

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import datetime

df=pd.read_csv('marketing_data.csv')
```

```
In [51]: df

Out[51]:
```

	ID	Year_Birth	Education	Marital_Status	Income	Kidhome	Teenhome	DT_Customer	Recency	MntWines	...	NumStorePurchases	NumWebVisitsMonth	AcceptedCmp3	AcceptedCmp4	AcceptedCmp5	Response	Complain	Country
0	1826	1970	Graduation	Divorced	\$84,835.00	0	0	6/16/14	0	189	...	6	1	0	0	0	0	0	US
1	1	1961	Graduation	Single	\$57,991.00	0	0	6/15/14	0	464	...	7	5	0	0	0	0	0	CA
2	10476	1959	Graduation	Married	\$67,267.00	0	1	5/13/14	0	134	...	5	2	0	0	0	0	0	US
3	1386	1967	Graduation	Together	\$32,474.00	1	1	5/11/14	0	10	...	2	7	0	0	0	0	0	US
4	5371	1989	Graduation	Single	\$21,474.00	1	0	4/18/14	0	6	...	2	7	1	0	0	0	0	US
...
2235	10142	1976	PHD	Divorced	\$66,476.00	0	1	3/7/13	99	372	...	11	4	0	0	0	0	0	US
2236	5263	1977	2n Cycle	Married	\$31,056.00	1	0	12/21/13	99	5	...	3	8	0	0	0	0	0	US
2237	22	1976	Graduation	Divorced	\$46,310.00	1	0	12/31/12	99	185	...	5	8	0	0	0	0	0	US
2238	528	1978	Graduation	Married	\$65,819.00	0	0	11/28/12	99	267	...	10	3	0	0	0	0	0	US
2239	4070	1969	PHD	Married	\$84,871.00	0	2	9/1/12	99	169	...	4	7	0	0	0	0	0	US

```
In [52]: df.shape

Out[52]: (2240, 28)
```

```
In [53]: df.info()

Out[53]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2240 entries, 0 to 2239
Data columns (total 28 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  NumStorePurchases       2240 non-null   int64
17  NumCatalogPurchases    2240 non-null   int64
18  NumWebVisitsMonth       2240 non-null   int64
19  AcceptedCmp3            2240 non-null   int64
20  AcceptedCmp4            2240 non-null   int64
21  AcceptedCmp5            2240 non-null   int64
22  AcceptedCmp6            2240 non-null   int64
23  AcceptedCmp7            2240 non-null   int64
24  AcceptedCmp8            2240 non-null   int64
25  Response                2240 non-null   int64
26  Complain                2240 non-null   int64
27  Country                 2240 non-null   object
dtypes: int64(23), object(5)
memory usage: 480.1 KB
```

```
In [54]: df.rename({'Income':'Income'},axis=1,inplace=True)
```

Check for null values

```
In [55]: df.isna().sum()

Out[55]:
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
NumStorePurchases      0
NumCatalogPurchases    0
NumWebVisitsMonth      0
AcceptedCmp3           0
AcceptedCmp4           0
AcceptedCmp5           0
AcceptedCmp6           0
AcceptedCmp7           0
AcceptedCmp8           0
Response               0
Complain               0
Country                0
dtype: int64
```

```
In [56]: df.dropna(inplace=True)
```

```
In [57]: df.info()

