

## OPTIMIZING CUSTOMER SEGMENTATION WITH CLUSTERING AND DATA MINING METHODS

**Dr. P. Rajavardhan Reddy<sup>1</sup>, Dr. K. Raghu Naga Prabhakar<sup>2</sup>**

<sup>1</sup> Professor, Department of Business Management, Aurora's PG College (MBA), Uppal, Hyderabad  
Email: [rajavardhan32@gmail.com](mailto:rajavardhan32@gmail.com)

<sup>2</sup> Professor, Department of Business Management, Aurora's PG College (MBA), Uppal, Hyderabad  
Email: [prabhakalepu@gmail.com](mailto:prabhakalepu@gmail.com)

### ABSTRACT

Understanding consumer behavior and classifying customers based on their demographics and purchase patterns is essential in today's competitive industry. This is a crucial component of customer segmentation since it helps marketers better target various audience segments with promotional, marketing, and product development strategies. Customer segmentation is the process of grouping customers according to shared traits, spending preferences, and shopping trends.

The process of figuring out how to engage with clients in different categories to optimize each client's value to the business is known as customer segmentation. Marketers may connect with each consumer in the most efficient manner by using customer segmentation. Utilizing a vast amount of client data, customer segmentation studies accurately identify discrete customer groups according to behavioral, demographic, and other factors. Unlike supervised machine learning techniques, K-means clustering is an unsupervised approach.

This approach is used when the dataset contains unlabeled data. Information that hasn't been categorized or grouped into any of the accessible groups is known as unlabeled data. In customer segmentation, various methods are used to find the ideal number of clusters; however, each method has limitations of its own. For example, the Density-based Spatial Clustering of Applications with Noise (DBSCAN) algorithm fails when the density of the clusters changes. Research on Regency, Frequency, Monetary Value (RFM) is not predicated on forecasts for the future, but rather on past facts. The usage of labeled data is necessary since the Hierarchical Clustering approach cannot undo previous work. On the other hand, the K-means approach guarantees convergence, initiates the centroid's positions, and promptly adjusts to novel and ideal cluster numbers.

**Keywords:** Customer Segmentation, Data Mining, Density-Based Spatial Clustering of Applications With Noise (Dbscan), Regency, Frequency, Monetary Value (Rfm), K-Means

### 1. INTRODUCTION

#### Overview

When creating criteria for segmenting consumers into logical groups, such as customers who came from a given source, dwell in a specific location, or purchased a specific product or service, it's always easier to make assumptions and depend on "gut emotions." On the other hand, these high-level categorizations seldom produce the desired results. Certain clients will,

## OPTIMIZING CUSTOMER SEGMENTATION WITH CLUSTERING AND DATA MINING METHODS

without a doubt, spend more money with a firm than others.

Most loyal consumers will spend a considerable quantity of money over an extended period. Good clients will either spend a little overtime or a lot in a short amount of time. Others will not overspend or stay for a lengthy amount of time. The most successful form of customer segmentation analysis is to split consumers into groups based on projections of their full future worth to the company, to address each group or individual in the most effective way possible to maximize that future, or a lifetime, value.

1. Behavioral
2. Demographic
3. Geographic
4. Psychographic

### **Motivation for the project**

Because the marketer's goal is to maximize the value (revenue and/or profit) from each client, knowing how each marketing action will affect the customer ahead of time is critical. In an ideal world, such "activity-centric" customer segmentation would focus on the long-term impact of a marketing action on customer lifetime value (CLV) rather than the short-term value of marketing activities. As a result, customers must be divided into groups or segments based on their CLV. [2]

In a distinct kind of customer segmentation, machine learning algorithms are utilized to find new groups. In contrast to marketer-designed segmentation models like the ones mentioned above, machine learning consumer segmentation allows advanced algorithms to reveal insights and catches that marketers could overlook on their own. Marketers that can relate their segmentation model to campaign results will see their customer groups improve over time. In these situations, the machine learning model will be able to not only fine-tune its segment definition, but also assess whether one segment is outperforming the others, therefore optimizing marketing performance.

