

SRM INSTITUTE OF SCIENCE AND TECHNOLOGY

CUSTOMER SEGMENTATION AND PROFILING USING CLUSTERING TECHNIQUES

A DATA ANALYSIS AND VISUALIZATION PROJECT

**PRESENTED TO
Dr. A. Shobanadevi
Vaskar Deka**

**PRESENTED BY
Pijush Pathak
RA2011027010152**

11-May-2023

Agenda

- 1 Abstract
- 2 Objectives
- 3 Importing Libraries
- 4 Loading Data
- 5 Data Cleaning
- 6 Data Preprocessing
- 7 Dimensionality Reduction
- 8 Clustering
- 9 Evaluating Models
- 10 Profiling
- 9 Conclusion
- 10 Future Enhancements

ABSTRACT

The project '**Customer Segmentation and Profiling using Clustering Technique**' focuses on utilizing clustering techniques to group customers based on their shared characteristics. The project follows a systematic approach, starting with data collection from various sources, including demographic information, purchase history, and online behavior. The collected data is then preprocessed, including handling missing values, removing outliers, and standardizing variables. Next, relevant variables are selected to capture the key aspects of customer behavior. Clustering techniques, such as k-means, hierarchical clustering, or DBSCAN, are applied to identify distinct customer segments. Finally, customer profiles are created based on the characteristics of each segment, enabling businesses to tailor their marketing efforts, enhance customer experiences, and improve overall business performance.

Objectives

- 1/ Importing Libraries
- 2/ Loading Data
- 3/ Data Cleaning
- 4/ Data Preprocessing
- 5/ Dimensionality Reduction
- 6/ Clustering
- 7/ Evaluating Models
- 8/ Profiling

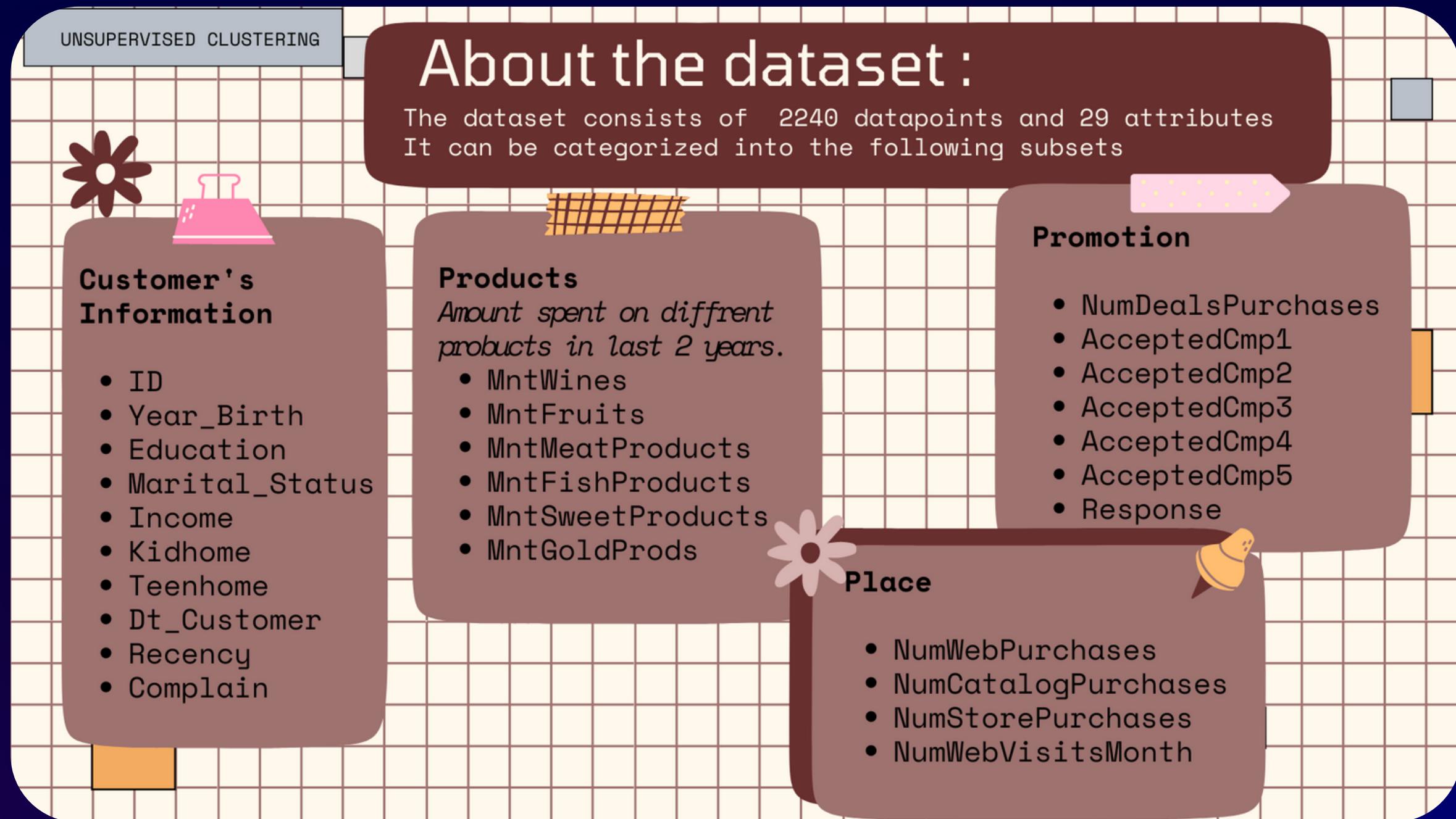
Importing Libraries

- 1 **Numpy:** It is a library for numerical computing in Python.
- 2 **Pandas:** It is a library for data manipulation and analysis
- 3 **Matplotlib:** It is a library for data visualization in Python
- 4 **Seaborn:** It is a library for creating more visually appealing and informative statistical graphics in Python.
- 5 **Label Encoder, Standard Library:** They are preprocessing techniques from the scikit-learn library for preparing data for machine learning models.
- 6 **KElbow Visualizer:** It helps to determine the optimal number of clusters for KMeans clustering algorithm.
- 7 **Axes3D:** It is a 3D plotting tool from matplotlib.

Loading Data

We shall load the dataset and find the number of datapoints.

