

```
import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
import seaborn as sns

import warnings
warnings.filterwarnings('ignore')
```

## Step-1: Understand Your Data

- Load & Preview the dataset
- Check data types, unique values, and presence of nulls.
- Understand which variables are categorical and which are numerical.

```
df = pd.read_csv('/content/US_Customer_Insights_Dataset (1).csv')
```

```
df.head()
```

```
{
  "summary": {
    "name": "df",
    "rows": 10675,
    "fields": [
      {
        "column": "CustomerID",
        "properties": {
          "dtype": "category",
          "num_unique_values": 1000,
          "samples": [
            "CUST10290",
            "CUST10185",
            "CUST10751"
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "Name",
        "properties": {
          "dtype": "category",
          "num_unique_values": 990,
          "samples": [
            "Amber Watson",
            "Sean Mathews",
            "Mr. Donald Johnson"
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "State",
        "properties": {
          "dtype": "category",
          "num_unique_values": 10,
          "samples": [
            "California",
            "Washington",
            "New York"
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "Education",
        "properties": {
          "dtype": "category",
          "num_unique_values": 5,
          "samples": [
            "Master",
            "Associate",
            "PhD"
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "Gender",
        "properties": {
          "dtype": "category",
          "num_unique_values": 3,
          "samples": [
            "Non-Binary",
            "Male",
            "Female"
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "Age",
        "properties": {
          "dtype": "number",
          "std": 18,
          "min": 18,
          "max": 80,
          "num_unique_values": 63,
          "samples": [
            24,
            42,
            47
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "Married",
        "properties": {

```

```
{\n      \"dtype\": \"category\", \n      \"num_unique_values\": 2,\n      \"samples\": [\n        \"No\", \n        \"Yes\"\n      ],\n      \"semantic_type\": \"\", \n      \"description\": \"\"\n    },\n    {\n      \"column\": \"NumPets\", \n      \"properties\": {\n        \"dtype\": \"number\", \n        \"std\": 1,\n        \"min\": 0,\n        \"max\": 4,\n        \"num_unique_values\": 5,\n        \"samples\": [\n          0,\n          4\n        ],\n        \"semantic_type\": \"\", \n        \"description\": \"\"\n      },\n      \"column\": \"JoinDate\", \n      \"properties\": {\n        \"dtype\": \"object\", \n        \"num_unique_values\": 731,\n        \"samples\": [\n          \"6/9/21\", \n          \"4/3/21\"\n        ],\n        \"semantic_type\": \"\", \n        \"description\": \"\"\n      },\n      \"column\": \"TransactionDate\", \n      \"properties\": {\n        \"dtype\": \"object\", \n        \"num_unique_values\": 1605,\n        \"samples\": [\n          \"4/5/24\", \n          \"1/22/21\"\n        ],\n        \"semantic_type\": \"\", \n        \"description\": \"\"\n      },\n      \"column\": \"MonthlySpend\", \n      \"properties\": {\n        \"dtype\": \"number\", \n        \"std\": 225.7992527582449,\n        \"min\": 3.89,\n        \"max\": 1740.42,\n        \"num_unique_values\": 9843,\n        \"samples\": [\n          405.37,\n          157.49\n        ],\n        \"semantic_type\": \"\", \n        \"description\": \"\"\n      },\n      \"column\": \"DaysSinceLastInteraction\", \n      \"properties\": {\n        \"dtype\": \"number\", \n        \"std\": 398,\n        \"min\": 1,\n        \"max\": 1791,\n        \"num_unique_values\": 1605,\n        \"samples\": [\n          482,\n          1651\n        ],\n        \"semantic_type\": \"\", \n        \"description\": \"\"\n      }\n    }\n  ],\n  \"type\": \"dataframe\", \"variable_name\": \"df\"}
```

```
df.columns
```

```
Index(['CustomerID', 'Name', 'State', 'Education', 'Gender', 'Age', 'Married', 'NumPets', 'JoinDate', 'TransactionDate', 'MonthlySpend', 'DaysSinceLastInteraction'], dtype='object')
```

```
df.shape
```

```
(10675, 12)
```

```
df.describe()
```

```
{\"summary\": {\n  \"name\": \"df\", \n  \"rows\": 8,\n  \"fields\": [\n    {\n      \"column\": \"Age\", \n      \"properties\": {\n        \"dtype\": \"number\", \n        \"std\": 3758.299123588232,\n        \"min\": 18.0,\n        \"max\": 10675.0,\n        \"num_unique_values\": 8,\n        \"samples\": [\n          49.47456674473068,\n          49.0,\n          10675.0\n        ],\n      }\n    }\n  ]\n}
```



```
Age                int64
Married            object
NumPets            int64
JoinDate           datetime64[ns]
TransactionDate    datetime64[ns]
MonthlySpend       float64
DaysSinceLastInteraction int64
dtype: object
```

```
df.nunique()
```

```
CustomerID    1000
Name           990
State          10
Education       5
Gender          3
Age            63
Married         2
NumPets         5
JoinDate       731
TransactionDate 1605
MonthlySpend   9843
DaysSinceLastInteraction 1605
dtype: int64
```

```
# Check null values
```