Out[57]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2216 entries, 0 to 2239
Data columns (total 28 columns):
#   Column                Non-Null Count  Dtype
---  --
0   ID                     2216 non-null   int64
1   Year_Birth             2216 non-null   int64
2   Education              2216 non-null   object
3   Marital_Status         2216 non-null   object
4   Income                 2216 non-null   float64
5   Kidhome                2216 non-null   int64
6   Teenhome               2216 non-null   int64
7   DT_Customer            2216 non-null   object
8   Recency                2216 non-null   int64
9   MntWines                2216 non-null   int64
10  MntFruits               2216 non-null   int64
11  MntMeatProducts         2216 non-null   int64
12  MntFishProducts         2216 non-null   int64
13  MntSweetProducts        2216 non-null   int64
14  MntGoldProds            2216 non-null   int64
15  NumDealsPurchases       2216 non-null   int64
16  NumStorePurchases       2216 non-null   int64
17  NumCatalogPurchases    2216 non-null   int64
18  NumWebVisitsMonth       2216 non-null   int64
19  AcceptedCmp3            2216 non-null   int64
20  AcceptedCmp4            2216 non-null   int64
21  AcceptedCmp5            2216 non-null   int64
22  AcceptedCmp6            2216 non-null   int64
23  AcceptedCmp7            2216 non-null   int64
24  AcceptedCmp8            2216 non-null   int64
25  Response                2216 non-null   int64
26  Complain                2216 non-null   int64
27  Country                 2216 non-null   object
dtypes: int64(23), object(5)
memory usage: 582.1 KB
```

Variable transformation

```
In [58]: # change income to float
df['Income'] = df['Income'].str.replace('$','',).astype(float)

C:\Users\DELL\AppData\Local\Temp\ipykernel_25980\968687995.py:3: FutureWarning: The default value of regex will change from True to False in a future version.
df['Income'] = df['Income'].str.replace('$','',).astype(float)

In [59]: df['DT_Customer'].head()

Out[59]:
0    6/16/14
1    6/15/14
2    5/13/14
3    5/11/14
4    4/8/14
Name: DT_Customer, dtype: object
```

```
In [60]: df['DT_Customer'].add.to_datetime(df['DT_Customer'])
```

```
In [61]: df.info()

Out[61]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2216 entries, 0 to 2239
Data columns (total 28 columns):
#   Column                Non-Null Count  Dtype
---  --
0   ID                     2216 non-null   int64
1   Year_Birth             2216 non-null   int64
2   Education              2216 non-null   object
3   Marital_Status         2216 non-null   object
4   Income                 2216 non-null   float64
5   Kidhome                2216 non-null   int64
6   Teenhome               2216 non-null   int64
7   DT_Customer            2216 non-null   datetime64[ns]
8   Recency                2216 non-null   int64
9   MntWines                2216 non-null   int64
10  MntFruits               2216 non-null   int64
11  MntMeatProducts         2216 non-null   int64
12  MntFishProducts         2216 non-null   int64
13  MntSweetProducts        2216 non-null   int64
14  MntGoldProds            2216 non-null   int64
15  NumDealsPurchases       2216 non-null   int64
16  NumStorePurchases       2216 non-null   int64
17  NumCatalogPurchases    2216 non-null   int64
18  NumWebVisitsMonth       2216 non-null   int64
19  AcceptedCmp3            2216 non-null   int64
20  AcceptedCmp4            2216 non-null   int64
21  AcceptedCmp5            2216 non-null   int64
22  AcceptedCmp6            2216 non-null   int64
23  AcceptedCmp7            2216 non-null   int64
24  AcceptedCmp8            2216 non-null   int64
25  Response                2216 non-null   int64
26  Complain                2216 non-null   int64
27  Country                 2216 non-null   object
dtypes: datetime64[ns](1), float64(1), int64(23), object(3)
memory usage: 581.1 KB
```

```
In [62]: df.head()

Out[62]:
```

	ID	Year_Birth	Education	Marital_Status	Income	Kidhome	Teenhome	DT_Customer	Recency	MntWines	...	NumStorePurchases	NumWebVisitsMonth	AcceptedCmp3	AcceptedCmp4	AcceptedCmp5	Response	Complain	Country
0	1826	1970	Graduation	Divorced	\$84835.0	0	0	2014-06-16	0	189	...	6	1	0	0	0	0	0	US
1	1	1961	Graduation	Single	\$57991.0	0	0	2014-06-15	0	464	...	7	5	0	0	0	0	0	CA
2	10476	1959	Graduation	Married	\$67267.0	0	1	2014-05-13	0	134	...	5	2	0	0	0	0	0	US
3	1386	1967	Graduation	Together	\$32474.0	1	1	2014-05-11	0	10	...	2	7	0	0	0	0	0	US
4	5371	1989	Graduation	Single	\$21474.0	1	0	2014-04-08	0	6	...	2	7	1	0	0	0	0	US

5 rows x 28 columns

Check for outliers

