

1. Introduction

In this project, I will be performing an unsupervised clustering of data on the customer's records from a groceries firm's database. Customer segmentation is the practice of separating customers into groups that reflect similarities among customers in each cluster. I will divide customers into segments to optimize the significance of each customer to the business. To modify products according to distinct needs and behaviours of the customers. It also helps the business to cater to the concerns of different types of customers.

IMPORTING LIBRARIES

```
#Importing the Libraries
import numpy as np
import pandas as pd
import datetime
import matplotlib
import matplotlib.pyplot as plt
from matplotlib import colors
import seaborn as sns
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from yellowbrick.cluster import KElbowVisualizer
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt, numpy as np
from mpl_toolkits.mplot3d import Axes3D
from sklearn.cluster import AgglomerativeClustering
from matplotlib.colors import ListedColormap
from sklearn import metrics
import warnings
import sys
if not sys.warnoptions:
    warnings.simplefilter("ignore")
np.random.seed(42)
```

LOADING DATA

```
#Loading the dataset
data = pd.read_csv("/content/marketing_campaign.csv", sep="\t")
print("Number of datapoints:", len(data))
data.head()
```

Number of datapoints: 2240

	ID	Year_Birth	Education	Marital_Status	Income	Kidhome	Teenhome	Dt_Customer
0	5524	1957	Graduation	Single	58138.0	0	0	04-09-201
1	2174	1954	Graduation	Single	46344.0	1	1	08-03-201
2	4141	1965	Graduation	Together	71613.0	0	0	21-08-201
3	6182	1984	Graduation	Together	26646.0	1	0	10-02-201
4	5324	1981	PhD	Married	58293.0	1	0	19-01-201

5 rows × 29 columns

DATA CLEANING

In this section

Data Cleaning Feature Engineering

```
#Information on features
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2240 entries, 0 to 2239
Data columns (total 29 columns):
#   Column                Non-Null Count  Dtype
---  -
0   ID                    2240 non-null  int64
1   Year_Birth            2240 non-null  int64
2   Education             2240 non-null  object
```

```

3  Marital_Status      2240 non-null  object
4  Income              2216 non-null  float64
5  Kidhome             2240 non-null  int64
6  Teenhome            2240 non-null  int64
7  Dt_Customer         2240 non-null  object
8  Recency              2240 non-null  int64
9  MntWines            2240 non-null  int64
10 MntFruits            2240 non-null  int64
11 MntMeatProducts     2240 non-null  int64
12 MntFishProducts     2240 non-null  int64
13 MntSweetProducts    2240 non-null  int64
14 MntGoldProds        2240 non-null  int64
15 NumDealsPurchases   2240 non-null  int64
16 NumWebPurchases     2240 non-null  int64
17 NumCatalogPurchases 2240 non-null  int64
18 NumStorePurchases   2240 non-null  int64
19 NumWebVisitsMonth    2240 non-null  int64
20 AcceptedCmp3        2240 non-null  int64
21 AcceptedCmp4        2240 non-null  int64
22 AcceptedCmp5        2240 non-null  int64
23 AcceptedCmp1        2240 non-null  int64
24 AcceptedCmp2        2240 non-null  int64
25 Complain            2240 non-null  int64
26 Z_CostContact        2240 non-null  int64
27 Z_Revenue           2240 non-null  int64
28 Response            2240 non-null  int64
dtypes: float64(1), int64(25), object(3)
memory usage: 507.6+ KB

```

From the above output, we can conclude and note that: There are missing values in income Dt_Customer that indicates the date a customer joined the database is not parsed as DateTime There are some categorical features in our data frame; as there are some features in dtype: object). So we will need to encode them into numeric forms later. First of all, for the missing values, I am simply going to drop the rows that have missing income values.

```

#To remove the NA values
data = data.dropna()
print("The total number of data-points after removing the rows with missing values are:", len(data))

```

→ The total number of data-points after removing the rows with missing values are: 2216

In the next step, I am going to create a feature out of "Dt_Customer" that indicates the number of days a customer is registered in the firm's database. However, in order to keep it simple, I am taking this value relative to the most recent customer in the record.

