

Walmart - Confidence Interval and CLT

Walmart is an American multinational retail corporation that operates a chain of supercenters, discount departmental stores, and grocery stores in 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 behaviour (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).

1. Import the dataset and do the usual data analysis steps like checking the structure & characteristics of the dataset.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
df = pd.read_csv("/user/input/walmart/walmart_data.csv")
df.head()
df.info()
df.memory_usage()
df.describe()
```

Observations

- There are no missing values in the dataset.
- Purchase amount might have outliers.

2 Detect Null values & Outliers (using boxplot, 'describe' method by checking the difference between mean and median, isnull etc.)

How many users are there in the dataset?

```
df['User_ID'].nunique() In [ ]:
```

How many products are there?

```
df['Product_ID'].nunique() In [ ]:
```

Value_counts for the following:

- Gender
- Age

Walmart - Confidence Interval and CLT

- Occupation
- City_Category
- Stay_In_Current_City_Years
- Marital_Status
- Product_Category

```
In [ ]:  
categorical_cols = ['Gender', 'Age', 'Occupation', 'City_Category', 'Stay_In_Current_City_Years', 'Marital_Status', 'Product_Category']  
df[categorical_cols].melt().groupby(['variable', 'value'])['value'].count()/len(df)
```

Observations

- ~ 80% of the users are between the age 18-50 (40%: 26-35, 18%: 18-25, 20%: 36-45)
- 75% of the users are male and 25% are female
- 60% Single, 40% Married
- 35% Staying in the city from 1 year, 18% from 2 years, 17% from 3 years
- Total of 20 product categories are there
- There are 20 different types of occupations in the city

Univariate Analysis

Understanding the distribution of data and detecting outliers for continuous variables

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

Observation

- Purchase is having outliers

Understanding the distribution of data for the categorical variables

- Gender
- Age
- Occupation
- City_Category
- Stay_In_Current_City_Years

Walmart - Confidence Interval and CLT

- Marital_Status
- Product_Category

```
In [ ]:  
categorical_cols = ['Gender', 'Occupation', 'City_Category', 'Marital_Status', 'Product_Category']
```

```
fig, axs = plt.subplots(nrows=2, ncols=2, figsize=(16, 12))  
sns.countplot(data=df, x='Gender', ax=axs[0,0])  
sns.countplot(data=df, x='Occupation', ax=axs[0,1])  
sns.countplot(data=df, x='City_Category', ax=axs[1,0])  
sns.countplot(data=df, x='Marital_Status', ax=axs[1,1])  
plt.show()
```

```
plt.figure(figsize=(10, 8))  
sns.countplot(data=df, x='Product_Category')  
plt.show()
```

Observations

- Most of the users are Male
- There are 20 different types of Occupation and Product_Category
- More users belong to B City_Category
- More users are Single as compare to Married
- Product_Category - 1, 5, 8, & 11 have highest purchasing frequency.

```
In [ ]:  
fig, axs = plt.subplots(nrows=1, ncols=2, figsize=(12, 8))  
  
data = df['Age'].value_counts(normalize=True)*100  
palette_color = sns.color_palette('BrBG_r')  
axs[0].pie(x=data.values, labels=data.index, autopct='%.0f%%', colors=palette_color)  
axs[0].set_title("Age")  
  
data = df['Stay_In_Current_City_Years'].value_counts(normalize=True)*100  
palette_color = sns.color_palette('YlOrRd_r')  
axs[1].pie(x=data.values, labels=data.index, autopct='%.0f%%', colors=palette_color)  
axs[1].set_title("Stay_In_Current_City_Years")
```

```
plt.show()
```

Upper two graphs are self-explanatory.

Bi-variate Analysis

```
In [ ]:  
attrs = ['Gender', 'Age', 'Occupation', 'City_Category', 'Stay_In_Current_City_Years', 'Marital_Status', 'Product_Category']  
sns.set_style("white")
```

```
fig, axs = plt.subplots(nrows=3, ncols=2, figsize=(20, 16))  
fig.subplots_adjust(top=1.3)
```

Walmart - Confidence Interval and CLT

