

In [1570]:

```
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)
```

In [1571]:

```
read_csv("C:\\Users\\User\\Desktop\\samsung\\project-samsung\\Final project\\ma
```

In [1572]:

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

In [1573]:

data

Out[1573]:

	ID	Year_Birth	Education	Marital_Status	Income	Kidhome	Teer
0	5524	1957	Graduation	Single	58138.0	0	
1	2174	1954	Graduation	Single	46344.0	1	
2	4141	1965	Graduation	Together	71613.0	0	
3	6182	1984	Graduation	Together	26646.0	1	
4	5324	1981	PhD	Married	58293.0	1	
...	
2235	10870	1967	Graduation	Married	61223.0	0	
2236	4001	1946	PhD	Together	64014.0	2	
2237	7270	1981	Graduation	Divorced	56981.0	0	
2238	8235	1956	Master	Together	69245.0	0	
2239	9405	1954	PhD	Married	52869.0	1	

2240 rows × 29 columns

In [1574]:

data['Z_Revenue'].value_counts()

Out[1574]:

```
11      2240
Name: Z_Revenue, dtype: int64
```

In [1575]:

```
data['MntFruits'].describe()
```

Out[1575]:

```
count    2240.000000
mean      26.302232
std       39.773434
min        0.000000
25%        1.000000
50%        8.000000
75%       33.000000
max      199.000000
Name: MntFruits, dtype: float64
```

Content

ID: Customer's unique identifier

Year_Birth: Customer's birth year

Education: Customer's education level

Marital_Status: Customer's marital status

Income: Customer's yearly household income

Kidhome: Number of children in customer's household

Teenhome: Number of teenagers in customer's household

Dt_Customer: Date of customer's enrollment with the company

Recency: Number of days since customer's last purchase

Complain: 1 if the customer complained in the last 2 years, 0 otherwise

MntWines: Amount spent on wine in last 2 years

MntFruits: Amount spent on fruits in last 2 years

MntMeatProducts: Amount spent on meat in last 2 years

MntFishProducts: Amount spent on fish in last 2 years

MntSweetProducts: Amount spent on sweets in last 2 years

MntGoldProds: Amount spent on gold in last 2 years

NumDealsPurchases: Number of purchases made with a discount

AcceptedCmp1: 1 if customer accepted the offer in the 1st campaign, 0 otherwise

AcceptedCmp2: 1 if customer accepted the offer in the 2nd campaign, 0 otherwise

AcceptedCmp3: 1 if customer accepted the offer in the 3rd campaign, 0 otherwise

AcceptedCmp4: 1 if customer accepted the offer in the 4th campaign, 0 otherwise

AcceptedCmp5: 1 if customer accepted the offer in the 5th campaign, 0 otherwise

Response: 1 if customer accepted the offer in the last campaign, 0 otherwise

NumWebPurchases: Number of purchases made through the company's website

NumCatalogPurchases: Number of purchases made using a catalogue

NumStorePurchases: Number of purchases made directly in stores

NumWebVisitsMonth: Number of visits to company's website in the last month

Questions:

1-Relationship between Date of customer's enrollment and marital status?

2-The relationship between the number of purchases and marital status?

3-The relationship between the number of purchases and the number of children and the family size?

4-What does age have to do with the number of purchases?

5-What is the relationship between education and income?

6-What is the relationship between income and the number of children?

7-What is the relationship between income and the number of purchases?

8-What is the relationship between the number of purchases from the website and the number of website visits?

9-What is the relationship between the number of purchases from a Deal with the number of purchases from the website, the number of purchases from the catalog, and the number of purchases from the store?

10-What is the relationship between the number of purchases from a Deal with accepted cmp 1 ,accepted cmp 2,accepted cmp 3 ,accepted cmp 4 ,accepted cmp 5 and Response?

11-What is the relationship between the complaint and Date of customer's enrollment?

-

-

-

-

In [1576]:

data

Out[1576]:

	ID	Year_Birth	Education	Marital_Status	Income	Kidhome	Teenhor
0	5524	1957	Graduation	Single	58138.0	0	
1	2174	1954	Graduation	Single	46344.0	1	
2	4141	1965	Graduation	Together	71613.0	0	
3	6182	1984	Graduation	Together	26646.0	1	
4	5324	1981	PhD	Married	58293.0	1	
...	
2235	10870	1967	Graduation	Married	61223.0	0	
2236	4001	1946	PhD	Together	64014.0	2	
2237	7270	1981	Graduation	Divorced	56981.0	0	
2238	8235	1956	Master	Together	69245.0	0	
2239	9405	1954	PhD	Married	52869.0	1	

2240 rows × 29 columns



Add new features or modify features to better clarify the data

In [1577]:

```
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("Date of registration of the company's newest client:",max(dates))
print("Date of registration of the company's oldest client:",min(dates))
```

Date of registration of the company's newest client: 2014-12-06

Date of registration of the company's oldest client: 2012-01-08

In [1578]:

```
d1 = max(dates) #taking it to be the newest customer
for i in dates:
    t=d1-i
    print(t)
```

```
362 days, 0:00:00
415 days, 0:00:00
473 days, 0:00:00
698 days, 0:00:00
939 days, 0:00:00
219 days, 0:00:00
700 days, 0:00:00
337 days, 0:00:00
388 days, 0:00:00
383 days, 0:00:00
465 days, 0:00:00
663 days, 0:00:00
284 days, 0:00:00
229 days, 0:00:00
511 days, 0:00:00
229 days, 0:00:00
310 days, 0:00:00
246 days, 0:00:00
727 days, 0:00:00
720 days, 0:00:00
```


In [1579]:

```
#Created a feature "Customer_From_days"
days = []
for i in dates:
    delta = d1 - i
    days.append(delta)
data["Customer_From_days"] = days
data["Customer_From_days"] = pd.to_numeric(data["Customer_From_days"], errors='coerce')
for i in range(len(data['Customer_From_days'])):
    t=0
    t=data['Customer_From_days'][i]
    data['Customer_From_days'][i]=t/60/60/24/1000000000
```

In []:

In [1580]:

```
# Create a feature that shows the age of the customer based on the date of birth
data["Age"] = 2021-data["Year_Birth"]
```

In [1581]:

data

Out[1581]:

	ID	Year_Birth	Education	Marital_Status	Income	Kidhome	Teenhor
0	5524	1957	Graduation	Single	58138.0	0	
1	2174	1954	Graduation	Single	46344.0	1	
2	4141	1965	Graduation	Together	71613.0	0	
3	6182	1984	Graduation	Together	26646.0	1	
4	5324	1981	PhD	Married	58293.0	1	
...	
2235	10870	1967	Graduation	Married	61223.0	0	
2236	4001	1946	PhD	Together	64014.0	2	
2237	7270	1981	Graduation	Divorced	56981.0	0	
2238	8235	1956	Master	Together	69245.0	0	
2239	9405	1954	PhD	Married	52869.0	1	

2240 rows × 31 columns

In [1582]:

data['Marital_Status'].value_counts()

Out[1582]:

```

Married      864
Together     580
Single       480
Divorced     232
Widow        77
Alone         3
Absurd        2
YOLO          2
Name: Marital_Status, dtype: int64

```

In [1583]:

```
# We will change the values[Alone,Absurd,YOLO] because there are few of them a
data['Marital_Status'].replace('Alone','Single',inplace=True)
data['Marital_Status'].replace('Absurd','Single',inplace=True)
data['Marital_Status'].replace('YOLO','Single',inplace=True)

data["Living_With"]=data["Marital_Status"].replace({"Married":"Partner", "Toge

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

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

In []:

In []:

Q1: Relationship between Date of customer's enrollment and marital status ?

In [1584]:

```
data['Marital_Status'].value_counts()
```

Out[1584]:

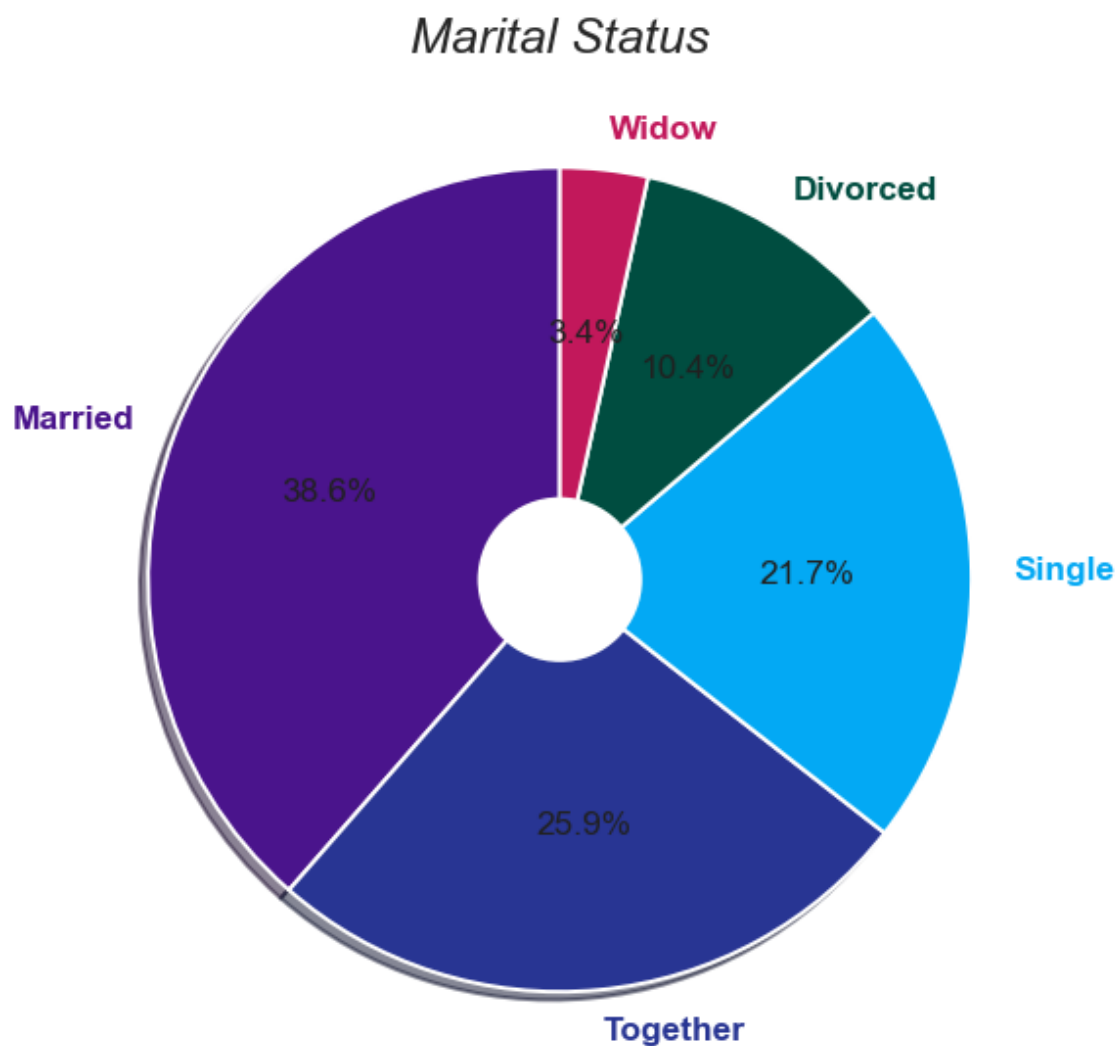
```
Married      864
Together     580
Single       487
Divorced      232
Widow         77
Name: Marital_Status, dtype: int64
```

In [1585]:

```
fig, ax = plt.subplots(figsize=(15, 8))
colors2=['#4a148c','#283593','#03a9f4','#004d40','#c2185b']
patches, texts, pcts = ax.pie(
    data['Marital_Status'].value_counts(), labels=[*data['Marital_Status'].val
    ,wedgeprops={'linewidth': 2.0, 'edgecolor': 'white'},
    textprops={'size': 'x-large'},
    startangle=90)

for i, patch in enumerate(patches):
    texts[i].set_color(patch.get_facecolor())
plt.setp(pcts, color='#212121')
plt.setp(texts, fontweight=600)
centre_circle = plt.Circle((0,0),0.20,fc='white')
plt.gcf().gca().add_artist(centre_circle)
plt.tight_layout()
plt.title(label='Marital Status',fontsize=25,fontstyle='italic')

plt.tight_layout()
```



The largest proportion of the company's customers are people who live with partners

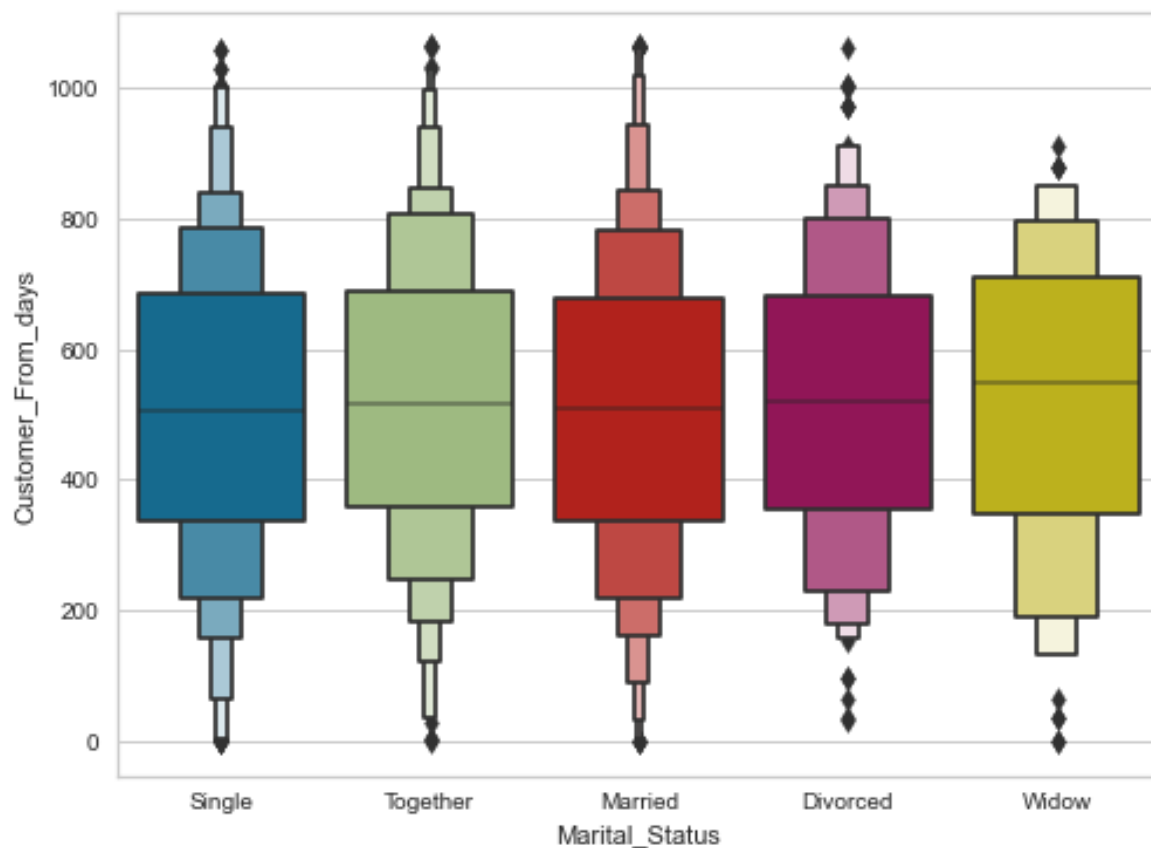
In []:

In [1586]:

```
plt.figure(figsize=(8,6))  
sns.boxenplot(data=data,x='Marital_Status',y='Customer_From_days')
```

Out[1586]:

<AxesSubplot:xlabel='Marital_Status', ylabel='Customer_From_days'>



There is no significant correlation between the marital status of customers and the date of their joining the company

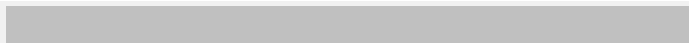
In []:

Q2: The relationship between the number of purchases and marital status?

In []:

In [1587]:

```
# Create a feature that shows the number of purchases for customers  
data['total_purchases']=data['MntFishProducts']+data["MntFruits"]+data['MntGol
```

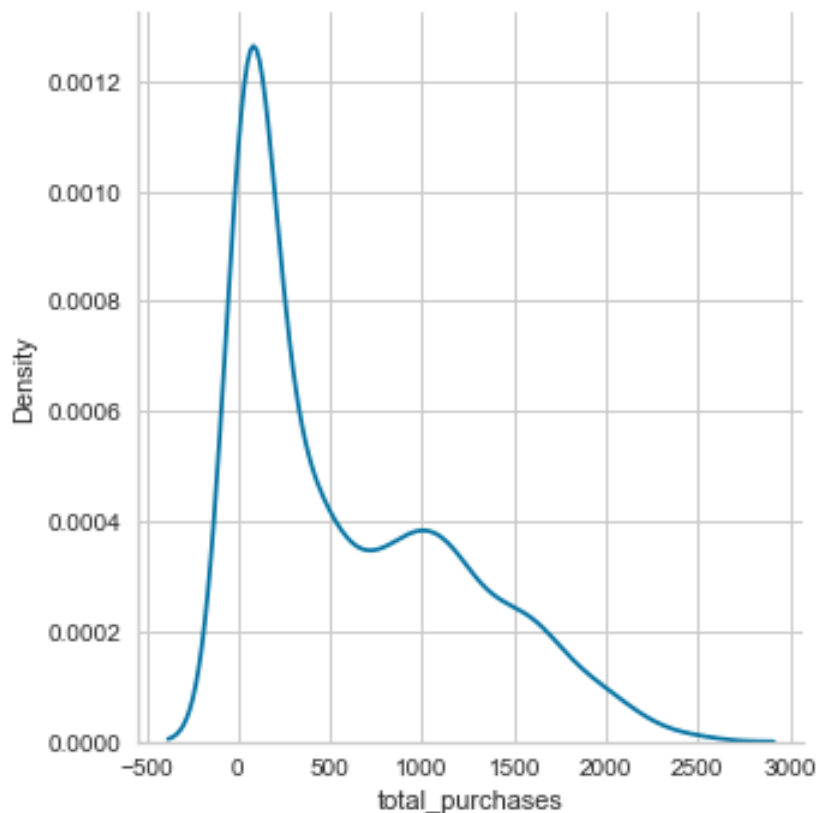


In [1588]:

```
sns.displot(data,x='total_purchases',kind='kde')
```

Out[1588]:

<seaborn.axisgrid.FacetGrid at 0x1c5d38624c0>



In [1589]:

```
print('Min:'+str(min(data['total_purchases'])),'| Max: '+str(max(data['total_p
```

Min:5 | Max: 2525

In [1590]:

```
inter=pd.interval_range(start=5,freq=210, end= 2525)  
inter
```

Out[1590]:

IntervalIndex([(5, 215], (215, 425], (425, 635], (635, 845], (845, 1055] ... (1475, 1685], (1685, 1895], (1895, 2105], (2105, 2315], (2315, 2525]], dtype='interval[int64, right]')

In [1591]:

```
# We will classify the number of purchases into more than one category
s=5
name_class=[]
for i in range(12):
    t='class ' + str(i) + " : (" +str(s)+ ", " +str(s+210) +')'
    name_class.append(t)
    s=s+210
inter=[5,215,425,635,845,1055,1265,1475,1685,1895,2105,2315,2525]
data['purchase_quantity']=pd.cut(data['total_purchases'],bins=inter,labels=name_class)
```

In [1592]:

```
data['purchase_quantity'].value_counts()
```

Out[1592]:

```
class 0 : (5, 215)          920
class 1 : (215, 425)        246
class 4 : (845, 1055)       184
class 2 : (425, 635)        177
class 5 : (1055, 1265)      172
class 3 : (635, 845)        161
class 6 : (1265, 1475)      117
class 7 : (1475, 1685)      117
class 8 : (1685, 1895)       67
class 9 : (1895, 2105)       50
class 10 : (2105, 2315)      20
class 11 : (2315, 2525)       8
Name: purchase_quantity, dtype: int64
```

In [1593]:

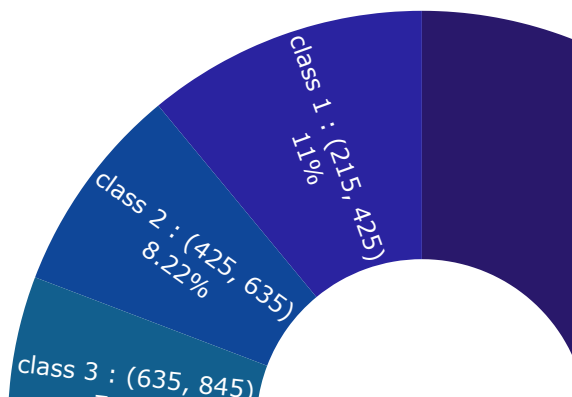
In [1594]:

```
d={'num_clas':data['purchase_quantity'].value_counts(),'clas':name_class}

import plotly.express as px

fig = px.pie(d,values='num_clas', names='clas',labels='clas',color_discrete_se
fig.update_traces(textposition='inside', textinfo='percent+label')
fig.update_traces(textposition='inside', hole=.4, hoverinfo="label+percent+nam
fig.show()
```

Buyer Categories



That the largest percentage of the number of purchases made by customers and up to 41.1% was between 5 to 215 purchases and the more purchases the less the percentage

In []:

In [1595]:

```
pre=[]
total=d['num_clas'].sum()
for i in d['num_clas']:
    n=i/total
    pre.append(n.round(3)*100)
```

In [1596]:

```
d_1={'num_clas':pre,'clas':name_class}
df=pd.DataFrame(data=d_1)
```

In [1597]:

df

Out[1597]:

	num_clas	clas
0	41.1	class 0 : (5, 215)
1	11.0	class 1 : (215, 425)
2	8.2	class 2 : (425, 635)
3	7.9	class 3 : (635, 845)
4	7.7	class 4 : (845, 1055)
5	7.2	class 5 : (1055, 1265)
6	5.2	class 6 : (1265, 1475)
7	5.2	class 7 : (1475, 1685)
8	3.0	class 8 : (1685, 1895)
9	2.2	class 9 : (1895, 2105)
10	0.9	class 10 : (2105, 2315)
11	0.4	class 11 : (2315, 2525)

In []:

In [1598]:

data.head()

Out[1598]:

	ID	Year_Birth	Education	Marital_Status	Income	Kidhome	Teenhome
0	5524	1957	Graduation	Single	58138.0	0	0
1	2174	1954	Graduation	Single	46344.0	1	1
2	4141	1965	Graduation	Together	71613.0	0	0
3	6182	1984	Graduation	Together	26646.0	1	0
4	5324	1981	PhD	Married	58293.0	1	0

5 rows × 36 columns

In [1599]:

name_class

Out[1599]:

```
[ 'class 0 : (5, 215)',
  'class 1 : (215, 425)',
  'class 2 : (425, 635)',
  'class 3 : (635, 845)',
  'class 4 : (845, 1055)',
  'class 5 : (1055, 1265)',
  'class 6 : (1265, 1475)',
  'class 7 : (1475, 1685)',
  'class 8 : (1685, 1895)',
  'class 9 : (1895, 2105)',
  'class 10 : (2105, 2315)',
  'class 11 : (2315, 2525)']
```

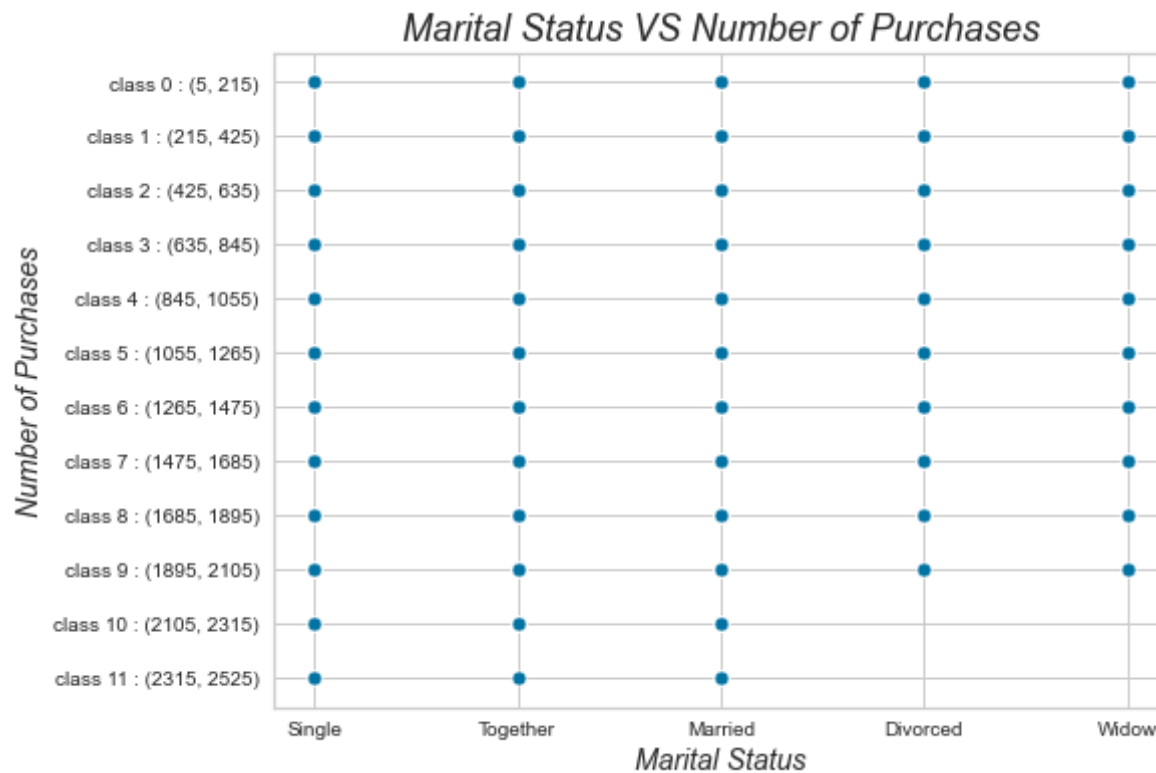
In []:

In [1600]:

```

plt.figure(figsize=(8,6))
sns.scatterplot(data=data,x='Marital_Status',y='purchase_quantity')
plt.xlabel(fontsize=14,xlabel='Marital Status',fontstyle='italic')
plt.ylabel(fontsize=14,ylabel='Number of Purchases',fontstyle='italic')
plt.title(label='Marital Status VS Number of Purchases',fontsize=18,fontstyle='italic')
plt.show()

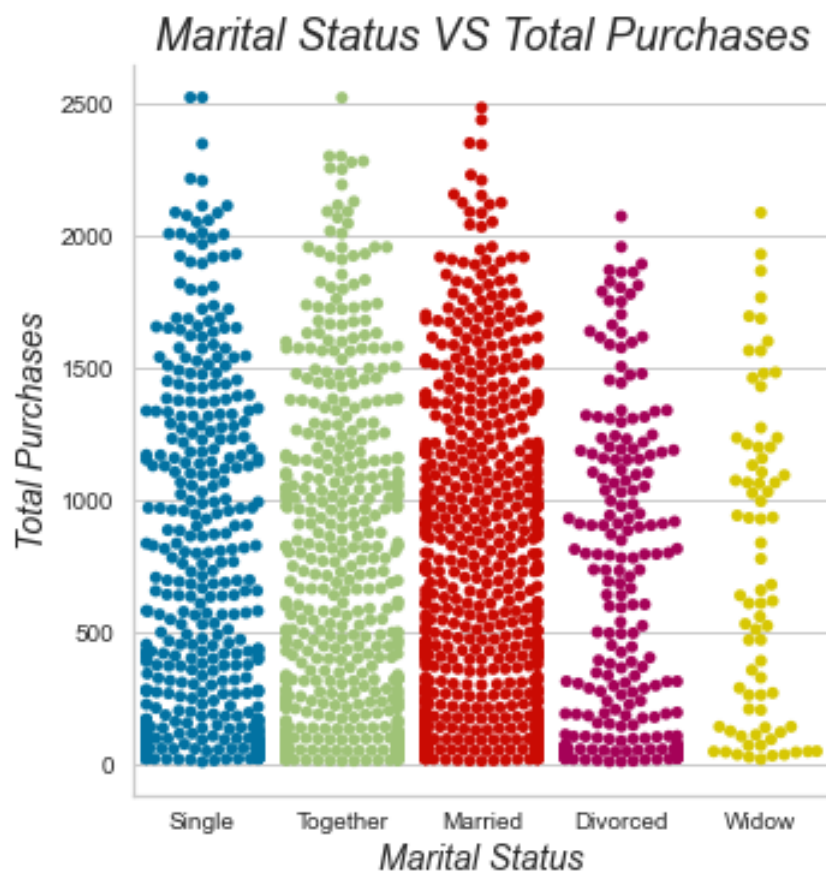
```



Divorced and widowed clients are not included in the categories 10 and 11

In [1601]:

```
sns.catplot(data=data,x='Marital_Status',y='total_purchases',kind='swarm')
plt.xlabel(fontsize=14,xlabel='Marital Status',fontstyle='italic')
plt.ylabel(fontsize=14,ylabel='Total Purchases',fontstyle='italic')
plt.title(label='Marital Status VS Total Purchases',fontsize=18,fontstyle='italic')
plt.show()
```



In []:

In []:

Q3: The relationship between the number of purchases and the number of children and the age them?

In [1602]:

```
data['Kidhome'].value_counts()
```

Out[1602]:

```
0    1293
1     899
2      48
Name: Kidhome, dtype: int64
```

In [1603]:

```
data['Teenhome'].value_counts()
```

Out[1603]:

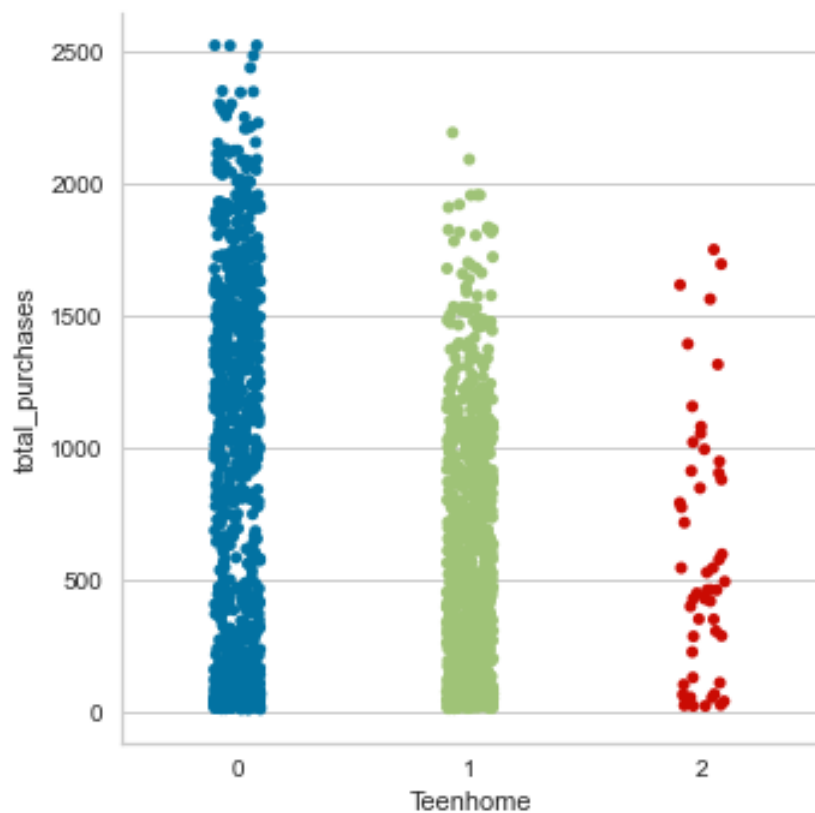
```
0    1158
1   1030
2     52
Name: Teenhome, dtype: int64
```

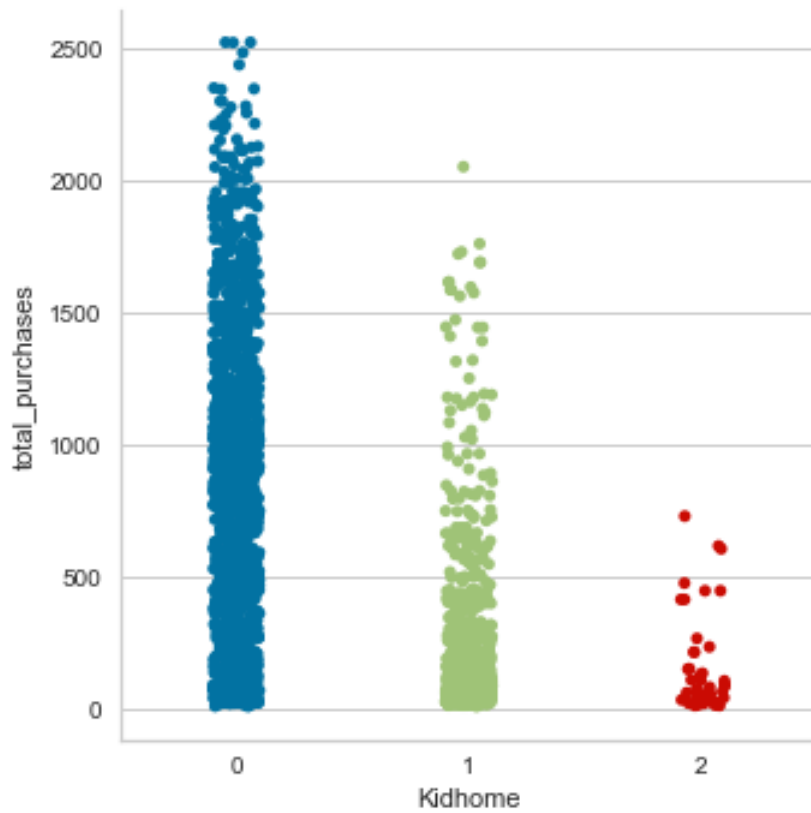
In [1604]:

```
sns.catplot(data=data,x='Teenhome',y='total_purchases',kind='strip')  
sns.catplot(data=data,x='Kidhome',y='total_purchases',kind='strip')  
plt.subplot()
```

Out[1604]:

<AxesSubplot:xlabel='Kidhome', ylabel='total_purchases'>





We see that customers with teenagers have more purchases than customers with children

In [1605]:

```
data['Family_Size'].value_counts()
```

Out[1605]:

3 889

2 764

4 301

1 254

5 32

Name: Family_Size, dtype: int64

In [1606]:

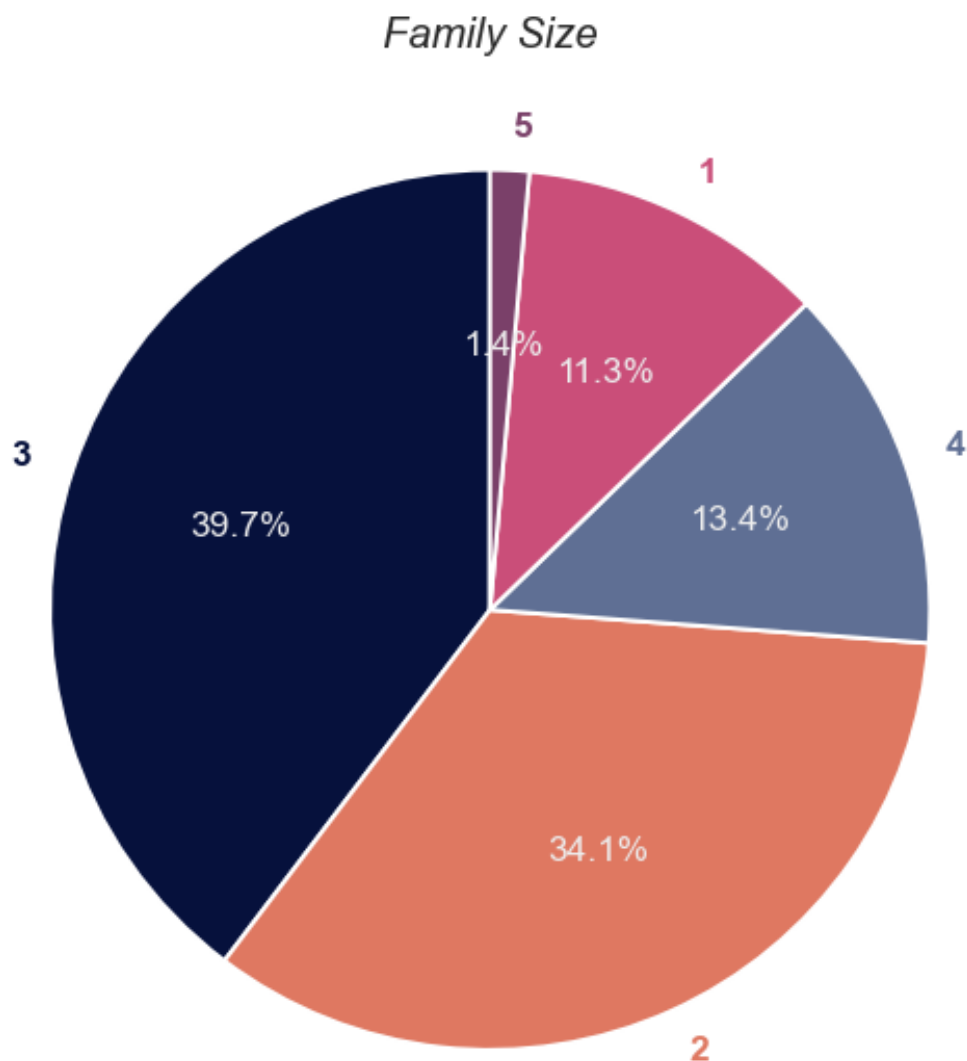
```
color_3=['#06113C', '#DF7861', '#5F6F94', '#CA4E79', '#7A4069']
```

In [1607]:

```
fig, ax = plt.subplots(figsize=(15, 8))
patches, texts, pcts = ax.pie(
    data['Family_Size'].value_counts(), labels=[*data['Family_Size'].value_cou
    ,wedgeprops={'linewidth': 2.0, 'edgecolor': 'white'},
    textprops={'size': 'x-large'},
    startangle=90)

for i, patch in enumerate(patches):
    texts[i].set_color(patch.get_facecolor())
plt.setp(pcts, color='#EEEEEE')
plt.setp(texts, fontweight=600)
plt.tight_layout()
plt.title(label='Family Size',fontsize=20,fontstyle='italic')

plt.tight_layout()
```



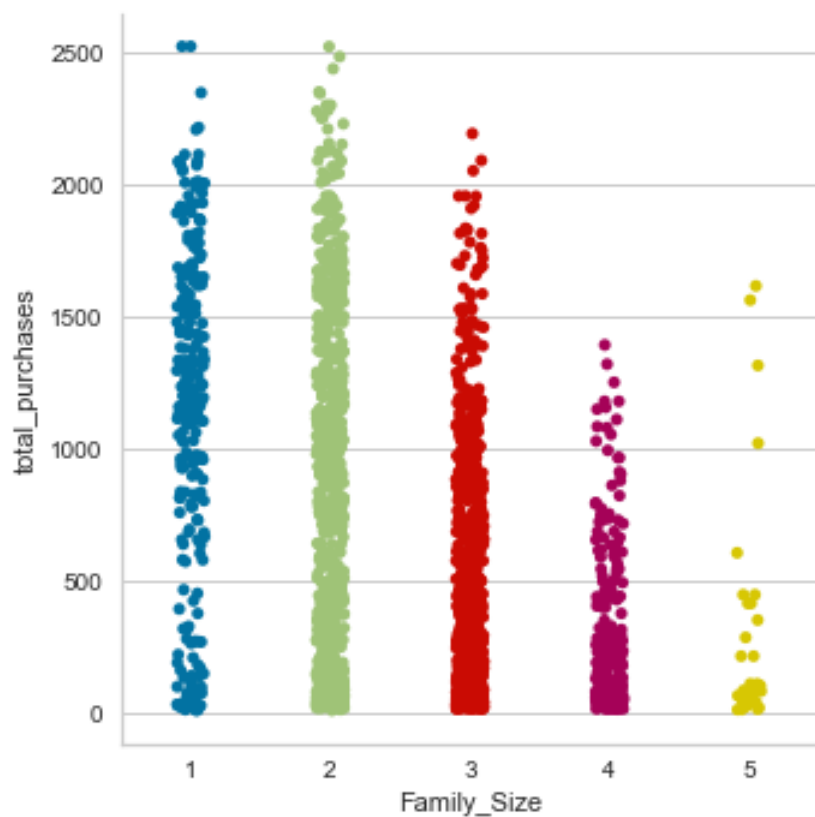
Most of the company's clients are families of two or 3

In [1608]:

```
sns.catplot(data=data,x='Family_Size',y='total_purchases',kind='strip')
```

Out[1608]:

<seaborn.axisgrid.FacetGrid at 0x1c5c56ae490>



The figure shows that the higher the number of family members, the fewer purchases

In []:

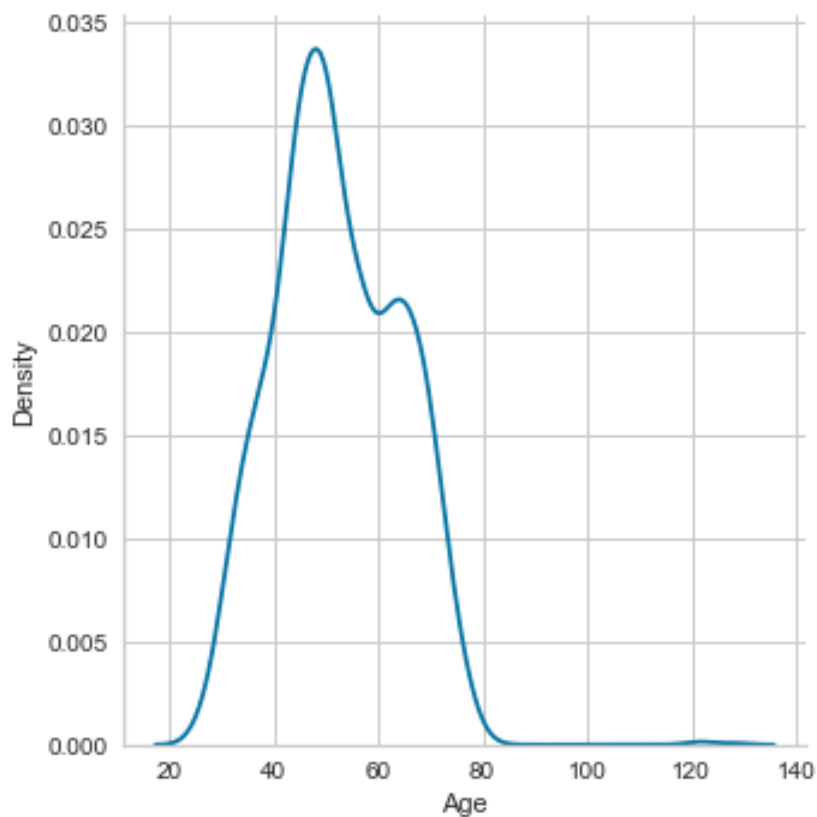
Q4: What does age have to do with the number of purchases?

In [1609]:

```
sns.displot(data,x='Age',kind='kde')
```

Out[1609]:

<seaborn.axisgrid.FacetGrid at 0x1c5c4fe9e80>



In [1610]:

```
print("Oldest customer:",max(data['Age']))  
print("Youngest customer:",min(data['Age']))
```

Oldest customer: 128

Youngest customer: 25

In [1611]:

```
data['Age'].mean()
```

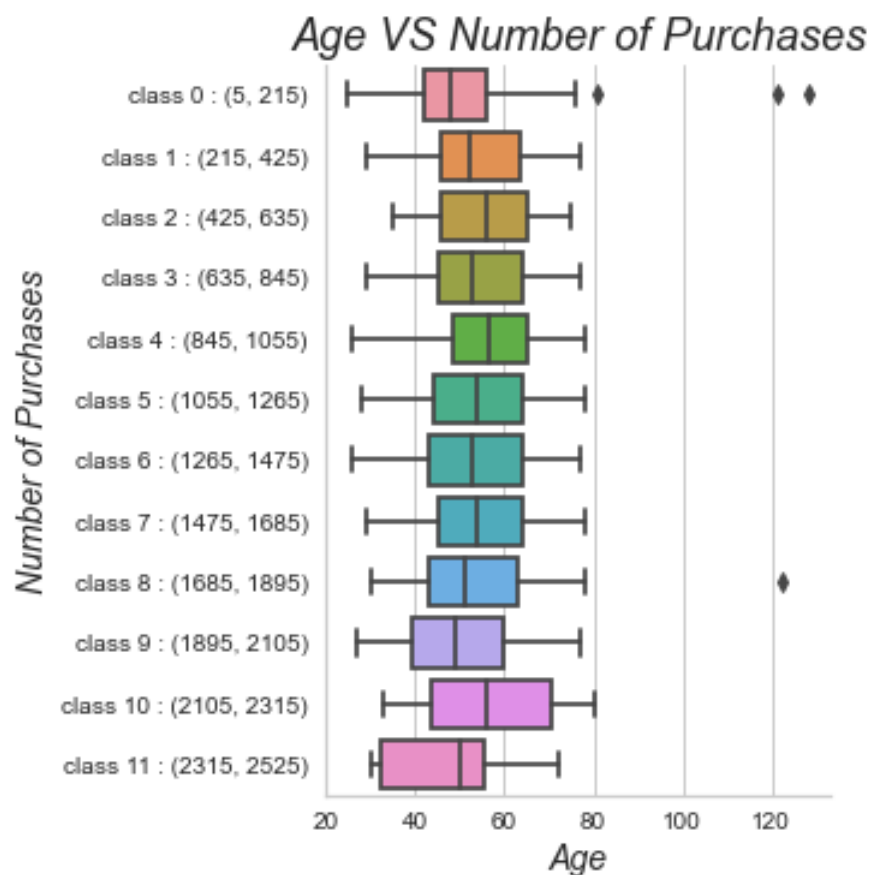
Out[1611]:

52.19419642857143

Average customer age 52

In [1612]:

```
sns.catplot(data=data, x='Age', y='purchase_quantity', kind='box')
plt.xlabel(fontsize=14, xlabel='Age', fontstyle='italic')
plt.ylabel(fontsize=14, ylabel='Number of Purchases', fontstyle='italic')
plt.title(label='Age VS Number of Purchases', fontsize=18, fontstyle='italic')
plt.show()
```



There is no relationship between age and number of purchases

In []:

Q5: What is the relationship between education and income?

In [1613]:

```
data['Education'].value_counts()
```

Out[1613]:

Graduation	1127
PhD	486
Master	370
2n Cycle	203
Basic	54

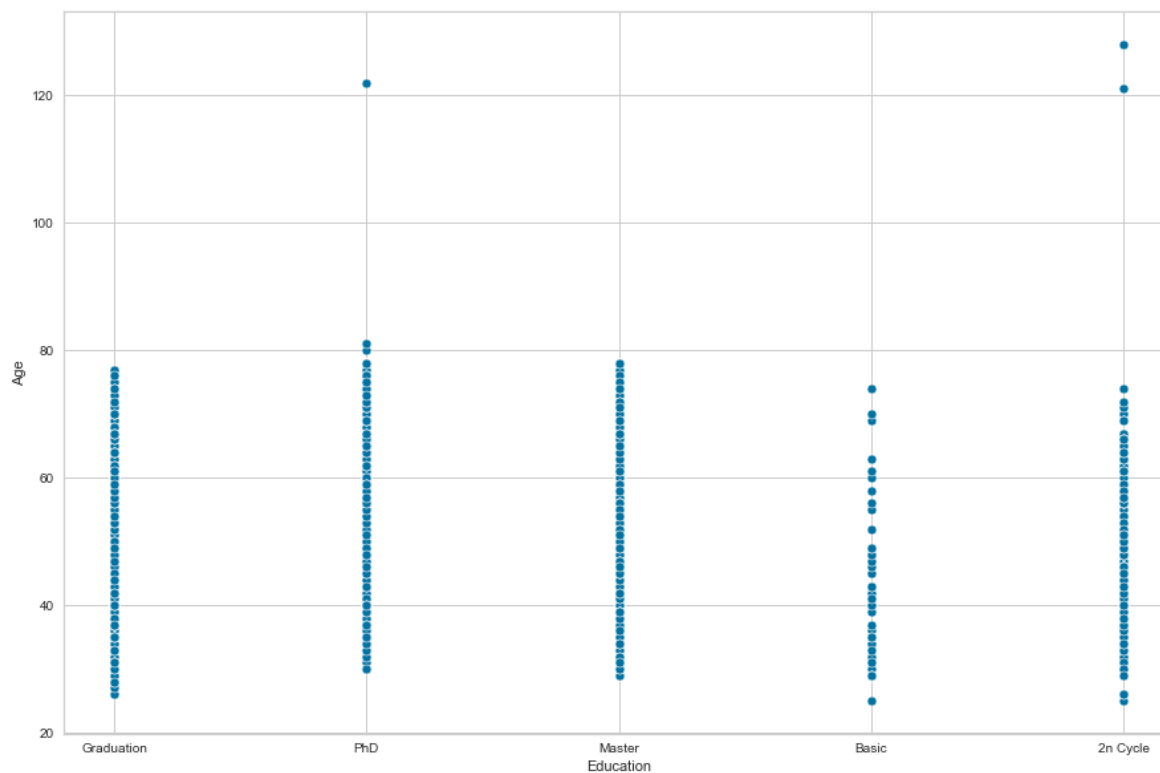
Name: Education, dtype: int64

In [1614]:

```
sns.scatterplot(data=data,x='Education',y='Age')
```

Out[1614]:

<AxesSubplot:xlabel='Education', ylabel='Age'>

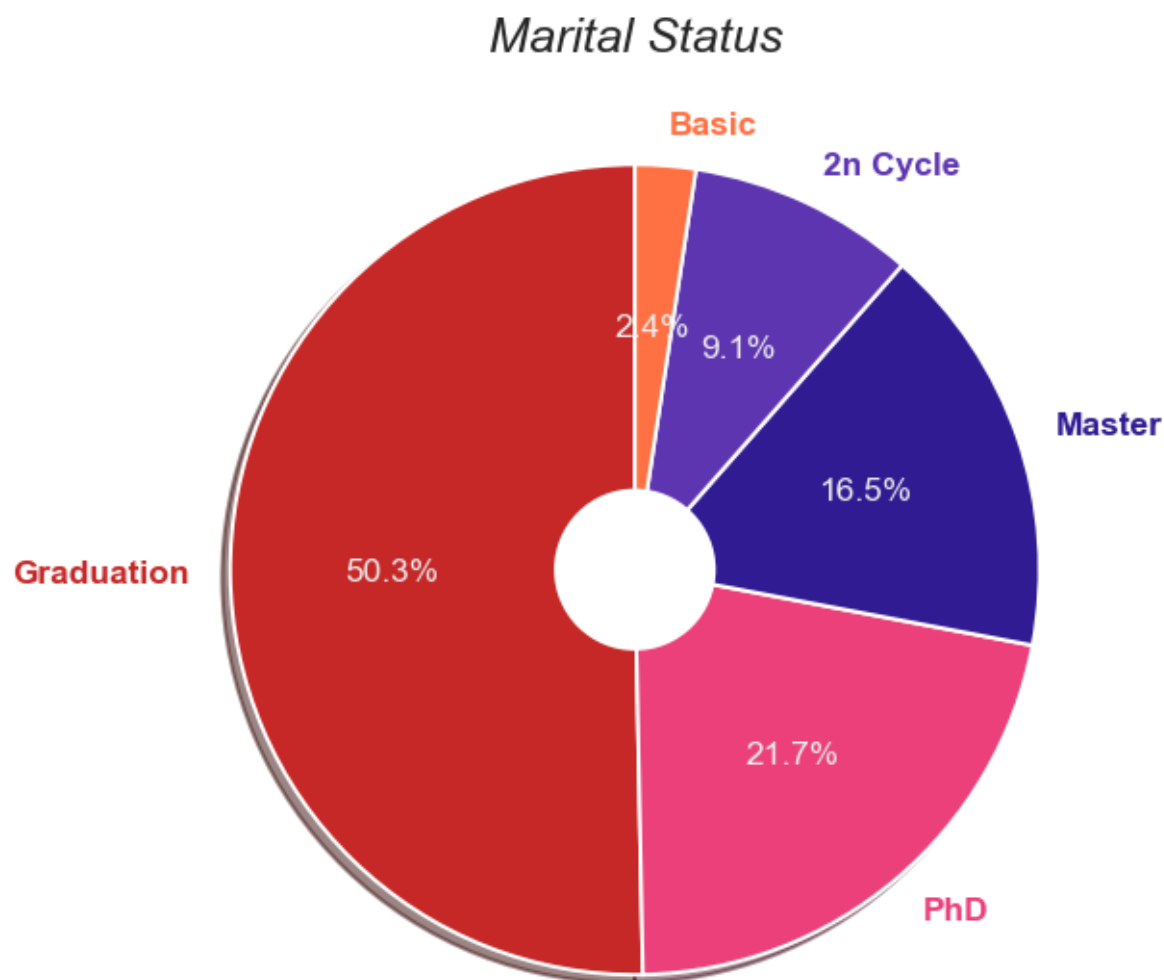


In [1615]:

```
fig, ax = plt.subplots(figsize=(15, 8))
patches, texts, pcts = ax.pie(
    data['Education'].value_counts(), labels=[*data['Education'].value_counts(
    ,wedgeprops={'linewidth': 2.0, 'edgecolor': 'white'},
    textprops={'size': 'x-large'},
    startangle=90)

for i, patch in enumerate(patches):
    texts[i].set_color(patch.get_facecolor())
plt.setp(pcts, color='#EEEEEE')
plt.setp(texts, fontweight=600)
centre_circle = plt.Circle((0,0),0.20,fc='white')
plt.gcf().gca().add_artist(centre_circle)
plt.tight_layout()
plt.title(label='Marital Status', fontsize=25, fontstyle='italic')

plt.tight_layout()
```



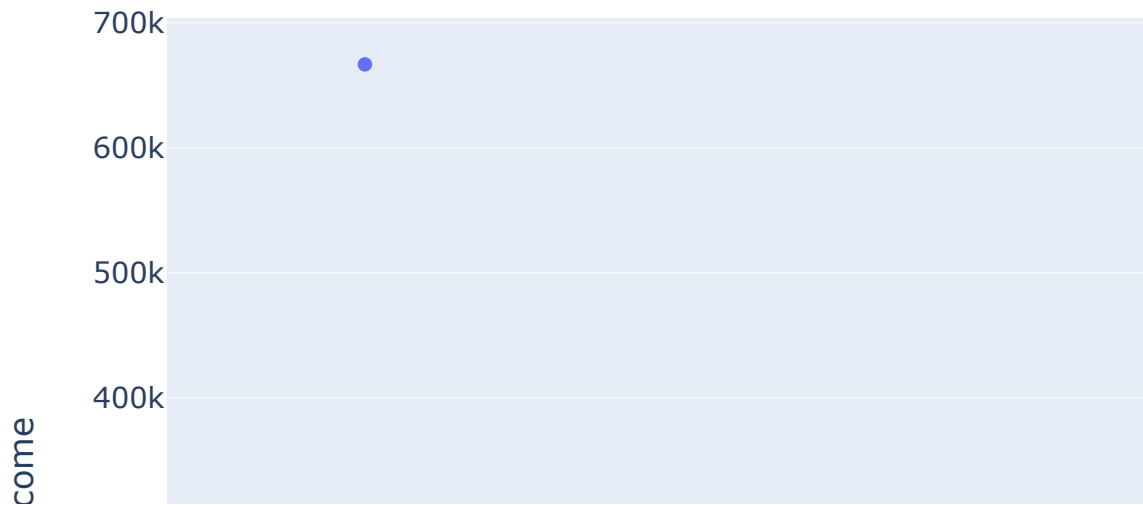
The percentage of clients with university degrees reaches 97.6%

In [1616]:

```
cmap=['RdGy','twilight_shifted','BuGn','PuBuGn']
```

In [1617]:

```
fig = px.box(data, x="Education", y="Income",)  
fig.show()
```



The average income of all clients is very similar except for Basic

In [1618]:

```
data['total_purchases']
```

Out[1618]:

0	1617
1	27
2	776
3	53
4	422

...

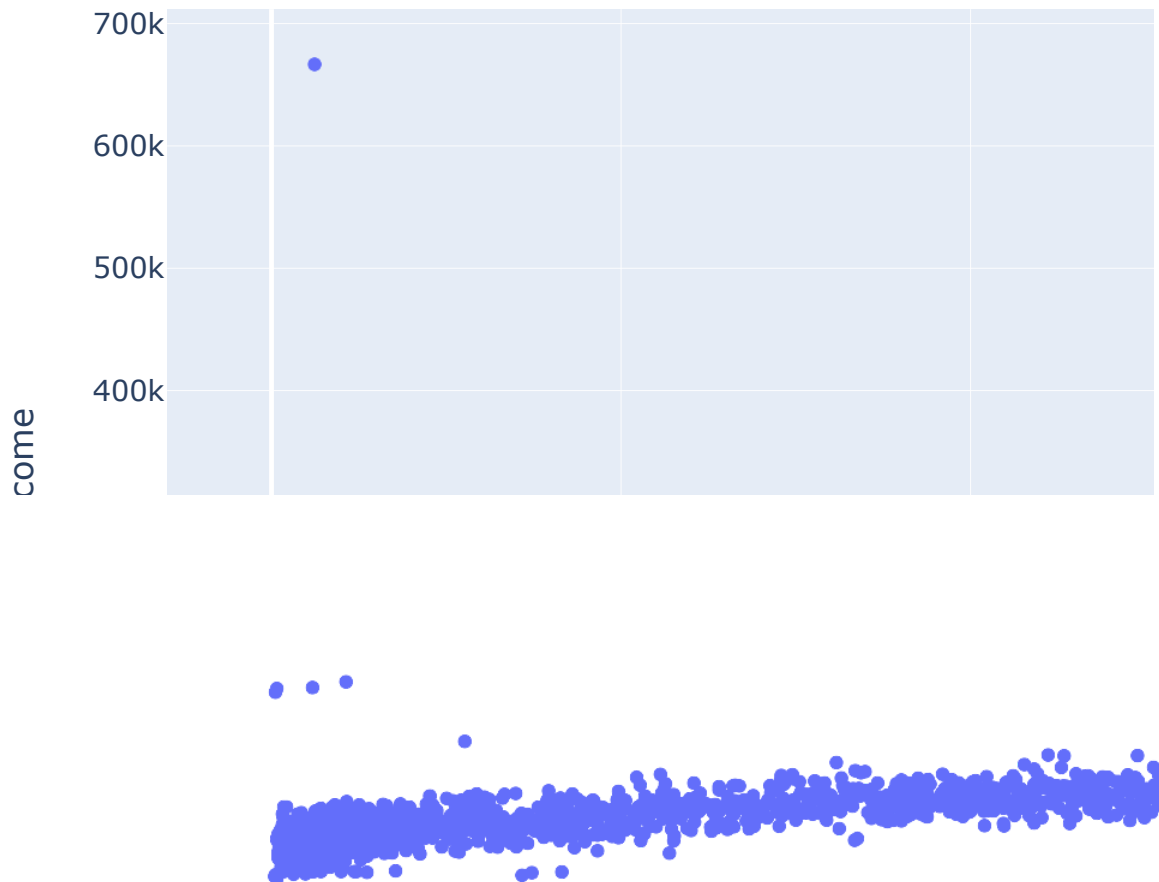
2235	1341
2236	444
2237	1241
2238	843
2239	172

Name: total_purchases, Length: 2240, dtype: int64

Q7: What is the relationship between income and the number of purchases?

In [1619]:

```
fig = px.scatter(data, x="total_purchases", y="Income")  
fig.show()
```



There is a near-linear relationship between income and the number of purchases

العملاء الذين بلغ عدد عمليات الشراء الخاصة بهم ٥٠٠ أغلبهم دخلهم السنوي يقل عن الـ ٥٠ ألف أما باقي العملاء الذين تزيد عدد عمليات الشراء الخاصة بهم الـ ٥٠٠ أغلبهم دخلهم يزيد عن ٥٠ ألف وقد يصل إلى الـ ١٠٠ ألف

أحكيها بالبريزنتيشن

In []:

Q8: What is the relationship between the number of purchases from the website and the number of website visits?

In [1620]:

```
data[ 'NumWebPurchases' ].value_counts()
```

Out[1620]:

2	373
1	354
3	336
4	280
5	220
6	205
7	155
8	102
9	75
0	49
11	44
10	43
27	2
23	1
25	1

Name: NumWebPurchases, dtype: int64

In [1621]:

```
data['NumWebVisitsMonth'].value_counts()
```

Out[1621]:

7	393
8	342
6	340
5	281
4	218
3	205
2	202
1	153
9	83
0	11
20	3
10	3
14	2
19	2
17	1
13	1

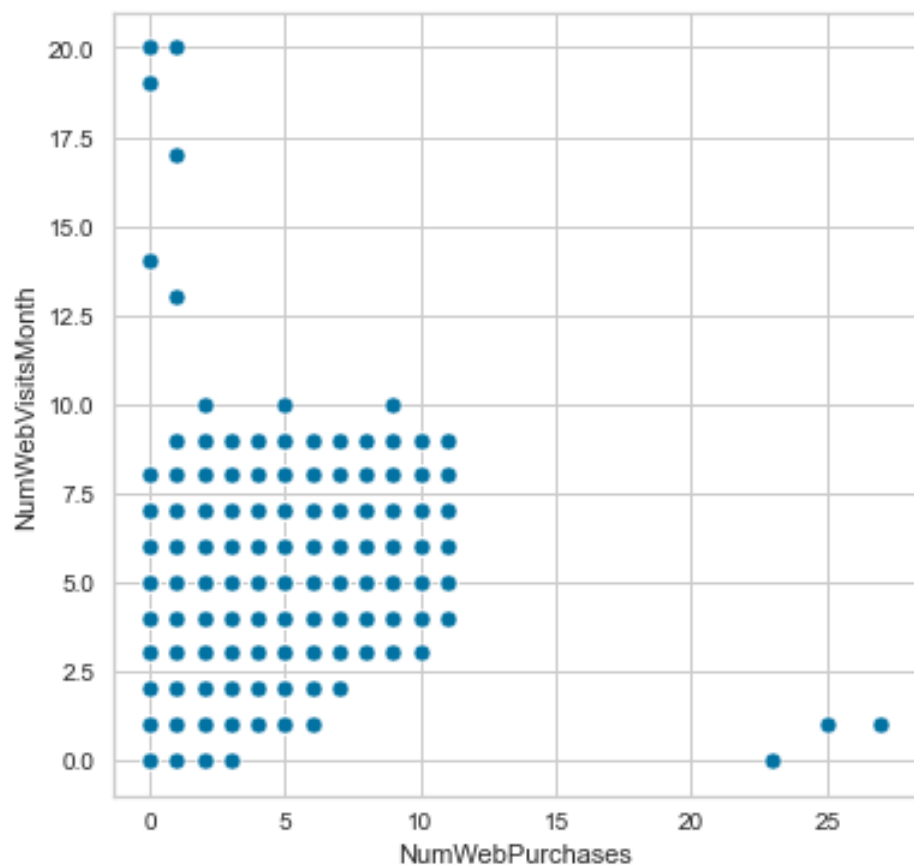
Name: NumWebVisitsMonth, dtype: int64

In [1622]:

```
plt.figure(figsize=(6,6))  
sns.scatterplot(data=data, x="NumWebPurchases", y="NumWebVisitsMonth")
```

Out[1622]:

<AxesSubplot:xlabel='NumWebPurchases', ylabel='NumWebVisitsMonth'>

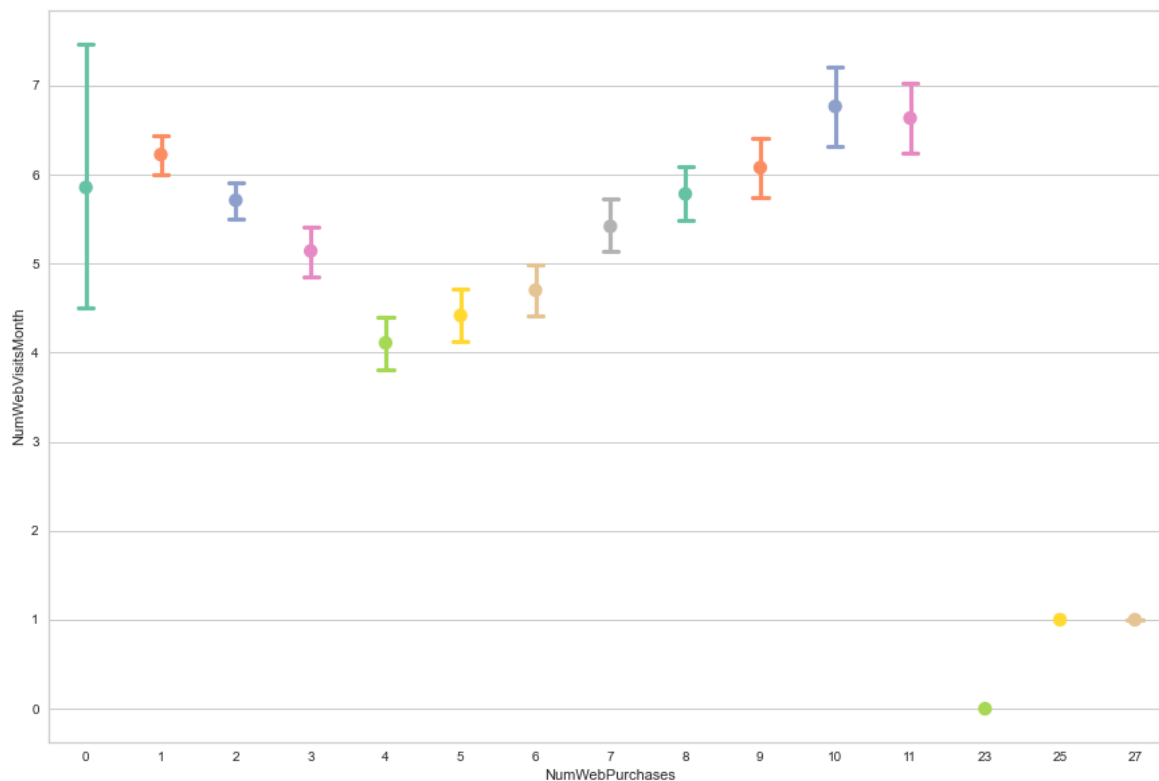


In [1623]:

```
sns.pointplot(data=data,x='NumWebPurchases',y='NumWebVisitsMonth',palette='Set
```

Out[1623]:

```
<AxesSubplot:xlabel='NumWebPurchases', ylabel='NumWebVisitsMonth'>
```



In general, there is a linear relationship between the number of visits to the site and the number of purchases from it

طبعاً باستثناء بعض الحالات الغريبة مثل أنه يوجد عميل اشترى من الموقع 23 مرة دون زيارة الموقع

In []:

Q9: What is the relationship between the number of purchases from a Deal with the number of purchases from the website, the number

of purchases from the catalog, and the number of purchases from the store?

In [1624]:

```
data['NumDealsPurchases'].value_counts()
```

Out[1624]:

1	970
2	497
3	297
4	189
5	94
6	61
0	46
7	40
8	14
9	8
15	7
10	5
11	5
12	4
13	3

Name: NumDealsPurchases, dtype: int64

In [1625]:

```
data['NumWebPurchases'].value_counts()
```

Out[1625]:

2	373
1	354
3	336
4	280
5	220
6	205
7	155
8	102
9	75
0	49
11	44
10	43
27	2
23	1
25	1

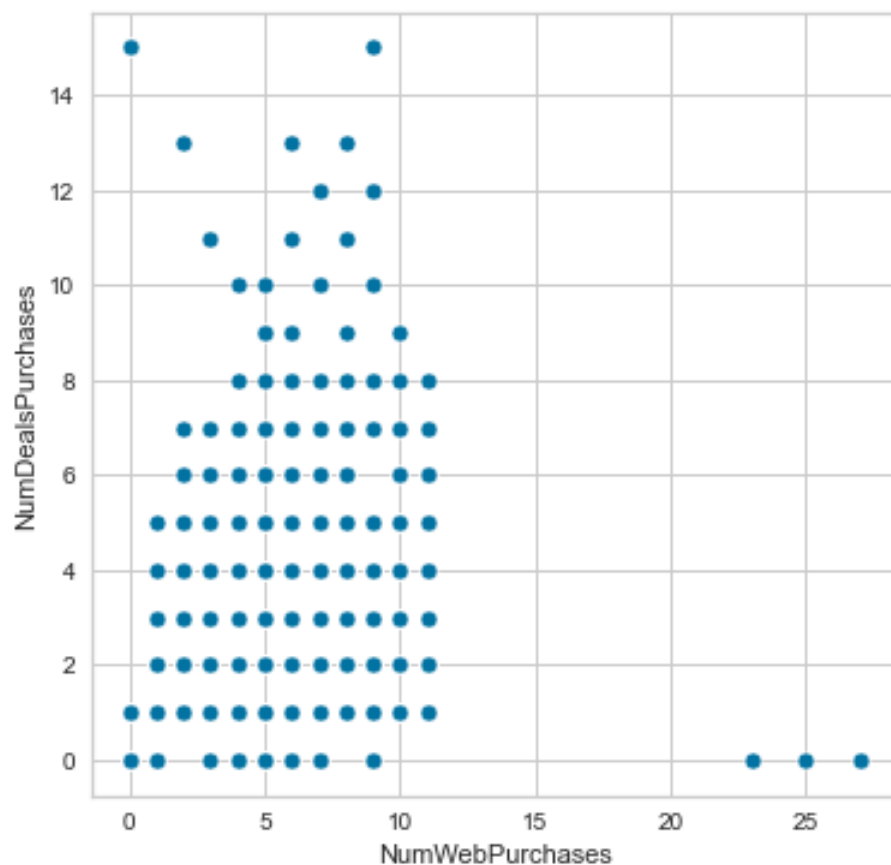
Name: NumWebPurchases, dtype: int64

In [1626]:

```
plt.figure(figsize=(6,6))  
sns.scatterplot(data=data,x='NumWebPurchases',y='NumDealsPurchases')
```

Out[1626]:

<AxesSubplot:xlabel='NumWebPurchases', ylabel='NumDealsPurchases'>



In [1627]:

```
data['NumCatalogPurchases'].value_counts()
```

Out[1627]:

0	586
1	497
2	276
3	184
4	182
5	140
6	128
7	79
8	55
10	48
9	42
11	19
28	3
22	1

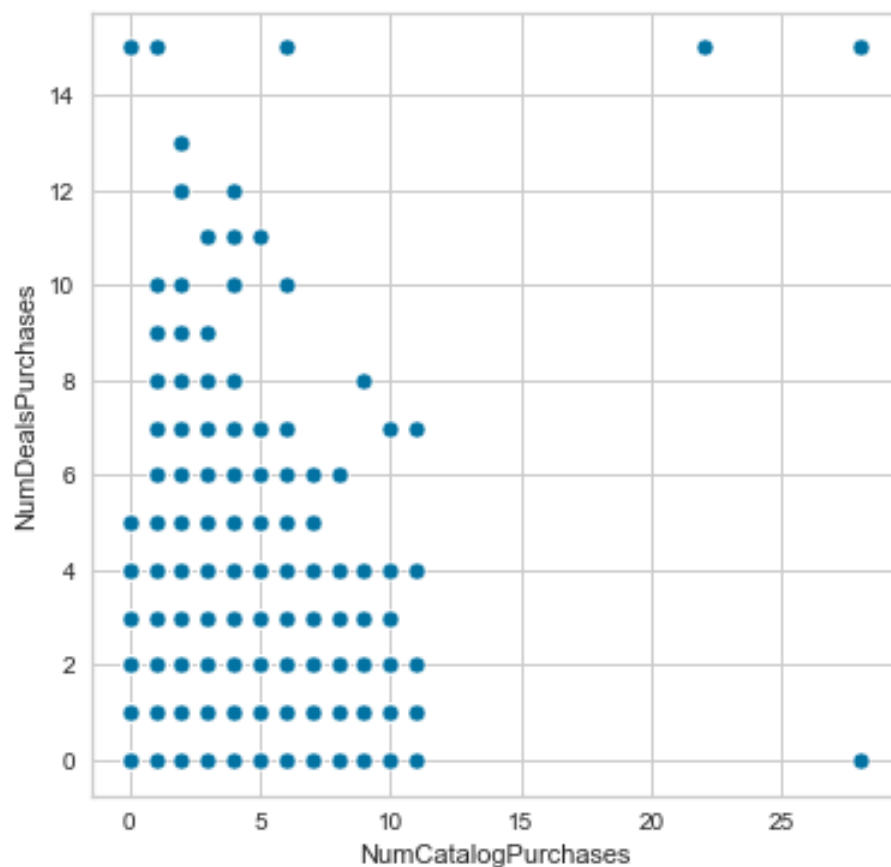
Name: NumCatalogPurchases, dtype: int64

In [1628]:

```
plt.figure(figsize=(6,6))  
sns.scatterplot(data=data,x='NumCatalogPurchases',y='NumDealsPurchases')
```

Out[1628]:

<AxesSubplot:xlabel='NumCatalogPurchases', ylabel='NumDealsPurchases'>



In [1629]:

```
data['NumStorePurchases'].value_counts()
```

Out[1629]:

3	490
4	323
2	223
5	212
6	178
8	149
7	143
10	125
9	106
12	105
13	83
11	81
0	15
1	7

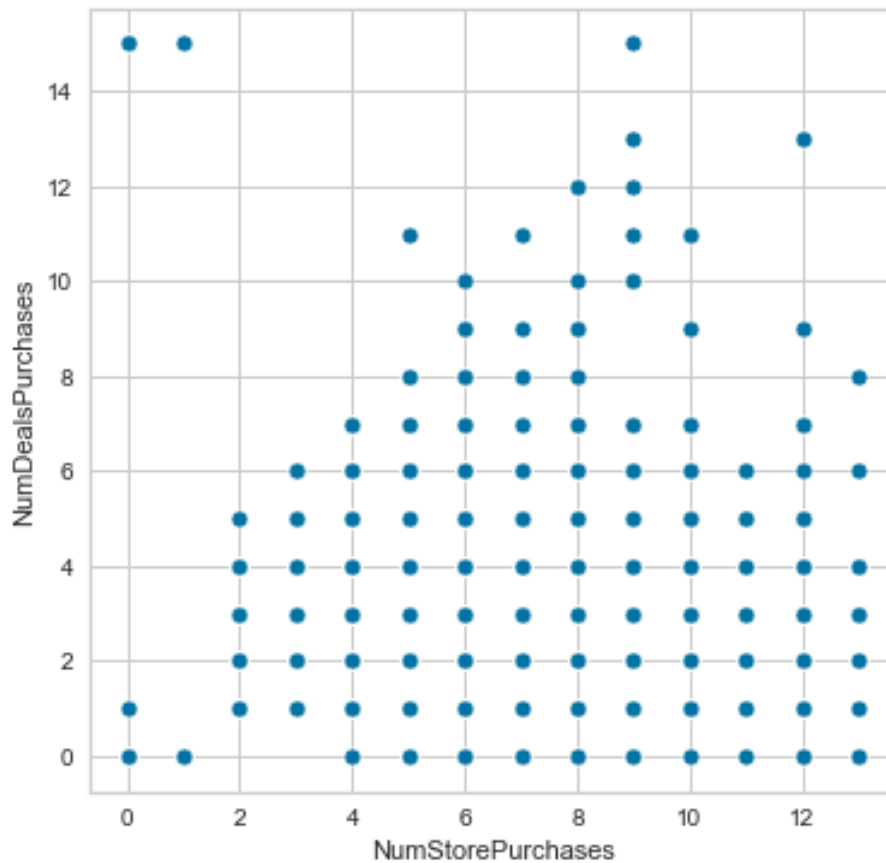
Name: NumStorePurchases, dtype: int64

In [1630]:

```
plt.figure(figsize=(6,6))  
sns.scatterplot(data=data,x='NumStorePurchases',y='NumDealsPurchases')
```

Out[1630]:

<AxesSubplot:xlabel='NumStorePurchases', ylabel='NumDealsPurchases'>



In []:

Q10: What is the relationship between the number of purchases from a Deal with accepted cmp 1 ,accepted cmp 2,accepted cmp 3 ,accepted cmp 4 ,accepted cmp 5 and Response?

In [1631]:

```
data['NumDealsPurchases'].value_counts()
```

Out[1631]:

1	970
2	497
3	297
4	189
5	94
6	61
0	46
7	40
8	14
9	8
15	7
10	5
11	5
12	4
13	3

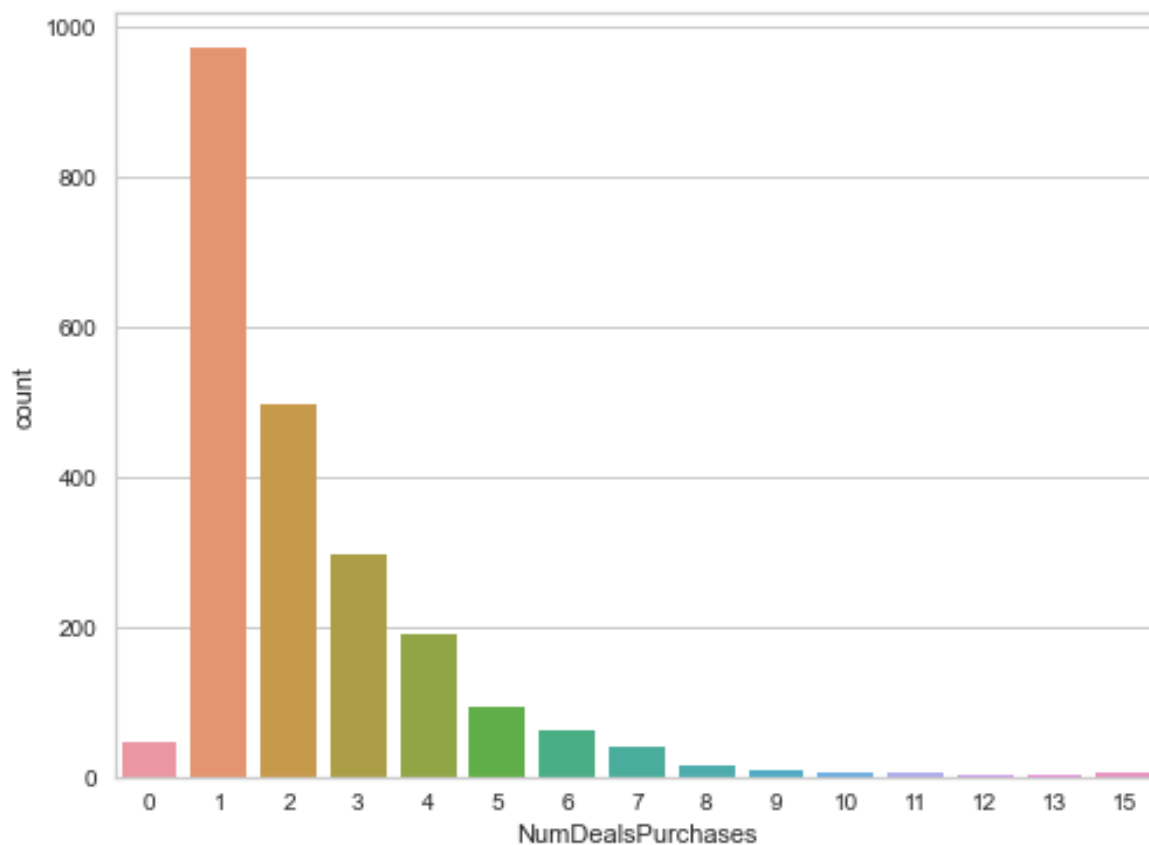
Name: NumDealsPurchases, dtype: int64

In [1632]:

```
plt.figure(figsize=(8,6))  
sns.countplot(data=data,x='NumDealsPurchases')
```

Out[1632]:

<AxesSubplot:xlabel='NumDealsPurchases', ylabel='count'>



Almost all customers have made a purchase at least once

In [1633]:

reate a feature that collects all offers

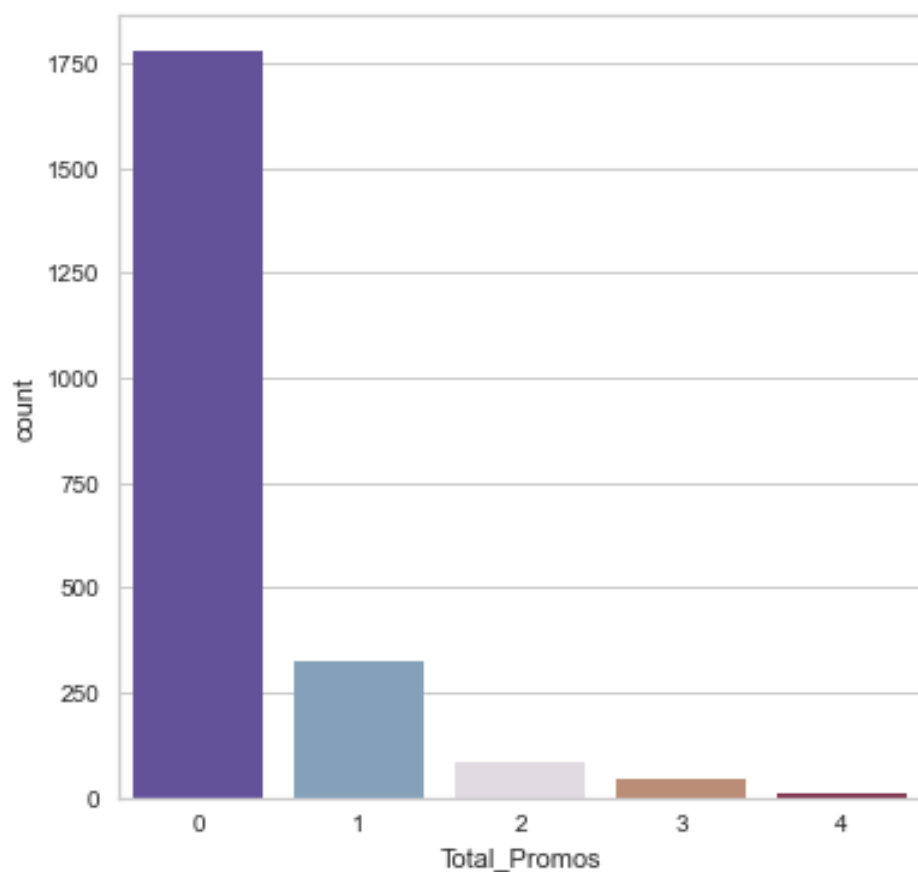
```
data["Total_Promos"] = data["AcceptedCmp1"] + data["AcceptedCmp2"] + data["AcceptedCmp3"]
```

In [1634]:

```
plt.figure(figsize=(6,6))  
sns.countplot(data=data, x='Total_Promos', palette='twilight_shifted')
```

Out[1634]:

<AxesSubplot:xlabel='Total_Promos', ylabel='count'>



The number of customers who did not accept the offers is very large, up to 80%

In []:

Q11: What is the relationship between the complaint and Date of customer's enrollment?

In [1635]:

```
data['Complain'].value_counts()
```

Out[1635]:

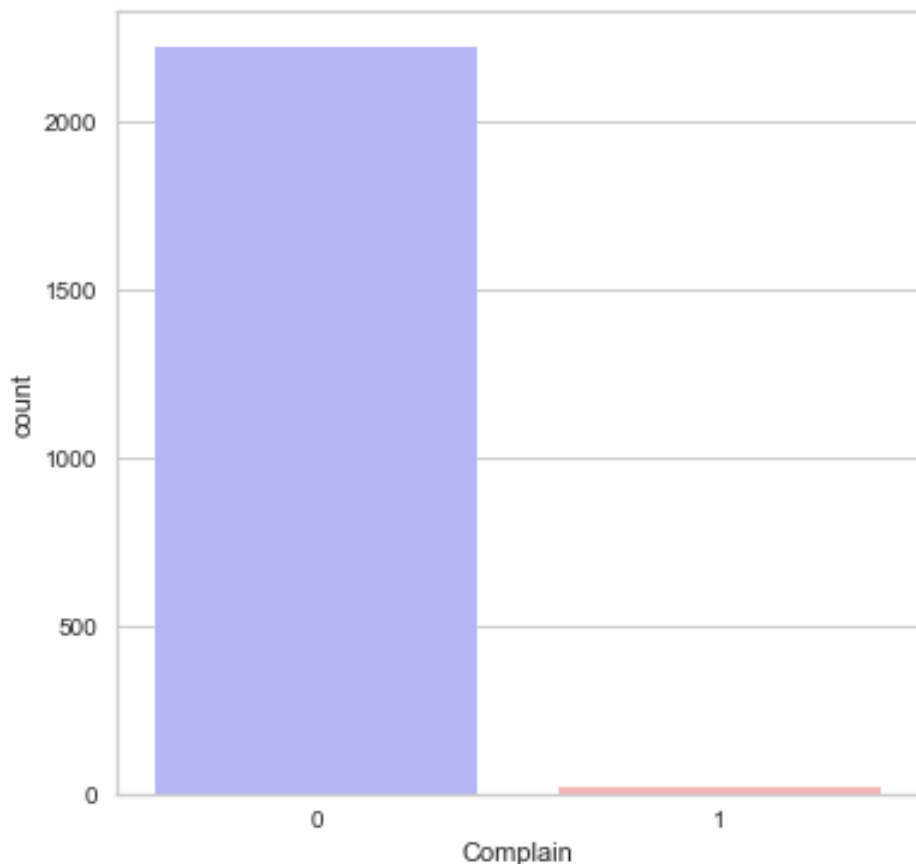
```
0    2219
1      21
Name: Complain, dtype: int64
```

In [1636]:

```
plt.figure(figsize=(6,6))
sns.countplot(data=data,x='Complain',palette='bwr')
```

Out[1636]:

```
<AxesSubplot:xlabel='Complain', ylabel='count'>
```



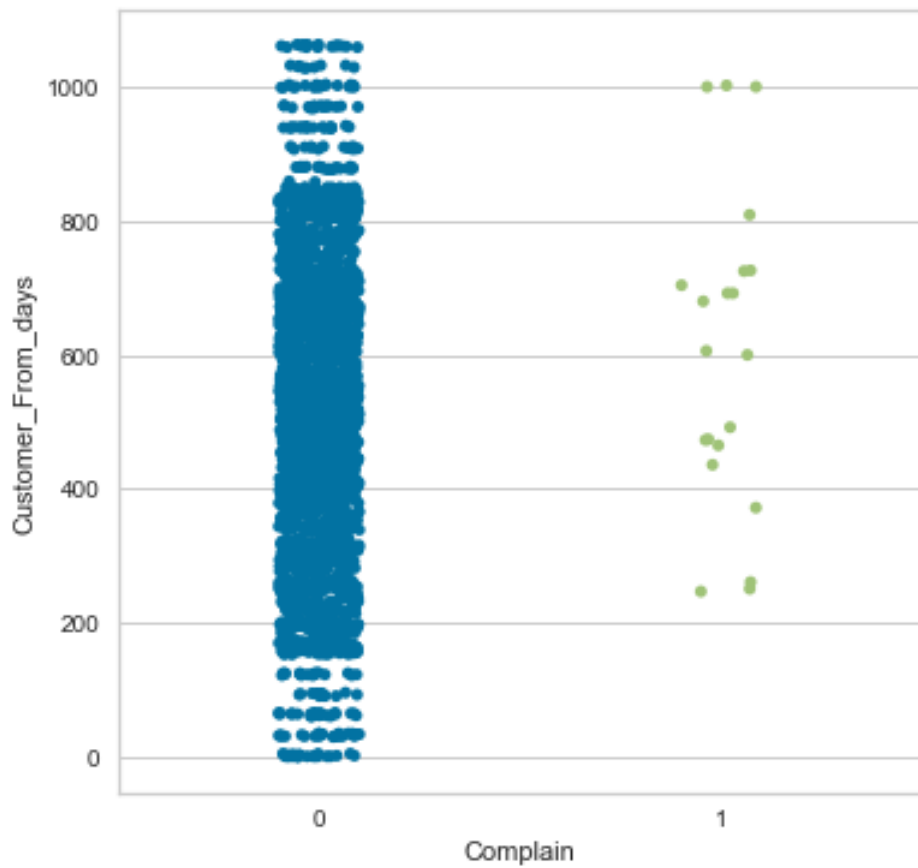
Very few customers who complained

In [1637]:

```
plt.figure(figsize=(6,6))  
sns.stripplot(data=data,x='Complain',y='Customer_From_days')
```

Out[1637]:

<AxesSubplot:xlabel='Complain', ylabel='Customer_From_days'>



All customers who filed a complaint as if they were with the company more than 200 days ago

In []:

Data Preprocessing

In [1638]:

```
data.isnull().sum()
```

Out[1638]:

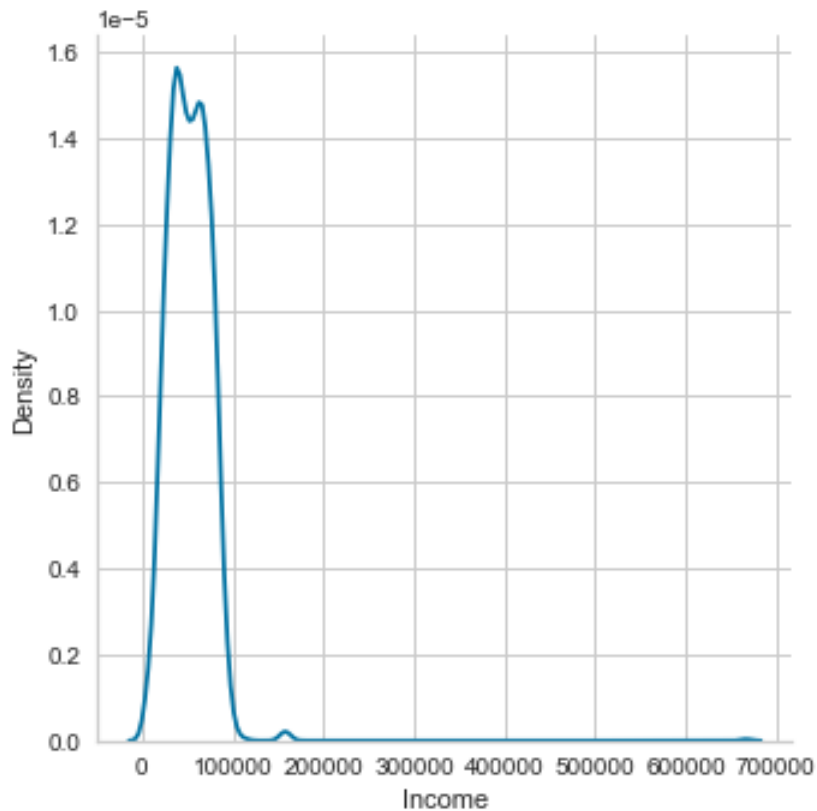
ID	0
Year_Birth	0
Education	0
Marital_Status	0
Income	24
Kidhome	0
Teenhome	0
Dt_Customer	0
Recency	0
MntWines	0
MntFruits	0
MntMeatProducts	0
MntFishProducts	0
MntSweetProducts	0
MntGoldProds	0
NumDealsPurchases	0
NumWebPurchases	0
NumCatalogPurchases	0
NumStorePurchases	0
NumWebVisitsMonth	0
AcceptedCmp3	0
AcceptedCmp4	0
AcceptedCmp5	0
AcceptedCmp1	0
AcceptedCmp2	0
Complain	0
Z_CostContact	0
Z_Revenue	0
Response	0
Customer_From_days	0
Age	0
Living_With	0
Num_Children	0
Family_Size	0
total_purchases	0
purchase_quantity	1
Total_Promos	0
dtype:	int64

In [1639]:

```
sns.displot(data,x='Income',kind="kde")
```

Out[1639]:

<seaborn.axisgrid.FacetGrid at 0x1c5e274bd30>



In [1640]:

```
fill_tobed=data['Income'].dropna()  
data['Income']=data['Income'].fillna(pd.Series(np.random.choice(fill_tobed,size=
```

In [1641]:

```
fill_tobed=data['purchase_quantity'].dropna()  
data['purchase_quantity']=data['purchase_quantity'].fillna(pd.Series(np.random
```


In [1642]:

```
data.isnull().sum()
```

Out[1642]:

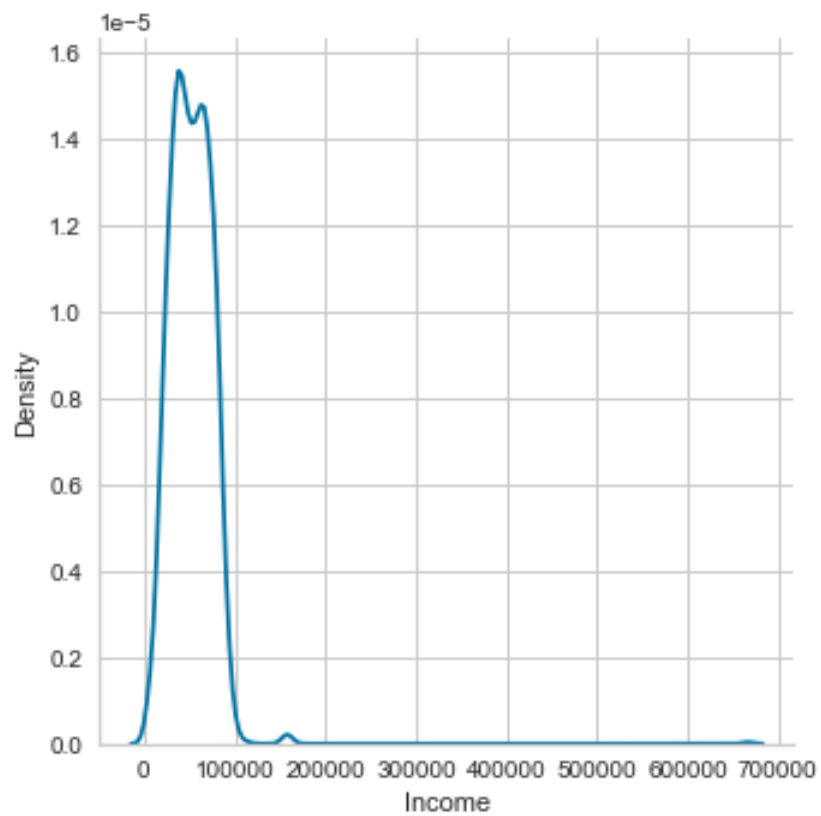
ID	0
Year_Birth	0
Education	0
Marital_Status	0
Income	0
Kidhome	0
Teenhome	0
Dt_Customer	0
Recency	0
MntWines	0
MntFruits	0
MntMeatProducts	0
MntFishProducts	0
MntSweetProducts	0
MntGoldProds	0
NumDealsPurchases	0
NumWebPurchases	0
NumCatalogPurchases	0
NumStorePurchases	0
NumWebVisitsMonth	0
AcceptedCmp3	0
AcceptedCmp4	0
AcceptedCmp5	0
AcceptedCmp1	0
AcceptedCmp2	0
Complain	0
Z_CostContact	0
Z_Revenue	0
Response	0
Customer_From_days	0
Age	0
Living_With	0
Num_Children	0
Family_Size	0
total_purchases	0
purchase_quantity	0
Total_Promos	0
dtype:	int64

In [1643]:

```
sns.displot(data,x='Income',kind="kde")
```

Out[1643]:

<seaborn.axisgrid.FacetGrid at 0x1c5cb272820>

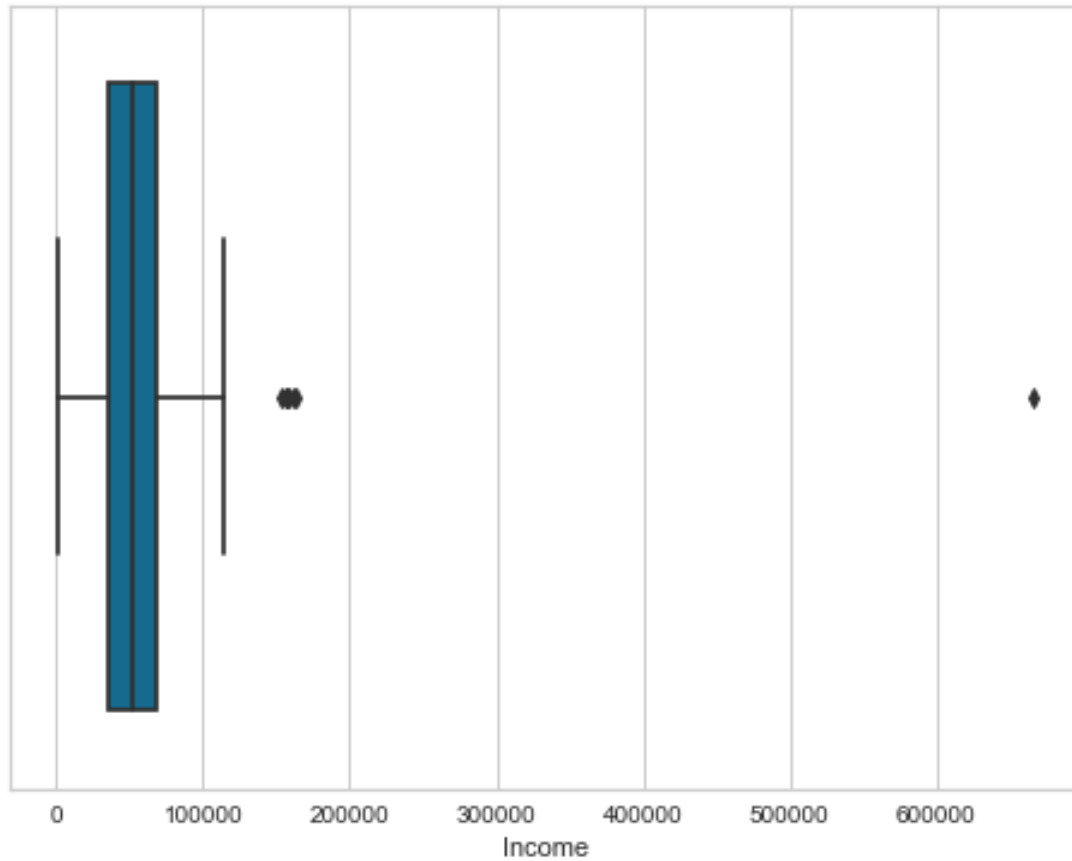


In [1644]:

```
plt.figure(figsize=(8,6))  
sns.boxplot(data=data,x='Income')
```

Out[1644]:

<AxesSubplot:xlabel='Income'>



In [1645]:

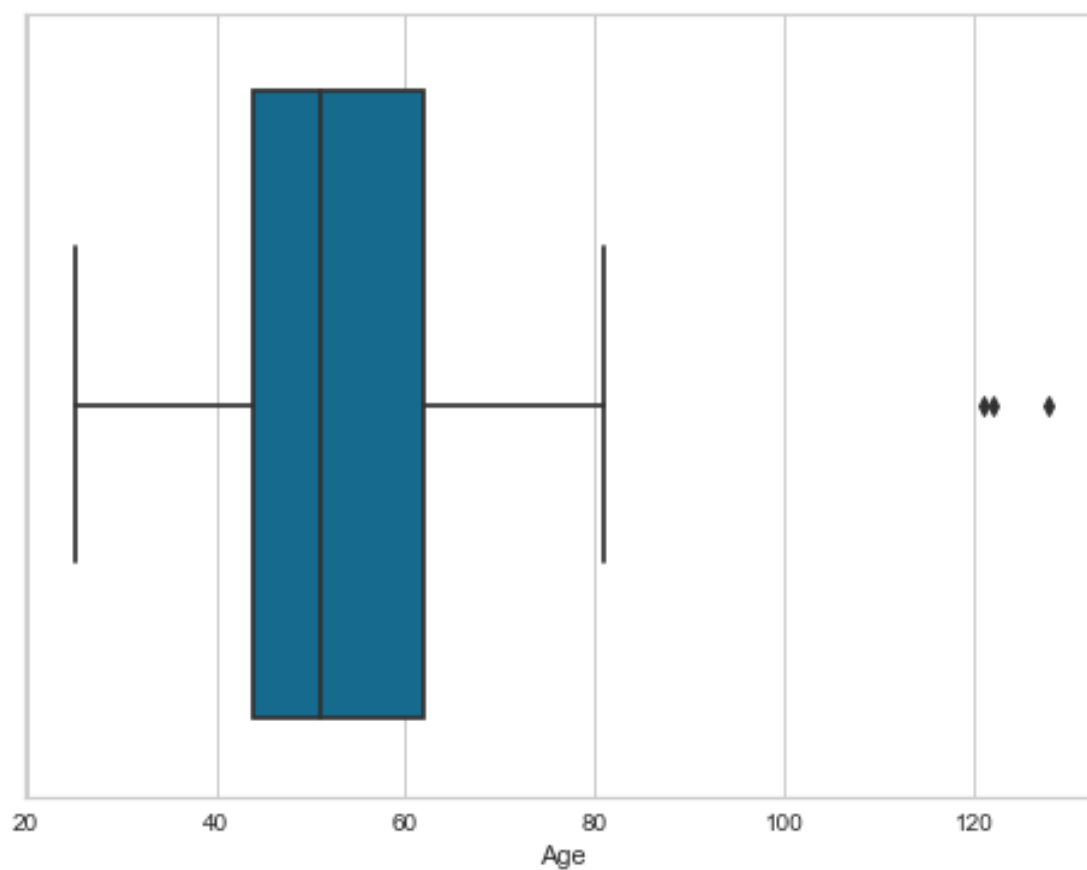
```
# Remove outliers  
data=data[data['Income']<150000]
```

In [1646]:

```
plt.figure(figsize=(8,6))  
sns.boxplot(data=data,x='Age')
```

Out[1646]:

<AxesSubplot: xlabel= 'Age' >



In [1647]:

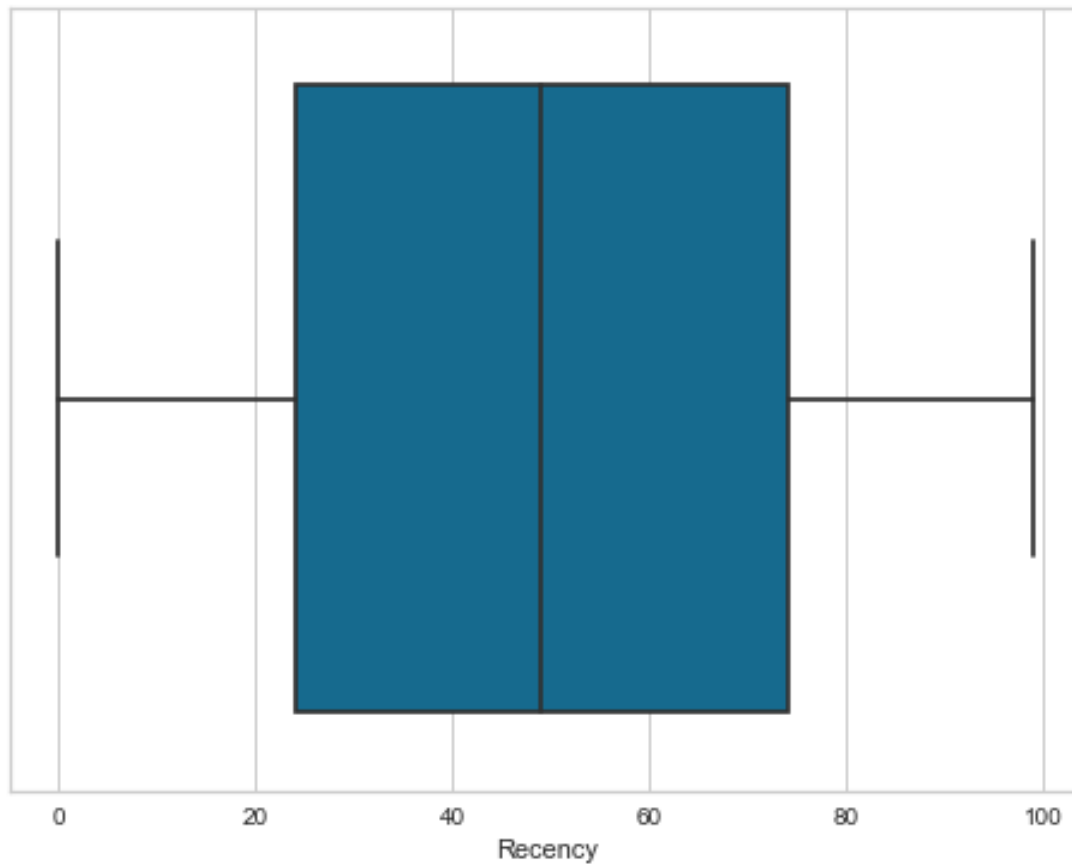
```
# Remove outliers  
data=data[data['Age']<100]
```

In [1648]:

```
plt.figure(figsize=(8,6))  
sns.boxplot(data=data,x='Recency')
```

Out[1648]:

<AxesSubplot:xlabel='Recency'>



In [1649]:

```
data.columns
```

Out[1649]:

```
Index(['ID', 'Year_Birth', 'Education', 'Marital_Status', 'Income', 'Kidhome',
      'Teenhome', 'Dt_Customer', 'Recency', 'MntWines', 'MntFruits',
      'MntMeatProducts', 'MntFishProducts', 'MntSweetProducts',
      'MntGoldProds', 'NumDealsPurchases', 'NumWebPurchases',
      'NumCatalogPurchases', 'NumStorePurchases', 'NumWebVisitsMonth',
      'AcceptedCmp3', 'AcceptedCmp4', 'AcceptedCmp5', 'AcceptedCmp1',
      'AcceptedCmp2', 'Complain', 'Z_CostContact', 'Z_Revenue', 'Response',
      'Customer_From_days', 'Age', 'Living_With', 'Num_Children',
      'Family_Size', 'total_purchases', 'purchase_quantity',
      'Total_Promos'],
      dtype='object')
```

In [1650]:

```
# Drop columns we don't need them
```

```
data=data[['Education', 'Marital_Status', 'Income', 'Recency', 'MntWines', 'MntMeatProducts', 'MntFishProducts', 'MntSweetProducts', 'Kidhome', 'Teenhome', 'Num_Children', 'MntGoldProds', 'NumDealsPurchases', 'NumWebPurchases', 'NumCatalogPurchases', 'NumStorePurchases', 'NumWebVisitsMonth', 'Customer_From_days', 'Age', 'Family_Size', 'total_purchases', 'purchase_quantity', 'Total_Promos']]
```

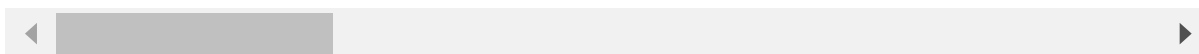
In [1651]:

```
data
```

Out[1651]:

	Education	Marital_Status	Income	Recency	MntWines	MntFruits	MntM
0	Graduation	Single	58138.0	58	635	88	
1	Graduation	Single	46344.0	38	11	1	
2	Graduation	Together	71613.0	26	426	49	
3	Graduation	Together	26646.0	26	11	4	
4	PhD	Married	58293.0	94	173	43	
...	
2235	Graduation	Married	61223.0	46	709	43	
2236	PhD	Together	64014.0	56	406	0	
2237	Graduation	Divorced	56981.0	91	908	48	
2238	Master	Together	69245.0	8	428	30	
2239	PhD	Married	52869.0	40	84	3	

2229 rows × 24 columns



In [1652]:

data.info()

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 2229 entries, 0 to 2239
Data columns (total 24 columns):
 #   Column                                Non-Null Count  Dtype  
---  -
 0   Education                            2229 non-null   object 
 1   Marital_Status                       2229 non-null   object 
 2   Income                               2229 non-null   float64
 3   Recency                             2229 non-null   int64  
 4   MntWines                             2229 non-null   int64  
 5   MntFruits                           2229 non-null   int64  
 6   MntMeatProducts                     2229 non-null   int64  
 7   MntFishProducts                     2229 non-null   int64  
 8   MntSweetProducts                    2229 non-null   int64  
 9   Kidhome                             2229 non-null   int64  
10  Teenhome                             2229 non-null   int64  
11  Num_Children                         2229 non-null   int64  
12  MntGoldProds                        2229 non-null   int64  
13  NumDealsPurchases                   2229 non-null   int64  
14  NumWebPurchases                     2229 non-null   int64  
15  NumCatalogPurchases                 2229 non-null   int64  
16  NumStorePurchases                   2229 non-null   int64  
17  NumWebVisitsMonth                   2229 non-null   int64  
18  Customer_From_days                  2229 non-null   int64  
19  Age                                 2229 non-null   int64  
20  Family_Size                         2229 non-null   int64  
21  total_purchases                     2229 non-null   int64  
22  purchase_quantity                   2229 non-null   category
23  Total_Promos                        2229 non-null   int64  
dtypes: category(1), float64(1), int64(20), object(2)
memory usage: 420.5+ KB

```


In [1653]:

```
s = (data.dtypes == 'object')
n = (data.dtypes == 'category')
object_cols = list(s[s].index)
category_col = list(n[n].index)
print("Categorical variables in the dataset:", object_cols, category_col)
```

Categorical variables in the dataset: ['Education', 'Marital_Status'] ['purchase_quantity']

In [1654]:

```
LE=LabelEncoder()
for i in object_cols:
    data[i]=data[[i]].apply(LE.fit_transform)
for i in category_col:
    data[i]=data[[i]].apply(LE.fit_transform)
print("All features are now numerical")
```

All features are now numerical

In [1655]:

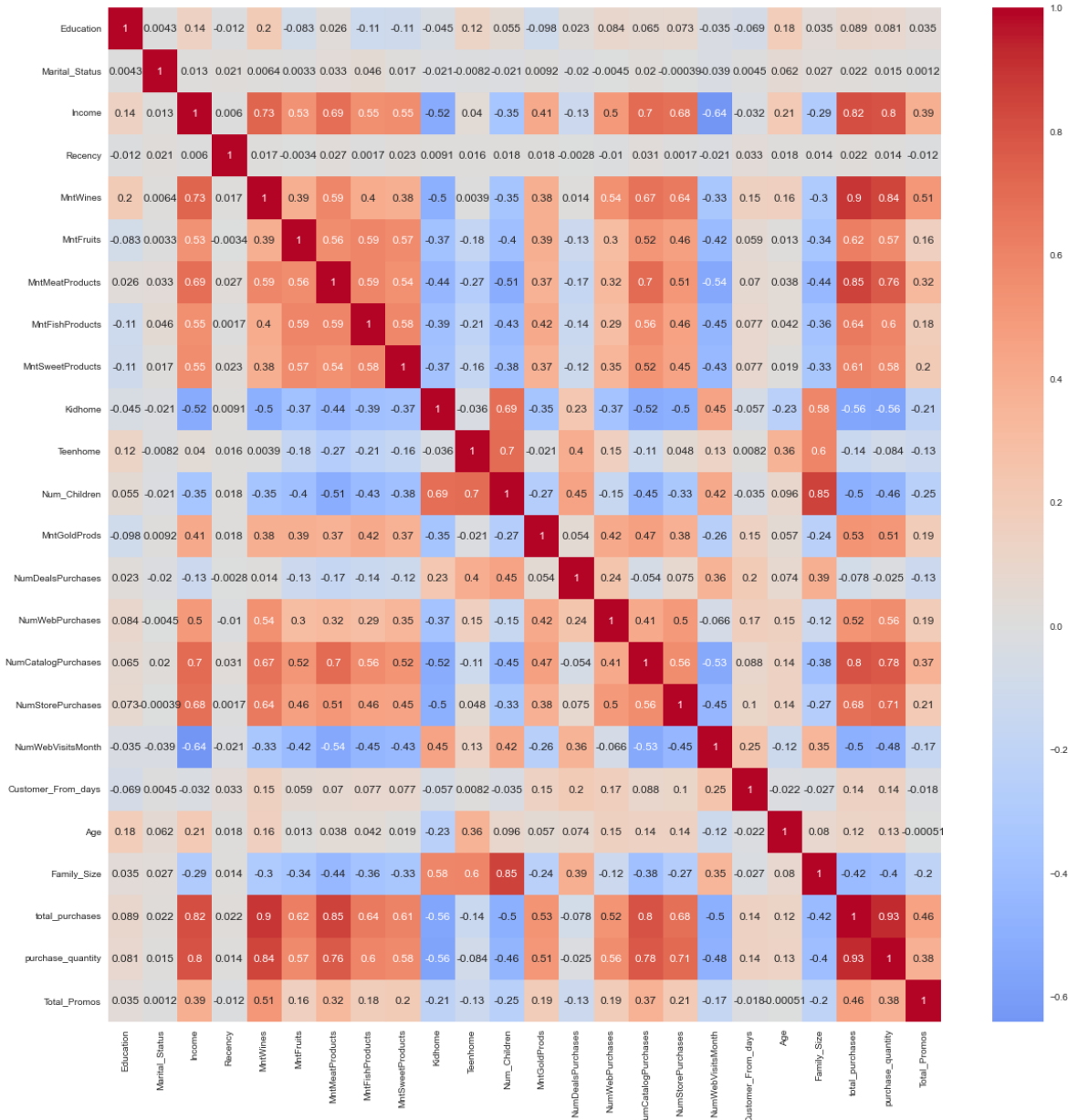
```

corrmat= data.corr()
plt.figure(figsize=(20,20))
sns.heatmap(corrmat,annot=True, cmap='coolwarm', center=0)

```

Out[1655]:

<AxesSubplot:>



In []:

In [1656]:

```
#Creating a copy of data
ds = data.copy()
#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

In [1657]:

```
ds.isnull().sum()
```

Out[1657]:

Education	0
Marital_Status	0
Income	0
Recency	0
MntWines	0
MntFruits	0
MntMeatProducts	0
MntFishProducts	0
MntSweetProducts	0
Kidhome	0
Teenhome	0
Num_Children	0
MntGoldProds	0
NumDealsPurchases	0
NumWebPurchases	0
NumCatalogPurchases	0
NumStorePurchases	0
NumWebVisitsMonth	0
Customer_From_days	0
Age	0
Family_Size	0
total_purchases	0
purchase_quantity	0
Total_Promos	0
dtype:	int64

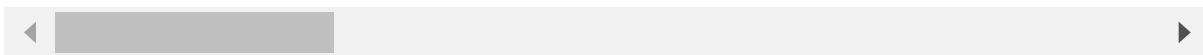
In [1658]:

```
ds
```

Out[1658]:

	Education	Marital_Status	Income	Recency	MntWines	MntFruits	MntM
0	2	2	58138.0	58	635	88	
1	2	2	46344.0	38	11	1	
2	2	3	71613.0	26	426	49	
3	2	3	26646.0	26	11	4	
4	4	1	58293.0	94	173	43	
...	
2235	2	1	61223.0	46	709	43	
2236	4	3	64014.0	56	406	0	
2237	2	0	56981.0	91	908	48	
2238	3	3	69245.0	8	428	30	
2239	4	1	52869.0	40	84	3	

2229 rows × 24 columns



In [1659]:

```
# Using PCA to reduce the dimensions of the data to 3 dimensions
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
```

Out[1659]:

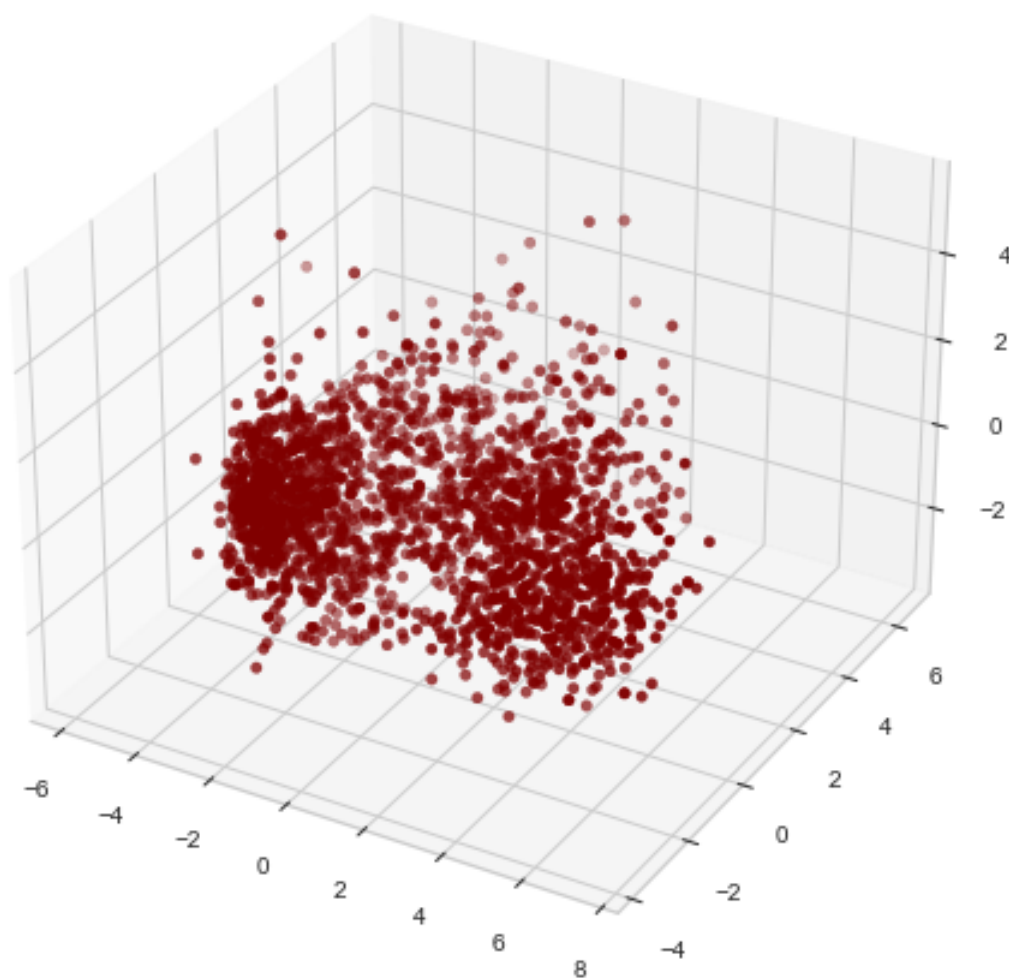
	count	mean	std	min	25%	50%	75%
col1	2229.0	-1.715392e-16	2.985744	-5.843514	-2.707964	-0.886203	2.595992
col2	2229.0	1.002139e-16	1.648486	-3.727375	-1.352172	-0.163528	1.160058
col3	2229.0	2.171634e-17	1.246159	-3.423161	-0.861863	0.011535	0.798212



In [1660]:

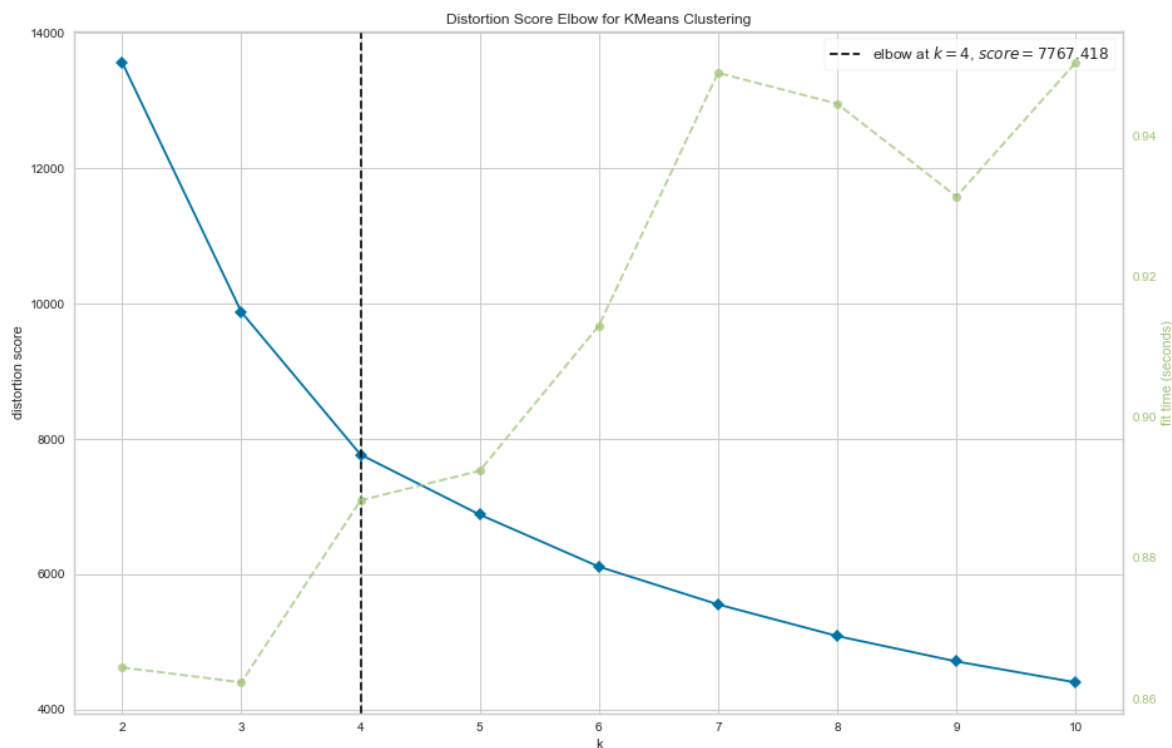
```
x =PCA_ds["col1"]
y =PCA_ds["col2"]
z =PCA_ds["col3"]
fig = plt.figure(figsize=(10,8))
ax = fig.add_subplot(111, projection="3d")
ax.scatter(x,y,z, c="maroon", marker="o" )
ax.set_title("A 3D Projection of Data After Dimensional Reduction")
plt.show()
```

A 3D Projection of Data After Dimensional Reduction



In [1661]:

```
# Using the elbow method to find the optimal number of clusters'  
Elbow_M = KElbowVisualizer(KMeans(), k=10)  
Elbow_M.fit(PCA_ds)  
Elbow_M.show()
```



Out[1661]:

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


In []:

In [1662]:

```
#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
```

In [1663]:

```
#Plotting the clusters
```

```
fig = plt.figure(figsize=(10,8))
```

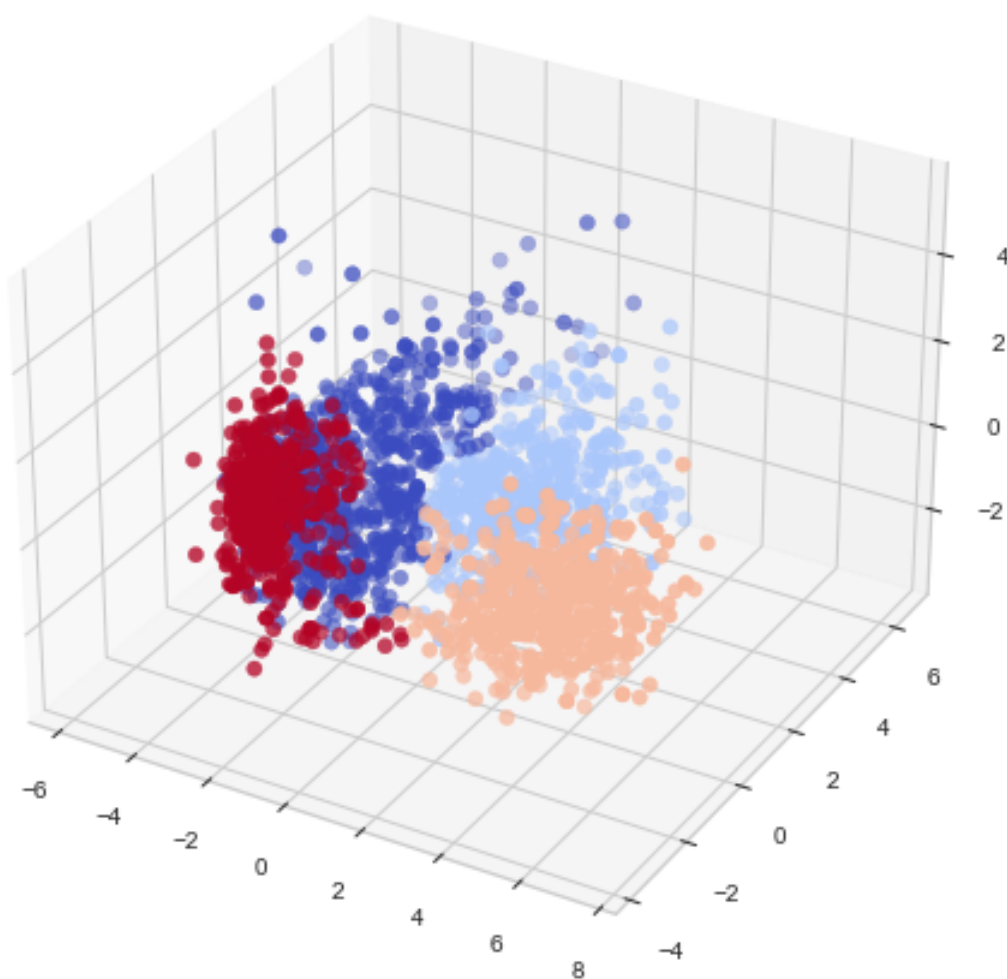
```
ax = plt.subplot(111, projection='3d', label="bla")
```

```
ax.scatter(x, y, z, s=40, c=PCA_ds["Clusters"], marker='o', cmap = 'coolwarm')
```

```
ax.set_title("The Plot Of The Clusters")
```

```
plt.show()
```

The Plot Of The Clusters

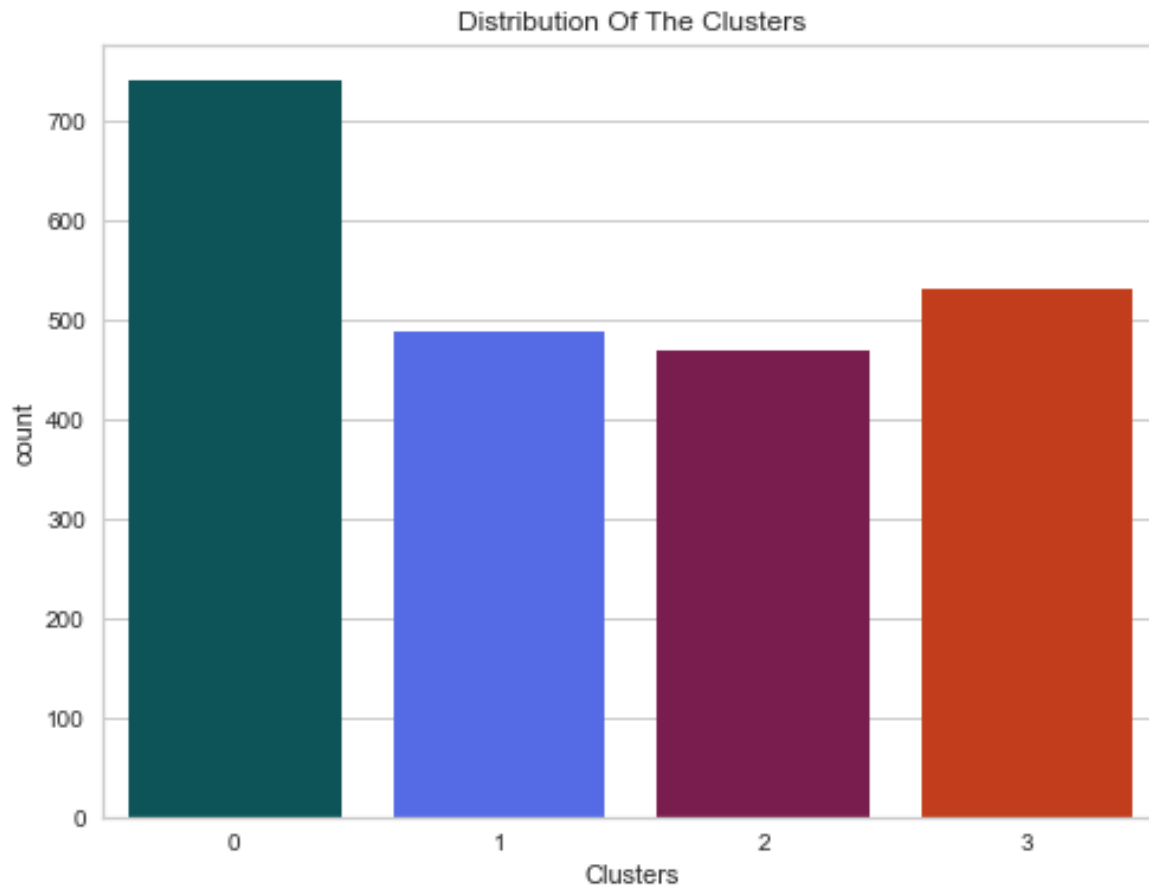


In [1664]:

```
color_2=[ '#006064', '#3d5afe', '#880e4f', '#dd2c00']
```

In [1665]:

```
#Plotting countplot of clusters  
plt.figure(figsize=(8,6))  
pl = sns.countplot(x=data["Clusters"], palette=color_2 )  
pl.set_title("Distribution Of The Clusters")  
plt.show()
```



In [1666]:

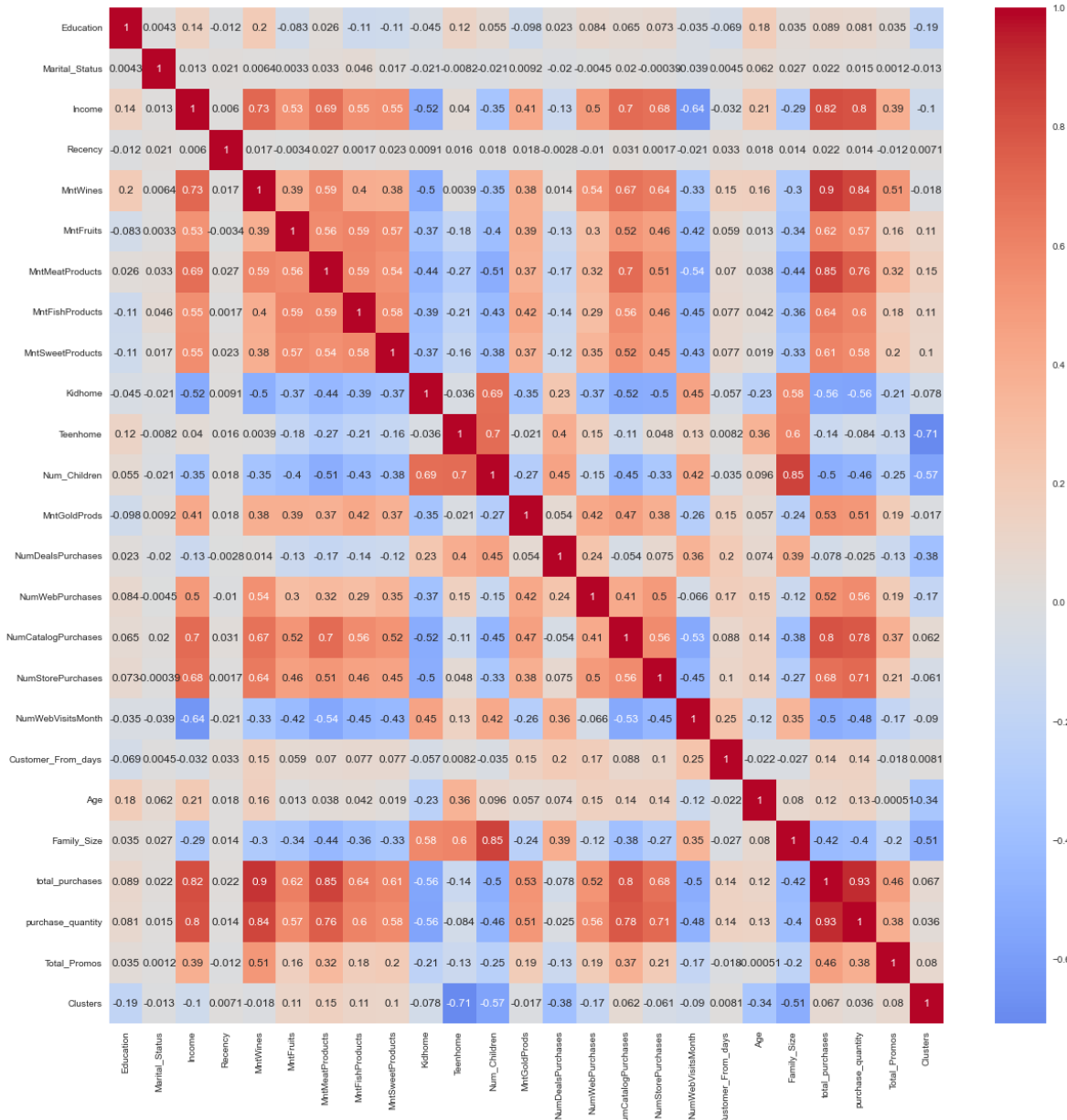
```

