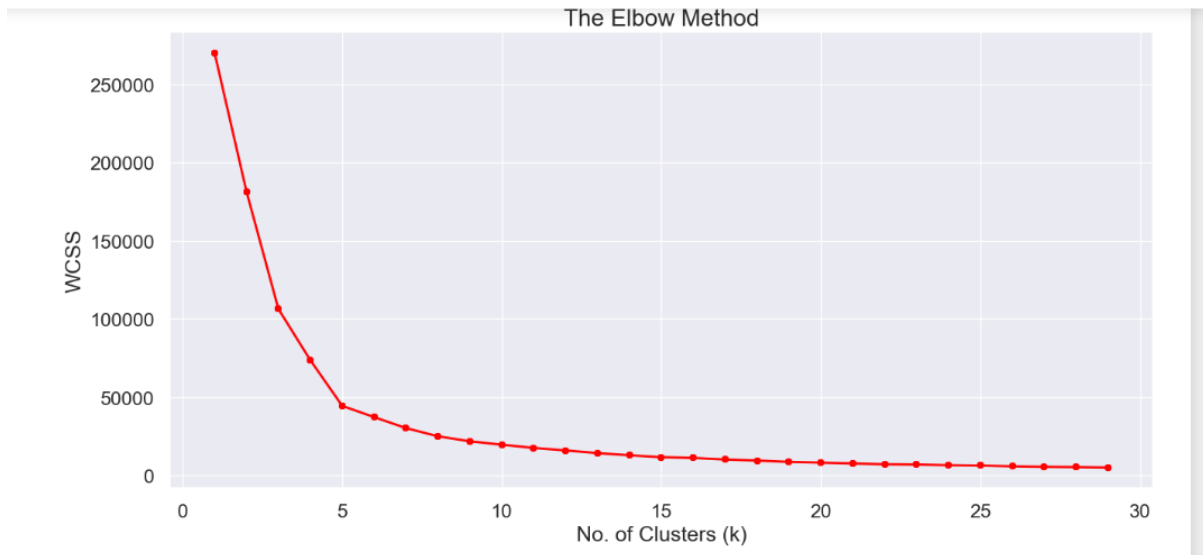


Figure 33. Screenshot 2

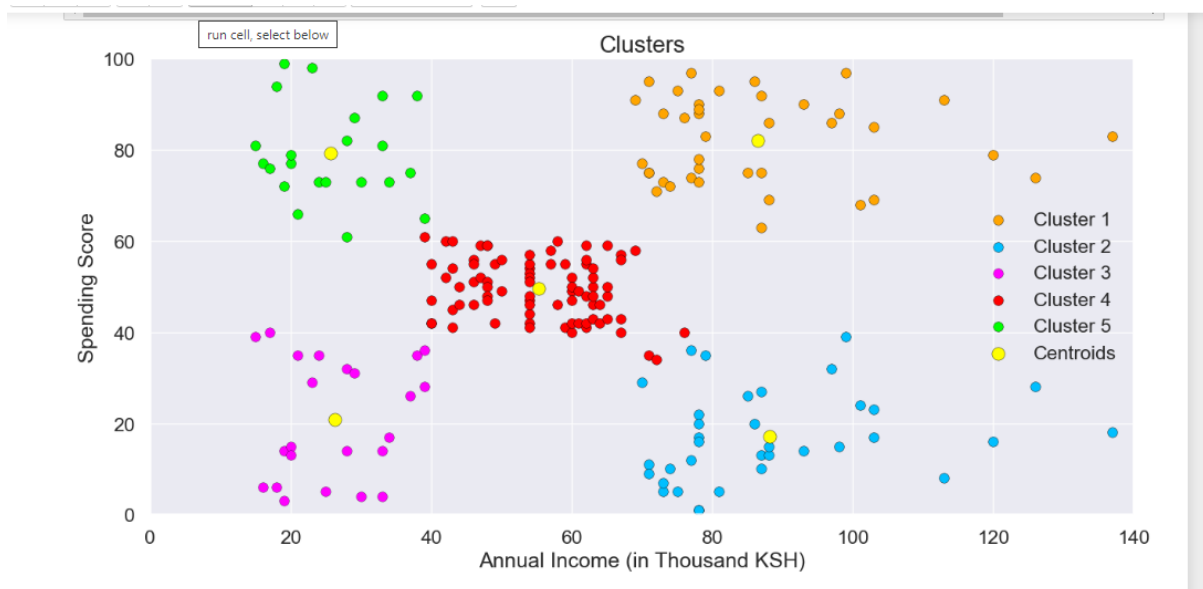


Figure 34. Screenshot 3

## 6.4 CONCLUSION

This chapter delineates the architectural design of a customer segmentation system in e-commerce, leveraging k-clustering for precise categorization of diverse consumer behaviors. The elucidation covers the key components, their interactions, and the user-system dynamics. Modeled through Data Flow Diagrams (DFD), the system comprises three pivotal components: segmentation, profiling, and targeting. K-clustering algorithms drive the segmentation process, offering a robust framework for grouping customers based on shared characteristics. The DFD illustrates the seamless flow of information within the system, showcasing how data moves through the segmentation, profiling, and targeting stages. Additionally, flow charts detailing the training and testing/predicting phases of the system are

presented, elucidating the step-by-step processes. The system's sequence diagram offers a comprehensive view of the interactions between users and the segmentation system. This customer segmentation system holds immense potential for applications in personalized marketing, enhanced user experience, and targeted campaigns. Although the system is currently in the developmental phase, preliminary results indicate its efficacy in optimizing marketing strategies and tailoring services to diverse customer segments, underscoring its promise in the evolving landscape of e-commerce.

## **CHAPTER 7**

### **CONCLUSION AND RECCOMENDATION**

#### **7.1 CONCLUSION**

This project introduces an innovative approach for customer segmentation in e-commerce using k-clustering techniques. The objective is to enhance targeted marketing strategies and improve personalized user experiences by categorizing diverse consumer behaviors. The project utilizes a dataset comprising various customer attributes relevant to online shopping behavior. The primary focus is on implementing k-clustering algorithms for efficient segmentation, allowing for the identification of distinct customer segments. The methodology involves preprocessing the customer dataset, applying k-clustering for segmentation, and refining the classification of customer segments.

