

## Data Collection and Preprocessing Phase

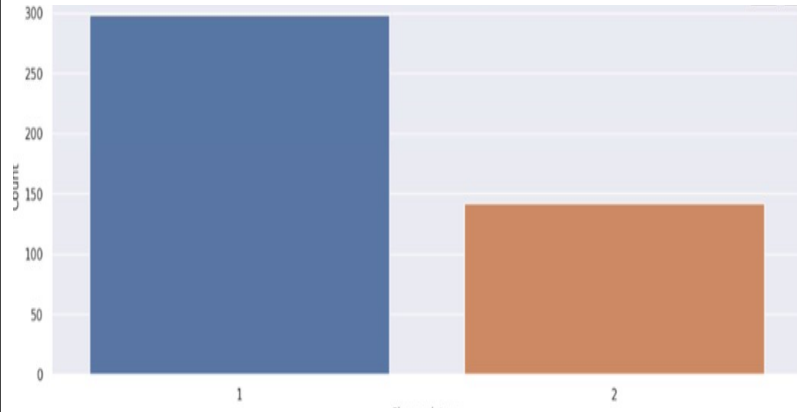
Date	July 5, 2024
Team ID	739683
Project Title	Customer Segmentation using Machine Learning
Maximum Marks	6 Marks

### Data Exploration and Preprocessing Template

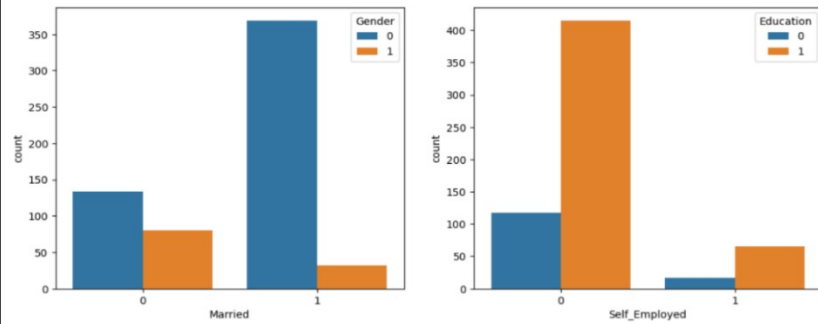
Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

Section	Description
Data Overview	<div> <div>Sex</div> <div>Marital status</div> <div>Age</div> <div>Education</div> <div>Income</div> <div>Occupation</div> <div>Settlement size</div> </div>
	<div>count</div> <div>2000.000000</div> <div>2000.000000</div> <div>2000.000000</div> <div>2000.000000</div> <div>2000.000000</div> <div>2000.000000</div> <div>2000.000000</div>
	<div>mean</div> <div>0.457000</div> <div>0.496500</div> <div>35.909000</div> <div>1.03800</div> <div>120954.419000</div> <div>0.810500</div> <div>0.739000</div>
	<div>std</div> <div>0.498272</div> <div>0.500113</div> <div>11.719402</div> <div>0.59978</div> <div>38108.824679</div> <div>0.638587</div> <div>0.812533</div>
	<div>min</div> <div>0.000000</div> <div>0.000000</div> <div>18.000000</div> <div>0.00000</div> <div>35832.000000</div> <div>0.000000</div> <div>0.000000</div>
	<div>25%</div> <div>0.000000</div> <div>0.000000</div> <div>27.000000</div> <div>1.00000</div> <div>97663.250000</div> <div>0.000000</div> <div>0.000000</div>
	<div>50%</div> <div>0.000000</div> <div>0.000000</div> <div>33.000000</div> <div>1.00000</div> <div>115548.500000</div> <div>1.000000</div> <div>1.000000</div>
	<div>75%</div> <div>1.000000</div> <div>1.000000</div> <div>42.000000</div> <div>1.00000</div> <div>138072.250000</div> <div>1.000000</div> <div>1.000000</div>
	<div>max</div> <div>1.000000</div> <div>1.000000</div> <div>76.000000</div> <div>3.00000</div> <div>309364.000000</div> <div>2.000000</div> <div>2.000000</div>