### **Problem Definition And Scenarios**

- To implement the customer segmentation analysis and dividing consumers into groups based on common characteristics, spending habits, and purchasing patterns
- Product Recommendation based on Customer Segmentation Engine.
- To develop a success prediction models for customer segmentation.
- Customer segmentation is a technique for determining how to interact with customers in various categories to maximize the value of each client to the company.

### **Organization of the report**

The chapter 2 is about the literature review on the various projects that are studied and understanding of the project in more detailed way and analysis of the system that are described.

The chapter 3 is discussed on the project description in which the analysis of the existed system and how it is differ from the proposed system and analysis of the model using spiral diagram.

The chapter 4 is all about the system design where the system is explained using block, flow etc. how they are developed on the basis of simplification this gives the clear understanding of project.

The chapter 5 is mainly discussed on the system requirements of the project such as hardware and software components and their specifications.

The chapter 6 is about all modules of the project and its description and explanation that the project is divided in to module to make every module completed successfully.

The chapter 7 is mainly discussed on the implementation of the project such as hardware and

software components and their specifications.

The chapter 8 is all about the result and analysis of parts that are used in our project.

The chapter 9 is all about the conclusion of the project in which how the project is concluded and what we have done in this project and also about the future enhancement of the project what will be the extension of this project for later development This the organization of thesis what we have discussed in the following chapter this brief understanding on each and every chapter.

## **2.LITERATURE SURVEY AND RELATED WORK**

### **Introduction**

This chapter is about analysis of the different projects that are published up to now and about the project the details that given and discussed in the analysis of project in detail.

#### **Explore Patterns of Customer Segments**

Customer behavior evolves and changes over time in real-world settings. In order to establish efficient marketing strategies, this dynamism must be considered while conducting consumer segmentation analyses and other business-related tasks. The study's major goal is to look at the patterns of structural changes in consumer groups. There hasn't been any study done on this subject yet. This is the first research to look at the influence of consumer dynamics on structural changes in segments. The goal of this study is to create a way to describe and explain this problem. Using clustering and sequential rule mining approaches, a novel strategy is suggested. A new concept and methodology for identifying distinct sequential rules are also established. The proposed strategy is tested using data from customers.[1]

#### **Discovering New Business Opportunities**

For sales and marketing departments inside major organizations, identifying and understanding new markets, clients, and partners is a vital task. Intel's Sales and Marketing Group (SMG) is experiencing similar challenges as it expands into new markets and industries and evolves its present business. To aid SMGs in sorting through millions of organizations across several locations and languages to find relevant routes, sophisticated automation that enables a fine-grained knowledge of enterprises is essential in today's challenging technological and commercial context. We show a system developed in our company that mines millions of public business

Web pages to provide a faceted client representation. We focus on two key characteristics of the customer that are important for discovering suitable opportunities: Industry sectors (which range from huge verticals like banking and insurance to smaller niches like manufacturing). [7]

#### **RFM Customer Segmentation**

The RFM model, which is utilized in the traditional retail business for consumer segmentation, is not ideal for an industry with different social group qualities; hence the RFMC model is established by adding the social relations parameter C. For this empirical investigation, educational e-commerce firm M was chosen, and the k-means algorithm was utilized to cluster legitimate clients of enterprise M, resulting in five unique customer groups and proving the model's usefulness. [8]

#### **High-Dimensional Customer Segmentation**

The Omni channel becomes a hot topic as a result of the rapid rise of e-commerce and clients' growing familiarity with multichannel shopping. To satisfy the present trend of client demand, several firm organizations to work on the Omni channel business issue and are progressively

## OPTIMIZING CUSTOMER SEGMENTATION WITH CLUSTERING AND DATA MINING METHODS

devoting their efforts to both online and offline business. As a result, there's no doubting that understanding online customers' purchasing patterns is critical to Omni channel success. Using the RFM (regency, frequency, monetary) model and the k-means clustering technique, consumers' information is retrieved and customers are segmented. To expand the RFM model, we divide total frequency and monetary data into weekly level data, resulting in a reduction in the number of variables associated with one week. [13]

