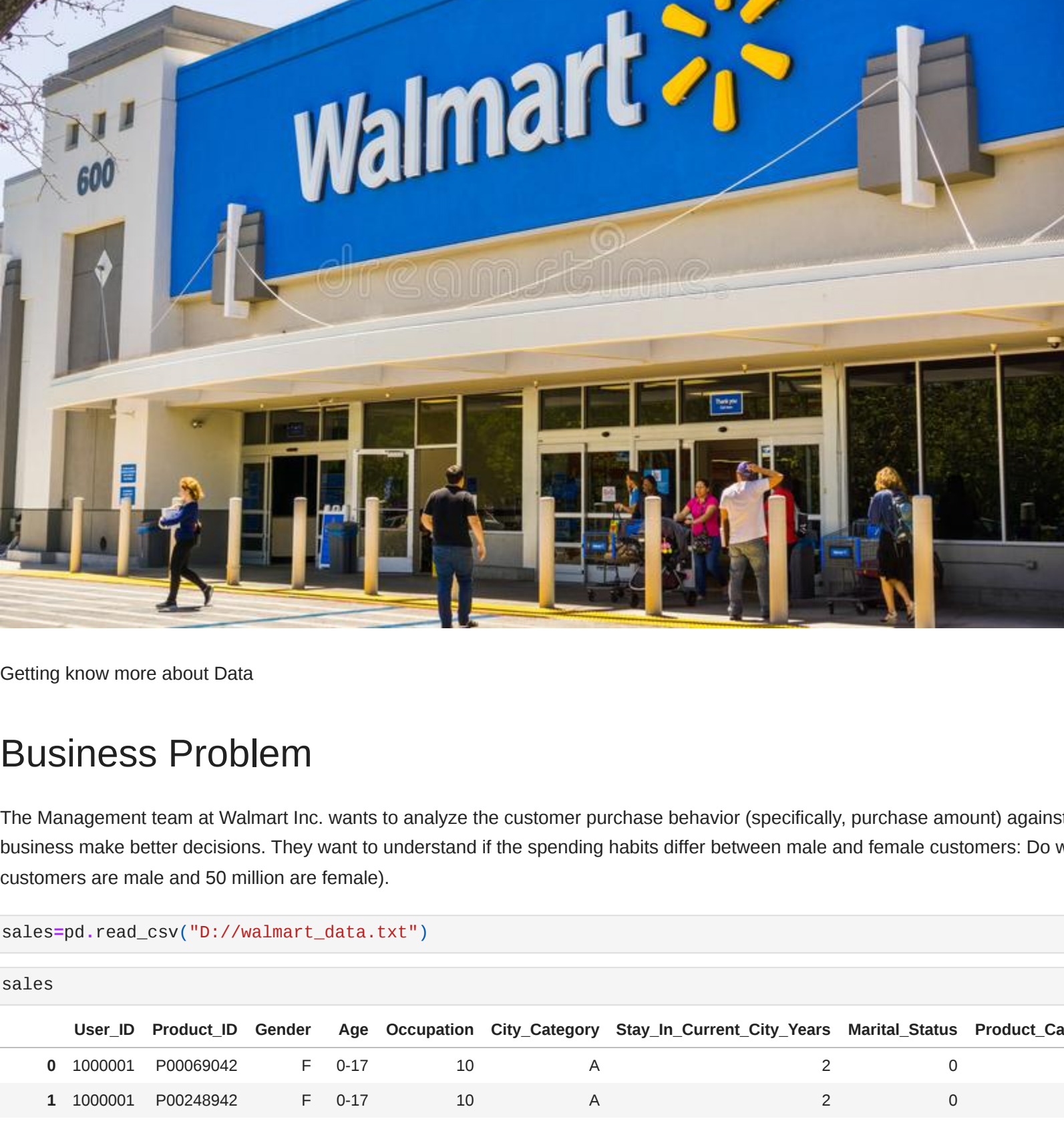


Business Case: Walmart - Confidence Interval and CLT

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import scipy.stats as stats
import statistics
```



Getting know more about Data

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).

```
In [1]: sales=pd.read_csv("0:\walmart_data.txt")

In [3]: sales

Out[3]:
```

User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category	Purchase	
0	1000001	P0009042	F	0-17	10	A	2	0	3	8370
1	1000001	P00248942	F	0-17	10	A	2	0	1	15200
2	1000001	P00087842	F	0-17	10	A	2	0	12	1422
3	1000001	P00085442	F	0-17	10	A	2	0	12	1057
...
550068	1006033	P00372445	M	51-55	13	B	1	1	20	368
550064	1006035	P00375436	F	26-35	1	C	3	0	20	371
550065	1006036	P00375436	F	26-35	15	B	4+	1	20	137
550066	1006038	P00375436	F	55+	1	C	2	0	20	365
550067	1006039	P00371644	F	46-50	0	B	4+	1	20	490

550068 rows × 10 columns

```
In [4]: sales.shape
Out[4]: (550068, 10)

In [5]: sales.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 550068 entries, 0 to 550067
Data columns (total 10 columns):
 # Column                Non-Null Count  Dtype
---  --
 0 User_ID                550068 non-null    int64
 1 Product_ID             550068 non-null    object
 2 Gender                 550068 non-null    object
 3 Age                    550068 non-null    object
 4 Occupation              550068 non-null    int64
 5 City_Category           550068 non-null    object
 6 Stay_In_Current_City_Years  550068 non-null    int64
 7 Marital_Status          550068 non-null    int64
 8 Product_Category        550068 non-null    int64
 9 Purchase                550068 non-null    int64
dtypes: int64(5), object(5)
memory usage: 42.8+ MB

In [6]: sales.describe()

Out[6]:
```

	User_ID	Occupation	Marital_Status	Product_Category	Purchase
count	5.50068e+05	550068.000000	550068.000000	550068.000000	550068.000000
mean	1.003292e+06	8.076707	0.408653	5.404270	9263.968713
std	1.727592e+03	6.522660	0.491770	3.930211	5023.065394
min	1.000000e+06	0.000000	0.000000	1.000000	12.000000
25%	1.001516e+06	2.000000	0.000000	1.000000	5823.000000
50%	1.003077e+06	7.000000	0.000000	5.000000	8947.000000
75%	1.004478e+06	14.000000	1.000000	8.000000	12054.000000
max	1.005604e+06	20.000000	1.000000	20.000000	23961.000000

```
In [7]: sales.isnull()

Out[7]:
```

User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category	Purchase
0	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False
...
550063	False	False	False	False	False	False	False	False	False
550064	False	False	False	False	False	False	False	False	False
550065	False	False	False	False	False	False	False	False	False
550066	False	False	False	False	False	False	False	False	False
550067	False	False	False	False	False	False	False	False	False

550068 rows × 10 columns

```
In [8]: sales.isnull().describe()

Out[8]:
```

	User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category	Purchase
count	550068	550068	550068	550068	550068	550068	550068	550068	550068	550068
unique	1	1	1	1	1	1	1	1	1	1
top	False	False	False	False	False	False	False	False	False	False
freq	550068	550068	550068	550068	550068	550068	550068	550068	550068	550068

```
In [9]: sales.isnull().sum()
Out[9]: 0

In [10]: sales

Out[10]:
```

User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category	Purchase	
0	1000001	P0009042	F	0-17	10	A	2	0	3	8370
1	1000001	P00248942	F	0-17	10	A	2	0	1	15200
2	1000001	P00087842	F	0-17	10	A	2	0	12	1422
...
550063	1006033	P00372445	M	51-55	13	B	1	1	20	368
550064	1006035	P00375436	F	26-35	1	C	3	0	20	371
550065	1006036	P00375436	F	26-35	15	B	4+	1	20	137
550066	1006038	P00375436	F	55+	1	C	2	0	20	365
550067	1006039	P00371644	F	46-50	0	B	4+	1	20	490

550068 rows × 10 columns

Outlier (using boxplot, "describe" method by checking the difference between mean and median, isnull etc.)

```
In [11]: P25 = sales['Purchase'].quantile(0.25)
P50 = sales['Purchase'].quantile(0.50)
P75 = sales['Purchase'].quantile(0.75)

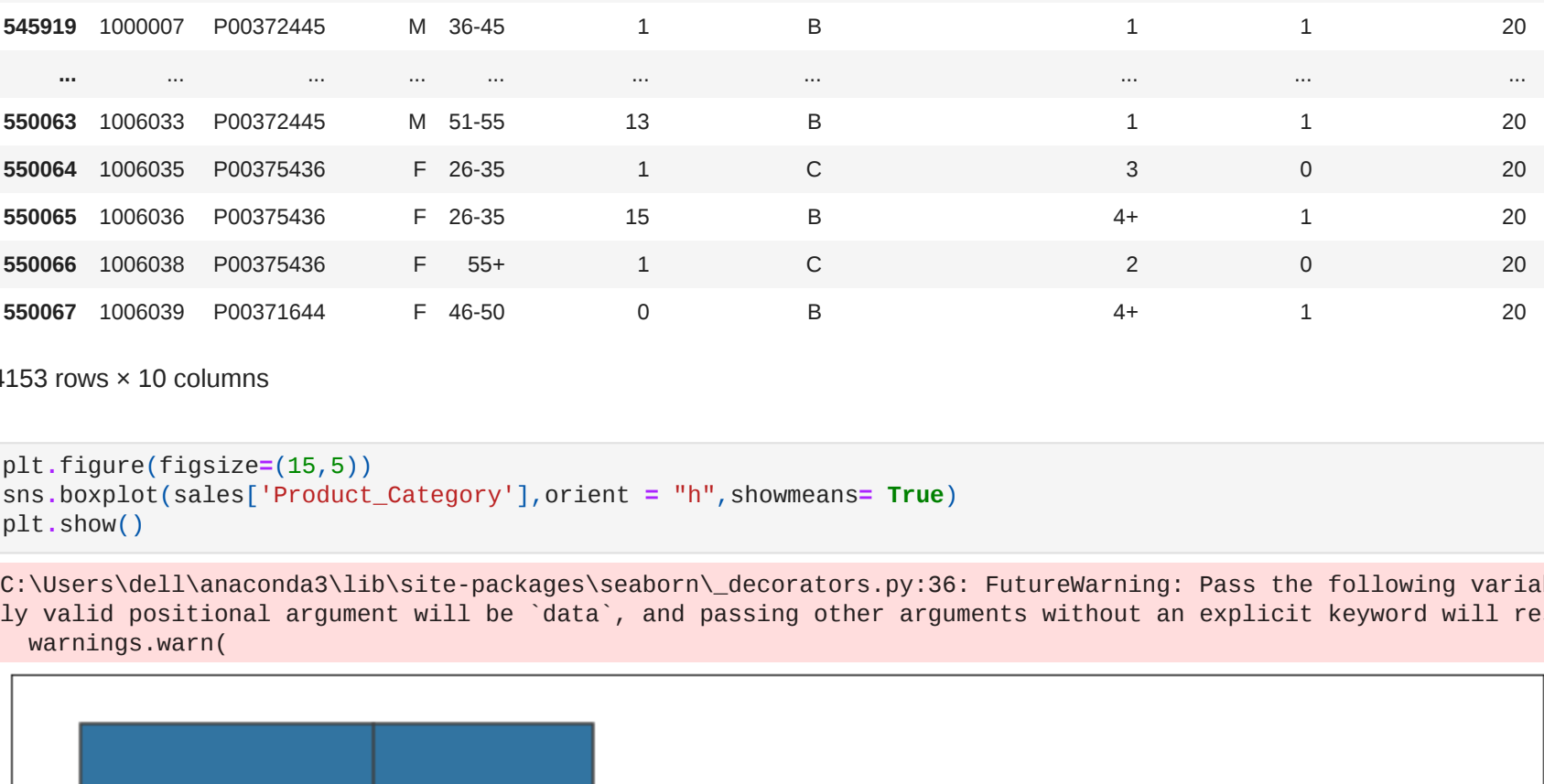
In [12]: IQR = P75 - P25
IQR
Out[12]: 6231.0

In [13]: lower = max(0, P25 - (1.5*IQR))
lower
Out[13]: 0

In [14]: upper = P75 + (1.5*IQR)
upper
Out[14]: 21490.5
```

```
In [15]: plt.figure(figsize=(15,5))
sns.boxplot(sales['Purchase'],orient="h",showmeans=True)
plt.show()

C:\Users\de11\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the on ly valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
warnings.warn()
```



```
In [16]: outliers = sales[sales['Purchase']>upper]
outliers

Out[16]:
```

User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category	Purchase
343	1000058	M	26-35	2	B	3	0	10	23603
672	1000062	F	36-45	3	A	1	0	10	23792
385	1000126	M	18-25	9	B	1	0	10	23233
173	1000139	F	26-35	20	C	2	0	10	23595
154	1000175	F	26-35	2	B	1	0	10	23341
...
54488	1005815	M	26-35	20	B	1	0	10	23752
54478	1005867	F	18-25	4	B	2	0	10	23724
54474	1005852	F	26-35	1	A	0	1	10	23529
54563	1006142	M	51-55	0	C	1	1	10	23663
54577	1006018	F	36-45	1	C	3	0	10	23496

2677 rows × 10 columns

```
In [17]: statistics.mean(sales['Purchase'])
Out[17]: 9263.968712959126

In [18]: statistics.median(sales['Purchase'])
Out[18]: 8947.0

In [19]: product25=sales['Product_Category'].quantile(0.25)
product50=sales['Product_Category'].quantile(0.50)
product75=sales['Product_Category'].quantile(0.75)
product_IQR = product75 - product25
product_lower = max(0, product25 - (1.5*product_IQR))
product_upper = product75 + (1.5*product_IQR)

In [20]: statistics.mean(sales['Product_Category'])
Out[20]: 5.404270817252186

In [21]: statistics.median(sales['Product_Category'])
Out[21]: 5.0

In [22]: outlier_product = sales[sales['Product_Category']>product_upper]
outlier_product

Out[22]:
```

User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category	Purchase
545915	1000001	F	0-17	10	A	2	0	20	612
545916	1000002	M	55+	16	C	4+	0	20	119
545917	1000004	M	46-50	7	B	2	1	20	481
545918	1000006	F	51-55	9	A	1	0	20	480
545919	1000007	M	36-45	1	B	1	1	20	241
...
550063	1006033	M	51-55	13	B	1	1	20	368
550064	1006035	F	26-35	1	C	3	0	20	371
550065	1006036	F	26-35	15	B	4+	1	20	137
550066	1006038	F	55+	1	C	2	0	20	365
550067	1006039	F	46-50	0	B	4+	1	20	490

4153 rows × 10 columns

```
In [23]: plt.figure(figsize=(15,5))
sns.boxplot(sales['Product_Category'],orient="h",showmeans=True)
plt.show()

C:\Users\de11\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the on ly valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
warnings.warn()
```



```
In [24]: df_1=sales.groupby(['Gender','Age','Product_Category'])['Purchase'].sum().reset_index()
df_2=df_1[df_1['Purchase']>5500000]

In [25]: df_2.mean(axis='Index')
Out[25]:
City_Category    10.0
Product_Category 17.5
dtype: float64

In [26]: sales['Purchase'].unique()
Out[26]: array([ 8370, 15280, 1422, ..., 135, 123, 613], dtype=int64)

In [27]: sales['Product_Category'].unique()
Out[27]: array([ 3, 4, 12, 8, 5, 4, 2, 6, 14, 13, 15, 15, 7, 16, 18, 10, 17, 19, 20, 29], dtype=int64)

In [28]: sales['Marital_Status'].unique()
Out[28]: array([0, 1], dtype=int64)

In [29]: df_2=sales.groupby(['User_ID'])['Gender'].unique()
df_2.value_counts()/len(df_2)
Out[29]:
[F]    0.717196
[M]    0.282804
Name: Gender, dtype: float64

In [30]: df_2=sales.groupby(['User_ID'])['Marital_Status'].unique()
df_2.value_counts()/len(df_2)
Out[30]:
[0]    0.588037
[1]    0.411963
Name: Marital_Status, dtype: float64

In [31]: df_2=sales.groupby(['User_ID'])['City_Category'].unique()
df_2.value_counts()/len(df_2)
Out[31]:
[C]    0.532847
[B]    0.289764
[A]    0.177389
Name: City_Category, dtype: float64

In [32]: df_2=pd.DataFrame(sales.groupby(['Gender'])[['Purchase']].sum())
df_2['percent'] = (df_2['Purchase'] / df_2['Purchase'].sum()) * 100
df_2

Out[32]:
```

	Purchase	percent
F	118622642	23.78576
M	3909580200	76.21424

```
In [33]: df_2=pd.DataFrame(sales.groupby(['Age'])[['Purchase']].sum())
df_2['percent'] = (df_2['Purchase'] / df_2['Purchase'].sum()) * 100
df_2

Out[33]:
```

	Purchase	percent
0-17	13491283	2.64730
18-25	913848675	17.93325
26-35	203187758	39.81374
36-45	1025698984	20.14501
46-50	423668058	8.25612
51-55	367095844	7.20561
55+	206767375	3.92850

```
In [34]: plt.figure(figsize=(10, 6))
sns.histplot(data=sales, x='Purchase', kde=True)
plt.show()

Out[34]:
```

```
In [35]: plt.figure(figsize=(2, 4))
sns.boxplot(data=sales, y='Purchase')
plt.show()

Out[35]:
```

```
In [36]: fig, axs = plt.subplots(nrows=2, ncols=2, figsize=(12, 10))
sns.countplot(data=sales, y='Purchase', x='Gender', ax=axs[0,0])
sns.countplot(data=sales, y='Purchase', x='Occupation', ax=axs[0,1])
sns.countplot(data=sales, y='Purchase', x='City_Category', ax=axs[1,0])
sns.countplot(data=sales, y='Purchase', x='Marital_Status', ax=axs[1,1])
plt.show()

Out[36]:
```

```
In [37]: plt.figure(figsize=(12, 5))
sns.countplot(data=sales, x='Product_Category')
plt.show()

Out[37]:
```

```
def validate
fig, axs = plt.subplots(nrows=1, ncols=2, figsize=(13,5))
sns.histplot(data=sales[sales['Gender']=='M'], ax=axs[0], set_title('Male Spending'))
sns.histplot(data=sales[sales['Gender']=='F'], ax=axs[1], set_title('Female Spending'))
plt.show()

Out[38]:
```

```
In [39]: fig, axs = plt.subplots(nrows=2, ncols=2, figsize=(15, 6))
fig.suptitle('adjust (top=1.5)')
sns.boxplot(data=sales, y='Purchase', x='Gender', hue='Age', ax=axs[0,0])
sns.boxplot(data=sales, y='Purchase', x='Gender', hue='City_Category', ax=axs[0,1])
sns.boxplot(data=sales, y='Purchase', x='Gender', hue='Marital_Status', ax=axs[1,0])
sns.boxplot(data=sales, y='Purchase', x='Gender', hue='Stay_In_Current_City_Years', ax=axs[1,1])
plt.show()

Out[39]:
```

```
In [40]: sns.pairplot(sales.corr(), annot=True, cmap='Greens', linewidth=0.5)

Out[40]:
```

```
In [41]: avg_gender = sales.groupby(['User_ID', 'Gender'])[['Purchase']].sum()
avg_gender = avg_gender.reset_index()

Out[41]:
```

User_ID	Gender	Purchase	
0	1000001	F	334093
1	1000002	M	810472
2	1000003	M	341635
3	1000004	M	206468
4	1000005	M	821001
...
5086	1000036	F	4110508
5087	1000037	F	1119538
5088	1000038	F	90034
5089	1000039	F	590319
5090	1000040	M	1853299

5091 rows × 3 columns

```
In [42]: fig, axs = plt.subplots(nrows=1, ncols=2, figsize=(16,5))
sns.histplot(data=avg_gender.groupby('Gender')[['Purchase']].mean())
sns.histplot(data=avg_gender.groupby('Gender')[['Purchase']].mean(), ax=axs[0], set_title('Females Avg Spend'))
sns.histplot(data=avg_gender.groupby('Gender')[['Purchase']].mean(), ax=axs[1], set_title('Males Avg Spend'))
Text(0.5, 1.0, 'Males Avg Spend')

Out[42]:
```

```
In [43]: avg_male = avg_gender.groupby('Gender')[['Purchase']].mean()
avg_female = avg_gender.groupby('Gender')[['Purchase']].mean()

In [44]: #Finding the sample (sample size=1000) for avg purchase amount for males and females
genders = ["M", "F"]
sample_size = 1000
num_reptitions = 1000
male_means = []
female_means = []

for i in range(num_reptitions):
    male_mean = avg_male.sample(sample_size, replace=True)['Purchase'].mean()
    female_mean = avg_female.sample(sample_size, replace=True)['Purchase'].mean()
    male_means.append(male_mean)
    female_means.append(female_mean)

In [45]: #final insight
z99=1.645 #95% Confidence Interval
z95=1.960 #95% Confidence Interval
z99=2.576 #99% Confidence Interval

print("Population avg spend amount for Male: {:.2f}".format(avg_male['Purchase'].mean()))
print("Population avg spend amount for Female: {:.2f}\n".format(avg_female['Purchase'].mean()))
print("Sample avg spend amount for Male: {:.2f}".format(np.mean(male_means)))
print("Sample avg spend amount for Female: {:.2f}\n".format(np.mean(female_means)))
print("Sample std for Male: {:.2f}".format(pd.Series(male_means).std()))
print("Sample std for Female: {:.2f}\n".format(pd.Series(female_means).std()))
print("Sample std error for Male: {:.2f}\n".format(pd.Series(male_means).std()/np.sqrt(1000)))
print("Sample std error for Female: {:.2f}\n".format(pd.Series(female_means).std()/np.sqrt(1000)))
sample_mean=np.mean(male_means)
sample_std=np.mean(female_means)
sample_std_male=pd.Series(male_means).std()
sample_std_female=pd.Series(female_means).std()
sample_std_error_male=sample_std_male/np.sqrt(1000)
sample_std_error_female=sample_std_female/np.sqrt(1000)
Upper_Limit_male=sample_std_error_male + sample_mean_male
Lower_Limit_male=sample_mean_male - sample_std_error_male
Upper_Limit_female=sample_std_error_female + sample_mean_female
Lower_Limit_female=sample_mean_female - sample_std_error_female
print("Male.CI: [{}Lower_Limit_male,Upper_Limit_male]".format(' '.join([str(x) for x in [Lower_Limit_male, Upper_Limit_male])]))
print("Female.CI: [{}Lower_Limit_female,Upper_Limit_female]".format(' '.join([str(x) for x in [Lower_Limit_female, Upper_Limit_female])]))
Population avg spend amount for Male: 925344.40
Population avg spend amount for Female: 712824.39

Sample avg spend amount for Male: 925714.29
Sample avg spend amount for Female: 711438.62
Sample std for Male: 32102.82
Sample std for Female: 23252.69

Sample std error for Male: 1015.39
Sample std error for Female: 799.82

Male.CI: [924944.3196454166, 927384.7628745634]
Female.CI: [710122.912604797, 712754.3361485283]
```

```
In [46]: fig, axis = plt.subplots(nrows=1, ncols=2, figsize=(20, 6))
axis[0].hist(male_means, bins=35)
axis[1].hist(female_means, bins=35)
axis[0].set_title('Male distribution of means, Sample size: 1000')
axis[1].set_title('Female distribution of means, Sample size: 1000')
plt.show()

Out[46]:
```

Question:

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

Ans: No. CIs of male and female do not overcome and uppers of female purchase CI are lesser than lower limits of male purchase CI. This proves that men usually spend more than women. Males might be doing the purchase for females. Salary can be a factor in lesser female purchase. We also need to see whether male-based products were sold more than women products to clearly identify difference in spending pattern.

1. Confidence intervals and distribution of the mean of the expenses by female and male customers. Ans: At 99% Confidence Interval with sample size 1000 Average amount spend by male customers lie in the lower parameters Average amount spend by female customers lie in higher parameter

2. 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 notover come. Attracted more female customers are chances to increase female purchases

3. Results when the same activity is performed for Married or Unmarried Ans: Average amount spend by married customers lie in the range: [841059.630978392, 845078.140167503] Average amount spend by unmarried customers lie in the range: [879093.3492016713, 884078.6762803260]

4. Results when the same activity is performed for Age Ans:its really based up on the different factors and the ratios that are more act on trends and personal preference and taste patterns

Recommendations