The number of Data Points in the dataset is 2240



Data Cleaning

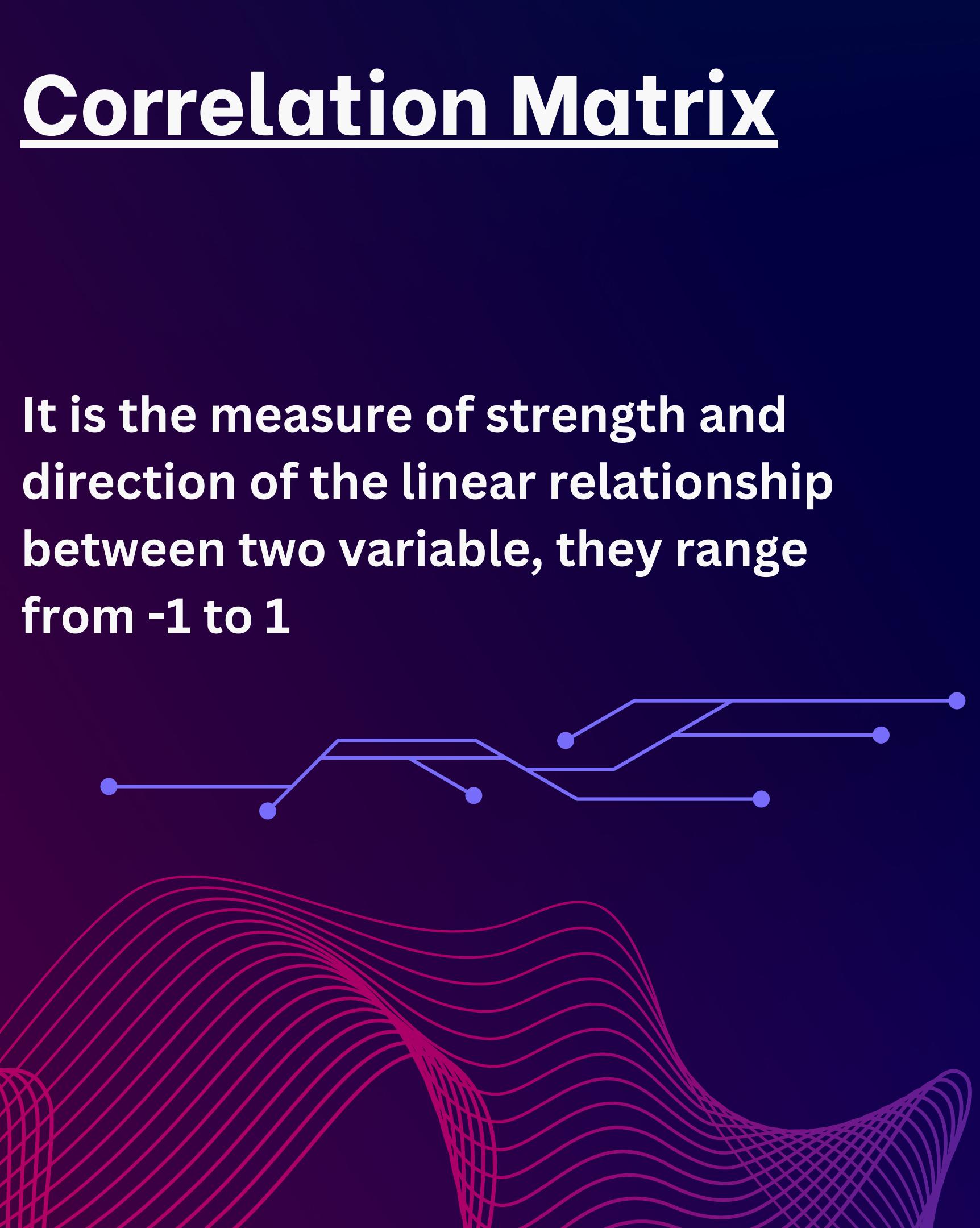
In order to, get a full grasp of what steps should be taking to clean the dataset. We shall have a look at the information in data. We found that there are missing values in income. First of all, for the missing values, I had dropped the rows that have missing income values.

In the next bit, I had performed the following steps to engineer some new features:

1. Extracted the "Age" of a customer by the "Year_Birth" that indicated the birth year of the respective person.
2. Created another feature "Spent" indicating the total amount spent by the customer in various categories over the span of two years.
3. Created another feature "Living_With" out of "Marital_Status" to extract the living situation of couples.
4. Created a feature "Children" to indicate total children in a household that is, kids and teenagers.
5. To get further clarity of household, I had created a feature indicating "Family_Size"
6. Created a feature "Is_Parent" to indicate parenthood status
7. Lastly, I had created three categories in the "Education" by simplifying its value counts. Dropping some of the redundant features

Correlation Matrix

It is the measure of strength and direction of the linear relationship between two variable, they range from -1 to 1



