

In [1]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import warnings
6 warnings.filterwarnings('ignore')
```

In [2]:

```
1 df=pd.read_csv("customer_data.csv")
```

In [3]:

```
1 df.head()
```

Out[3]:

ID	Year_Birth	Education	Marital_Status	Income	Kidhome	Teenhome	Dt_Customer	Recenc
0	5524	1957	Graduation	Single	58138.0	0	04/09/12	5
1	2174	1954	Graduation	Single	46344.0	1	08/03/14	3
2	4141	1965	Graduation	Together	71613.0	0	21/08/13	2
3	6182	1984	Graduation	Together	26646.0	1	10/02/14	2
4	5324	1981	PhD	Married	58293.0	1	19/01/14	9

5 rows × 9 columns



In [4]:

```
1 df.columns
```

Out[4]:

```
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'],
      dtype='object')
```

EDA

In [5]:

```
1 #Separating Out Date,Month & Year of Customer Enrollment
2 df["Date_enroll"]=df["Dt_Customer"].str.split("/").str[0]
3 df["Month_enroll"]=df["Dt_Customer"].str.split("/").str[1]
4 df["Year_enroll"]=df["Dt_Customer"].str.split("/").str[2]
```

In [6]:

```
1 df.head()
```

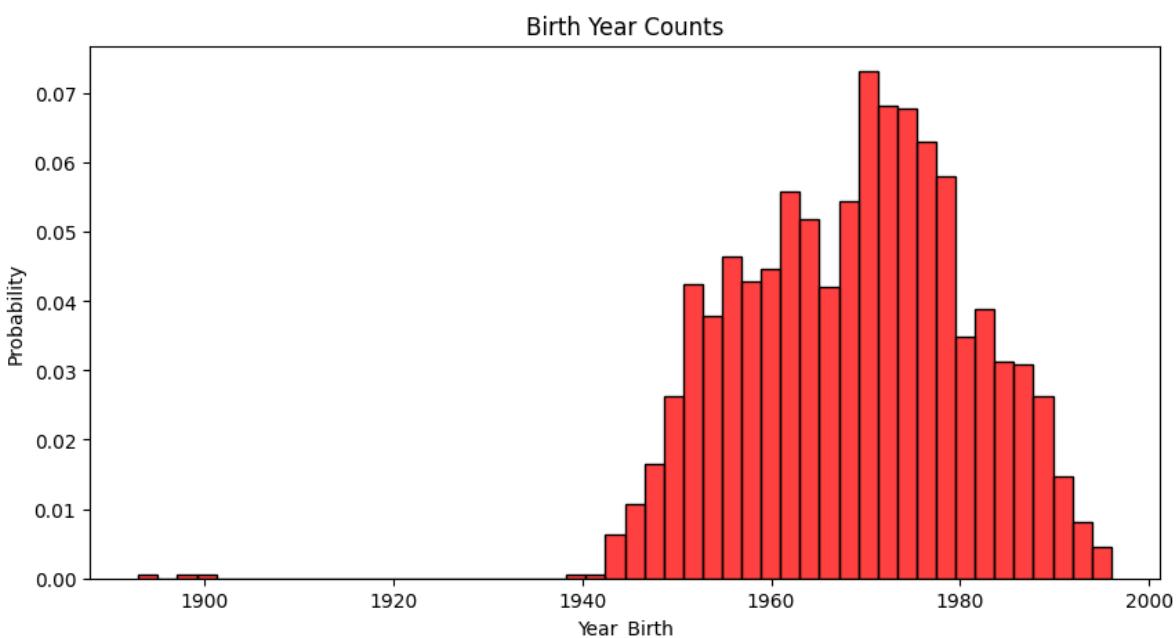
Out[6]:

ID	Year_Birth	Education	Marital_Status	Income	Kidhome	Teenhome	Dt_Customer	Recenc
0	5524	1957	Graduation	Single	58138.0	0	04/09/12	5
1	2174	1954	Graduation	Single	46344.0	1	08/03/14	3
2	4141	1965	Graduation	Together	71613.0	0	21/08/13	2
3	6182	1984	Graduation	Together	26646.0	1	10/02/14	2
4	5324	1981	PhD	Married	58293.0	1	19/01/14	9

5 rows × 32 columns

In [7]:

```
1 #Birth Year Values
2 plt.figure(figsize=(10,5))
3 sns.histplot(df["Year_Birth"],bins=50,color="red",stat="probability")
4 plt.title("Birth Year Counts")
5 plt.show()
```

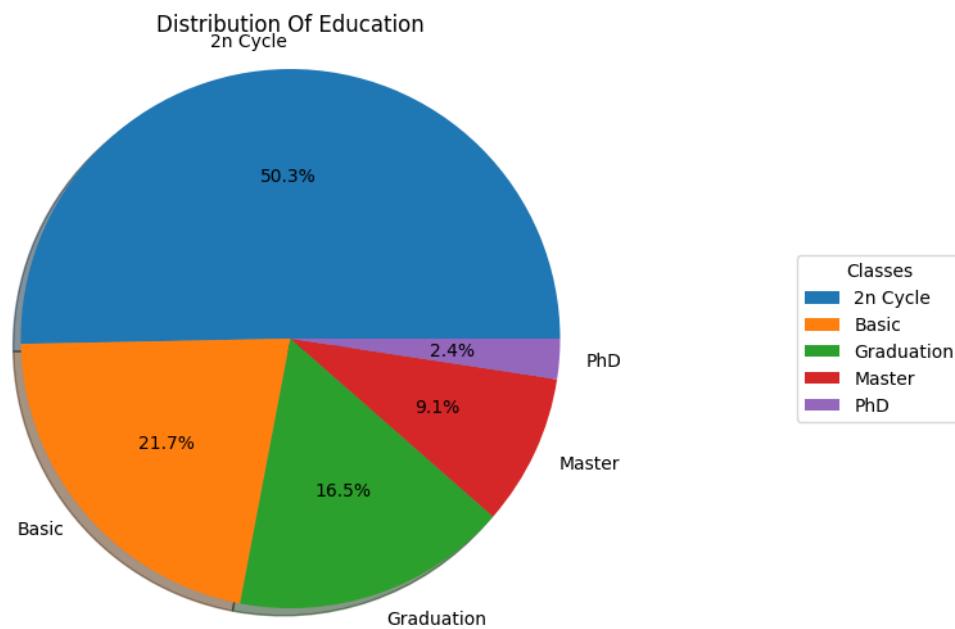


1. Dataset is of people born between 1940-2000.

2 or 3 variables of data are from year 1900 or earlier it can be result of wrong entry or an outlier.

In [8]:

```
1 #Types Of Education
2 plt.figure(figsize=(10,6))
3 plt.pie(df["Education"].value_counts(), autopct="%1.1f%%", labels=np.unique(df["Education"]))
4 plt.title("Distribution Of Education")
5 plt.legend(title="Classes", loc="center left", bbox_to_anchor=(1, 0.5))
6 plt.axis("equal")
7 plt.show()
```

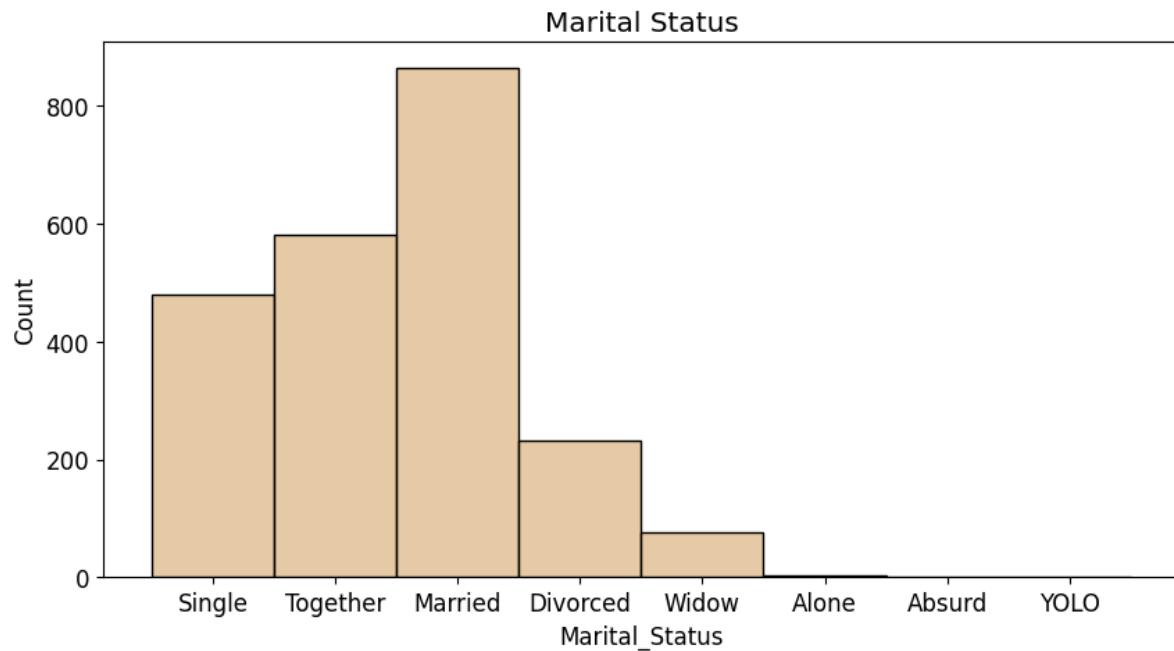


1.Education is divided into 5 types:

- i.Graduation ---16.5%
- ii.Basic ---21.7%
- iii.Master ---9.1%
- iv.PhD ---2.4%
- v.2n Cycle ---50.3%

In [29]:

```
1 #Marital Status Values
2 plt.figure(figsize=(10,5))
3 sns.histplot(df["Marital_Status"],bins=8,color="burlywood",element="bars")
4 plt.title("Marital Status")
5 plt.show()
```

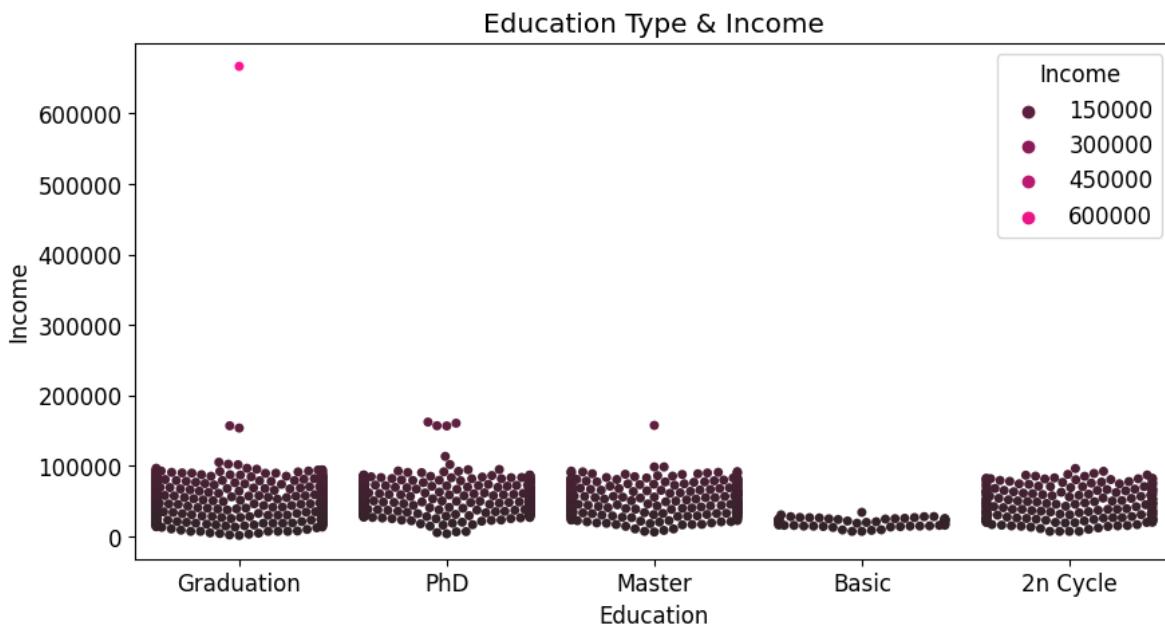


1. Marital Status is divided into 8 types:

'Single', 'Together', 'Married', 'Divorced', 'Widow', 'Alone', 'Absurd', 'YOLO'

In [30]:

```
1 #Education Type & Income
2 plt.figure(figsize=(10,5))
3 sns.swarmplot(x=df["Education"],y=df["Income"],hue=df["Income"],color="deeppink")
4 plt.title("Education Type & Income")
5 plt.show()
```

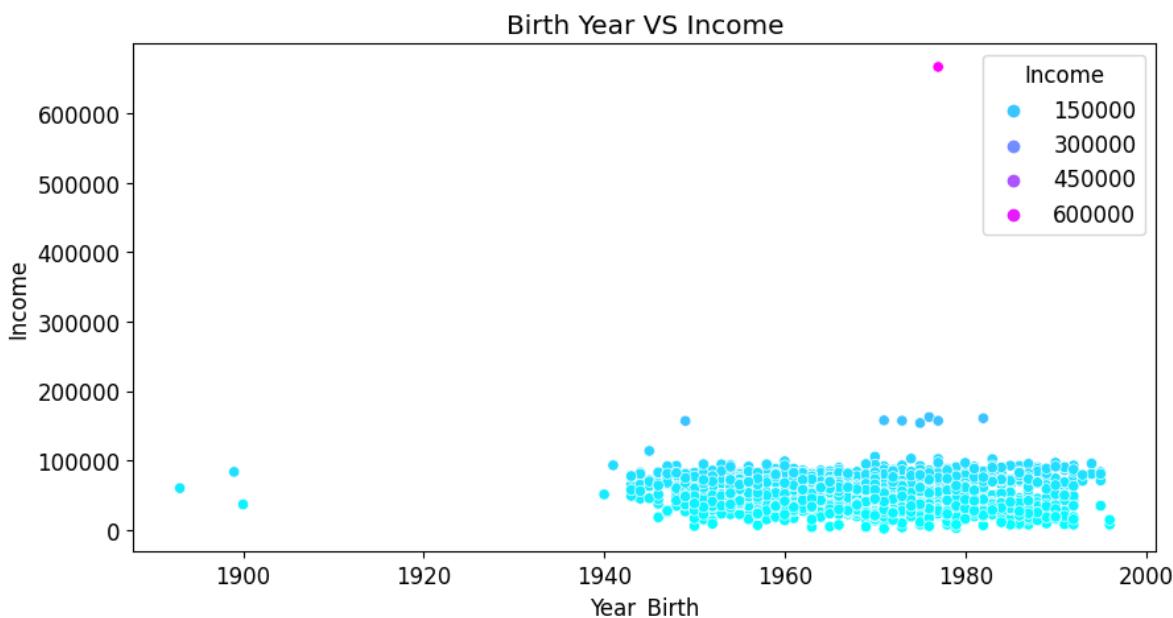


1.Income of Basic education is less compared to others.

2.Highest Income generated is from Graduation but it can be wrong entry or outlier as well.

In [31]:

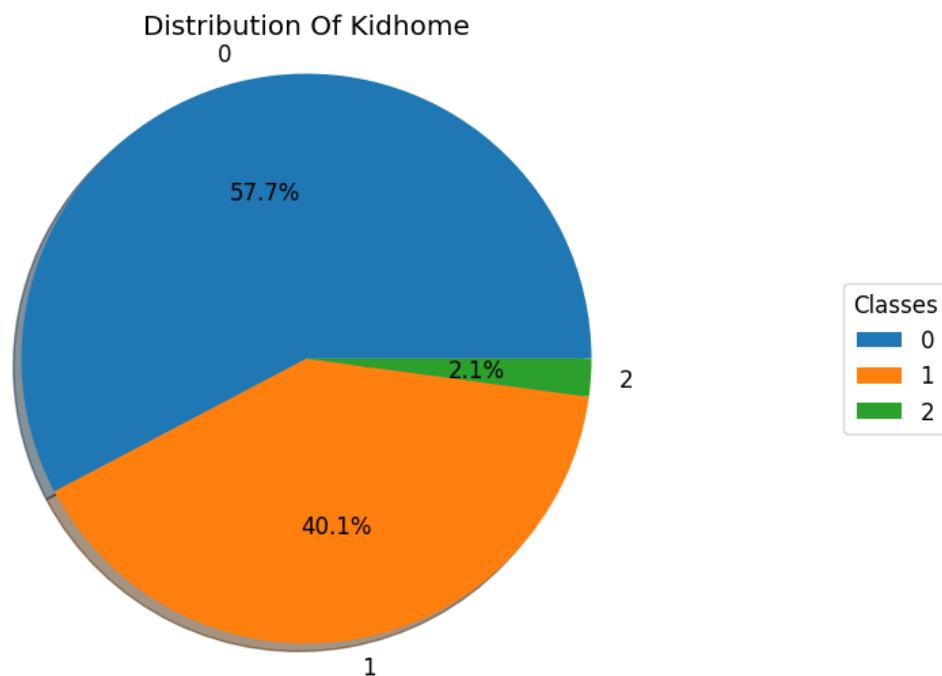
```
1 #Birth Year VS Income
2 plt.figure(figsize=(10,5))
3 sns.scatterplot(data=df, x="Year_Birth", y="Income",hue="Income",palette="cool")
4 plt.title("Birth Year VS Income")
5 plt.show()
```



1..Birth year wise income is distributed equally but some group of people from birth year 1970-1980 are earning more compared to others.

In [32]:

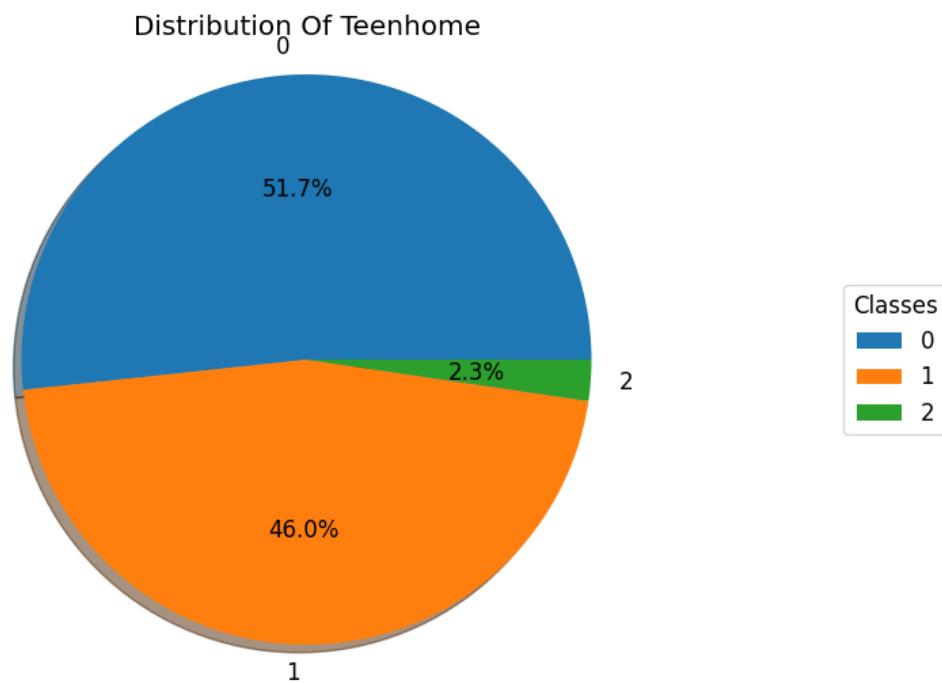
```
1 #Distribution Of Kidhome
2 plt.figure(figsize=(10,6))
3 plt.pie(df["Kidhome"].value_counts(), autopct="%1.1f%%", labels=np.unique(df["Kidhome"]), sha
4 plt.title("Distribution Of Kidhome")
5 plt.legend(title="Classes", loc="center left", bbox_to_anchor=(1, 0.5))
6 plt.axis("equal")
7 plt.show()
```



1.50.7% of people do not have kids,40.1% people have 1 child whereas 2.1% people have 2 childs.

In [33]:

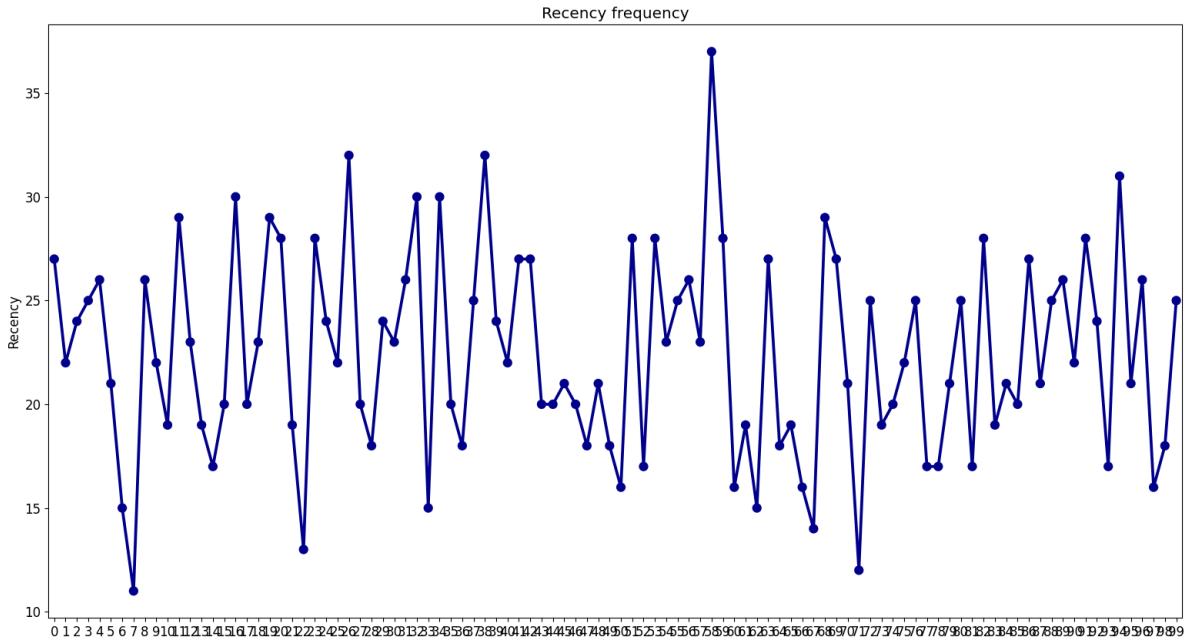
```
1 #Distribution Of Teenhome
2 plt.figure(figsize=(10,6))
3 plt.pie(df["Teenhome"].value_counts(), autopct="%1.1f%%", labels=np.unique(df["Teenhome"]),
4 plt.title("Distribution Of Teenhome")
5 plt.legend(title="Classes", loc="center left", bbox_to_anchor=(1, 0.5))
6 plt.axis("equal")
7 plt.show()
```



1.51.7% of people do not have teens at their home,46.0% of people have 1 teen wheras 2.3% of people have 2 teens at home.

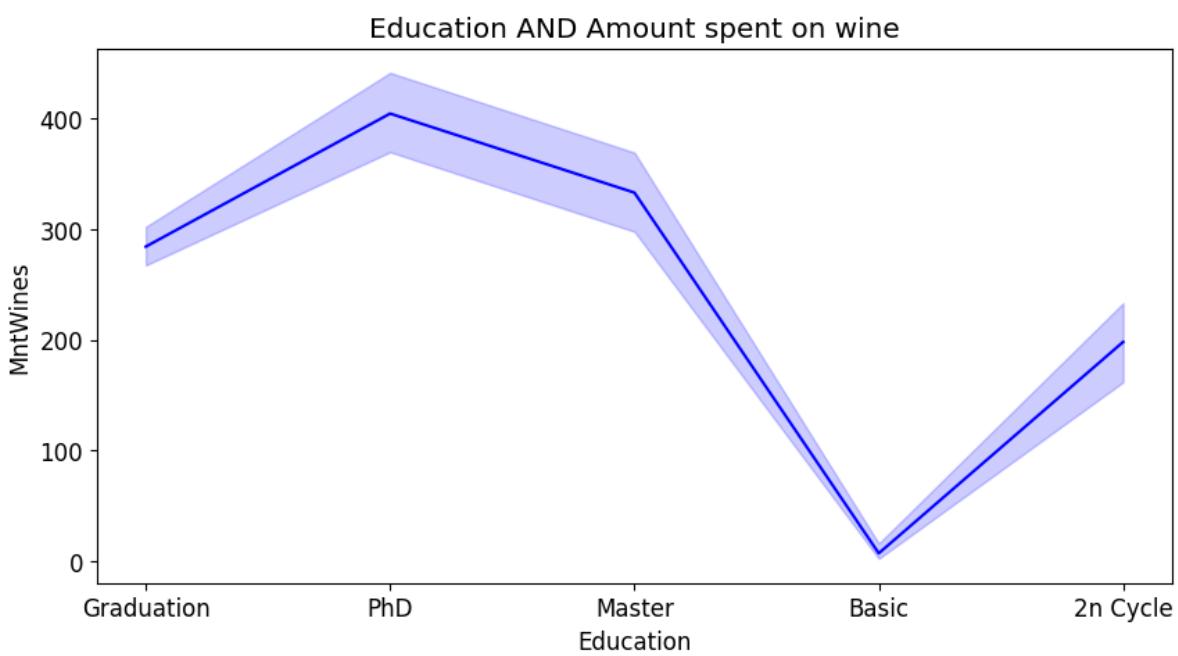
In [34]:

```
1 #Recency
2 plt.figure(figsize=(19,10))
3 sns.pointplot(x=df["Recency"].unique(),y=df["Recency"].value_counts(),color="darkblue")
4 plt.title("Recency frequency")
5 plt.show()
```



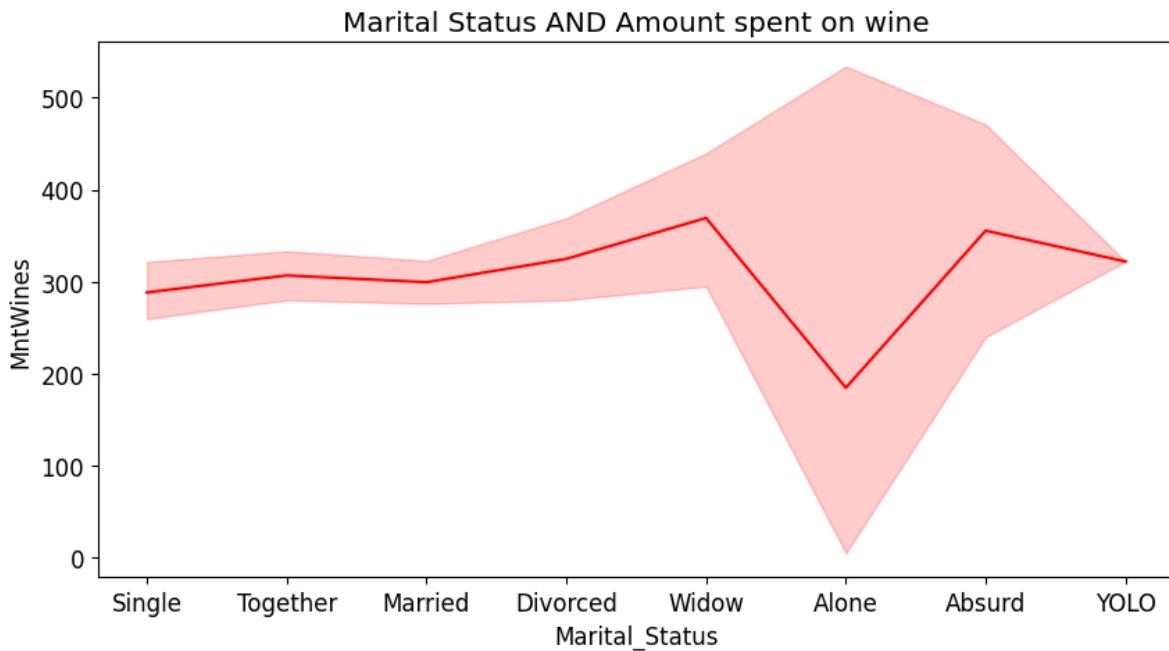
In [35]:

```
1 #MntWines:Amount Spent On wines
2 #Education & Wines
3 plt.figure(figsize=(10,5))
4 sns.lineplot(data=df,x="Education",y="MntWines",color="blue")
5 plt.title("Education AND Amount spent on wine")
6 plt.show()
```



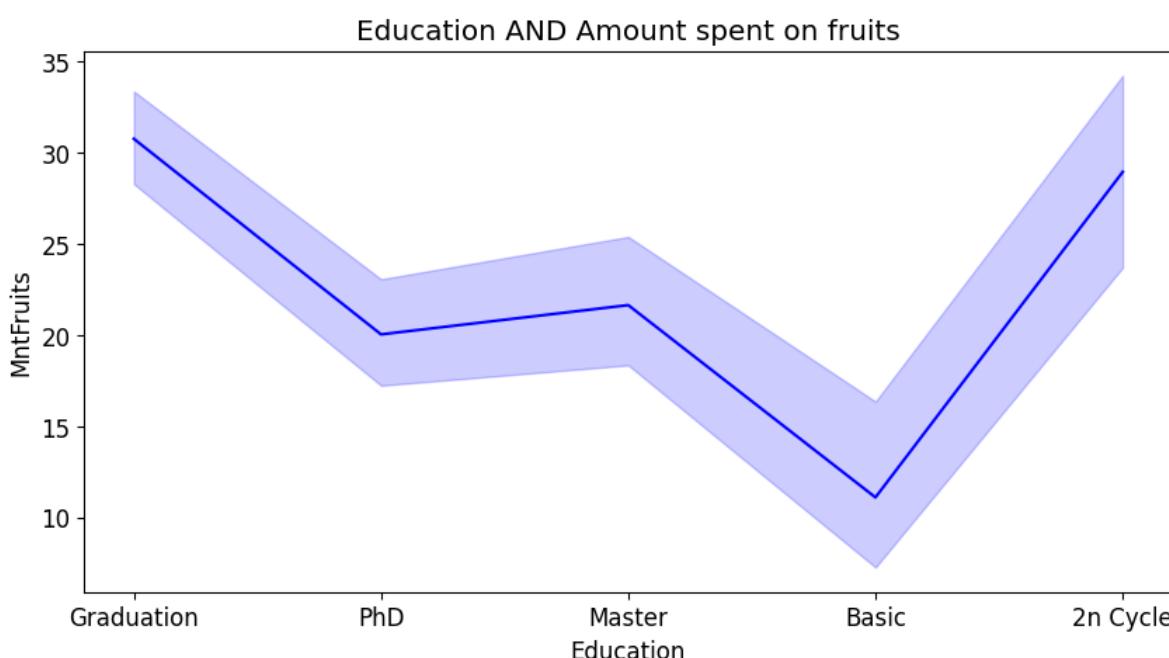
In [36]:

```
1 #Marital Status & Wines
2 plt.figure(figsize=(10,5))
3 sns.lineplot(data=df,x="Marital_Status",y="MntWines",color="red")
4 plt.title("Marital Status AND Amount spent on wine")
5 plt.show()
```



In [37]:

```
1 #MntFruits:Amount Spent On Fruits
2 #Education & Fruits
3 plt.figure(figsize=(10,5))
4 sns.lineplot(data=df,x="Education",y="MntFruits",color="blue")
5 plt.title("Education AND Amount spent on fruits")
6 plt.show()
```

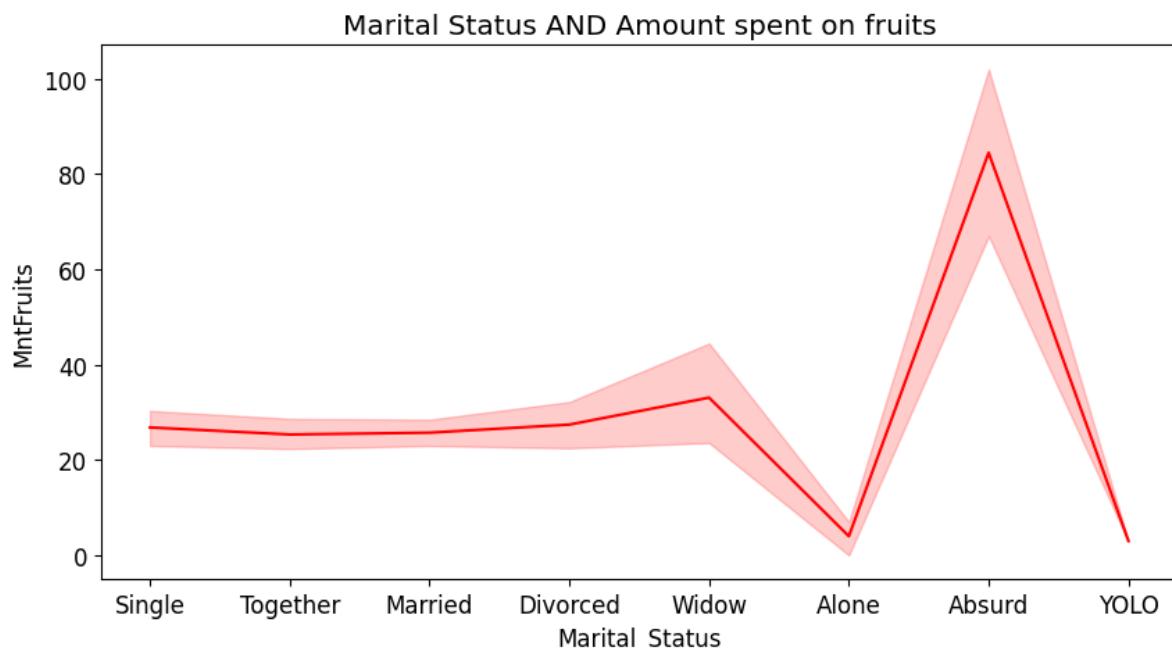


1. Maximum/Highest amount is spend on wine.

2. People with basic education are not spending any amount on wines.

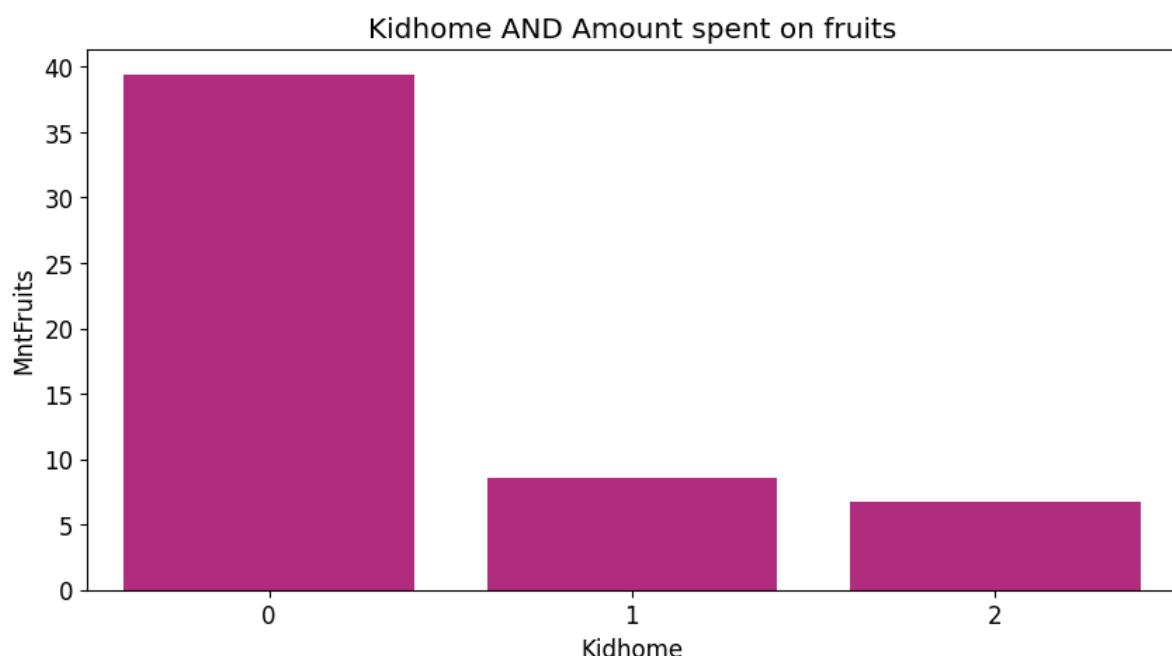
In [38]:

```
1 #Marital Status & fruits
2 plt.figure(figsize=(10,5))
3 sns.lineplot(data=df,x="Marital_Status",y="MntFruits",color="red")
4 plt.title("Marital Status AND Amount spent on fruits")
5 plt.show()
```



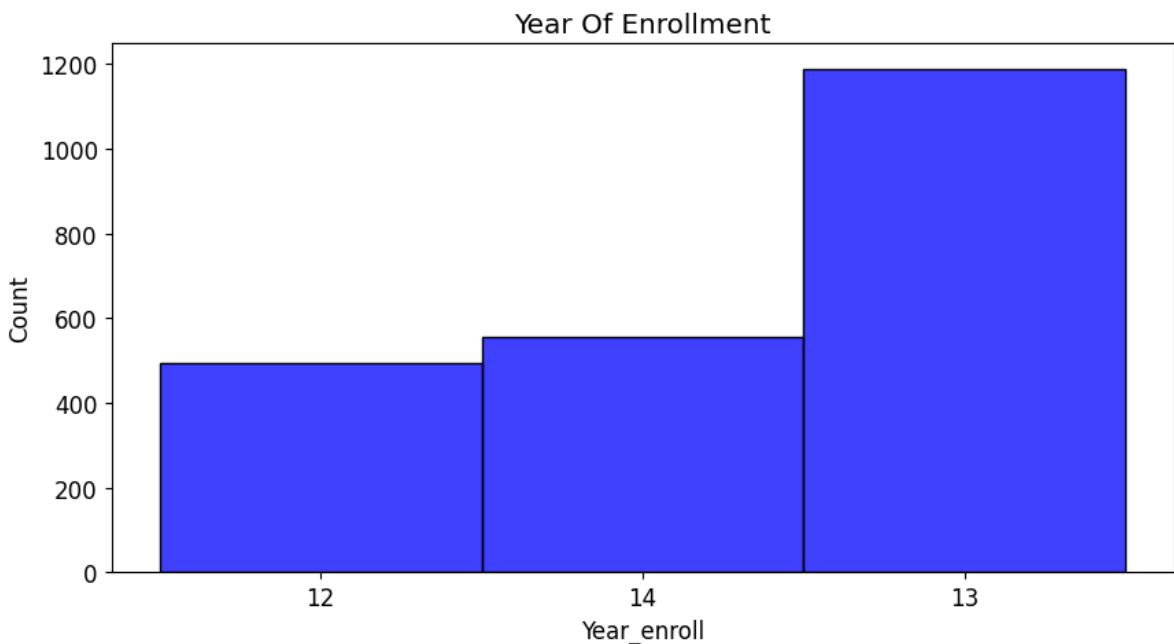
In [39]:

```
1 #Kidhome & fruits
2 plt.figure(figsize=(10,5))
3 sns.barplot(data=df,x="Kidhome",y="MntFruits",color="mediumvioletred",ci=None)
4 plt.title("Kidhome AND Amount spent on fruits")
5 plt.show()
```



In [40]:

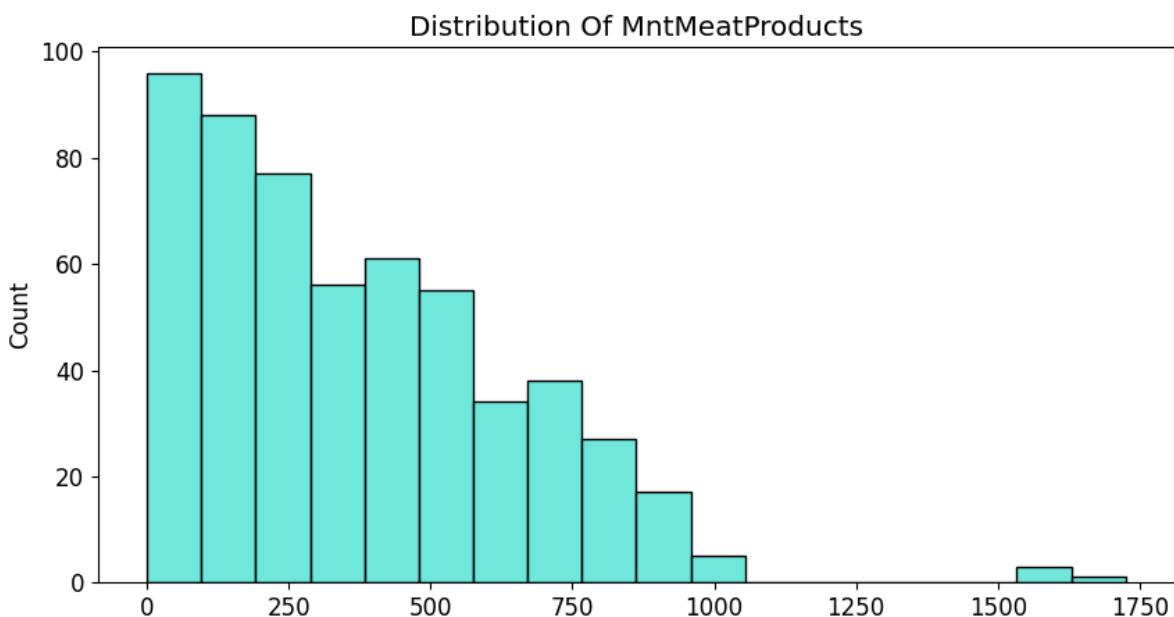
```
1 #Enrollment year wise
2 plt.figure(figsize=(10,5))
3 sns.histplot(df["Year_enroll"],color="blue")
4 plt.title("Year Of Enrollment")
5 plt.show()
```



1. Most of the registrations are from year 2013 in dataset provided of 2 years.

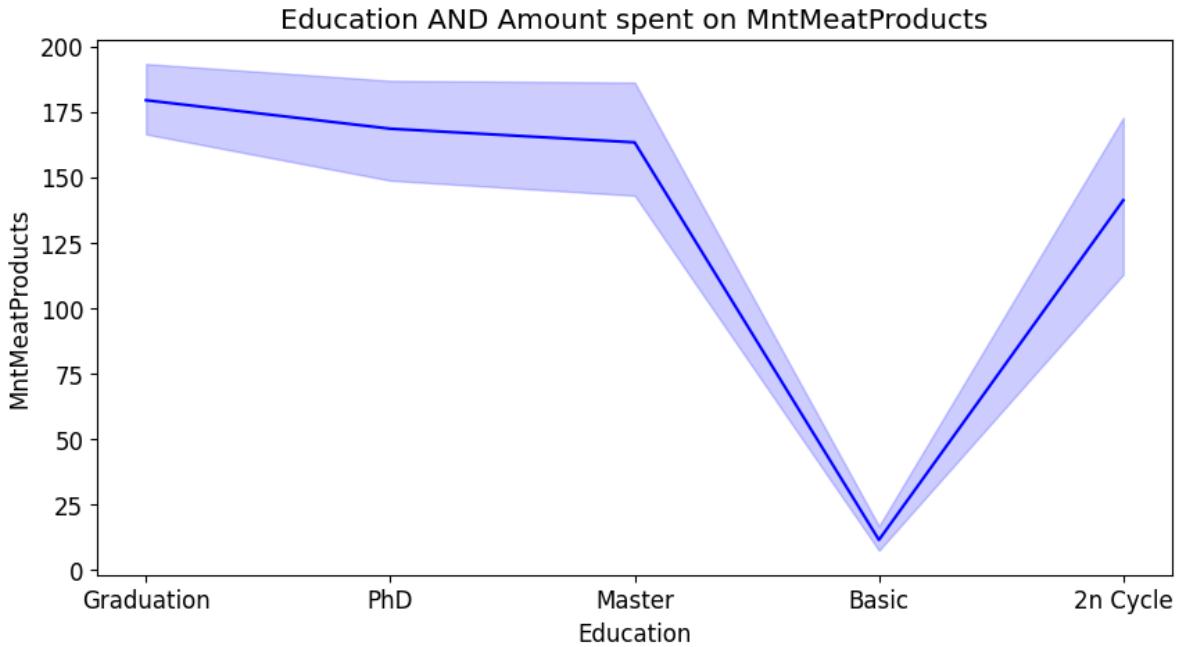
In [41]:

```
1 #Distribution Of MntMeatProducts
2 plt.figure(figsize=(10,5))
3 sns.histplot(df["MntMeatProducts"].unique(),color="turquoise")
4 plt.title("Distribution Of MntMeatProducts")
5 plt.show()
```



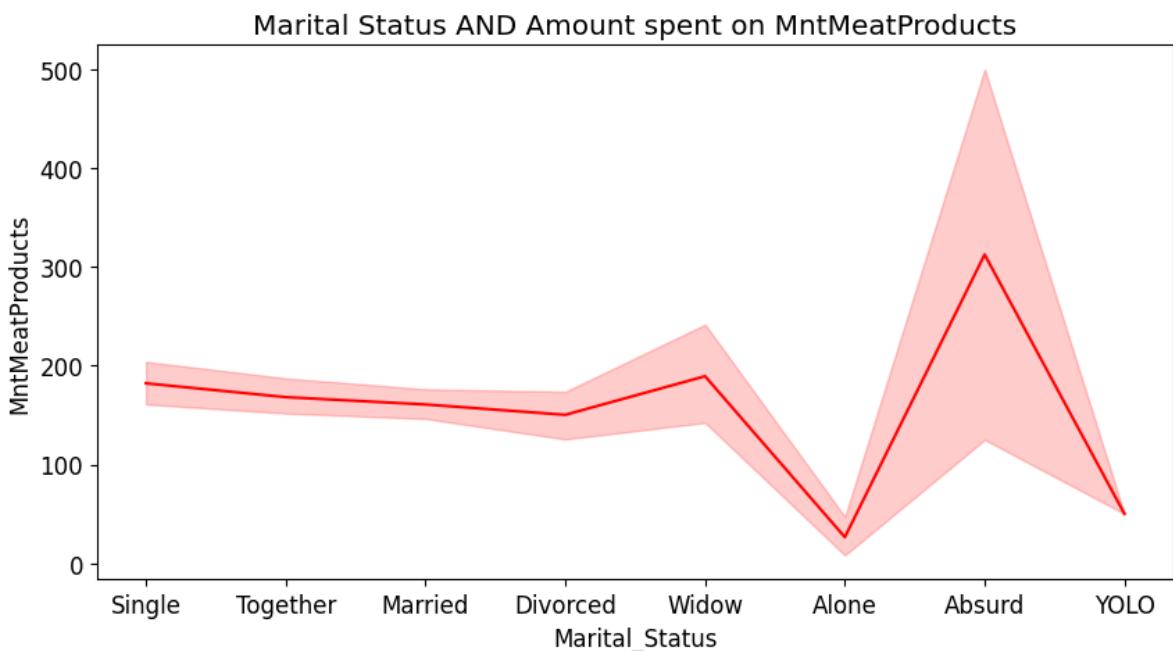
In [42]:

```
1 #MntFruits:Amount Spent On MntMeatProducts
2 #Education & MntMeatProducts
3 plt.figure(figsize=(10,5))
4 sns.lineplot(data=df,x="Education",y="MntMeatProducts",color="blue")
5 plt.title("Education AND Amount spent on MntMeatProducts")
6 plt.show()
```



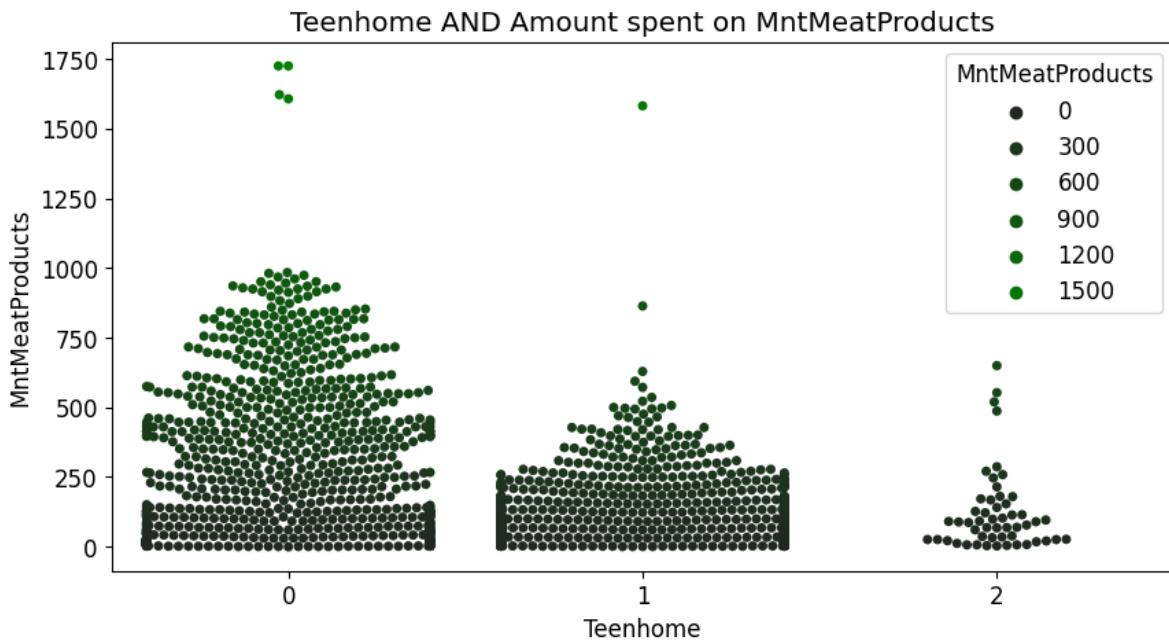
In [43]:

```
1 #Marital Status & MntMeatProducts
2 plt.figure(figsize=(10,5))
3 sns.lineplot(data=df,x="Marital_Status",y="MntMeatProducts",color="red")
4 plt.title("Marital Status AND Amount spent on MntMeatProducts")
5 plt.show()
```



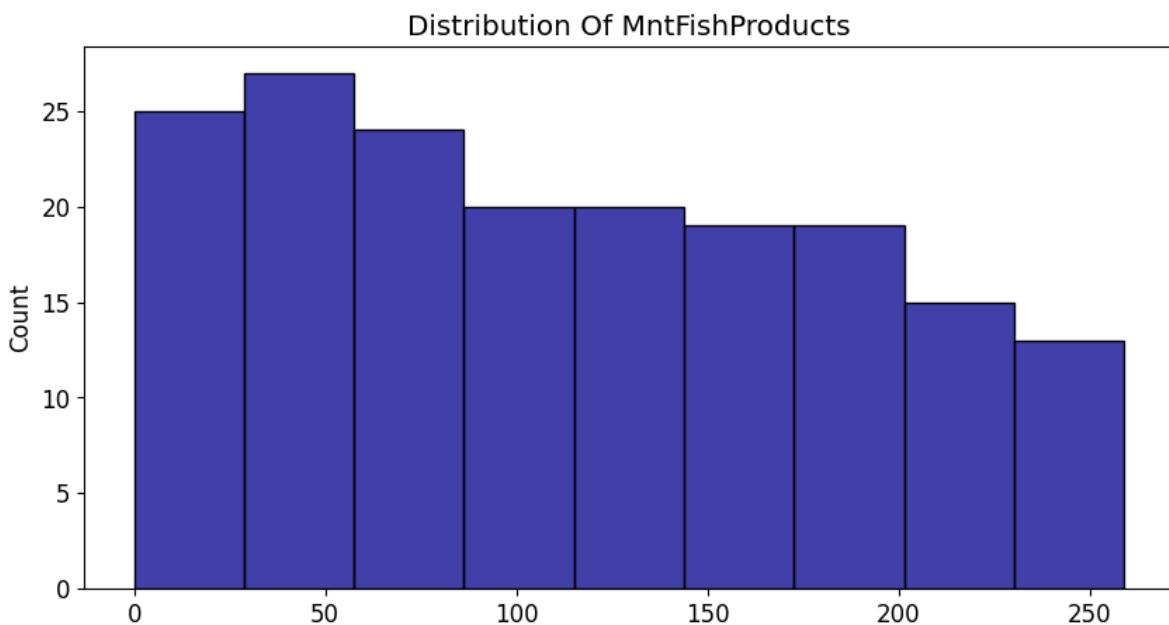
In [44]:

```
1 # Teenhome & MntMeatProducts
2 plt.figure(figsize=(10,5))
3 sns.swarmplot(data=df,x="Teenhome",y="MntMeatProducts",color="green",hue="MntMeatProducts")
4 plt.title("Teenhome AND Amount spent on MntMeatProducts")
5 plt.show()
```



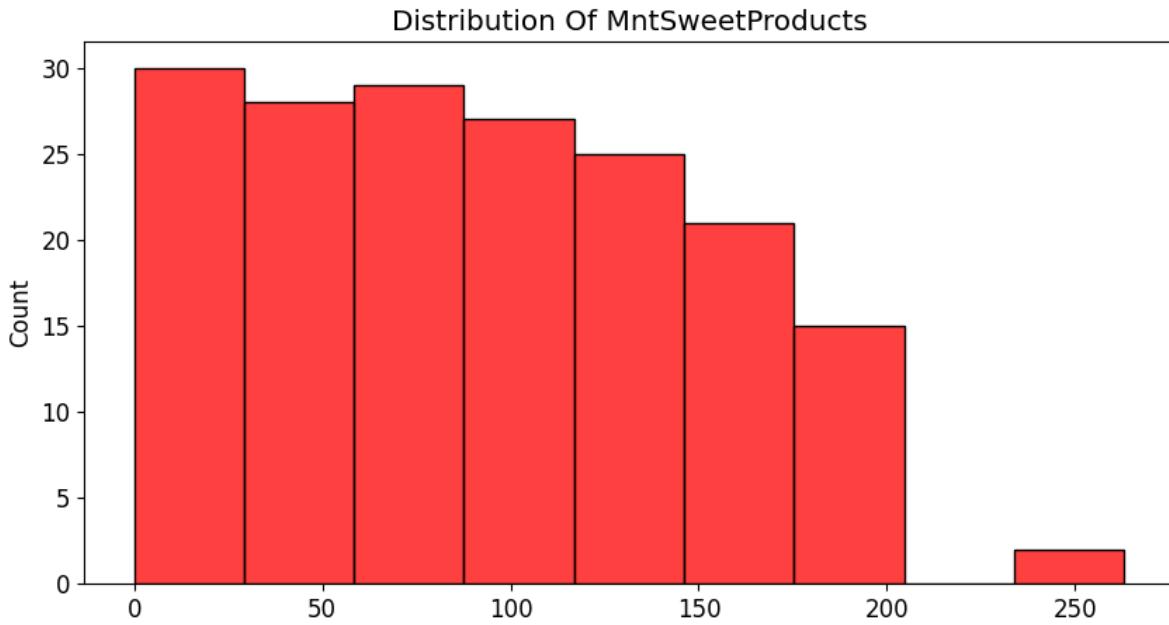
In [45]:

```
1 #Distribution Of MntMeatProducts
2 plt.figure(figsize=(10,5))
3 sns.histplot(df["MntFishProducts"].unique(),color="darkblue")
4 plt.title("Distribution Of MntFishProducts")
5 plt.show()
```



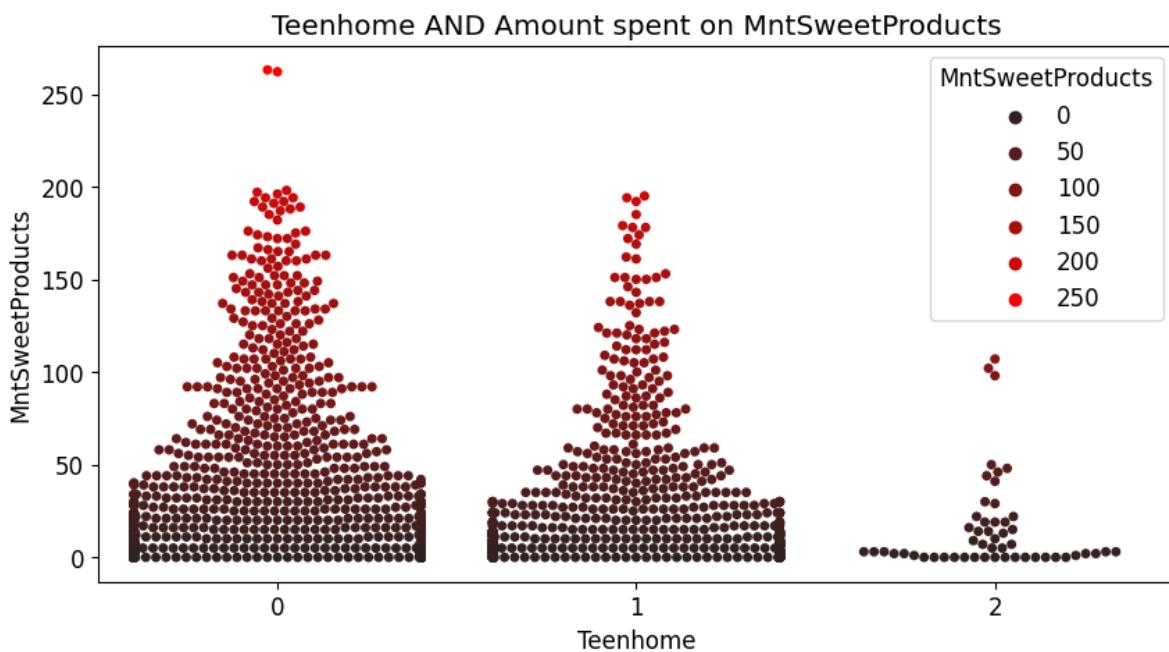
In [46]:

```
1 #Distribution Of MntSweetProducts
2 plt.figure(figsize=(10,5))
3 sns.histplot(df["MntSweetProducts"].unique(),color="red")
4 plt.title("Distribution Of MntSweetProducts")
5 plt.show()
```



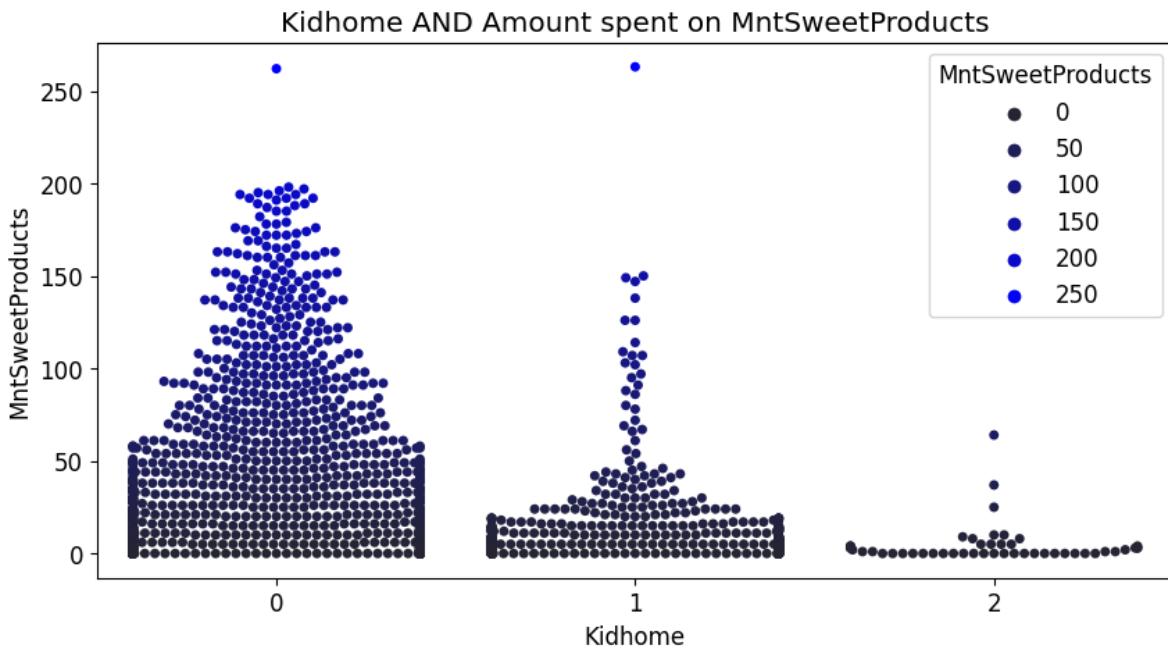
In [47]:

```
1 # Teenhome & MntSweetProducts
2 plt.figure(figsize=(10,5))
3 sns.swarmplot(data=df,x="Teenhome",y="MntSweetProducts",color="red",hue="MntSweetProducts")
4 plt.title("Teenhome AND Amount spent on MntSweetProducts")
5 plt.show()
```



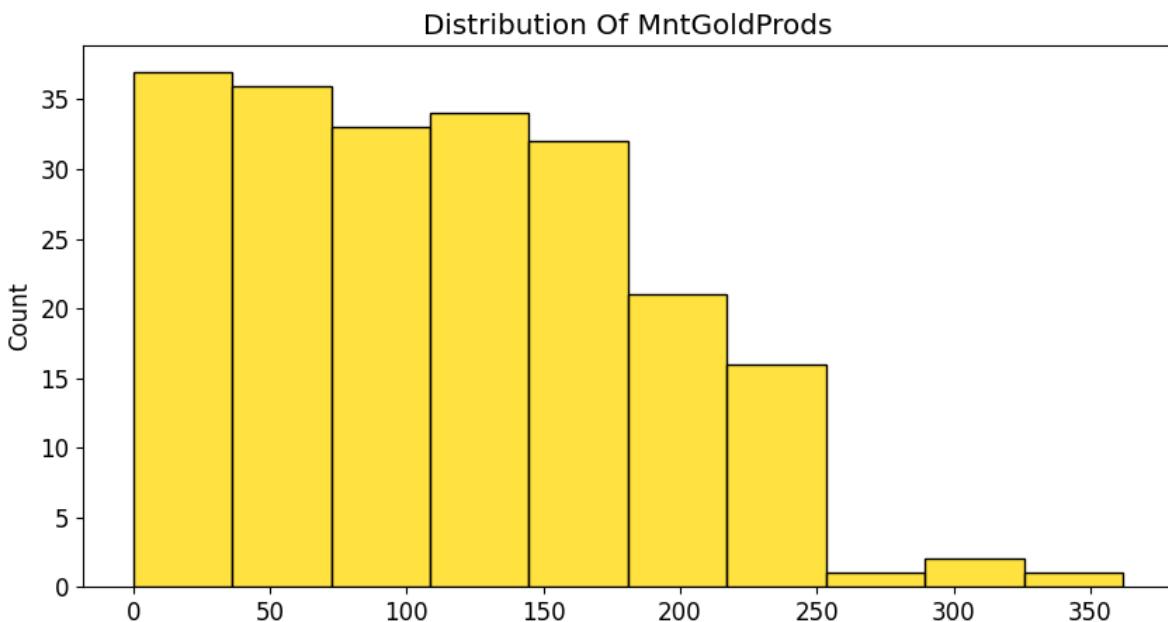
In [48]:

```
1 # Kidhome & MntSweetProducts
2 plt.figure(figsize=(10,5))
3 sns.swarmplot(data=df,x="Kidhome",y="MntSweetProducts",color="blue",hue="MntSweetProducts")
4 plt.title("Kidhome AND Amount spent on MntSweetProducts")
5 plt.show()
```



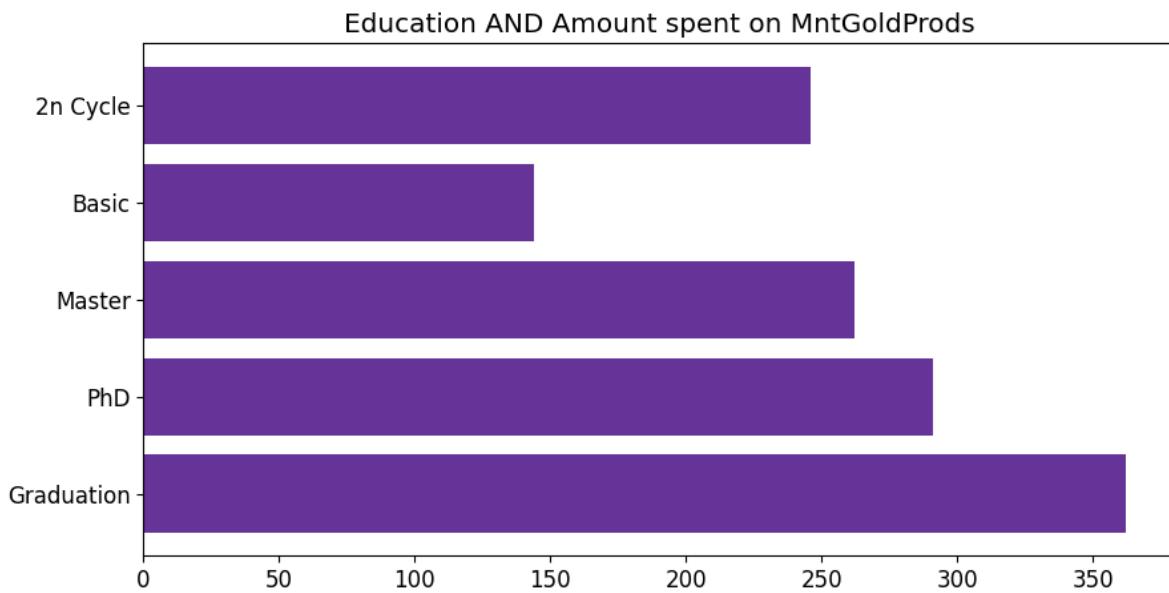
In [49]:

```
1 #Distribution Of MntGoldProds
2 plt.figure(figsize=(10,5))
3 sns.histplot(df["MntGoldProds"].unique(),color="gold")
4 plt.title("Distribution Of MntGoldProds")
5 plt.show()
```



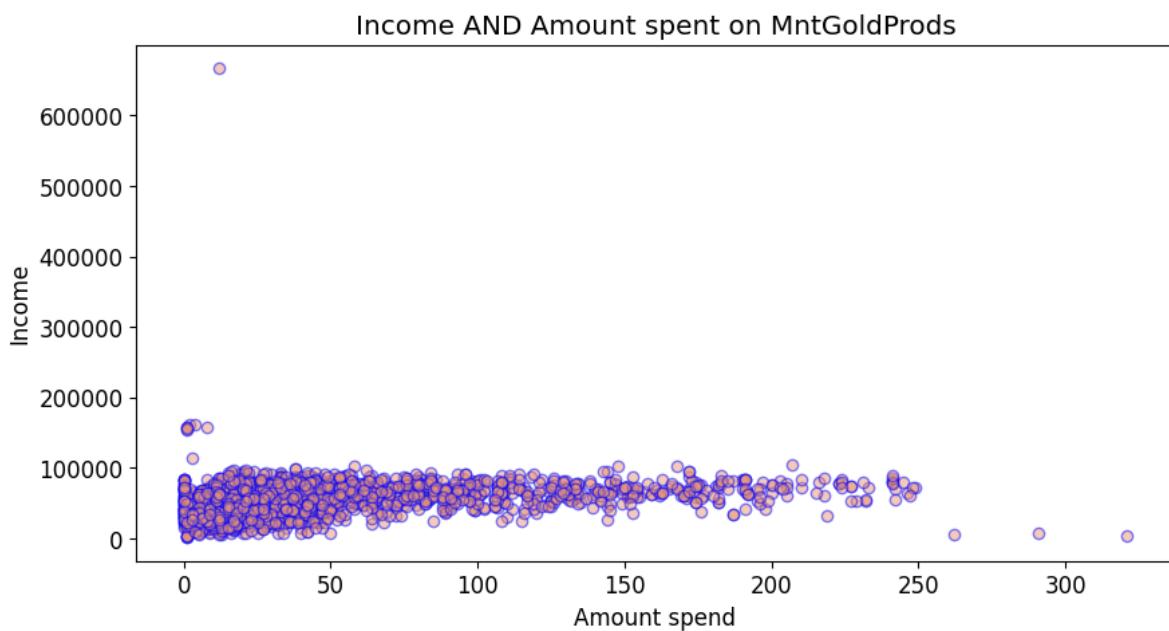
In [50]:

```
1 # MntGoldProds & Education
2 plt.figure(figsize=(10,5))
3 plt.barh(df["Education"],df["MntGoldProds"],color="rebeccapurple")
4 plt.title("Education AND Amount spent on MntGoldProds")
5 plt.show()
```



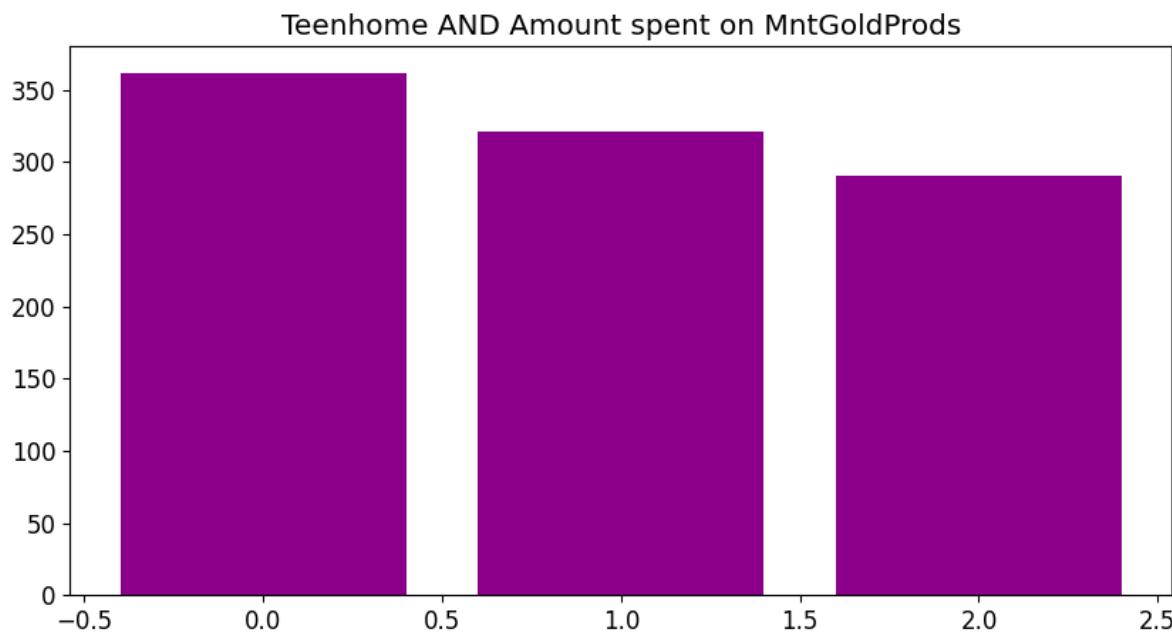
In [51]:

```
1 # Income & MntGoldProds
2 plt.figure(figsize=(10,5))
3 plt.scatter(df["MntGoldProds"],df["Income"],color="darksalmon",alpha=0.5,edgecolors="blue")
4 plt.xlabel("Amount spend")
5 plt.ylabel("Income")
6 plt.title("Income AND Amount spent on MntGoldProds")
7 plt.show()
```



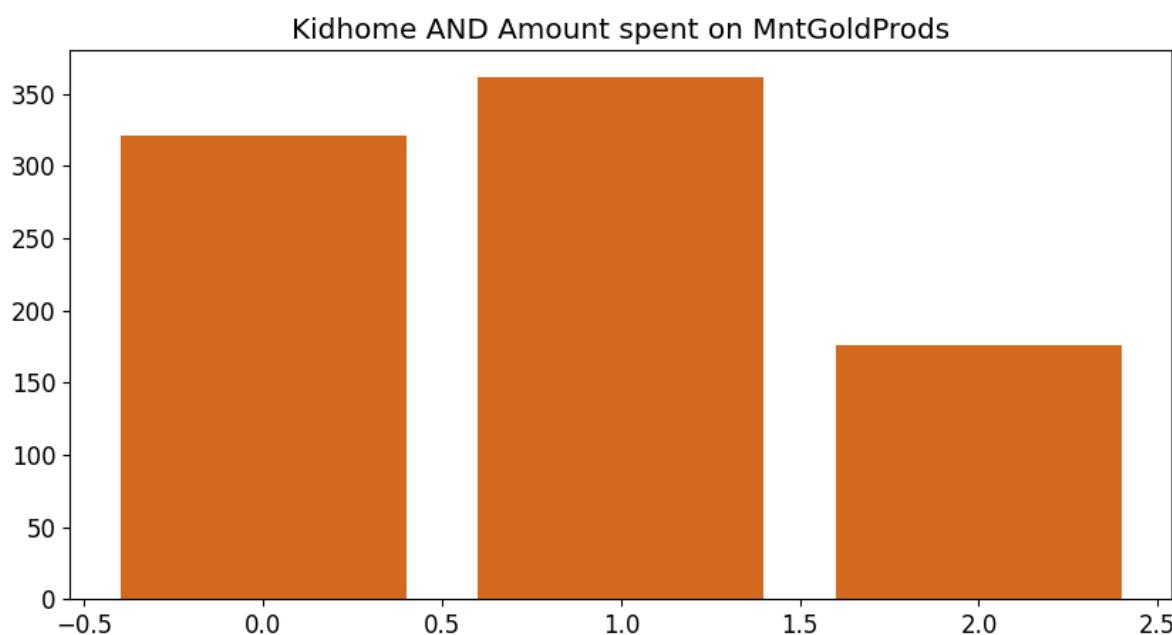
In [52]:

```
1 # Teenhome & MntGoldProds
2 plt.figure(figsize=(10,5))
3 plt.bar(df["Teenhome"],df["MntGoldProds"],color="darkmagenta")
4 plt.title( "Teenhome AND Amount spent on MntGoldProds")
5 plt.show()
```



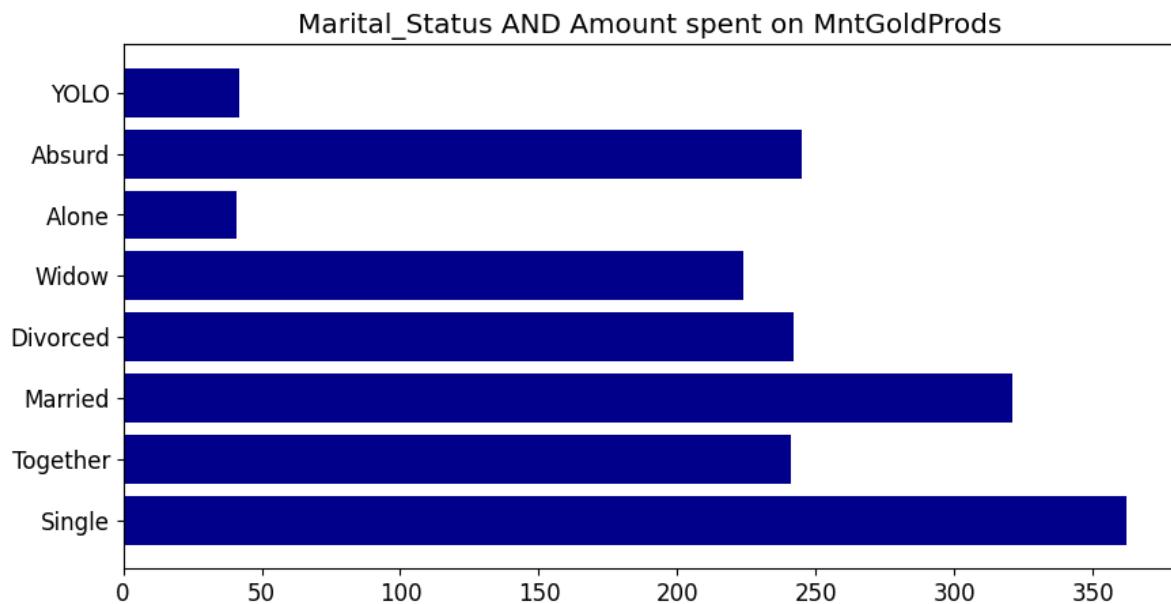
In [53]:

```
1 # Kidhome & MntGoldProds
2 plt.figure(figsize=(10,5))
3 plt.bar(df["Kidhome"],df["MntGoldProds"],color="chocolate")
4 plt.title( "Kidhome AND Amount spent on MntGoldProds")
5 plt.show()
```



In [54]:

```
1 # MntGoldProds & Marital_Status
2 plt.figure(figsize=(10,5))
3 plt.barh(df["Marital_Status"],df["MntGoldProds"],color="darkblue")
4 plt.title("Marital_Status AND Amount spent on MntGoldProds")
5 plt.show()
```



1. People having 0 teens are spending more on gold followed by 1 and 2.

2. People having 2 kids are spending less on gold compared with having 1 kid or none.

3.. Single & married are spending more on gold than compared to others.

In [55]:

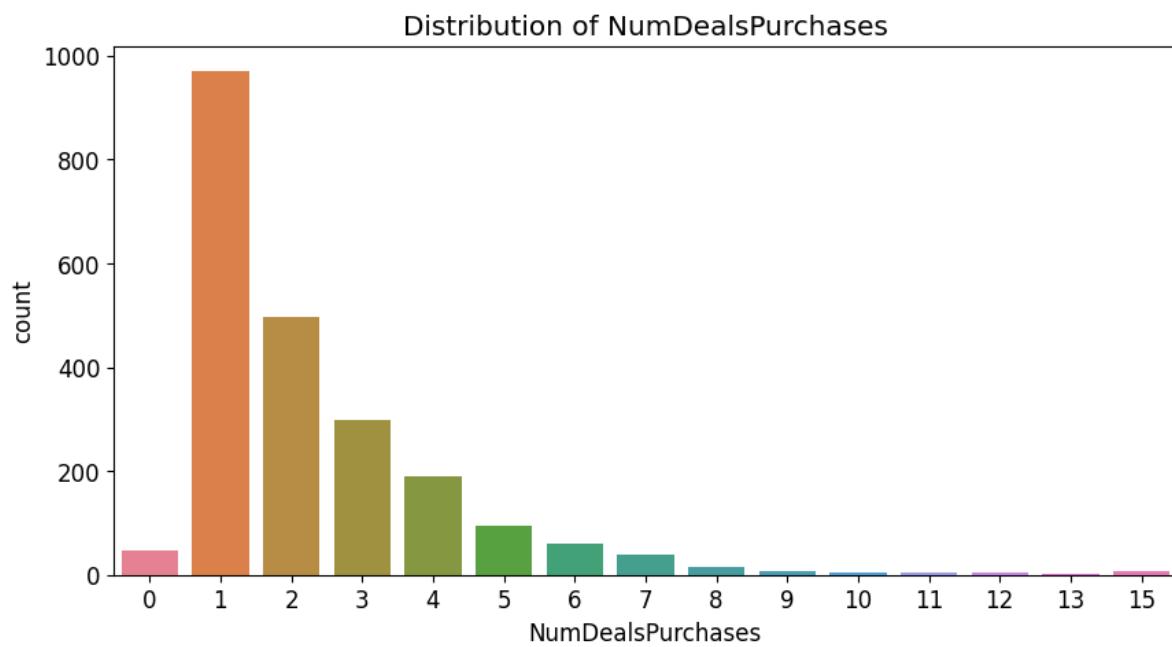
```
1 df["NumDealsPurchases"].values
```

Out[55]:

```
array([3, 2, 1, ..., 1, 2, 3], dtype=int64)
```

In [56]:

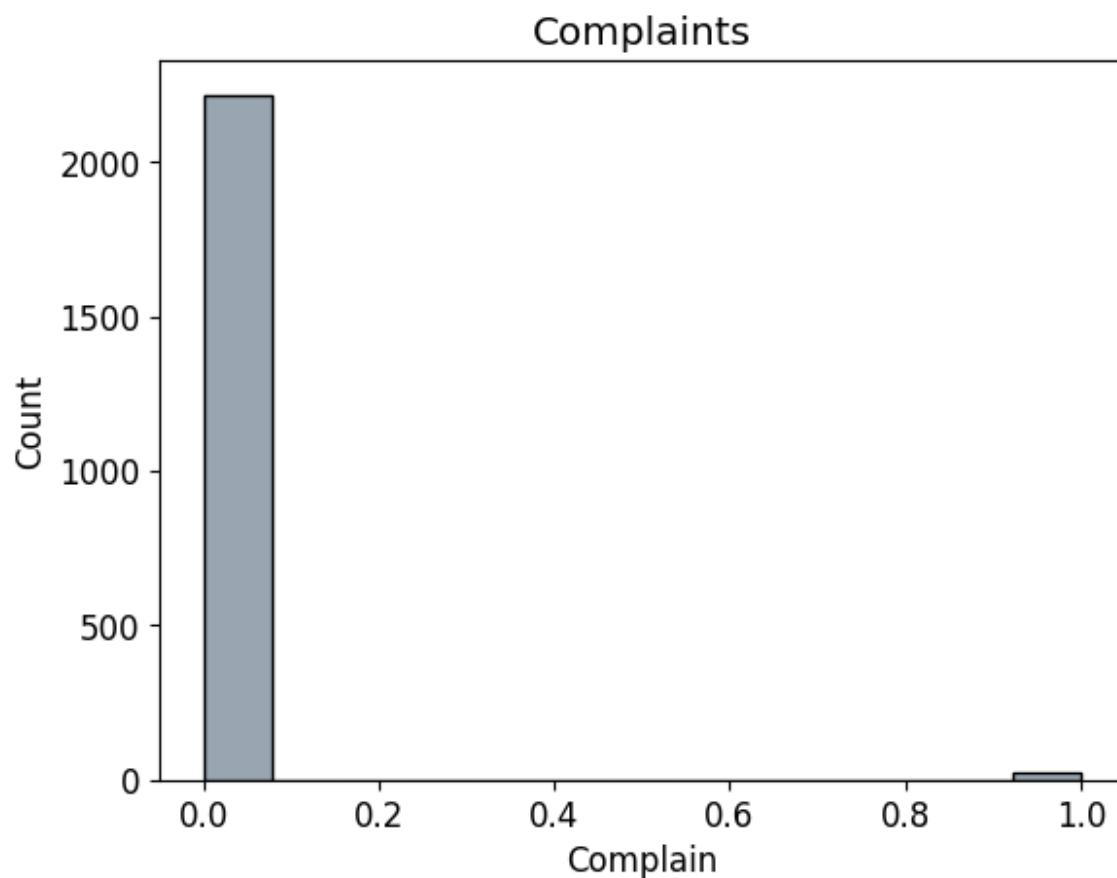
```
1 #NumDealsPurchases
2 plt.figure(figsize=(10,5))
3 sns.countplot(x=df["NumDealsPurchases"], palette="husl")
4 plt.ylabel("count")
5 plt.title("Distribution of NumDealsPurchases")
6 plt.show()
```



1. People Like to purchase when discount is there can be clearly seen that when discount was there many people have brought the items.

In [57]:

```
1 sns.histplot(df["Complain"],color="lightslategrey")
2 plt.title("Complaints")
3 plt.show()
```



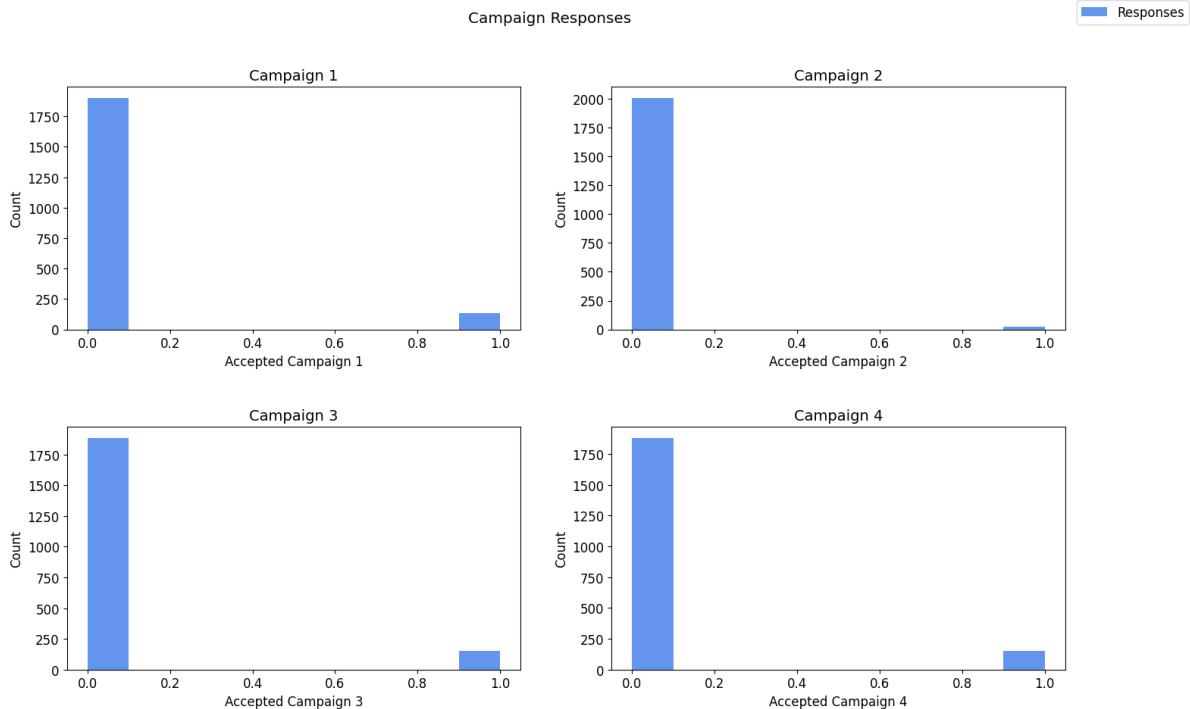
1. There are not much complaints of customer only 21 complaints are there

In [59]:

```
1 #plotting responses of campaigns
```

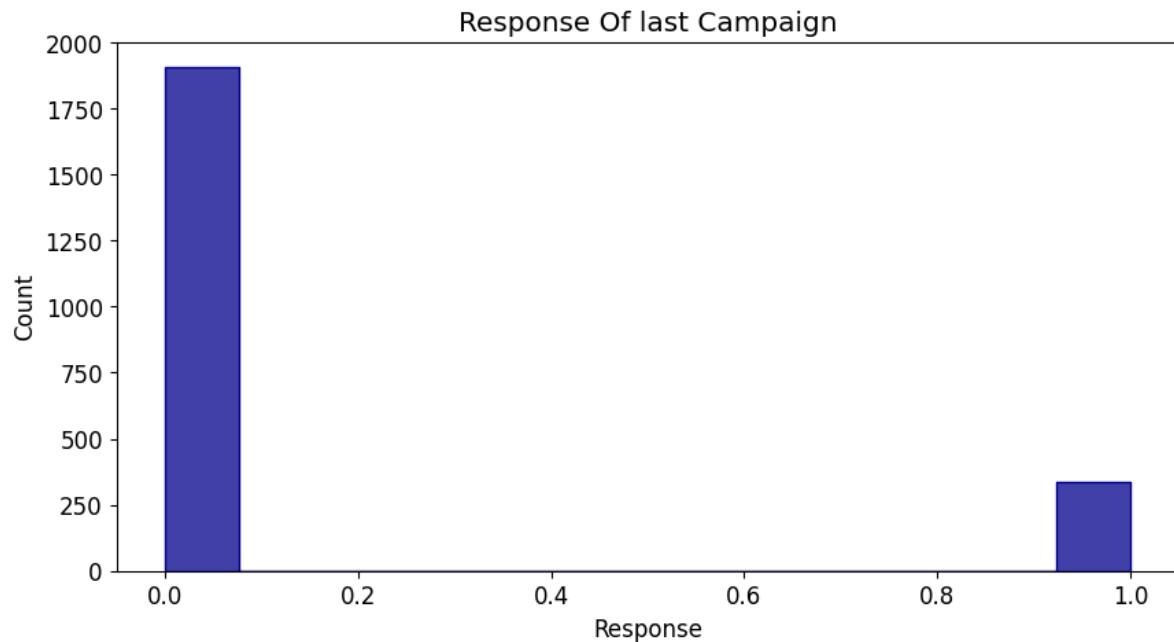
In [134]:

```
1 color = "cornflowerblue"
2
3 fig, axs = plt.subplots(nrows=2, ncols=2, figsize=(17,10))
4
5 axs[0, 0].hist(df["AcceptedCmp1"], color=color)
6 axs[0, 0].set_title("Campaign 1", fontsize=14)
7 axs[0, 0].set_xlabel("Accepted Campaign 1", fontsize=12)
8 axs[0, 0].set_ylabel("Count", fontsize=12)
9
10 axs[0, 1].hist(df["AcceptedCmp2"], color=color)
11 axs[0, 1].set_title("Campaign 2", fontsize=14)
12 axs[0, 1].set_xlabel("Accepted Campaign 2", fontsize=12)
13 axs[0, 1].set_ylabel("Count", fontsize=12)
14
15 axs[1, 0].hist(df["AcceptedCmp3"], color=color)
16 axs[1, 0].set_title("Campaign 3", fontsize=14)
17 axs[1, 0].set_xlabel("Accepted Campaign 3", fontsize=12)
18 axs[1, 0].set_ylabel("Count", fontsize=12)
19
20 axs[1, 1].hist(df["AcceptedCmp4"], color=color)
21 axs[1, 1].set_title("Campaign 4", fontsize=14)
22 axs[1, 1].set_xlabel("Accepted Campaign 4", fontsize=12)
23 axs[1, 1].set_ylabel("Count", fontsize=12)
24
25 fig.subplots_adjust(hspace=0.4, wspace=0.2)
26 fig.legend(labels=['Responses'], loc='upper right', fontsize=12)
27 fig.suptitle("Campaign Responses")
28 plt.show()
```



In [61]:

```
1 #Plotting response i.e of last campaign
2 plt.figure(figsize=(10,5))
3 sns.histplot(df["Response"],color="darkblue",element="step")
4 plt.title("Response Of last Campaign")
5 plt.show()
```



1. It can be clearly seen people doesn't accept the offer in 2nd campaign ,the ratio is almost same in other 3 campaigns.

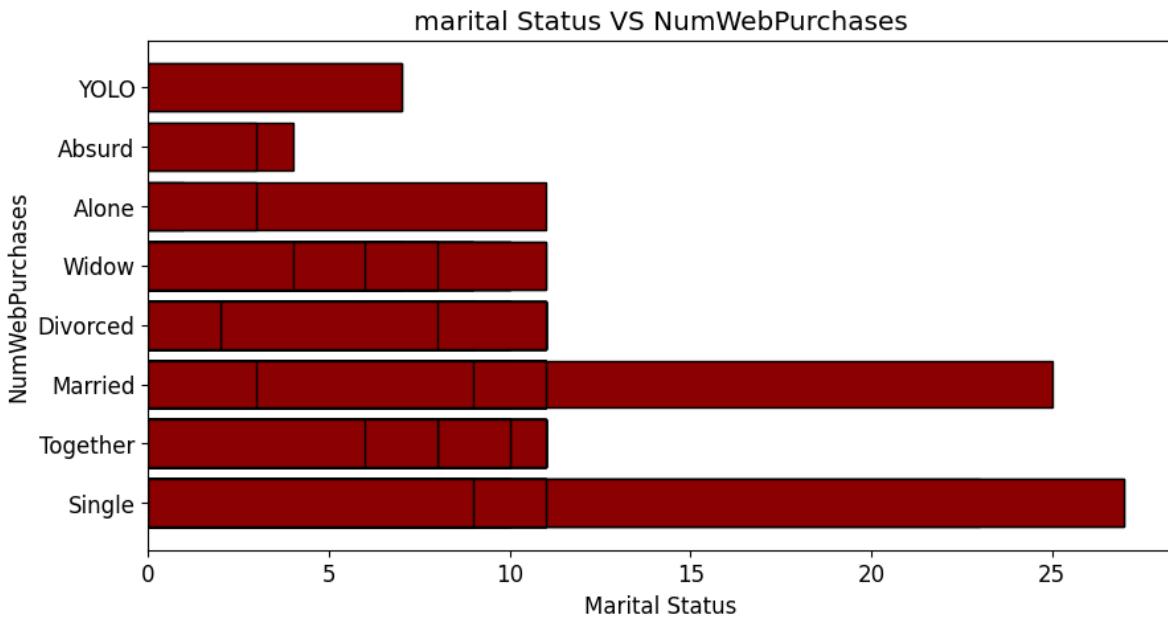
2. Also majority of people have not accepted offer in last campaign i.e.they are least affected by campaigns.

In [62]:

```
1 #NumWebPurchases, NumStorePurchases & NumCatalogPurchases
```

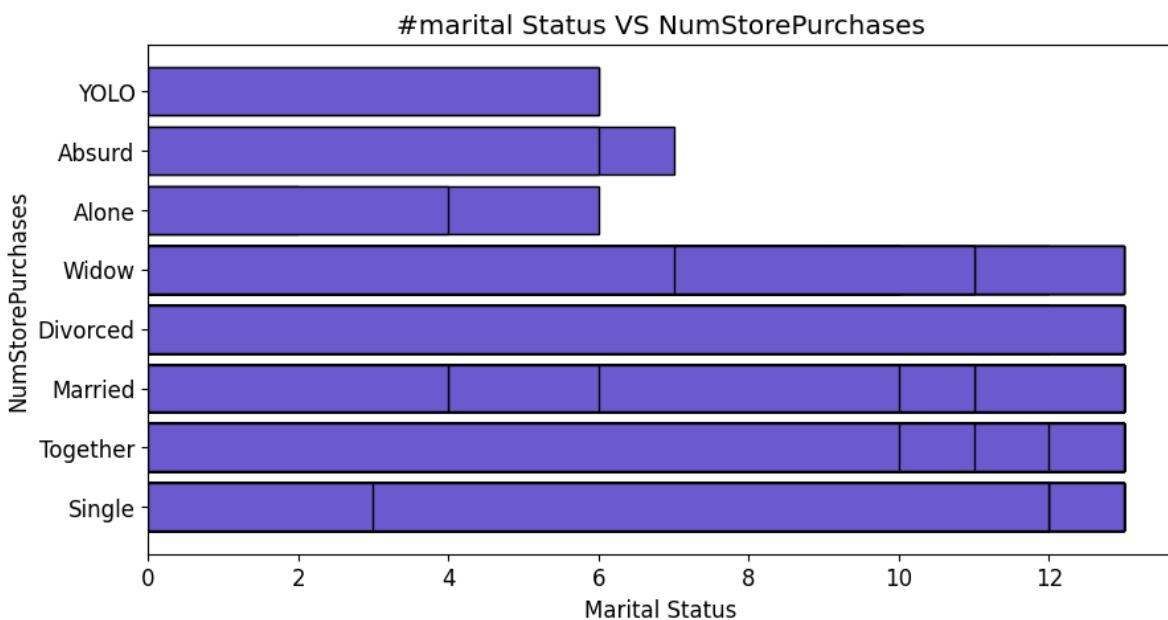
In [63]:

```
1 #marital Status VS NumWebPurchases
2 plt.figure(figsize=(10,5))
3 plt.barh(df["Marital_Status"],df[ "NumWebPurchases"],color="darkred",edgecolor="black")
4 plt.xlabel("Marital Status")
5 plt.ylabel("NumWebPurchases")
6 plt.title("marital Status VS NumWebPurchases")
7 plt.show()
```



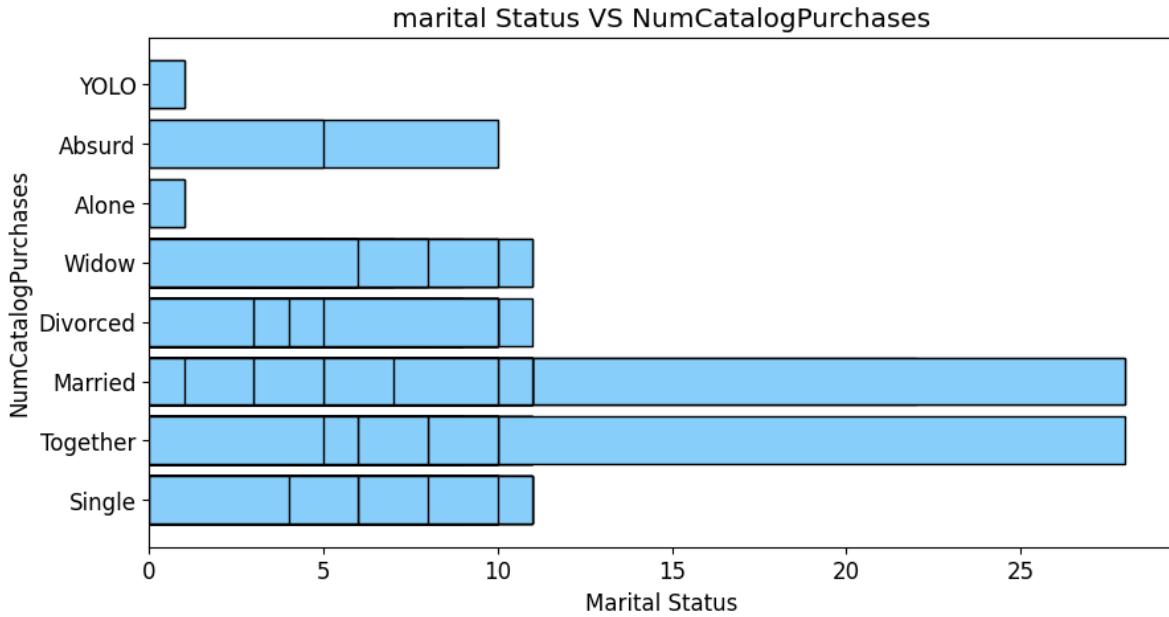
In [64]:

```
1 #marital Status VS NumStorePurchases
2 plt.figure(figsize=(10,5))
3 plt.barh(df["Marital_Status"],df[ "NumStorePurchases"],color="slateblue",edgecolor="black")
4 plt.xlabel("Marital Status")
5 plt.ylabel("NumStorePurchases")
6 plt.title("#marital Status VS NumStorePurchases")
7 plt.show()
```



In [65]:

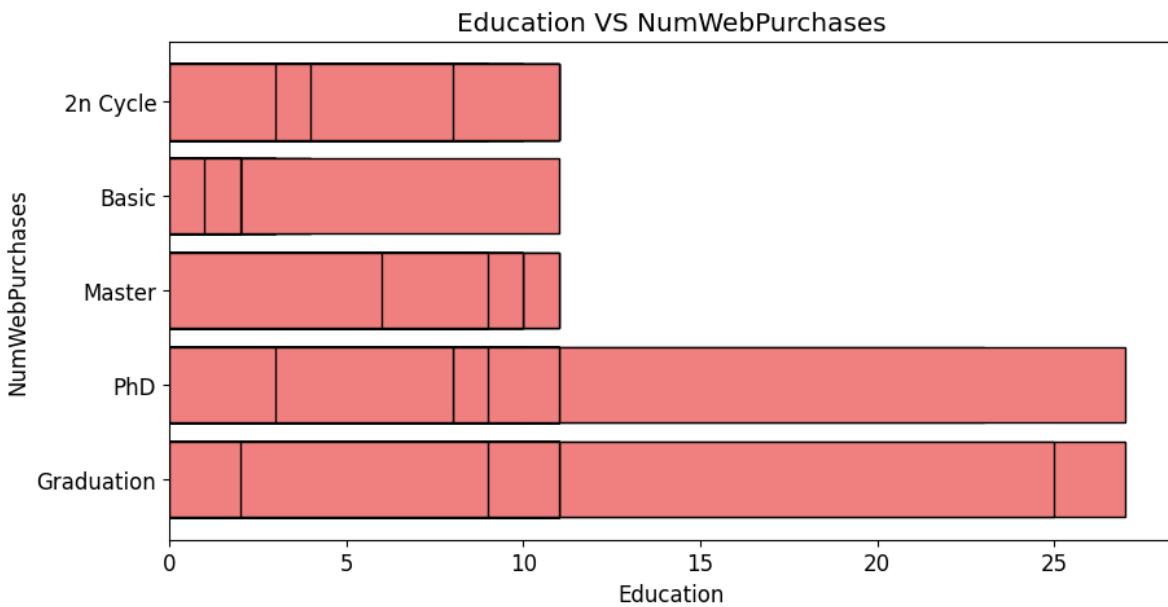
```
1 #marital Status VS NumCatalogPurchases
2 plt.figure(figsize=(10,5))
3 plt.barh(df["Marital_Status"],df["NumCatalogPurchases"],color="lightskyblue",edgecolor="black")
4 plt.xlabel("Marital Status")
5 plt.ylabel("NumCatalogPurchases")
6 plt.title("marital Status VS NumCatalogPurchases")
7 plt.show()
```



1. Web purchases are mostly done in majority by single followed by married persons.
2. Store purchases are done equally by people on marital status those are alone,absurd,YOLO are not in large number anyways.
3. Catalog purchases are done mostly by married & together.

In [66]:

```
1 #Education VS NumWebPurchases
2 plt.figure(figsize=(10,5))
3 plt.barh(df["Education"],df["NumWebPurchases"],color="lightcoral",edgecolor="black")
4 plt.xlabel("Education")
5 plt.ylabel("NumWebPurchases")
6 plt.title("Education VS NumWebPurchases")
7 plt.show()
```



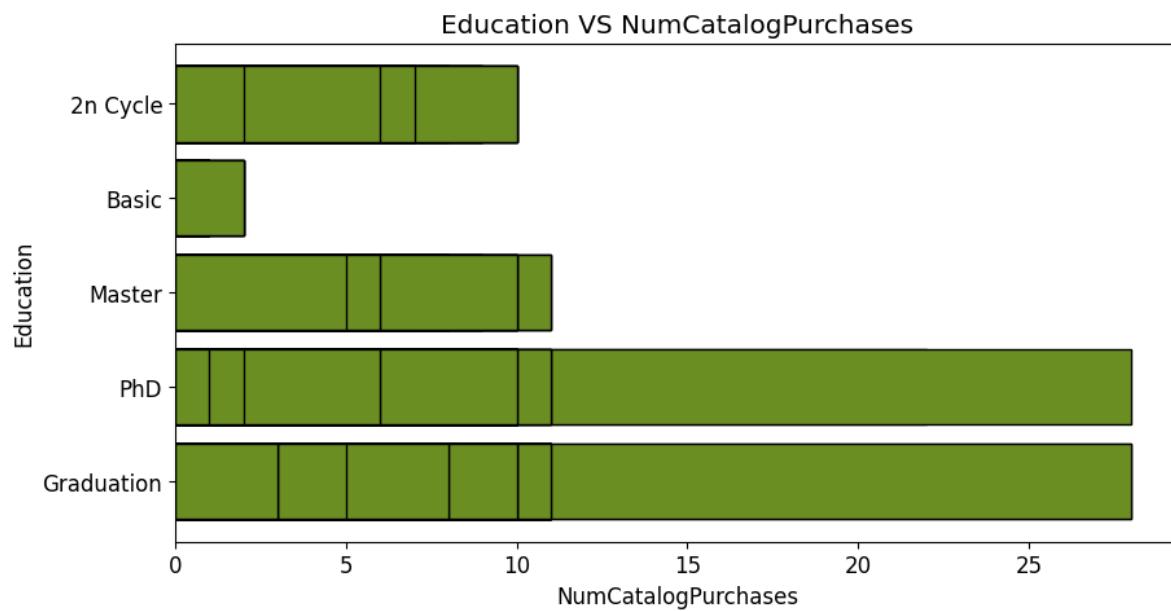
In [67]:

```
1 #Education VS NumStorePurchases
2 plt.figure(figsize=(10,5))
3 plt.barh(df["Education"],df["NumStorePurchases"],color="mediumorchid",edgecolor="black")
4 plt.xlabel("Number of Store Purchases")
5 plt.ylabel("Education")
6 plt.title("#Education VS NumStorePurchases")
7 plt.show()
```



In [68]:

```
1 #Education VS NumCatalogPurchases
2 plt.figure(figsize=(10,5))
3 plt.barh(df["Education"], df["NumCatalogPurchases"], color="olivedrab", edgecolor="black")
4 plt.xlabel("NumCatalogPurchases")
5 plt.ylabel("Education")
6 plt.title("Education VS NumCatalogPurchases")
7 plt.show()
```



1. People with basic education are least in catalog purchase.

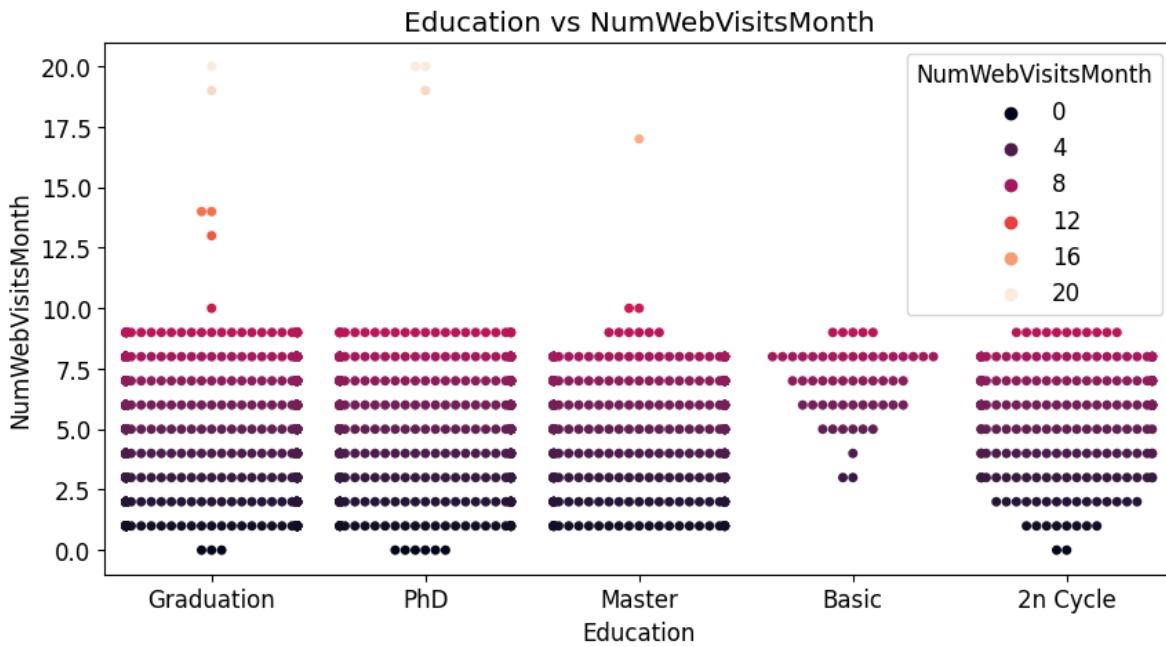
2. People with all sort of education visit store but the no.of people with basic are less.

In [69]:

```
1 #NumWebVisitsMonth
```

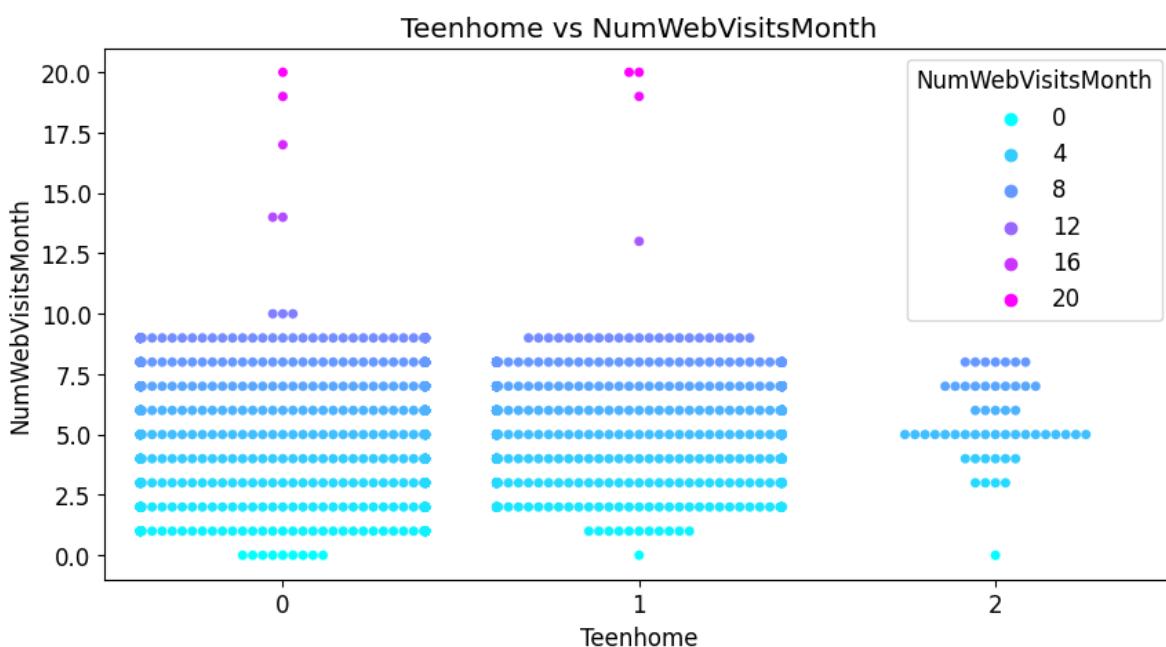
In [70]:

```
1 #Education vs NumWebVisitsMonth
2 plt.figure(figsize=(10,5))
3 sns.swarmplot(x=df["Education"],y=df["NumWebVisitsMonth"],hue=df["NumWebVisitsMonth"],pal
4 plt.title("Education vs NumWebVisitsMonth")
5 plt.show()
```



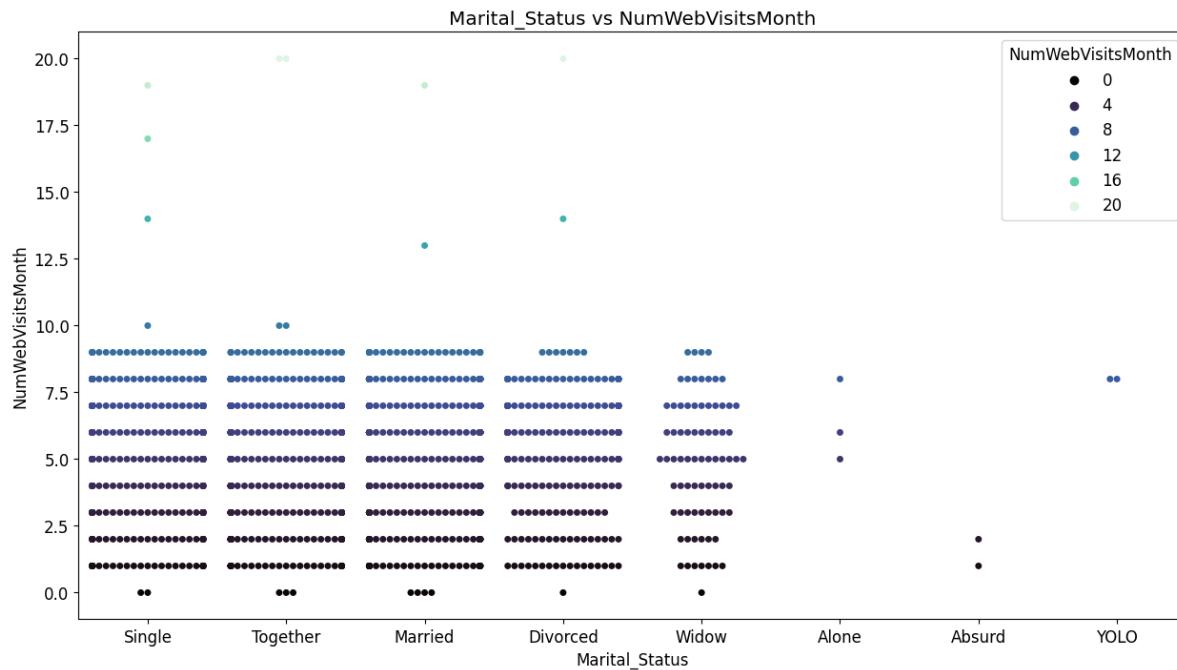
In [71]:

```
1 #Teenhome vs NumWebVisitsMonth
2
3 plt.figure(figsize=(10,5))
4 sns.swarmplot(x=df["Teenhome"],y=df["NumWebVisitsMonth"],hue=df["NumWebVisitsMonth"],pal
5 plt.title("Teenhome vs NumWebVisitsMonth")
6 plt.show()
```



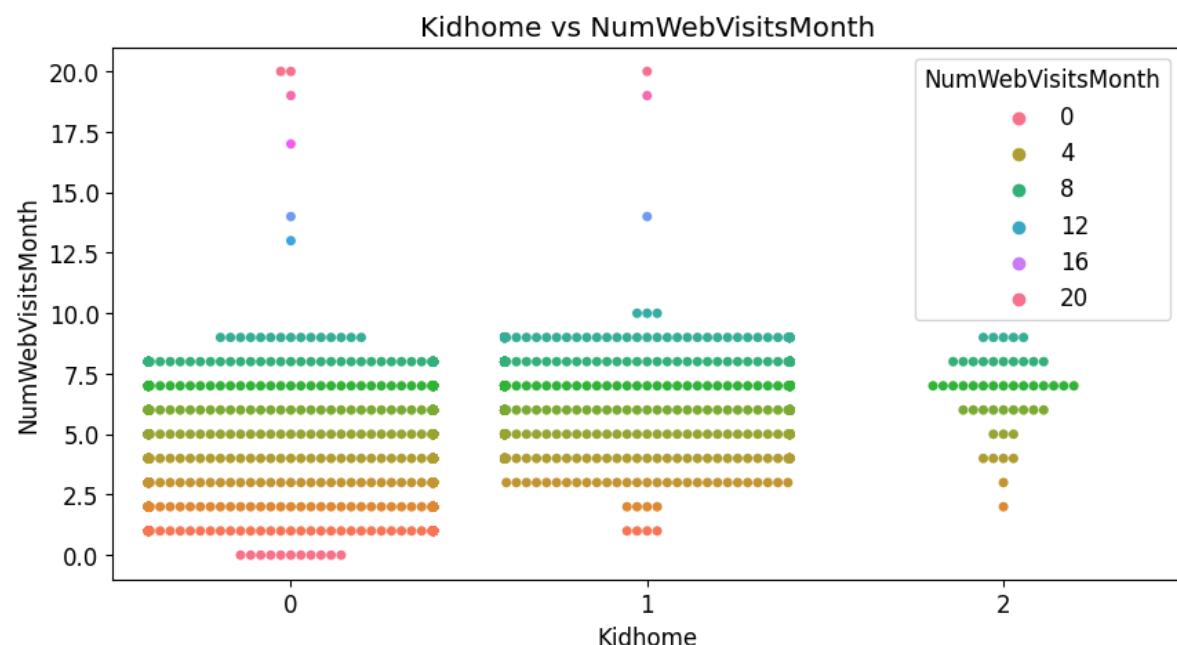
In [72]:

```
1 #Marital_Status vs NumWebVisitsMonth
2 plt.figure(figsize=(15,8))
3
4 sns.swarmplot(x=df["Marital_Status"],y=df["NumWebVisitsMonth"],hue=df["NumWebVisitsMonth"])
5 plt.title("Marital_Status vs NumWebVisitsMonth")
6 plt.show()
```



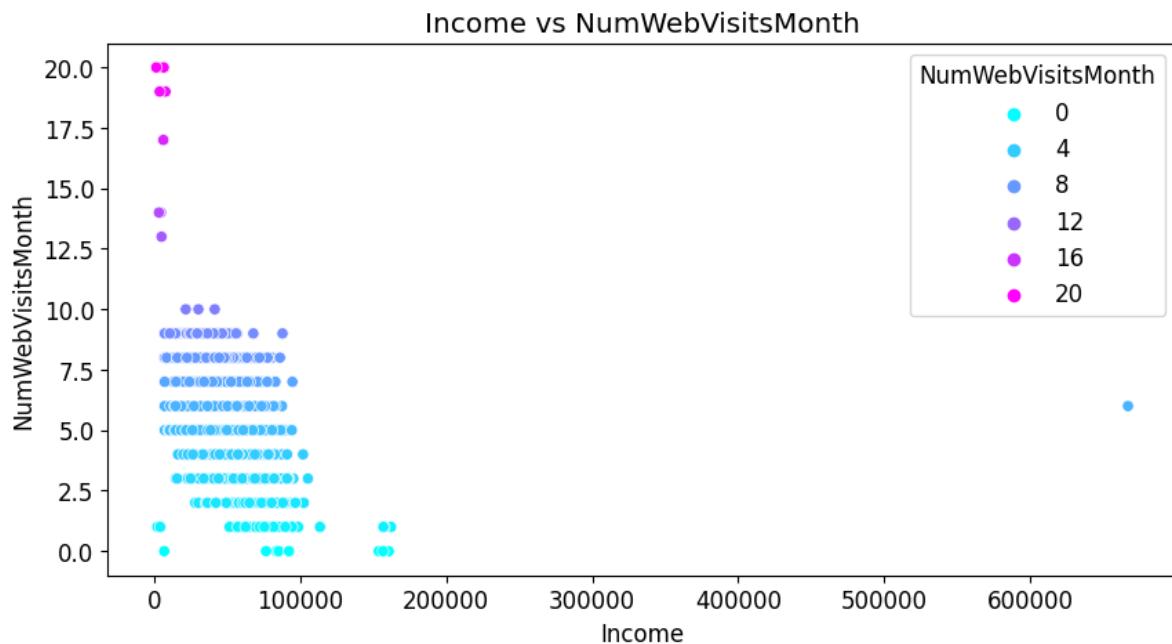
In [73]:

```
1 # Kidhome vs NumWebVisitsMonth
2 plt.figure(figsize=(10,5))
3
4 sns.swarmplot(data=df,x="Kidhome",y="NumWebVisitsMonth",hue=df["NumWebVisitsMonth"],palette="viridis")
5 plt.title("Kidhome vs NumWebVisitsMonth")
6 plt.show()
```



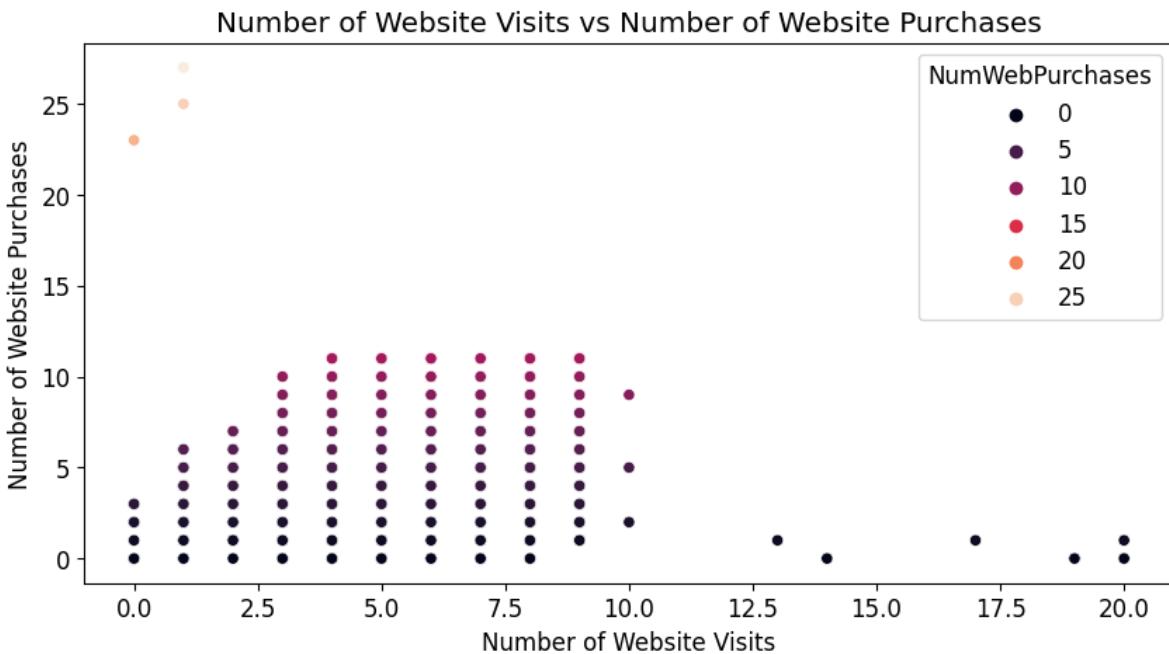
In [74]:

```
1 # Income vs NumWebVisitsMonth
2 plt.figure(figsize=(10,5))
3
4 sns.scatterplot(data=df,x="Income",y="NumWebVisitsMonth",palette="cool",hue=df[ "NumWebVisitsMonth"])
5 plt.title("Income vs NumWebVisitsMonth")
6 plt.show()
```



In [75]:

```
1 #Number of Website Visits vs Number of Website Purchases
2 plt.figure(figsize=(10,5))
3 sns.scatterplot(x=df["NumWebVisitsMonth"], y=df["NumWebPurchases"], hue=df["NumWebPurchase"])
4 plt.title("Number of Website Visits vs Number of Website Purchases")
5 plt.xlabel("Number of Website Visits")
6 plt.ylabel("Number of Website Purchases")
7 plt.show()
```



Conclusions/Observations:

Birth Year

1. Dataset is of people born between 1940-2000.

2. 2 or 3 variables of data are from year 1900 or earlier it can be result of wrong entry or an outlier.

Education:

1. Education is divided into 5 types:

- i.Graduation ---16.5%
- ii.Basic ---21.7%
- iii.Master ---9.1%
- iv.PhD ---2.4%
- v.2n Cycle ---50.3%

Marital Status:

1. Marital Status is divided into 8 types:

- i.Single
 - ii.Togther
 - iii.Married
 - iv.Divorced
 - v.Widow
 - vi.Alone
 - vii.Absurd
-

Income:

- 1.Income of Basic education is less comapred to others.
- 2.Highest Income generated is from Graduation but it is too high,can be wrong entry or outlier as well.
- 3.Birth year wise income is distributed equally but some group of people from birth year 1970-1980 are earning more compared to others.

Kids & Teens:

- 1.50.7% of people do not have kids,40.1% people have 1 child whereas 2.1% people have 2chilids.
- 2.51.7% of people do not have teens at their home,46.0% of people have 1 teen wheras 2.3% of people have 2 teens at home.

Amount Spend:

- 1.Maximum/Highest amount is spend on wine.
- 2.People with basic education are not spending any amount on wines.
- 3.People with basic education are not spending much on fruits & meat as comapared to others
- 4.All types of people differentiated on basis of martial status are spending nearby same on fruits& meat products.

Num Deals Purchase (Discount)

- 1.People Like to purchase when discount is there can be clearly seen that when discount was there many people have brought the items.

Registrations:

Most of the registrations are from year 2013 in dataset provided of 2 years.

Complaints:

There are not much complaints of customer only 21 complaints are there

Amount Spent On Gold:

- 1.People having 0 teens are spending more on gold followed by 1 and 2.
- 2.People having 2 kids are spending less on gold comapred with having 1 kid or none.

3..Single & married are spending more on gold than compared to others.

Offer Campaigns:

1.It can be clearly seen people doesn't accept the offer in 2nd campaign ,the ratio is almost same in other 3 campaigns.

2.Also majority of people have not accepted offer in last campaign i.e.they are least affected by campaigns.

Purchases:

1.Web purchases are mostly done in majority by single followed by married persons.

2.Store purchases are done equally by people on marital status those are alone,absurd,YOLO are not in large number anyways.

3.Catalog purchases are done mostly by married & together.

4.People with basic education are least in catalog purchase.

Store & Website Visit

1.All sort of people visit website but no.of people with basic education are less visiting comparatively to others.

2.Was trying to find if no.of visits depends on Teenhome & Kidhome but that was not the case,No sort of relation found between them.

3.Marital Status is also a not a consideration for visiting website people with all sort of marital status visit it,count of widows and others is less because we have very small no.of data regarding to them in the dataset.

4.People with all sort of education visit store but the no.of people with basic are less.

5.People with all sort of income are visiting website frequently no pattern identified on one class.

6.Most of the purchases are done in 1-9 visits.

DATA PREPROCESSING

In [76]:

```
1 #Dropping Date,Month & Year Columns Separated for EDA purpose
2 #Dropping ID As well because of high cardinality
3 #Dropping Dt_Customer as its in str format & even its not of that important
```

In [77]:

```
1 df.drop(["ID","Date_enroll","Month_enroll","Year_enroll","Dt_Customer"],axis=1,inplace=True)
```

In [78]:

```
1 #Checking for duplicates
2 df.duplicated().sum()
```

Out[78]:

182

In [79]:

```
1 #Dropping Duplicates
2 df.drop_duplicates(inplace=True)
```

In [80]:

```
1 df.duplicated().sum()
```

Out[80]:

0

In [81]:

```
1 #ALL duplicates values have been removed
```

In [82]:

```
1 #Checking For null values
2 df.isna().sum()
```

Out[82]:

Year_Birth	0
Education	0
Marital_Status	0
Income	24
Kidhome	0
Teenhome	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

dtype: int64

In [83]:

```
1 #As Income have 24 null values we are going to drop them
2 df.dropna(inplace=True)
```

In [84]:

```
1 df.isna().sum()
```

Out[84]:

```
Year_Birth          0
Education          0
Marital_Status     0
Income             0
Kidhome            0
Teenhome            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
dtype: int64
```

In [85]:

```
1 df.describe()
```

Out[85]:

	Year_Birth	Income	Kidhome	Teenhome	Recency	MntWines	MntFruit
count	2034.000000	2034.000000	2034.000000	2034.000000	2034.000000	2034.000000	2034.000000
mean	1968.802852	52357.791544	0.442970	0.508358	48.844641	305.180924	26.29252
std	11.975420	25526.956988	0.535914	0.546130	28.983678	337.323274	39.77288
min	1893.000000	1730.000000	0.000000	0.000000	0.000000	0.000000	0.00000
25%	1959.000000	35528.250000	0.000000	0.000000	24.000000	23.000000	2.00000
50%	1970.000000	51533.000000	0.000000	0.000000	49.000000	175.500000	8.00000
75%	1977.000000	68480.750000	1.000000	1.000000	74.000000	505.000000	33.00000
max	1996.000000	666666.000000	2.000000	2.000000	99.000000	1493.000000	199.00000

8 rows × 25 columns

In [86]:

```
1 df.head()
```

Out[86]:

	Year_Birth	Education	Marital_Status	Income	Kidhome	Teenhome	Recency	MntWines	MntFru
0	1957	Graduation	Single	58138.0	0	0	58	635	
1	1954	Graduation	Single	46344.0	1	1	38	11	
2	1965	Graduation	Together	71613.0	0	0	26	426	
3	1984	Graduation	Together	26646.0	1	0	26	11	
4	1981	PhD	Married	58293.0	1	0	94	173	

5 rows × 27 columns

In [87]:

```
1 #Separating categorical & numeric Data
2 df_Numeric=df.select_dtypes(include="int")
3 df_categorical=df.select_dtypes(include="object")
```

In [88]:

```
1 df_categorical.head()
```

Out[88]:

	Education	Marital_Status
0	Graduation	Single
1	Graduation	Single
2	Graduation	Together
3	Graduation	Together
4	PhD	Married

In [89]:

```
1 #Encoding
```

In [90]:

```
1 enc_categorical=pd.get_dummies(df_categorical)
2 enc_categorical.head()
```

Out[90]:

	Education_2n Cycle	Education_Basic	Education_Graduation	Education_Master	Education_PhD	Marital
0	0	0	1	0	0	0
1	0	0	1	0	0	0
2	0	0	1	0	0	0
3	0	0	1	0	0	0
4	0	0	0	0	0	1

In [91]:

```
1 #Concatenating Numeric & Encoded Categorical Data
```

In [92]:

```
1 final_df=pd.concat([enc_categorical,df_Numeric],axis=1)
```

In [93]:

```
1 final_df.head()
```

Out[93]:

	Education_2n Cycle	Education_Basic	Education_Graduation	Education_Master	Education_PhD	Marital
0	0	0	1	0	0	0
1	0	0	1	0	0	0
2	0	0	1	0	0	0
3	0	0	1	0	0	0
4	0	0	0	0	0	1

5 rows × 37 columns



In [94]:

```
1 final_df.columns
```

Out[94]:

```
Index(['Education_2n Cycle', 'Education_Basic', 'Education_Graduation',
       'Education_Master', 'Education_PhD', 'Marital_Status_Absurd',
       'Marital_Status_Alone', 'Marital_Status_Divorced',
       'Marital_Status_Married', 'Marital_Status_Single',
       'Marital_Status_Together', 'Marital_Status_Widow',
       'Marital_Status_YOLO', 'Year_Birth', 'Kidhome', 'Teenhome', 'Recency',
       'MntWines', 'MntFruits', 'MntMeatProducts', 'MntFishProducts',
       'MntSweetProducts', 'MntGoldProds', 'NumDealsPurchases',
       'NumWebPurchases', 'NumCatalogPurchases', 'NumStorePurchases',
       'NumWebVisitsMonth', 'AcceptedCmp3', 'AcceptedCmp4', 'AcceptedCmp5',
       'AcceptedCmp1', 'AcceptedCmp2', 'Complain', 'Z_CostContact',
       'Z_Revenue', 'Response'],
      dtype='object')
```

PCA

In [137]:

```
1 from sklearn.preprocessing import StandardScaler
2 sc=StandardScaler()
3 scaled_df=sc.fit_transform(final_df)
4 df2=pd.DataFrame(scaled_df,columns=final_df.columns)
```

In [138]:

```
1 scaled_df.shape
```

Out[138]:

(2034, 37)

In [139]:

```
1 df2.head()
```

Out[139]:

	Education_2n Cycle	Education_Basic	Education_Graduation	Education_Master	Education_PhD	Marital
0	-0.316313	-0.157115	0.998035	-0.444837	-0.529198	
1	-0.316313	-0.157115	0.998035	-0.444837	-0.529198	
2	-0.316313	-0.157115	0.998035	-0.444837	-0.529198	
3	-0.316313	-0.157115	0.998035	-0.444837	-0.529198	
4	-0.316313	-0.157115	-1.001969	-0.444837	1.889652	

5 rows × 37 columns



In [140]:

```
1 from sklearn.decomposition import PCA
```

In [141]:

```
1 pca=PCA(n_components=37)
2 pca.fit_transform(df2)
```

Out[141]:

```
array([[ 3.92126114e+00, -9.61448340e-01, -4.18362311e-01, ...,
       6.81726420e-15,  4.70969941e-16,  1.84690340e-17],
       [-2.36140150e+00, -5.12049484e-01, -3.04982935e-01, ...,
        3.94533695e-15,  1.38987152e-17, -3.82208087e-17],
       [ 1.50092229e+00, -4.27692395e-01, -1.38093928e+00, ...,
       -1.53724166e-14,  7.59451335e-17,  3.11820417e-17],
       ...,
       [ 1.50103128e+00,  2.68854853e-01,  4.08566507e-01, ...,
       -3.13574979e-16,  1.56782122e-17, -9.09263745e-18],
       [ 1.56188422e+00,  1.38250713e+00, -1.05948585e+00, ...,
        7.86248061e-16,  1.99674898e-17, -6.51998320e-18],
       [-1.78295060e+00,  1.65644879e+00,  1.32880184e+00, ...,
       1.41045263e-15, -4.02049749e-18]])
```

In [142]:

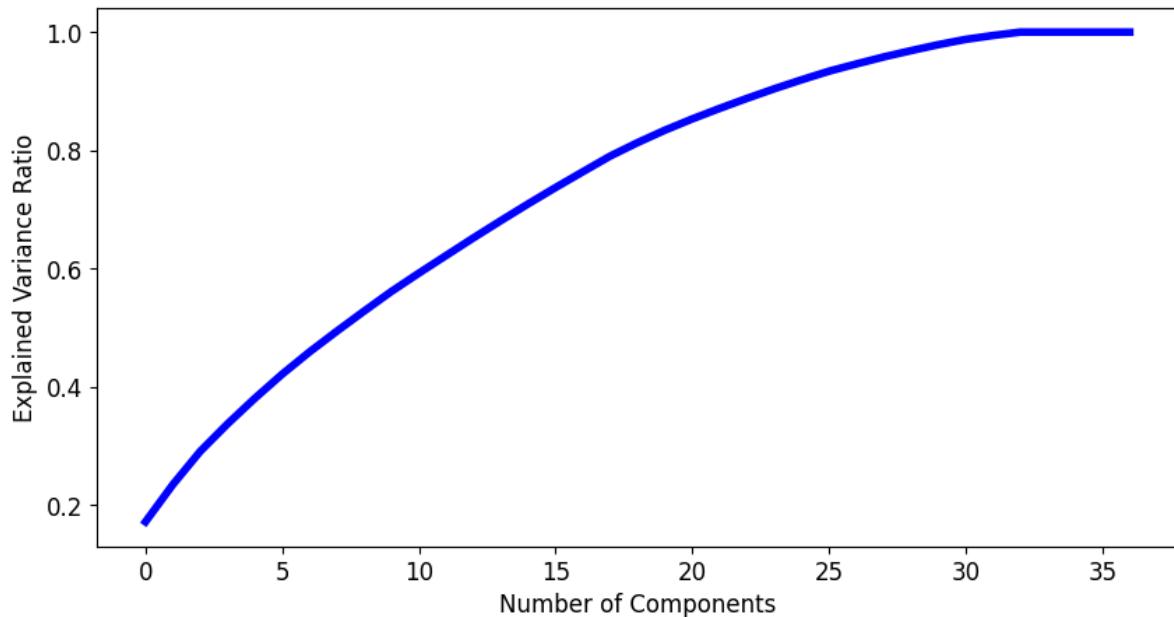
```
1 #explained variance ratio
2 pca.explained_variance_ratio_
```

Out[142]:

```
array([1.72192328e-01, 6.28764599e-02, 5.61055137e-02, 4.64831620e-02,
       4.32938627e-02, 4.07701572e-02, 3.76059249e-02, 3.50710137e-02,
       3.40452270e-02, 3.32582107e-02, 3.07888139e-02, 2.99418829e-02,
       2.98386064e-02, 2.90265643e-02, 2.84050295e-02, 2.75235307e-02,
       2.67648958e-02, 2.64538757e-02, 2.27422479e-02, 2.09773868e-02,
       1.90604822e-02, 1.77636535e-02, 1.69051260e-02, 1.62791683e-02,
       1.52097793e-02, 1.44799030e-02, 1.22756402e-02, 1.18222605e-02,
       1.06760730e-02, 1.00919538e-02, 8.94729489e-03, 6.60407766e-03,
       5.71989396e-03, 8.92816949e-30, 2.40570590e-30, 7.05358138e-34,
       7.05358138e-34])
```

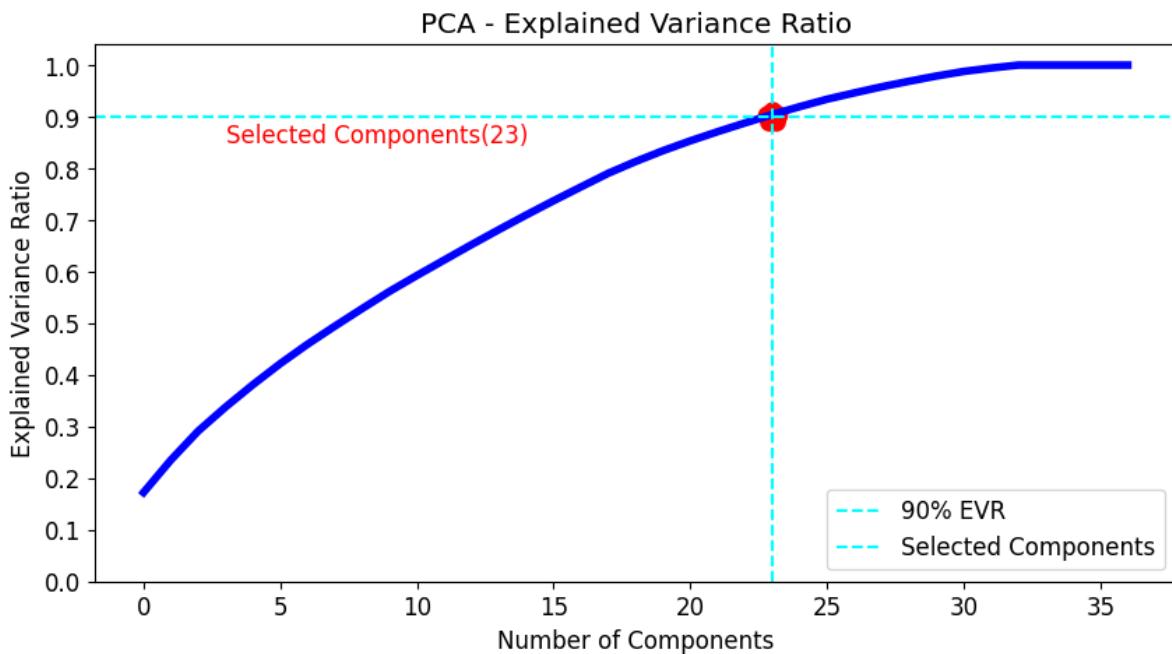
In [154]:

```
1 plt.figure(figsize=(10,5))
2 plt.plot(np.cumsum(pca.explained_variance_ratio_), linewidth=4, markersize=10, color="blue")
3 plt.xlabel('Number of Components')
4 plt.ylabel('Explained Variance Ratio')
5 plt.show()
```



In [163]:

```
1 plt.figure(figsize=(10,5))
2 plt.plot(np.cumsum(pca.explained_variance_ratio_),color='blue', linewidth=4, )
3 plt.axhline(y=0.9, color='cyan', linestyle='--', label='90% EVR')
4 plt.axvline(x=23, color='cyan', linestyle='--', label='Selected Components')
5 plt.scatter(23,0.9,marker="*",linewidths=10,color="red")
6 plt.title('PCA - Explained Variance Ratio')
7 plt.xlabel('Number of Components')
8 plt.ylabel('Explained Variance Ratio')
9 plt.yticks(np.arange(0, 1.1, step=0.1))
10 plt.text(3, 0.85, 'Selected Components(23)', color='r')
11 plt.legend(loc="lower right")
12 plt.show()
```



In [164]:

```
1 #23 components are showing 90% of the variance in our data, so we'll build our PCA again w
```

In [165]:

```
1 pca=PCA(n_components=23)
```

In [166]:

```
1 df3=pd.DataFrame(pca.fit_transform(df2))
```

In [167]:

```
1 df3.head()
```

Out[167]:

	0	1	2	3	4	5	6	7	8
0	3.921261	-0.961448	-0.418362	1.552713	-0.998744	2.012662	-1.186018	0.352056	-0.528134
1	-2.361402	-0.512049	-0.304983	0.592715	-1.036934	-0.347183	-1.857775	0.830816	-0.269280
2	1.500922	-0.427692	-1.380939	0.160132	-1.317587	-0.722257	0.748517	-0.629505	-0.117337
3	-2.259714	-1.557122	-0.086753	0.581038	-1.127287	-0.954844	1.062398	-0.710251	-0.124947
4	-0.590543	0.331769	-0.096791	-0.775769	1.749521	0.541293	-0.048923	-1.439071	-1.174442

5 rows × 23 columns



Clustering

1.Kmeans Clustering

In [168]:

```
1 from sklearn.cluster import KMeans
```

In [169]:

```
1 wcss=[]
```

In [170]:

```
1 for i in range(1,5):
2     kmean=KMeans(n_clusters=i,init="random",random_state=42)
3     kmean.fit(df3)
4     wcss.append(kmean.inertia_)
```

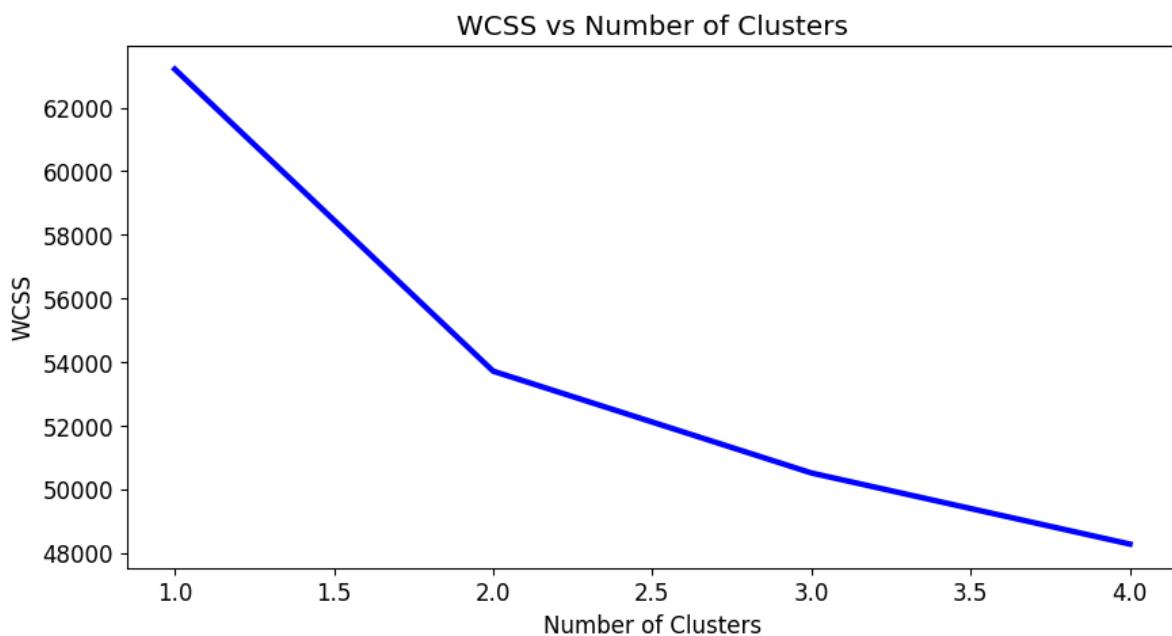
In [171]:

```
1 print(wcss)
```

[63209.17068842075, 53715.05325913177, 50513.20588634547, 48277.99630918754]

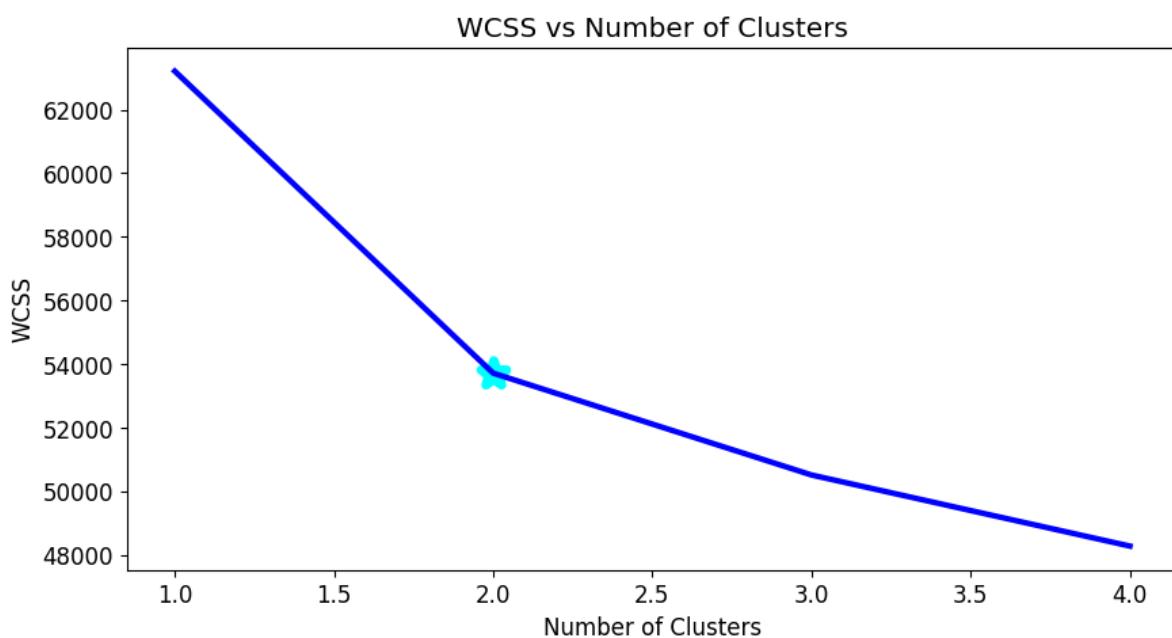
In [172]:

```
1 plt.figure(figsize=(10,5))
2 plt.plot(range(1,5),wcss,color="blue",linewidth=3)
3 plt.title('WCSS vs Number of Clusters')
4 plt.xlabel('Number of Clusters')
5 plt.ylabel('WCSS')
6 plt.show()
```



In [209]:

```
1 plt.figure(figsize=(10,5))
2 plt.plot(range(1,5),wcss,color="blue",linewidth=3)
3 plt.title('WCSS vs Number of Clusters')
4 plt.scatter(2,wcss[1],s=200,c="cyan",marker="*",linewidth=5)
5 plt.xlabel('Number of Clusters')
6 plt.ylabel('WCSS')
7 plt.show()
```



In [174]:

```
1 from sklearn.metrics import silhouette_score
```

In [175]:

```
1 silhouette_score_lst=[]
2
3 for i in range(2,6):
4     silhouette_score_lst.append(silhouette_score(df3,(KMeans(n_clusters=i).fit_predict(df3)
```

In [176]:

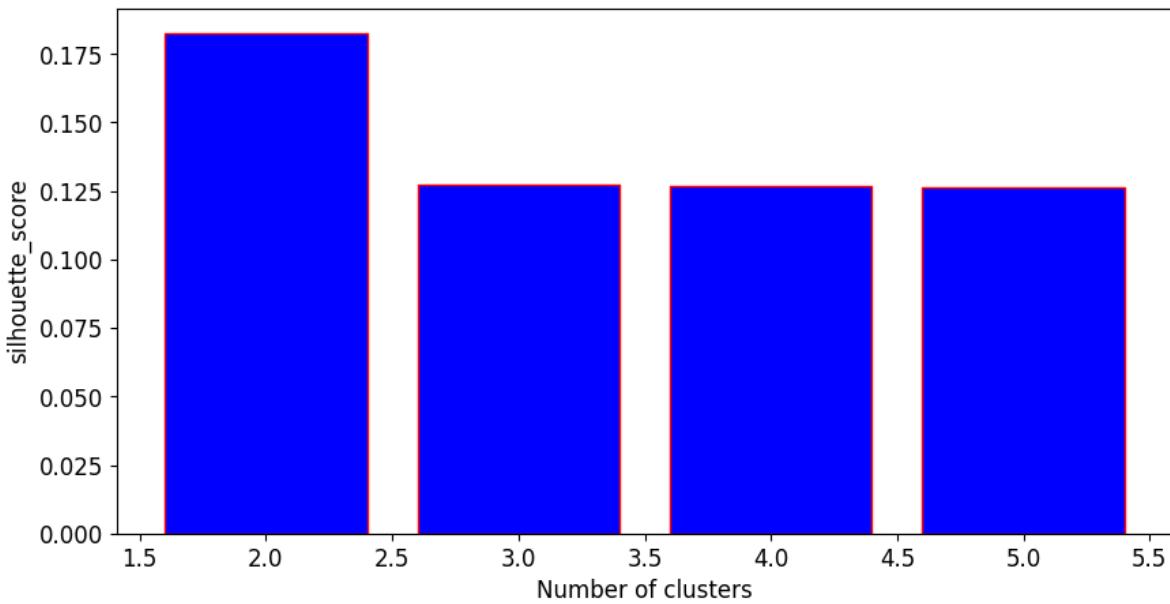
```
1 silhouette_score_lst
```

Out[176]:

```
[0.18256678718476013,
 0.12731079376982762,
 0.12685890301831806,
 0.1261638680247387]
```

In [212]:

```
1 #plotting
2 plt.figure(figsize=(10,5))
3
4 k=[2,3,4,5]
5 plt.bar(k,silhouette_score_lst,color="blue",edgecolor="red")
6
7 plt.xlabel("Number of clusters")
8 plt.ylabel("silhouette_score")
9 plt.show()
```



Elbow Graph shows that we should select 2 clusters and also silhouette scores is maximum when clusters are 2

In [178]:

```
1 #Making model again with 2 clusters
```

In [179]:

```
1 kmeans=KMeans(n_clusters=2,random_state=42)
```

In [180]:

```
1 kmeans.fit(df3)
```

Out[180]:

```
KMeans(n_clusters=2, random_state=42)
```

In [181]:

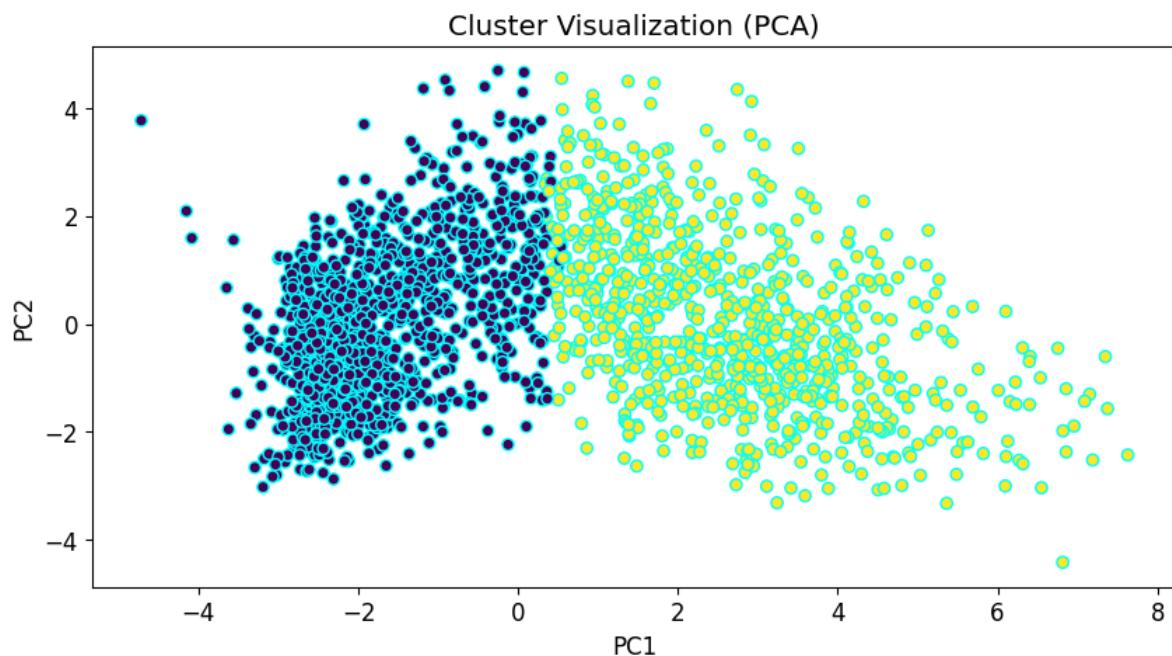
```
1 pred=kmeans.predict(df3)
2 pred
```

Out[181]:

```
array([1, 0, 1, ..., 1, 1, 0])
```

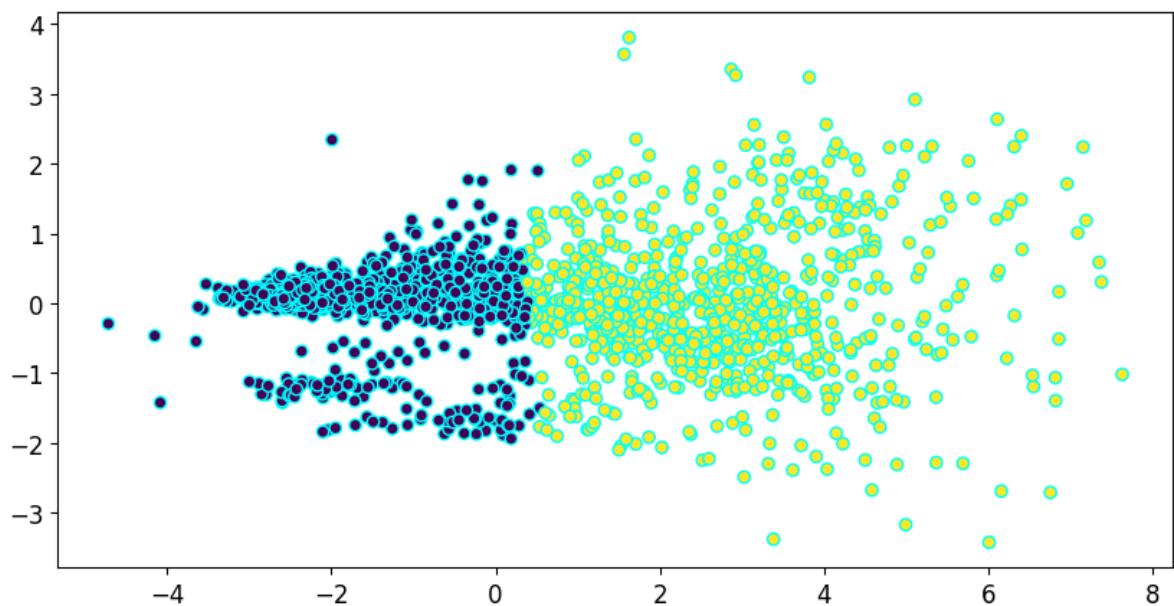
In [213]:

```
1 x = df3.iloc[:, 0]
2 y = df3.iloc[:, 1]
3 plt.figure(figsize=(10,5))
4 plt.scatter(x, y, c=pred, edgecolors="cyan")
5 plt.xlabel('PC1')
6 plt.ylabel('PC2')
7 plt.title('Cluster Visualization (PCA)')
8 plt.show()
```



In [220]:

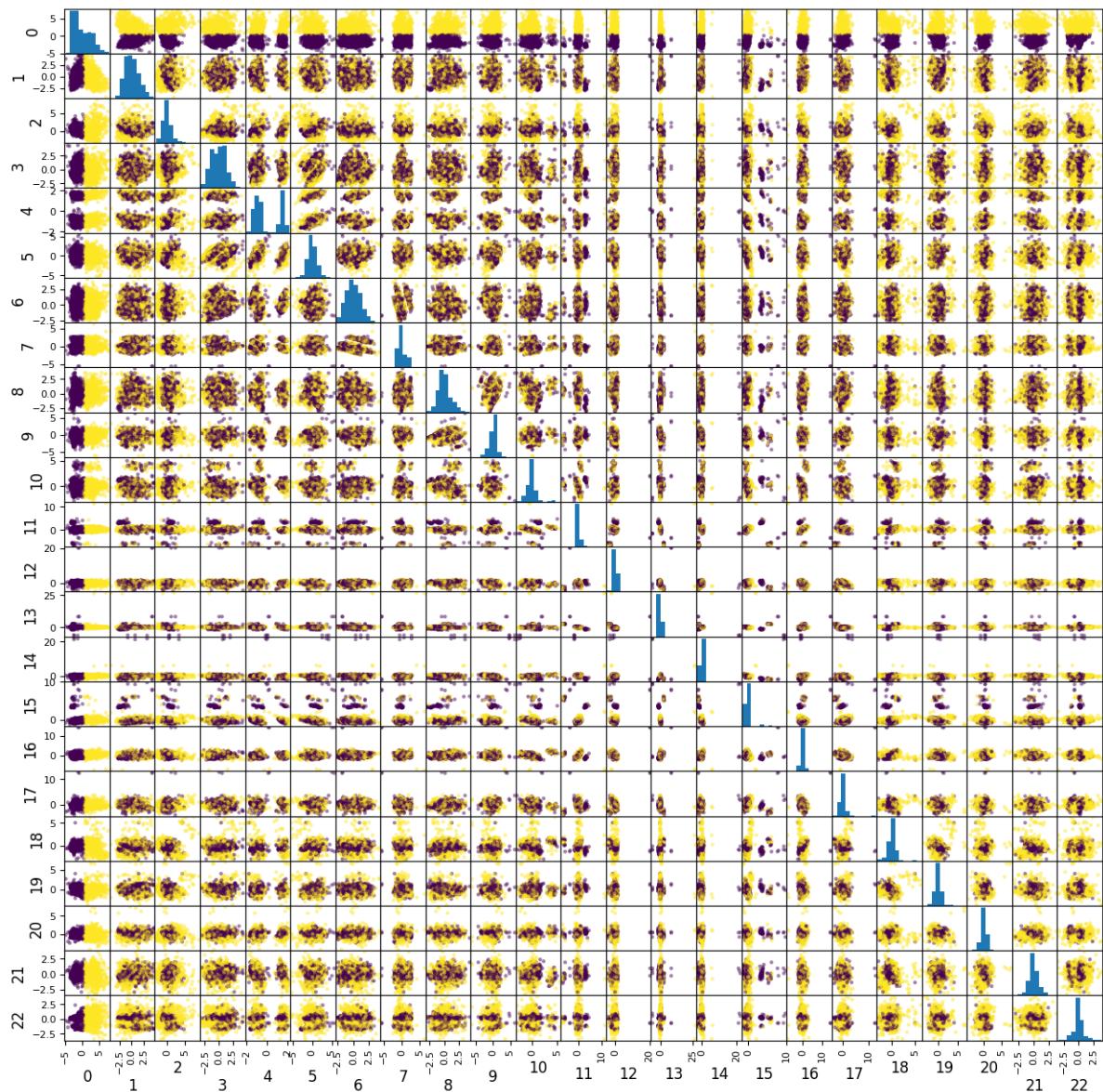
```
1 x = df3.iloc[:, 0]
2 y = df3.iloc[:, 22]
3 plt.figure(figsize=(10,5))
4 plt.scatter(x, y, c=pred, edgecolors="cyan")
5 plt.show()
```



In [184]:

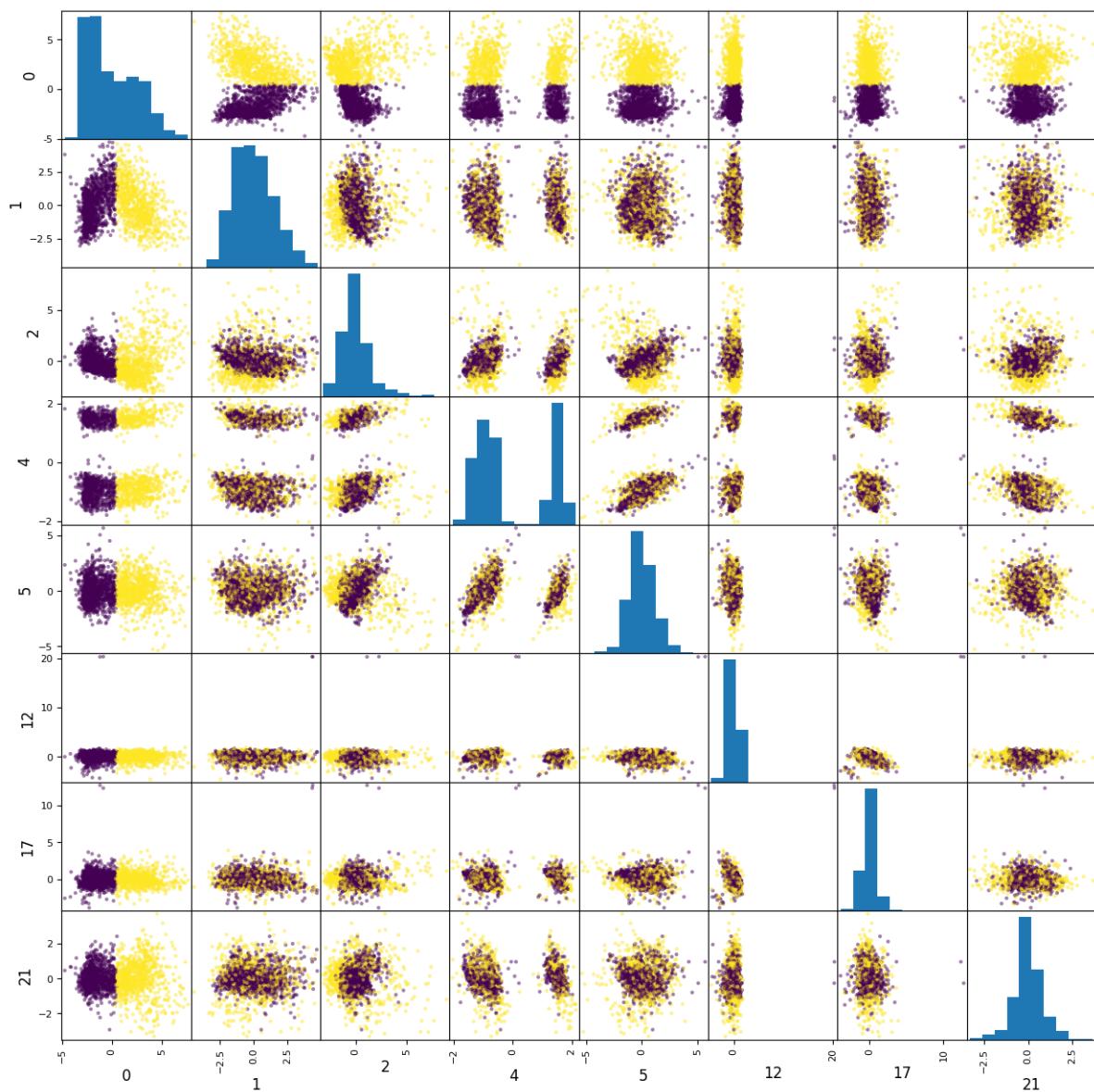
```
1 # visualizing all
2 plt.figure(figsize=(10,5))
3
4 pd.plotting.scatter_matrix(pd.DataFrame(df3), c=pred, figsize=(15,15), diagonal='hist')
5 plt.show()
```

<Figure size 1000x500 with 0 Axes>



In [221]:

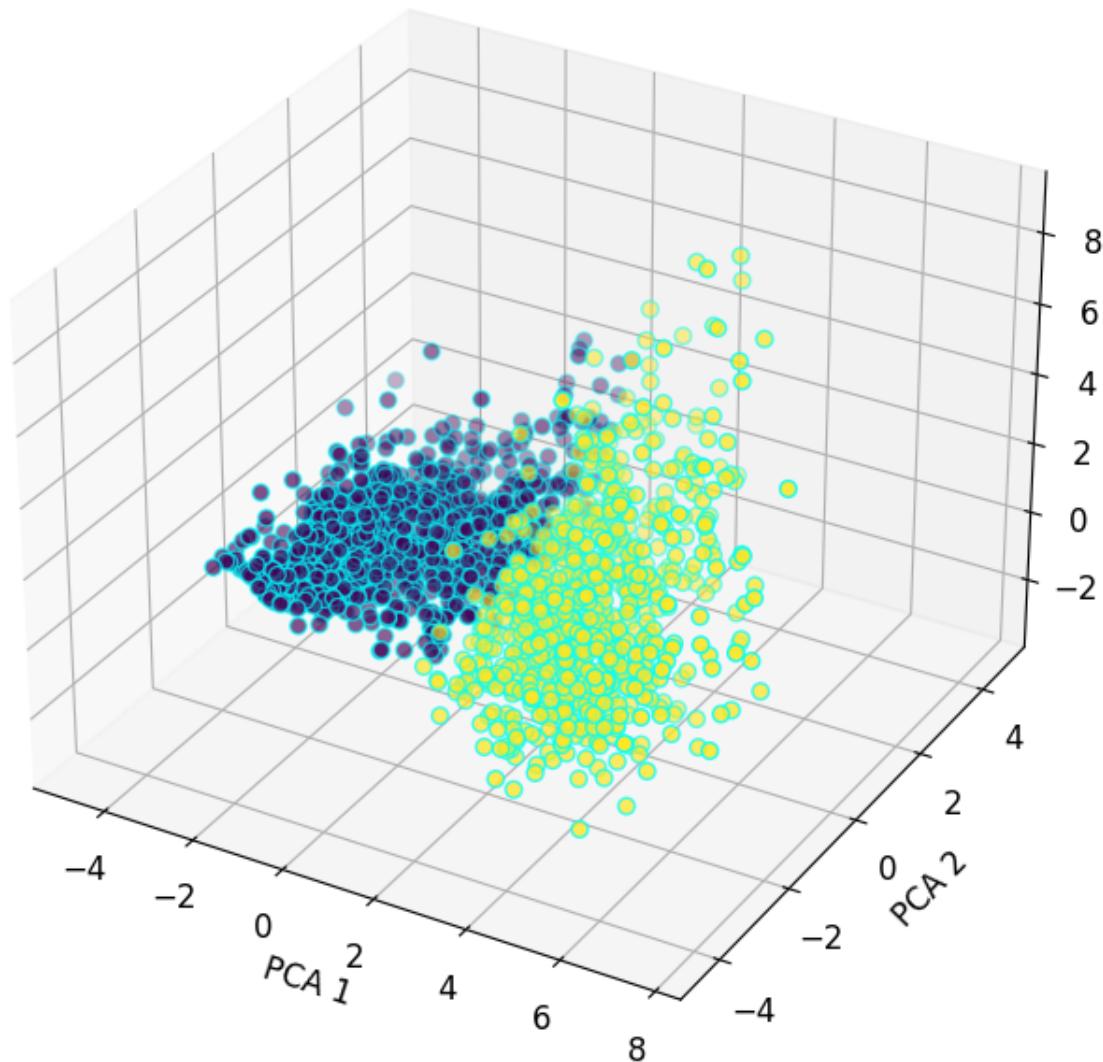
```
1 #Viewing subset of it
2 dim_subset = [0, 1, 2, 4, 5, 12, 17, 21]
3 pd.plotting.scatter_matrix(pd.DataFrame(df3.iloc[:, dim_subset]), c=pred, figsize=(15,15))
4 plt.show()
```



In [186]:

```
1 from mpl_toolkits.mplot3d import Axes3D
2
3 # Create a 3D plot with different colors for each cluster
4 fig = plt.figure(figsize=(8,8))
5 ax = fig.add_subplot(111, projection='3d')
6 ax.scatter(df3.iloc[:,0], df3.iloc[:,1], df3.iloc[:,2], c=pred, s=40, edgecolor="cyan")
7 ax.set_xlabel('PCA 1')
8 ax.set_ylabel('PCA 2')
9 ax.set_zlabel('PCA 3')
10 ax.set_title('K-Means Clustering Results (3D)')
11 plt.show()
```

K-Means Clustering Results (3D)



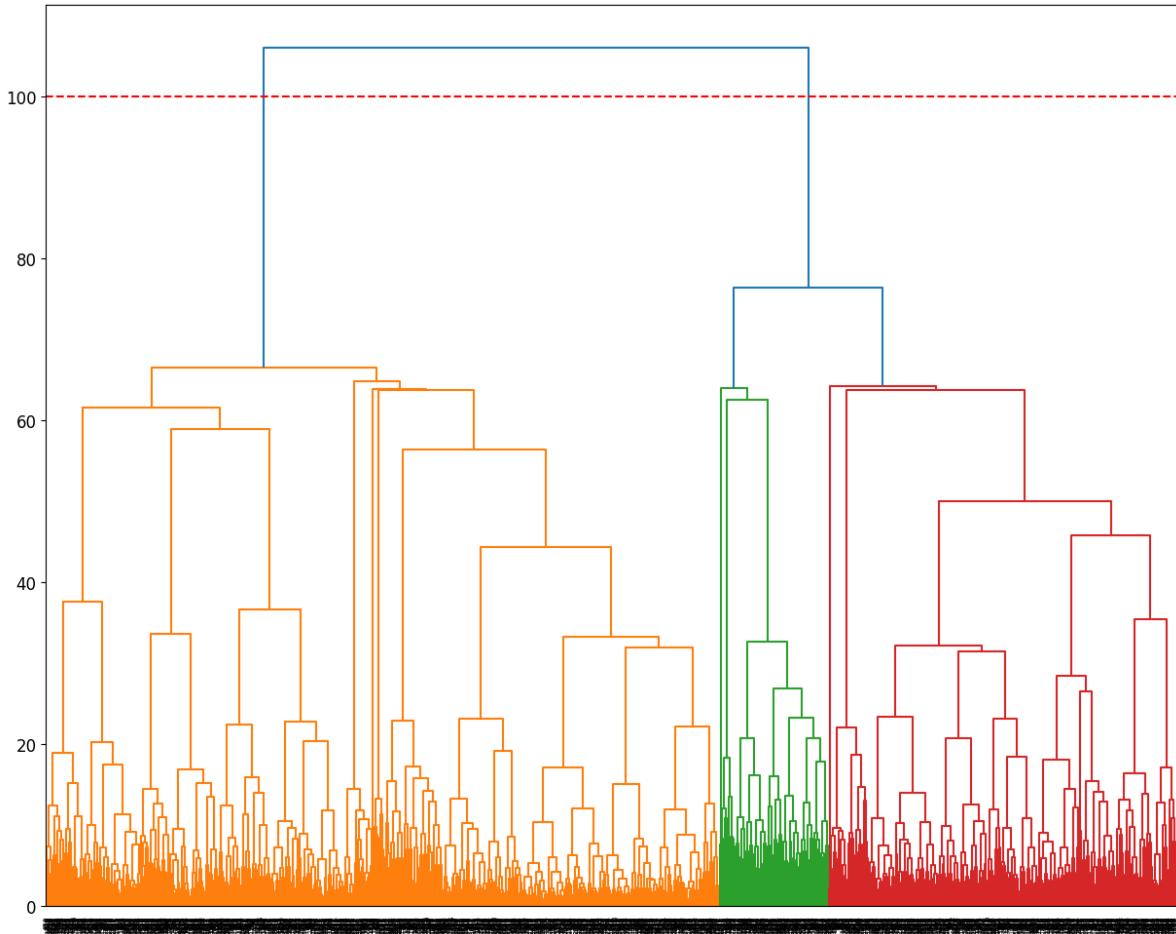
2.Hierachial Clustering

In [187]:

```
1 from sklearn.cluster import AgglomerativeClustering  
2 import scipy.cluster.hierarchy as shc
```

In [223]:

```
1 plt.figure(figsize=(15,12))  
2 den=shc.dendrogram(shc.linkage(df3,method="ward"))  
3 plt.axhline(y=100,color="r",linestyle="--")  
4 plt.show()
```



Hierachial is showing with 100 clusters so its definately not a good approach.

3.DBSCAN

In [189]:

```
1 from sklearn.cluster import DBSCAN
```

In [190]:

```
1 dbSCAN=DBSCAN(eps=27,min_samples=70)  
2 dbSCAN.fit(df3)
```

Out[190]:

```
DBSCAN(eps=27, min_samples=70)
```

In [191]:

```
1 #Checking labels
2 dbscan.labels_
```

Out[191]:

```
array([0, 0, 0, ..., 0, 0, 0], dtype=int64)
```

In [192]:

```
1 #Total no of clusters
2 len(set(dbscan.labels_))
```

Out[192]:

```
2
```

In [193]:

```
1 from sklearn.metrics import silhouette_score
2 silhouette_score(df3,dbscan.labels_)
```

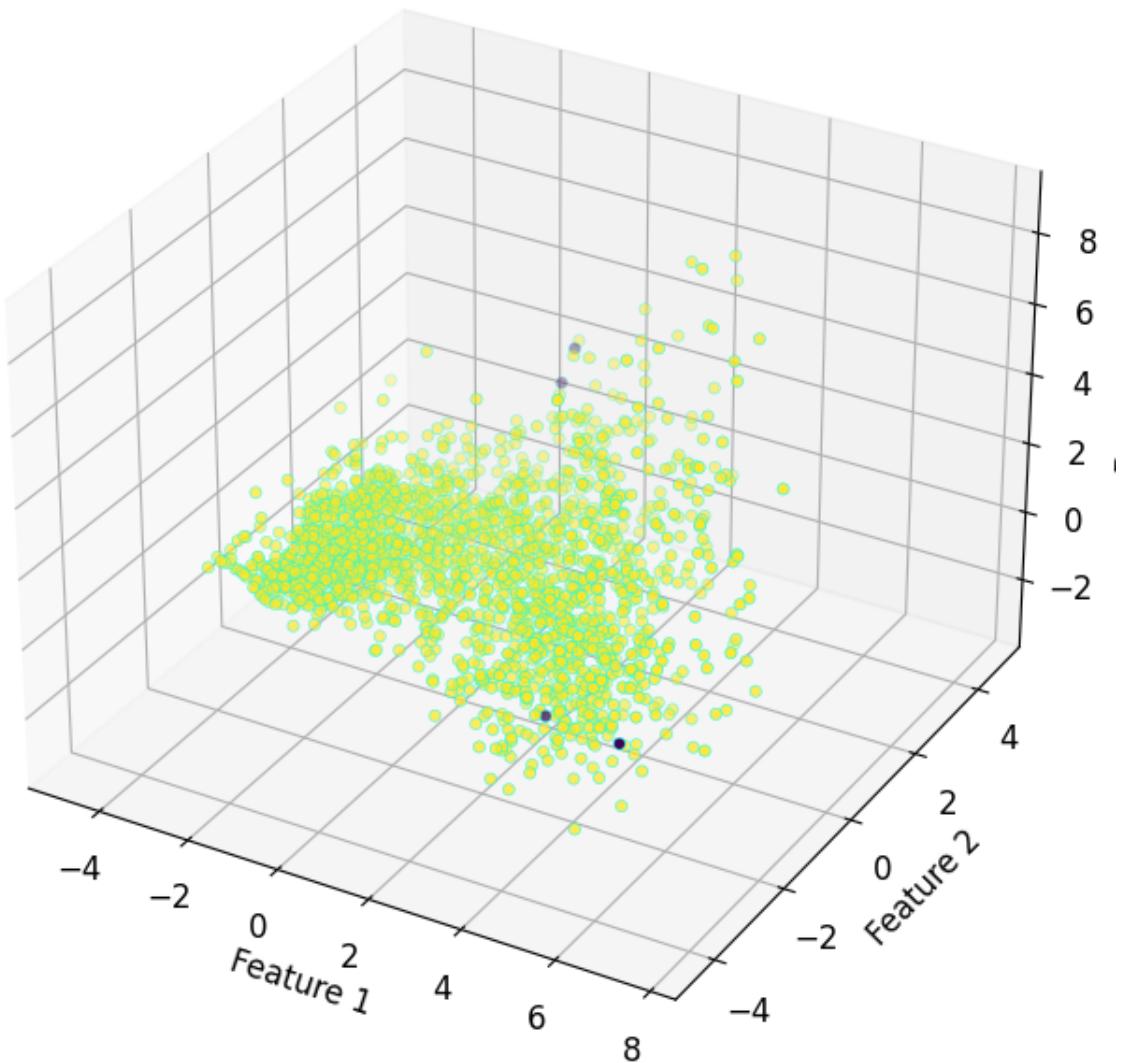
Out[193]:

```
0.7784382051017664
```

In [258]:

```
1 fig = plt.figure(figsize=(10, 8))
2 ax = fig.add_subplot(111, projection='3d')
3 ax.scatter(df3.iloc[:, 0], df3.iloc[:, 1], df3.iloc[:, 2], c=dbSCAN.labels_, edgecolor="cyan")
4 ax.set_xlabel('Feature 1')
5 ax.set_ylabel('Feature 2')
6 ax.set_zlabel('Feature 3')
7 ax.set_title('DBSCAN Clustering Results (3D)')
8 plt.show()
```

DBSCAN Clustering Results (3D)



Conclusions:

1. Both KMeans and DBSCAN were able to efficiently separate out the clusters in the dataset.
2. However, when comparing the two algorithms using the silhouette score, DBSCAN outperformed KMeans with a score of 0.7784, while KMeans had a score of 0.182. This suggests that, despite the fact that KMeans appeared to perform better visually, the clusters produced by DBSCAN were more distinct and better separated.
3. Therefore, based on the silhouette score, DBSCAN is a better clustering algorithm for this dataset.

Project By:

Rohit Vyawahare

Linkedin <https://www.linkedin.com/in/rohitvyawahare2001/>
[\(https://www.linkedin.com/in/rohitvyawahare2001/\)](https://www.linkedin.com/in/rohitvyawahare2001/)

Github : <https://github.com/RohitVyawahare2001>
[\(https://github.com/RohitVyawahare2001\)](https://github.com/RohitVyawahare2001)