Thus to get the values I must check the newest and oldest recorded dates.

```
data["Dt_Customer"].head()
```

```

→ 0    04-09-2012
   1    08-03-2014
   2    21-08-2013
   3    10-02-2014
   4    19-01-2014
   Name: Dt_Customer, dtype: object

```

```
data["Dt_Customer"] = pd.to_datetime(data["Dt_Customer"], format="%d-%m-%Y")
```

```
data["Dt_Customer"] = pd.to_datetime(data["Dt_Customer"], format="%d-%m-%Y", dayfirst=True)
```

```
data["Dt_Customer"] = pd.to_datetime(data["Dt_Customer"], infer_datetime_format=True)
```

```

→ <ipython-input-9-0ca5f9b53590>:1: UserWarning: The argument 'infer_datetime_format' is deprecated and will be removed in a future version
data["Dt_Customer"] = pd.to_datetime(data["Dt_Customer"], infer_datetime_format=True)

```

```

data["Dt_Customer"] = pd.to_datetime(data["Dt_Customer"])
dates = []
for i in data["Dt_Customer"]:
    i = i.date()
    dates.append(i)
#Dates of the newest and oldest recorded customer
print("The newest customer's enrolment date in therecords:",max(dates))
print("The oldest customer's enrolment date in the records:",min(dates))

```

→ The newest customer's enrolment date in therecords: 2014-06-29
The oldest customer's enrolment date in the records: 2012-07-30

Creating a feature ("Customer_For") of the number of days the customers started to shop in the store relative to the last recorded date

```
#Created a feature "Customer_For"
days = []
d1 = max(dates) #taking it to be the newest customer
for i in dates:
    delta = d1 - i
    days.append(delta)
data["Customer_For"] = days
data["Customer_For"] = pd.to_numeric(data["Customer_For"], errors="coerce")
```

Now we will be exploring the unique values in the categorical features to get a clear idea of the data.

```
print("Total categories in the feature Marital_Status:\n", data["Marital_Status"].value_counts(), "\n")
print("Total categories in the feature Education:\n", data["Education"].value_counts())
```

```

Total categories in the feature Marital_Status:
Marital_Status
Married      857
Together     573
Single       471
Divorced     232
Widow        76
Alone         3
Absurd        2
YOLO         2
Name: count, dtype: int64

```

```
Total categories in the feature Education:
Education
Graduation    1116
PhD            481
Master        365
2n Cycle      200
Basic         54
Name: count, dtype: int64
```

Extract the "Age" of a customer by the "Year_Birth" indicating the birth year of the respective person. Create another feature "Spent" indicating the total amount spent by the customer in various categories over the span of two years. Create another feature "Living_With" out of "Marital_Status" to extract the living situation of couples. Create a feature "Children" to indicate total children in a household that is, kids and teenagers. To get further clarity of household, Creating feature indicating "Family_Size" Create a feature "Is_Parent" to indicate parenthood status Lastly, I will create three categories in the "Education" by simplifying its value counts. Dropping some of the redundant features

```
#Feature Engineering
#Age of customer today
data["Age"] = 2021-data["Year_Birth"]

#Total spendings on various items
data["Spent"] = data["MntWines"]+ data["MntFruits"]+ data["MntMeatProducts"]+ data["MntFishProducts"]+ data["MntSweetProducts"]+ data["MntOtherProducts"]

#Deriving living situation by marital status"Alone"
data["Living_With"]=data["Marital_Status"].replace({"Married":"Partner", "Together":"Partner", "Absurd":"Alone", "Widow":"Alone", "YOLO":"Alone"})

#Feature indicating total children living in the household
data["Children"]=data["Kidhome"]+data["Teenhome"]

#Feature for total members in the household
data["Family_Size"] = data["Living_With"].replace({"Alone": 1, "Partner":2})+ data["Children"]

#Feature pertaining parenthood
data["Is_Parent"] = np.where(data.Children> 0, 1, 0)

#Segmenting education levels in three groups
data["Education"]=data["Education"].replace({"Basic":"Undergraduate", "2n Cycle":"Undergraduate", "Graduation":"Graduate", "Master":"Postgraduate", "PhD":"Postgraduate"})

#For clarity
data=data.rename(columns={"MntWines": "Wines", "MntFruits":"Fruits", "MntMeatProducts":"Meat", "MntFishProducts":"Fish", "MntSweetProducts":"Sweet", "MntOtherProducts":"Other"})

#Dropping some of the redundant features
to_drop = ["Marital_Status", "Dt_Customer", "Z_CostContact", "Z_Revenue", "Year_Birth", "ID"]
data = data.drop(to_drop, axis=1)
```