	Income	Kidhome	Teenhome	Recency	Wines	Fruits	Meat	Fish	Sweets	Gold	NumDealsPurchases	NumWebPurchases	NumCatalogPurchases	NumStorePurchases	NumWebVisitsMonth	AcceptedCmp3	AcceptedCmp4	AcceptedCmp5	AcceptedCmp1	AcceptedCmp2	Complain	Response	Customer_For	Age	Spent	Children	Family_Size	Is_Parent			
Income	1	-0.51	0.035	0.008	0.69	0.51	0.69	0.52	0.52	0.39	-0.11	0.46	0.7	0.63	-0.65	0.015	0.22	0.4	0.33	0.1	-0.028	0.16	-0.028	0.2	0.79	-0.34	-0.29	-0.4			
Kidhome	-0.51	1	-0.039	0.011	-0.5	-0.37	-0.44	-0.39	-0.38	-0.35	0.22	-0.37	-0.5	-0.5	0.45	0.016	-0.16	-0.2	-0.17	-0.082	0.037	0.078	0.058	-0.24	-0.56	0.69	0.58	0.52			
Teenhome	0.035	0.039	1	0.014	0.003	0.90	0.18	-0.26	-0.21	-0.16	0.019	0.39	0.16	-0.11	0.049	0.13	-0.043	0.038	-0.19	-0.15	-0.016	0.007	0.15	0.009	0.36	-0.14	0.7	0.59	0.59		
Recency	0.008	0.010	0.014	1	0.01	0.005	0.028	0.007	0.025	0.018	0.002	0.005	0.020	0.0004	0.019	0.032	0.018	0.0002	0.02	0.001	0.005	0.057	-0.2	0.03	0.10	0.16	0.02	0.018	0.015	0.0022	
Wines	0.69	-0.50	0.003	0.016	1	0.39	0.57	0.4	0.39	0.39	0.009	0.105	0.63	0.64	-0.32	0.061	0.37	0.47	0.35	0.21	-0.036	0.25	0.15	0.16	0.89	-0.35	-0.3	-0.34			
Fruits	0.51	-0.37	-0.18	0.005	0.39	1	0.55	0.59	0.57	0.39	-0.13	0.3	0.49	0.46	-0.42	0.015	0.066	0.21	0.19	0.009	0.003	0.12	0.06	0.013	0.61	-0.4	-0.34	-0.41			
Meat	0.69	-0.44	-0.26	0.023	0.57	0.55	1	0.57	0.53	0.36	-0.12	0.31	0.73	0.49	-0.54	0.018	0.092	0.38	0.31	0.044	0.021	0.24	0.071	0.034	0.85	-0.5	-0.43	-0.57			
Fish	0.52	-0.39	-0.21	0.0007	0.90	0.59	0.57	1	0.58	0.43	-0.14	0.3	0.53	0.46	-0.45	0.002	0.016	0.19	0.26	0.002	0.019	0.11	0.078	0.041	0.64	-0.43	-0.36	-0.45			
Sweets	0.52	-0.38	-0.16	0.025	0.39	0.57	0.53	0.58	1	0.36	-0.12	0.33	0.49	0.46	-0.42	0.001	0.029	0.26	0.25	0.01	-0.021	0.12	0.076	0.022	0.61	-0.39	-0.33	-0.4			
Gold	0.39	-0.35	0.019	0.018	0.39	0.39	0.36	0.43	0.36	1	0.053	0.41	0.44	0.39	-0.25	0.13	0.024	0.18	0.17	0.051	-0.03	0.14	0.15	0.06	0.53	-0.27	-0.24	-0.25			
NumDealsPurchases	-0.11	0.22	0.39	0.002	0.6	0.009	0.10	0.13	-0.12	-0.14	-0.12	0.053	1	0.24	-0.012	0.066	0.35	-0.023	0.016	-0.18	-0.13	0.038	0.003	0.032	0.2	0.066	0.066	0.44	0.37	0.39	
NumWebPurchases	0.46	-0.37	0.16	0.005	0.57	0.55	0.3	0.31	0.3	0.33	0.41	0.24	1	0.39	0.52	-0.052	0.043	0.16	0.14	0.16	0.035	0.014	0.15	0.17	0.16	0.53	-0.15	-0.12	-0.073		
NumCatalogPurchases	0.7	-0.5	-0.11	0.024	0.63	0.49	0.73	0.53	0.49	0.44	-0.012	0.39	1	0.52	-0.52	0.1	0.14	0.32	0.31	0.1	-0.019	0.22	0.091	0.13	0.78	-0.44	-0.37	-0.45			
NumStorePurchases	0.63	-0.5	0.049	0.004	0.64	0.46	0.49	0.46	0.46	0.39	0.066	0.52	0.52	1	-0.43	0.069	0.18	0.21	0.18	0.085	0.012	0.036	0.1	0.14	0.68	-0.32	-0.27	-0.28			
NumWebVisitsMonth	-0.65	0.45	0.13	-0.019	0.32	-0.42	-0.54	-0.45	-0.42	-0.25	0.35	-0.052	0.52	-0.43	1	0.061	0.029	0.28	-0.2	-0.007	0.021	0.002	0.26	-0.12	0.5	0.42	0.35	0.48			
AcceptedCmp3	-0.015	0.016	0.043	0.032	0.061	0.10	0.015	0.018	0.000	0.001	0.017	0.13	-0.023	0.043	0.1	-0.069	0.061	1	-0.080	0.081	0.096	0.072	0.009	0.60	0.250	0.006	0.061	0.053	-0.02	0.024	0.0055
AcceptedCmp4	0.22	-0.16	0.038	0.018	0.37	0.006	0.092	0.016	0.029	0.024	0.016	0.16	0.14	0.18	-0.029	0.08	1	0.31	0.24	0.3	-0.027	0.18	0.014	0.07	0.25	-0.088	0.077	0.077			
AcceptedCmp5	0.4	-0.2	-0.19	0.002	0.47	0.21	0.38	0.19	0.26	0.18	-0.18	0.14	0.32	0.21	-0.28	0.081	0.31	1	0.41	0.220	0.008	0.432	-0.023	0.019	0.47	-0.28	-0.23	-0.35			
AcceptedCmp1	0.33	-0.17	-0.15	0.021	0.35	0.19	0.31	0.26	0.25	0.17	-0.13	0.16	0.31	0.18	-0.2	0.096	0.24	0.41	1	0.18	-0.025	0.3	-0.037	0.012	0.38	-0.23	-0.19	-0.28			
AcceptedCmp2	0.1	-0.08	0.01	0.001	0.4	0.21	0.009	0.044	0.002	0.030	0.01	0.051	0.038	0.035	0.1	0.085	0.007	0.072	0.3	0.22	0.18	1	-0.011	0.17	0.006	0.007	0.80	0.14	-0.07	-0.06	0.082
Complain	-0.028	0.037	0.007	0.005	0.36	0.003	0.021	0.019	0.021	0.03	0.014	0.012	0.021	0.009	0.02	0.008	0.025	0.011	1	-0.001	0.04	0.004	0.034	0.032	0.027	0.018					
Response	0.16	-0.078	0.15	-0.2	0.25	0.12	0.24	0.11	0.12	0.14	0.0032	0.15	0.22	0.03	0.026	0.25	0.18	0.32	0.3	0.17	0.001	0.141	0.18	-0.021	0.26	-0.17	-0.22	-0.2			
Customer_For	-0.028	0.058	0.009	0.031	0.15	0.06	0.071	0.078	0.076	0.15	0.2	0.17	0.091	0.1	0.26	0.006	0.014	0.023	0.006	0.042	0.18	1	-0.021	0.14	-0.035	0.028	0.0052				
Age	0.2	-0.24	0.36	0.016	0.16	0.013	0.034	0.041	0.022	0.06	0.066	0.16	0.13	0.14	-0.12	0.061	0.07	-0.019	0.012	0.007	0.004	0.021	0.021	1	0.12	0.093	0.079	0.012			
Spent	0.79	-0.56	-0.14	0.02	0.89	0.61	0.85	0.64	0.61	0.53	-0.066	0.53	0.78	0.68	-0.5	0.053	0.25	0.47	0.38	0.14	-0.034	0.26	0.14	0.12	1	-0.5	-0.42	-0.52			
Children	-0.34	0.69	0.7	0.018	0.35	-0.4	-0.5	-0.43	-0.39	-0.27	0.44	-0.15	-0.44	-0.32	0.42	-0.02															

Data Preprocessing

Label Encoding the categorial features

Label encoding is a technique used to transform categorical data into numerical data. This is done to represent the data in a way that can be easily used by machine learning algorithms.

Scaling the features using the standard scalar

Scaling is a common data preprocessing step in machine learning that involves transforming features to have the same scale, typically with the goal of improving the performance of machine learning algorithms.