```
count = 0
for row in range(3):
    for col in range(2):
        sns.boxplot(data=df, y='Purchase', x=attrs[count], ax=axes[row, col], palette='Set3')
        axes[row, col].set_title(f"Purchase vs {attrs[count]}", pad=12, fontsize=13)
        count += 1
plt.show()

plt.figure(figsize=(10, 8))
sns.boxplot(data=df, y='Purchase', x=attrs[-1], palette='Set3')
plt.show()
```

Multivariate Analysis

```
In [ ]:
fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(20, 6))
fig.subplots_adjust(top=1.5)
sns.boxplot(data=df, y='Purchase', x='Gender', hue='Age', palette='Set3', ax=axes[0, 0])
sns.boxplot(data=df, y='Purchase', x='Gender', hue='City_Category', palette='Set3', ax=axes[0, 1])

sns.boxplot(data=df, y='Purchase', x='Gender', hue='Marital_Status', palette='Set3', ax=axes[1, 0])
sns.boxplot(data=df, y='Purchase', x='Gender', hue='Stay_In_Current_City_Years', palette='Set3', ax=axes[1, 1])
axes[1, 1].legend(loc='upper left')

plt.show()
```

```
In [ ]:
df.head(10)
```

Average amount spend per customer for Male and Female

```
In [ ]:
amt_df = df.groupby(['User_ID', 'Gender'])[['Purchase']].sum()
amt_df = amt_df.reset_index()
amt_df
```

```
In [ ]:
# Gender wise value counts in avg_amt_df
avg_amt_df['Gender'].value_counts()
```

```
In [ ]:
# histogram of average amount spend for each customer - Male & Female
amt_df[amt_df['Gender']=='M']['Purchase'].hist(bins=35)
plt.show()
```

```
amt_df[amt_df['Gender']=='F']['Purchase'].hist(bins=35)
plt.show()
```

```
In [ ]:
male_avg = amt_df[amt_df['Gender']=='M']['Purchase'].mean()
female_avg = amt_df[amt_df['Gender']=='F']['Purchase'].mean()

print("Average amount spend by Male customers: {:.2f}".format(male_avg))
print("Average amount spend by Female customers: {:.2f}".format(female_avg))
```

Walmart - Confidence Interval and CLT

Observation

1. Male customers spend more money than female customers

```
In [ ]:
male_df = amt_df[amt_df['Gender']=='M']
female_df = amt_df[amt_df['Gender']=='F']

In [ ]:
genders = ["M", "F"]

male_sample_size = 3000
female_sample_size = 1500
num_repitions = 1000
male_means = []
female_means = []

for _ in range(num_repitions):
    male_mean = male_df.sample(male_sample_size, replace=True)['Purchase'].mean()
    female_mean = female_df.sample(female_sample_size, replace=True)['Purchase'].mean()

    male_means.append(male_mean)
    female_means.append(female_mean)

In [ ]:
fig, axis = plt.subplots(nrows=1, ncols=2, figsize=(20, 6))

axis[0].hist(male_means, bins=35)
axis[1].hist(female_means, bins=35)
axis[0].set_title("Male - Distribution of means, Sample size: 3000")
axis[1].set_title("Female - Distribution of means, Sample size: 1500")

plt.show()

In [ ]:
print("Population mean - Mean of sample means of amount spend for Male: {:.2f}".format(np.mean(male_means)))
print("Population mean - Mean of sample means of amount spend for Female: {:.2f}".format(np.mean(female_means)))

print("\nMale - Sample mean: {:.2f} Sample std: {:.2f}".format(male_df['Purchase'].mean(), male_df['Purchase'].std()))
print("Female - Sample mean: {:.2f} Sample std: {:.2f}".format(female_df['Purchase'].mean(), female_df['Purchase'].std()))
```

Observation

Now using the **Central Limit Theorem** for the **population** we can say that:

1. Average amount spend by male customers is **9,26,341.86**
2. Average amount spend by female customers is **7,11,704.09**