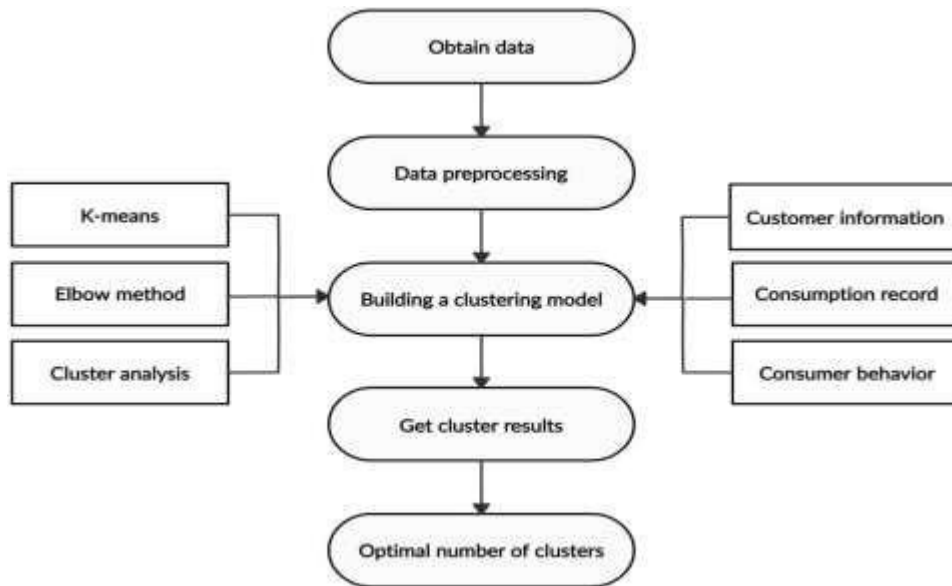
### 3. EXISTING SYSTEM

- Outliers that are isolated in low-density regions are identified as outliers by DBSCAN, which merges points that are close together. Two important factors create the model's 'density': the minimum number of points required to generate dense regions min samples and the distance required to establish a neighborhood eps. Larger min samples or lower eps demand a higher density to construct a cluster.
- RFM Analysis is a marketing method that combines the three aspects of regency,, and monetary value to analyze and understand customer behavior. The RFM frequency Analysis will help businesses divide their customer base into a variety of homogenous groups, allowing them to connect with each group using a variety of targeted marketing strategies.
- A way of grouping comparable components is hierarchical clustering, often known as hierarchical cluster analysis. The endpoint is made up of several clusters, each of which is unique yet shares a lot of similarities.[2]

### 4. PROPOSED SYSTEM

The number of segmentation choices is practically limitless, and they are mostly determined by the quantity of consumer data we have available. It starts with the basics, such as gender, hobby, or age, and progresses to elements such as the amount of time since the user last visited the business or the amount of time spent on website.

- i. Geographic: The concept of geographic consumer segmentation is straightforward; it all boils down to the user's geographic location. This may be done in several different ways. You may sort your results by city, zip code, nation, or state.
- ii. Demographics: Demographic segmentation takes into account the structure, size, and movement patterns of consumers over time and space. Many businesses create and market products that are based on gender differences. Another important factor is the status of the parents. This type of information may be obtained through customer surveys.
- iii. Behavioral: Customer segmentation based on past behavior can be used to forecast future behavior. Customers' favorite brands, for example, or the times of the year when they make the most purchases. The behavioral component of consumer segmentation seeks to understand not just why people buy things, but also how those reasons change over time.
- iv. Psychological: The psychological segmentation of customers includes personality traits, attitudes, and beliefs. Consumer surveys are used to collect this information, which may then be used to assess customer sentiment.



**FIG 1 – SYSTEM ARCHITECTURE**

## 5. METHODOLOGIES SEGMENTATION

Let's have a look at how the consumer segmentation challenge is being solved.