Creating a subset dataframe for dimentionality reduction

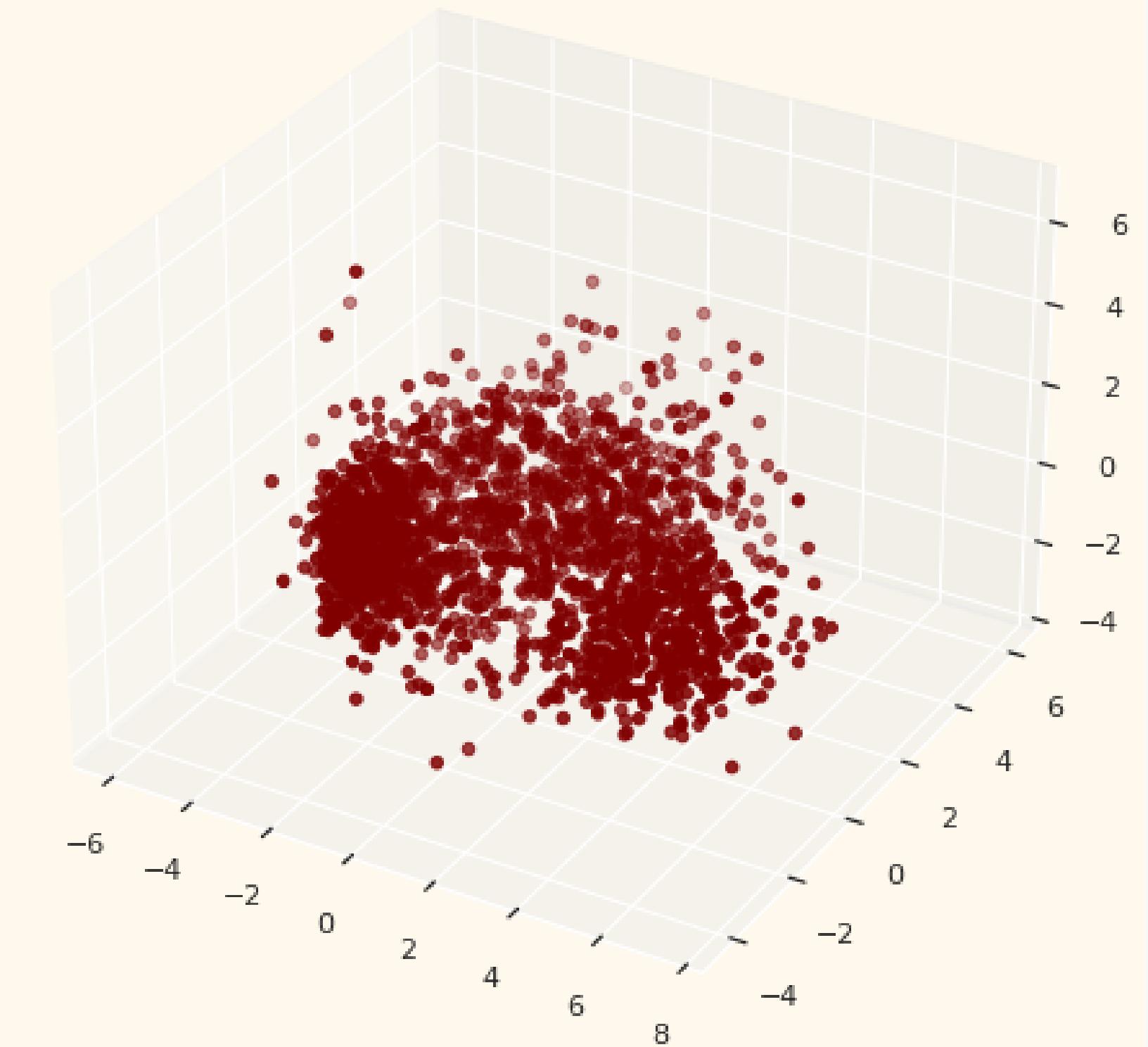
We create a subset of a dataset for dimensionality reduction when we want to reduce the number of features in the dataset while retaining as much of the original information as possible.

Dimentionality Reduction

We shall apply dimensionality reduction using PCA(Principle Component Analysis)

After applying PCA, we can visualize the data in the new reduced-dimension shape to gain insights into the structure of the data. In PCA, the original high dimension dataset is projected into the lower - dimension sub space

A 3D Projection Of Data In The Reduced Dimension



Data Preprocessing

Elbow method to determine the number of clusters to be formed

Elbow method is a common technique that is used to determine the optimal number of cluster for a given dataset. As the number of cluster increases, the within cluster Sum of Squared Error(SSE) decrease

Clustering via Agglomerative Clustering

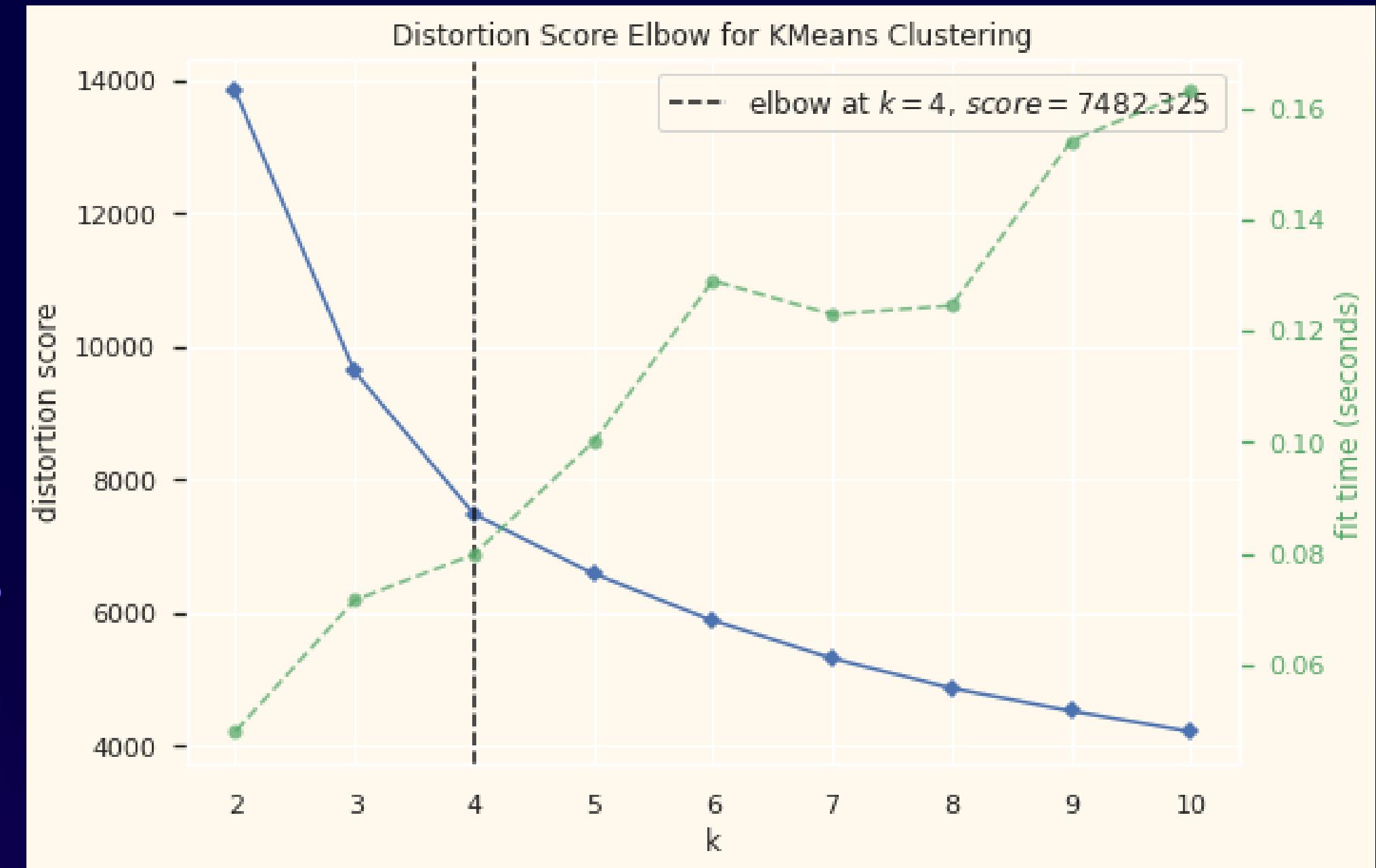
Agglomerative Clustering is a hierarchical clustering algorithm the starts by treating each data point as a separate cluster and then iteratively merges the closest pairs of clusters until all the data points belong to a single cluster.

Examining the cluster formed via Scatter Plot

I had examined the clusters formed by the 3-D distribution of the clusters.

