

## - Import libraries and dataset

### Business Problem

The Management team at Walmart Inc. wants to analyze the customer purchase behavior (specifically, purchase amount) against the customer's gender and the various other factors to help the business make better decisions. They want to understand if the spending habits differ between male and female customers: Do women spend more on Black Friday than men? (Assume 50 million customers are male and 50 million are female).

### What good looks like?

- Import the dataset and do usual data analysis steps like checking the structure & characteristics of the dataset.
- ok
- Detect Null values & Outliers (using boxplot, "describe" method by checking the difference between mean and median, isnan etc.)
- ok
- Do some data exploration steps like:
  - Tracking the amount spent per transaction of all the 50 million female customers, and all the 50 million male customers, calculate the average, and conclude the results.
  - ok
  - Inference after computing the average female and male expenses.
  - ok
  - Use the sample average to find out an interval within which the population average will lie. Using the sample of female customers you will calculate the interval within which the average spending of 50 million male and female customers may lie.
  - ok
  - Use the Central limit theorem to compute the interval. Change the sample size to observe the distribution of the mean of the expenses by female and male customers.
  - ok
  - The interval that you calculated is called Confidence Interval. The width of the interval is mostly decided by the business: Typically 90%, 95%, or 99%. Play around with the width parameter and report the observations.
  - ok
  - Conclude the results and check if the confidence intervals of average male and female spends are overlapping or not overlapping. How can Walmart leverage this conclusion to make changes or improvements?
  - ok
  - Perform the same activity for Married vs Unmarried and Age
  - ok
  - For Age, you can try bins based on life stages: 0-17, 18-25, 26-35, 36-50, 51+ years.
  - ok
  - Give recommendations and action items to Walmart.
  - ok

## - Analysing basic metrics

- Observations on shape of data, data types of all the attributes, conversion of categorical attributes to 'category' (If required), statistical summary

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from scipy.stats import norm
6 df = pd.read_csv('/content/walmart_data.txt')

1 df
```

|   | User_ID | Product_ID | Gender | Age  | Occupation | City_Category | Stay_In_Current_City_Years | Marital_Status |
|---|---------|------------|--------|------|------------|---------------|----------------------------|----------------|
| 0 | 1000001 | P00069042  | F      | 0-17 | 10         | A             |                            | 2              |
| 1 | 1000001 | P00248942  | F      | 0-17 | 10         | A             |                            | 2              |

```
1 type(df)
2 #type of dataset is Pandas DataFrame
```

```
pandas.core.frame.DataFrame
```

```
1 df.shape
2 #(rows, columns)
```

```
(550068, 10)
```

```
1 df.info()
```

| # | Column                     | Non-Null Count | Dtype  |
|---|----------------------------|----------------|--------|
| 0 | User_ID                    | 550068         | int64  |
| 1 | Product_ID                 | 550068         | object |
| 2 | Gender                     | 550068         | object |
| 3 | Age                        | 550068         | object |
| 4 | Occupation                 | 550068         | int64  |
| 5 | City_Category              | 550068         | object |
| 6 | Stay_In_Current_City_Years | 550068         | object |
| 7 | Marital_Status             | 550068         | int64  |
| 8 | Product_Category           | 550068         | int64  |
| 9 | Purchase                   | 550068         | int64  |

dtypes: int64(5), object(5)  
memory usage: 42.0+ MB

```
1 columns=['Occupation', 'Marital_Status', 'Product_Category']
2 df[columns]=df[columns].astype('object')
3
4 # Now we are converting categoric columns from int datatype to object datatype.
5 # That will be helpful, when we will analyse description for categoric columns.
```

## Non Graphical Analysis

```
1 df.nunique()
2
3 # no column has count of unique values same as len of all rows,
4 # so we have the only default index values as unique identifier.
5 # Column 'Purchase' is continuous and we kept it as int datatype,
6 # Columns like 'Gender', 'Age', 'Occupation', 'City_Category', 'Marital_Status', 'Product_Category' are categorical.
```

|                            |              |
|----------------------------|--------------|
| User_ID                    | 5891         |
| Product_ID                 | 3631         |
| Gender                     | 2            |
| Age                        | 7            |
| Occupation                 | 21           |
| City_Category              | 3            |
| Stay_In_Current_City_Years | 5            |
| Marital_Status             | 2            |
| Product_Category           | 20           |
| Purchase                   | 18105        |
|                            | dtype: int64 |

```
1 # Checking name and count of unique attributes in categoric variables
2
3 for col in df.columns:
4     print('{} :{} = {} '.format(col,df[col].unique(),df[col].nunique()))
5 print()
6
7 # There are 20 product categories in total.
8 # There are 21 different types of occupations in the city.
```

|   |         |
|---|---------|
| User_ID :[1000001 1000002 1000003 ... 1004113 1005391 1001529]                            | = 5891  |
| Product_ID :['P00069042' 'P00248942' 'P00087842' ... 'P00370293' 'P00371644' 'P00370853'] | = 3631  |
| Gender :['F' 'M']   | = 2     |
| Age :['0-17' '55+' '26-35' '46-50' '51-55' '36-45' '18-25']                               | = 7     |
| Occupation :[10 16 15 7 20 9 1 12 17 0 3 4 11 8 19 2 18 5 14 13 6]                        | = 21    |
| City_Category :['A' 'C' 'B']  | = 3     |
| Stay_In_Current_City_Years :['2' '4+' '3' '1' '0']  | = 5     |
| Marital_Status :[0 1]   | = 2     |
| Product_Category :[3 1 12 8 5 4 2 6 14 11 13 15 7 16 18 10 17 9 20 19]                    | = 20    |
| Purchase :[ 8370 15200 1422 ... 135 123 613]  | = 18105 |

```

1 #Checking how the data is spread on basis of distinct users (customer analysis)
2
3 df2=df.groupby(['User_ID']).first()
4 df2
5 categ_cols = ['Gender', 'Age', 'City_Category', 'Stay_In_Current_City_Years', 'Marital_Status']
6 cat_count = df2[categ_cols].melt().groupby(['variable', 'value'])[['value']].size().reset_index(name='counts')
7 s = df2[categ_cols].melt().variable.value_counts()
8 cat_count['Percent'] = cat_count['counts'].div(cat_count['variable'].map(s)).mul(100).round().astype('int')
9 cat_count.groupby(['variable', 'value']).first()
10
11 # 35% customers aged between the age 26-35
12 # 73% customers aged between the age 18-45
13 # 72% customers are Male and 28% are Female
14 # 58% customers are Single, 40% are Married
15 # 35% customers staying in city from 1 year

```

|                                   |              | counts | Percent |
|-----------------------------------|--------------|--------|---------|
|                                   | variable     | value  |         |
| <b>Age</b>                        | <b>0-17</b>  | 218    | 4       |
|                                   | <b>18-25</b> | 1069   | 18      |
|                                   | <b>26-35</b> | 2053   | 35      |
|                                   | <b>36-45</b> | 1167   | 20      |
|                                   | <b>46-50</b> | 531    | 9       |
|                                   | <b>51-55</b> | 481    | 8       |
|                                   | <b>55+</b>   | 372    | 6       |
| <b>City_Category</b>              | <b>A</b>     | 1045   | 18      |
|                                   | <b>B</b>     | 1707   | 29      |
|                                   | <b>C</b>     | 3139   | 53      |
| <b>Gender</b>                     | <b>F</b>     | 1666   | 28      |
|                                   | <b>M</b>     | 4225   | 72      |
| <b>Marital_Status</b>             | <b>0</b>     | 3417   | 58      |
|                                   | <b>1</b>     | 2474   | 42      |
| <b>Stay_In_Current_City_Years</b> | <b>0</b>     | 772    | 13      |
|                                   | <b>1</b>     | 2086   | 35      |
|                                   | <b>2</b>     | 1145   | 19      |
|                                   | <b>3</b>     | 979    | 17      |
|                                   | <b>4+</b>    | 909    | 15      |