Now that we have some new features let's have a look at the data's stats.

```
data.describe()
```



	Income	Kidhome	Teenhome	Recency	Wines	Fruits
count	2216.000000	2216.000000	2216.000000	2216.000000	2216.000000	2216.000000
mean	52247.251354	0.441787	0.505415	49.012635	305.091606	26.356047
std	25173.076661	0.536896	0.544181	28.948352	337.327920	39.793917
min	1730.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	35303.000000	0.000000	0.000000	24.000000	24.000000	2.000000
50%	51381.500000	0.000000	0.000000	49.000000	174.500000	8.000000
75%	68522.000000	1.000000	1.000000	74.000000	505.000000	33.000000
max	666666.000000	2.000000	2.000000	99.000000	1493.000000	199.000000

8 rows × 28 columns

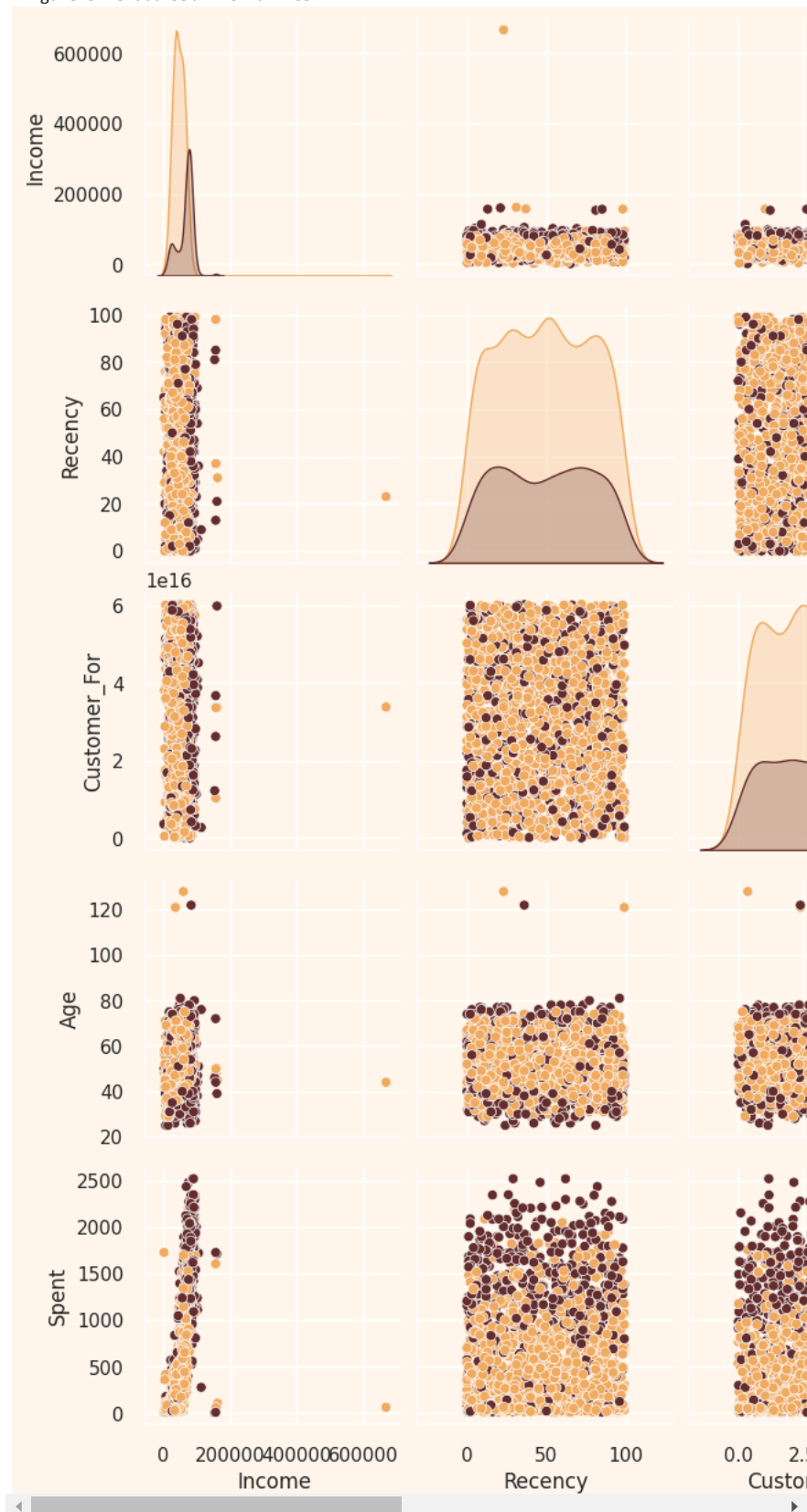
The above stats show some discrepancies in mean Income and Age and max Income and age.