```
In [63]: df.columns

Out[63]:
Index(['ID', 'Year_Birth', 'Education', 'Marital_Status', 'Income', 'Kidhome',
       'Teenhome', 'DT_Customer', 'Recency', 'MntWines', 'MntFruits',
       'MntMeatProducts', 'MntFishProducts', 'MntSweetProducts',
       'MntGoldProds', 'NumDealsPurchases', 'NumStorePurchases', 'NumCatalogPurchases',
       'NumWebVisitsMonth', 'AcceptedCmp3', 'AcceptedCmp4', 'AcceptedCmp5',
       'AcceptedCmp6', 'AcceptedCmp7', 'AcceptedCmp8', 'Response', 'Complain',
       'Country'],
      dtype='object')
```

```
In [64]: numeric_cols = df.select_dtypes(include=['number']).columns
numeric_cols = numeric_cols.drop(['ID', 'AcceptedCmp3', 'AcceptedCmp4', 'AcceptedCmp5', 'AcceptedCmp6',
                                  'AcceptedCmp7', 'AcceptedCmp8', 'Response', 'Complain'])
df[numeric_cols].plot(subplots=True, layout=(4,4), kind='box', figsize=(18,15))

plt.show()
```



Outliers are seen in year_birth and income columns

```
In [65]: df = df.drop(df[df['Year_Birth'] <= 1900].index)

Out[65]: df['Year_Birth'].sort_values().head()

Out[65]:
2171    1949
2140    1941
1350    1943
1444    1943
1289    1943
Name: Year_Birth, dtype: int64
```

```
In [67]: df[df['Income'] >= 600000]

Out[67]:
```

	ID	Year_Birth	Education	Marital_Status	Income	Kidhome	Teenhome	DT_Customer	Recency	MntWines	...	NumStorePurchases	NumWebVisitsMonth	AcceptedCmp3	AcceptedCmp4	AcceptedCmp5	Response	Complain	Country
527	9432	1977	Graduation	Together	\$66666.0	1	0	2013-06-02	23	9	...	3	6	0	0	0	0	0	US

1 rows x 28 columns

```
In [68]: df = df.drop(df[df['Income'] >= 600000].index)
```

```
In [69]: df['Income'].sort_values()