```
In [ ]:
male_margin_of_error_clt = 1.96*male_df['Purchase'].std()/np.sqrt(len(male_df))
male_sample_mean = male_df['Purchase'].mean()
male_lower_lim = male_sample_mean - male_margin_of_error_clt
male_upper_lim = male_sample_mean + male_margin_of_error_clt
```

Walmart - Confidence Interval and CLT

```
female_margin_of_error_clt = 1.96*female_df['Purchase'].std()/np.sqrt(len(female_d
f))
female_sample_mean = female_df['Purchase'].mean()
female_lower_lim = female_sample_mean - female_margin_of_error_clt
female_upper_lim = female_sample_mean + female_margin_of_error_clt

print("Male confidence interval of means: ({:.2f}, {:.2f})".format(male_lower_lim,
male_upper_lim))
print("Female confidence interval of means: ({:.2f}, {:.2f})".format(female_lower_
lim, female_upper_lim))
```

Now we can infer about the population that, **95% of the times**:

1. Average amount spend by male customer will lie in between: **(895617.83, 955070.97)**
2. Average amount spend by female customer will lie in between: **(673254.77, 750794.02)**

Doing the same activity for married vs unmarried

```
amt_df                                                    In [ ]:

amt_df = df.groupby(['User_ID', 'Marital_Status'])[['Purchase']].sum()
amt_df = amt_df.reset_index()
amt_df                                                    In [ ]:

amt_df['Marital_Status'].value_counts()                  In [ ]:

marid_samp_size = 3000
unmarid_sample_size = 2000
num_repitions = 1000
marid_means = []
unmarid_means = []

for _ in range(num_repitions):
    marid_mean = amt_df[amt_df['Marital_Status']==1].sample(marid_samp_size, repla
ce=True)['Purchase'].mean()
    unmarid_mean = amt_df[amt_df['Marital_Status']==0].sample(unmarid_sample_size,
replace=True)['Purchase'].mean()

    marid_means.append(marid_mean)
    unmarid_means.append(unmarid_mean)

fig, axis = plt.subplots(nrows=1, ncols=2, figsize=(20, 6))

axis[0].hist(marid_means, bins=35)
axis[1].hist(unmarid_means, bins=35)
axis[0].set_title("Married - Distribution of means, Sample size: 3000")
axis[1].set_title("Unmarried - Distribution of means, Sample size: 2000")

plt.show()
```

Walmart - Confidence Interval and CLT

```
print("Population mean - Mean of sample means of amount spend for Married: {:.2f}"
      .format(np.mean(marid_means)))
print("Population mean - Mean of sample means of amount spend for Unmarried: {:.2f}"
      .format(np.mean(unmarid_means)))

print("\nMarried - Sample mean: {:.2f} Sample std: {:.2f}".format(amt_df[amt_df['M
arital_Status']==1]['Purchase'].mean(), amt_df[amt_df['Marital_Status']==1]['Purch
ase'].std()))
print("Unmarried - Sample mean: {:.2f} Sample std: {:.2f}".format(amt_df[amt_df['M
arital_Status']==0]['Purchase'].mean(), amt_df[amt_df['Marital_Status']==0]['Purch
ase'].std()))

In [ ]:

for val in ["Married", "Unmarried"]:

    new_val = 1 if val == "Married" else 0

    new_df = amt_df[amt_df['Marital_Status']==new_val]

    margin_of_error_clt = 1.96*new_df['Purchase'].std()/np.sqrt(len(new_df))
    sample_mean = new_df['Purchase'].mean()
    lower_lim = sample_mean - margin_of_error_clt
    upper_lim = sample_mean + margin_of_error_clt

    print("{} confidence interval of means: ({:.2f}, {:.2f})".format(val, lower_li
m, upper_lim))
```

Calculating the average amount spent by Age

```
In [ ]:

amt_df = df.groupby(['User_ID', 'Age'])[['Purchase']].sum()
amt_df = amt_df.reset_index()
amt_df