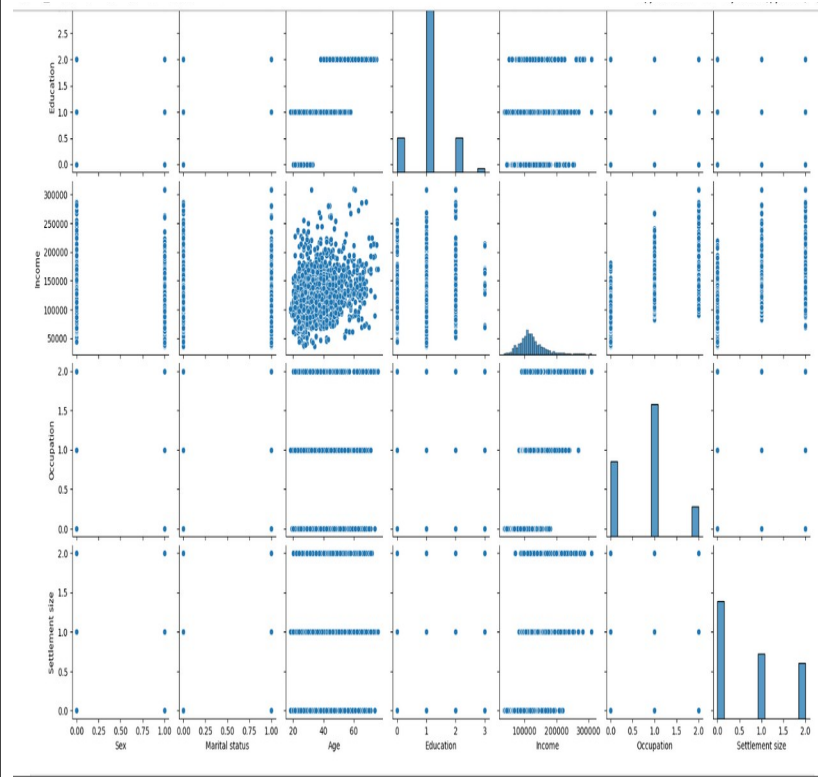
## Univariate Analysis



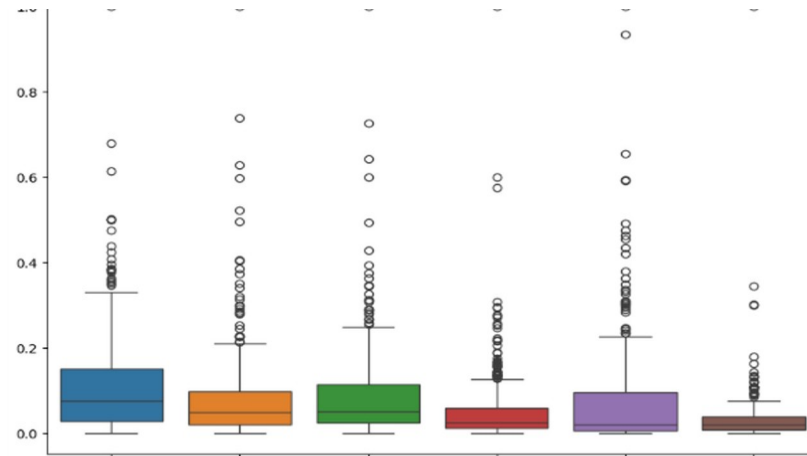
## Bivariate Analysis



## Multivariate Analysis



## Outliers and Anomalies



## Data Preprocessing Code Screenshots

### Loading Data

```
3]: os.chdir(r"C:/Users/kusur/Downloads")
data = pd.read_csv('segmentation data.csv', header='infer')
data.head()

3]:
```

	ID	Sex	Marital status	Age	Education	Income	Occupation	Settlement size
0	1000000001	0	0	67	2	124670	1	2
1	1000000002	1	1	22	1	150773	1	2
2	1000000003	0	0	49	1	89210	0	0
3	1000000004	0	0	45	1	171565	1	1
4	1000000005	0	0	53	1	149031	1	1

### Handling Missing Data

```
: data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2000 entries, 0 to 1999
Data columns (total 7 columns):
#   Column              Non-Null Count  Dtype
---  ---             
0   Sex                  2000 non-null   int64
1   Marital status       2000 non-null   int64
2   Age                  2000 non-null   int64
3   Education             2000 non-null   int64
4   Income                2000 non-null   int64
5   Occupation            2000 non-null   int64
6   Settlement size       2000 non-null   int64
dtypes: int64(7)
memory usage: 109.5 KB
```

Data Transformation	<pre>data = minmax_scale(data,feature_range=(0,1))  import pickle  pickle.dump(data,open("scale.pk2", 'wb'))  names = ['Sex','Marital status','Age','Education','Income','Occupation','Settlement size'] data = pd.DataFrame(data,columns=names)  wcss = [] for i in range(1, 11):     kmeans = sk.cluster.KMeans(n_clusters=i, init='k-means++', random_state=0)     kmeans.fit(data)     wcss.append(kmeans.inertia_)</pre>
Feature Engineering	Attached the codes in final submission
Save Processed Data	-