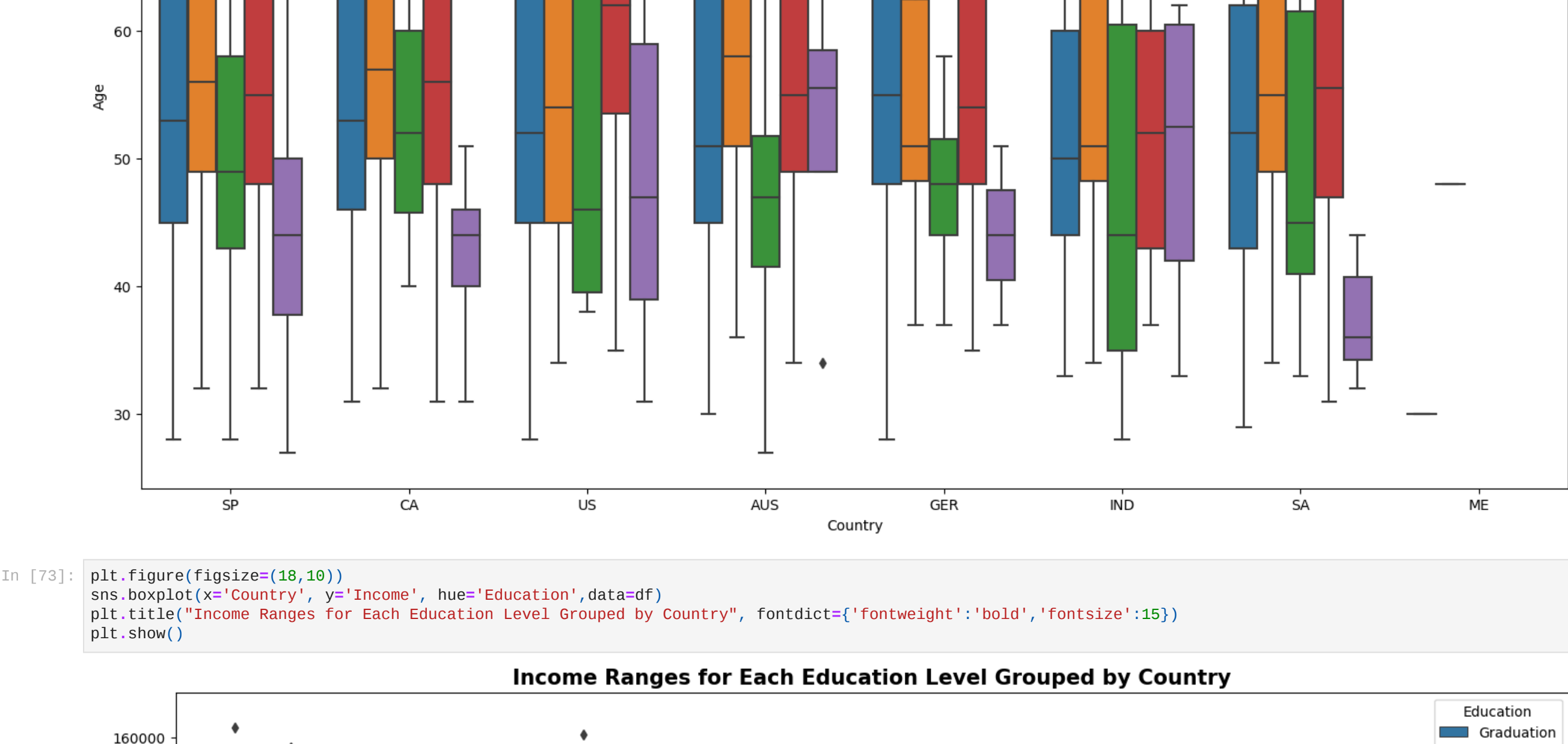
Out[69]:
1456    1739.0
961     2447.0
1291     3502.0
663     4923.0
14     4428.0
Name: Income, Length: 2212, dtype: float64
```

```
In [70]: # lets check education and marital status columns
df[Education].value_counts()

Out[70]:
Graduation    1115
PHD            480
Master         365
2n Cycle       198
Basic           54
Name: Education, dtype: int64
```

```
In [71]: df['Age'] = datetime.datetime.now().year - df['Year_Birth']
```

```
In [72]: plt.figure(figsize=(18,10))
sns.boxplot(x='Country', y='Age', hue='Education', data=df)
plt.title('Age Ranges for Each Education Level Grouped by Country', fontdict={'fontweight':'bold', 'fontsize':15})
plt.show()
```



```
In [73]: plt.figure(figsize=(18,10))
sns.boxplot(x='Country', y='Income', hue='Education', data=df)
plt.title('Income Ranges for Each Education Level Grouped by Country', fontdict={'fontweight':'bold', 'fontsize':15})
plt.show()
```



```
In [74]: df['New_Education'] = df[['Graduation','Undergraduate','2n Cycle','Master','PHD','PHD', 'Basic','Basic', 'Master','Master']]
```

```
In [75]: df['New_Education'].value_counts(dropna=False)
```

```
Out[75]:
Undergraduate    1115
Master           480
PHD              365
Basic            198
Name: New_Education, dtype: int64
```

```
In [76]: df['New_Education'].value_counts(dropna=False)
```

```
Out[76]:
Graduation    1115
PHD            480
Master         365
2n Cycle       198
Basic           54
Name: Education, dtype: int64
```