Five and Three different segments were identified based on the clustering analysis, representing distinct patterns in online shopping behavior. The evaluation of the proposed algorithm's performance employed metrics such as Precision. The dataset was divided into training and testing samples to assess the algorithm's robustness. Precision results demonstrated the effectiveness of the k-clustering approach in achieving accurate and meaningful customer segmentation.

While the project has shown encouraging results, acknowledging the complexity of customer segmentation in e-commerce, future efforts will focus on refining the algorithm and incorporating additional features to enhance segmentation accuracy. The ultimate aim is to deploy a robust customer segmentation system that significantly contributes to optimizing marketing strategies and tailoring services for diverse customer segments in the dynamic landscape of e-commerce.

#### **7.2 FUTURE RECCOMENDATIONS**

To begin, investing in advanced machine learning techniques, such as deep learning, can improve segmentation accuracy and predictability for specific client preferences. Real-time segmentation is critical for quickly adapting to changing client behaviors, necessitating solutions that dynamically modify segments based on the most recent data. Predictive analytics integration is critical for anticipating future client needs and exploiting past data and market trends for proactive marketing tactics. An omni-channel segmentation strategy guarantees that the consumer journey is consistent and personalized across all touchpoints. To sustain customer trust, ethical data usage should be prioritized, including transparent policies

and compliance with data protection rules. Integrating customer feedback into segmentation methods yields significant insights, and ongoing improvement of behavioral segmentation with new factors is required.

Implementing an experimentation and A/B testing culture enables iterative optimization of segmentation models. Prioritizing customer segments based on lifetime value encourages long-term partnerships and revenue growth. Finally, encouraging collaboration between AI and marketing specialists enables a comprehensive approach that combines technological expertise with business acumen for the effective deployment of advanced segmentation strategies. These ideas aim to move client segmentation in e-commerce into a more intelligent, adaptable, and morally sound future.

## REFERENCES

1. Li, Zeying. "Research on customer segmentation in retailing based on clustering model." In 2011 International Conference on Computer Science and Service System (CSSS), pp. 3437-3440. IEEE, 2011.
2. Kamakura, W. A., Mela, C. F., A. Bodapatti, Fader, P. S., Iyengar, R., Naik, P. A., Neslin, S. A., Sun, B., Verhoef, P. C., Wedel, M., & Wilcox, R. T. (2005). Choice Models and Customer Relationship Management. *Marketing Letters*, 16(3-4), 279–291.
3. Mishra, S.K.; Dwivedi, V.; Sarvanan, C.K.; Pathak, K. Pattern Discovery in Hydrological Time Series Data Mining during the Monsoon Period of the High Flood Years in Brahmaputra River Basin. *Int. J. Comput. Appl.* 2013, 67, 7–14
4. Rajas Sanjay Ubhare. (2020, August 11). Telecom Industry Customer Churn Prediction with K Nearest Neighbor. Medium; Chatbots Life. <https://chatbotslife.com/telecom-industry-customer-churn-prediction-with-k-nearest-neighbor-1d5784952c45>
5. (Frank, Massy, Wind 1972, McDonald & Dunbar 2004, Anna-Lena 2001, Jiang and Tuzhilin, 2006). Predicting customer quality in e-commerce social networks: a machine learning approach. *Review of Managerial Science*, 13(3), 589–603.  
<https://doi.org/10.1007/s11846-018-0316-x>
6. Chien, T. W., & Tsaur, S. H. (2007). Exploring the impact of customer relationship management on customer satisfaction in e-commerce. *Total Quality Management & Business Excellence*, 18(4), 387-398.
7. Xie, C., Li, X., Bao, Z., & Huang, L. (2018). A novel customer segmentation model based on RFM and LRFMC methods. *Sustainability*, 10(5), 1476.
8. Bertrand, J. W., & Franses, P. H. (2003). The service paradox and endogenous asymmetric quality. *Journal of Service Research*, 5(4), 302-311.
9. Nguyen, T. H. (2018). Customer segmentation: A review. *Journal of Retailing and Consumer Services*, 43, 166-177.
10. Borle, S., Boatwright, P., Kadane, J. B., & Manchanda, P. (2005). The effect of customer returns on retailers' pricing and ordering strategies. *Quantitative Marketing and Economics*, 3(4), 355-376.
11. Verhoef, P. C., Franses, P. H., & Hoekstra, J. C. (2002). The effect of relational constructs on customer referrals and number of services purchased from a

- multiservice provider: Does age of relationship matter? *Journal of the Academy of Marketing Science*, 30(3), 202-216.
12. Hennig-Thurau, T., Gwinner, K. P., Walsh, G., & Gremler, D. D. (2004). Electronic word-of-mouth via consumer-opinion platforms: What motivates consumers to articulate themselves on the internet? *Journal of Interactive Marketing*, 18(1), 38-52.
  13. Wedel, M., & Kamakura, W. A. (2000). *Market segmentation: Conceptual and methodological foundations* (Vol. 8). Springer Science & Business Media.
  14. Han, J., Kamber, M., & Pei, J. (2012). *Data mining: Concepts and techniques* (3rd ed.). Elsevier.
  15. Schmitt, P., Böckenholt, U., & Van den Bulte, C. (2010). The analysis of segmentation tables: Why latent class analysis outperforms better known cluster analysis procedures. *Journal of Business Research*, 63(7), 714-724.
  16. Gupta, S., Lehmann, D. R., & Stuart, J. A. (2004). Valuing customers. *Journal of Marketing Research*, 41(1), 7-18.
  17. Panigrahi, R., & Sahoo, B. (2015). A systematic literature review on customer relationship management: A comprehensive update. *Decision Science Letters*, 4(2), 195-210.
  18. Brink, A., & Berndt, A. (2017). Customer segmentation in an online business-to-business network. *Southern African Business Review*, 21(1), 98-119.
  19. Kim, D., & Yoo, C. (2015). Segmenting customers based on values and behaviors in mobile commerce: A clustering approach. *Telematics and Informatics*, 32(1), 45-55.
  20. Sun, T., & Kim, S. (2013). Customer segmentation and strategy development based on customer lifetime value: A case study. *Expert Systems with Applications*, 40(14), 5621-5628.
  21. Rabinovich, E., & Bailey, R. (2015). Examining the role of business analytics in firm performance: Evidence from the field. *Journal of Business Research*, 68(12), 2668-2675.
  22. Han, J., Kamber, M., & Pei, J. (2011). *Data mining: Concepts and techniques*. Morgan Kaufmann.
  23. Jain, A. K., Murty, M. N., & Flynn, P. J. (1999). Data clustering: A review. *ACM Computing Surveys (CSUR)*, 31(3), 264-323.
  24. Hair, J. F., Black, W. C., Babin, B. J., & Anderson, R. E. (2019). *Multivariate data analysis* (8th ed.). Cengage Learning.

25. Wu, P. F., & Li, H. (2018). The impact of customer satisfaction on customer loyalty and firm profitability: Findings from China's telecommunications industry. *Management Decision*, 56(2), 443-464.
26. Choudhury, M. M., & Harrigan, P. (2014). CRM to social CRM: The integration of new technologies into customer relationship management. *Journal of Strategic Marketing*, 22(2), 149-176.
27. Ryzhkova, I., Polosukhina, I., & Tikhomirova, D. (2018). Social media for customer relationship management: Benefits and risks. *Management Decision*, 56(6), 1264-1280.
28. Varadarajan, R., & Yadav, M. S. (2002). Marketing strategy and the internet: An organizing framework. *Journal of the Academy of Marketing Science*, 30(4), 296-312.
29. Jain, S., & Sharma, A. (2018). Customer segmentation in e-commerce: A review and future directions. *Management Decision*, 56(10), 2206-2230.
30. Coussement, K., De Bock, K. W., & De Pelsmacker, P. (2014). Customer churn prediction in the online gambling industry: The beneficial effect of ensemble learning. *Decision Support Systems*, 60, 21-31.
31. Kumar, V., & Reinartz, W. (2016). Creating enduring customer value. *Journal of Marketing*, 80(6), 36-68.
32. Yoo, B., Donthu, N., & Lee, S. (2000). An examination of selected marketing mix elements and brand equity. *Journal of the Academy of Marketing Science*, 28(2), 195-211.
33. Jain, V., Dixit, A., & Khare, A. (2018). Customer segmentation in e-commerce using machine learning: A comprehensive review. *Journal of King Saud University-Computer and Information Sciences*.
34. Verhoef, P. C., Neslin, S. A., & Vroomen, B. (2007). Multichannel customer management: Understanding the research-shopper phenomenon. *International Journal of Research in Marketing*, 24(2), 129-148.
35. Saarijärvi, H., Kannan, P. K., & Kuusela, H. (2013). Value co-creation: theoretical approaches and practical implications. *European Business Review*, 25(1), 6-19.
- 36.

## **APPENDICES**

### **APPENDICE A: SAMPLE CODE FOR CLUSTER PREDICTION**

```
data.isna().sum()

data.head()

clustering_data = data.iloc[:,[2,4]]

clustering_data.head()

fig, ax = plt.subplots(figsize=(15,7))

sns.set(font_scale=1.5)

ax = sns.scatterplot(y=clustering_data['Electronics'],x=clustering_data['Annual_Income'],
s=70, color='#f73434', edgecolor='black', linewidth=0.3)

ax.set_ylabel('Electronics bought')

ax.set_xlabel('Annual Income (in X 1000 KSH )')

plt.title('Electronics bought per Annual Income', fontsize = 20)

plt.show()

import warnings

warnings.filterwarnings('ignore')

from sklearn.cluster import KMeans
```



```

wcss=[]

for i in range(1,30):

    km = KMeans(i)

    km.fit(clustering_data)

    wcss.append(km.inertia_)

np.array(wcss)

fig, ax = plt.subplots(figsize=(15,7))

ax = plt.plot(range(1,30),wcss, linewidth=2, color="red", marker ="8")

plt.ylabel('WCSS')

plt.xlabel('No. of Clusters (k)')

plt.title('The Elbow Method', fontsize = 20)

plt.show()

from sklearn.cluster import KMeans

kms = KMeans(n_clusters=3, init='k-means++')

kms.fit(clustering_data)

clusters = clustering_data.copy()

clusters['Cluster_Prediction'] = kms.fit_predict(clustering_data)

clusters.head()

kms.cluster_centers_

fig, ax = plt.subplots(figsize=(15,7))

plt.scatter(x=clusters[clusters['Cluster_Prediction'] == 1]['Annual_Income'],

            y=clusters[clusters['Cluster_Prediction'] == 1]['Electronics'],

            s=70,edgecolor='black', linewidth=0.3, c='red', label='Cluster 1')

plt.scatter(x=clusters[clusters['Cluster_Prediction'] == 2]['Annual_Income'],

```

```

y=clusters[clusters['Cluster_Prediction'] == 2]['Electronics'],

s=70,edgecolor='black', linewidth=0.2, c='Magenta', label='Cluster 2')

plt.scatter(x=clusters[clusters['Cluster_Prediction'] == 0]['Annual_Income'],

y=clusters[clusters['Cluster_Prediction'] == 0]['Electronics'],

s=70,edgecolor='black', linewidth=0.3, c='deepskyblue', label='Cluster 3')


plt.scatter(x=kms.cluster_centers_[0], y=kms.cluster_centers_[0], s = 120, c = 'yellow',
label = 'Centroids',edgecolor='black', linewidth=0.3)

plt.legend(loc='right')

plt.xlim(0,140)

plt.ylim(0,30)

plt.xlabel('Annual Income (in Thousand KSH)')

plt.ylabel('Electronics bought')

plt.title('Clusters', fontsize = 20)

plt.show()

fig, ax = plt.subplots(nrows=2, ncols=2, figsize=(15,20))

ax[0,0].scatter(x=clusters[clusters['Cluster_Prediction'] == 1]['Annual_Income'],

y=clusters[clusters['Cluster_Prediction'] == 1]['Electronics'],

s=40,edgecolor='black', linewidth=0.3, c='red', label='Cluster 1')

ax[0,0].scatter(x=kms.cluster_centers_[1,0], y=kms.cluster_centers_[1,1],

s = 120, c = 'yellow',edgecolor='black', linewidth=0.3)

ax[0,0].set(xlim=(0,140), ylim=(0,30), xlabel='Annual Income', ylabel='Spending Score',
title='Cluster 1')

ax[0,1].scatter(x=clusters[clusters['Cluster_Prediction'] == 2]['Annual_Income'],

```

```

y=clusters[clusters['Cluster_Prediction'] == 2]['Electronics'],

s=40,edgecolor='black', linewidth=0.3, c='magenta', label='Cluster 2')

ax[0,1].scatter(x=kms.cluster_centers_[2,0], y=kms.cluster_centers_[2,1],

s = 120, c = 'yellow',edgecolor='black', linewidth=0.3)

ax[0,1].set(xlim=(0,140), ylim=(0,30), xlabel='Annual Income', ylabel='Spending Score',
title='Cluster 2')


ax[1,0].scatter(x=clusters[clusters['Cluster_Prediction'] == 0]['Annual_Income'],

y=clusters[clusters['Cluster_Prediction'] == 0]['Electronics'],

s=40,edgecolor='black', linewidth=0.3, c='deepskyblue', label='Cluster 3')

ax[1,0].scatter(x=kms.cluster_centers_[0,0], y=kms.cluster_centers_[0,1],

s = 120, c = 'yellow',edgecolor='black', linewidth=0.3)

ax[1,0].set(xlim=(0,140), ylim=(0,30), xlabel='Annual Income', ylabel='Spending Score',
title='Cluster 3')


fig.delaxes(ax[1,1])

fig.legend(loc='right')

fig.suptitle('Individual Clusters')

plt.show()

```