- i. Create a business case: The first stage is to make a business case for the project. There must be an aim for everything. You don't want to become involved with this blindly. Otherwise, the result will be unorganized and untidy. You'll need a business case instead. It's identifying the most lucrative consumer groupings within the total pool of customers in this scenario. [9]
- ii. Prepare the data: The data must then be prepared in the second stage. What is the size of your data set? In this instance, a hundred, thousand or ten thousand client data is preferable. This is due to the fact that you will be able to see more patterns and trends. You'll also need a set of features based on the most relevant business indicators inside the data collection. The data is then preprocessed to reduce discrepancies, which helps in better data analysis.
- iii. Data analysis and exploration: Data analysis and investigation is the third stage. This is an important stage since it will allow you to discover some intriguing relationships and trends in your data. You will be able to better comprehend the client's interests and purchase patterns as a result of this, and you will be able to determine which traits are most directly associated with the consumer and, obviously, the business. [14]
- iv. Clustering analysis: Clustering analysis in the context of a client is the fourth phase in the process. The use of a mathematical model to uncover groupings of similar consumers based on the tiniest differences between customers is known as segmentation clustering analysis. The purpose of cluster analysis within each group is to precisely categorize consumers so that personalization may be used to create more successful customer marketing. A mathematical approach called k-means clustering analysis is a popular cluster analysis tool. The resulting clusters aid in improved consumer modeling and analysis. There are no pre-set thresholds or standards in place for this procedure. The data, on the other hand, exposes the customer prototypes that exist intrinsically inside the consumer base.

## OPTIMIZING CUSTOMER SEGMENTATION WITH CLUSTERING AND DATA MINING METHODS

- v. Choosing optimal hyperparameters: The fifth step is to select the best hyperparameters. "Tuning" or "hyperparameter optimization" is the process of selecting the optimal set of hyperparameters for an algorithm. Based on our past work, this is the next stage since it assists us in identifying the most accurate and rewarding client groups.[10]
- vi. Visualization and interpretation: Visualization
- vii. and interpretation are the final steps in the process. Now it's time to visualize and interpret your findings. Businesses can improve marketing campaigns by targeting features, launches, and product roadmaps when they have profitable customer profiles at their fingertips. This gives the company a much clearer picture of which customers are most likely to stick around.[15]

### K-MEANS ALGORITHM

Let's look at the K-means method, which we'll use in this project. The k-mean technique is an iterative method for splitting a data collection into K distinct, non-overlapping subgroups, each of which contains just one data point. It sought to make the data points in the intracluster clusters as comparable as feasible. while preserving as much space between the clusters as possible. It distributes data points to clusters in such a way that the squared distances add up to one. Inside clusters, the cluster centroid is the location with the least variation, and the most homogeneous data points are found within the same cluster. It allows us to cluster the data into different groups and a convenient way to discover the categories of groups in the unlabeled dataset on its own without the need for any training.

It is a centroid-based algorithm, where each cluster is associated with a centroid. The main aim of this algorithm is to minimize the sum of distances between the data point and their corresponding clusters.

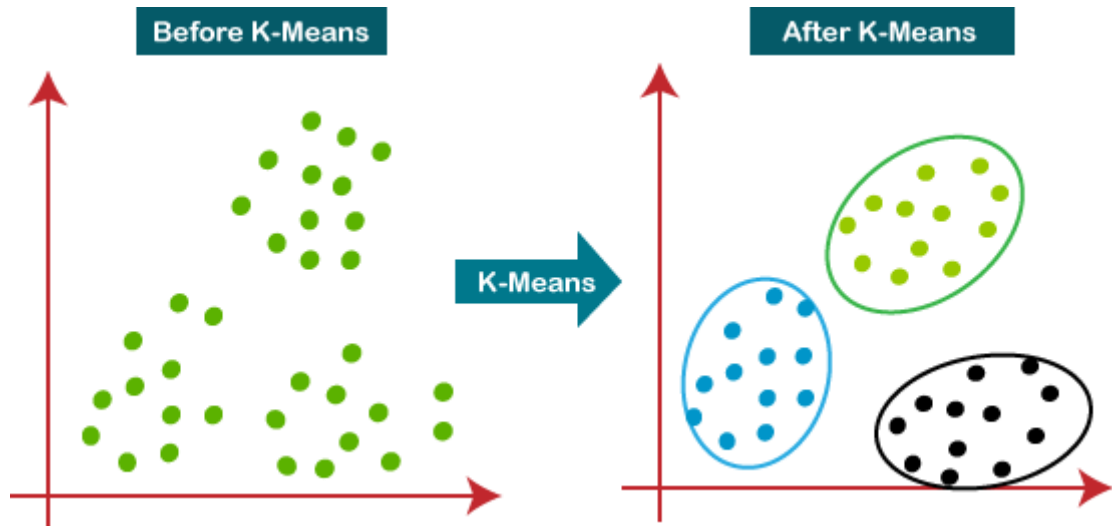
The algorithm takes the unlabeled dataset as input, divides the dataset into k-number of clusters, and repeats the process until it does not find the best clusters. The value of k should be predetermined in this algorithm.

