# Walmart Business case study
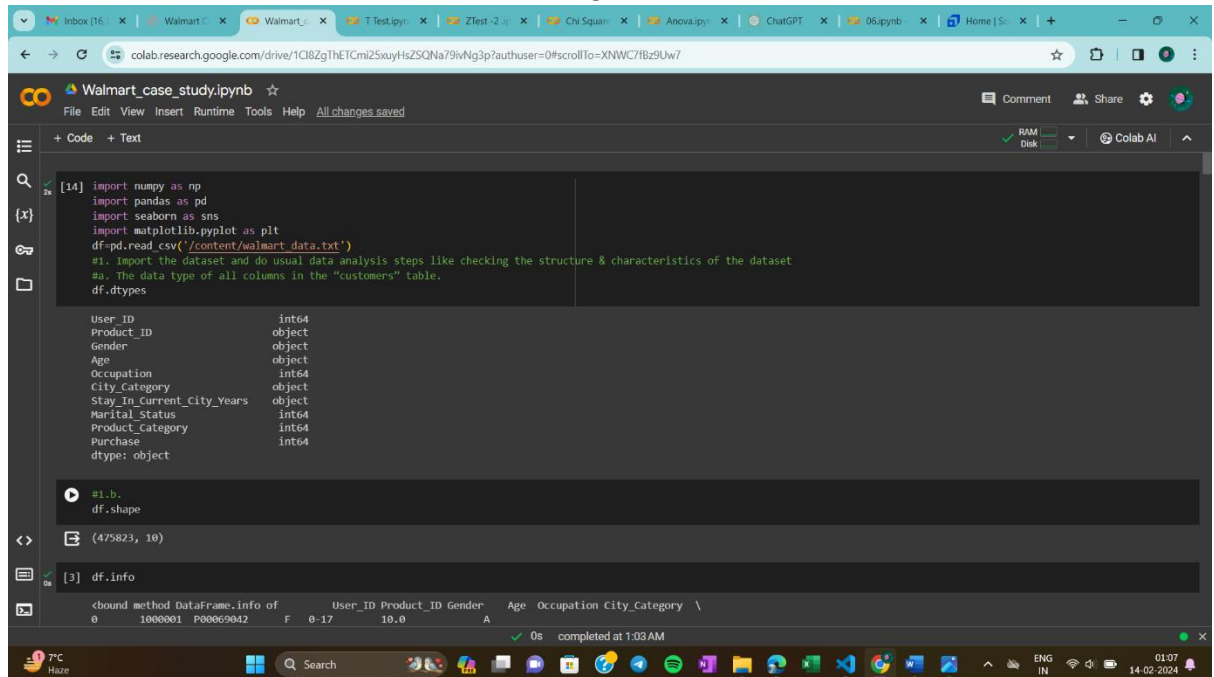
1. Import the dataset and do usual data analysis steps like checking the structure & characteristics of the dataset.
   a. The data type of all columns in the "customers" table.
   b. b. You can find the number of rows and columns given in the dataset



   c. Check for the missing values and find the number of missing values in each column.



2. Detect Null values and outliers

   a. Find the outliers for every continuous variable in the dataset.

colab.research.google.com/drive/1CI8ZgThETCmi25xuyHsZSQNa79ivNg3p?authuser=0#scrollTo=5w7zxutP3Ozq

**Walmart_case_study.ipynb** ☆

File Edit View Insert Runtime Tools Help  All changes saved

+ Code  + Text

```python
q_1=df['Purchase'].quantile(0.25)
q_3=df['Purchase'].quantile(0.75)
IQR=q_3-q_1
lower_bound=IQR-IQR*1.5
upper_bound=IQR+IQR*1.5
outliers =[value for value in df['Purchase'] if value < lower_bound or value >upper_bound]
sns.boxplot(outliers)
```

<Axes: >



completed at 5:45 AM

---

colab.research.google.com/drive/1CI8ZgThETCmi25xuyHsZSQNa79ivNg3p?authuser=0#scrollTo=5w7zxutP3Ozq

**Walmart_case_study.ipynb** ☆

File Edit View Insert Runtime Tools Help  All changes saved

+ Code  + Text

```python
q_1=df['Product_Category'].quantile(0.25)
q_3=df['Product_Category'].quantile(0.75)
IQR=q_3-q_1
lower_bound=IQR-IQR*1.5
upper_bound=IQR+IQR*1.5
outliers =[value for value in df['Product_Category'] if value < lower_bound or value >upper_bound]
sns.boxplot(outliers)
```

<Axes: >



completed at 5:45 AM

b. Remove/clip the data between the 5 percentile and 95 percentiles.