Do note that max-age is 128 years, As I calculated the age that would be today (i.e. 2021) and the data is old.

I must take a look at the broader view of the data. I will plot some of the selected features.

```
#To plot some selected features
#Setting up colors preferences
sns.set(rc={"axes.facecolor":"#FFF9ED","figure.facecolor":"#FFF9ED"})
pallet = ["#682F2F", "#9E726F", "#D6B2B1", "#B9C0C9", "#9F8A78", "#F3AB60"]
cmap = colors.ListedColormap(["#682F2F", "#9E726F", "#D6B2B1", "#B9C0C9", "#9F8A78", "#F3AB60"])
#Plotting following features
To_Plot = [ "Income", "Recency", "Customer_For", "Age", "Spent", "Is_Parent"]
print("Relative Plot Of Some Selected Features: A Data Subset")
plt.figure()
sns.pairplot(data[To_Plot], hue= "Is_Parent",palette= (["#682F2F","#F3AB60"]))
#Taking hue
plt.show()
```

Relative Plot Of Some Selected Features: A Data Subset
<Figure size 800x550 with 0 Axes>



Clearly, there are a few outliers in the Income and Age features. I will be deleting the outliers in the data.

```
#Dropping the outliers by setting a cap on Age and income.
data = data[(data["Age"]<90)]
data = data[(data["Income"]<600000)]
print("The total number of data-points after removing the outliers are:", len(data))
```

The total number of data-points after removing the outliers are: 2212

Next, let us look at the correlation amongst the features. (Excluding the categorical attributes at this point)

```
data["Education"] = data["Education"].astype("category")
```

```
data = data.drop("Education", axis=1)
```