The k-means [clustering](#) algorithm mainly performs two tasks

- Determines the best value for K center points or centroids by an iterative process.
- Assigns each data point to its closest k-center. Those data points which are near to the particular k-center, create a cluster.

Hence each cluster has data points with some commonalities, and it is away from other clusters.

The below diagram explains the working of the K-means Clustering Algorithm:



**Fig.2 K-Means Clustering Algorithm**

The working of the K-Means algorithm is explained in the below steps:

**Step-1:** Select the number K to decide the number of clusters.

Step-2: Select random K points or centroids. (It can be other from the input dataset).

Step-3: Assign each data point to their closest centroid, which will form the predefined K clusters.

Step-4: Calculate the variance and place a new centroid of each cluster.

Step-5: Repeat the third steps, which means reassign each datapoint to the new closest centroid of each cluster.

Step-6: If any reassignment occurs, then go to step-4 else go to FINISH.

Step-7: The model is ready.

The performance of the K-means clustering algorithm depends upon highly efficient clusters that it forms. But choosing the optimal number of clusters is a big task. There are some different ways to find the optimal number of clusters, but here we are discussing the most appropriate method to find the number of clusters or value of K

### **ELBOW METHOD**

The Elbow method is one of the most popular ways to find the optimal number of clusters. This method uses the concept of WCSS value. WCSS stands for Within Cluster Sum of Squares, which defines the total variations within a cluster. The formula to calculate the value of WCSS (for 3 clusters) is given below:

$$WCSS = \sum_{P_i \text{ in Cluster1}} \text{distance}(P_i, C_1)^2 + \sum_{P_i \text{ in Cluster2}} \text{distance}(P_i, C_2)^2 + \sum_{P_i \text{ in Cluster3}} \text{distance}(P_i, C_3)^2$$

In the above formula of WCSS,

$\sum_{P_i \text{ in Cluster1}} \text{distance}(P_i, C_1)^2$ : It is the sum of the square of the distances between each data point and its centroid within a cluster1 and the same for the other two terms.

To measure the distance between data points and centroid, we can use any method such as Euclidean distance or Manhattan distance.

To find the optimal value of clusters, the elbow method follows the below steps:

- It executes the K-means clustering on a given dataset for different K values (ranges from 1-

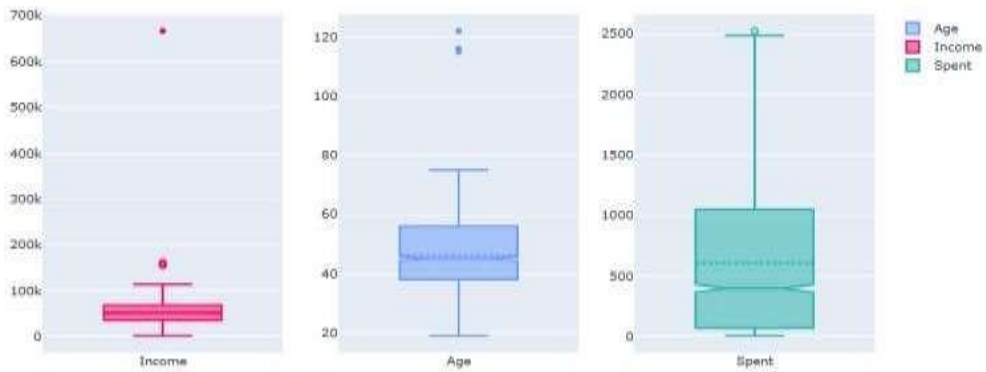
# OPTIMIZING CUSTOMER SEGMENTATION WITH CLUSTERING AND DATA MINING METHODS

10).