3. Data Exploration

a. What products are different age groups buying?

```python
#3. Data Exploration
#a. What products are different age groups buying?
import seaborn as sns
import matplotlib.pyplot as plt

# Plot a countplot of product purchases for each age group
plt.figure(figsize=(10, 8))
sns.countplot(data=df, x='Product_Category', hue='Age')
plt.title('Product Purchases by Age Group')
plt.xlabel('Product Category')
plt.ylabel('Count')
plt.xticks(rotation=45)  # Rotate x-axis labels for better readability
plt.legend(title='Age Group')
plt.show()
```



```python
#3. Data Exploration
#a. What products are different age groups buying?(another way)
plt.figure(figsize=(10, 6))
sns.countplot(data=df, x='Age', hue='Product_Category')
plt.title('Products Purchased by Age Groups')
plt.xlabel('Age Group')
plt.ylabel('Count')
plt.legend(title='Product Category')
plt.show()
```

b. Is there a relationship between age, marital status, and the amount spent?





c. Are there preferred product categories for different genders?

```
#3.c. Are there preferred product categories for different genders?
# Separate data by gender
male_data = df[df['Gender'] == 'M']
female_data = df[df['Gender'] == 'F']

plt.figure(figsize=(12, 6))
plt.subplot(1, 2, 1)
sns.countplot(data=male_data, x='Product_Category', palette='Reds')
plt.title('Product Category Distribution for Males')
plt.xlabel('Product Category')
plt.ylabel('Count')

plt.subplot(1, 2, 2)
sns.countplot(data=female_data, x='Product_Category', palette='Purples')
plt.title('Product Category Distribution for Females')
plt.xlabel('Product Category')
plt.ylabel('Count')

plt.tight_layout()
plt.show()
```

```
<ipython-input-11-2a83270179d5>:8: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

  sns.countplot(data=male_data, x='Product_Category', palette='Reds')
<ipython-input-11-2a83270179d5>:14: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

  sns.countplot(data=female_data, x='Product_Category', palette='Purples')
```



4. How does gender affect the amount spent?

```python
#4. How does gender affect the amount spent?
mean_purchase_by_gender = df.groupby('Gender')['Purchase'].mean()
std_purchase_by_gender = df.groupby('Gender')['Purchase'].std()
# Calculate confidence interval for the entire dataset
# Calculate confidence interval for each gender separately
confidence_intervals = {}
for gender in df['Gender'].unique():
    gender_data = df[df['Gender'] == gender]['Purchase']
    n = gender_data.count()
    mean_purchase = gender_data.mean()
    std_purchase = gender_data.std()
    z_score = stats.norm.ppf(0.975)  # 95% confidence interval
    margin_of_error = z_score * (std_purchase / np.sqrt(n))
    confidence_intervals[gender] = (mean_purchase - margin_of_error, mean_purchase + margin_of_error)
# Print confidence intervals
for gender, interval in confidence_intervals.items():
    print(f"{gender}: {interval}")
# Bootstrap function to calculate confidence intervals
def bootstrap_confidence_interval(data, num_samples, alpha=0.05):
    bootstrapped_means = []
    for _ in range(num_samples):
        bootstrap_sample = np.random.choice(data, size=len(data), replace=True)
        bootstrapped_means.append(np.mean(bootstrap_sample))
    lower_percentile = (alpha / 2) * 100
    upper_percentile = (1 - alpha / 2) * 100
    lower_bound = np.percentile(bootstrapped_means, lower_percentile)
    upper_bound = np.percentile(bootstrapped_means, upper_percentile)
    return lower_bound, upper_bound

# Perform bootstrapping and calculate confidence intervals for different sample sizes
sample_sizes = [300, 3000, 30000]
```