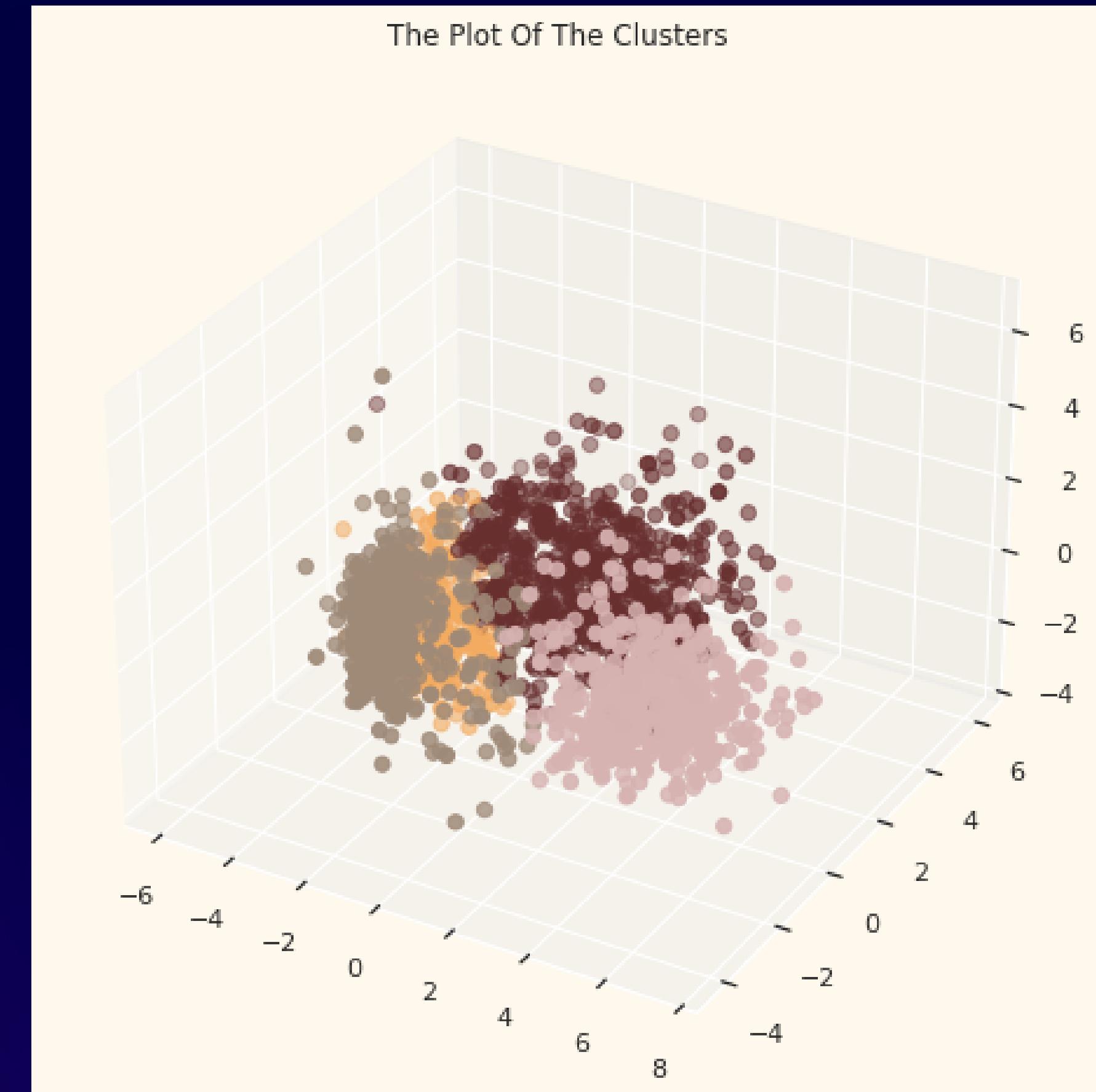
KElbowVisualizer

Elbow Method to determine the number of clusters to be formed.



