

Title: Final Report

Author: Dennis

Date: Jan 2022

Import Your Data

```
In [1]:  
import numpy as np  
from numpy import count_nonzero  
from numpy import median  
from pandas import DataFrame  
import pandas as pd  
import matplotlib.pyplot as plt  
import seaborn as sns  
import plotly.express as px  
import random  
  
import statsmodels.api as sm  
import statsmodels.formula.api as smf  
from statsmodels.formula.api import ols  
  
import datetime  
from datetime import timedelta  
  
import scipy.stats  
from collections import Counter  
  
matplotlib_inline  
#set_ipython.terminal.autosave frequency in seconds  
autosave = 60  
sns.set_style('dark')  
sns.set(font_scale=1.2)  
  
plt.rcParams['axes.titlesize']=9  
plt.rcParams['labelsize']=14  
plt.rcParams['xtick', 'labelsize']=12  
plt.rcParams['ytick', 'labelsize']=12  
  
import warnings  
warnings.filterwarnings('ignore')  
  
pd.set_option('display.max_columns',None)  
pd.set_option('display.max_rows',None)  
pd.set_option('display.width', 1000)  
pd.set_option('display.float_format','{:.2f}'.format)  
  
random.seed(0)  
np.random.seed(0)  
np.set_printoptions(suppress=True)  
  
Autosaving every 60 seconds
```