DATA PREPROCESSING

```
#Get list of categorical variables
```

```
s = (data.dtypes == 'object')
```

```
object_cols = list(s[s].index)
```

```
print("Categorical variables in the dataset:", object_cols)
```

```
↗ Categorical variables in the dataset: ['Living_With']
```

```
#Label Encoding the object dtypes.
```

```
LE=LabelEncoder()
```

```
for i in object_cols:
```

```
    data[i]=data[[i]].apply(LE.fit_transform)
```

```
print("All features are now numerical")
```

```
↗ All features are now numerical
```

```
#Creating a copy of data
```

```
ds = data.copy()
```

```
# creating a subset of dataframe by dropping the features on deals accepted and promotions
```

```
cols_del = ['AcceptedCmp3', 'AcceptedCmp4', 'AcceptedCmp5', 'AcceptedCmp1', 'AcceptedCmp2', 'Complain', 'Response']
```

```
ds = ds.drop(cols_del, axis=1)
```

```
#Scaling
```

```
scaler = StandardScaler()
```

```
scaler.fit(ds)
```

```
scaled_ds = pd.DataFrame(scaler.transform(ds), columns= ds.columns )
```

```
print("All features are now scaled")
```

```
↗ All features are now scaled
```

```
#Scaled data to be used for reducing the dimensionality
```

```
print("Dataframe to be used for further modelling:")
```

```
scaled_ds.head()
```

```
↗ Dataframe to be used for further modelling:
```

	Income	Kidhome	Teenhome	Recency	Wines	Fruits	Meat	Fish
0	0.287105	-0.822754	-0.929699	0.310353	0.977660	1.552041	1.690293	2.453472
1	-0.260882	1.040021	0.908097	-0.380813	-0.872618	-0.637461	-0.718230	-0.651004
2	0.913196	-0.822754	-0.929699	-0.795514	0.357935	0.570540	-0.178542	1.339513
3	-1.176114	1.040021	-0.929699	-0.795514	-0.872618	-0.561961	-0.655787	-0.504911
4	0.294307	1.040021	-0.929699	1.554453	-0.392257	0.419540	-0.218684	0.152508

```
5 rows × 22 columns
```

DIMENSIONALITY REDUCTION

In this problem, there are many factors on the basis of which the final classification will be done. These factors are basically attributes or features. The higher the number of features, the harder it is to work with it. Many of these features are correlated, and hence redundant. This is why I will be performing dimensionality reduction on the selected features before putting them through a classifier. Dimensionality reduction is the process of reducing the number of random variables under consideration, by obtaining a set of principal variables.

```
#Initiating PCA to reduce dimentions aka features to 3
```

```
pca = PCA(n_components=3)
```

```
pca.fit(scaled_ds)
```

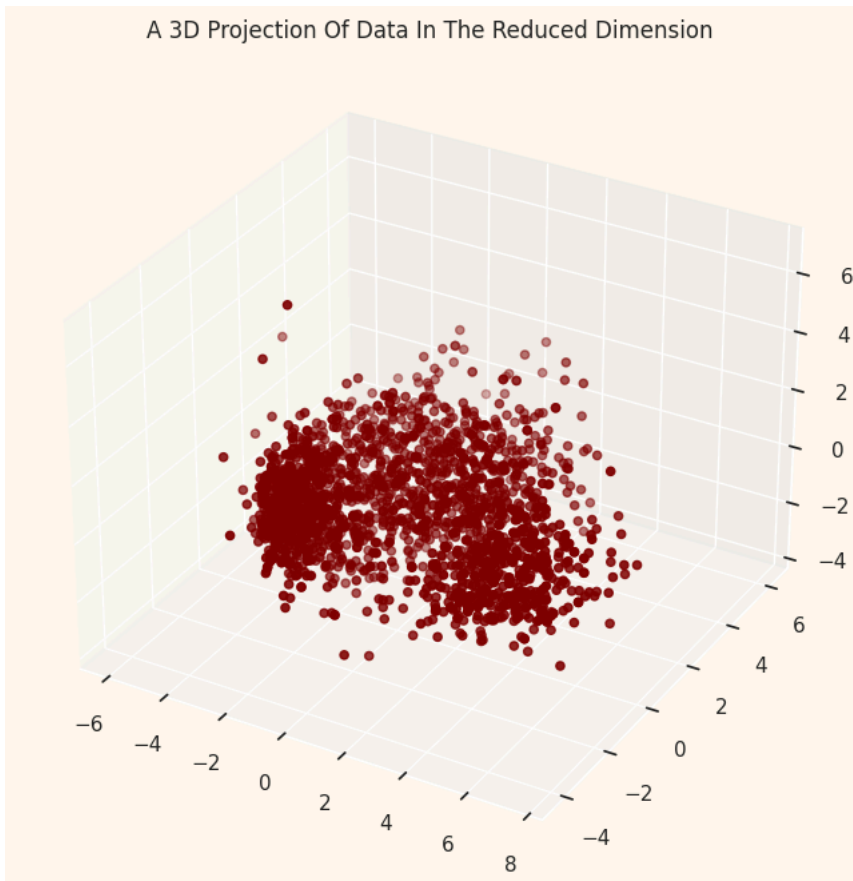
```
PCA_ds = pd.DataFrame(pca.transform(scaled_ds), columns=["col1", "col2", "col3"])
```

```
PCA_ds.describe().T
```