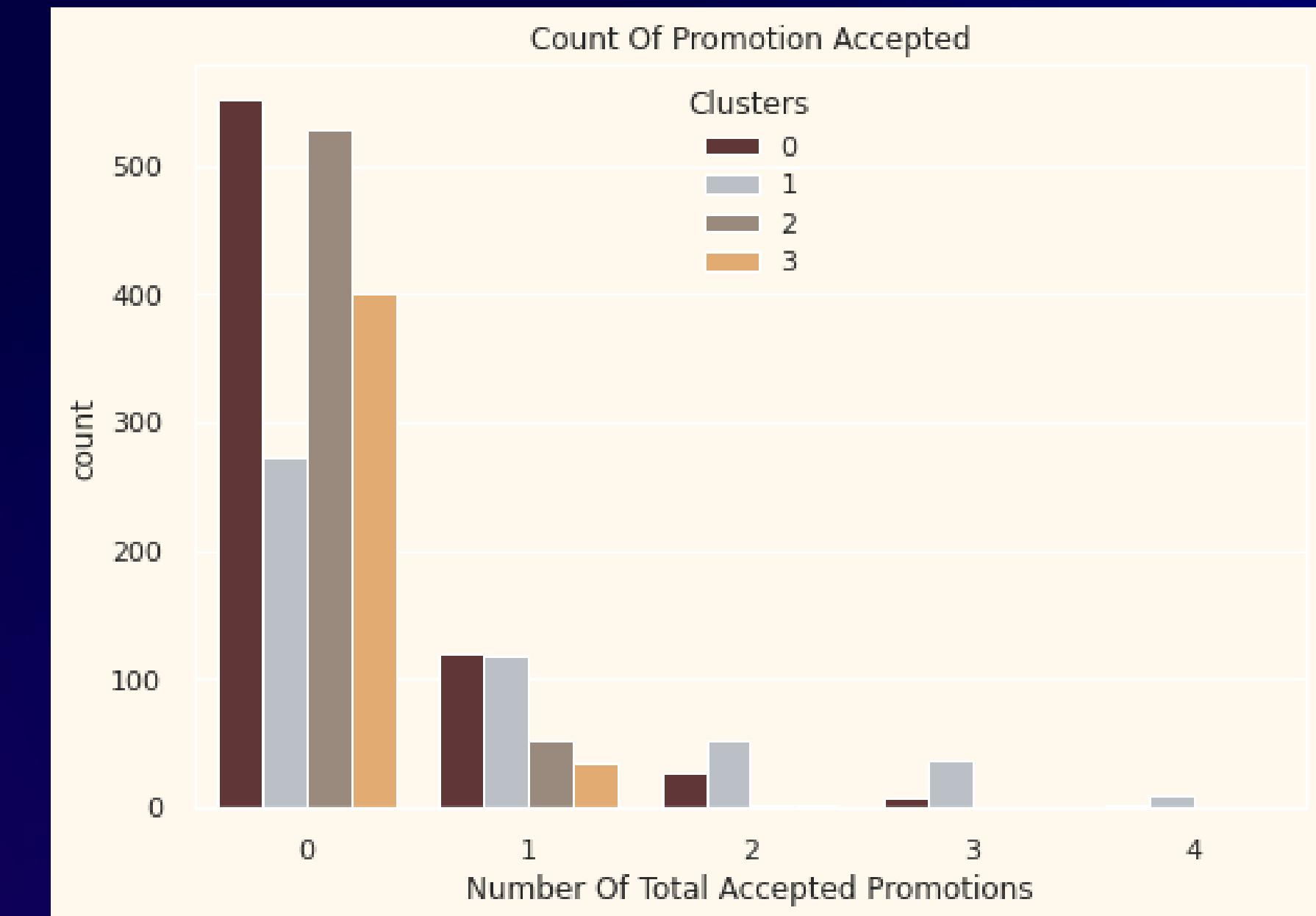
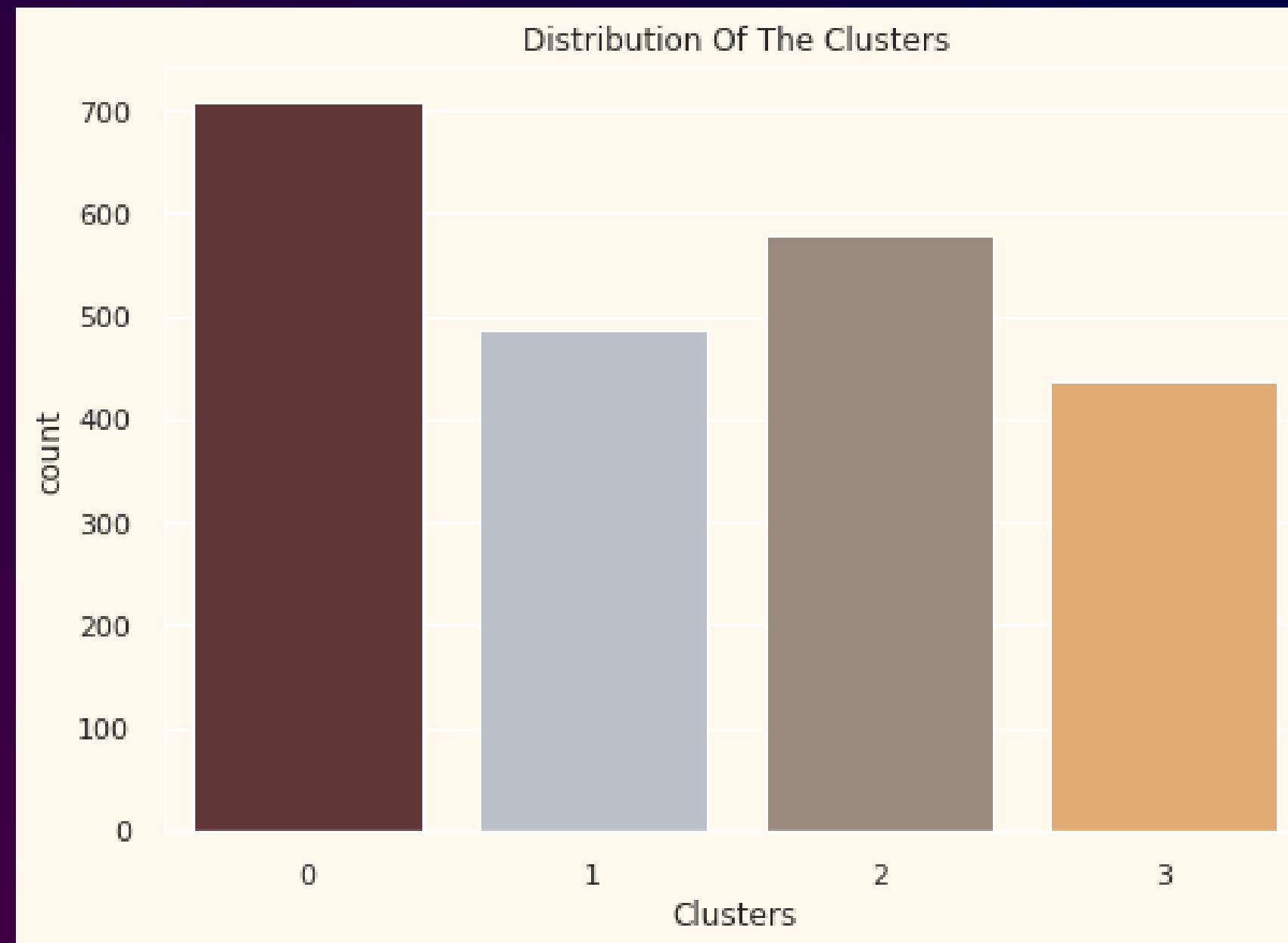
Scatter Plot

After applying clustering algorithms to a dataset, it is common to visualize the resulting clusters using scatter plots. A scatter plot can help to visualize how well the data points are separated into different clusters.