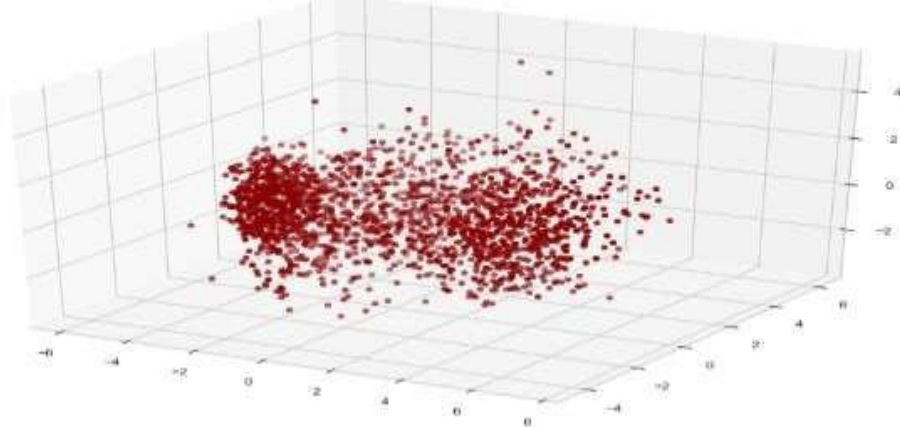
- For each value of K, calculates the WCSS value.
- Plots a curve between calculated WCSS values and the number of clusters K.
- The sharp point of bend or a point of the plot looks like an arm, then that point is considered as the best value of K.