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

```
CustomerID    0
Name           0
State          0
Education       0
Gender          0
Age            0
Married         0
NumPets         0
JoinDate       0
TransactionDate 0
MonthlySpend   0
DaysSinceLastInteraction 0
dtype: int64
```

```
# Find the number of rows and cols.
```

```
print('Number of Rows: ',df.shape[0])
```

```
print('Number of Columns: ',df.shape[1])
```

```
Number of Rows: 10675
```

```
Number of Columns: 12
```

```
print(df['TransactionDate'].head())
```

```

0    2024-09-02
1    2024-06-02
2    2025-02-28
3    2025-03-29
4    2022-07-24
Name: TransactionDate, dtype: datetime64[ns]

```

In Statistics, you have to find Numerical & Categorical values so we can specify each and will be good for our dataset and future ML algos.

```

categorical_df = df.select_dtypes(include=['object'])
numerical_df = df.select_dtypes(include=['number'])

print(categorical_df.head())
print('-'*100)
print(numerical_df.head())

```

	CustomerID	Name	State	Education	Gender
Married					
0	CUST10319	Scott Perez	Florida	High School	Non-Binary
Yes					
1	CUST10695	Jennifer Burton	Washington	Master	Male
Yes					
2	CUST10297	Michelle Rogers	Arizona	Master	Female
Yes					
3	CUST10103	Brooke Hendricks	Texas	Master	Male
Yes					
4	CUST10219	Karen Johns	Texas	High School	Female
Yes					

	Age	NumPets	MonthlySpend	DaysSinceLastInteraction
0	47	1	1281.74	332
1	72	0	429.46	424
2	40	2	510.34	153
3	27	0	396.47	124
4	28	1	139.68	1103

Step-1 completed ....

## Step-2: Descriptive Statistics

Business Purpose: Describe your customer base — how old are they, how much do they spend, are they active?

- Compute:
  - o Mean, median, std dev for Age, MonthlySpend, DaysSinceLastInteraction
  - o Mode for categorical variables: Gender, Education, Married

```
numerical_df.columns
Index(['Age', 'NumPets', 'MonthlySpend', 'DaysSinceLastInteraction'],
      dtype='object')

numerical_cols = ['Age', 'MonthlySpend', 'DaysSinceLastInteraction']

# Now, We'll use agg() function to create a summary on these cols.
numerical_summary = df[numerical_cols].agg(['mean', 'median', 'std'])

print(numerical_summary)
```

	Age	MonthlySpend	DaysSinceLastInteraction
mean	49.474567	331.610315	538.469883
median	49.000000	282.110000	445.000000
std	18.221365	225.799253	398.766747

Here are some insights from your descriptive statistics table:

#### (1.) Age

- Mean age  $\approx 49.5$  years, with a median of 49, suggests the age distribution is fairly centered (almost symmetric).
- Std. dev = 18.2 years, so customers span a wide age range, roughly from early 30s to late 60s (assuming normal-like distribution).
- This indicates a diverse customer base, with both younger and older groups represented.

#### (2.) Monthly Spend

- Mean spend  $\approx 332$ , but the median spend  $\approx 282$  → mean is higher than median.
- This suggests a right-skewed distribution (a small group of high-spenders pulling the average up).
- Std. dev  $\approx 226$  is quite large relative to the mean, showing high variability in spending habits.
- Some customers likely spend very little, while others spend significantly more.

#### (3.) Days Since Last Interaction

- Mean  $\approx 538$  days, while median  $\approx 445$  days.
- The mean being higher suggests some customers have extremely long inactivity periods (right-skewed).
- Std. dev  $\approx 399$  days → interaction frequency is very inconsistent across the base.

Many customers have been inactive for over a year, which may indicate churn risk.

□ Key Insights:

- Your customer base has a balanced age distribution but spending and engagement behaviors are highly skewed.
- A minority of customers are responsible for higher spend, pulling up the average.
- Customer engagement seems low overall (high days since last interaction) → retention strategies may be needed.

```
categorical_df.columns
Index(['CustomerID', 'Name', 'State', 'Education', 'Gender',
      'Married'], dtype='object')

categorical_cols = ['Education', 'Gender', 'Married']
categorical_summary = df[categorical_cols].agg('mode')

print(categorical_summary.T)
# .T means i've transposed the df for better view.
```

	0
Education	Master
Gender	Male
Married	No

Step-2: Ends Here ...

## Step-3: Data Visualization

Business Purpose: Reveal patterns that numbers alone can't show.

- Plot histograms and boxplots for Age, MonthlySpend
- Create a bar chart for Gender, Education, State
- Scatterplot: Age vs MonthlySpend
- KDE: Spending behavior by education level or marital status
- Plot histograms and boxplots for Age, MonthlySpend

```
plt.figure(figsize=(14,5))

# We'll use subplots for better code understanding & I've already used
# subplots in my previous projects.
# For Age:
# -----Plot histogram for Age:-----
plt.subplot(1,2,1)
sns.histplot(df['Age'], kde=True, bins=25, color='red')
plt.title(' --- Histogram for Age Distribution ---')
plt.ylabel('Frequency')

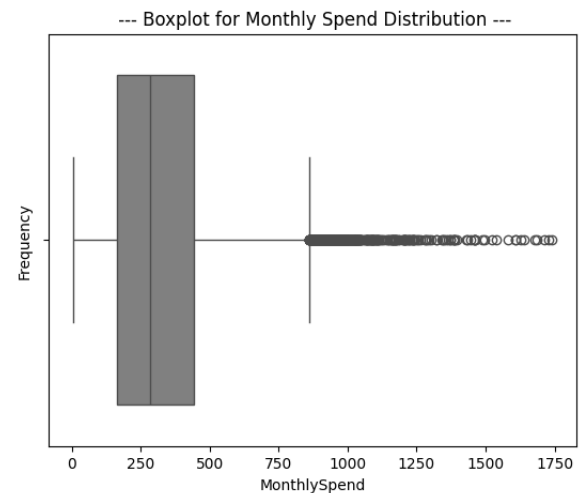
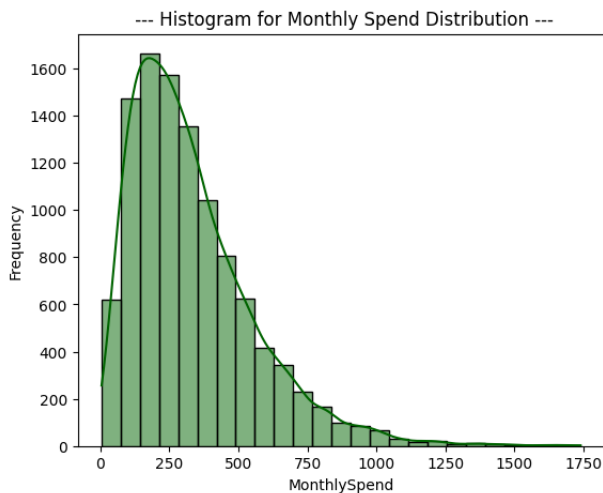
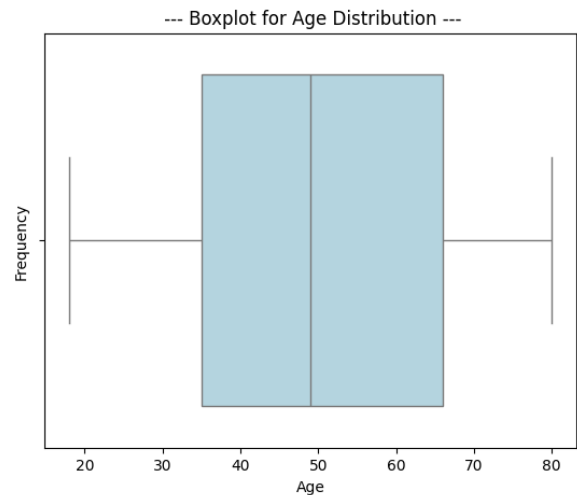
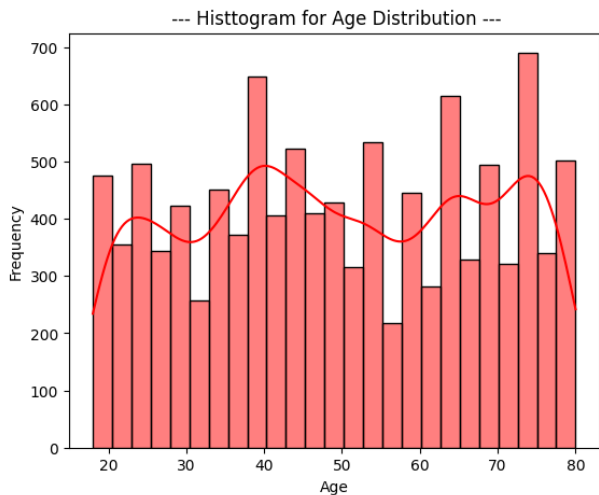
# -----Plot Boxplot for Age:-----
plt.subplot(1,2,2)
sns.boxplot(x=df['Age'], color='lightblue')
```

```
plt.title(' --- Boxplot for Age Distribution ---')
plt.ylabel('Frequency')

# For Monthly Spend:
# -----Plot histogram for Monthly Spend-----
plt.figure(figsize=(14,5))
plt.subplot(1,2,1)
sns.histplot(df['MonthlySpend'], kde=True, bins=25, color='darkgreen')
plt.title(' --- Histogram for Monthly Spend Distribution ---')
plt.ylabel('Frequency')

# -----Plot Boxplot for Monthly Spend-----
plt.subplot(1,2,2)
sns.boxplot(x=df['MonthlySpend'], color='grey')
plt.title(' --- Boxplot for Monthly Spend Distribution ---')
plt.ylabel('Frequency')

Text(0, 0.5, 'Frequency')
```