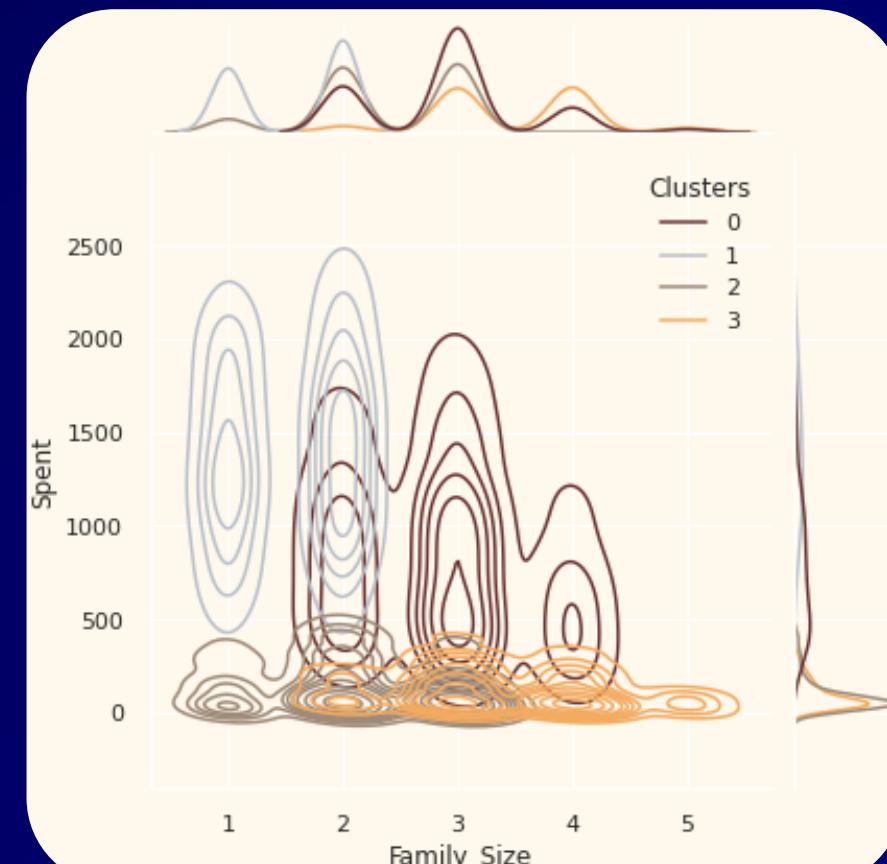
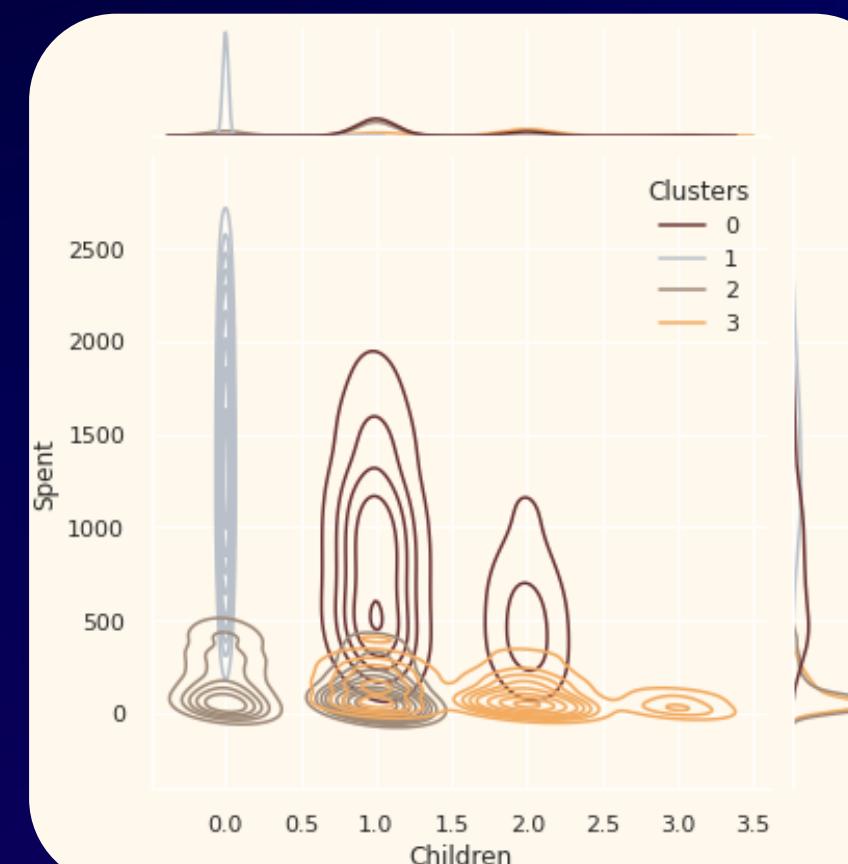
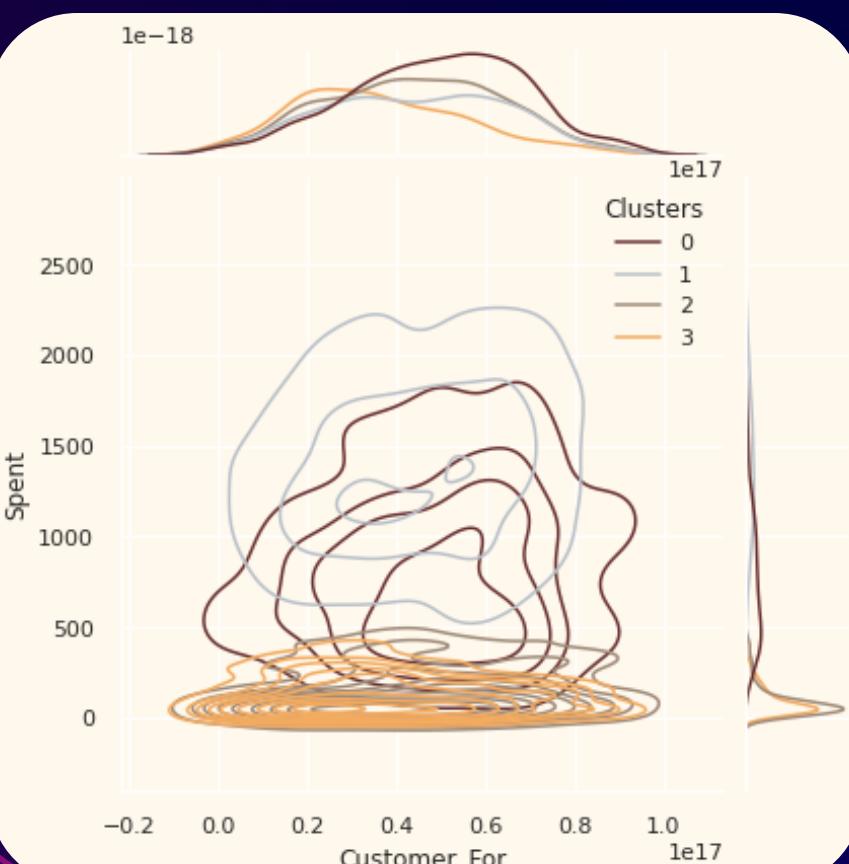
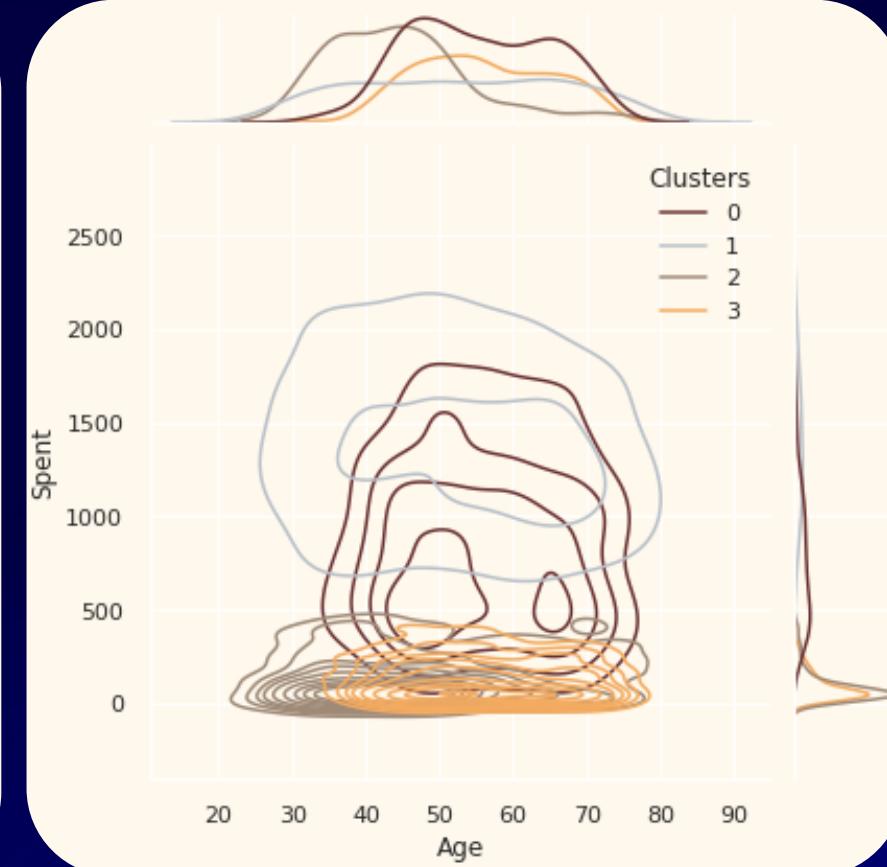
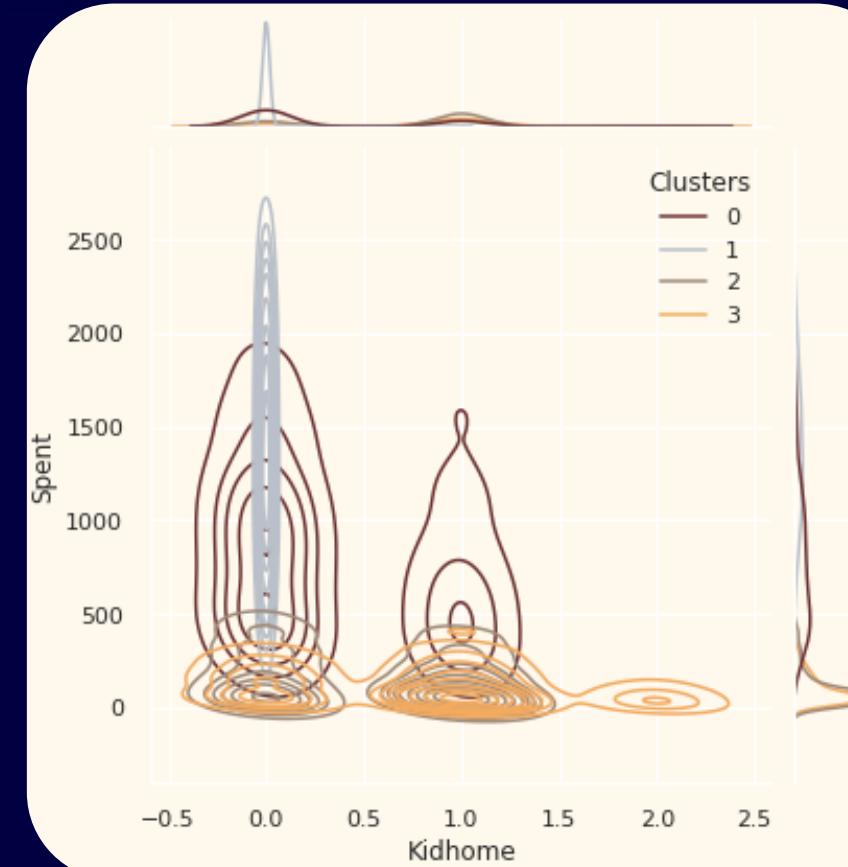
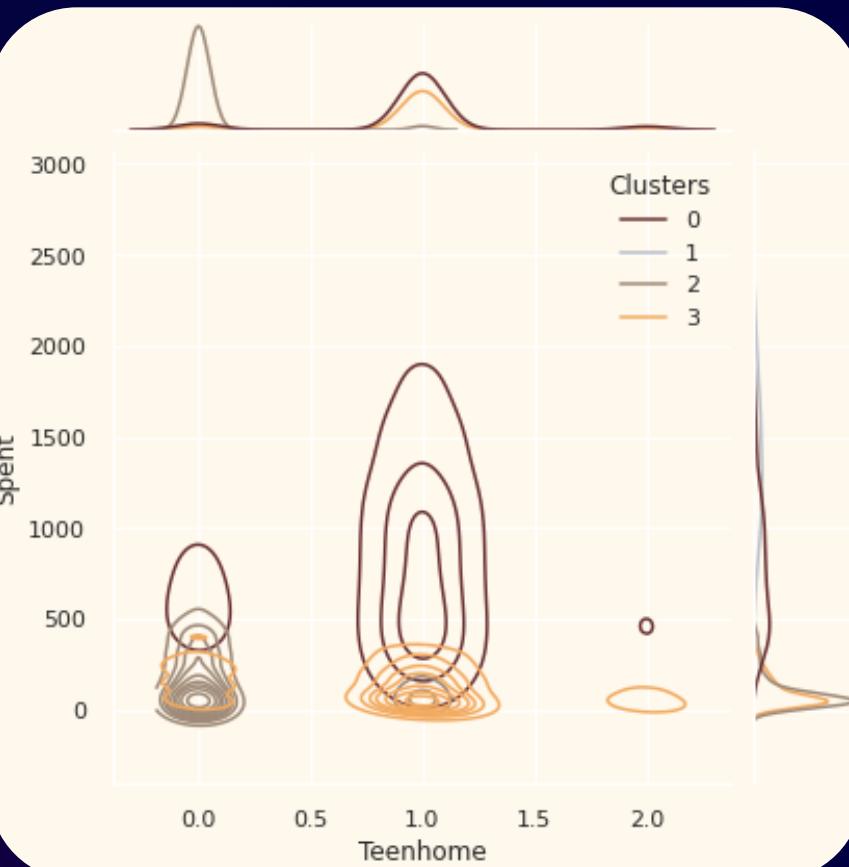
Evaluating Models