```

1 #Checking how the data is spread on basis of distinct transection
2
3 df2=df.copy()
4 df2
5 categ_cols = ['Gender', 'Age', 'City_Category', 'Stay_In_Current_City_Years', 'Marital_Status']
6 cat_count = df2[categ_cols].melt().groupby(['variable', 'value'])[['value']].size().reset_index(name='counts')
7 s = df2[categ_cols].melt().variable.value_counts()
8 cat_count['Percent'] = cat_count['counts'].div(cat_count['variable'].map(s)).mul(100).round().astype('int')
9 cat_count.groupby(['variable', 'value']).first()
10
11 # There are 35% customers are in age group 26-35 but they have purchased 40% products
12 # There are 72% customers are male but they have purchased 75% products

```

|          |       | counts | Percent |
|----------|-------|--------|---------|
| variable | value |        |         |
| Age      | 0-17  | 15102  | 3       |
|          | 18-25 | 99660  | 18      |
|          | 26-35 | 219587 | 40      |
|          | 36-45 | 110013 | 20      |
|          | 46-50 | 45701  | 8       |
|          | 51-55 | 38501  | 7       |

```

1 # Checking how product category contributes to the entire data (product analysis)
2
3 categ_cols = ['Product_Category']
4 cat_count = df[categ_cols].melt().groupby(['variable', 'value'])[['value']].size().reset_index(name='counts')
5 s = df[categ_cols].melt().variable.value_counts()
6 cat_count['Percent'] = cat_count['counts'].div(cat_count['variable'].map(s)).mul(100).round().astype('int')
7 cat_count.groupby(['variable', 'value']).first().sort_values('Percent', ascending=False)
8
9 # 27% sold products belong to category 5
10 # 26% sold products belong to category 1
11 # 21% sold products belong to category 8
12 # product category 9 and 17 are least sold.

```

|                  |       | counts | Percent |
|------------------|-------|--------|---------|
| variable         | value |        |         |
| Product_Category | 5     | 150933 | 27      |
|                  | 1     | 140378 | 26      |
|                  | 8     | 113925 | 21      |
|                  | 6     | 20466  | 4       |
|                  | 2     | 23864  | 4       |
|                  | 11    | 24287  | 4       |
|                  | 3     | 20213  | 4       |
|                  | 4     | 11753  | 2       |
|                  | 16    | 9828   | 2       |
|                  | 7     | 3721   | 1       |
|                  | 10    | 5125   | 1       |
|                  | 12    | 3947   | 1       |
|                  | 13    | 5549   | 1       |
|                  | 15    | 6290   | 1       |
|                  | 18    | 3125   | 1       |
|                  | 9     | 410    | 0       |
|                  | 14    | 1523   | 0       |
|                  | 17    | 578    | 0       |
|                  | 19    | 1603   | 0       |
|                  | 20    | 2550   | 0       |

```

1 # Checking how product category contributes to the entire data (product analysis)
2 # Checking how categorical variables contributes to the entire data
3 categ_cols = ['Product_ID']
4 cat_count = df[categ_cols].melt().groupby(['variable', 'value'])[['value']].count()
5 cat_count.rename(columns={'value':'counts'}, inplace=True)
6 cat_count.sort_values('counts', ascending=False)
7 # 3631 types of products are present in data
8 # P00265242 is most purchased product

```

**counts**

| variable   | value     |
|------------|-----------|
| Product_ID | P00265242 |
|            | 1880      |
|            | P00025442 |
|            | 1615      |
|            | P00110742 |
|            | 1612      |

```
1 # Users with highest number of purchases
2 df.groupby(['User_ID'])['Purchase'].count().nlargest(10)
```

```
User_ID
1001680    1026
1004277    979
1001941    898
1001181    862
1000889    823
1003618    767
1001150    752
1001015    740
1005795    729
1005831    727
Name: Purchase, dtype: int64
```

```
1 #Users with highest purchases amount
2 df.groupby(['User_ID'])['Purchase'].sum().nlargest(10)
```

```
User_ID
1004277    10536909
1001680    8699596
1002909    7577756
1001941    6817493
1000424    6573609
1004448    6566245
1005831    6512433
1001015    6511314
1003391    6477160
1001181    6387961
Name: Purchase, dtype: int64
```

```
1 describe_obj=df.describe(include='object')
2 describe_obj.rename(index={'top':'Mode'}, inplace=True)
3 describe_obj.loc['Mode %'] = describe_obj.apply(lambda x: int(round(x[3]/x[0]*100)))
4 describe_obj.T
5 # Customer with User_ID 1001680 is most loyal customer.
6 # Product with Product_ID P00265242 is most popular product.
7 # Products from Product_Category 5 are most purchased.
8 # 40% of the purchase done by aged 26-35.
9 # 59% Single contributes to the purchase count.
10 # 35% Staying in the city from 1 year, 18% from 2 years, 17% from 3 years
11 # There are 20 product categories in total.
12 # There are 20 different types of occupations in the city.
```

|                            | count  | unique | Mode      | freq   | Mode % |
|----------------------------|--------|--------|-----------|--------|--------|
| User_ID                    | 550068 | 5891   | 1001680   | 1026   | 0      |
| Product_ID                 | 550068 | 3631   | P00265242 | 1880   | 0      |
| Gender                     | 550068 | 2      | M         | 414259 | 75     |
| Age                        | 550068 | 7      | 26-35     | 219587 | 40     |
| Occupation                 | 550068 | 21     | 4         | 72308  | 13     |
| City_Category              | 550068 | 3      | B         | 231173 | 42     |
| Stay_In_Current_City_Years | 550068 | 5      | 1         | 193821 | 35     |
| Marital_Status             | 550068 | 2      | 0         | 324731 | 59     |
| Product_Category           | 550068 | 20     | 5         | 150933 | 27     |

```
1 df.describe().T
2 #descriptive statistics of continuous variables - mean, std, quartile, median, range
3 #
```

|          | count    | mean        | std         | min  | 25%    | 50%    | 75%     | max     |
|----------|----------|-------------|-------------|------|--------|--------|---------|---------|
| Purchase | 550068.0 | 9263.968713 | 5023.065394 | 12.0 | 5823.0 | 8047.0 | 12054.0 | 23961.0 |

```
1 df3=df.groupby(['City_Category'])['Purchase'].sum()
2 df3
```

```
City_Category
A      1316471661
B      2115533605
C      1663807476
Name: Purchase, dtype: int64
```

```

1 df2=df.groupby(['User_ID'])['City_Category'].unique()
2 df2=pd.DataFrame(df2)
3 df2.reset_index().astype('str').City_Category.value_counts()/len(df2)

['C']    0.532847
['B']    0.289764
['A']    0.177389
Name: City_Category, dtype: float64

```

```

1 #Checking the age group distribution in different city categories
2 data_crosstab = pd.crosstab(index=df["City_Category"],columns=df["Age"],margins=True,normalize="columns")
3 data_crosstab.round(2)
4
5 # We have seen earlier that city category B and A constitutes less percentage of total population
6 # but they contribute more towards purchase counts for customers aged 26-35 for B(40%) and A (50%)
7 # which can be the reason for these city categories to be more actively purchasing

```

|               | Age | 0-17 | 18-25 | 26-35 | 36-45 | 46-50 | 51-55 | 55+  | All  |
|---------------|-----|------|-------|-------|-------|-------|-------|------|------|
| City_Category | A   | 0.17 | 0.28  | 0.34  | 0.24  | 0.17  | 0.16  | 0.17 | 0.27 |
|               | B   | 0.36 | 0.43  | 0.42  | 0.43  | 0.45  | 0.46  | 0.24 | 0.42 |
|               | C   | 0.47 | 0.29  | 0.25  | 0.33  | 0.39  | 0.38  | 0.59 | 0.31 |

```

1 data_crosstab = pd.crosstab(df['Marital_Status'],df['Gender'],margins = True, normalize="all")
2 data_crosstab.round(2)
3
4 #normalize over all values. Probability of single female customer is .17 over all customers.
5
6 # Marginal probability: Probability of customers who are single is = 0.59
7 # Conditional probabitity: Probability of customers who are single and male is =0.45

```

|                | Gender | F    | M    | All  |
|----------------|--------|------|------|------|
| Marital_Status | 0      | 0.14 | 0.45 | 0.59 |
|                | 1      | 0.10 | 0.31 | 0.41 |
|                | All    | 0.25 | 0.75 | 1.00 |

```

1 df.isna().sum()
2 #dataframe has no missing values in any column.
3 #no need for data cleaning or imputation.

```

```

User_ID              0
Product_ID           0
Gender               0
Age                  0
Occupation           0
City_Category         0
Stay_In_Current_City_Years  0
Marital_Status        0
Product_Category      0
Purchase              0
dtype: int64

```

```

1 df[df.duplicated()]
2
3 # No duplicate row present in dataset

```

```
User_ID  Product_ID  Gender  Age  Occupation  City_Category  Stay_In_Current_City_Years  Marital_Status
```

◀ ▶

## Visual Analysis

```

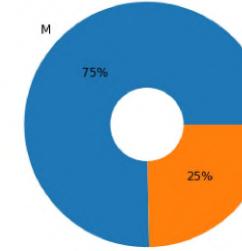
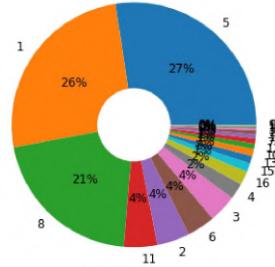
1 fig=plt.figure(figsize = [18,5])
2
3 plt.subplot(1,2,1)
4 x_bar = df.Product_Category.value_counts().index
5 y_bar = df.Product_Category.value_counts()
6 plt.pie(y_bar, labels=x_bar, autopct="%0.0f%%", textprops={"fontsize":12})
7 centre_circle= plt.Circle((0,0),0.3,color='black', fc ='white', linewidth=0, )
8 fig = plt.gcf()
9 fig.gca().add_artist(centre_circle)
10 plt.axis('equal')
11
12 plt.subplot(1,2,2)
13 x_bar = df.Gender.value_counts().index
14 y_bar = df.Gender.value_counts()

```

```

15 plt.pie(y_bar, labels = x_bar, autopct="%0.0f%%", textprops={"fontsize":12})
16 centre_circle= plt.Circle((0,0),0.3,color='black', fc ='white', linewidth=0, )
17 fig = plt.gcf()
18 fig.gca().add_artist(centre_circle)
19 plt.axis('equal')
20
21 plt.show()
22
23 # In dataset there are 20 Product categories.
24 # Product category 5,1,8 are most-popular product catagories.
25 # 27% of all products purchased from category 5.
26 # 26% of all products purchased from category 1.
27 # 21% of all products purchased from category 8.
28 # other categories are liked by less then 4% customers for each categories.
29
30 # 75% products purchased by male
31 # 25% products purchased by female

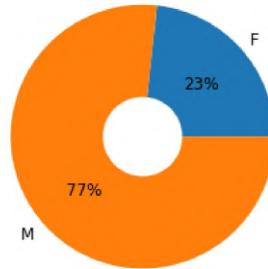
```



```

1
2 plt.show()
3
4 # Customers from City_Catagory (B) purchased most amount of products.
5 # Customers from Age group (26-35) purchased most amount of products.

```



```

1 top_cat_sale = df.groupby(['Product_Catagory'])['Purchase'].sum().head(3).sum()
2 total_sale = df['Purchase'].sum()
3 print("sale percentage :",round(top_cat_sale/total_sale,2)*100)
4
5 # nearly half of of total sales generated from only 3 product catagories out of 20 product catagories.

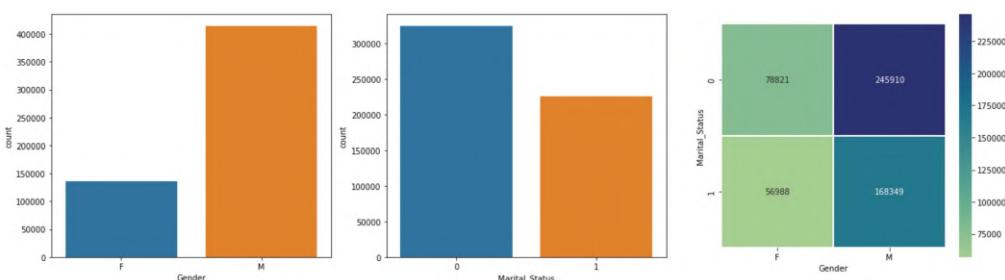
```

sale percentage : 47.0

```

1 plt.figure(figsize = [18,5])
2
3 plt.subplot(1, 3, 1)
4 sns.countplot(data=df, x="Gender")
5
6 plt.subplot(1, 3, 2)
7 sns.countplot(data=df, x="Marital_Status")
8
9 plt.subplot(1, 3, 3)
10 ct_counts = df.groupby(['Gender', 'Marital_Status']).size()
11 ct_counts = ct_counts.reset_index(name = 'count')
12 ct_counts = ct_counts.pivot(index = 'Marital_Status', columns = 'Gender', values = 'count')
13 sns.heatmap(ct_counts, cmap= "crest", annot = True, fmt = 'd', square=1, linewidth=1.)
14 plt.tight_layout()
15 plt.show()
16
17 # Major quantity of products are purchased by Single Male.
18 # Least quantity of products purchased by Married Female.

```



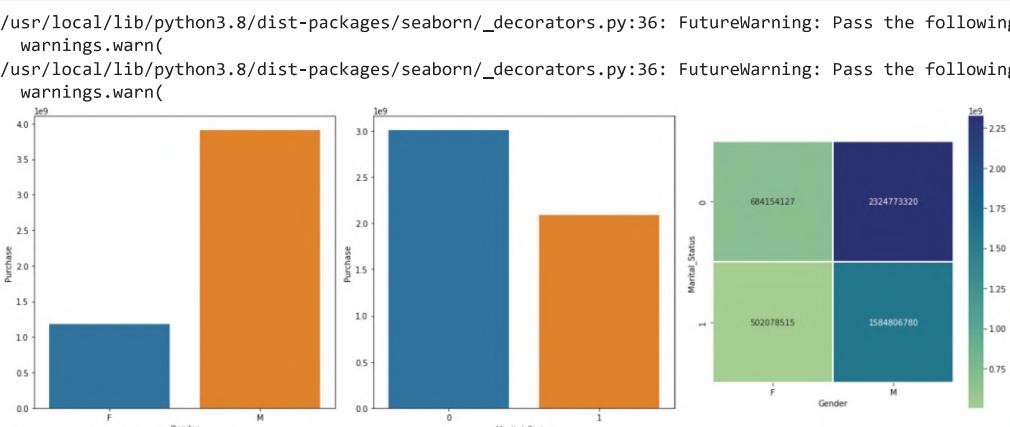
```
1 df.groupby(['Gender'])['Purchase'].count()
2 # Total products purchased by males is around 144 thousand whereas for females it's just 135 thousand.
```

```
Gender
F    135809
M    414259
Name: Purchase, dtype: int64
```

```
1 df.groupby(['Gender'])['Purchase'].sum()
2 # Total amount spent by males is around 4 billion whereas for females it's 1.2 billion.
```

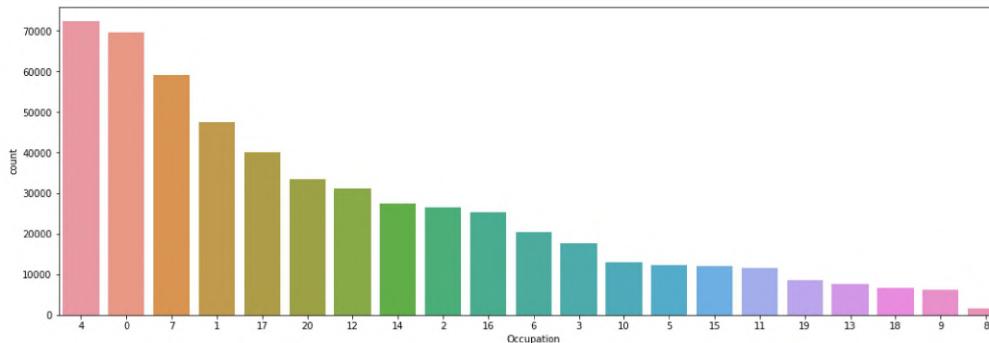
```
Gender
F    1186232642
M    3909580100
Name: Purchase, dtype: int64
```

```
1 df4= df.groupby(['Gender'])['Purchase'].sum().reset_index()
2 df5= df.groupby(['Marital_Status'])['Purchase'].sum().reset_index()
3 plt.figure(figsize=[18,6])
4 plt.subplot(1,3,1)
5 sns.barplot(df4['Gender'],df4['Purchase'])
6 plt.subplot(1,3,2)
7 sns.barplot(df5['Marital_Status'],df5['Purchase'])
8 plt.subplot(1, 3, 3)
9 ct_counts = df.groupby(['Gender', 'Marital_Status'])['Purchase'].sum()
10 ct_counts = ct_counts.reset_index(name = 'count')
11 ct_counts = ct_counts.pivot(index = 'Marital_Status', columns = 'Gender', values = 'count')
12 sns.heatmap(ct_counts, cmap= "crest", annot = True, fmt = 'd', square=1, linewidth=1.)
13 plt.tight_layout()
14 plt.show()
15 plt.show()
16
17 # Customers from City_Catagory (B) purchased most amount of products.
18 # Customers from Age group (26-35) purchased most amount of products.
```



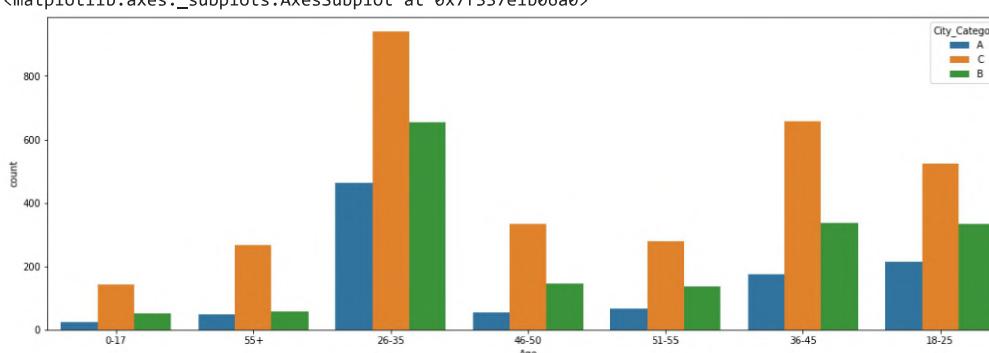
```
1 plt.figure(figsize=[18,6])
2 sns.countplot(df['Occupation'], order = df['Occupation'].value_counts().index)
3 plt.show()
4
5 # In dataset there are 21 occupations categories.
6 # Occupation category 4, 0, and 7 are with higher number of purchases.
7 # Occupation category 8 has least number of purchases.
```

```
/usr/local/lib/python3.8/dist-packages/seaborn/_decorators.py:36: FutureWarning: Pass the following
warnings.warn(
```



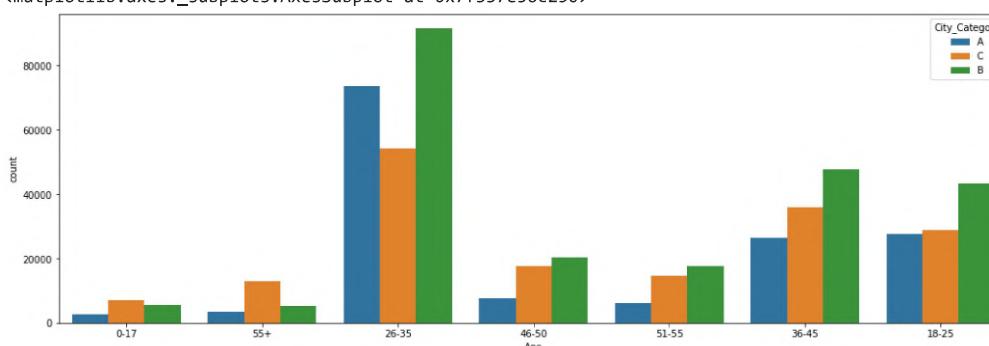
```
1 df2=df.groupby(['User_ID']).first()
2 plt.figure(figsize=[18,6])
3 sns.countplot(df2['Age'], hue=df2['City_Category'])
4
5 # for every age category there are least count of costomer those belong to city category A, except age category '26-35'
6 # In age category '26-35' the most count of customers belongs to city category C.
```

```
/usr/local/lib/python3.8/dist-packages/seaborn/_decorators.py:36: FutureWarning: Pass the following
warnings.warn(
<matplotlib.axes._subplots.AxesSubplot at 0x7f337e1b06a0>
```



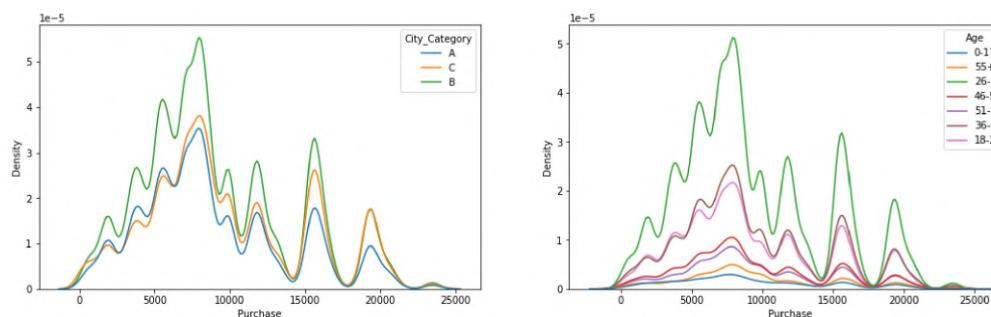
```
1 plt.figure(figsize=[18,6])
2 sns.countplot(df['Age'], hue=df['City_Category'])
3 # In age category '26-35' the most count of customers belongs to city category C.
4 # But they have least count count of purchase than city category A and B.
```

```
/usr/local/lib/python3.8/dist-packages/seaborn/_decorators.py:36: FutureWarning: Pass the following
warnings.warn(
<matplotlib.axes._subplots.AxesSubplot at 0x7f337e36c250>
```

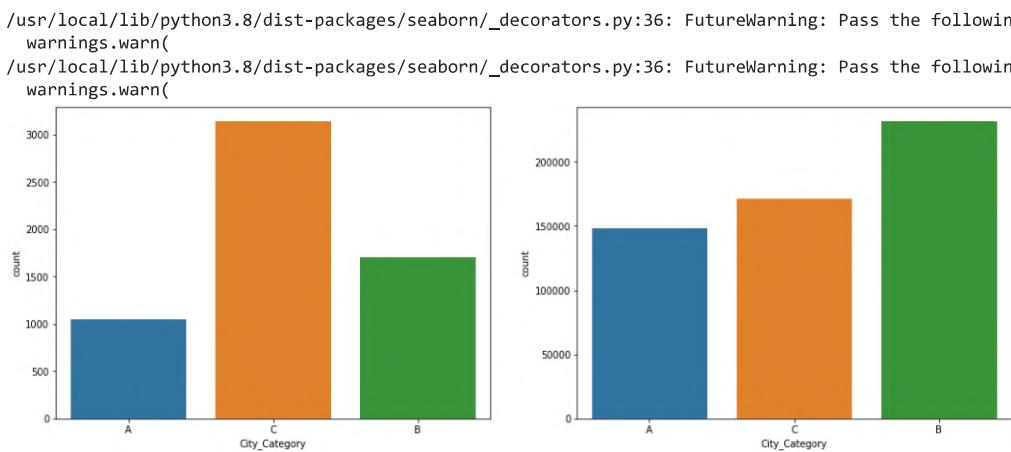


```
1 plt.figure(figsize=[18,5])
2 plt.subplot(1,2,1)
3 sns.kdeplot(df['Purchase'],hue=df['City_Category'])
4 plt.subplot(1,2,2)
5 sns.kdeplot(df['Purchase'],hue=df['Age'])
6 plt.show()
7
```

```
8 # There is not significant difference in purchase trend  
9 # City B and age group 26-35 have higher density because of higher customer counts from those sections.
```



```
1 df2=df.groupby(['User_ID']).first()  
2 plt.figure(figsize=[18,6])  
3 plt.subplot(1,2,1)  
4 sns.countplot(df2['City_Category'])  
5 plt.subplot(1,2,2)  
6 sns.countplot(df['City_Category'])  
7 plt.show()  
8  
9 # Mostly customers lives in City_Catagory C.  
10 # customers from City_Categroy (B) purchased most quantity of products.
```



```
1 df4= df.groupby(['City_Category'])['Purchase'].sum().reset_index()  
2 df5= df.groupby(['Age'])['Purchase'].sum().reset_index()  
3 plt.figure(figsize=[18,6])  
4 plt.subplot(1,2,1)  
5 sns.barplot(df4['City_Category'],df4['Purchase'])  
6 plt.subplot(1,2,2)  
7 sns.barplot(df5['Age'],df5['Purchase'])  
8 plt.show()  
9  
10 # Customers from City_Catagory (B) purchased most amount of products.  
11 # Customers from Age group (26-35) purchased most amount of products.
```

```

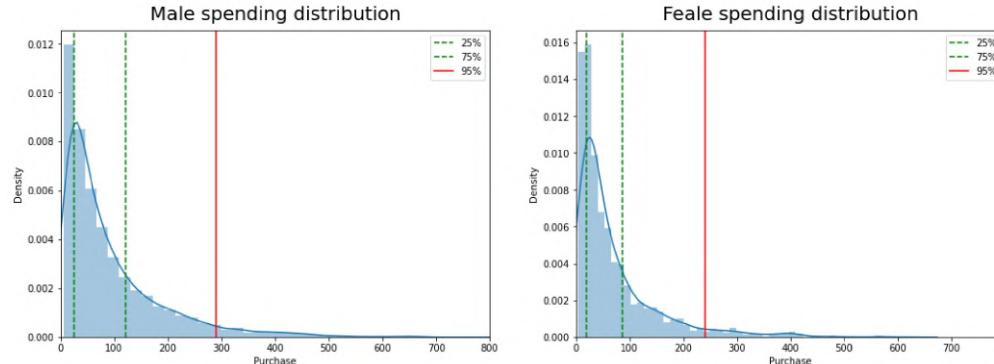
/usr/local/lib/python3.8/dist-packages/seaborn/_decorators.py:36: FutureWarning: Pass the following
warnings.warn(
/usr/local/lib/python3.8/dist-packages/seaborn/_decorators.py:36: FutureWarning: Pass the following
warnings.warn(
    1e9
  2.00
  1.75
  1e9
  2.00
  1.75
1 df2 = df.groupby(['User_ID', 'Gender'])[['Purchase']].sum().reset_index()
2
3 plt.figure(figsize=[18,6])
4 plt.subplot(1,2,1)
5 sns.distplot(df2[df2['Gender']=='M']['Purchase']/10000)
6 plt.axvline(np.percentile(df2[df2['Gender']=='M']['Purchase']/10000, 25), linestyle='--', color='green',label="25%")
7 plt.legend()
8 plt.axvline(np.percentile(df2[df2['Gender']=='M']['Purchase']/10000, 75), linestyle='--', color='green',label="75%")
9 plt.legend()
10 plt.axvline(np.percentile(df2[df2['Gender']=='M']['Purchase']/10000, 95), linestyle='-', color='red',label="95%")
11 plt.legend()
12 plt.title("Male spending distribution", y=1.015, fontsize=20)
13 plt.xlim(0,800)
14 plt.subplot(1,2,2)
15 sns.distplot(df2[df2['Gender']=='F']['Purchase']/10000)
16 plt.axvline(np.percentile(df2[df2['Gender']=='F']['Purchase']/10000, 25), linestyle='--', color='green',label="25%")
17 plt.legend()
18 plt.axvline(np.percentile(df2[df2['Gender']=='F']['Purchase']/10000, 75), linestyle='--', color='green',label="75%")
19 plt.legend()
20 plt.axvline(np.percentile(df2[df2['Gender']=='F']['Purchase']/10000, 95), linestyle='-', color='red',label="95%")
21 plt.legend()
22 plt.title("Female spending distribution", y=1.015, fontsize=20)
23 plt.xlim(0,800)
24 plt.show()
25
26 # distributions for Male spending and Female spending are highly right skewed (positive distribution).
27 # so in order to measure of variability or spread of the distribution,
28 # we are using first and third quartiles range, instead of standard deviation.
29 # There is more variance in Male spending distribution than female spending distribution.
30 # data points in Female spending distribution tend to be very close to the median.
31 # 5% of male customers purchased more than 2.8 million usd amount of products.
32 # 5% of female customers purchased more than 2.4 million usd amount of products.

```

```

/usr/local/lib/python3.8/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.8/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is
warnings.warn(msg, FutureWarning)

```



```

1 df2[df2['Gender']=='M']['Purchase'].skew(axis = 0),df2[df2['Gender']=='F']['Purchase'].skew(axis = 0)
(2.37378034880871, 2.4802021716417357)

```

```

1 (np.percentile(df2[df2['Gender']=='M']['Purchase'], 95),(np.percentile(df2[df2['Gender']=='F']['Purchase'], 95)))
(2898857.4, 2409077.75)

```

```

1 plt.figure(figsize=[18,14])
2
3 df2 = df.groupby(['User_ID', 'Gender'])[['Purchase']].sum().reset_index()
4 plt.subplot(2,2,1)
5 sns.boxplot(y=df2['Purchase'], x=df2['Gender'])
6
7 df2 = df.groupby(['User_ID', 'Marital_Status'])[['Purchase']].sum().reset_index()
8 plt.subplot(2,2,2)
9 sns.boxplot(y=df2['Purchase'], x=df2['Marital_Status'])
10

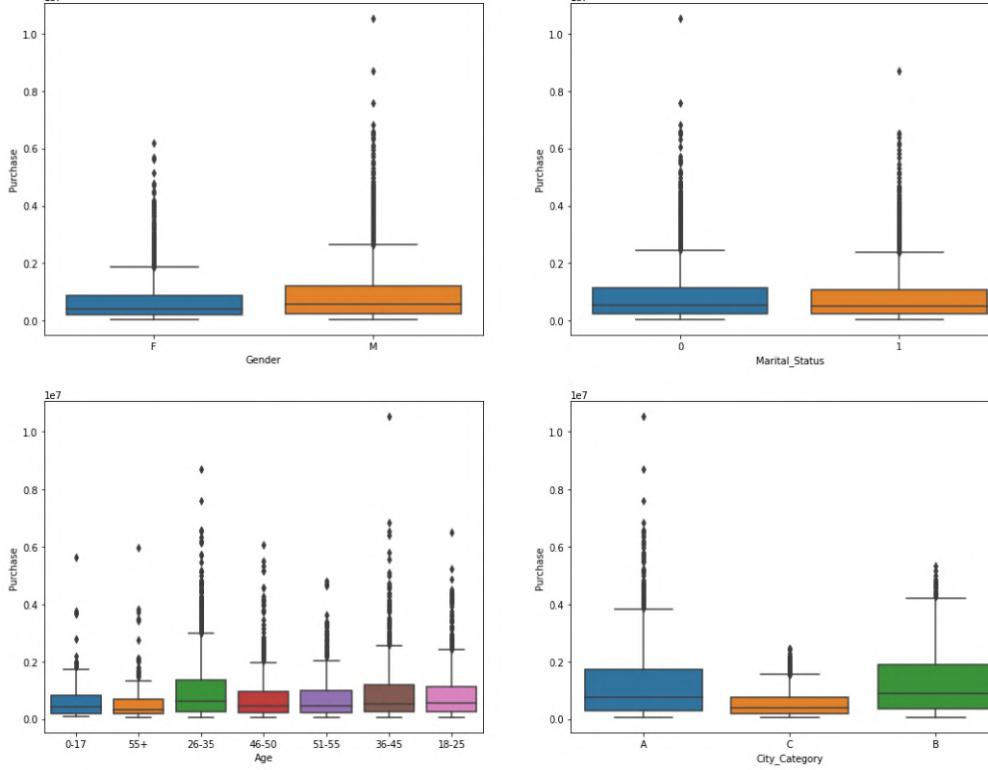
```

```

11 df2 = df.groupby(['User_ID', 'Age'])[['Purchase']].sum().reset_index()
12 plt.subplot(2,2,3)
13 sns.boxplot(y=df2['Purchase'], x=df2['Age'])
14
15 df2 = df.groupby(['User_ID', 'City_Category'])[['Purchase']].sum().reset_index()
16 plt.subplot(2,2,4)
17 sns.boxplot(y=df2['Purchase'], x=df2['City_Category'])
18
19 plt.suptitle('Outliers Detections (customer based)', y=.95, fontsize=20)
20 plt.show()
21
22 # There are more outliers in Male, few male customer purchased extreme amount of products.
23 # There are more outliers in Single Marital_Status, few single customer purchased extreme amount of products.

```

Outliers Detections (customer based)

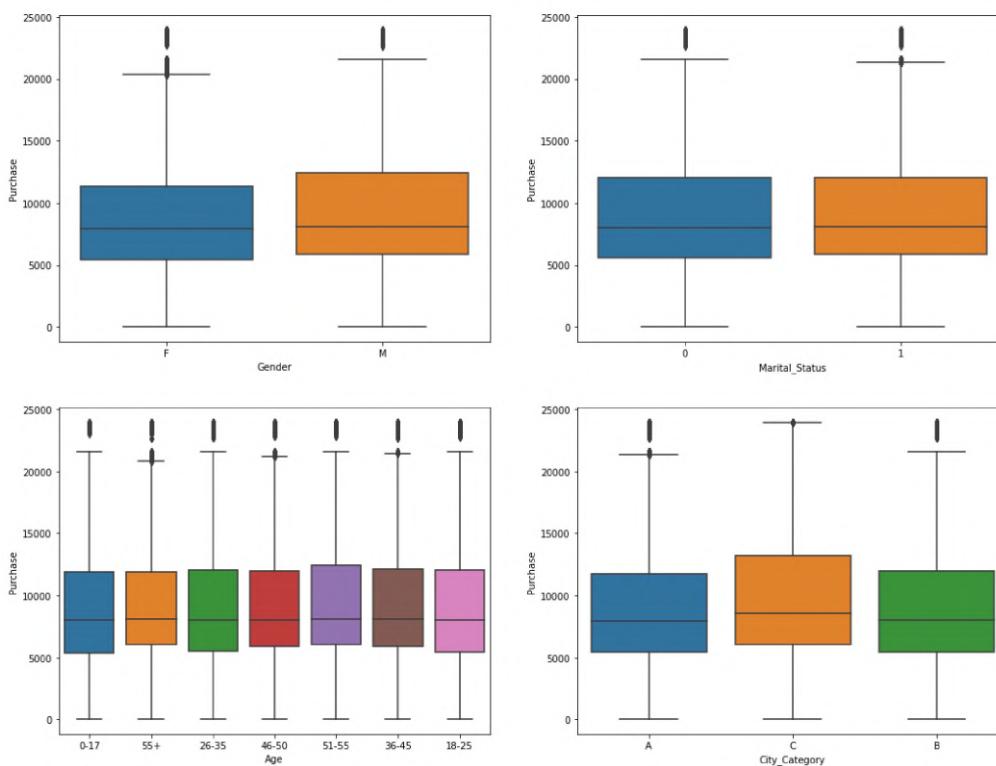


```

1 plt.figure(figsize=[18,14])
2 plt.subplot(2,2,1)
3 sns.boxplot(y=df['Purchase'], x=df['Gender'])
4
5 plt.subplot(2,2,2)
6 sns.boxplot(y=df['Purchase'], x=df['Marital_Status'])
7
8 plt.subplot(2,2,3)
9 sns.boxplot(y=df['Purchase'], x=df['Age'])
10
11 plt.subplot(2,2,4)
12 sns.boxplot(y=df['Purchase'], x=df['City_Category'])
13
14 plt.suptitle('Outliers Detections (product based)', y=.95, fontsize=20)
15 plt.show()
16
17 # Male customer purchase some more expansive products.
18 # Customers from City_Category ('C') purchase some more expansive products.

```

## Outliers Detections (product based)



```

1 def find_outliers_IQR(df):
2     q1=df['Purchase'].quantile(0.25)
3     q3=df['Purchase'].quantile(0.75)
4     IQR=q3-q1
5     outliers = df[df['Purchase']>(q3+1.5*IQR)]
6     return outliers.sort_values('Purchase',ascending=False)
7

```

```

1 find_outliers_IQR(df.groupby(['User_ID', 'Gender'])[['Purchase']].sum().reset_index())
2
3 # Below we can see the customers who purchased with extreme total amount of products.

```

|      | User_ID | Gender | Purchase |
|------|---------|--------|----------|
| 4166 | 1004277 | M      | 10536909 |
| 1634 | 1001680 | M      | 8699596  |
| 2831 | 1002909 | M      | 7577756  |
| 1885 | 1001941 | M      | 6817493  |
| 416  | 1000424 | M      | 6573609  |
| ...  | ...     | ...    | ...      |
| 3956 | 1004058 | M      | 2451245  |
| 5400 | 1005539 | F      | 2450068  |
| 4079 | 1004186 | M      | 2448401  |
| 3905 | 1004007 | M      | 2447282  |
| 1697 | 1001746 | M      | 2445649  |

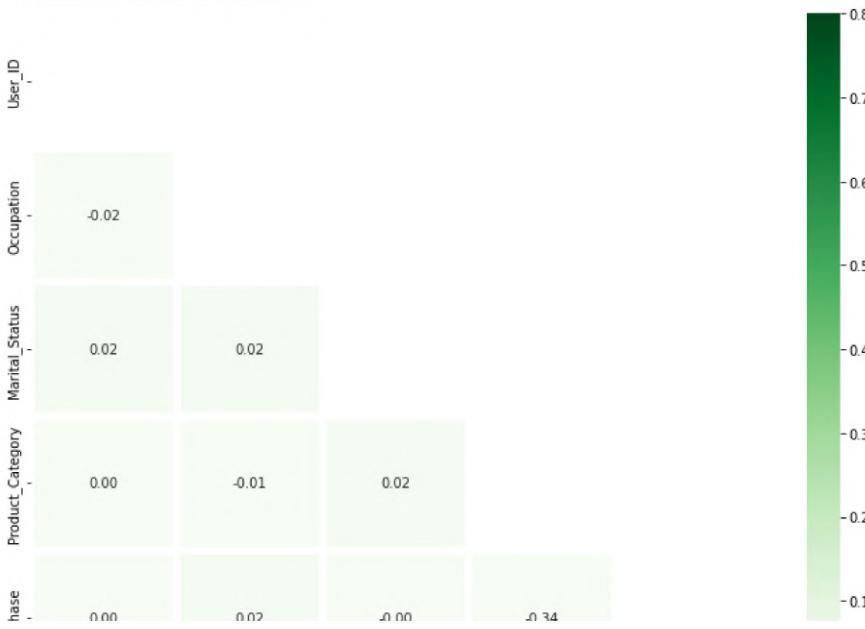
409 rows × 3 columns

```

1 plt.figure(figsize=[10,8])
2 dfc = pd.read_csv('/content/walmart_data.txt')
3 mask = np.triu(np.ones_like(dfc.corr()))
4 sns.heatmap(dfc.corr().round(2), cmap= "Greens", annot=True, mask=mask, fmt=".2f", linewidths=5, vmin=0, vmax=0.8)
5 plt.title('correlation heatmap', color='w', fontsize=20, fontweight = 'normal', backgroundcolor = 'g', pad = 20, loc='left')
6 plt.tight_layout()
7
8 # From the correlation plot, we can see the correlation is not significant between any pair of variables.

```

## correlation heatmap



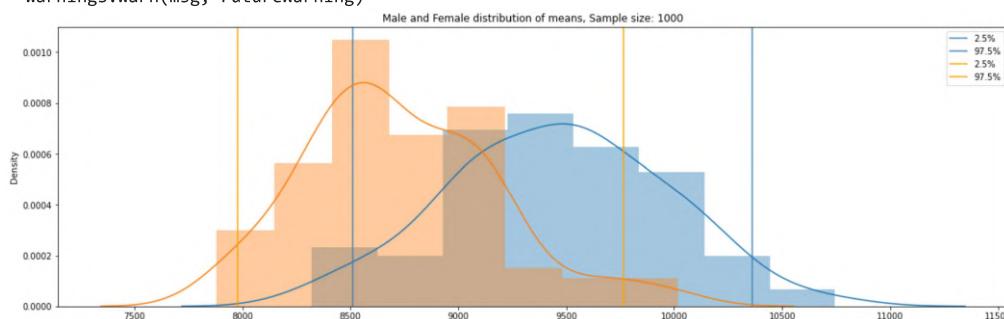
## ▼ CLT and CI Visual illustrations

```

1 avgamt_gender = df[['User_ID','Gender','Purchase']].copy()
2
3 avgamt_male = avgamt_gender[avgamt_gender['Gender']=='M'][['Purchase']]
4 avgamt_female = avgamt_gender[avgamt_gender['Gender']=='F'][['Purchase']]
5
6 sample_size = 100
7 num_repetitions = 100
8
9 male_means = [avgamt_male.sample(sample_size).mean() for i in range (num_repetitions)]
10 female_means = [avgamt_female.sample(sample_size).mean() for i in range (num_repetitions)]
11
12 plt.figure(figsize = (20,6))
13
14 sns.distplot(male_means).set_title("Male and Female distribution of means, Sample size: 1000")
15 plt.axvline(np.percentile(male_means, 2.5), linestyle='--', color='steelblue',label="2.5%")
16 plt.axvline(np.percentile(male_means, 97.5), linestyle='--', color='steelblue',label="97.5%")
17
18 sns.distplot(female_means)
19 plt.axvline(np.percentile(female_means, 2.5), linestyle='--', color='orange',label="2.5%")
20 plt.axvline(np.percentile(female_means, 97.5), linestyle='--', color='orange',label="97.5%")
21
22 plt.legend()
23 plt.show()

```

/usr/local/lib/python3.8/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is warnings.warn(msg, FutureWarning)  
 /usr/local/lib/python3.8/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is warnings.warn(msg, FutureWarning)



With taking small size samples and less repetitions. Here we can see that the confidence interval overlap and distribution is not normally distributed. Thus, now we cannot conclude that man spend more money per transaction than women. We will increase the size of samples and repetitions to get overwhelming evidence.

```

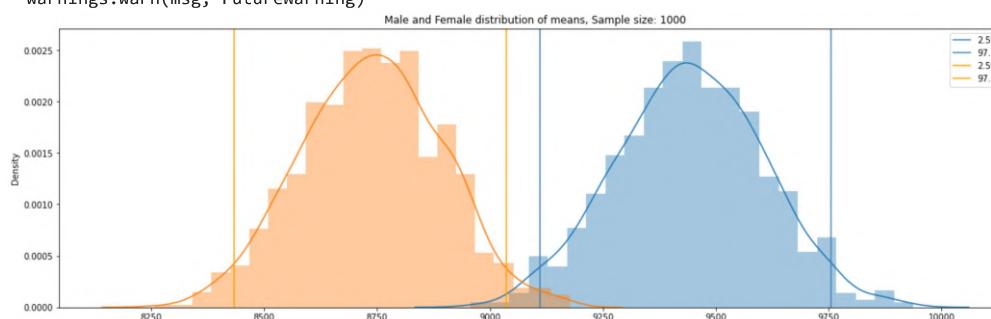
1 avgamt_gender = df[['User_ID','Gender','Purchase']].copy()
2
3 avgamt_male = avgamt_gender[avgamt_gender['Gender']=='M'][['Purchase']]
4 avgamt_female = avgamt_gender[avgamt_gender['Gender']=='F'][['Purchase']]
5
6 sample_size = 1000
7 num_repetions = 1000
8
9 male_means = [avgamt_male.sample(sample_size).mean() for i in range (num_repetions)]
10 female_means = [avgamt_female.sample(sample_size).mean() for i in range (num_repetions)]
11
12 plt.figure(figsize = (20,6))
13
14 sns.distplot(male_means).set_title("Male and Female distribution of means, Sample size: 1000")
15 plt.axvline(np.percentile(male_means, 2.5), linestyle='-', color='steelblue',label="2.5%")
16 plt.axvline(np.percentile(male_means, 97.5), linestyle='-', color='steelblue',label="97.5%")
17
18 sns.distplot(female_means)
19 plt.axvline(np.percentile(female_means, 2.5), linestyle='-', color='orange',label="2.5%")
20 plt.axvline(np.percentile(female_means, 97.5), linestyle='-', color='orange',label="97.5%")
21
22 plt.legend()
23 plt.show()

```

```

/usr/local/lib/python3.8/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is
  warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.8/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is
  warnings.warn(msg, FutureWarning)

```



confidence intervals of average male and female spending are not overlapping. By increasing the size of samples, narrow confidence intervals in relation to the point estimate tell us that the estimated value is relatively stable!

When the sample size increases, the standard error decreases. Hence, as we increase the sample size, the difference between the sample mean and the population mean tends to decrease. In other words, larger sample sizes produce more precise estimates!

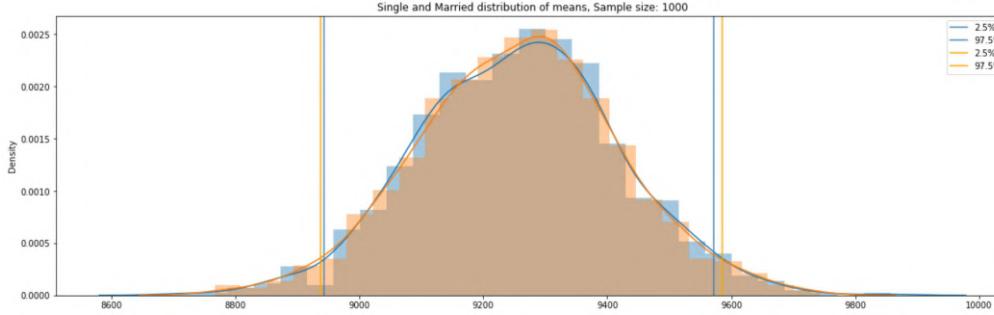
As number of repetition increase, the sampling distributions more closely approximate the normal distribution and become more tightly clustered around the population mean even for skewed, nonnormal data!

```

1 avgamt_status = df[['User_ID','Marital_Status','Purchase']].copy()
2
3 avgamt_single = avgamt_status[avgamt_status['Marital_Status']==0][['Purchase']]
4 avgamt_married = avgamt_status[avgamt_status['Marital_Status']==1][['Purchase']]
5
6 sample_size = 1000
7 num_repetions = 1000
8
9 single_means = [avgamt_single.sample(sample_size).mean() for k in range (num_repetions)]
10 married_means = [avgamt_married.sample(sample_size).mean() for k in range (num_repetitions)]
11
12 plt.figure(figsize = (20,6))
13
14 sns.distplot(single_means).set_title("Single and Married distribution of means, Sample size: 1000")
15 plt.axvline(np.percentile(single_means, 2.5), linestyle='-', color='orange',label="2.5%")
16 plt.axvline(np.percentile(single_means, 97.5), linestyle='-', color='orange',label="97.5%")
17
18 sns.distplot(married_means)
19 plt.axvline(np.percentile(married_means, 2.5), linestyle='-', color='orange',label="2.5%")
20 plt.axvline(np.percentile(married_means, 97.5), linestyle='-', color='orange',label="97.5%")
21
22 plt.legend()
23 plt.show()

```

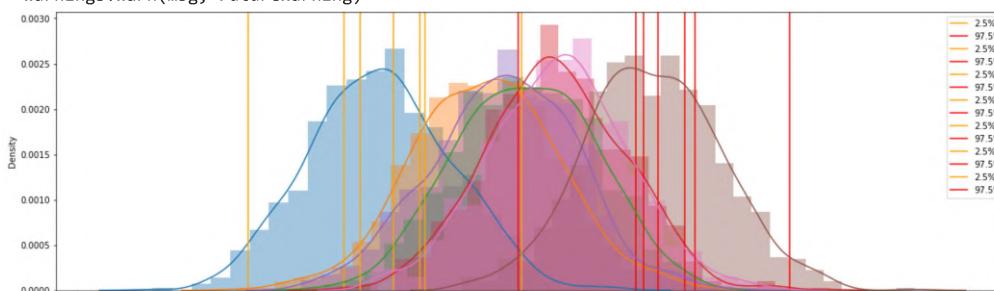
```
/usr/local/lib/python3.8/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.8/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is
warnings.warn(msg, FutureWarning)
```



It's observed that the avg spending on purchase in between married and single is same and overlap almost completely.

```
1 plt.figure(figsize = (20,6))
2
3
4 avgamt_age = df[['User_ID','Age','Purchase']].copy()
5
6 for i in (sorted(list(avgamt_age['Age'].unique()))):
7     avgamt_single = avgamt_age[avgamt_age['Age']==i][['Purchase']]
8
9 sample_size = 1000
10 num_repetitions = 1000
11
12 means = [avgamt_single.sample(sample_size).mean() for k in range (num_repetitions)]
13
14 sns.distplot(means)
15 plt.axvline(np.percentile(means, 2.5), linestyle='--', color='orange',label="2.5%")
16 plt.axvline(np.percentile(means, 97.5), linestyle='--', color='red',label="97.5%")
17 plt.legend()
18 plt.show()
```

```
/usr/local/lib/python3.8/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is
warnings.warn(msg, FutureWarning)
```



It's observed that the avg spending on purchase in among age groups is almost similer and overlaping eachother. There is not significant difference.

## ▼ Confidence intervals

```
1 for i in(['Gender','Marital_Status']):
2     cat = list(df[i].unique())
3     for j in (cat):
4         sample_size = 1000
5         num_repetitions = 1000
6         pop = df[df[i]==j][['Purchase']]
7         means = [pop.sample(sample_size).mean() for k in range (num_repetitions)]
8
9     def CI(p):
10        sample_n = 1000
11        std_dev = np.std(df['Purchase'][df[i]==j])
12        margin_error = norm.ppf((p+100-p)/2)*std_dev/np.sqrt(sample_n)
13        Lower_Limit = round(np.mean(means)-margin_error)
14        Upper_Limit = round(np.mean(means)+margin_error)
15        print(f"{p}% CI for {i} {j}: [{Lower_Limit},{Upper_Limit}]")
16 CI(90)
17 CI(95)
18 CI(99)
19 print()
```

90% CI for Gender F: [8495,8991]  
95% CI for Gender F: [8448,9039]  
99% CI for Gender F: [8355,9131]

90% CI for Gender M: [9169,9698]  
95% CI for Gender M: [9118,9749]  
99% CI for Gender M: [9019,9848]

90% CI for Marital\_Status 0: [9007,9530]  
95% CI for Marital\_Status 0: [8957,9580]  
99% CI for Marital\_Status 0: [8859,9678]

90% CI for Marital\_Status 1: [8995,9517]  
95% CI for Marital\_Status 1: [8945,9567]  
99% CI for Marital\_Status 1: [8847,9665]

```
1 for i in(['Age']):
2     cat = sorted(list(df[i].unique()))
3     for j in (cat):
4         sample_size = 1000
5         num_repetitions = 1000
6         pop = df[df[i]==j][['Purchase']]
7         means = [pop.sample(sample_size).mean() for k in range (num_repetitions)]
8
9     def CI(p):
10        sample_n = 1000
11        std_dev = np.std(df['Purchase'][df[i]==j])
12        margin_error = norm.ppf((p+100-p)/2)*std_dev/np.sqrt(sample_n)
13        Lower_Limit = round(np.mean(means)-margin_error)
14        Upper_Limit = round(np.mean(means)+margin_error)
15        print(f"{p}% CI for {i} {j}: [{Lower_Limit},{Upper_Limit}]")
16 CI(90)
17 CI(95)
18 CI(99)
19 print()
```

90% CI for Age 0-17: [8673,9204]  
95% CI for Age 0-17: [8622,9255]  
99% CI for Age 0-17: [8522,9355]

90% CI for Age 18-25: [8907,9430]  
95% CI for Age 18-25: [8856,9481]  
99% CI for Age 18-25: [8758,9579]

90% CI for Age 26-35: [8992,9513]  
95% CI for Age 26-35: [8942,9563]  
99% CI for Age 26-35: [8844,9660]

90% CI for Age 36-45: [9064,9586]  
95% CI for Age 36-45: [9014,9636]  
99% CI for Age 36-45: [8916,9734]

90% CI for Age 46-50: [8961,9478]  
95% CI for Age 46-50: [8911,9527]  
99% CI for Age 46-50: [8815,9624]

90% CI for Age 51-55: [9269,9798]  
95% CI for Age 51-55: [9219,9849]  
99% CI for Age 51-55: [9119,9948]

90% CI for Age 55+: [9074,9595]  
95% CI for Age 55+: [9024,9645]  
99% CI for Age 55+: [8927,9743]

The larger confidence levels lead to wider confidence intervals but lower precision.

With a 95 percent confidence interval, we have a 5 percent possibility of being wrong. With a 99 percent confidence interval will be wider than a 95 percent confidence interval and precision has to be lower. That's the intuitive way to understand it. We have seen it mathematically by using the confidence interval formula as well as Visual Analysis.

## Questions and Answers:

(1) Are women spending more money per transaction than men? Why or Why not?

Ans: No. Sampling mean of per transaction from male customer and female customers are different. CI's of male and female do not overlap and upper limits of female purchase CI are lesser than lower limits of male purchase CI. This proves that men usually spend more than women. The reason for less purchase by women could have several factors: Walmart store might have less products or less categories for females. Males might be doing the purchase for females. Salary can be a factor in less purchase for females.

(2) Confidence intervals and distribution of the mean of the expenses by female and male customers.

At 95% Confidence Interval with sample size 1000

Average amount spent by male customers lie in the range: [8448,9039]

Average amount spent by female customers lie in the range: [9118,9749]

(3) Are confidence intervals of average male and female spending overlapping? How can Walmart leverage this conclusion to make changes or improvements?

Ans: No. Confidence intervals of average male and female spending are not overlapping. This trend can be changed via introducing female centric marketing strategies by Walmart so that more female customers are attracted to increase female purchases to achieve comparable statistics close to 50%.

(4) Results when the same activity is performed for Married vs Unmarried

At 95% Confidence Interval with sample size 1000

Average amount spent by single customers lie in the range: [8957,9580]

Average amount spent by married customers lie in the range: [8945,9567]

It's observed that the avg spending on purchase in between married and single is same and overlap almost completely.

Walmart can focus these two groups as a single cohort to target them as a single group in order to increase the sale.

But when it comes to specific category of items, Single male customer purchase more products and Married female customers purchase less product. Walmart can increase items in store that are made specially for married female customers or give discounts to attract those customers in store.

(5) Results when the same activity is performed for Age

with sample size 1000

95% CI for Age 0-17: [8622,9255]

95% CI for Age 18-25: [8856,9481]

95% CI for Age 26-35: [8942,9563]

95% CI for Age 36-45: [9014,9636]

95% CI for Age 46-50: [8911,9527]

95% CI for Age 51-55: [9219,9849]

95% CI for Age 55+: [9024,9645]

It's observed that the avg spending on purchase in among age groups is almost similer and CI values overlapping. So we can say that There is not significant difference in avg purchase among age groups. All the upper limits are greater than lower limits of different Age groups, that means overlapping.

## Note -

- I would like to inform that, All the non visual analysis results are statistics of given data (550068 customers), that results are not 100% true for 50 million customers.
- All the appropriate basic insights I have mentioned with relevant code, and calculations based on results.
- There is no significant relationship between variables.
- Comments for univariate and bivariate plots I have mentioned with relevant plots.
- I have illustrated the insights of CLT while exploration.
- We could consider the whole data set of 550068 rows as one sample, and we could do resampling bootstrapping of 550068 sample size to predict statistics for 50 million customers population.

- In all analysis I have considered given data of 550068 rows as population then I have taken many samples of size = 1000 from given dataset.

## Insights and Recommendations

1. Men spent more money than women, company can focus on retaining the male customers and getting more male customers. we can say like women spend less on Black Friday than men. Company should also take initiatives to attract female customers, company can increase products in store that made for women, like makeup items, household items, kitchen utensils, accessories .
2. Product\_Category - 1, 5, 8 have highest purchasing frequency. it means these are the products in these categories are in more demand. Company can focus on selling more of these products.
3. Unmarried customers spend more money than married customers, So company should focus on acquisition of Unmarried customers.
4. Customers in the age 26-35 spend more money than the others, So company should focus on acquisition of customers who are in the age 26-35.
5. We have more customers aged 26-35 in the city category B and A, company can focus more on these customers for these cities to increase the business.
6. Male customers living in City\_Category C spend more money than other male customers living in B or C, Selling more products in the City\_Category C will help the company increase the revenue.
7. Product category 9, 17 have very less purchase. Company can think of dropping it.
8. The top 10 users who have purchased more company should give more offers and discounts so that they can be retained and can be helpful for companies business.
9. The occupation which are contributing more company can think of offering credit cards or other benefits to those customers by liaising with some financial partners to increase the sales.
10. The top products should be given focus in order to maintain the quality in order to further increase the sales of those products.
11. People who are staying in city for an year have contributed to 35% of the total purchase amount. Company can focus on such customer base who are neither too old nor too new residents in the city.
12. We have highest frequency of purchase order between 5k and 10k, company can focus more on these mid range products to increase the sales.