	count	mean	std	min	25%	50%	75%	max
col1	2212.0	2.120064e-16	2.877672	-5.965864	-2.540765	-0.771030	2.402031	7.470700
col2	2212.0	-1.927331e-17	1.707872	-4.249285	-1.314183	-0.181514	1.230197	6.205023

```
#A 3D Projection Of Data In The Reduced Dimension
x =PCA_ds["col1"]
y =PCA_ds["col2"]
z =PCA_ds["col3"]
#To plot
fig = plt.figure(figsize=(10,8))
ax = fig.add_subplot(111, projection="3d")
ax.scatter(x,y,z, c="maroon", marker="o" )
ax.set_title("A 3D Projection Of Data In The Reduced Dimension")
plt.show()
```



CLUSTERING

Now that I have reduced the attributes to three dimensions, I will be performing clustering via Agglomerative clustering. Agglomerative clustering is a hierarchical clustering method. It involves merging examples until the desired number of clusters is achieved.

Steps involved in the Clustering

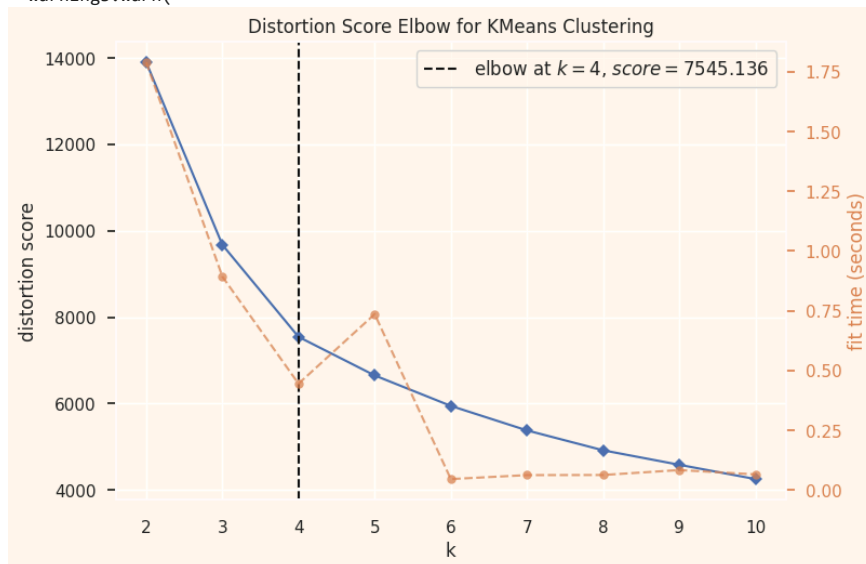
Elbow Method to determine the number of clusters to be formed Clustering via Agglomerative Clustering Examining the clusters formed via scatter plot

```
# Quick examination of elbow method to find numbers of clusters to make.
print('Elbow Method to determine the number of clusters to be formed:')
Elbow_M = KElbowVisualizer(KMeans(), k=10)
Elbow_M.fit(PCA_ds)
Elbow_M.show()
```

```

Elbow Method to determine the number of clusters to be formed:
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning
warnings.warn(

