```
1 #importing libraries
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from scipy.stats import t
7 import warnings
8 warnings.filterwarnings('ignore')
9 import copy
```

```
1 !gdown https://d2beiqkhq929f0.cloudfront.net/public_assets/assets/000/001/293/original/walmart_d
```

```
Downloading...
From: https://d2beiqkhq929f0.cloudfront.net/public_assets/assets/000/001/293/original/walmart_d
To: /content/walmart_data.csv?1641285094
100% 23.0M/23.0M [00:00<00:00, 166MB/s]
```

```
1 df = pd.read_csv('walmart_data.csv?1641285094')
```

```
1 df.head()
```

|   | User_ID | Product_ID | Gender | Age | Occupation | City_Category | Stay_In_Current_City_Years | Mar: |
|---|---------|-----------|--------|-----|-----------|---------------|---------------------------|------|
| 0 | 1000001 | P00069042 | F | 0-17 | 10 | A | 2 | |
| 1 | 1000001 | P00248942 | F | 0-17 | 10 | A | 2 | |
| 2 | 1000001 | P00087842 | F | 0-17 | 10 | A | 2 | |

```
1 df.tail()
```

|   | User_ID | Product_ID | Gender | Age | Occupation | City_Category | Stay_In_Current_City_Years |
|---|---------|-----------|--------|-----|-----------|---------------|---------------------------|
| 550063 | 1006033 | P00372445 | M | 51-55 | 13 | B | 1 |
| 550064 | 1006035 | P00375436 | F | 26-35 | 1 | C | 3 |
| 550065 | 1006036 | P00375436 | F | 26-35 | 15 | B | 4+ |

```
1 df.shape
```

```
(550068, 10)
```

```
1 df.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  object
  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: 42.0+ MB
```

## 🔍 Insights

- From the above analysis, it is clear that, data has total of 10 features with lots of mixed alpha numeric data.

- Apart from Purchase Column, all the other data types are of categorical type. We will change the datatypes of all such columns to category

## ∨ 🔁 Changing the Datatype of Columns

```
1 for i in df.columns[:-1]:
2     df[i] = df[i].astype('category')
3
4 df.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  category
  1   Product_ID                    550068 non-null  category
  2   Gender                        550068 non-null  category
  3   Age                           550068 non-null  category
  4   Occupation                    550068 non-null  category
  5   City_Category                 550068 non-null  category
  6   Stay_In_Current_City_Years    550068 non-null  category
  7   Marital_Status                550068 non-null  category
  8   Product_Category              550068 non-null  category
  9   Purchase                      550068 non-null  int64
dtypes: category(9), int64(1)
memory usage: 10.3 MB
```

## ∨ 📝 Statistical Summary

### Satistical summary of object type columns

```
1 df.describe(include = 'category')
```

| | User_ID | Product_ID | Gender | Age | Occupation | City_Category | Stay_In_Current_City_Ye |
|---|---|---|---|---|---|---|---|
| count | 550068 | 550068 | 550068 | 550068 | 550068 | 550068 | 550 |
| unique | 5891 | 3631 | 2 | 7 | 21 | 3 | |
| top | 1001680 | P00265242 | M | 26-35 | 4 | B | |
| freq | 1026 | 1880 | 414259 | 219587 | 72308 | 231173 | 193 |

## 🔍 Insights

1. User_ID - Among 5,50,068 transactions there are 5891 unique user_id, indicating same customers buying multiple products.

2. Product_ID - Among 5,50,068 transactions there are 3631 unique products,with the product having the code P00265242 being the highest seller, with a maximum of 1,880 units sold.

3. Gender - Out of 5,50,068 transactions, 4,14,259 (nearly 75%) were done by male gender indicating a significant disparity in purchase behavior between males and females during the Black Friday event.

4. Age - We have 7 unique age groups in the dataset. 26 - 35 Age group has maximum of 2,19,587 transactions. We will analyse this feature in detail in future

5. Stay_In_Current_City_Years - Customers with 1 year of stay in current city accounted to maximum of 1,93,821 transactions among all the other customers with (0,2,3,4+) years of stay in current city

6. Marital_Status - 59% of the total transactions were done by Unmarried Customers and 41% by Married Customers.

## Satistical summary of numerical data type columns

```
1 df.describe()
```

| | Purchase |
|---|---|
| count | 550068.000000 |
| mean | 9263.968713 |
| std | 5023.065394 |
| min | 12.000000 |
| 25% | 5823.000000 |
| 50% | 8047.000000 |
| 75% | 12054.000000 |
| max | 23961.000000 |

## 🔍 Insights

The purchase amounts vary widely, with the minimum recorded purchase being $12$ and the maximum reaching $23961$. The median purchase amount of $8047$ is notably lower than the mean purchase amount of $9264$, indicating a right-skewed distribution where a few high-value purchases pull up the mean

## 👥 Duplicate Detection

```
1 df.duplicated().value_counts()
```

```
False    550068
dtype: int64
```

## 🔍 Insights

- There are no duplicate entries in the dataset

## ✅ Sanity Check for columns

```
1 # checking the unique values for columns
2 for i in df.columns:
3     print('Unique Values in',i,'column are :-')
4     print(df[i].unique())
5     print('-'*70)
```

```
Unique Values in User_ID column are :-
[1000001, 1000002, 1000003, 1000004, 1000005, ..., 1004588, 1004871, 1004113, 1005391, 1001529]
Length: 5891
Categories (5891, int64): [1000001, 1000002, 1000003, 1000004, ..., 1006037, 1006038, 1006039,
----------------------------------------------------------------------
Unique Values in Product_ID column are :-
['P00069042', 'P00248942', 'P00087842', 'P00085442', 'P00285442', ..., 'P00375436', 'P00372445'
Length: 3631
Categories (3631, object): ['P00000142', 'P00000242', 'P00000342', 'P00000442', ..., 'P0099642'
                                'P0099742', 'P0099842', 'P0099942']
----------------------------------------------------------------------
Unique Values in Gender column are :-
['F', 'M']
Categories (2, object): ['F', 'M']
----------------------------------------------------------------------
Unique Values in Age column are :-
['0-17', '55+', '26-35', '46-50', '51-55', '36-45', '18-25']
Categories (7, object): ['0-17', '18-25', '26-35', '36-45', '46-50', '51-55', '55+']
----------------------------------------------------------------------
Unique Values in Occupation column are :-
[10, 16, 15, 7, 20, ..., 18, 5, 14, 13, 6]
Length: 21
Categories (21, int64): [0, 1, 2, 3, ..., 17, 18, 19, 20]
----------------------------------------------------------------------
Unique Values in City_Category column are :-
['A', 'C', 'B']
Categories (3, object): ['A', 'B', 'C']
----------------------------------------------------------------------
Unique Values in Stay_In_Current_City_Years column are :-
['2', '4+', '3', '1', '0']
```

```
Categories (5, object): ['0', '1', '2', '3', '4+']
-----------------------------------------------------------------
Unique Values in Marital_Status column are :-
[0, 1]
Categories (2, int64): [0, 1]
-----------------------------------------------------------------
Unique Values in Product_Category column are :-
[3, 1, 12, 8, 5, ..., 10, 17, 9, 20, 19]
Length: 20
Categories (20, int64): [1, 2, 3, 4, ..., 17, 18, 19, 20]
-----------------------------------------------------------------
Unique Values in Purchase column are :-
[ 8370 15200  1422 ...   135   123   613]
-----------------------------------------------------------------
```

## ✓ 🔍 Insights

- The dataset does not contain any abnormal values.
- We will convert the 0,1 in Marital Status column as married and unmarried

```
1 #replacing the values in marital_status column
2
3 df['Marital_Status'] = df['Marital_Status'].replace({0:'Unmarried',1:'Married'})
4 df['Marital_Status'].unique()
```

```
['Unmarried', 'Married']
Categories (2, object): ['Unmarried', 'Married']
```

## ✓ 👨🏽‍💻 Missing Value Analysis

```
1 df.isnull().sum()
```

```
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
```

## 🔍 Insights

- The dataset does not contain any missing values.

## ✓ Univariate Analysis

Numerical Variables

💰 Purchase Amount Distribution

```python
1  #setting the plot style
2
3  fig = plt.figure(figsize = (15,10))
4  gs = fig.add_gridspec(2,1,height_ratios=[0.65, 0.35])
5
6                                    #creating purchase amount histogram
7
8  ax0 = fig.add_subplot(gs[0,0])
9
10 ax0.hist(df['Purchase'],color= '#5C8374',linewidth=0.5,edgecolor='black',bins = 20)
11 ax0.set_xlabel('Purchase Amount',fontsize = 12,fontweight = 'bold')
12 ax0.set_ylabel('Frequency',fontsize = 12,fontweight = 'bold')
13
14 #removing the axis lines
15 for s in ['top','left','right']:
16     ax0.spines[s].set_visible(False)
17
18 #setting title for visual
19 ax0.set_title('Purchase Amount Distribution',{'font':'serif', 'size':15,'weight':'bold'})
20
21
22                                    #creating box plot for purchase amount
23
24 ax1 = fig.add_subplot(gs[1,0])
25 boxplot = ax1.boxplot(x = df['Purchase'],vert = False,patch_artist = True,widths = 0.5)
26
27 # Customize box and whisker colors
28 boxplot['boxes'][0].set(facecolor='#5C8374')
29
30 # Customize median line
31 boxplot['medians'][0].set(color='red')
32
33 # Customize outlier markers
34 for flier in boxplot['fliers']:
35     flier.set(marker='o', markersize=8, markerfacecolor= "#4b4b4c")
36
37 #removing the axis lines
38 for s in ['top','left','right']:
39     ax1.spines[s].set_visible(False)
40
41 #adding 5 point summary annotations
42 info = [i.get_xdata() for i in boxplot['whiskers']] #getting the upperlimit,Q1,Q3 and lowerlimit
43
44 median = df['Purchase'].quantile(0.5) #getting Q2
45
46 for i,j in info: #using i,j here because of the output type of info list comprehension
47
48     ax1.annotate(text = f"{i:.1f}", xy = (i,1), xytext = (i,1.4),fontsize = 12,
49                 arrowprops= dict(arrowstyle="<-", lw=1, connectionstyle="arc,rad=0"))
50
51     ax1.annotate(text = f"{j:.1f}", xy = (j,1), xytext = (j,1.4),fontsize = 12,
52                 arrowprops= dict(arrowstyle="<-", lw=1, connectionstyle="arc,rad=0"))
53
54 #adding the median separately because it was included in info list
55 ax1.annotate(text = f"{median:.1f}",xy = (median,1),xytext = (median + 1,1.4),fontsize = 12,
56             arrowprops= dict(arrowstyle="<-", lw=1, connectionstyle="arc,rad=0"))
57
58 #removing y-axis ticks
59 ax1.set_yticks([])
60
61 #adding axis label
```

```
62 ax1.set_xlabel('Purchase Amount',fontweight = 'bold',fontsize = 12)
63
64 plt.show()
```

**Purchase Amount Distribution**



## Calculating the Number of Outliers

- As seen above, Purchase amount over 21399 is considered as outlier. We will count the number of outliers as below

```
1 len(df.loc[df['Purchase'] > 21399,'Purchase'])
```

2677

# 🔍 Insights

## Outliers

- There are total of 2677 outliers which is roughly 0.48% of the total data present in purchase amount. We will not remove them as it indicates a broad range of spending behaviors during the sale, highlighting the importance of tailoring marketing strategies to both regular and high-value customers to maximize revenue.

## Distribution

- Data suggests that the majority of customers spent between 5,823 USD and 12,054 USD, with the median purchase amount being 8,047 USD.

- The lower limit of 12 USD while the upper limit of 21,399 USD reveal significant variability in customer spending

## ⌄ Categorical Variables

👩🏻 🧑🏻 Gender, 👫🏻 Marital Status and 🌍 City Category Distribution

```
1  #setting the plot style
2  fig = plt.figure(figsize = (15,12))
3  gs = fig.add_gridspec(1,3)
4
5                                          # creating pie chart for gender disribution
6  ax0 = fig.add_subplot(gs[0,0])
7
8  color_map = ["#3A7089", "#4b4b4c"]
9  ax0.pie(df['Gender'].value_counts().values,labels = df['Gender'].value_counts().index,autopct =
10         shadow = True,colors = color_map,textprops={'fontsize': 13, 'color': 'black'})
11
12 #setting title for visual
13 ax0.set_title('Gender Distribution',{'font':'serif', 'size':15,'weight':'bold'})
14
15                                          # creating pie chart for marital status
16 ax1 = fig.add_subplot(gs[0,1])
17
18 color_map = ["#3A7089", "#4b4b4c"]
19 ax1.pie(df['Marital_Status'].value_counts().values,labels = df['Marital_Status'].value_counts().
20         shadow = True,colors = color_map,textprops={'fontsize': 13, 'color': 'black'})
21
22 #setting title for visual
23 ax1.set_title('Marital Status Distribution',{'font':'serif', 'size':15,'weight':'bold'})
24
25                                          # creating pie chart for city category
26 ax1 = fig.add_subplot(gs[0,2])
27
28 color_map = ["#3A7089", "#4b4b4c",'#99AEBB']
29 ax1.pie(df['City_Category'].value_counts().values,labels = df['City_Category'].value_counts().in
30         shadow = True,colors = color_map,textprops={'fontsize': 13, 'color': 'black'})
31
32 #setting title for visual
33 ax1.set_title('City Category Distribution',{'font':'serif', 'size':15,'weight':'bold'})
34 plt.show()
```



🔍 Insights

1. Gender Distribution - Data indicates a significant disparity in purchase behavior between males and females during the Black Friday event.

2. Marital Status - Given that unmarried customers account for a higher percentage of transactions, it may be worthwhile to consider specific marketing campaigns or promotions that appeal to this group.

3. City Category - City B saw the most number of transactions followed by City C and City A respectively

## Confidence Interval Construction: Estimating Average Purchase Amount per Transaction

1. **Step 1 - Building CLT Curve**

As seen above, the purchase amount distribution is not Normal. So we need to use Central Limit Theorem. It states the distribution of sample means will approximate a normal distribution, regardless of the underlying population distribution 2.** Step 2 - Building Confidence Interval**

After building CLT curve, we will create a confidence interval predicting population mean at 99%,95% and 90% Confidence level.

```
1 #creating a function to calculate confidence interval
2
3 def confidence_interval(data,ci):
4     #converting the list to series
5     l_ci = (100-ci)/2
6     u_ci = (100+ci)/2
7
8     #calculating lower limit and upper limit of confidence interval
9     interval = np.percentile(data,[l_ci,u_ci]).round(0)
10
11    return interval
```

```python
1  #defining a function for plotting the visual for given confidence interval
2
3  def plot(ci):
4
5      #setting the plot style
6      fig = plt.figure(figsize = (15,8))
7      gs = fig.add_gridspec(2,2)
8
9      #creating separate data frames for each gender
10     df_male = df.loc[df['Gender'] == 'M','Purchase']
11     df_female = df.loc[df['Gender'] == 'F','Purchase']
12
13     #sample sizes and corresponding plot positions
14     sample_sizes = [(100,0,0),(1000,0,1),(5000,1,0),(50000,1,1)]
15
16     #number of samples to be taken from purchase amount
17     bootstrap_samples = 20000
18
19     male_samples = {}
20     female_samples = {}
21
22     for i,x,y in sample_sizes:
23         male_means = [] #list for collecting the means of male sample
24         female_means = [] #list for collecting the means of female sample
25
26         for j in range(bootstrap_samples):
27
28             #creating random 5000 samples of i sample size
29             male_bootstrapped_samples = np.random.choice(df_male,size = i)
30             female_bootstrapped_samples = np.random.choice(df_female,size = i)
31
32             #calculating mean of those samples
33             male_sample_mean = np.mean(male_bootstrapped_samples)
34             female_sample_mean = np.mean(female_bootstrapped_samples)
35
36             #appending the mean to the list
37             male_means.append(male_sample_mean)
38             female_means.append(female_sample_mean)
39
40         #storing the above sample generated
41         male_samples[f'{ci}%_{i}'] = male_means
42         female_samples[f'{ci}%_{i}'] = female_means
43
44         #creating a temporary dataframe for creating kdeplot
45         temp_df = pd.DataFrame(data = {'male_means':male_means,'female_means':female_means})
46
47                                                  #plotting kdeplots
48         #plot position
49         ax = fig.add_subplot(gs[x,y])
50
51         #plots for male and female
52         sns.kdeplot(data = temp_df,x = 'male_means',color ="#3A7089" ,fill = True, alpha = 0.5,a
53         sns.kdeplot(data = temp_df,x = 'female_means',color ="#4b4b4c" ,fill = True, alpha = 0.5
54
55         #calculating confidence intervals for given confidence level(ci)
56         m_range = confidence_interval(male_means,ci)
57         f_range = confidence_interval(female_means,ci)
58
59         #plotting confidence interval on the distribution
60         for k in m_range:
61             ax.axvline(x = k,ymax = 0.9, color ="#3A7089",linestyle = '--')
```
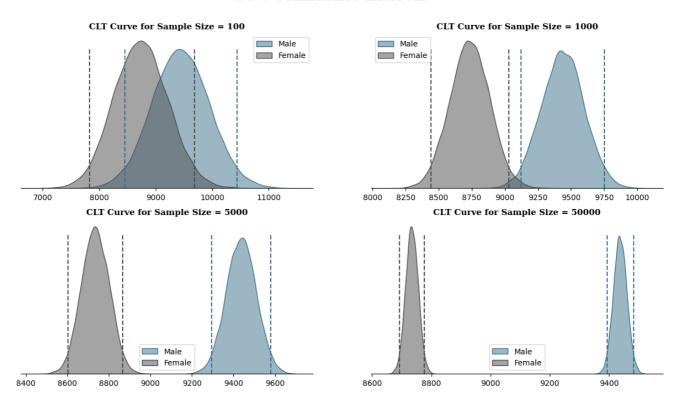
```
62
63            for k in f_range:
64                ax.axvline(x = k,ymax = 0.9, color ="#4b4b4c",linestyle = '--')
65
66
67            #removing the axis lines
68            for s in ['top','left','right']:
69                ax.spines[s].set_visible(False)
70
71            # adjusting axis labels
72            ax.set_yticks([])
73            ax.set_ylabel('')
74            ax.set_xlabel('')
75
76            #setting title for visual
77            ax.set_title(f'CLT Curve for Sample Size = {i}',{'font':'serif', 'size':11,'weight':'bol
78
79            plt.legend()
80
81        #setting title for visual
82        fig.suptitle(f'{ci}% Confidence Interval',font = 'serif', size = 18, weight = 'bold')
83
84        plt.show()
85
86        return male_samples,female_samples
```

```
1 m_samp_90,f_samp_90 = plot(90)
```

# 90% Confidence Interval



```
1 m_samp_95,f_samp_95 = plot(95)
```

## 95% Confidence Interval



# 🔍 Insights

   1. Sample Size

The analysis highlights the importance of sample size in estimating population parameters. It suggests that as the sample size increases, the confidence intervals become narrower and more precise. In business, this implies that larger sample sizes can provide more reliable insights and estimates.

   2. Confidence Intervals

From the above analysis, we can see that except for the Sample Size of 100, the confidence interval do not overlap as the sample size increases. This means that there is a statistically significant difference between the average spending per transaction for men and women within the given samples.

   3. Population Average

We are 95% confident that the true population average for males falls between $9,393$ and $9,483$, and for females, it falls between $8,692$ and $8,777$.

4. Women spend less

Men tend to spend more money per transaction on average than women, as the upper bounds of the confidence intervals for men are consistently higher than those for women across different sample sizes.

## 👫 Marital Status VS 💰 Purchase Amount

### 📊 Data Visualization

```
1  #creating a df for purchase amount vs marital status
2  temp = df.groupby('Marital_Status')['Purchase'].agg(['sum','count']).reset_index()
3
4  #calculating the amount in billions
5  temp['sum_in_billions'] = round(temp['sum'] / 10**9,2)
6
7  #calculationg percentage distribution of purchase amount
8  temp['%sum'] = round(temp['sum']/temp['sum'].sum(),3)
9
10 #calculationg per purchase amount
11 temp['per_purchase'] = round(temp['sum']/temp['count'])
12
13 temp
```
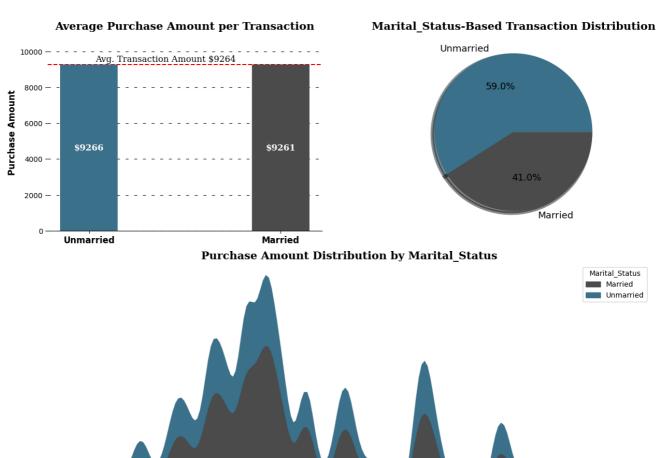
| | Marital_Status | sum | count | sum_in_billions | %sum | per_purchase |
|---|---|---|---|---|---|---|
| 0 | Unmarried | 3008927447 | 324731 | 3.01 | 0.59 | 9266.0 |
| 1 | Married | 2086885295 | 225337 | 2.09 | 0.41 | 9261.0 |

```python
1  #setting the plot style
2  fig = plt.figure(figsize = (15,14))
3  gs = fig.add_gridspec(3,2,height_ratios =[0.10,0.4,0.5])
4
5
6                                              #Distribution of Purchase Amount
7  ax = fig.add_subplot(gs[0,:])
8
9  #plotting the visual
10 ax.barh(temp.loc[0,'Marital_Status'],width = temp.loc[0,'%sum'],color = "#3A7089",label = 'Unmar
11 ax.barh(temp.loc[0,'Marital_Status'],width = temp.loc[1,'%sum'],left =temp.loc[0,'%sum'], color
12
13 #inserting the text
14 txt = [0.0] #for left parameter in ax.text()
15
16 for i in temp.index:
17     #for amount
18     ax.text(temp.loc[i,'%sum']/2 + txt[0],0.15,f"${temp.loc[i,'sum_in_billions']} Billion",
19             va = 'center', ha='center',fontsize=18, color='white')
20
21     #for marital status
22     ax.text(temp.loc[i,'%sum']/2 + txt[0],- 0.20 ,f"{temp.loc[i,'Marital_Status']}",
23             va = 'center', ha='center',fontsize=14, color='white')
24
25     txt += temp.loc[i,'%sum']
26
27 #removing the axis lines
28 for s in ['top','left','right','bottom']:
29     ax.spines[s].set_visible(False)
30
31 #customizing ticks
32 ax.set_xticks([])
33 ax.set_yticks([])
34 ax.set_xlim(0,1)
35
36 #plot title
37 ax.set_title('Marital_Status-Based Purchase Amount Distribution',{'font':'serif', 'size':15,'wei
38
39
40                                              #Distribution of Purchase Amount per Transaction
41
42 ax1 = fig.add_subplot(gs[1,0])
43
44 color_map = ["#3A7089", "#4b4b4c"]
45
46 #plotting the visual
47 ax1.bar(temp['Marital_Status'],temp['per_purchase'],color = color_map,zorder = 2,width = 0.3)
48
49 #adding average transaction line
50 avg = round(df['Purchase'].mean())
51
52 ax1.axhline(y = avg, color ='red', zorder = 0,linestyle = '--')
53
54 #adding text for the line
55 ax1.text(0.4,avg + 300, f"Avg. Transaction Amount ${avg:.0f}",
56         {'font':'serif','size' : 12},ha = 'center',va = 'center')
57
58 #adjusting the ylimits
59 ax1.set_ylim(0,11000)
60
61 #adding the value_counts
```

```python
62 for i in temp.index:
63     ax1.text(temp.loc[i,'Marital_Status'],temp.loc[i,'per_purchase']/2,f"${temp.loc[i,'per_purch
64             {'font':'serif','size' : 12,'color':'white','weight':'bold' },ha = 'center',va = 'c
65
66 #adding grid lines
67 ax1.grid(color = 'black',linestyle = '--', axis = 'y', zorder = 0, dashes = (5,10))
68
69 #removing the axis lines
70 for s in ['top','left','right']:
71     ax1.spines[s].set_visible(False)
72
73 #adding axis label
74 ax1.set_ylabel('Purchase Amount',fontweight = 'bold',fontsize = 12)
75 ax1.set_xticklabels(temp['Marital_Status'],fontweight = 'bold',fontsize = 12)
76
77 #setting title for visual
78 ax1.set_title('Average Purchase Amount per Transaction',{'font':'serif', 'size':15,'weight':'bol
79
80                                          # creating pie chart for Marital_Status disribution
81 ax2 = fig.add_subplot(gs[1,1])
82
83 color_map = ["#3A7089", "#4b4b4c"]
84 ax2.pie(temp['count'],labels = temp['Marital_Status'],autopct = '%.1f%%',
85         shadow = True,colors = color_map,wedgeprops = {'linewidth': 5},textprops={'fontsize': 13
86
87 #setting title for visual
88 ax2.set_title('Marital_Status-Based Transaction Distribution',{'font':'serif', 'size':15,'weight
89
90                                          # creating kdeplot for purchase amount distribution
91
92 ax3 = fig.add_subplot(gs[2,:])
93 color_map = [ "#4b4b4c","#3A7089"]
94
95 #plotting the kdeplot
96 sns.kdeplot(data = df, x = 'Purchase', hue = 'Marital_Status', palette = color_map,fill = True,
97             ax = ax3,hue_order = ['Married','Unmarried'])
98
99 #removing the axis lines
100 for s in ['top','left','right']:
101     ax3.spines[s].set_visible(False)
102
103 # adjusting axis labels
104 ax3.set_yticks([])
105 ax3.set_ylabel('')
106 ax3.set_xlabel('Purchase Amount',fontweight = 'bold',fontsize = 12)
107
108 #setting title for visual
109 ax3.set_title('Purchase Amount Distribution by Marital_Status',{'font':'serif', 'size':15,'weigh
110
111 plt.show()
```

## Marital_Status-Based Purchase Amount Distribution

| $3.01 Billion | $2.09 Billion |
|---|---|
| Unmarried | Married |

### Average Purchase Amount per Transaction



### Marital_Status-Based Transaction Distribution



### Purchase Amount Distribution by Marital_Status



## 🔍 Insights

### 1. Total Sales and Transactions Comparison

The total purchase amount and number of transactions by Unmarried customers was more than 20% the amount and transactions by married customers indicating that they had a more significant impact on the Black Friday sales.

> 2. Average Transaction Value

The average purchase amount per transaction was almost similar for married and unmarried customers ( $9261 vs 9266$ ).

> 3. Distribution of Purchase Amount

As seen above, the purchase amount for both married and unmarried customers is not normally distributed

## ˅ Confidence Interval Construction: Estimating Average Purchase Amount per Transaction

> 1. Step 1 - Building CLT Curve

As seen above, the purchase amount distribution is not Normal. So we need to use Central Limit Theorem. It states the distribution of sample means will approximate a normal distribution, regardless of the underlying population distribution
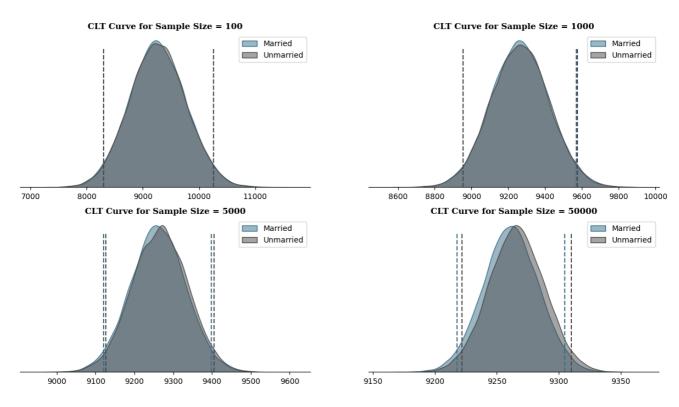
> 2. Step 2 - Building Confidence Interval

After building CLT curve, we will create a confidence interval predicting population mean at 95% Confidence level. Note - We will use different sample sizes of [100,1000,5000,50000]

```python
1  #defining a function for plotting the visual for given confidence interval
2
3  def plot(ci):
4
5      #setting the plot style
6      fig = plt.figure(figsize = (15,8))
7      gs = fig.add_gridspec(2,2)
8
9      #creating separate data frames
10     df_married = df.loc[df['Marital_Status'] == 'Married','Purchase']
11     df_unmarried = df.loc[df['Marital_Status'] == 'Unmarried','Purchase']
12
13     #sample sizes and corresponding plot positions
14     sample_sizes = [(100,0,0),(1000,0,1),(5000,1,0),(50000,1,1)]
15
16     #number of samples to be taken from purchase amount
17     bootstrap_samples = 20000
18
19     married_samples = {}
20     unmarried_samples = {}
21
22     for i,x,y in sample_sizes:
23         married_means = [] #list for collecting the means of married sample
24         unmarried_means = [] #list for collecting the means of unmarried sample
25
26         for j in range(bootstrap_samples):
27
28             #creating random 5000 samples of i sample size
29             married_bootstrapped_samples = np.random.choice(df_married,size = i)
30             unmarried_bootstrapped_samples = np.random.choice(df_unmarried,size = i)
31
32             #calculating mean of those samples
33             married_sample_mean = np.mean(married_bootstrapped_samples)
34             unmarried_sample_mean = np.mean(unmarried_bootstrapped_samples)
35
36             #appending the mean to the list
37             married_means.append(married_sample_mean)
38             unmarried_means.append(unmarried_sample_mean)
39
40         #storing the above sample generated
41         married_samples[f'{ci}%_{i}'] = married_means
42         unmarried_samples[f'{ci}%_{i}'] = unmarried_means
43
44         #creating a temporary dataframe for creating kdeplot
45         temp_df = pd.DataFrame(data = {'married_means':married_means,'unmarried_means':unmarried
46
47                                                      #plotting kdeplots
48         #plot position
49         ax = fig.add_subplot(gs[x,y])
50
51         #plots for married and unmarried
52         sns.kdeplot(data = temp_df,x = 'married_means',color ="#3A7089" ,fill = True, alpha = 0.
53         sns.kdeplot(data = temp_df,x = 'unmarried_means',color ="#4b4b4c" ,fill = True, alpha =
54
55         #calculating confidence intervals for given confidence level(ci)
56         m_range = confidence_interval(married_means,ci)
57         u_range = confidence_interval(unmarried_means,ci)
58
59         #plotting confidence interval on the distribution
60         for k in m_range:
61             ax.axvline(x = k,ymax = 0.9, color ="#3A7089",linestyle = '--')
```

```python
62
63          for k in u_range:
64              ax.axvline(x = k,ymax = 0.9, color ="#4b4b4c",linestyle = '--')
65
66
67          #removing the axis lines
68          for s in ['top','left','right']:
69              ax.spines[s].set_visible(False)
70
71          # adjusting axis labels
72          ax.set_yticks([])
73          ax.set_ylabel('')
74          ax.set_xlabel('')
75
76          #setting title for visual
77          ax.set_title(f'CLT Curve for Sample Size = {i}',{'font':'serif', 'size':11,'weight':'bol
78
79          plt.legend()
80
81      #setting title for visual
82      fig.suptitle(f'{ci}% Confidence Interval',font = 'serif', size = 18, weight = 'bold')
83
84      plt.show()
85
86      return married_samples,unmarried_samples
```

```python
1 m_samp_95,u_samp_95 = plot(95)
2
```

## 95% Confidence Interval



🔍 Insights

1. Sample Size

The analysis highlights the importance of sample size in estimating population parameters. It suggests that as the sample size increases, the confidence intervals become narrower and more precise. In business, this implies that larger sample sizes can provide more reliable insights and estimates.

2. Confidence Intervals

From the above analysis, we can see that the confidence interval overlap for all the sample sizes. This means that there is no statistically significant difference between the average spending per transaction for married and unmarried customers within the given samples.

3. Population Average

We are 95% confident that the true population average for married customers falls between $9,217$ and $9,305$, and for unmarried customers, it falls between $9,222$ and $9,311$.

4. Both the customers spend equal

The overlapping confidence intervals of average spending for married and unmarried customers indicate that both married and unmarried customers spend a similar amount per transaction. This implies a resemblance in spending behavior between the two groups

## 📅 Customer Age VS 💰 Purchase Amount.

## 📊 Data Visualization

```
1  #creating a df for purchase amount vs age group
2  temp = df.groupby('Age')['Purchase'].agg(['sum','count']).reset_index()
3
4  #calculating the amount in billions
5  temp['sum_in_billions'] = round(temp['sum'] / 10**9,2)
6
7  #calculationg percentage distribution of purchase amount
8  temp['%sum'] = round(temp['sum']/temp['sum'].sum(),3)
9
10 #calculationg per purchase amount
11 temp['per_purchase'] = round(temp['sum']/temp['count'])
12
13 temp
```

|   | Age | sum | count | sum_in_billions | %sum | per_purchase |
|---|-----|-----|-------|-----------------|------|--------------|
| 0 | 0-17 | 134913183 | 15102 | 0.13 | 0.026 | 8933.0 |
| 1 | 18-25 | 913848675 | 99660 | 0.91 | 0.179 | 9170.0 |
| 2 | 26-35 | 2031770578 | 219587 | 2.03 | 0.399 | 9253.0 |
| 3 | 36-45 | 1026569884 | 110013 | 1.03 | 0.201 | 9331.0 |
| 4 | 46-50 | 420843403 | 45701 | 0.42 | 0.083 | 9209.0 |
| 5 | 51-55 | 367099644 | 38501 | 0.37 | 0.072 | 9535.0 |
| 6 | 55+ | 200767375 | 21504 | 0.20 | 0.039 | 9336.0 |

```python
1 #setting the plot style
2 fig = plt.figure(figsize = (20,14))
3 gs = fig.add_gridspec(3,1,height_ratios =[0.10,0.4,0.5])
4
5
```