```
In [77]: df['Marital_Status'].value_counts()
```

```
Out[77]:
Married      857
Together     571
Single        470
Divorced      221
Widow         76
Alone          3
YOLO           2
Absurd         2
Name: Marital_Status, dtype: int64
```

```
In [78]: fcs = 4
cols = 2
```



```
In [79]: df['Marital_Status'].mode()
```

```
Out[79]:
0    Married
Name: Marital_Status, dtype: object
```

```
In [80]: df['New_Marital_Status'] = df[['Divorced':'Divorced', 'Single':'Single', 'Married':'Married', 'Together':'Single',
                                     'Widow':'Widow', 'YOLO':'Married', 'Alone':'Single', 'Absurd':'Married']]
df['New_Marital_Status'] = df['New_Marital_Status'].map(New_Marital_Status)
```

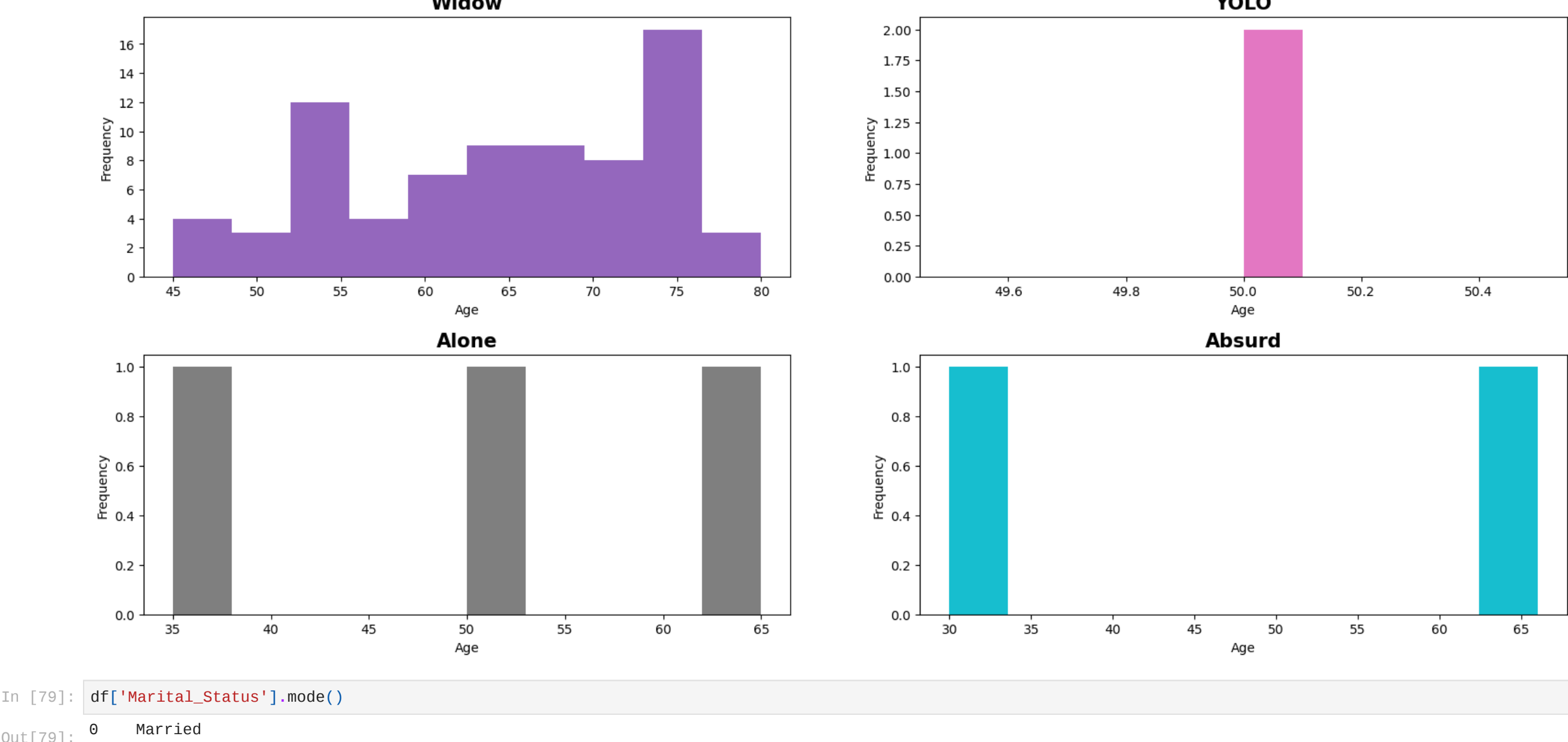
```
In [81]: df['New_Marital_Status'].value_counts(dropna=False)
```