## 6. RESULTS AND SCREEN SHOTS:

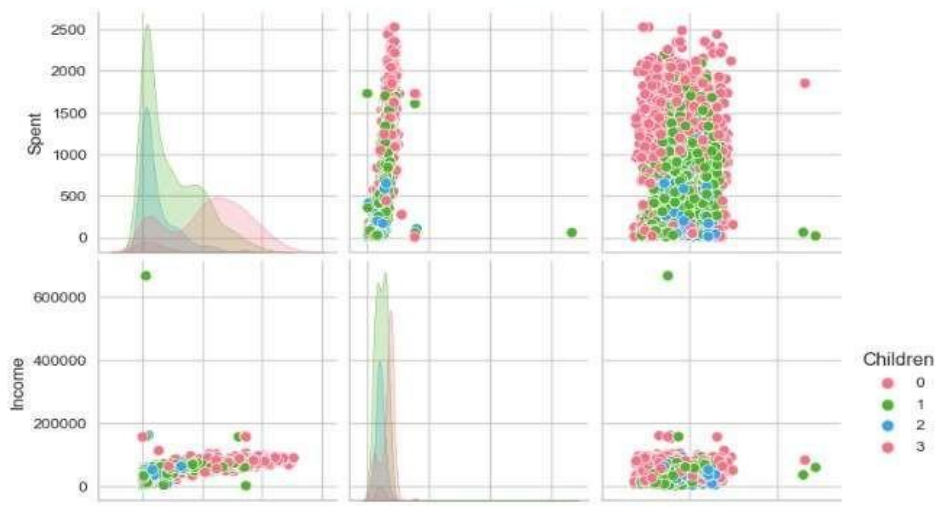
Box Plots for Numerical Variables



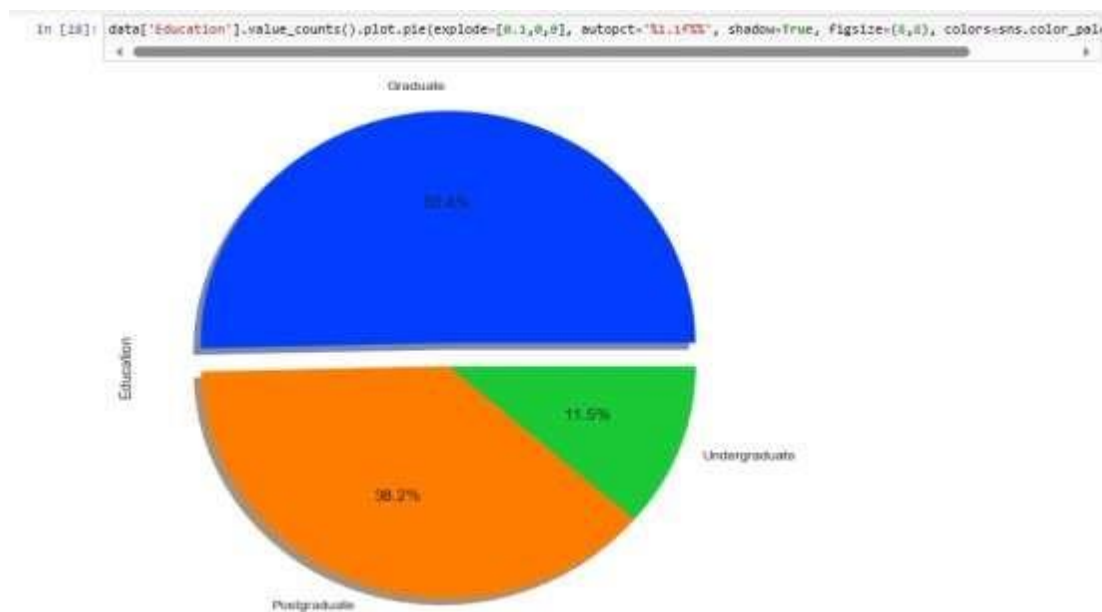
A 3D Projection of Data in the Reduced Dimension



```
In [24]: sns.pairplot(data , vars=['Spent','Income','Age'] , hue='Children' , palette='husl');
```







## 7. CONCLUSION AND FUTURE SCOPE

### CONCLUSION

Our contribution to this research paper is to highlight and solve the challenges of product recommendation based on the customer segmentation engine. Data science can use Customer segmentation to build a stronger relationship with their customers. It enables them to make educated retention choices, develop new features, and position their product strategically in the market.

### FUTURE SCOPE

As part of their customer segmentation strategy, retailers with extensive, multi-category offers must show their things in such a way that target customers may search and pick from those offerings. In the first piece of this dissertation, a product segmentation approach is described. The proposed approach offers merchants a methodology for identifying customer-centric, cross-category product segments from large numbers of goods across multiple categories, where products within a segment are bought by the same type of customers. The research also looks at the relationship between the recommended product segmentation approach and a customer segmentation method. Because the approaches are so closely linked, the product and customer groups inferred by each are likely to be identical.

## 8. REFERENCES

- [1] ElhamAkhondZadehNoughabi;AmirAlbadvi;BehrouzHomayoun Far "How Can We Explore Patterns of Customer Segments' Structural Changes? A Sequential Rule Mining Approach" In Proceedings of the 2015 IEEE International Conference on Information Reuse and Integration, 26 October 2015
- [2] Wayne XinZhao;SuiLi;YulanHe;Edward Y. Chang;Ji-RongWen;XiaomingLi"Connecting Social Media to E-Commerce: Cold-Start Product Recommendation Using Microblogging Information" IEEE Transactions on Knowledge and Data Engineering,volume: 28, pages: 1147-

## OPTIMIZING CUSTOMER SEGMENTATION WITH CLUSTERING AND DATA MINING METHODS

1159, 17 December 2015

[3] Xiaojun Chen; Yixiang Fang; Min Yang; FeipingNie; Zhou Zhao; Joshua Zhexue Huang "PurTreeClust: A Clustering Algorithm for Customer Segmentation from Massive Customer Transaction Data" IEEE Transactions on Knowledge and Data Engineering, volume: 30, page 559-572,26 October 2017

[4] QingluGao;DeXia;YangyanShi;JiQuan"Policies Adoption for Supply Disruption Mitigation Based on Customer Segmentation" IEEE Access, volume: 7, page 47329 - 47338, 29 March 2019

[5] Fuxiang Liu"3D Block Matching Algorithm in Concealed Image Recognition and E-Commerce Customer Segmentation" IEEE Sensors Journal, volume: 20, pp. 11761 - 11769, 19 August 2019

[6] Caroline GobboSáCavalcante;Diego Castro Fettermann"Recommendations for Product Development of Intelligent Products" IEEE Latin America Transactions, volume: 17, pages: 1645-1652, October 2019

[7] Itay Lieder; MeiravSegal;EranAvidan;AsafCohen;Tom Hope "Learning a Faceted Customer Segmentation for Discovering new Business Opportunities at Intel", In Proceedings of 2019 IEEE International Conference on Big Data (Big Data),24 Dec. 2019