```



```

<Axes: title={'center': 'Distortion Score Elbow for KMeans Clustering'}, xlabel='k',
ylabel='distortion score'>

```

The above cell indicates that four will be an optimal number of clusters for this data. Next, we will be fitting the Agglomerative Clustering Model to get the final clusters.

```

#Initiating the Agglomerative Clustering model
AC = AgglomerativeClustering(n_clusters=4)
# fit model and predict clusters
yhat_AC = AC.fit_predict(PCA_ds)
PCA_ds["Clusters"] = yhat_AC
#Adding the Clusters feature to the original dataframe.
data["Clusters"] = yhat_AC

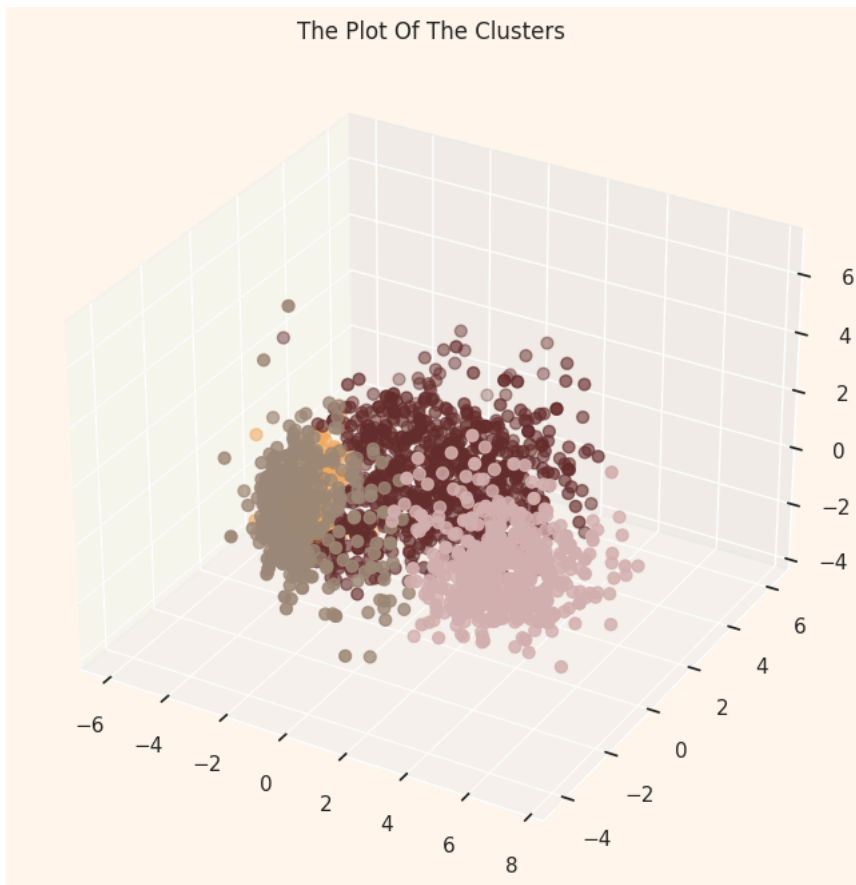
```

To examine the clusters formed let's have a look at the 3-D distribution of the clusters.

```

#Plotting the clusters
fig = plt.figure(figsize=(10,8))
ax = plt.subplot(111, projection='3d', label="bla")
ax.scatter(x, y, z, s=40, c=PCA_ds["Clusters"], marker='o', cmap = cmap )
ax.set_title("The Plot Of The Clusters")
plt.show()

```

✓ 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.

For that, we will be having a look at the data in light of clusters via exploratory data analysis and drawing conclusions.