corrmat= data.corr()
plt.figure(figsize=(20,20))
sns.heatmap(corrmat,annot=True, cmap='coolwarm', center=0)

```

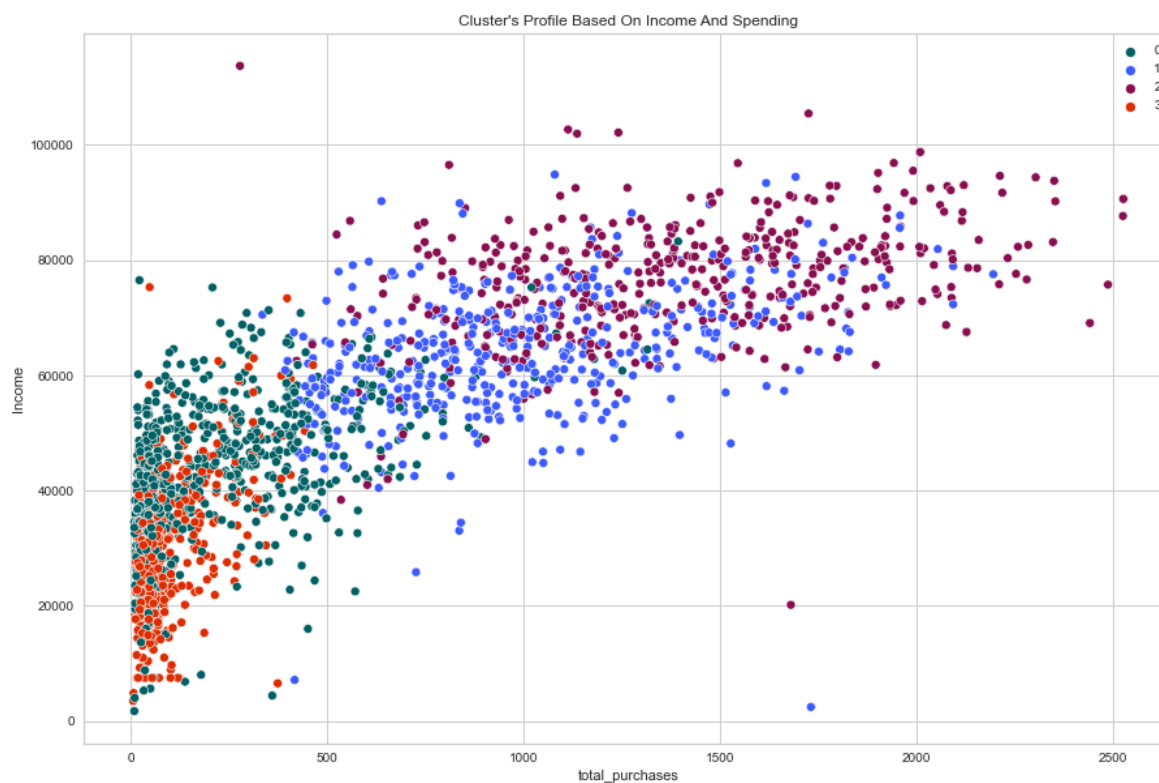
Out[1666]:

<AxesSubplot:>



In [1667]:

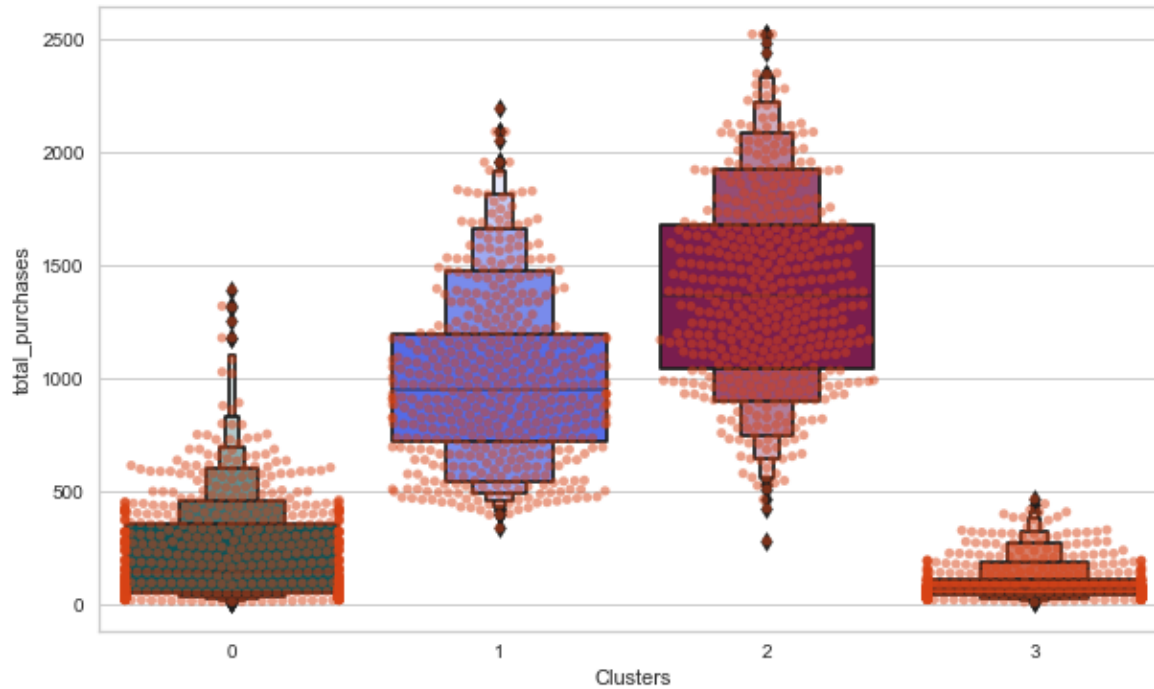
```
pl = sns.scatterplot(data = data,x=data["total_purchases"], y=data["Income"],h  
pl.set_title("Cluster's Profile Based On Income And Spending")  
plt.legend()  
plt.show()
```



The relationship between income and the number of purchases by clusters

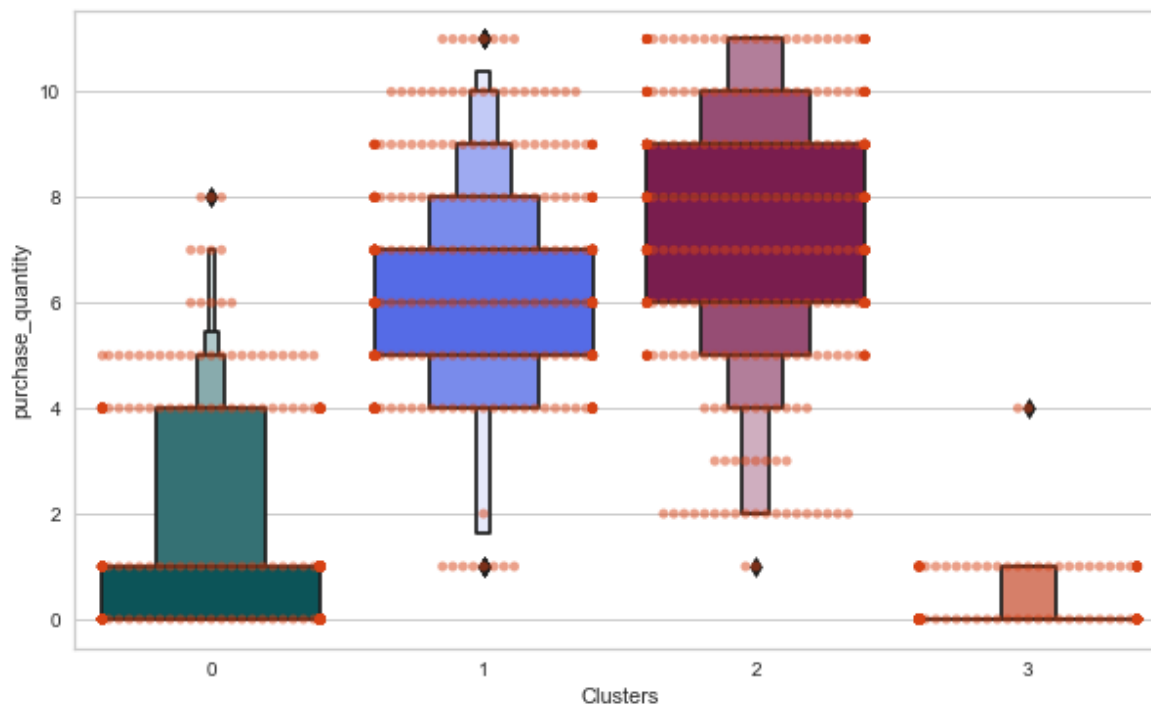
In [1668]:

```
plt.figure(figsize=(10,6))  
pl=sns.swarmplot(x=data["Clusters"], y=data["total_purchases"], color= '#d8431d')  
pl=sns.boxenplot(x=data["Clusters"], y=data["total_purchases"], palette=color_  
plt.show()
```



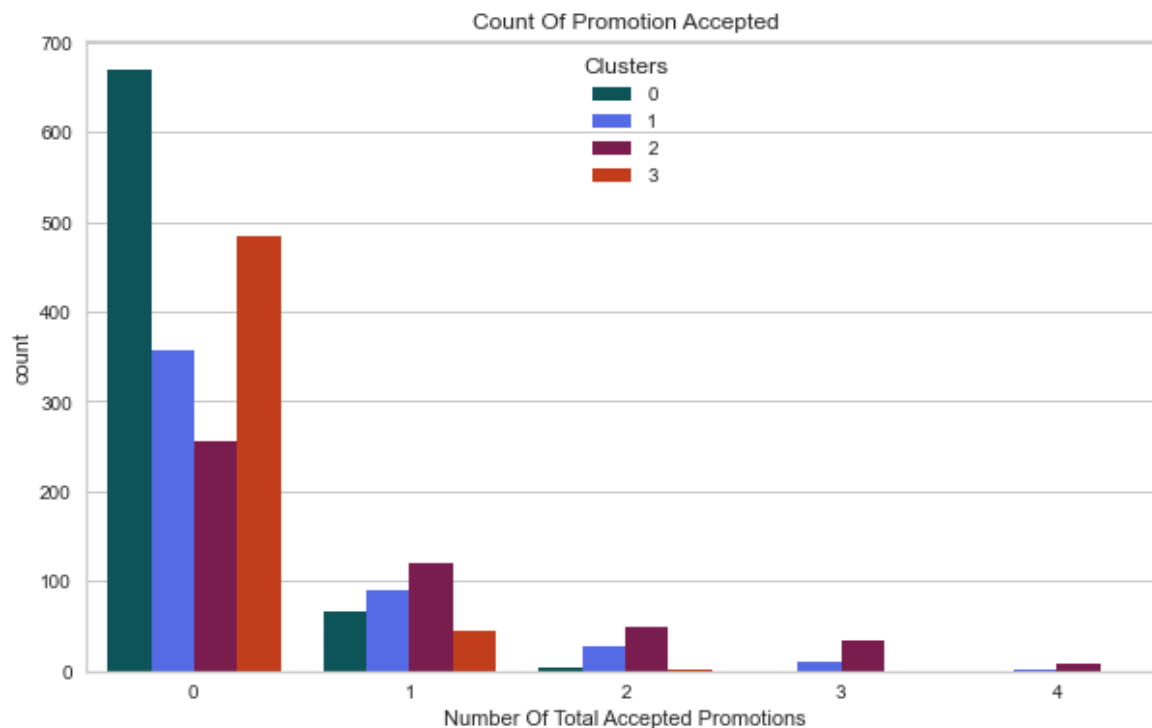
In [1669]:

```
plt.figure(figsize=(10,6))
pl=sns.swarmplot(x=data["Clusters"], y=data["purchase_quantity"], color= '#d84
pl=sns.boxenplot(x=data["Clusters"], y=data["purchase_quantity"], palette=colc
plt.show()
```



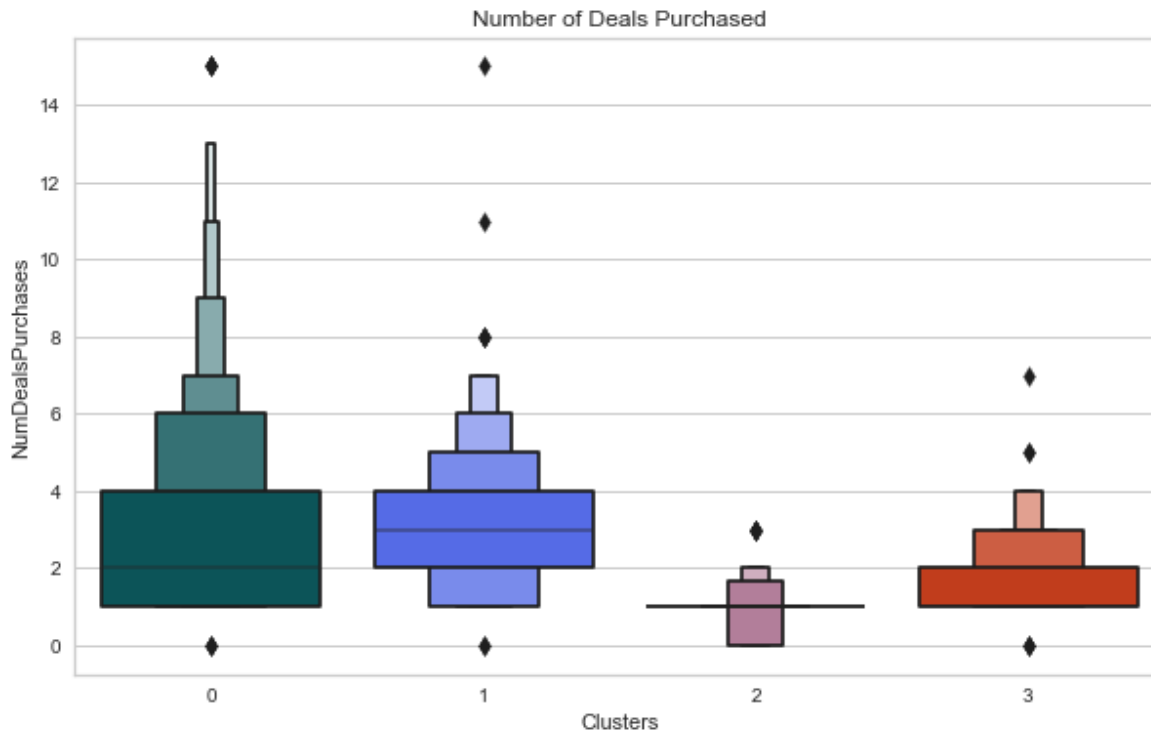
In [1670]:

```
plt.figure(figsize=(10,6))
pl = sns.countplot(x=data["Total_Promos"],hue=data["Clusters"], palette= color
pl.set_title("Count Of Promotion Accepted")
pl.set_xlabel("Number Of Total Accepted Promotions")
plt.show()
```



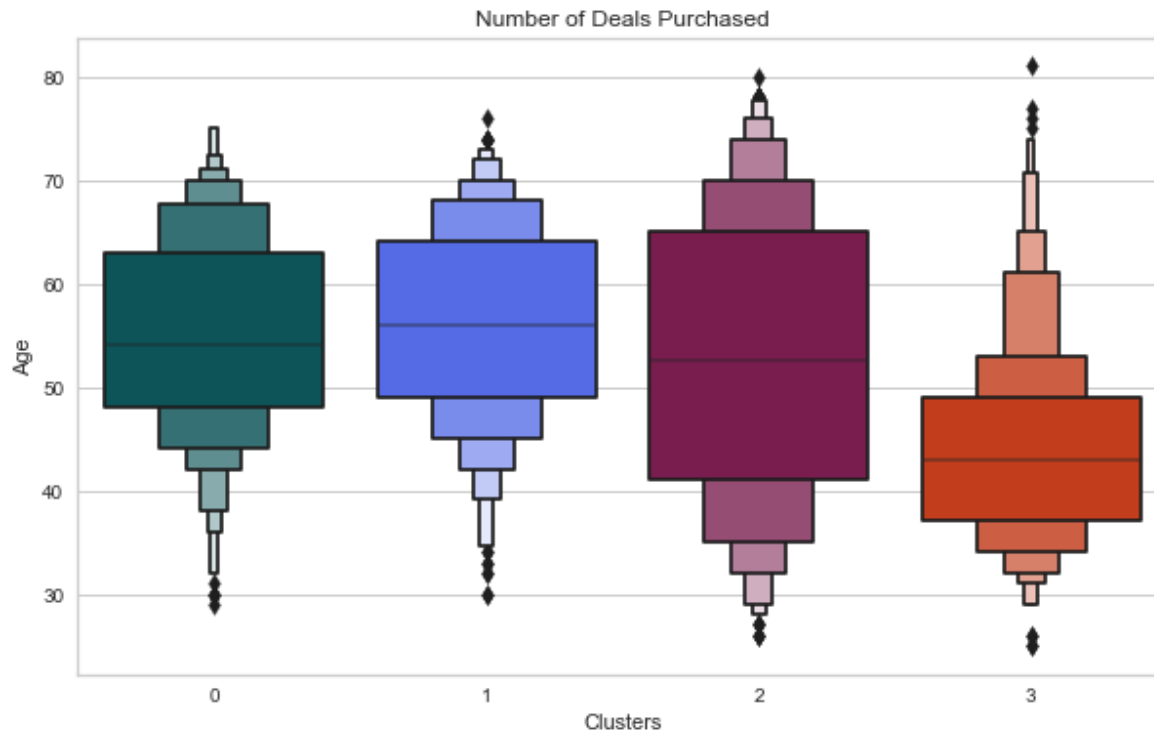
In [1671]:

```
plt.figure(figsize=(10,6))
pl=sns.boxenplot(y=data["NumDealsPurchases"],x=data["Clusters"], palette= colc
pl.set_title("Number of Deals Purchased")
plt.show()
```



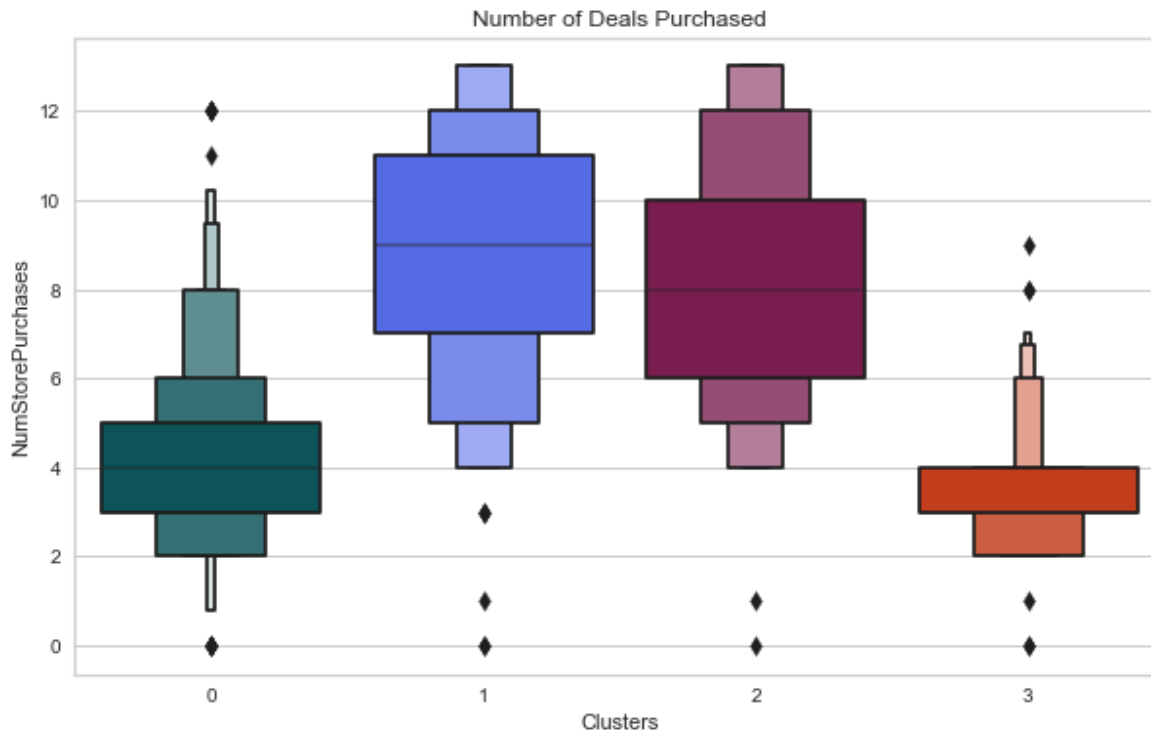
In [1672]:

```
plt.figure(figsize=(10,6))  
pl=sns.boxenplot(y=data["Age"],x=data["Clusters"], palette= color_2)  
pl.set_title("Number of Deals Purchased")  
plt.show()
```



In [1673]:

```
plt.figure(figsize=(10,6))
pl=sns.boxenplot(y=data["NumStorePurchases"],x=data["Clusters"], palette= colc
pl.set_title("Number of Deals Purchased")
plt.show()
```



In [1674]:

```
cmaps
```

Out[1674]:

```
[('Perceptually Uniform Sequential',  
  ['viridis', 'plasma', 'inferno', 'magma', 'cividis']),  
( 'Sequential',  
  ['Greys',  
   'Purples',  
   'Blues',  
   'Greens',  
   'Oranges',  
   'Reds',  
   'YlOrBr',  
   'YlOrRd',  
   'OrRd',  
   'PuRd',  
   'RdPu',  
   'BuPu',  
   'GnBu',  
   'PuBu',  
   'YlGnBu',  
   'PuBuGn',  
   'BuGn',  
   'YlGn']),  
( 'Sequential (2)',  
  ['binary',  
   'gist_yarg',  
   'gist_gray',  
   'gray',  
   'bone',  
   'pink',  
   'spring',  
   'summer',  
   'autumn',  
   'winter',  
   'cool',  
   'Wistia',  
   'hot',  
   'afmhot',  
   'gist_heat',  
   'copper']),  
( 'Diverging',  
  ['PiYG',  
   'PRGn',
```

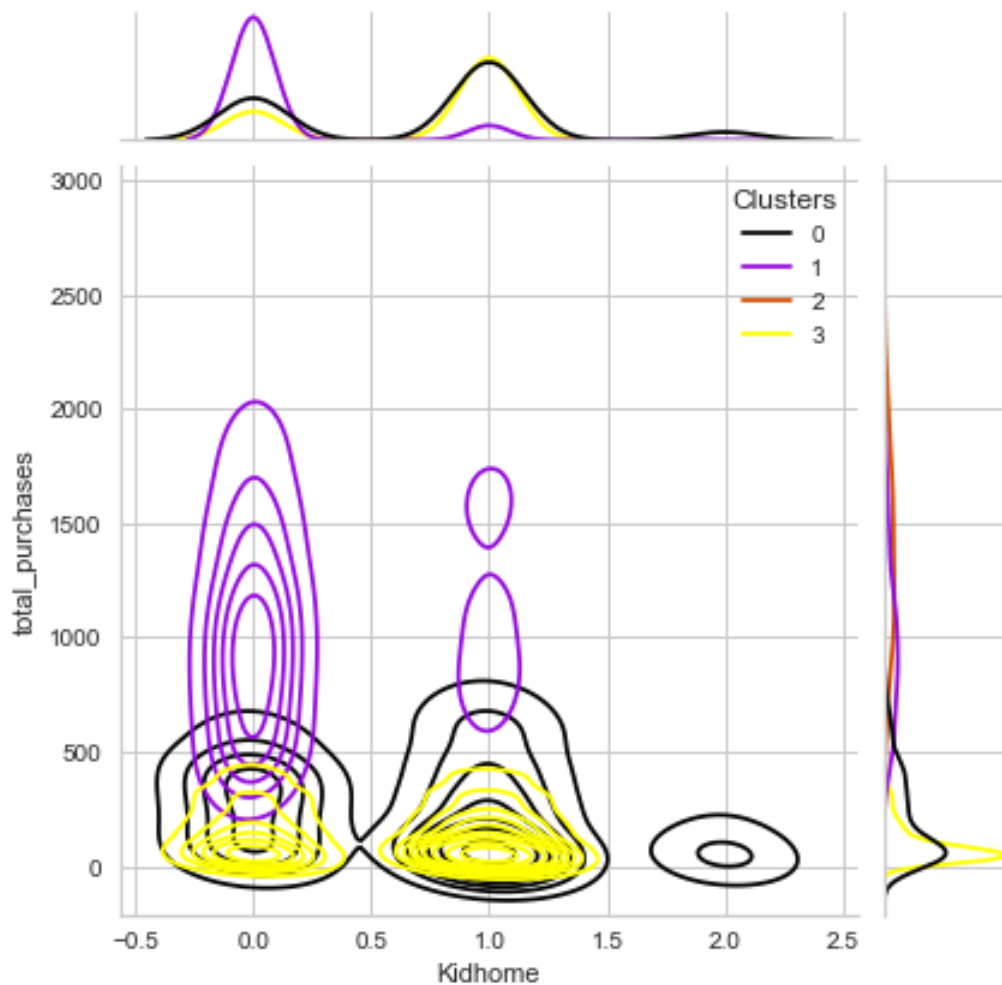
```
'BrBG',
'PuOr',
'RdGy',
'RdBu',
'RdYlBu',
'RdYlGn',
'Spectral',
'coolwarm',
'bwr',
'seismic']]),
('Cyclic', ['twilight', 'twilight_shifted', 'hsv']),
('Qualitative',
 ['Pastel1',
  'Pastel2',
  'Paired',
  'Accent',
  'Dark2',
  'Set1',
  'Set2',
  'Set3',
  'tab10',
  'tab20',
  'tab20b',
  'tab20c']]),
('Miscellaneous',
 ['flag',
  'prism',
  'ocean',
  'gist_earth',
  'terrain',
  'gist_stern',
  'gnuplot',
  'gnuplot2',
  'CMRmap',
  'cubehelix',
  'brg',
  'gist_rainbow',
  'rainbow',
  'jet',
  'turbo',
  'nipy_spectral',
  'gist_ncar']])]
```

In [1676]:

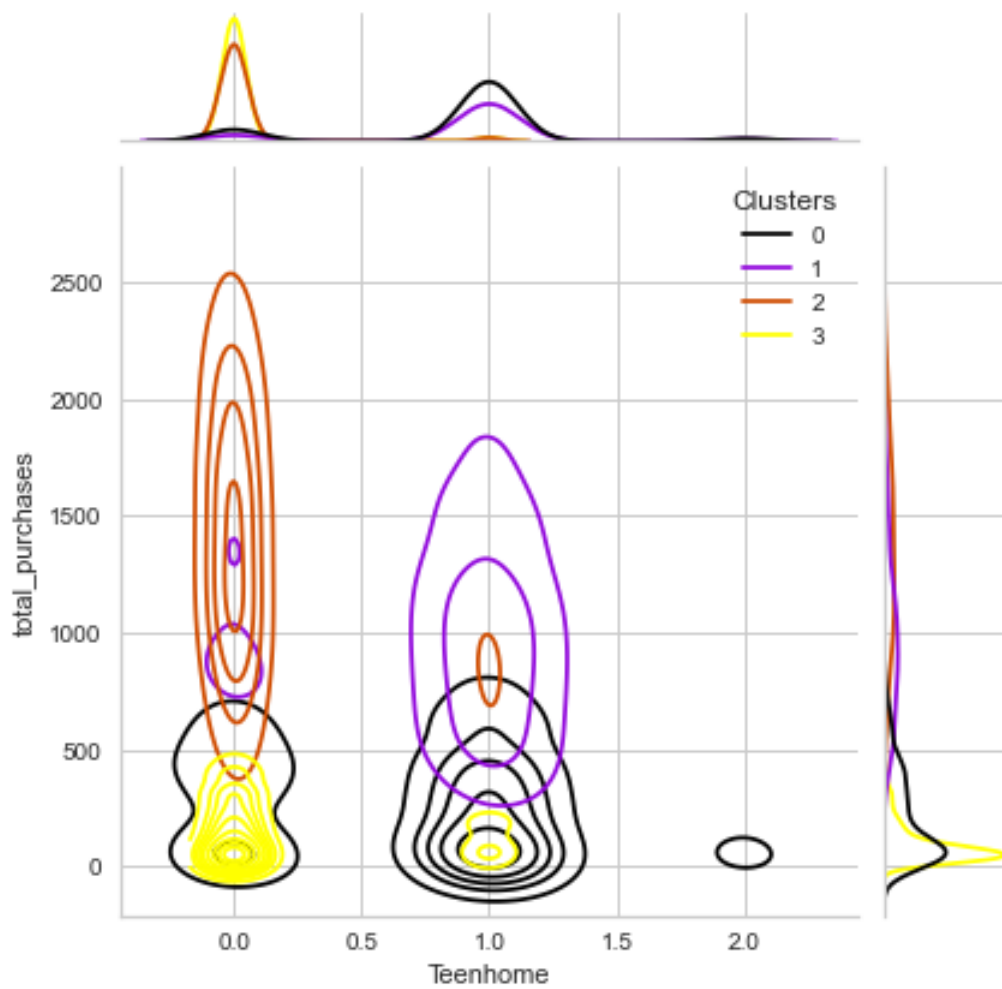
```
Personal = [ "Kidhome", "Teenhome", "Customer_From_days", "Age", "Num_Children",

for i in Personal:
    plt.figure()
    sns.jointplot(x=data[i], y=data["total_purchases"], hue =data["Clusters"],
    plt.show()
```

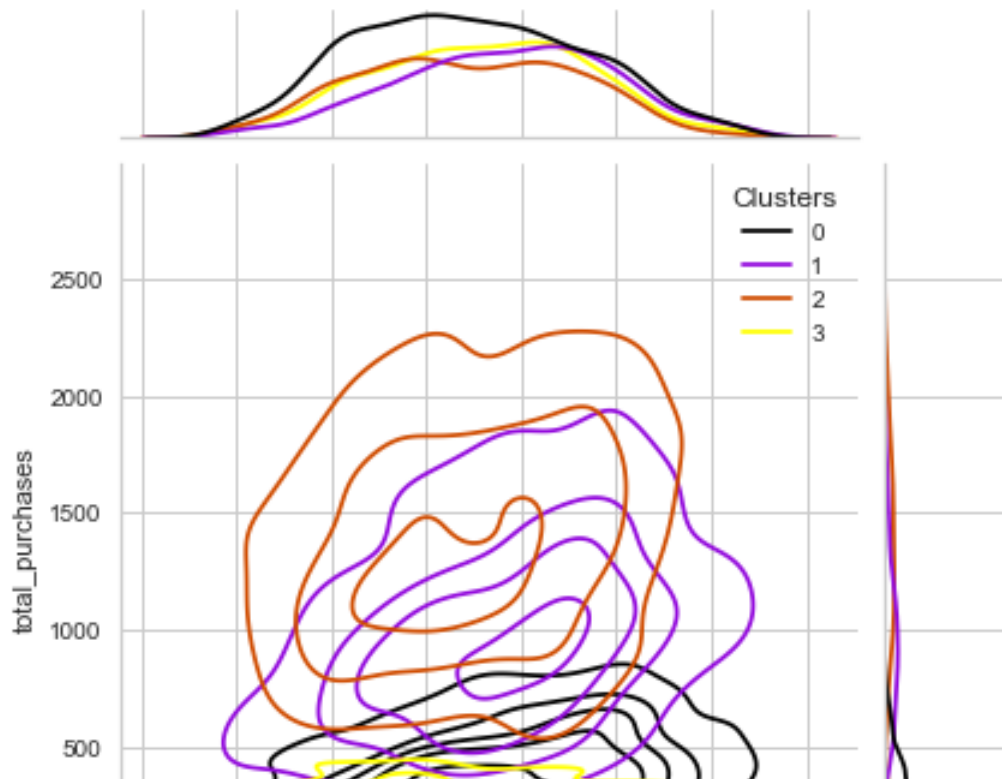
<Figure size 1080x720 with 0 Axes>



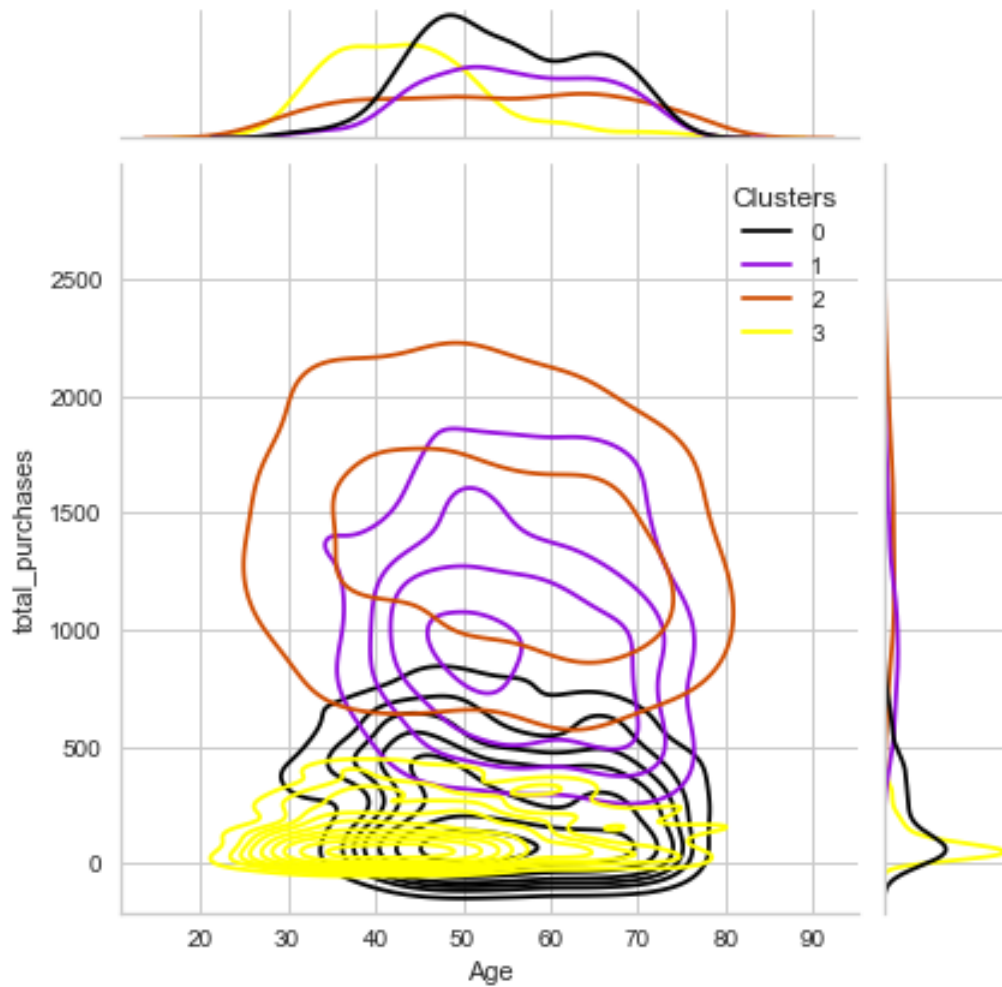
<Figure size 1080x720 with 0 Axes>



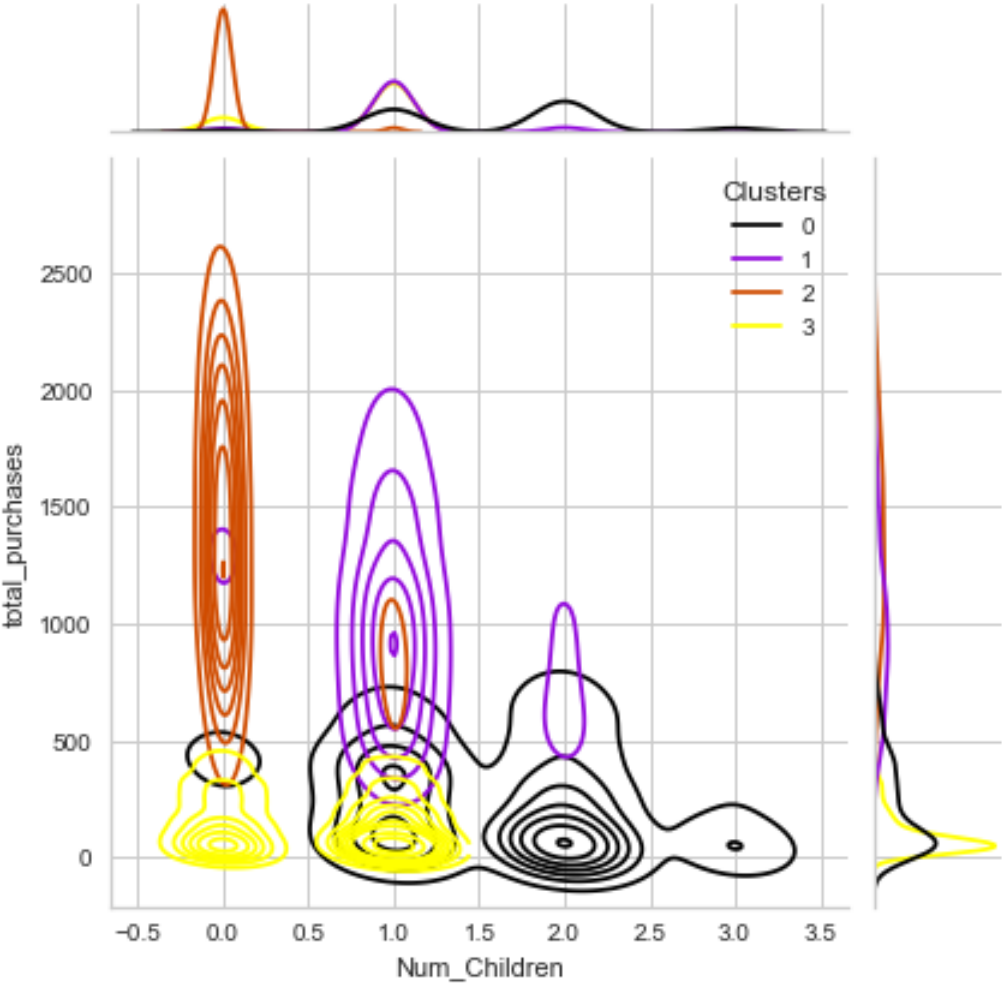
<Figure size 1080x720 with 0 Axes>



<Figure size 1080x720 with 0 Axes>

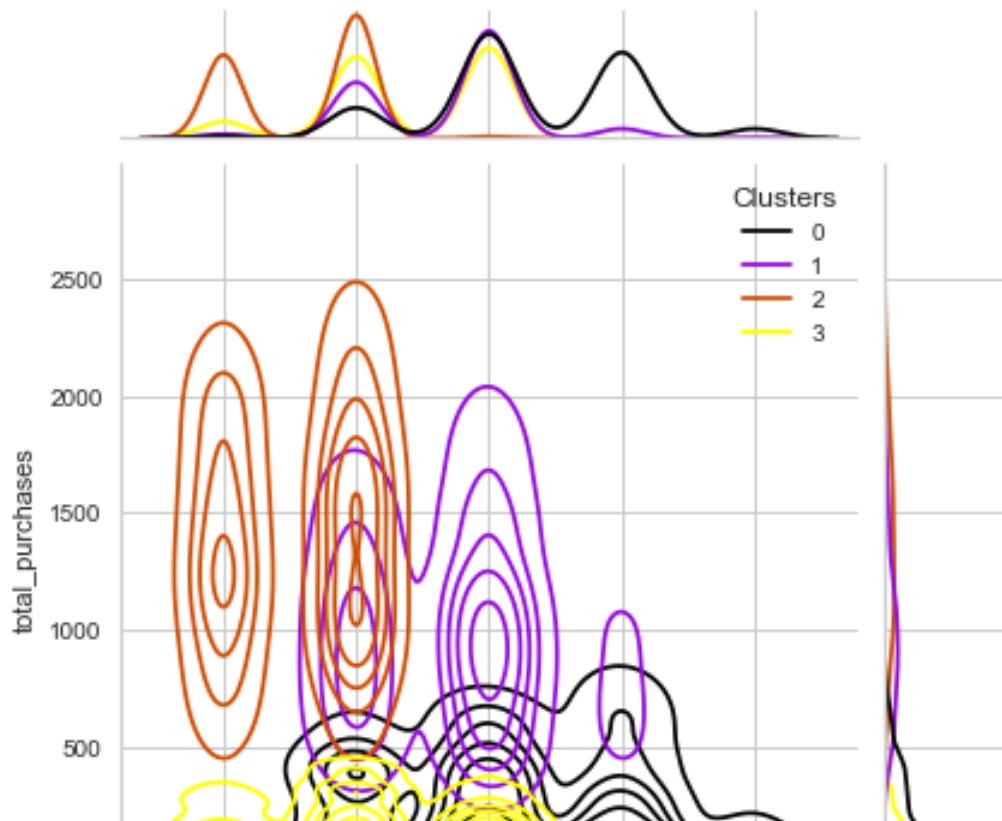


<Figure size 1080x720 with 0 Axes>

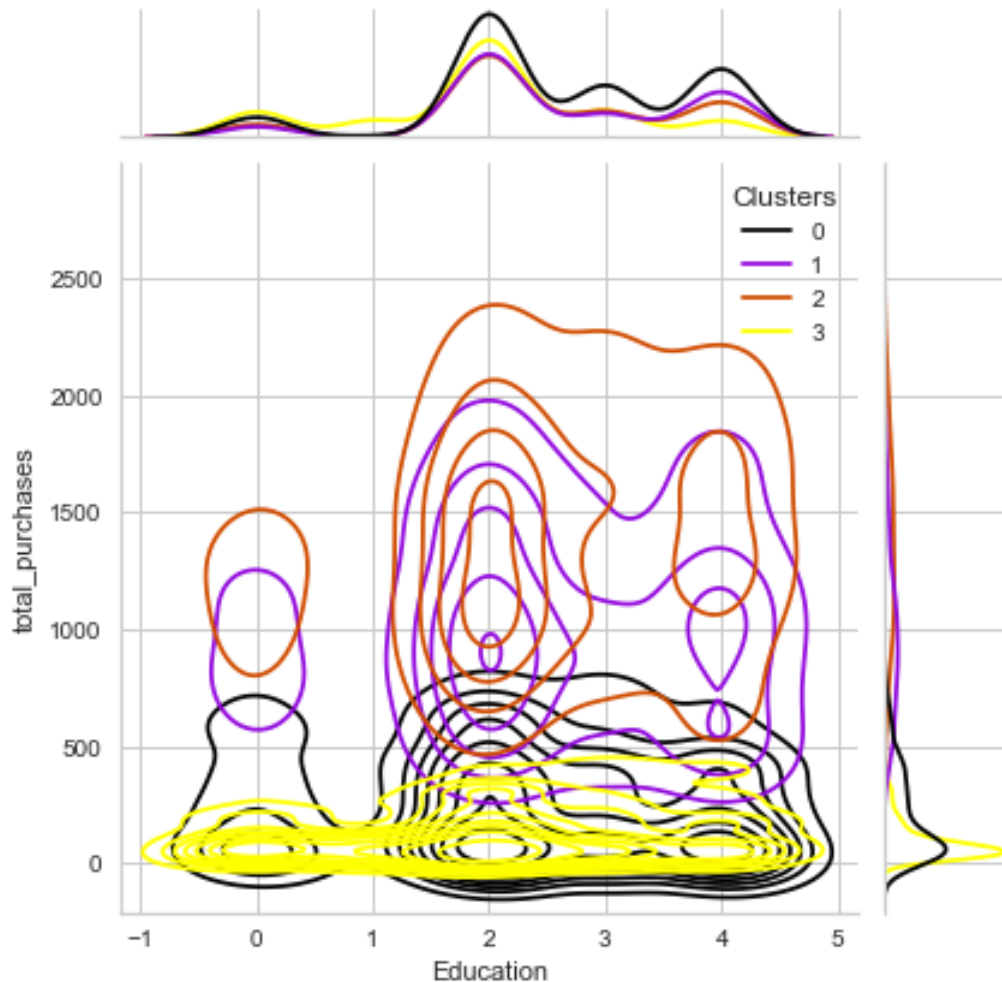


<Figure size 1080x720 with 0 Axes>





<Figure size 1080x720 with 0 Axes>



In []:

In []:

In []:

In []: