

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5
6

```

## ✓ Table of Contents

### 1. [Step 1: Understanding the Problem and the Data](#)

### 2. [Step 2: Import and Inspect the Data](#)

- [Statistical summary for numerical columns and Unique values for categorical data](#)
- [Step 3: Handling Missing Values](#)

### 3. [Step 4: Explore Data Characteristics](#)

- [Univariate and Bivariate Analysis](#)
  - [Univariate analysis of User\\_ID and Purchase column](#)
  - [Univariate analysis of Product\\_ID column](#)
  - [Bivariate analysis of Product and Gender](#)
  - [Bivariate analysis of Product and Marital Status](#)
  - [Bivariate analysis of Product Category and Gender, Purchase](#)
  - [Univariate analysis of Gender column](#)
  - [Bivariate analysis of Gender and Purchase column](#)
  - [Univariate analysis of Age column](#)
  - [Bivariate analysis of Age and Gender column](#)
  - [Univariate Analysis of Occupation](#)
  - [Bivariate Analysis of Occupation and Age Group](#)
  - [Bivariate Analysis of City\\_Category and Purchase](#)
  - [Bivariate Analysis of City\\_Category and User\\_ID](#)
  - [Bivariate Analysis of City\\_Category and Gender](#)
  - [Bivariate Analysis of City\\_Category and Age Group](#)
  - [Bivariate Analysis of City\\_Category and Occupation](#)
  - [Bivariate Analysis of City\\_Category and Product Category](#)
  - [Bivariate Analysis of Stay\\_In\\_Current\\_City\\_Years and User\\_ID](#)
  - [Bivariate Analysis of Stay\\_In\\_Current\\_City\\_Years and Stay Duration](#)
  - [Bivariate Analysis of Marital\\_Status and Gender](#)
  - [Bivariate Analysis of Marital\\_Status and Age Category](#)
  - [Bivariate Analysis of Marital\\_Status and Occupation](#)
  - [Bivariate Analysis of Marital\\_Status and City](#)
  - [Bivariate Analysis of Marital\\_Status and Stay Duration](#)
  - [Univariate Analysis of Purchase column](#)
  -

### 4. [Answering Questions asked in Case Study](#)

- [Are women spending more money per transaction than men? Why or Why not](#)
- [Confidence intervals and distribution of the mean of the expenses by female and male customers](#)
  - [Effect of sample size on the confidence interval](#)
  - [Effect of significance level on width of the confidence interval](#)
- [Are confidence intervals of average male and female spending overlapping?](#)
- [Are confidence intervals of Married vs Unmarried spending overlapping?](#)
- [Are confidence intervals of spending habits of different Age Groups overlapping?](#)

### 5. [Business Insights](#)

## 6. [Recommendations](#)

### LEGEND

Analysis Steps in Red color

Sub Headings in orange color

These are insights got from analysing the data, in Blue color

""" code sections are notes to self, some quick code snippets, or some notes while doing analysis

```
1 ''' Notes to Self
2 my markdown templates
3 <span style="font-size:50px; font-family:Arial;color:red">Use this for heading/span><li>
4 <span style="font-size:30px; font-family:Arial;color:orange">Sub Headings in orange color</span><li>
5 <span style="font-size:25px; font-family:Arial;color:#8DDBF2">Use this for insight </span><li>
6 '''
```

```
➦ ' Notes to Self\nmy markdown templates\n<span style="font-size:50px; font-family:Arial;color:red">Use this for heading/span>
<li>\n<span style="font-size:30px; font-family:Arial;color:orange">Sub Headings in orange color</span><li>\n<span style="font-
size:25px; font-family:Arial;color:#8DDBF2">Use this for insight </span><li>\n'
```

## ✓ 1.0

### Step 1: Understanding the Problem and the Data

#### About Walmart

Walmart is an American multinational retail corporation that operates a chain of supercenters, discount departmental stores, and grocery stores from the United States. Walmart has more than 100 million customers worldwide.

#### 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).

#### Key points:

1. WHAT: Analyse customer purchase behavior, focussing on the purchase amount wrt gender and other factors
2. WHY: Identify the spending habits of Men and Women to see if we can gain any insights out of the analysis

Dataset The company collected the transactional data of customers who purchased products from the Walmart Stores during Black Friday.

The dataset has the following features: Dataset link: [Walmart\\_data.csv](#)

- User\_ID: User ID
- Product\_ID: Product ID
- Gender: Sex of User
- Age: Age in bins
- Occupation: Occupation(Masked)
- City\_Category: Category of the City (A,B,C)
- StayInCurrentCityYears: Number of years stay in current city
- Marital\_Status: Marital Status
- ProductCategory: Product Category (Masked)
- Purchase: Purchase Amount

## ✓ 2.0

### Step 2: Import and Inspect the Data

```

1 '''
2 Notes to Self
3 During this step, looking into the statistics is critical to gain initial know-how of its structure, variable kinds, and capabi
4
5 '''

```

```

↗ ' \nNotes to Self\nDuring this step, looking into the statistics is critical to gain initial know-how of its structure,
variable kinds, and capability issues.\n\n'

```

```

1 df =pd.read_csv("walmart_data.csv")
2 df.head()

```

```
↗
```

	User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category	Pi
0	1000001	P00069042	F	0-17	10	A	2	0		3
1	1000001	P00248942	F	0-17	10	A	2	0		1
2	1000001	P00087842	F	0-17	10	A	2	0		12

```
1 df.columns
```

```

↗ Index(['User_ID', 'Product_ID', 'Gender', 'Age', 'Occupation', 'City_Category',
       'Stay_In_Current_City_Years', 'Marital_Status', 'Product_Category',
       'Purchase'],
       dtype='object')

```

## 2.1

Statistical summary for numerical columns and Unique values for categorical data

```

1 #Instead of manually creating tables and collecting values
2 #lets automate it.
3 #We will generate a table of all columns in Markdown format
4 #We will paste this table into a markdown cell and then we may add additional analysis
5 #!pip install tabulate
6 from tabulate import tabulate
7 def generate_markdown_table_with_details(df):
8     column_details = []
9     for col in df.columns:
10         dtype = str(df[col].dtype)
11         if df[col].dtype == 'object':
12             cat_or_num = 'Categorical'
13             details = df[col].unique().tolist()
14             if(len(details)>10):
15                 details= f'{df[col].nunique()} unique values'
16         else:
17             cat_or_num = 'Numerical'
18             details = f"Min: {df[col].min()}, Max: {df[col].max()}"
19         column_details.append([col, dtype, cat_or_num, details, ''])
20
21 # Create a DataFrame for the table
22 table_df = pd.DataFrame(column_details, columns=['Column Name', 'Data Type', 'Categorical or Numerical', 'Details', 'Busine
23
24 # Convert DataFrame to Markdown
25 markdown_table = tabulate(table_df, headers='keys', tablefmt='pipe', showindex=False)
26 return markdown_table
27
28 markdown_table = generate_markdown_table_with_details(df)
29 print(f'There are a total of {df.shape[0]} rows and {df.shape[1]} columns')
30 print(markdown_table)
31

```

```

↗ There are a total of 550068 rows and 10 columns
| Column Name | Data Type | Categorical or Numerical | Details |
|:-----|:-----|:-----|:-----|

```

User_ID	int64	Numerical	Min: 1000001, Max: 1006040
Product_ID	object	Categorical	3631 unique values
Gender	object	Categorical	['F', 'M']
Age	object	Categorical	['0-17', '55+', '26-35', '46-50', '51-55', '36-45', '18-25']
Occupation	int64	Numerical	Min: 0, Max: 20
City_Category	object	Categorical	['A', 'C', 'B']
Stay_In_Current_City_Years	object	Categorical	['2', '4+', '3', '1', '0']
Marital_Status	int64	Numerical	Min: 0, Max: 1
Product_Category	int64	Numerical	Min: 1, Max: 20
Purchase	int64	Numerical	Min: 12, Max: 23961

```
1 '''
2 Notes to Self
3 Adding insights about columns and data size.
4 We identify which columns are numerical, which are categorical
5 Within categorical, we identify if it is nominal or ordinal data. This helps us to decide best graph to use
6 Also ordinal data can be converted into numeric using codes thereby making it available for some kind of graphs.
7 '''
```

```
'''
Notes to Self
Adding insights about columns and data size.
We identify which columns are numerical, which are categorical
Within categorical, we identify if it is nominal or ordinal data. This helps us to decide best graph to use
Also ordinal data can be converted into numeric using codes thereby making it available for some kind of graphs.'''
```

There are a total of 550068 rows and 10 columns

Column Name	Data Type	Categorical or Numerical	Details	Business Description(fill manually)
User_ID	int64	Cat(Nominal)	Min: 1000001, Max: 1006040	
Product_ID	object	Cat(Nominal)	3631 unique values	
Gender	object	Cat(Nominal)	['F', 'M']	
Age	object	Cat(Ordinal)	['0-17', '55+', '26-35', '46-50', '51-55', '36-45', '18-25']	Age in Bins
Occupation	int64	Cat(Nominal)	Min: 0, Max: 20	Has been encoded into numerical values
City_Category	object	Cat(Nominal)	['A', 'C', 'B']	Category of the City (A,B,C)
Stay_In_Current_City_Years	object	Cat(Ordinal)	['2', '4+', '3', '1', '0']	Number of years stay in current city
Marital_Status	int64	Cat(Nominal)	Min: 0, Max: 1	Has been encoded into numerical values
Product_Category	int64	Cat(Nominal)	Min: 1, Max: 20	Has been encoded into numerical values
Purchase	int64	Numerical	Min: 12, Max: 23961	

There is only ONE column that has real numerical data, the purchase column, rest all columns are categorical in nature

```
1 '''
2 Notes to Self
3 Let's get a quick summary of the dataset using the pandas describe() method.
4 The describe() function applies basic statistical computations on the dataset like extreme values, count of data points standar
5 Any missing value or NaN value is automatically skipped. describe() function gives a good picture of the distribution of data.
6 '''
7 df.describe(include ='all')
```

	User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category
count	5.500680e+05	550068	550068	550068	550068.000000	550068	550068	550068.000000	550068
unique	NaN	3631	2	7	NaN	3	5	NaN	5
top	NaN	P00265242	M	26-35	NaN	B	1	NaN	5
freq	NaN	1880	414259	219587	NaN	231173	193821	NaN	5
mean	1.003029e+06	NaN	NaN	NaN	8.076707	NaN	NaN	0.409653	5
std	1.727592e+03	NaN	NaN	NaN	6.522660	NaN	NaN	0.491770	5
min	1.000001e+06	NaN	NaN	NaN	0.000000	NaN	NaN	0.000000	5
25%	1.001516e+06	NaN	NaN	NaN	2.000000	NaN	NaN	0.000000	5
50%	1.003077e+06	NaN	NaN	NaN	7.000000	NaN	NaN	0.000000	5
75%	1.004478e+06	NaN	NaN	NaN	14.000000	NaN	NaN	1.000000	5
max	1.006040e+06	NaN	NaN	NaN	20.000000	NaN	NaN	1.000000	5

## 2.2

### Step 3: Handling Missing Values

```

1 '''
2 Notes to Self
3 Missing Data can also refer to as NA(Not Available) values in pandas. There are several useful functions for detecting, removing, and replacing null values in Pandas DataFrame
4
5 isnull()
6 notnull()
7 dropna()
8 fillna()
9 replace()
10 interpolate()
11 '''

```

```

'''
Notes to Self
Missing Data can also refer to as NA(Not Available) values in pandas. There are several useful functions for detecting, removing, and replacing null values in Pandas DataFrame
:
isnull()
notnull()
dropna()
fillna()
replace()
interpolate()
'''

```

```

1 #Now let's check if there are any missing values in our dataset or not.
2 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

```

Awesome! the data does not have any missing values

Lets add some columns for analysing the categorical and numerical data better

- For categorical data- we will add some codes
- For numerical data -we will add some labels after binning the data

```

1 #REMOVE IF NOT NEEDED
2
3 # Create a new column with the binned values

```

```
4 # df['IncomeGroup'] = pd.cut(df['Income'], bins=income_bin_edges, labels=["Less than 60k", "More than 60K"])
5 # df['AgeGroup'] = pd.cut(df['Age'], bins=age_bin_edges, labels=["Less than 24", "24-35", "36 and above"])
```

## 3.0

### Step 4: Explore Data Characteristics

```
1 '''
2 Notes to self
3 By exploring the characteristics of your information very well, you can gain treasured insights into its structure, pick out ca
4 Let's start by exploring the data according to the dataset.
5
6 for categorical columns, use
7 -----
8 sns.countplot(data=df, x='Product')
9 df.Product.value_counts().plot.bar()
10 df['Product'].value_counts().plot.pie(autopct='%1.1f%%', startangle=90, colors=sns.color_palette('pastel'))
11
12
13
14 for numerical columns use
15 -----
16 df.hist() # Quickly create histograms for all numeric columns
17 sns.histplot(data=df["Age"], kde=True)
18 df.Age.plot.hist() #hist for Specific column
19 df.Age.plot.kde()
20 df.Age.plot.box()
21
22 df.plot.box()#Print boxplots of all numerical columns in one go
23
24 sns.violinplot(data=df , x='Age')
25
26
27
28
29 '''

'''
Notes to self
By exploring the characteristics of your information very well, you can gain treasured insights into its
structure, pick out capability problems or anomalies, and inform your subsequent evaluation and modeling choices. Documenting
any findings or observations from this step is critical, as they may be relevant for destiny reference or communication with
stakeholders.
Let's start by exploring the data according to the dataset.

for categorical columns, use
-----
sns.countplot(data=df,
x='Product')
df.Product.value_counts().plot.bar()
df['Product'].value_counts().plot.pie(autopct='%1.1f%%',
startangle=90, colors=sns.color_palette('pastel'))

for numerical columns use
-----
df.hist()
# Quickly create histograms for all numeric columns
sns.histplot(data=df["Age"], kde=True)
df.Age.plot.hist() #hist for
Specific column
df.Age.plot.kde()
df.Age.plot.box()
df.plot.box()#Print boxplots of all numerical columns in one
go
sns.violinplot(data=df , x='Age')
'''
```

## 3.1

Univariate and Bivariate Analysis (Analysis of one attribute at a time.)

### 3.1.1

Analysis of User\_ID and Purchase column

```
1 frequentUsers = df.groupby('User_ID').Purchase.count().reset_index()
2 frequentUsers.describe()
3
4
```



	User_ID	Purchase
count	5.891000e+03	5891.000000
mean	1.003025e+06	93.374300
std	1.743379e+03	107.190049
min	1.000001e+06	6.000000
25%	1.001518e+06	26.000000
50%	1.003026e+06	54.000000
75%	1.004532e+06	117.000000
max	1.006040e+06	1026.000000

We have data of 5891 customers


- Each user has shopped multiple times, from a minimum of 6 to a maximum of 1026 times
- On average, each user has shopped for 93 times.
- This is a very decent number to be able to quantify some individual user's buying habit

```
1 # The user data in the original dataset can be separated out as a separate dataframe so that we can run analysis on users indep
2 users = df[['User_ID','Gender', 'Age', 'Occupation', 'City_Category', 'Stay_In_Current_City_Years', 'Marital_Status','Purchase']]
3 users
```



	User_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Purchase		
								count	sum	mean
0	1000001	F	0-17	10	A	2	0	35	334093	9545.514286
1	1000002	M	55+	16	C	4+	0	77	810472	10525.610390
2	1000003	M	26-35	15	A	3	0	29	341635	11780.517241
3	1000004	M	46-50	7	B	2	1	14	206468	14747.714286
4	1000005	M	26-35	20	A	1	1	106	821001	7745.292453
...	...	...	...	...	...	...	...	...	...	...
5886	1006036	F	26-35	15	B	4+	1	514	4116058	8007.894942
5887	1006037	F	46-50	1	C	4+	0	122	1119538	9176.540984
5888	1006038	F	55+	1	C	0	0	10	88884	7500.000000


```
1 # Flatten the column names
2 users.columns = [''.join(col).strip() for col in users.columns]
3 users
```



	User_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Purchasecount	Purchasesum	I
0	1000001	F	0-17	10	A	2	0	35	334093	
1	1000002	M	55+	16	C	4+	0	77	810472	
2	1000003	M	26-35	15	A	3	0	29	341635	
3	1000004	M	46-50	7	B	2	1	14	206468	
4	1000005	M	26-35	20	A	1	1	106	821001	
...	...	...	...	...	...	...	...	...	...	...
5886	1006036	F	26-35	15	B	4+	1	514	4116058	
5887	1006037	F	46-50	1	C	4+	0	122	1119538	
5888	1006038	F	55+	1	C	2	0	12	90034	
5889	1006039	F	46-50	0	B	4+	1	74	590319	
5890	1006040	M	26-35	6	B	2	0	180	1653299	

5891 rows × 10 columns

```
1 users.describe()
```



	User_ID	Occupation	Marital_Status	Purchasecount	Purchasesum	Purchasemean
count	5.891000e+03	5891.000000	5891.000000	5891.000000	5.891000e+03	5891.000000
mean	1.003025e+06	8.153285	0.419963	93.374300	8.650166e+05	9568.839914
std	1.743379e+03	6.323140	0.493594	107.190049	9.436445e+05	1890.087105
min	1.000001e+06	0.000000	0.000000	6.000000	4.668100e+04	2318.733333
25%	1.001518e+06	3.000000	0.000000	26.000000	2.376780e+05	8287.212366
50%	1.003026e+06	7.000000	0.000000	54.000000	5.212130e+05	9386.208333
75%	1.004532e+06	14.000000	1.000000	117.000000	1.119250e+06	10654.633199
max	1.006040e+06	20.000000	1.000000	1026.000000	1.053691e+07	18577.893617

Average spend per user is 9568, with maximum spend at 18577 and minimum at 2318

3.1.2

Analysis of Product\_ID column

```
1 products = df.groupby(['Product_Category', 'Product_ID'], as_index=False).agg(  
2     Revenue=('Purchase', 'sum'),  
3     Quantity=('Purchase', 'count'),  
4     Male_count=('Gender', lambda x: (x == 'M').sum()),  
5     Female_count=('Gender', lambda x: (x == 'F').sum()),  
6     Married_count=('Marital_Status', lambda x: (x == 1).sum()),  
7     Unmarried_count=('Marital_Status', lambda x: (x == 0).sum())  
8 )  
9 products
```





	Product_Category	Product_ID	Revenue	Quantity	Male_count	Female_count	Married_count	Unmarried_count
0	1	P00000642	7635578	512	441	71	203	309
1	1	P00000942	581125	55	48	7	17	38
2	1	P00001042	6922380	503	427	76	193	310
3	1	P00001542	640313	69	58	11	35	34
4	1	P00002042	889731	93	86	7	33	60
...	...	...	...	...	...	...	...	...
3626	19	P00370293	28790	785	543	242	331	454
3627	19	P00370853	30588	818	609	209	326	492
3628	20	P00371644	326257	899	643	256	375	524
3629	20	P00372445	313817	837	603	234	345	492
3630	20	P00375436	304653	814	581	233	350	464

3631 rows × 8 columns

1 Start coding or generate with AI.


```
1 products['Price'] = round(products.Revenue/products.Quantity,0)
2 products['MalePercentage'] = round(products.Male_count*100/(products.Male_count+products.Female_count),0)
3 products['UnmarriedPercentage'] = round(products.Unmarried_count*100/(products.Unmarried_count+products.Married_count),0)
4 products
```



	Product_Category	Product_ID	Revenue	Quantity	Male_count	Female_count	Married_count	Unmarried_count	Price	MalePe
0	1	P00000642	7635578	512	441	71	203	309	14913.0	
1	1	P00000942	581125	55	48	7	17	38	10566.0	
2	1	P00001042	6922380	503	427	76	193	310	13762.0	
3	1	P00001542	640313	69	58	11	35	34	9280.0	
4	1	P00002042	889731	93	86	7	33	60	9567.0	
...	...	...	...	...	...	...	...	...	...	...
3626	19	P00370293	28790	785	543	242	331	454	37.0	
3627	19	P00370853	30588	818	609	209	326	492	37.0	
3628	20	P00371644	326257	899	643	256	375	524	363.0	
3629	20	P00372445	313817	837	603	234	345	492	375.0	
3630	20	P00375436	304653	814	581	233	350	464	374.0	

3631 rows × 11 columns


```
1 products.sort_values(by='Revenue', ascending=False)
```



	Product_Category	Product_ID	Revenue	Quantity	Male_count	Female_count	Married_count	Unmarried_count	Price	MaleP
31	1	P00025442	27995166	1615	1267	348	652	963	17334.0	
129	1	P00110742	26722309	1612	1247	365	648	964	16577.0	
3545	16	P00255842	25168963	1383	1008	375	535	848	18199.0	
1807	6	P00059442	24338343	1406	1048	358	577	829	17310.0	
253	1	P00184942	24334887	1440	1141	299	597	843	16899.0	
...	...	...	...	...	...	...	...	...	...	...
853	5	P00012942	1717	1	1	0	0	1	1717.0	
3287	11	P00325342	1656	1	1	0	1	0	1656.0	
3312	11	P00353042	1545	1	1	0	1	0	1545.0	
3355	12	P00309042	726	1	0	1	1	0	726.0	
3375	13	P00091742	405	1	1	0	0	1	405.0	

3631 rows × 11 columns

```
1 products.describe()
```



	Product_Category	Revenue	Quantity	Male_count	Female_count	Married_count	Unmarried_count	Price	Male
count	3631.000000	3.631000e+03	3631.000000	3631.000000	3631.000000	3631.000000	3631.000000	3631.000000	3631.000000
mean	6.484164	1.403419e+06	151.492151	114.089507	37.402644	62.059212	89.432939	7874.911319	7874.911319
std	3.801726	2.647138e+06	212.852932	163.547026	53.307847	86.409797	127.347006	3873.754005	3873.754005
min	1.000000	4.050000e+02	1.000000	0.000000	0.000000	0.000000	0.000000	37.000000	37.000000
25%	5.000000	1.152100e+05	19.500000	13.000000	5.000000	9.000000	11.000000	5386.000000	5386.000000
50%	6.000000	4.378210e+05	71.000000	51.000000	17.000000	29.000000	41.000000	6943.000000	6943.000000
75%	8.000000	1.517297e+06	194.000000	142.000000	49.000000	81.000000	113.000000	10154.500000	10154.500000
max	20.000000	2.799517e+07	1880.000000	1372.000000	508.000000	793.000000	1087.000000	21257.000000	21257.000000

- There are 3631 products
- Average product price is 7874, *with maximum at 21257* and minimum at \$37

```
1
2 top1PercentProducts = products[:int(round(len(products)*0.01,0)]]
3 top1PercentProducts
```

	Product_Category	Product_ID	Revenue	Quantity	Male_count	Female_count	Married_count	Unmarried_count	Price	MaleF
0	1	P00000642	7635578	512	441	71	203	309	14913.0	
1	1	P00000942	581125	55	48	7	17	38	10566.0	
2	1	P00001042	6922380	503	427	76	193	310	13762.0	
3	1	P00001542	640313	69	58	11	35	34	9280.0	
4	1	P00002042	889731	93	86	7	33	60	9567.0	
5	1	P00002142	10424160	735	578	157	315	420	14183.0	
6	1	P00004242	1269862	113	95	18	38	75	11238.0	
7	1	P00004442	1632933	147	127	20	53	94	11108.0	
8	1	P00006942	6165491	498	425	73	167	331	12381.0	
9	1	P00007042	351139	44	40	4	20	24	7980.0	
10	1	P00008842	1299569	109	93	16	38	71	11923.0	
11	1	P00009342	3964157	351	288	63	127	224	11294.0	
12	1	P00009642	24539	3	3	0	1	2	8180.0	
13	1	P00010742	22164153	1350	1069	281	574	776	16418.0	
14	1	P00010942	4073491	305	265	40	110	195	13356.0	
15	1	P00013542	15836	1	0	1	0	1	15836.0	
16	1	P00014042	2471491	194	175	19	75	119	12740.0	
17	1	P00014842	4255327	412	329	83	138	274	10328.0	
18	1	P00015542	2644080	295	247	48	107	188	8963.0	
19	1	P00015842	3579328	267	225	42	100	167	13406.0	
20	1	P00016042	6406520	447	384	63	168	279	14332.0	
21	1	P00016342	2168552	191	138	53	79	112	11354.0	
22	1	P00016542	1883681	159	132	27	52	107	11847.0	
23	1	P00016742	2710920	266	232	34	101	165	10191.0	
24	1	P00016842	2902550	315	280	35	122	193	9214.0	
25	1	P00017642	333018	51	42	9	17	34	6530.0	
26	1	P00019342	4330757	330	284	46	120	210	13124.0	
27	1	P00019942	1224440	107	95	12	42	65	11443.0	
28	1	P00020342	3649188	358	289	69	130	228	10193.0	
29	1	P00022142	785570	75	69	6	27	48	10474.0	
30	1	P00024542	136554	10	8	2	6	4	13655.0	
31	1	P00025442	27995166	1615	1267	348	652	963	17334.0	
32	1	P00027842	12128	1	0	1	1	0	12128.0	
33	1	P00028042	4597211	376	314	62	141	235	12227.0	
34	1	P00028442	5332834	418	335	83	159	259	12758.0	
35	1	P00029542	2545638	218	198	20	85	133	11677.0	

- The highest grossing Top 1% products are all from Product\_Category 1.

### 3.1.3

Bivariate analysis of Product and Gender

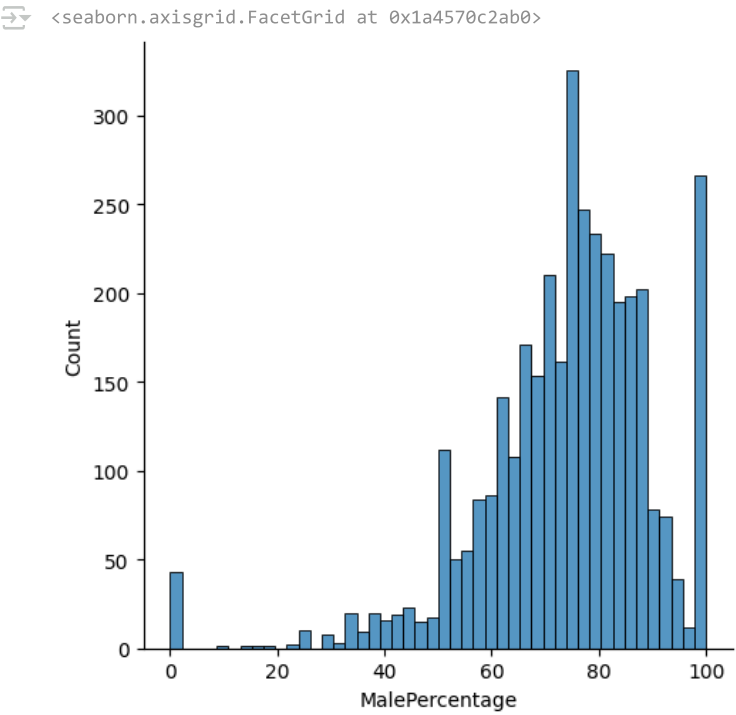
```
1 top1MalePreferredProducts = products.sort_values(by='MalePercentage', ascending=False)
2 top1MalePreferredProducts
```

↗

	Product_Category	Product_ID	Revenue	Quantity	Male_count	Female_count	Married_count	Unmarried_count	Price	MalePe
2734	8	P00262742	11882	2	2	0	1	1	5941.0	
1954	7	P00135942	16954	1	1	0	1	0	16954.0	
1966	7	P00194942	55208	4	4	0	3	1	13802.0	
1965	7	P00188142	25185	2	2	0	0	2	12592.0	
1964	7	P00172242	54146	3	3	0	2	1	18049.0	
...	...	...	...	...	...	...	...	...	...	...
1400	5	P00231642	7181	1	0	1	0	1	7181.0	
3419	14	P00152042	52338	4	0	4	1	3	13084.0	
1726	5	P00350742	6989	1	0	1	0	1	6989.0	
90	1	P00060842	15737	1	0	1	0	1	15737.0	
1583	5	P00301942	7147	1	0	1	1	0	7147.0	

3631 rows × 11 columns

```
1 sns.displot(data=top1MalePreferredProducts.MalePercentage)
```



The breakup in percentages per product for male-females largely is in line with the overall count of 75%-25% across all products

```
1 top1MalePreferredProducts.query('MalePercentage>90').Product_ID.count()
```

↗

429

```
1 top1MalePreferredProducts.query('MalePercentage>99').Product_ID.count()
```

↗

264

There are 264 products bought exclusively by males

```
1 top1MalePreferredProducts.query('MalePercentage<1').Product_ID.count()
```

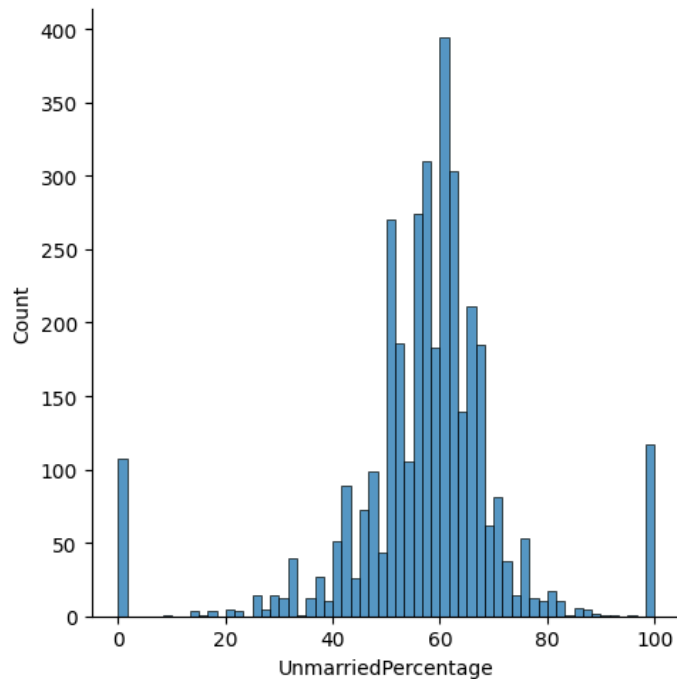
↗ 43

There are 43 products bought exclusively by females

analysis

```
1 sns.displot(data=products.UnmarriedPercentage)
```

↗ <seaborn.axisgrid.FacetGrid at 0x1a471f1a870>



### 3.1.4

Bivariate analysis of Product and Marital Status

```
1 products.query('UnmarriedPercentage>99').Product_ID.count()
```

↗ 117

There are 117 products bought exclusively by Unmarried customers

```
1 products.query('UnmarriedPercentage<1').Product_ID.count()
```

↗ 107

There are 107 products bought exclusively by Married customers

### 3.1.5

Bivariate analysis of Product Category and Gender. Purchase

```
1 product_categories = products.groupby('Product_Category', as_index=False).agg(
2     Revenue=('Revenue', 'sum'),
3     Male_count=('Male_count', 'sum'),
4     Female_count=('Female_count', 'sum'),
5     Married_count=('Married_count', 'sum'),
6     Unmarried_count=('Unmarried_count', 'sum')
```

```

7     )
8 product_categories.sort_values(by='Revenue', ascending=False, inplace=True)
9 product_categories
10

```




	Product_Category	Revenue	Male_count	Female_count	Married_count	Unmarried_count
0	1	1910013754	115547	24831	56003	84375
4	5	941835229	108972	41961	61277	89656
7	8	854318799	80367	33558	48514	65411
5	6	324150302	15907	4559	8327	12139
1	2	268516186	18206	5658	9726	14138
2	3	204084713	14207	6006	7854	12359
15	16	145120612	7426	2402	4115	5713
10	11	113791115	19548	4739	9619	14668
9	10	100837301	3963	1162	2347	2778
14	15	92969042	5244	1046	2667	3623
6	7	60896731	2778	943	1681	2040
3	4	27380488	8114	3639	4576	7177
13	14	20014696	900	623	677	846
17	18	9290201	2743	382	1484	1641
8	9	6370324	340	70	163	247
16	17	5878699	516	62	280	298
11	12	5331844	2415	1532	1913	2034
12	13	4008601	4087	1462	2387	3162
19	20	944727	1827	723	1070	1480
18	19	59378	1152	451	657	946

```

1 product_categories['MalePercentage'] = round(product_categories.Male_count*100/(product_categories.Male_count+product_categories.Female_count),2)
2 product_categories['UnmarriedPercentage'] = round(product_categories.Unmarried_count*100/(product_categories.Unmarried_count+product_categories.Married_count),2)
3 product_categories

```



	Product_Category	Revenue	Male_count	Female_count	Married_count	Unmarried_count	MalePercentage	UnmarriedPercentage
0	1	1910013754	115547	24831	56003	84375	82.0	60.0
4	5	941835229	108972	41961	61277	89656	72.0	59.0
7	8	854318799	80367	33558	48514	65411	71.0	57.0
5	6	324150302	15907	4559	8327	12139	78.0	59.0
1	2	268516186	18206	5658	9726	14138	76.0	59.0
2	3	204084713	14207	6006	7854	12359	70.0	61.0
15	16	145120612	7426	2402	4115	5713	76.0	58.0
10	11	113791115	19548	4739	9619	14668	80.0	60.0
9	10	100837301	3963	1162	2347	2778	77.0	54.0
14	15	92969042	5244	1046	2667	3623	83.0	58.0
6	7	60896731	2778	943	1681	2040	75.0	55.0
3	4	27380488	8114	3639	4576	7177	69.0	61.0
13	14	20014696	900	623	677	846	59.0	56.0
17	18	9290201	2743	382	1484	1641	88.0	53.0
8	9	6370324	340	70	163	247	83.0	60.0
16	17	5878699	516	62	280	298	89.0	52.0
11	12	5331844	2415	1532	1913	2034	61.0	52.0
12	13	4008601	4087	1462	2387	3162	74.0	57.0
19	20	944727	1827	723	1070	1480	72.0	58.0
18	19	59378	1152	451	657	946	72.0	59.0

- Product Category 1 is the highest grossing category with USD 1.91 Billion in revenue
- For this category, 82% buyers are male. However, as overall we have 75% male customers for Walmart, it would be prudent to check if the 82% is statistically significantly different from the average.
- This category contains ALL the Top 1% revenue-grossing products for Walmart
- Categories 5,8,6,2,3,16,11,10 are the other 1 Billion plus categories
- Category 19 is the lowest grossing category with just \$59378 revenue

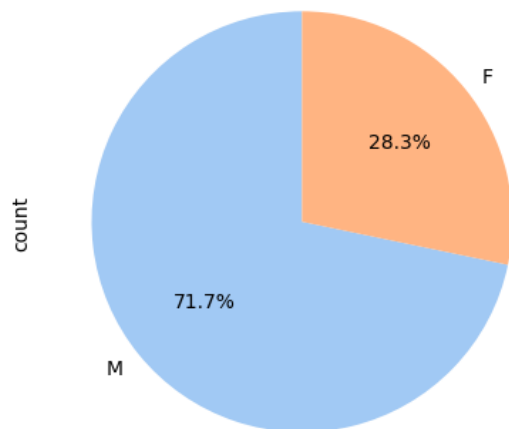
### 3.1.6

Univariate analysis of Gender column

```
1 users['Gender'].value_counts().plot.pie(
2 autopct='%1.1f%%',
3 startangle=90, colors=sns.color_palette('pastel'))
4 plt.title("Gender breakup in overall users")
```

```
Text(0.5, 1.0, 'Gender breakup in overall users')
```

Gender breakup in overall users



```
1 users['Gender'].value_counts()
```

```
Gender
M    4225
F    1666
Name: count, dtype: int64
```

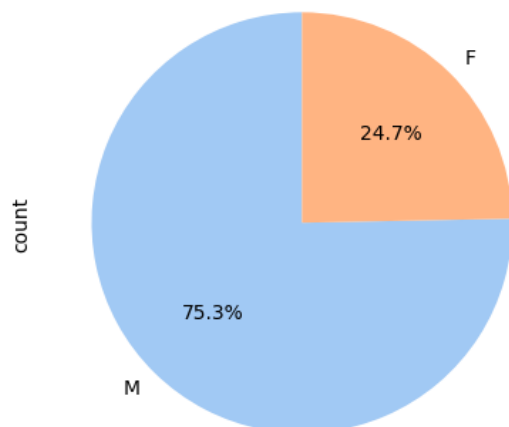
### 3.1.7

Bivariate analysis of Gender and Purchase column

```
1 df['Gender'].value_counts().plot.pie(
2 autopct='%1.1f%%',
3 startangle=90, colors=sns.color_palette('pastel'))
4 plt.title("Gender breakup in overall purchases")
```

```
Text(0.5, 1.0, 'Gender breakup in overall purchases')
```

Gender breakup in overall purchases



- There are 5891 unique users
- Out of these 4225(72%) are males, 1666 (28%) are females
- In terms of purchases too, 75% of purchases are by males and 25% are by females



```

1 males= df[df.Gender=='M']
2 females = df[df.Gender=='F']
3
4 PopMeanFemales = females.Purchase.mean()
5 PopMeanMales = males.Purchase.mean()
6 print(f'Population Mean purchase amount for females is ${PopMeanFemales:.2f}')
7 print(f'Population Mean purchase amount for males is ${PopMeanMales:.2f}')

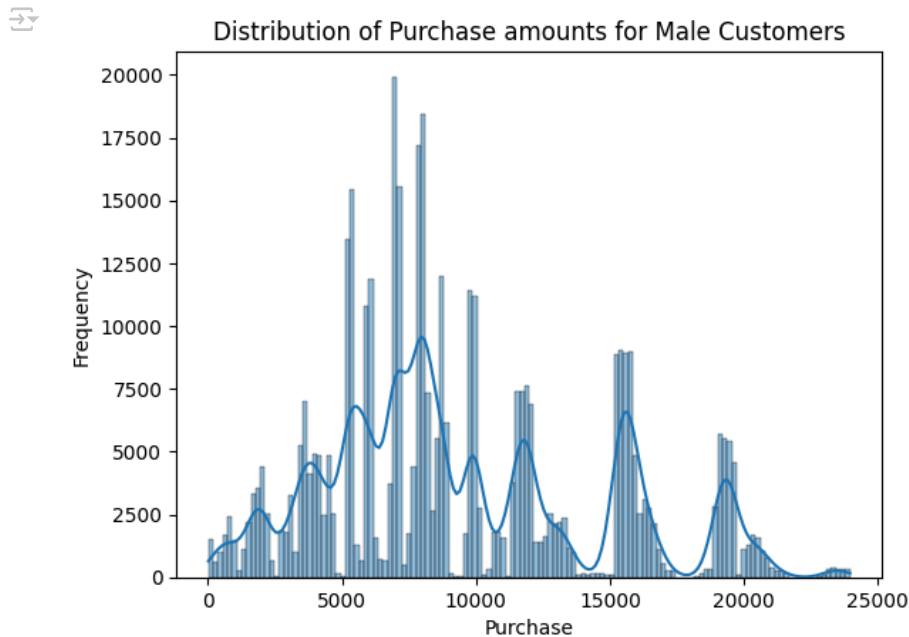
```

Population Mean purchase amount for females is \$8734.57  
Population Mean purchase amount for males is \$9437.53

```

1 # plot the sample_means_male to see the distribution
2 sns.histplot(males.Purchase, kde=True)
3 plt.xlabel('Purchase')
4 plt.ylabel('Frequency')
5 plt.title('Distribution of Purchase amounts for Male Customers')
6 plt.show()

```



- The distribution of purchase amounts for males is not really normal. However we can use CLT to estimate the mean value with some confidence.
- Instead of giving exact average, it is better to give an interval for the mean purchase amount.
- We can do this by using the Central limit theorem
- We will select a sample, and use the sample mean to calculate an interval

### 3.1.8

Univariate analysis of Age column

```

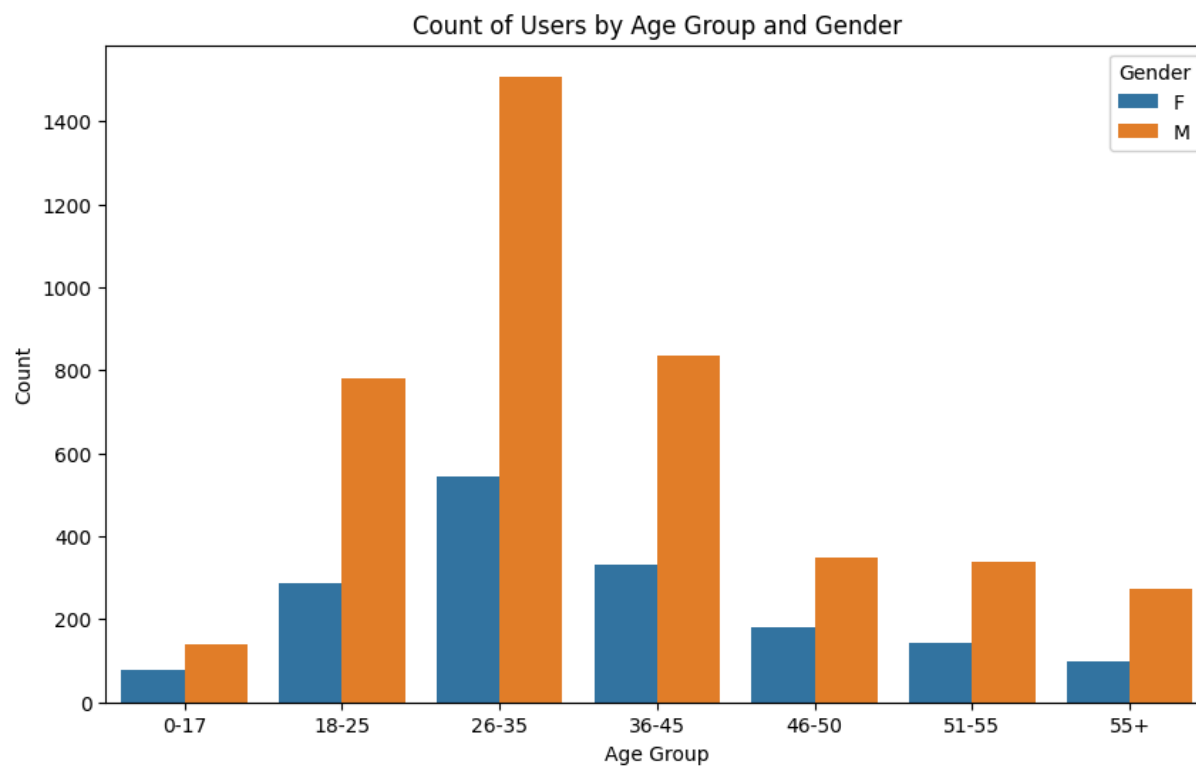
1 age = users.groupby(['Age', 'Gender'], as_index=False).size().sort_values(by="Age")
2 age

```



	Age	Gender	size
0	0-17	F	78
1	0-17	M	140
2	18-25	F	287
3	18-25	M	782
4	26-35	F	545
5	26-35	M	1508
6	36-45	F	333
7	36-45	M	834
8	46-50	F	182
9	46-50	M	349
10	51-55	F	142
11	51-55	M	339
12	55+	F	99
13	55+	M	273

```
1 # Univariate Analysis
2
3 plt.figure(figsize=(10, 6))
4 sns.barplot(data=age, x='Age', y='size', hue='Gender')
5 plt.title('Count of Users by Age Group and Gender')
6 plt.xlabel('Age Group')
7 plt.ylabel('Count')
8 plt.show()
```



### 3.1.9

Bivariate analysis of Age and Gender column

```

1 #select all records where gender is M
2
3 male_users = age[age['Gender'] == 'M']
4 male_users
5
6 # Insights from the age distribution plot
7 total_users = age['size'].sum()
8 male_users_count = male_users['size'].sum()
9 female_users_count = total_users - male_users_count
10
11 print(f"Total Users: {total_users}")
12 print(f"Male Users: {male_users_count} ({(male_users_count / total_users) * 100:.2f}%)")
13 print(f"Female Users: {female_users_count} ({(female_users_count / total_users) * 100:.2f}%)")
14
15 # Age group distribution
16 age_groups = age.groupby('Age')['size'].sum().reset_index()
17 age_groups['percentage'] = (age_groups['size'] / total_users) * 100
18 print("\nAge Group Distribution:")
19 print(age_groups)
20
21 # Gender distribution within each age group
22 age_gender_distribution = age.pivot(index='Age', columns='Gender', values='size').fillna(0)
23 age_gender_distribution['Total'] = age_gender_distribution.sum(axis=1)
24 age_gender_distribution['Male_Percentage'] = (age_gender_distribution['M'] / age_gender_distribution['Total']) * 100
25 age_gender_distribution['Female_Percentage'] = (age_gender_distribution['F'] / age_gender_distribution['Total']) * 100
26 print("\nAge and Gender Distribution:")
27 print(age_gender_distribution)

```

➡ Total Users: 5891  
 Male Users: 4225 (71.72%)  
 Female Users: 1666 (28.28%)

Age Group Distribution:

	Age	size	percentage
0	0-17	218	3.700560
1	18-25	1069	18.146325
2	26-35	2053	34.849771
3	36-45	1167	19.809879
4	46-50	531	9.013750
5	51-55	481	8.164997
6	55+	372	6.314717

Age and Gender Distribution:

Gender	F	M	Total	Male_Percentage	Female_Percentage
Age					
0-17	78	140	218	64.220183	35.779817
18-25	287	782	1069	73.152479	26.847521
26-35	545	1508	2053	73.453483	26.546517
36-45	333	834	1167	71.465296	28.534704
46-50	182	349	531	65.725047	34.274953
51-55	142	339	481	70.478170	29.521830
55+	99	273	372	73.387097	26.612903

- The highest number of users are in 26-35 age category and comprise 34% of overall users
- 19% users are in 36-45 age group and 18% of users belong to 18-25 age group
- In the 0-17 age group, while only 3.7% of users lie, 35% of them are females, which is highest female ratio in any age group.

### ✓ 3.1.10

#### Bivariate Analysis of Occupation and Gender

```

1 occ = users.groupby(['Occupation', 'Gender'], as_index=False).size().sort_values(by="Occupation")
2 occ

```



	Occupation	Gender	size
0	0	F	226
1	0	M	462
2	1	F	203
3	1	M	314
4	2	F	88
5	2	M	168
6	3	F	98
7	3	M	72
8	4	F	228
9	4	M	512
11	5	M	80
10	5	F	31
12	6	F	99
13	6	M	129
14	7	F	137
15	7	M	532
16	8	F	3
17	8	M	14
18	9	F	85
19	9	M	3
21	10	M	126
20	10	F	66
22	11	F	22
23	11	M	106
24	12	F	46
25	12	M	330
26	13	F	33
27	13	M	107
28	14	F	78
29	14	M	216
31	15	M	112
30	15	F	28
32	16	F	49
33	16	M	186
34	17	F	50
35	17	M	441
36	18	F	4
37	18	M	63
38	19	F	15
39	19	M	56
40	20	F	77
41	20	M	196

### 3.1.11

#### Univariate Analysis of Occupation

```
1 # Insights from the occupation distribution plot
2 occ_groups = occ.groupby('Occupation')['size'].sum().reset_index()
3 occ_groups['percentage'] = (occ_groups['size'] / total_users) * 100
4 print("Occupation Distribution:")
5 print(occ_groups)
```

↗ Occupation Distribution:

	Occupation	size	percentage
0	0	688	11.678832
1	1	517	8.776099
2	2	256	4.345612
3	3	170	2.885758
4	4	740	12.561535
5	5	111	1.884230
6	6	228	3.870311
7	7	669	11.356306
8	8	17	0.288576
9	9	88	1.493804
10	10	192	3.259209
11	11	128	2.172806
12	12	376	6.382618
13	13	140	2.376507
14	14	294	4.990664
15	15	140	2.376507
16	16	235	3.989136
17	17	491	8.334748
18	18	67	1.137328
19	19	71	1.205228
20	20	273	4.634188

Occupations 0,4 and 7 have the maximum number of users, each have 11-12% users

```
1 # Gender distribution within each Occupation
2 occ_gender_distribution = occ.pivot(index='Occupation', columns='Gender', values='size').fillna(0)
3 occ_gender_distribution['Total'] = occ_gender_distribution.sum(axis=1)
4 occ_gender_distribution['Male_Percentage'] = (occ_gender_distribution['M'] / occ_gender_distribution['Total']) * 100
5 occ_gender_distribution['Female_Percentage'] = (occ_gender_distribution['F'] / occ_gender_distribution['Total']) * 100
6 print("\nOccupation and Gender Distribution:")
7 print(occ_gender_distribution)
```

↗ \nOccupation and Gender Distribution:

Gender	F	M	Total	Male_Percentage	Female_Percentage
Occupation					
0	226	462	688	67.151163	32.848837
1	203	314	517	60.735010	39.264990
2	88	168	256	65.625000	34.375000
3	98	72	170	42.352941	57.647059
4	228	512	740	69.189189	30.810811
5	31	80	111	72.072072	27.927928
6	99	129	228	56.578947	43.421053
7	137	532	669	79.521674	20.478326
8	3	14	17	82.352941	17.647059
9	85	3	88	3.409091	96.590909
10	66	126	192	65.625000	34.375000
11	22	106	128	82.812500	17.187500
12	46	330	376	87.765957	12.234043
13	33	107	140	76.428571	23.571429
14	78	216	294	73.469388	26.530612
15	28	112	140	80.000000	20.000000
16	49	186	235	79.148936	20.851064
17	50	441	491	89.816701	10.183299
18	4	63	67	94.029851	5.970149
19	15	56	71	78.873239	21.126761
20	77	196	273	71.794872	28.205128

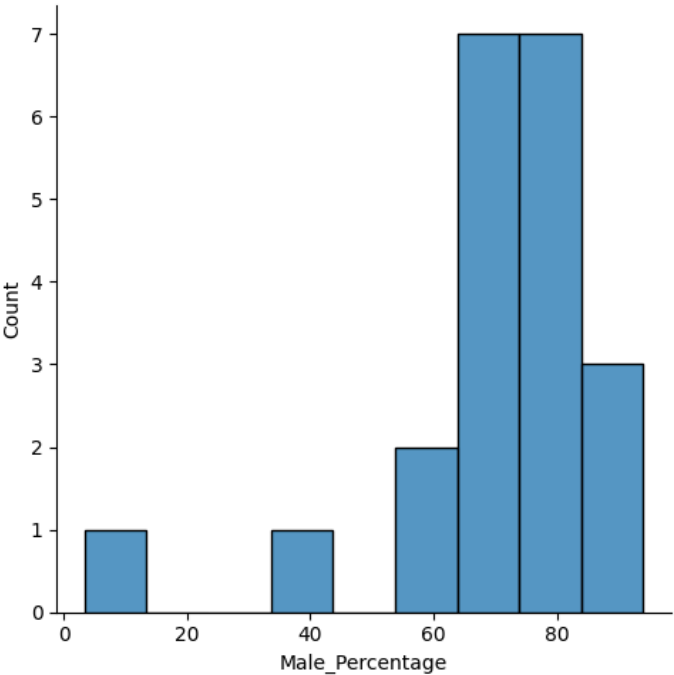
```
<>:6: SyntaxWarning: invalid escape sequence '\0'
```

```
<>:6: SyntaxWarning: invalid escape sequence '\0'
```

```
C:\Users\Admin\AppData\Local\Temp\ipykernel_37668\2747592413.py:6: SyntaxWarning: invalid escape sequence '\0'
print("\nOccupation and Gender Distribution:")
```

```
1 sns.displot(data=occ_gender_distribution['Male_Percentage'])
```

```
<seaborn.axisgrid.FacetGrid at 0x1a455b648c0>
```



The gender wise breakup for occupation is similar to overall gender breakup. One Way ANOVA on Occupation and Male Percentage columnn can be used to confirm this further Occupation 9 is a occupation where 96.5% of users are males Occupation 18 is a occupation where 94% of users are females

3.1.12

Bivariate Analysis of Occupation and Age Group

```
1 occAge = pd.crosstab(users.Age, users.Occupation)
2 occAge
```

Occupation	0	1	2	3	4	5	6	7	8	9	...	11	12	13	14	15	16	17	18	19	20
Age																					
0-17	26	4	3	0	3	0	0	2	1	0	...	1	1	1	2	0	0	2	0	11	0
18-25	101	44	42	18	518	18	9	17	1	6	...	7	58	0	37	12	17	50	10	30	46
26-35	292	152	109	67	197	52	75	251	6	28	...	50	177	0	129	71	56	197	22	15	106
36-45	131	115	46	41	14	24	57	213	3	34	...	31	83	5	67	25	56	129	16	9	67
46-50	62	80	29	16	3	10	29	71	1	9	...	17	32	7	25	17	37	56	8	2	20
51-55	43	68	17	16	5	6	38	76	2	7	...	16	19	27	19	8	45	35	9	2	23
55+	33	54	10	12	0	1	20	39	3	4	...	6	6	100	15	7	24	22	2	2	11

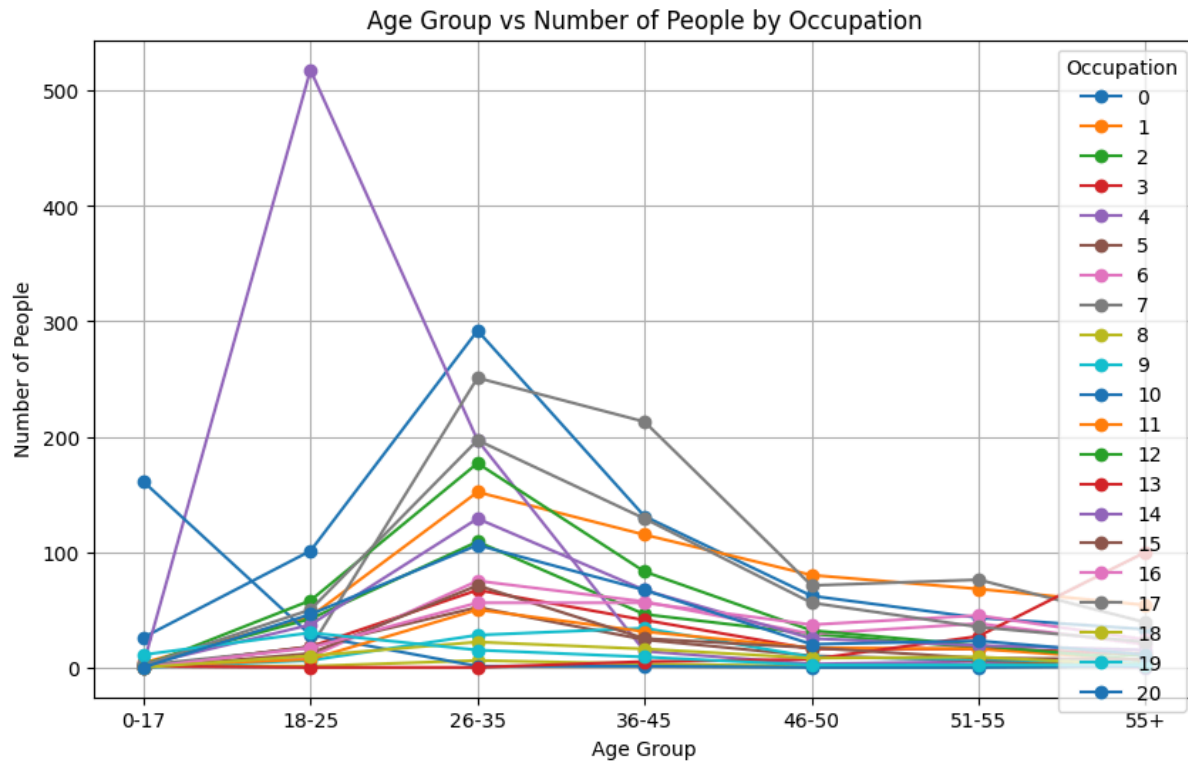
7 rows × 21 columns

```
1 # Plot the line chart
2 plt.figure(figsize=(10, 6))
3
4 for column in occAge.columns:
5     plt.plot(occAge.index, occAge[column], marker='o', label=column)
6
7 # Customize the plot
8 plt.title('Age Group vs Number of People by Occupation')
9 plt.xlabel('Age Group')
10 plt.ylabel('Number of People')
```

```

11 plt.legend(title='Occupation')
12 plt.grid(True)
13 plt.show()

```



```

1 #find which occupation has maximum value for 26-35 age group
2 # Total number of users in the 26-35 age group
3
4 def age_group_top3_occupations(age_group):
5     total_users_in_age_group = cTrans[age_group].sum()
6
7     top_3 = cTrans.nlargest(3, age_group)[age_group]
8     print(f'{age_group} occupations {top_3.index.values}')
9     print(f'Total users in {age_group} = {(top_3/total_users_26_35).sum()*100}')
10
11
12 age_group_top3_occupations('18-25')
13 age_group_top3_occupations('26-35')
14 age_group_top3_occupations('36-45')
15

```



```
18-25 occupations [ 4  0 12]
```

```

-----
NameError                                Traceback (most recent call last)
Cell In[181], line 12
      8 print(f'{age_group} occupations {top_3.index.values}')
      9 print(f'Total users in {age_group} = {(top_3/total_users_26_35).sum()*100}')
--> 12 age_group_top3_occupations('18-25')
      13 age_group_top3_occupations('26-35')
      14 age_group_top3_occupations('36-45')

Cell In[181], line 9, in age_group_top3_occupations(age_group)
      7 top_3 = cTrans.nlargest(3, age_group)[age_group]
      8 print(f'{age_group} occupations {top_3.index.values}')
--> 9 print(f'Total users in {age_group} = {(top_3/total_users_26_35).sum()*100}')

NameError: name 'total_users_26_35' is not defined

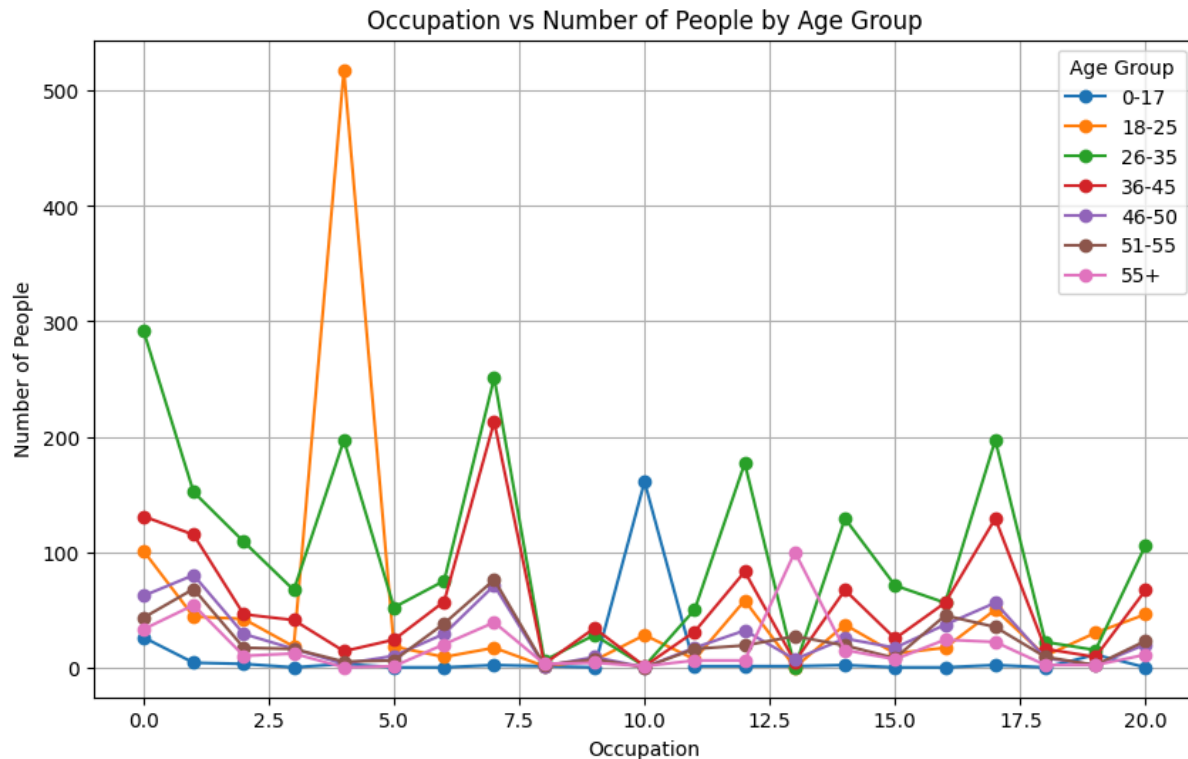
```

- For 26-35 age group, 36% of users are from occupations 0,4,7
- For 18-25 age group, 32% of users are from occupations 0,4,12
- For 36-45 age group, 23% of users are from occupations 0,7,17

```

1 # Plot the line chart
2 plt.figure(figsize=(10, 6))
3 cTrans = occAge.T
4 for column in cTrans.columns:
5     plt.plot(cTrans.index, cTrans[column], marker='o', label=column)
6
7 # Customize the plot
8 plt.title('Occupation vs Number of People by Age Group')
9 plt.xlabel('Occupation')
10 plt.ylabel('Number of People')
11 plt.legend(title='Age Group')
12 plt.grid(True)
13 plt.show()

```



### 3.1.13

Bivariate Analysis of City\_Category and Purchase

```

1 city = df.groupby('City_Category',as_index=False).size().sort_values(by="City_Category")
2 city['percent'] = round(city['size']*100/len(df),0)
3 city

```



	City_Category	size	percent
0	A	147720	27.0
1	B	231173	42.0
2	C	171175	31.0

42% of purchases are from City B, 31% from city C

### 3.1.14

Bivariate Analysis of City\_Category and User\_ID


```

1 cityUser = users.groupby('City_Category',as_index=False).size().sort_values(by="City_Category")
2 cityUser['percent'] = round(cityUser['size']*100/len(users),0)

```



3 cityUser




	City_Category	size	percent
0	A	1045	18.0
1	B	1707	29.0
2	C	3139	53.0

53% of users are from City C, 29% from City B

3.1.15

Bivariate Analysis of City\_Category and Gender

```
1 cityGender = pd.crosstab(users.City_Category, users.Gender, margins=True)
2 cityGender['percent'] = cityGender.M*100/cityGender.All
3 cityGender
```




	Gender	F	M	All	percent
City_Category					
A		295	750	1045	71.770335
B		503	1204	1707	70.533099
C		868	2271	3139	72.347881
All		1666	4225	5891	71.719572

The gender distribution across the three cities is same.

3.1.16


Bivariate Analysis of City\_Category and Age Group

```
1 cityAge= pd.crosstab(users.City_Category, users.Age,margins=True)
2
3 #cityAge['percent'] = cityAge.M*100/cityGender.All
4 cityAge
```



	Age	0-17	18-25	26-35	36-45	46-50	51-55	55+	All
City_Category									
A		25	214	461	176	53	67	49	1045
B		50	331	652	335	146	135	58	1707
C		143	524	940	656	332	279	265	3139
All		218	1069	2053	1167	531	481	372	5891

```
1 cityAge= pd.crosstab(users.City_Category, users.Age,normalize = 'index')
2 cityAge*100
```



	Age	0-17	18-25	26-35	36-45	46-50	51-55	55+
City_Category								
A		2.392344	20.478469	44.114833	16.842105	5.071770	6.411483	4.688995
B		2.929115	19.390744	38.195665	19.625073	8.553017	7.908612	3.397774
C		4.555591	16.693214	29.945843	20.898375	10.576617	8.888181	8.442179

```
1 cityAge= pd.crosstab(users.City_Category, users.Age,normalize = 'columns',margins= True)
2 round(cityAge*100,0)
```

↗

	Age	0-17	18-25	26-35	36-45	46-50	51-55	55+	All
City_Category									
	A	11.0	20.0	22.0	15.0	10.0	14.0	13.0	18.0
	B	23.0	31.0	32.0	29.0	27.0	28.0	16.0	29.0
	C	66.0	49.0	46.0	56.0	63.0	58.0	71.0	53.0

- In City A, 44% of users are from 26-35 age group
- While 53% of users are from City C, it has 71% of all 55+ aged users

3.1.17

Bivariate Analysis of City\_Category and Occupation

```
1 cityOcc= pd.crosstab(users.City_Category, users.Occupation,margins=True)
2
3 #cityAge['percent'] = cityAge.M*100/cityGender.All
4 cityOcc
```

↗

	Occupation	0	1	2	3	4	5	6	7	8	9	...	12	13	14	15	16	17	18	19	20	All
City_Category																						
	A	129	94	68	34	172	14	24	113	4	7	...	77	12	53	25	35	62	7	14	59	1045
	B	203	140	80	47	236	40	74	170	2	25	...	113	30	83	43	61	136	12	18	106	1707
	C	356	283	108	89	332	57	130	386	11	56	...	186	98	158	72	139	293	48	39	108	3139
	All	688	517	256	170	740	111	228	669	17	88	...	376	140	294	140	235	491	67	71	273	5891

4 rows × 22 columns

```
1 cityOcc= pd.crosstab(users.City_Category, users.Occupation,normalize = 'index')
2
3 round(cityOcc*100,0)
4
```

↗

	Occupation	0	1	2	3	4	5	6	7	8	9	...	11	12	13	14	15	16	17	18	19	20
City_Category																						
	A	12.0	9.0	7.0	3.0	16.0	1.0	2.0	11.0	0.0	1.0	...	2.0	7.0	1.0	5.0	2.0	3.0	6.0	1.0	1.0	6.0
	B	12.0	8.0	5.0	3.0	14.0	2.0	4.0	10.0	0.0	1.0	...	3.0	7.0	2.0	5.0	3.0	4.0	8.0	1.0	1.0	6.0
	C	11.0	9.0	3.0	3.0	11.0	2.0	4.0	12.0	0.0	2.0	...	2.0	6.0	3.0	5.0	2.0	4.0	9.0	2.0	1.0	3.0

3 rows × 21 columns

```
1 cityOcc.loc['A'].sum()
```

↗

0.9999999999999999

```
1 cityOccTrans= cityOcc.T
2 def city_group_top3_occupations(city_group):
3     top_3 = cityOccTrans.nlargest(3, city_group)[city_group]
4     print(f'{city_group} occupations {top_3.index.values}')
5     print(f'% users in top 3 occupations for {city_group} = {(top_3.sum()*100}')
```

```
A occupations [4 0 7]
% users in top 3 occupations for A = 39.61722488038278
B occupations [4 0 7]
% users in top 3 occupations for B = 35.67662565905097
C occupations [7 0 4]
% users in top 3 occupations for C = 34.21471806307741
```

3.1.18

Bivariate Analysis of City\_Category and Product Category

```
1 cityProdCat= pd.crosstab(df.City_Category, df.Product_Category,margins=True)
2 cityProdCat
```

Product_Category	1	2	3	4	5	6	7	8	9	10	...	12	13	14	15	16	17
City_Category																	
A	35081	6141	4943	3050	42211	5507	1226	32179	110	1333	...	1063	1614	481	1717	2848	121
B	58253	10444	8587	5226	64138	8526	1599	47553	174	2063	...	1675	2271	632	2638	4038	267
C	47044	7279	6683	3477	44584	6433	896	34193	126	1729	...	1209	1664	410	1935	2942	190
All	140378	23864	20213	11753	150933	20466	3721	113925	410	5125	...	3947	5549	1523	6290	9828	578

4 rows × 21 columns

3.1.19

Bivariate Analysis of Stay\_In\_Current\_City\_Years and UserID

```
1 stayDuration= users.groupby('Stay_In_Current_City_Years',as_index=False).size().sort_values('Stay_In_Current_City_Years')
2 stayDuration['percent'] = round(stayDuration['size']*100/stayDuration['size'].sum(),0)
3 stayDuration
4
```

Stay_In_Current_City_Years	size	percent
0	772	13.0
1	2086	35.0
2	1145	19.0
3	979	17.0
4	909	15.0

35% of users have stayed in same city for 1-2 yrs

3.1.20

Bivariate Analysis of Stay\_In\_Current\_City\_Years and Stay Duration

```
1 cityStayDuration = pd.crosstab(users.City_Category,users.Stay_In_Current_City_Years, normalize='index')
2 cityStayDuration
```

Stay_In_Current_City_Years	0	1	2	3	4+
City_Category					
A	0.140670	0.354067	0.175120	0.172249	0.157895
B	0.123609	0.356180	0.200351	0.172818	0.147042
C	0.131889	0.352979	0.197515	0.160561	0.157056

The duration of stay in the city is similar in all three cities across various duration categories.

```
1 genderStayDuration = pd.crosstab(users.Gender,users.Stay_In_Current_City_Years, normalize='index')
2 genderStayDuration
```

Stay_In_Current_City_Years	0	1	2	3	4+
Gender					
F	0.128451	0.362545	0.196879	0.171669	0.140456
M	0.132071	0.350769	0.193373	0.164024	0.159763

### 3.1.21

Bivariate Analysis of Marital\_Status and Gender

```
1 maritalStatusGender= pd.crosstab(users.Marital_Status, users.Gender,margins=True)
2 maritalStatusGender
```

Gender	F	M	All
Marital_Status			
0	947	2470	3417
1	719	1755	2474
All	1666	4225	5891

```
1 print('Unmarried % = ',round(100*maritalStatusGender.at[0,'All'] /maritalStatusGender.at['All','All'],1))
2 print('Unmarried male% = ',round(100*maritalStatusGender.at[0,'M'] /maritalStatusGender.at['All','M'],1))
3 print('Unmarried female% = ',round(100*maritalStatusGender.at[0,'F'] /maritalStatusGender.at['All','F'],1))
```

```
Unmarried % = 58.0
Unmarried male% = 58.5
Unmarried female% = 56.8
```

58% of users are unmarried Ratio of unmarried users is same across genders

### 3.1.22

Bivariate Analysis of Marital\_Status and Age Category

```
1 maritalStatusAge= pd.crosstab(users.Marital_Status, users.Age,margins=True)
2 maritalStatusAge
```

Age	0-17	18-25	26-35	36-45	46-50	51-55	55+	All
Marital_Status								
0	218	825	1244	705	156	136	133	3417
1	0	244	809	462	375	345	239	2474
All	218	1069	2053	1167	531	481	372	5891

```
1 # Calculate percentages
2 maritalStatusAge_percent = maritalStatusAge.div(maritalStatusAge['All'], axis=0) * 100
3 maritalStatusAge_percent
4
5
```



	Age	0-17	18-25	26-35	36-45	46-50	51-55	55+	All
Marital_Status									
0		6.379865	24.143986	36.406204	20.632133	4.565408	3.980100	3.892303	100.0
1		0.000000	9.862571	32.700081	18.674212	15.157639	13.945028	9.660469	100.0
All		3.700560	18.146325	34.849771	19.809879	9.013750	8.164997	6.314717	100.0

```
1 print(24/(9+24))
2 print(36.4/(36.4+32))
3 print(20.6/(20.6+18.6))
```



0.7272727272727273  
0.5321637426900584  
0.5255102040816326

- 72% of users in 18-25 age group are unmarried
- 53% of users in 26-35 age group are unmarried
- 52% of users in 36-45 age group are unmarried

3.1.23

Bivariate Analysis of Marital\_Status and Occupation

```
1 maritalStatusOcc= pd.crosstab(users.Marital_Status, users.Occupation,margins=True)
2 maritalStatusOcc
```



Occupation	0	1	2	3	4	5	6	7	8	9	...	12	13	14	15	16	17	18	19	20	All
Marital_Status																					
0	402	279	144	90	544	65	102	355	11	41	...	219	50	151	73	110	288	35	57	144	3417
1	286	238	112	80	196	46	126	314	6	47	...	157	90	143	67	125	203	32	14	129	2474
All	688	517	256	170	740	111	228	669	17	88	...	376	140	294	140	235	491	67	71	273	5891

3 rows × 22 columns


```
1 #Calculate the breakup of the married unmarried group acorss occupations
2 maritalStatusOcc_percent = maritalStatusOcc.div(maritalStatusOcc['All'], axis=0) * 100
3 maritalStatusOcc_percent
```




Occupation	0	1	2	3	4	5	6	7	8	9	...	12
Marital_Status												
0	11.764706	8.165057	4.214223	2.633889	15.920398	1.902253	2.985075	10.389230	0.321920	1.199883	...	6.409131
1	11.560226	9.620049	4.527082	3.233630	7.922393	1.859337	5.092967	12.691997	0.242522	1.899757	...	6.345998
All	11.678832	8.776099	4.345612	2.885758	12.561535	1.884230	3.870311	11.356306	0.288576	1.493804	...	6.382618

3 rows × 22 columns

```
1 #Calculate the breakup of the married unmarried group within each occupation
2 sum_first_two = maritalStatusOcc_percent.iloc[0:2].sum()
3 # Calculate the percentages for each column
4 percentage_row = round((maritalStatusOcc_percent.iloc[0:1] / sum_first_two) * 100,1)
5
6 # Add the percentage row to the DataFrame
7 maritalStatusOcc_percent = pd.concat([maritalStatusOcc_percent,percentage_row])
8 maritalStatusOcc_percent
```

 Index([0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 'All'], dtype='object', name='Occupation')

 Index([0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 'All'], dtype='object', name='Occupation')


Occupation	0	1	2	3	4	5	6	7	8	9	...
Marital_Status											
0	11.764706	8.165057	4.214223	2.633889	15.920398	1.902253	2.985075	10.389230	0.321920	1.199883	...
1	11.560226	9.620049	4.527082	3.233630	7.922393	1.859337	5.092967	12.691997	0.242522	1.899757	...
All	11.678832	8.776099	4.345612	2.885758	12.561535	1.884230	3.870311	11.356306	0.288576	1.493804	...
0	50.400000	45.900000	48.200000	44.900000	66.800000	50.600000	37.000000	45.000000	57.000000	38.700000	... 5


4 rows × 22 columns


- 66% of users in Occupation 4 are unmarried
- 50% of users in Occupation 0 are unmarried
- 45% of users in Occupation 7 are unmarried
- 75% of users in Occupation 19 are unmarried

3.1.24


Bivariate Analysis of Marital\_Status and City


 1 maritalStatusCity= pd.crosstab(users.Marital\_Status, users.City\_Category,margins=True)


 2 maritalStatusCity




City_Category	A	B	C	All
Marital_Status				
0	652	1004	1761	3417
1	393	703	1378	2474
All	1045	1707	3139	5891

 1 #Calculate the breakup of the married unmarried group across each city


 2 maritalStatusCity\_percent = maritalStatusCity.div(maritalStatusCity['All'], axis=0) \* 100


 3 maritalStatusCity\_percent





City_Category	A	B	C	All
Marital_Status				
0	19.081065	29.382499	51.536435	100.0
1	15.885206	28.415521	55.699272	100.0
All	17.738924	28.976405	53.284672	100.0


- The breakup of married-unmarried users is same as the general breakup of users across the cities


 1 #Calculate the breakup of the married unmarried group within each city


 2 sum\_first\_two = maritalStatusCity\_percent.iloc[0:2].sum()


 3 # Calculate the percentages for each column


 4 percentage\_row = round((maritalStatusCity\_percent.iloc[0:1] / sum\_first\_two) \* 100,1)

 5

 6 # Add the percentage row to the DataFrame

 7 maritalStatusCity\_percent = pd.concat([maritalStatusCity\_percent,percentage\_row])

 8 maritalStatusCity\_percent



City_Category	A	B	C	All
Marital_Status				
0	19.081065	29.382499	51.536435	100.0
1	15.885206	28.415521	55.699272	100.0
All	17.738924	28.976405	53.284672	100.0
0	54.600000	50.800000	48.100000	50.0

- While overall we have 58% unmarried users, in City C we have only 48% of unmarried users


### 3.1.25

Bivariate Analysis of Marital\_Status and Stay Duration

```

1 maritalStatusStayDuration= pd.crosstab(users.Marital_Status, users.Stay_In_Current_City_Years, margins=True)
2 maritalStatusStayDuration
3
4 #Calculate the breakup of the married-unmarried group across each city
5 maritalStatusStayDuration_percent = maritalStatusStayDuration.div(maritalStatusStayDuration['All'], axis=0) * 100
6 maritalStatusStayDuration_percent
7
8
9

```



Stay_In_Current_City_Years	0	1	2	3	4+	All
Marital_Status						
0	13.666959	34.182031	19.402985	16.944688	15.803336	100.0
1	12.328213	37.105901	19.482619	16.168149	14.915117	100.0
All	13.104736	35.409947	19.436428	16.618571	15.430317	100.0

- The breakup of married-unmarried users is same as the general breakup of duration of stay of users irrespective of marital status. So marital status does not have any impact on duration of stay in a city for a user

### 3.1.26

Univariate Analysis of Purchase column

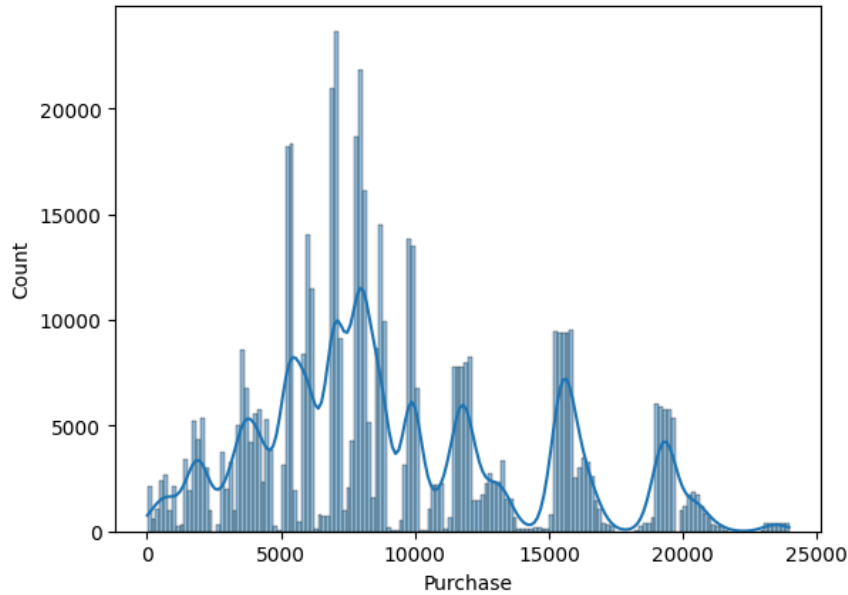
```
1 df.Purchase.mean()
```



```
9263.968712959126
```

```
1 sns.histplot(data=df["Purchase"], kde=True)
```

 <Axes: xlabel='Purchase', ylabel='Count'>

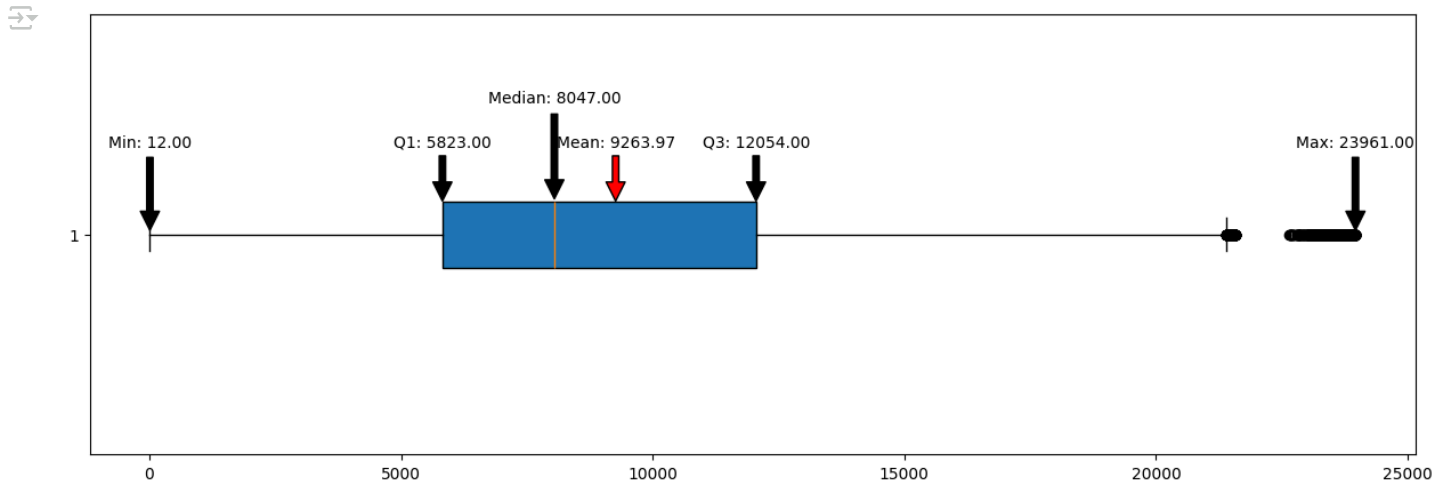


```

1 data=df.Purchase #Only this one line is to be updated
2
3 fig, ax = plt.subplots(figsize=(15, 5))
4 # Create boxplot
5 boxplot = ax.boxplot(data, vert=False, patch_artist=True)
6
7 # Calculate the five-number summary + mean
8 mean = np.mean(data)
9 min_val = np.min(data)
10 q1 = np.percentile(data, 25)
11 median = np.median(data)
12 q3 = np.percentile(data, 75)
13 max_val = np.max(data)
14
15 # Annotate the five-number summary on the boxplot
16 ax.annotate(f'Min: {min_val:.2f}', xy=(min_val, 1), xytext=(min_val, 1.2),
17             arrowprops=dict(facecolor='black', shrink=0.05), ha='center')
18
19 ax.annotate(f'Q1: {q1:.2f}', xy=(q1, 1.07), xytext=(q1, 1.2),
20             arrowprops=dict(facecolor='black', shrink=0.05), ha='center')
21
22 ax.annotate(f'Median: {median:.2f}', xy=(median, 1.07), xytext=(median, 1.3),
23             arrowprops=dict(facecolor='black', shrink=0.05), ha='center')
24
25 ax.annotate(f'Q3: {q3:.2f}', xy=(q3, 1.07), xytext=(q3, 1.2),
26             arrowprops=dict(facecolor='black', shrink=0.05), ha='center')
27
28 ax.annotate(f'Max: {max_val:.2f}', xy=(max_val, 1), xytext=(max_val, 1.2),
29             arrowprops=dict(facecolor='black', shrink=0.05), ha='center')
30
31 ax.annotate(f'Mean: {mean:.2f}', xy=(mean, 1.07), xytext=(mean, 1.2),
32             arrowprops=dict(facecolor='red', shrink=0.05), ha='center')
33
34
35 # Display plot
36 plt.show()
37
38

```





- 50% of customers shop between 5823 and 12054
- Overall range of purchase amount is from 12 to 23961

## ✓ 4.0

Answering Questions asked in Case Study

### 4.1

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

We have already established that Mean purchase amount for males is 9437.53 *Mean purchase amount for females is 8734.57*

So women are not spending more money per transaction than men as per this data. Later we have statistically established that the spending habits of the two groups is different

## ✓ Why are women spending less?

We observe that only 25% of customers are women as per given data.

### Possible reasons

- Walmart's marketing strategies might be more appealing to men, leading to higher male foot traffic and spending
- Cultural norms and societal expectations might influence shopping behaviors. For example, women might be more likely to shop for groceries and household items, while men might shop for higher-ticket items
- Gender Roles: Traditional gender roles might play a part in determining who does the shopping and what they buy.

### It could indicate that shopping experience of women needs improvement

This can be done by

- keeping more items relevant to women. The product mix at Walmart might cater more to male preferences, such as electronics and automotive products, which could result in higher spending by men
- making it easier to reach
- helpful and women friendly staff
- ensuring women safety
- Women are more likely to use coupons and look for discounts, which might lead them to shop at stores that offer better deals

## ✓ 4.2

## Confidence intervals and distribution of the mean of the expenses by female and male customers

```

1 # select a random sample of 1000 rows for each gender
2 sample_males = males.sample(n=1000)
3 sample_females = females.sample(n=1000)
4
5 # identify confidence interval for the mean of the purchase
6 from scipy import stats
7
8 def my_decorator(func):
9     def wrapper(data, confidence=0.95, text="Mean"):
10         mean, lower, upper = func(data, confidence)
11         return f'Using CLT and single sample of {len(data)} rows, {text} = {mean:.2f}, CI = [{lower:.2f}, {upper:.2f}] with {c
12     return wrapper
13
14 #We want to decorate the function mean_confidence_interval to add some text to the output
15 @my_decorator
16 def mean_confidence_interval_decorated(data, confidence=0.95):
17     return mean_confidence_interval(data, confidence)
18
19 #Function to calculate the confidence interval for the mean of the purchase amount
20 def mean_confidence_interval(data, confidence=0.95):
21     a = 1.0 * np.array(data)
22     n = len(a)
23     m, se = np.mean(a), stats.sem(a)
24     h = se * stats.t.ppf((1 + confidence) / 2., n-1) #dof is n-1
25     return m, m-h, m+h
26
27
28
29 print(f'Using actual population of {len(males)} rows, Population Mean purchase amount for males is ${PopMeanMales:.2f}')
30 print(f'{mean_confidence_interval_decorated(sample_males.Purchase, confidence=0.95, text="Mean Purchase amount for Males")}')
31 print()
32 print(f'Using actual population of {len(females)} rows, Population Mean purchase amount for females is ${PopMeanFemales:.2f}')
33 print(f'{mean_confidence_interval_decorated(sample_females.Purchase, confidence=0.95, text="Mean Purchase amount for Females")}')
34
35
36
37

```

Using actual population of 414259 rows, Population Mean purchase amount for males is \$9437.53  
 Using CLT and single sample of 1000 rows, Mean Purchase amount for Males = 9191.12, CI = [8889.37, 9492.87] with 95.0% confiden

Using actual population of 135809 rows, Population Mean purchase amount for females is \$8734.57  
 Using CLT and single sample of 1000 rows, Mean Purchase amount for Females = 8811.79, CI = [8514.98, 9108.60] with 95.0% confid

```

1 #Utility function to plot sample means distribution and confidence interval
2
3 import matplotlib.patches as mpatches
4
5 def plotConfidenceIntervals(sample_means, pop_mean, text, xlabel):
6     #Plot the second and third columns of sample_means which are the lower and upper bounds of the confidence interval
7
8     totalIntervalCount = len(sample_means)
9     plt.figure(figsize=(10, 5+round(totalIntervalCount/10,0)))
10    missed=0
11
12    #Plot 1000 confidence intervals to see how many contain the population mean
13    for i in range(totalIntervalCount):
14        # Plot the horizontal line for the confidence interval
15        mean, lower, upper = sample_means[i]
16        #Plot the confidence interval that miss the population mean in red
17        if(lower<=pop_mean<=upper):
18            plt.hlines(i*4, lower, upper, color='g', linewidth=2, label=f'95% CI for {mean:.2f}')
19        else:
20            missed=missed+1
21            plt.hlines(i*4, lower, upper, color='r', linewidth=2, label=f'95% CI for {mean:.2f}')
22        # Plot the sample mean as a dot
23        plt.plot(mean, i*4, 'ro', label='Mean')
24

```

```

25 # Plot the population mean
26 plt.axvline(pop_mean, color='b', linestyle='dashed', linewidth=2, label='Mean')
27
28 # Add custom legend
29 red_line = mpatches.Patch(color='red', label=f'{missed}({missed*100/totalIntervalCount}) Intervals that missed the populati
30 green_line = mpatches.Patch(color='green', label=f'{totalIntervalCount-missed}({(totalIntervalCount-missed)*100/totalInterv
31 blue_line = mpatches.Patch(color='blue', label=f'Population mean = {pop_mean:.2f}')
32 red_dot = plt.Line2D([0], [0], marker='o', color='w', label='Mean of the sample', markerfacecolor='red', markersize=10)
33
34 plt.legend(handles=[red_line, green_line, red_dot, blue_line], loc='upper left')
35
36
37 plt.title(f'Population Mean {pop_mean:.2f} and {totalIntervalCount} 95% Confidence Interval {text}')
38 plt.xlabel(xlabel)
39 plt.yticks([]) # Remove y-axis ticks
40
41
42 def plotSampleMeanDistributionAndConfidenceIntervals(data,
43             sample_size,
44             pop_mean,
45             confidence=0.95,
46             text="for"):
47
48     sample_means=[]
49     for _ in range(1000):
50         sample = data.sample(n=sample_size)
51         m,l,u = mean_confidence_interval(sample, confidence)
52         sample_means.append([m,l,u])
53     #select first column of the sample_means_male, whcih is the mean
54     sample_means = np.array(sample_means)
55     # plot the sample_means_male to see the distribution
56     sns.histplot(sample_means[:,0], kde=True)
57     plt.xlabel(data.name)
58     plt.ylabel('Frequency')
59     plt.title(f'Distribution of Sample means of {data.name} amounts {text}')
60     plt.show()
61
62     plotConfidenceIntervals(sample_means, pop_mean, text, data.name)
63     return sample_means
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

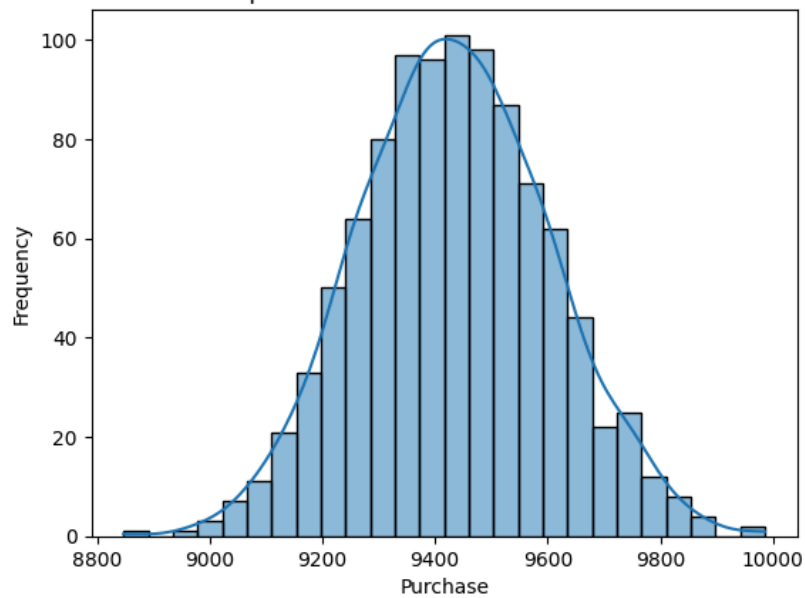
```

1 male_sample_means = plotSampleMeanDistributionAndConfidenceIntervals(males.Purchase, 1000, PopMeanMales, confidence=0.95, text=

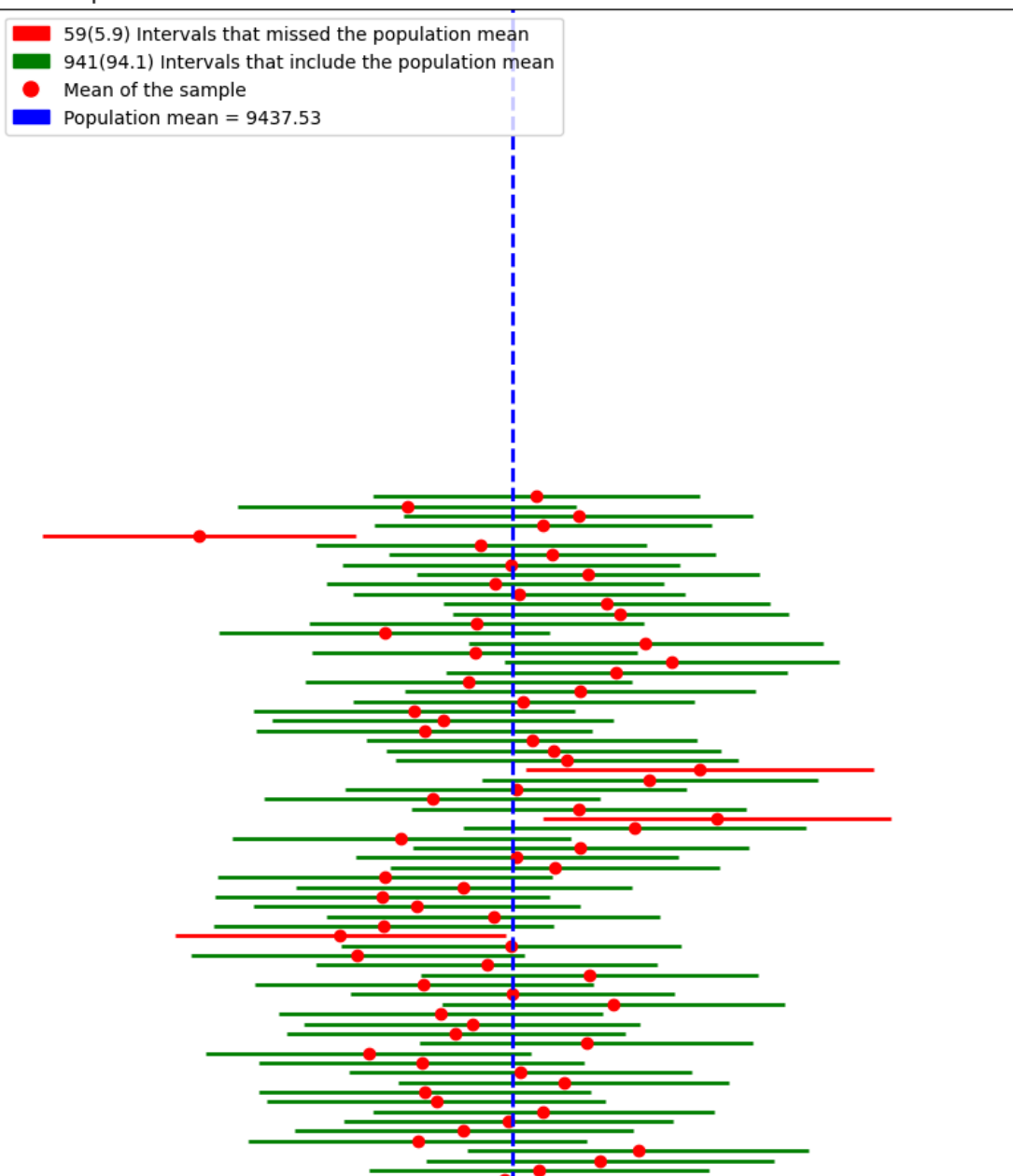
```

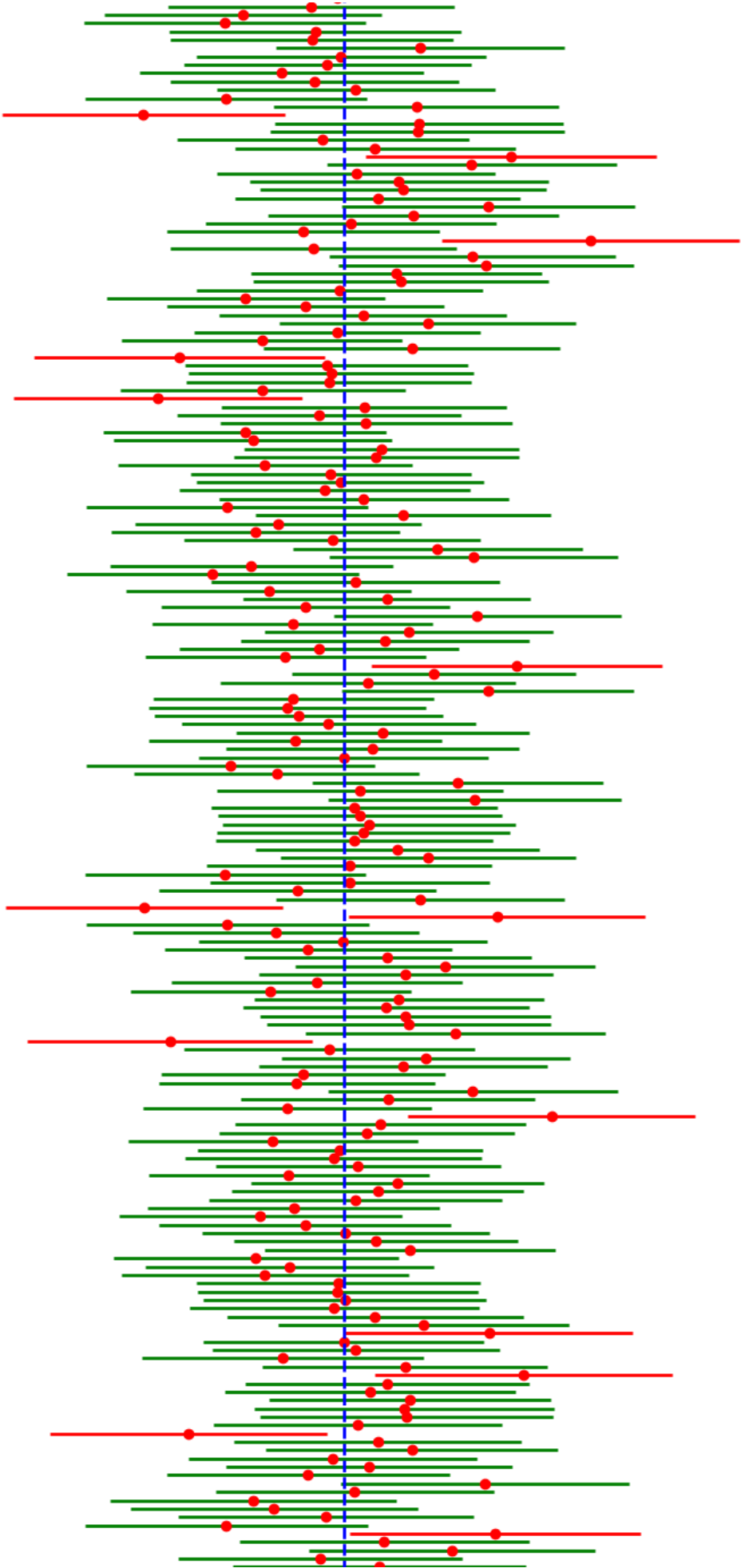


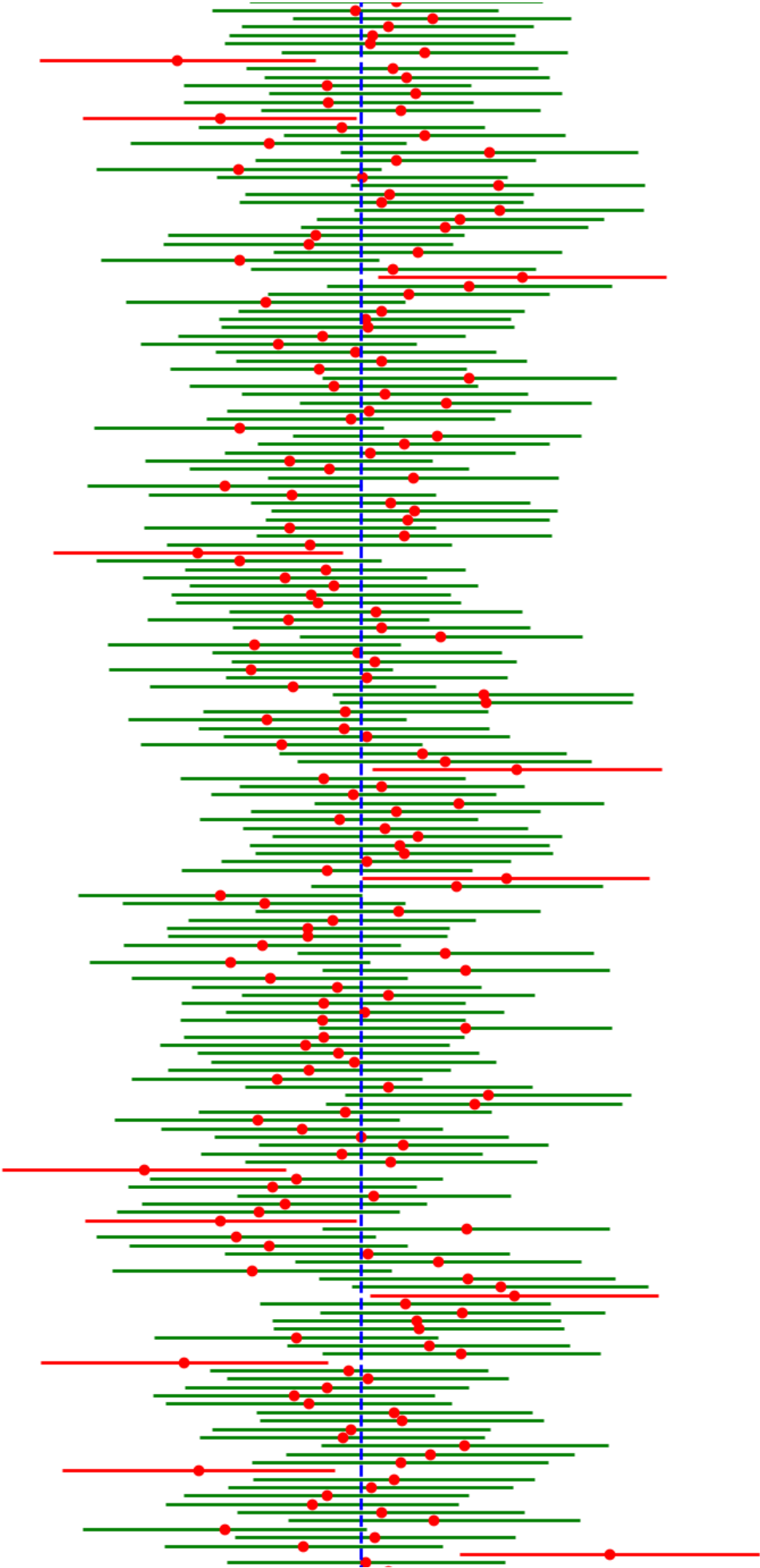
Distribution of Sample means of Purchase amounts for Male Customers

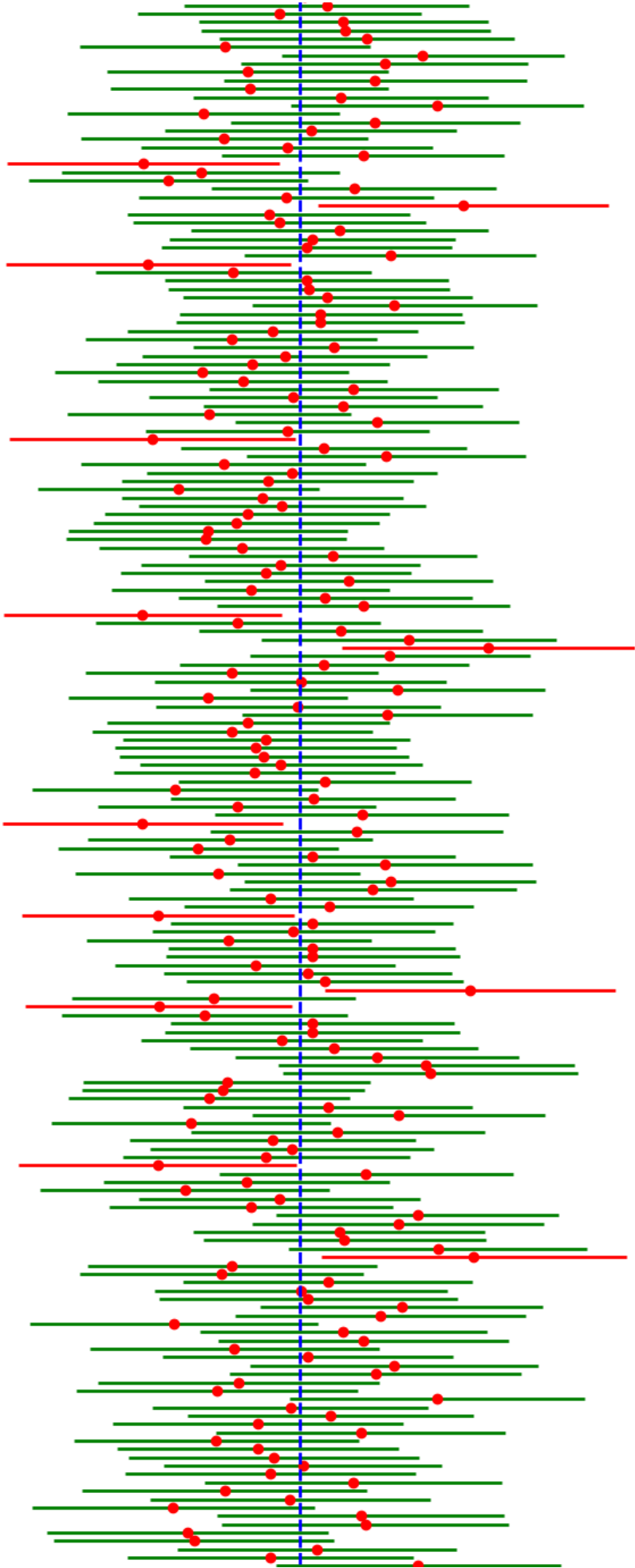


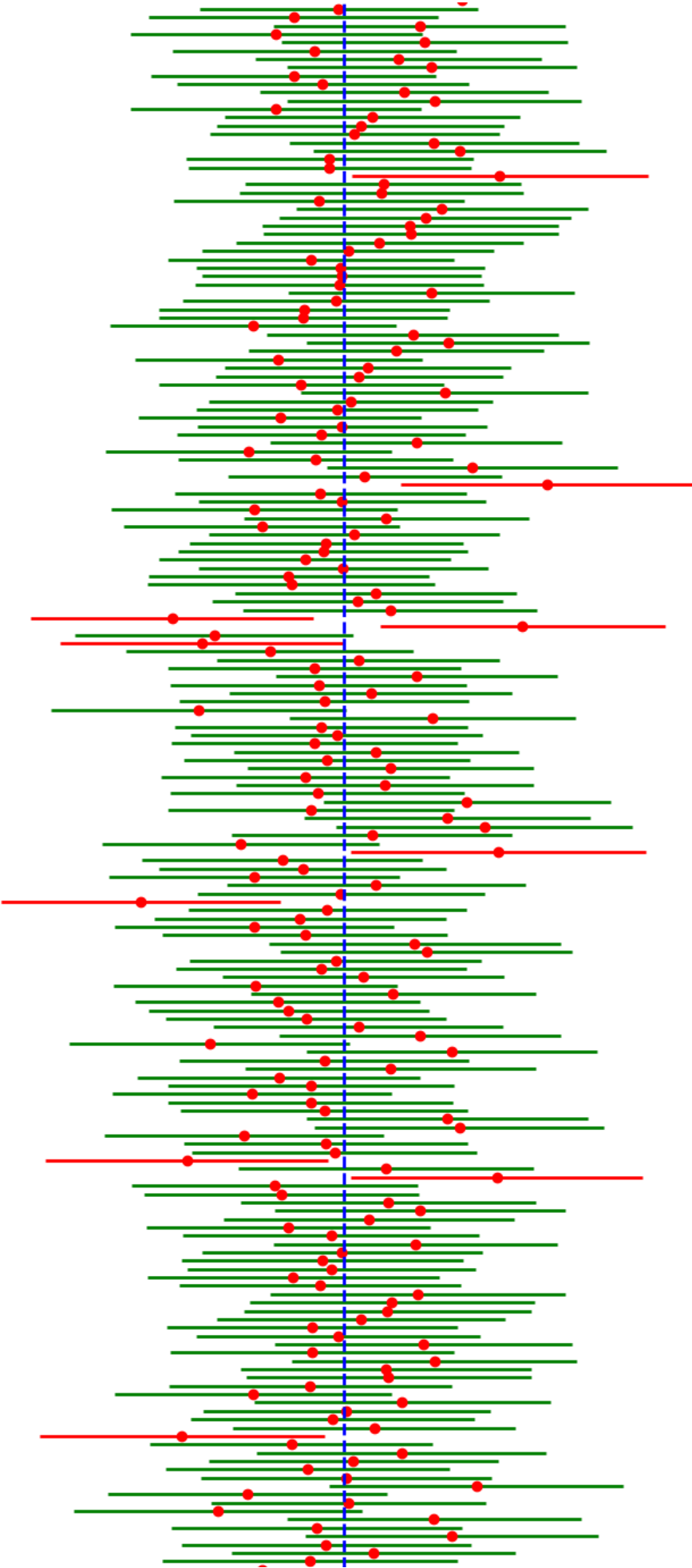
Population Mean 9437.53 and 1000 95% Confidence Interval for Male Customers



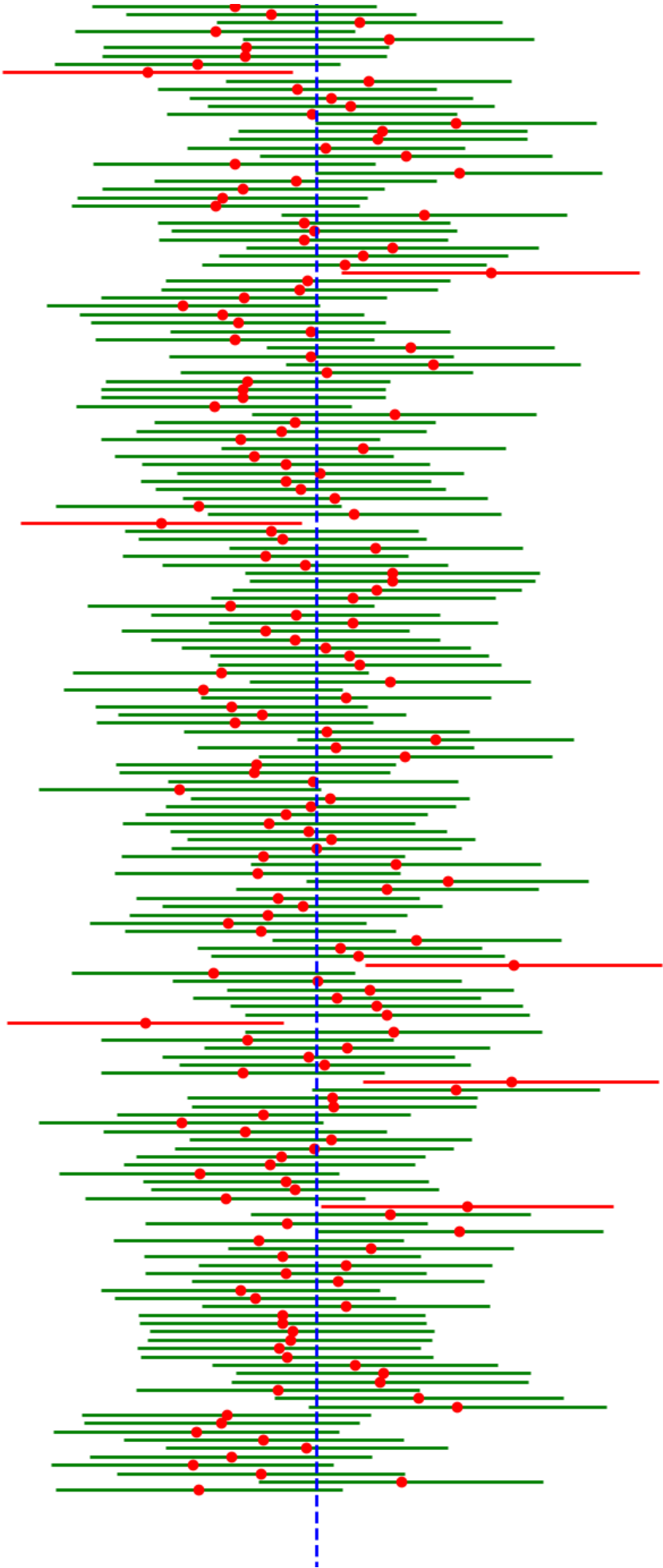


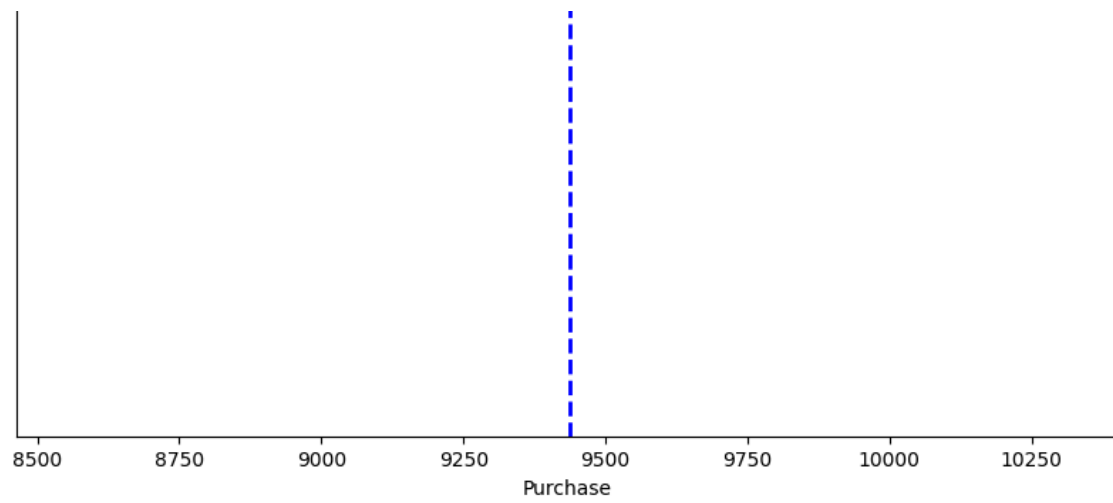








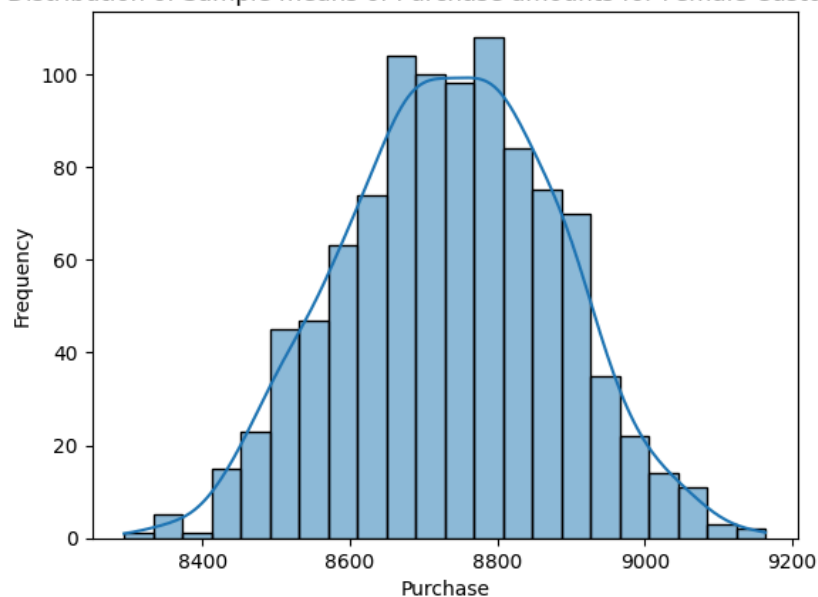




```
1 female_sample_means = plotSampleMeanDistributionAndConfidenceIntervals(females.Purchase, 1000, PopMeanFemales, confidence=0.95,  
2
```



Distribution of Sample means of Purchase amounts for Female Customers



Population Mean 8734.57 and 1000 95% Confidence Interval for Female Customers

- 35(3.5) Intervals that missed the population mean
- 965(96.5) Intervals that include the population mean
- Mean of the sample
- Population mean = 8734.57