Exploratory Data Analysis

```
In [2]:  
df = pd.read_csv("marketing_data.csv",parse_dates=['Dt_Customer'])  
  
In [3]:  
df  
  
Out[3]:  
ID Year_Birth Education Marital_Status Income Kidhome Teenhome Dt_Customer Recency MntWines MntFruits MntMeatProd  
0 1826 1970 Graduation Divorced 84835.00 0 0 2014-06-16 0 189 104  
1 1 1961 Graduation Single 57091.00 0 0 2014-06-15 0 464 5  
2 10476 1958 Graduation Married 67267.00 0 1 2014-05-13 0 134 11  
3 1386 1967 Graduation Together 32474.00 1 1 2014-11-05 0 10 0  
4 5371 1989 Graduation Single 21474.00 1 0 2014-08-04 0 6 16  
... ... ... ... ... ... ... ... ... ... ... ... ...  
2235 10142 1976 PhD Divorced 66476.00 0 1 2013-07-03 99 372 18  
2236 5263 1977 2n Cycle Married 31056.00 1 0 2013-01-22 99 5 10  
2237 22 1976 Graduation Divorced 46310.00 1 0 2012-03-12 99 185 2  
2238 528 1978 Graduation Married 65819.00 0 0 2012-11-29 99 267 38  
2239 4070 1969 PhD Married 94871.00 0 2 2012-01-09 99 169 24  
2240 rows × 28 columns  
  
In [4]:  
df.info()  
  
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 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          2240 non-null   float64  
 5   Kidhome         2240 non-null   int64    
 6   Teenhome        2240 non-null   int64    
 7   Dt_Customer     2240 non-null   datetime64[ns]  
 8   Recency         2240 non-null   int64    
 9   MntWines        2240 non-null   int64    
 10  MntFruits       2240 non-null   int64    
 11  MntMeatProd     2240 non-null   int64    
 12  MntFruitProducts 2240 non-null   int64    
 13  MntSweetProd   2240 non-null   int64    
 14  MntGoldProd    2240 non-null   int64    
 15  NumWebPurchases 2240 non-null   int64    
 16  NumTelPurchases 2240 non-null   int64    
 17  NumCatalogPurchases 2240 non-null   int64    
 18  NumWebVisitsMonth 2240 non-null   int64    
 19  NumWebVisitsMonth 2240 non-null   int64    
 20  AcceptedCmp1    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  AcceptedCmp3    2240 non-null   int64    
 26  AcceptedCmp4    2240 non-null   int64    
 27  Country          2240 non-null   object    
dtypes: datetime64[ns](1), float64(1), int64(23), object(3)  
memory usage: 490.1+ KB
```

```
In [5]:  
df.describe(include='all')
```

```
Out[5]:  
ID Year_Birth Education Marital_Status Income Kidhome Teenhome Dt_Customer Recency MntWines MntFruits MntMeatProd  
count 2240.00 2240.00 2240.00 2240.00 2240.00 2240.00 2240.00 2240.00 2240.00 2240.00 2240.00 2240.00  
unique NaN NaN NaN 5 8 NaN NaN NaN 663 NaN NaN NaN  
top NaN NaN NaN Graduation Married NaN NaN NaN 2012-08-31 00:00:00 NaN NaN NaN  
freq NaN NaN NaN 1127 864 NaN NaN NaN 2012-01-08 00:00:00 NaN NaN NaN  
first NaN NaN NaN NaN NaN NaN NaN NaN 2014-12-06 00:00:00 NaN NaN NaN  
last NaN NaN NaN NaN NaN NaN NaN NaN 2014-12-06 00:00:00 NaN NaN NaN  
mean 5592.16 1968.81 NaN 35303.00 0.44 0.51 9.11 303.94 26.30  
std 3246.66 11.98 NaN 25173.08 0.54 0.54 9.11 28.96 33.60 39.77  
min 0.00 1893.00 NaN 1730.00 0.00 0.00 NaN 0.00 0.00 0.00 0.00  
25% 5282.25 1959.00 NaN 35303.00 0.00 0.00 NaN 24.00 23.75 1.00  
50% 5458.50 1970.00 NaN 51381.50 0.00 0.00 NaN 49.00 173.50 8.00  
75% 5427.75 1977.00 NaN 68522.00 1.00 1.00 NaN 74.00 504.25 33.00  
max 11191.00 1996.00 NaN 66666.00 2.00 2.00 NaN 99.00 1493.00 199.00
```

```
In [6]:  
df.isnull().sum()
```

```
Out[6]:  
ID Year_Birth Education Marital_Status Income Kidhome Teenhome Dt_Customer Recency MntWines MntFruits MntMeatProd  
Year_Birth 0  
Education 0  
Marital_Status 0  
Income 24  
Kidhome 0  
Teenhome 0  
Dt_Customer 0  
Recency 0  
MntWines 0  
MntFruits 0  
MntMeatProd 0  
MntFruitProducts 0  
MntGoldProd 0  
NumWebPurchases 0  
NumTelPurchases 0  
NumCatalogPurchases 0  
NumWebVisitsMonth 0  
AcceptedCmp3 0  
AcceptedCmp4 0  
AcceptedCmp5 0  
AcceptedCmp1 0  
AcceptedCmp2 0  
Response 0  
Complain 0  
Country 0  
dtype: int64
```

```
In [7]:  
df.dropna(inplace=True)
```

```
In [8]:  
df.head()
```

```
Out[8]:  
ID Year_Birth Education Marital_Status Income Kidhome Teenhome Dt_Customer Recency MntWines MntFruits MntMeatProd  