Since this is an unsupervised clustering. We do not have a tagged feature to evaluate or score our model. The purpose of this section is to study the patterns in the clusters formed and determine the nature of the clusters' patterns.



Profiling

To decide the conclusion of the clustering I will be plotting some of the features that are indicative of the customer's personal traits in light of the cluster they are in. On the basis of the outcomes, I will be arriving at the conclusions.



Points to be noted after profiling



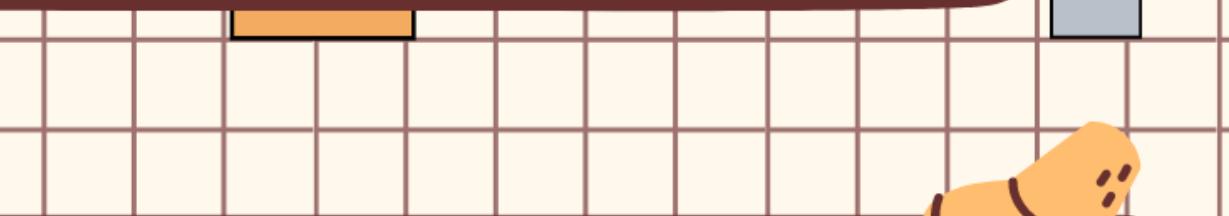
0

About Cluster Number : 0

- Are definitely a parent
- At the max have 4 members in the family and at least 2
- Single parents are a subset of this group
- Most have a teenager at home
- Relatively older

1

Profiling The Clusters



About Cluster Number : 1

- Are definitely not a parent
- At the max are only 2 members in the family
- A slight majority of couples over single people
- Span all ages
- A high income group

2

About Cluster Number : 2

- The majority of these people are parents
- At the max are 3 members in the family
- They majorly have one kid (and not teenagers, typically)
- Relatively younger



3

About Cluster Number : 3

- They are definitely a parent
- At the max are 5 members in the family and at least 2
- Majority of them have a teenager at home
- Relatively older
- A lower-income group



UNSUPERVISED CLUSTERING



Conclusion

In conclusion, the project 'Customer Segmentation and Profiling using Clustering Technique' is a valuable approach for businesses to gain insights into their customer base. By using unsupervised machine learning to group customers based on shared characteristics, businesses can tailor their marketing strategies and product offerings to better meet customer demands. This can lead to improved customer satisfaction and retention, which ultimately translates to increased revenue and a competitive advantage in the market. Overall, this project highlights the importance of utilizing data-driven techniques to understand and cater to customer needs in today's competitive business landscape.

Future Enhancement

- 1 One possibility is to incorporate additional data sources, such as social media activity, to gain a more comprehensive understanding of customer behavior and preferences. This could provide more nuanced insights into customer needs and allow for even more targeted marketing strategies.
- 2 Another potential enhancement is to explore more advanced clustering techniques, such as hierarchical clustering or density-based clustering. These techniques may be better suited for datasets with complex or non-linear relationships.
- 3 Additionally, incorporating predictive modeling techniques could help businesses anticipate future customer behavior and tailor their strategies accordingly. For example, businesses could use clustering to identify customers who are likely to churn and develop retention strategies to prevent this.