✓ Connected to Python 3 Google Compute Engine backend

```python
# Bootstrap function to calculate confidence intervals
def bootstrap_confidence_interval(data, num_samples, alpha=0.05):
    bootstrapped_means = []
    for _ in range(num_samples):
        bootstrap_sample = np.random.choice(data, size=len(data), replace=True)
        bootstrapped_means.append(np.mean(bootstrap_sample))
    lower_percentile = (alpha / 2) * 100
    upper_percentile = (1 - alpha / 2) * 100
    lower_bound = np.percentile(bootstrapped_means, lower_percentile)
    upper_bound = np.percentile(bootstrapped_means, upper_percentile)
    return lower_bound, upper_bound

# Perform bootstrapping and calculate confidence intervals for different sample sizes
sample_sizes = [300, 3000, 30000]
bootstrapped_intervals = {}
for gender in df['Gender'].unique():
    data = df[df['Gender'] == gender]['Purchase']
    bootstrapped_intervals[gender] = {}
    for sample_size in sample_sizes:
        lower_bound, upper_boun (variable) gender: Any _interval(data, num_samples=sample_size)
        bootstrapped_intervals[gender][sample_size] = (lower_bound, upper_bound)

print("Confidence Intervals for the Entire Dataset:")
for gender, interval in confidence_intervals.items():
    print(f"{gender}: {interval}")

print("\nConfidence Intervals for Bootstrapped Samples (Sample Size: 300, 3000, 30000):")
for gender, intervals in bootstrapped_intervals.items():
    print(f"{gender}:")
    for sample_size, interval in intervals.items():
        print(f"  Sample Size {sample_size}: {interval}")
```

✓ Connected to Python 3 Google Compute Engine backend

colab.research.google.com/drive/1CI8ZgThETCmi25xuyHsZSQNa79ivNg3p?authuser=0#scrollTo=_qfWpR6s1CAo

**Walmart_case_study.ipynb**

File  Edit  View  Insert  Runtime  Tools  Help

+ Code   + Text

```
print("\nConfidence Intervals for Bootstrapped Samples (Sample Size: 300, 3000, 30000):")
for gender, intervals in bootstrapped_intervals.items():
    print(f"{gender}:")
    for sample_size, interval in intervals.items():
        print(f"  Sample Size {sample_size}: {interval}")
```

```
F: (8611.598235475956, 8848.987146921872)
M: (9327.270601912533, 9469.283972516645)
nan: (nan, nan)
/usr/local/lib/python3.10/dist-packages/numpy/core/fromnumeric.py:3432: RuntimeWarning: Mean of empty slice.
  return _methods._mean(a, axis=axis, dtype=dtype,
/usr/local/lib/python3.10/dist-packages/numpy/core/_methods.py:190: RuntimeWarning: invalid value encountered in double_scalars
  ret = ret.dtype.type(ret / rcount)
Confidence Intervals for the Entire Dataset:
F: (8611.598235475956, 8848.987146921872)
M: (9327.270601912533, 9469.283972516645)
nan: (nan, nan)

Confidence Intervals for Bootstrapped Samples (Sample Size: 300, 3000, 30000):
F:
  Sample Size 300: (8617.379192809904, 8853.219370018653)
  Sample Size 3000: (8610.06420213668, 8848.664778701035)
  Sample Size 30000: (8612.529226725454, 8850.247702221468)
M:
  Sample Size 300: (9329.261308062072, 9476.465728355714)
  Sample Size 3000: (9328.312211975548, 9468.718837191076)
  Sample Size 30000: (9327.895379852658, 9469.15304091123)
nan:
  Sample Size 300: (nan, nan)
  Sample Size 3000: (nan, nan)
  Sample Size 30000: (nan, nan)
```

Connected to Python 3 Google Compute Engine backend

a.  From the above calculated CLT answer the following questions.

   i.    Is the confidence interval computed using the entire dataset wider for one of the
         genders? Why is this the case?
         Yes, the ci for entire dataset is wider for the male gender compared to the female
         gender. This is because the std dev of purchase amounts for males is generally higher
         than that for females, resulting in a wider interval.

   ii.   How is the width of the confidence interval affected by the sample size?
         As the sample size increases, the width of ci generally increases .(due to CLT)

   iii.  Do the confidence intervals for different sample sizes overlap?
         Yes, the ci for different sample size overlap

   iv.   How does the sample size affect the shape of the distributions of the means?
         As the sample size increase the shape of distribution of the means become more
         normally distributed.

## 5. How does Marital Status affect the amount spent?



```python
#5. How does Marital_Status affect the amount spent?
mean_purchase_by_Marital_Status = df.groupby('Marital_Status')['Purchase'].mean()
std_purchase_by_Marital_Status = df.groupby('Marital_Status')['Purchase'].std()

confidence_intervals = {}
for Marital_status in df['Marital_Status'].unique():
    Marital_Status_data = df[df['Marital_Status'] == Marital_Status]['Purchase']
    n = Marital_Status_data.count()
    mean_purchase = Marital_Status_data.mean()
    std_purchase = Marital_Status_data.std()
    z_score = stats.norm.ppf(0.975)  # 95% confidence interval
    margin_of_error = z_score * (std_purchase / np.sqrt(n))
    confidence_intervals[Marital_Status] = (mean_purchase - margin_of_error, mean_purchase + margin_of_error)

# Print confidence intervals
for Marital_Status, interval in confidence_intervals.items():
    print(f"{Marital_Status}: {interval}")

# Bootstrap function to calculate confidence intervals
def bootstrap_confidence_interval(data, num_samples, alpha=0.05):
    bootstrapped_means = []
    for _ in range(num_samples):
        bootstrap_sample = np.random.choice(data, size=len(data), replace=True)
        bootstrapped_means.append(np.mean(bootstrap_sample))
    lower_percentile = (alpha / 2) * 100
    upper_percentile = (1 - alpha / 2) * 100
    lower_bound = np.percentile(bootstrapped_means, lower_percentile)
    upper_bound = np.percentile(bootstrapped_means, upper_percentile)
    return lower_bound, upper_bound

# Perform bootstrapping and calculate confidence intervals for different sample sizes
```



```python
# Bootstrap function to calculate confidence intervals
def bootstrap_confidence_interval(data, num_samples, alpha=0.05):
    bootstrapped_means = []
    for _ in range(num_samples):
        bootstrap_sample = np.random.choice(data, size=len(data), replace=True)
        bootstrapped_means.append(np.mean(bootstrap_sample))
    lower_percentile = (alpha / 2) * 100
    upper_percentile = (1 - alpha / 2) * 100
    lower_bound = np.percentile(bootstrapped_means, lower_percentile)
    upper_bound = np.percentile(bootstrapped_means, upper_percentile)
    return lower_bound, upper_bound

# Perform bootstrapping and calculate confidence intervals for different sample sizes
sample_sizes = [300, 3000, 30000]
bootstrapped_intervals = {}
for Marital_Status in df['Marital_Status'].unique():
    data = df[df['Marital_Status'] == Marital_Status]['Purchase']
    bootstrapped_intervals[Marital_Status] = {}
    for sample_size in sample_sizes:
        lower_bound, upper_bound = bootstrap_confidence_interval(data, num_samples=sample_size)
        bootstrapped_intervals[Marital_Status][sample_size] = (lower_bound, upper_bound)

# Print confidence intervals for the entire dataset and bootstrapped samples
print("Confidence Intervals for the Entire Dataset:")
for Marital_Status, interval in confidence_intervals.items():
    print(f"{Marital_Status}: {interval}")

print("\nConfidence Intervals for Bootstrapped Samples (Sample Size: 300, 3000, 30000):")
for Marital_Status, intervals in bootstrapped_intervals.items():
    print(f"{Marital_Status}:")
    for sample_size, interval in intervals.items():
        print(f"    Sample Size {sample_size}: {interval}")
```

```
0.0: (9130.489781252667, 9289.254780711011)
1.0: (9191.726553149285, 9383.531869323971)
nan: (nan, nan)
/usr/local/lib/python3.10/dist-packages/numpy/core/fromnumeric.py:3432: RuntimeWarning: Mean of empty slice.
  return _methods._mean(a, axis=axis, dtype=dtype,
/usr/local/lib/python3.10/dist-packages/numpy/core/_methods.py:190: RuntimeWarning: invalid value encountered in double_scalars
  ret = ret.dtype.type(ret / rcount)
Confidence Intervals for the Entire Dataset:
0.0: (9130.489781252667, 9289.254780711011)
1.0: (9191.726553149285, 9383.531869323971)
nan: (nan, nan)

Confidence Intervals for Bootstrapped Samples (Sample Size: 300, 3000, 30000):
0.0:
  Sample Size 300: (9136.920010975853, 9291.862916583517)
  Sample Size 3000: (9125.323899088673, 9288.407949178474)
  Sample Size 30000: (9130.034698662941, 9288.890835162643)
1.0:
  Sample Size 300: (9193.18456213136, 9370.922273318005)
  Sample Size 3000: (9194.209277216834, 9380.34721583525)
  Sample Size 30000: (9193.597870638809, 9384.790775267418)
nan:
  Sample Size 300: (nan, nan)
  Sample Size 3000: (nan, nan)
  Sample Size 30000: (nan, nan)
```

   b.   From the above calculated CLT answer the following questions.

   1.Is the confidence interval computed using the entire dataset wider for one of the genders?     Why is this the case?

Yes, the ci for entire dataset is wider for the male gender compared to the female gender. This is because the std dev of purchase amounts for males is generally higher than that for females, resulting in a wider interval.

2. How is the width of the confidence interval affected by the sample size?

As the sample size increases, the width of ci generally increases .(due to CLT)

3.Do the confidence intervals for different sample sizes overlap?

Yes, the ci for different sample size overlap

4.How does the sample size affect the shape of the distributions of the means?

As the sample size increase the shape of distribution of the means become more normally distributed.

6. How does Age affect the amount spent?

```python
#How does Age affect the amount spent?
mean_age=df.groupby('Age')['Purchase'].mean()
std_age=df.groupby('Age')['Purchase'].std()

confidence_intervals = {}
for Age in df['Age'].unique():
    Age_data = df[df['Age'] == Age]['Purchase']
    n = Age_data.count()
    mean_purchase = Age_data.mean()
    std_purchase = Age_data.std()
    z_score = stats.norm.ppf(0.975)  # 95% confidence interval
    margin_of_error = z_score * (std_purchase / np.sqrt(n))
    confidence_intervals[Age] = (mean_purchase - margin_of_error, mean_purchase + margin_of_error)

# Print confidence intervals
for Age, interval in confidence_intervals.items():
    print(f"{Age}: {interval}")


# Bootstrap function to calculate confidence intervals
def bootstrap_confidence_interval(data, num_samples, alpha=0.05):
    bootstrapped_means = []
    for _ in range(num_samples):
        bootstrap_sample = np.random.choice(data, size=len(data), replace=True)
        bootstrapped_means.append(np.mean(bootstrap_sample))
    lower_percentile = (alpha / 2) * 100
    upper_percentile = (1 - alpha / 2) * 100
    lower_bound = np.percentile(bootstrapped_means, lower_percentile)
    upper_bound = np.percentile(bootstrapped_means, upper_percentile)
    return lower_bound, upper_bound
```

```python
# Perform bootstrapping and calculate confidence intervals for different sample sizes
sample_sizes = [300, 3000, 30000]
bootstrapped_intervals = {}
for Age in df['Age'].unique():
    data = df[df['Age'] == Age]['Purchase']
    bootstrapped_intervals[Age] = {}
    for sample_size in sample_sizes:
        lower_bound, upper_bound = bootstrap_confidence_interval(data, num_samples=sample_size)
        bootstrapped_intervals[Age][sample_size] = (lower_bound, upper_bound)

# Print confidence intervals for the entire dataset and bootstrapped samples
print("Confidence Intervals for the Entire Dataset:")
for Age, interval in confidence_intervals.items():
    print(f"{Age}: {interval}")

print("\nConfidence Intervals for Bootstrapped Samples (Sample Size: 300, 3000, 30000):")
for Age, intervals in bootstrapped_intervals.items():
    print(f"{Age}:")
    for sample_size, interval in intervals.items():
        print(f"  Sample Size {sample_size}: {interval}")
```

```
Confidence Intervals for the Entire Dataset:
0-17: (8600.747254268435, 9346.465907391079)
55+: (8905.737678498861, 9489.331008369825)
26-35: (9153.369299461547, 9349.377087674622)
46-50: (9069.720562284421, 9504.098309732932)
51-55: (9222.54753418532, 9696.106090670837)
36-45: (9118.640484429454, 9395.392166944797)
18-25: (9022.703949885052, 9296.507422663968)
nan: (nan, nan)

Confidence Intervals for Bootstrapped Samples (Sample Size: 300, 3000, 30000):
0-17:
    Sample Size 300: (8623.6613018598, 9343.610836909871)
    Sample Size 3000: (8596.907761087268, 9346.9833369098713)
    Sample Size 30000: (8602.069277539342, 9343.884334763949)
55+:
    Sample Size 300: (8906.304747474747, 9486.58659090909)
    Sample Size 3000: (8902.14747474747475, 9491.727020202019)
    Sample Size 30000: (8909.11143939394, 9489.052247474749)
26-35:
    Sample Size 300: (9157.347804294728, 9347.820537349888)
    Sample Size 3000: (9150.392644514553, 9348.705103806229)
    Sample Size 30000: (9152.081370852839, 9348.523514146143)
46-50:
    Sample Size 300: (9118.150664316703, 9498.43786605206)
    Sample Size 3000: (9075.621203904555, 9498.816133405639)
    Sample Size 30000: (9067.650311822124, 9506.423766268981)
51-55:
    Sample Size 300: (9204.185500575373, 9681.628552934408)
    Sample Size 3000: (9219.395368239355, 9698.198158803223)
    Sample Size 30000: (9225.964643268124, 9696.105681818182)
36-45:
    Sample Size 300: (9133.898982227733, 9381.440623062615)
    Sample Size 3000: (9120.56864021492, 9395.914770613763)
    Sample Size 30000: (9116.137605910311, 9395.413132878695)
```

c.From the above calculated CLT answer the following questions.

1.Is the confidence interval computed using the entire dataset wider for one of the genders?    Why is this the case?

Yes, the ci for entire dataset is wider for the male gender compared to the female gender. This is because the std dev of purchase amounts for males is generally higher than that for females, resulting in a wider interval.

2. How is the width of the confidence interval affected by the sample size?

As the sample size increases, the width of ci generally increases. (due to CLT)

3.Do the confidence intervals for different sample sizes overlap?

Yes, the ci for different sample size overlap

4.How does the sample size affect the shape of the distributions of the means?

As the sample size increase the shape of distribution of the means become more normally distributed.

7. Create a report

a. Report whether the confidence intervals for the average amount spent by males and females (computed using all the data) overlap. How can Walmart leverage this conclusion to make changes or improvements?

No, ci for male and female do not overlap. This suggests that there is significant difference in amount spent by male or female.

Walmart can make different marketing strategies for male and female which can be more efficient and can lead to higher sales.

c. Report whether the confidence intervals for the average amount spent by married and unmarried (computed using all the data) overlap. How can Walmart leverage this conclusion to make changes or improvements?

Yes, ci for married and unmarried do not overlap. This suggests that there significant difference in amount spent by male or female.

Walmart can make multiple marketing strategies for male and female which can be more efficient and can lead to higher sales.



d. Report whether the confidence intervals for the average amount spent by different age groups (computed using all the data) overlap. How can Walmart leverage this conclusion to make changes or improvements?

Yes, ci for different age groups do not overlap. This suggests that there is significant difference in amount spent by male or female.

Walmart can make different marketing strategies for male and female which can be more efficient and can lead to higher sales.