0 1826 1970 Graduation Divorced 84835.00 0 0 2014-06-16 0 189 104 37  
1 1 1961 Graduation Single 57091.00 0 0 2014-06-15 0 464 5 5  
2 10476 1958 Graduation Married 67267.00 0 1 2014-05-13 0 134 11 11  
3 1386 1967 Graduation Together 32474.00 1 1 2014-11-05 0 10 0 0  
4 5371 1989 Graduation Single 21474.00 1 0 2014-08-04 0 6 16 16
```

```
In [9]:  
df["Age"] = 2022 - df["Year_Birth"]
```

```
In [10]:  
df.head()
```

```
Out[10]:  
ID Year_Birth Education Marital_Status Income Kidhome Teenhome Dt_Customer Recency MntWines MntFruits MntMeatProd  
0 1826 1970 Graduation Divorced 84835.00 0 0 2014-06-16 0 189 104 37  
1 1 1961 Graduation Single 57091.00 0 0 2014-06-15 0 464 5 5  
2 10476 1958 Graduation Married 67267.00 0 1 2014-05-13 0 134 11 11  
3 1386 1967 Graduation Together 32474.00 1 1 2014-11-05 0 10 0 0  
4 5371 1989 Graduation Single 21474.00 1 0 2014-08-04 0 6 16 16
```

```
In [11]:  
df.NumDealsPurchases.unique()
```

```
Out[11]:  
array([1, 0, 1, 2, 3, 0, 4, 12, 7, 6, 4, 8, 3, 9, 0, 17, 13, 10, 14, 19, 20], dtype=int64)
```

```
In [12]:  
df.columns
```

```
Out[12]:  
Index(['ID', 'Year_Birth', 'Education', 'Marital_Status', 'Income', 'Kidhome', 'Teenhome', 'Dt_Customer', 'Recency', 'MntWines', 'MntFruits', 'MntMeatProd', 'MntFruitProducts', 'MntGoldProd', 'NumWebPurchases', 'NumTelPurchases', 'NumCatalogPurchases', 'NumWebVisitsMonth', 'AcceptedCmp1', 'AcceptedCmp2', 'AcceptedCmp3', 'AcceptedCmp4', 'AcceptedCmp5', 'Response', 'Complain', 'Country', 'Age'], dtype='object')
```

```
In [13]:  
plt.figure(figsize=(10,5))  
sns.barplot(x=df.groupby(["Education"])["NumWebPurchases"].sum())
```

```
Out[13]:  
NumWebPurchases  
Education  
2n Cycle 753  
Basic 102  
Graduation 4593  
Master 1473  
PhD 2132
```

```
In [14]:  
plt.figure(figsize=(10,5))  
sns.barplot(x=df.groupby(["Education Level"])["NumWebPurchases"], size=20)  
plt.xlabel("Education Level")  
plt.ylabel("Web Purchases Numbers")  
plt.show()
```

```
Web Purchases by Education Level
```



```
In [15]:  
lineplot = pd.DataFrame(df.groupby(["Dt_Customer"])["NumWebPurchases"].sum())
```

```
Out[15]:  
NumWebPurchases  
Dt_Customer  
2012-01-08 17  
2012-01-09 8  
2012-01-10 22  
2012-01-11 6  
2012-01-12 19  
2014-12-02 13  
2014-12-03 21  
2014-12-04 2  
2014-12-05 41  
2014-12-06 0  
662 rows × 1 columns
```

```
In [16]:  
fig = px.line(data_frame=lineplot, x=lineplot.index, y="NumWebPurchases",  
               labels={"x": "Dt_Customer", "y": "Dates",  
                      "color": "NumWebPurchases"},  
               title="Web Purchases by date")  
fig.show()
```

```
Web Purchases by date
```



```
In [17]:  
scatter = df[["Income", "NumWebPurchases"]]
```

```
Out[17]:  
Income NumWebPurchases  
0 84835.00 4  
1 57091.00 7  
2 67267.00 3  
3 32474.00 1  
4 21474.00 3  
... ... ...  
2235 66476.00 5  
2236 31056.00 1  
2237 46310.00 6  
2238 65819.00 5  
2239 94871.00 8
```

```
2216 rows × 2 columns
```

```
In [18]:  
plt.figure(figsize=(10,5))  
sns.scatterplot(x=df["Income"], y=df["NumWebPurchases"], data=scatter)
```

```
Out[18]:  
Scatter plot showing the relationship between Income and NumWebPurchases. The X-axis is Income and the Y-axis is NumWebPurchases.
```

```
In [19]:  
plt.title("Web Purchases by Income", size=20)
```

```
Out[19]:  
Web Purchases by Income
```



```
In [20]:  
df.NumDealsPurchases.unique()
```

```
Out[20]:  
array([1, 2, 3, 0, 4, 12, 7, 6, 4, 8, 3, 9, 0, 17, 13, 10, 14, 19, 20], dtype=int64)
```

```
In [21]:  
g = sns.catplot(x="Kidhome", col = "Marital_Status", col_wrap=4,  
                 kind="bar", data=df, height = 5, aspect = 1)  
g.set_x_labels("Number of Kids")  
g.set_y_labels("Count")  
#g.set_axis_labels("Tip", "Total bill(USD)").set_xlim(0,100).set_ylim(0,100)  
g.set_xlim(0, None)  
g.set_xticklabels(rotation=0)  
plt.suptitle("Number of Kids by Marital Status", x=0.5, y=1.1, ha="center", fontsize=20)  
plt.show()
```

```
<Figure size 1440x1440 with 0 Axes>
```

```
Number of Kids by Marital Status
```