In [ ]:

amt_df['Age'].value_counts()

In [ ]:

sample_size = 200
num_repitions = 1000

all_means = {}

age_intervals = ['26-35', '36-45', '18-25', '46-50', '51-55', '55+', '0-17']
for age_interval in age_intervals:
    all_means[age_interval] = []

for age_interval in age_intervals:
    for _ in range(num_repitions):
        mean = amt_df[amt_df['Age']==age_interval].sample(sample_size, replace=True
e)['Purchase'].mean()
        all_means[age_interval].append(mean)

In [ ]:

for val in ['26-35', '36-45', '18-25', '46-50', '51-55', '55+', '0-17']:

    new_df = amt_df[amt_df['Age']==val]

    margin_of_error_clt = 1.96*new_df['Purchase'].std()/np.sqrt(len(new_df))
```

Walmart - Confidence Interval and CLT

```
sample_mean = new_df['Purchase'].mean()
lower_lim = sample_mean - margin_of_error_clt
upper_lim = sample_mean + margin_of_error_clt

print("For age {} --> confidence interval of means: ({:.2f}, {:.2f})".format(v
al, lower_lim, upper_lim))
```

Insights

- ~ 80% of the users are between the age 18-50 (40%: 26-35, 18%: 18-25, 20%: 36-45)
 - 75% of the users are Male and 25% are Female
 - 60% Single, 40% Married
 - 35% Staying in the city from 1 year, 18% from 2 years, 17% from 3 years
 - Total of 20 product categories are there
 - There are 20 different types of occupations in the city
-
- Most of the users are Male
 - There are 20 different types of Occupation and Product_Category
 - More users belong to B City_Category
 - More users are Single as compare to Married
 - Product_Category - 1, 5, 8, & 11 have highest purchasing frequency.
-
- **Average amount** spend by **Male** customers: **925344.40**
 - **Average amount** spend by **Female** customers: **712024.39**

Confidence Interval by Gender

Now using the **Central Limit Theorem** for the **population**:

1. Average amount spend by **male** customers is **9,26,341.86**
2. Average amount spend by **female** customers is **7,11,704.09**

Now we can infer about the population that, **95% of the times**:

1. Average amount spend by **male** customer will lie in between: **(895617.83, 955070.97)**

Walmart - Confidence Interval and CLT

2. Average amount spend by **female** customer will lie in between: **(673254.77, 750794.02)**

Confidence Interval by Marital_Status

1. **Married** confidence interval of means: **(806668.83, 880384.76)**
2. **Unmarried** confidence interval of means: **(848741.18, 912410.38)**

Confidence Interval by Age

1. For **age 26-35** --> confidence interval of means: **(945034.42, 1034284.21)**
2. For **age 36-45** --> confidence interval of means: **(823347.80, 935983.62)**
3. For **age 18-25** --> confidence interval of means: **(801632.78, 908093.46)**
4. For **age 46-50** --> confidence interval of means: **(713505.63, 871591.93)**
5. For **age 51-55** --> confidence interval of means: **(692392.43, 834009.42)**
6. For **age 55+** --> confidence interval of means: **(476948.26, 602446.23)**
7. For **age 0-17** --> confidence interval of means: **(527662.46, 710073.17)**

linkcode

Recommendations

1. Men spent more money than women, So company should focus on retaining the male customers and getting more male customers.
2. **Product_Category - 1, 5, 8, & 11** have highest purchasing frequency. it means these are the products in these categories are liked more by customers. Company can focus on selling more of these products or selling more of the products which are purchased less.
3. **Unmarried** customers spend more money than married customers, So company should focus on acquisition of Unmarried customers.
4. Customers in the **age 18-45** spend more money than the others, So company should focus on acquisition of customers who are in the **age 18-45**
5. Male customers living in City_Category C spend more money than other male customers living in B or C, Selling more products in the

Walmart - Confidence Interval and CLT

City_Category C will help the company increase the revenue. **How many users are there in the dataset?**