```
Out[81]:
Single      1844
Married     861
Divorced    221
Widow        76
Name: New_Marital_Status, dtype: int64
```

```
In [82]: df['Marital_Status'].value_counts(dropna=False)
```

```
Out[82]:
Married      857
Together     571
Single        470
Divorced      221
Widow         76
Alone          3
YOLO           2
Absurd         2
Name: Marital_Status, dtype: int64
```

```
In [83]: df.corr()
```



ordinal encoding and one hot encoding

```
In [86]: categorical_df = df.select_dtypes(exclude='number')

dummy_df = pd.DataFrame()
for col in categorical_df.columns:
    dummy = pd.get_dummies(categorical_df[col], prefix=col)
    dummy_df = pd.concat([dummy_df, dummy], axis=1)

df = pd.concat([df, dummy_df], axis=1)
df.drop(categorical_df.columns, axis=1, inplace=True)
```

```
In [87]: df.info()

Out[87]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2212 entries, 0 to 2239
Data columns (total 43 columns):
#   Column                Non-Null Count  Dtype
---  --
0   ID                     2212 non-null   int64
1   Year_Birth             2212 non-null   int64
2   Income                 2212 non-null   float64
3   Kidhome                2212 non-null   int64
4   Teenhome               2212 non-null   int64
5   Recency                2212 non-null   int64
6   MntWines                2212 non-null   int64
7   MntFruits               2212 non-null   int64
8   MntMeatProducts         2212 non-null   int64
9   MntFishProducts         2212 non-null   int64
10  MntSweetProducts        2212 non-null   int64
11  MntGoldProds            2212 non-null   int64
12  NumDealsPurchases       2212 non-null   int64
13  NumCatalogPurchases    2212 non-null   int64
14  NumStorePurchases       2212 non-null   int64
15  NumWebVisitsMonth       2212 non-null   int64
16  AcceptedCmp3            2212 non-null   int64
17  AcceptedCmp4            2212 non-null   int64
18  AcceptedCmp5            2212 non-null   int64
19  AcceptedCmp6            2212 non-null   int64
20  AcceptedCmp7            2212 non-null   int64
21  AcceptedCmp8            2212 non-null   int64
22  Response                2212 non-null   int64
23  Complain                2212 non-null   int64
24  Education_2n Cycle      2212 non-null   int64
25  Education_Basic         2212 non-null   int64
26  Education_Graduation    2212 non-null   int64
27  Education_Master        2212 non-null   int64
28  Education_2n Cycle      2212 non-null   int64
29  Education_Basic         2212 non-null   int64
30  Education_Graduation    2212 non-null   int64
31  Education_Master        2212 non-null   int64
32  Education_Phd           2212 non-null   int64
33  Education_Master        2212 non-null   int64
34  Marital_Status_Absurd   2212 non-null   int64
35  Marital_Status_Absurd   2212 non-null   int64
36  Marital_Status_Absurd   2212 non-null   int64
37  Marital_Status_Single   2212 non-null   int64
38  Marital_Status_Single   2212 non-null   int64
39  Marital_Status_Single   2212 non-null   int64
40  Marital_Status_Single   2212 non-null   int64
41  Marital_Status_Single   2212 non-null   int64
42  Marital_Status_Single   2212 non-null   int64
dtypes: float64(1), int64(24), uint8(18)
memory usage: 488.2 KB
```

```
In [ ]:
```

```
In [ ]:
```