Firstly, let us have a look at the group distribution of clustering

```
#Plotting countplot of clusters
pal = ["#682F2F", "#B9C0C9", "#9F8A78", "#F3AB60"]
pl = sns.countplot(x=data["Clusters"], palette= pal)
pl.set_title("Distribution Of The Clusters")
plt.show()
```

```
<ipython-input-34-2f63248a592e>:3: FutureWarning:
```

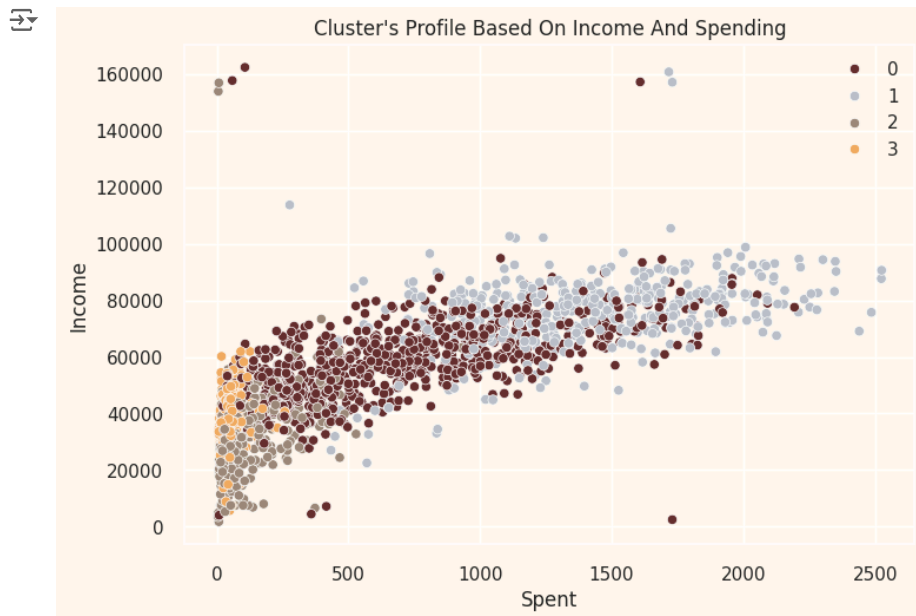
Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.

```
pl = sns.countplot(x=data["Clusters"], palette= pal)
```



The clusters seem to be fairly distributed.

```
pl = sns.scatterplot(data = data,x=data["Spent"], y=data["Income"],hue=data["Clusters"], palette= pal)
pl.set_title("Cluster's Profile Based On Income And Spending")
plt.legend()
plt.show()
```



Income vs spending plot shows the clusters pattern

group 0: high spending & average income group 1: high spending & high income group 2: low spending & low income group 3: high spending & low income Next, I will be looking at the detailed distribution of clusters as per the various products in the data. Namely: Wines, Fruits, Meat, Fish, Sweets and Gold

```
plt.figure()
pl=sns.swarmplot(x=data["Clusters"], y=data["Spent"], color= "#CBEDDD", alpha=0.5 )
pl=sns.boxenplot(x=data["Clusters"], y=data["Spent"], palette=pl)
plt.show()
```

```
↳ /usr/local/lib/python3.10/dist-packages/seaborn/categorical.py:3398: UserWarning: 25.  
  warnings.warn(msg, UserWarning)  
/usr/local/lib/python3.10/dist-packages/seaborn/categorical.py:3398: UserWarning: 70.  
  warnings.warn(msg, UserWarning)  
/usr/local/lib/python3.10/dist-packages/seaborn/categorical.py:3398: UserWarning: 69.  
  warnings.warn(msg, UserWarning)  
<ipython-input-36-b5a4c0fdc6d1>:3: FutureWarning:  
  
Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.  
  
  pl=sns.boxenplot(x=data["Clusters"], y=data["Spent"], palette=pal)  
/usr/local/lib/python3.10/dist-packages/seaborn/categorical.py:3398: UserWarning: 35.  
  warnings.warn(msg, UserWarning)  
/usr/local/lib/python3.10/dist-packages/seaborn/categorical.py:3398: UserWarning: 8.2  
  warnings.warn(msg, UserWarning)  
/usr/local/lib/python3.10/dist-packages/seaborn/categorical.py:3398: UserWarning: 72.  
  warnings.warn(msg, UserWarning)
```