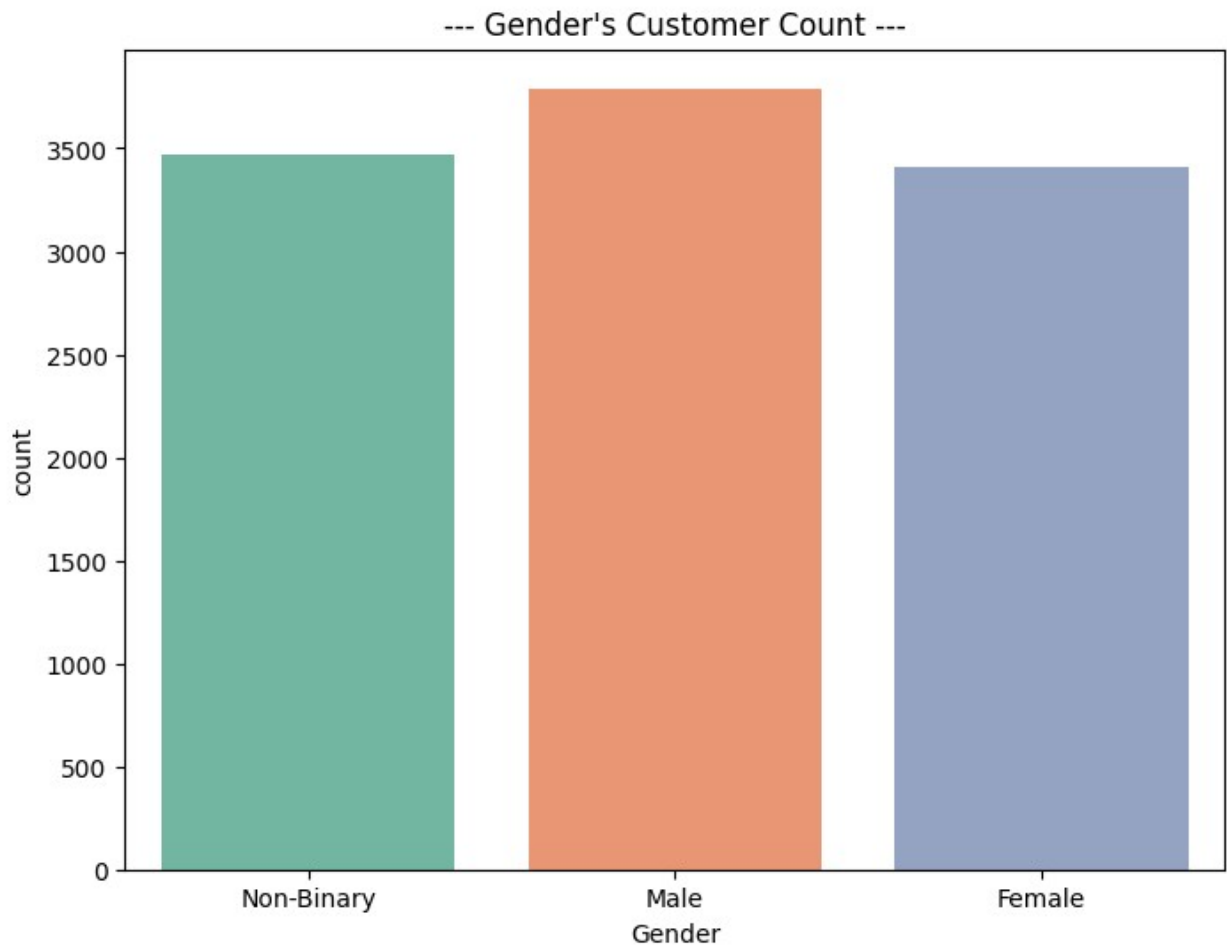


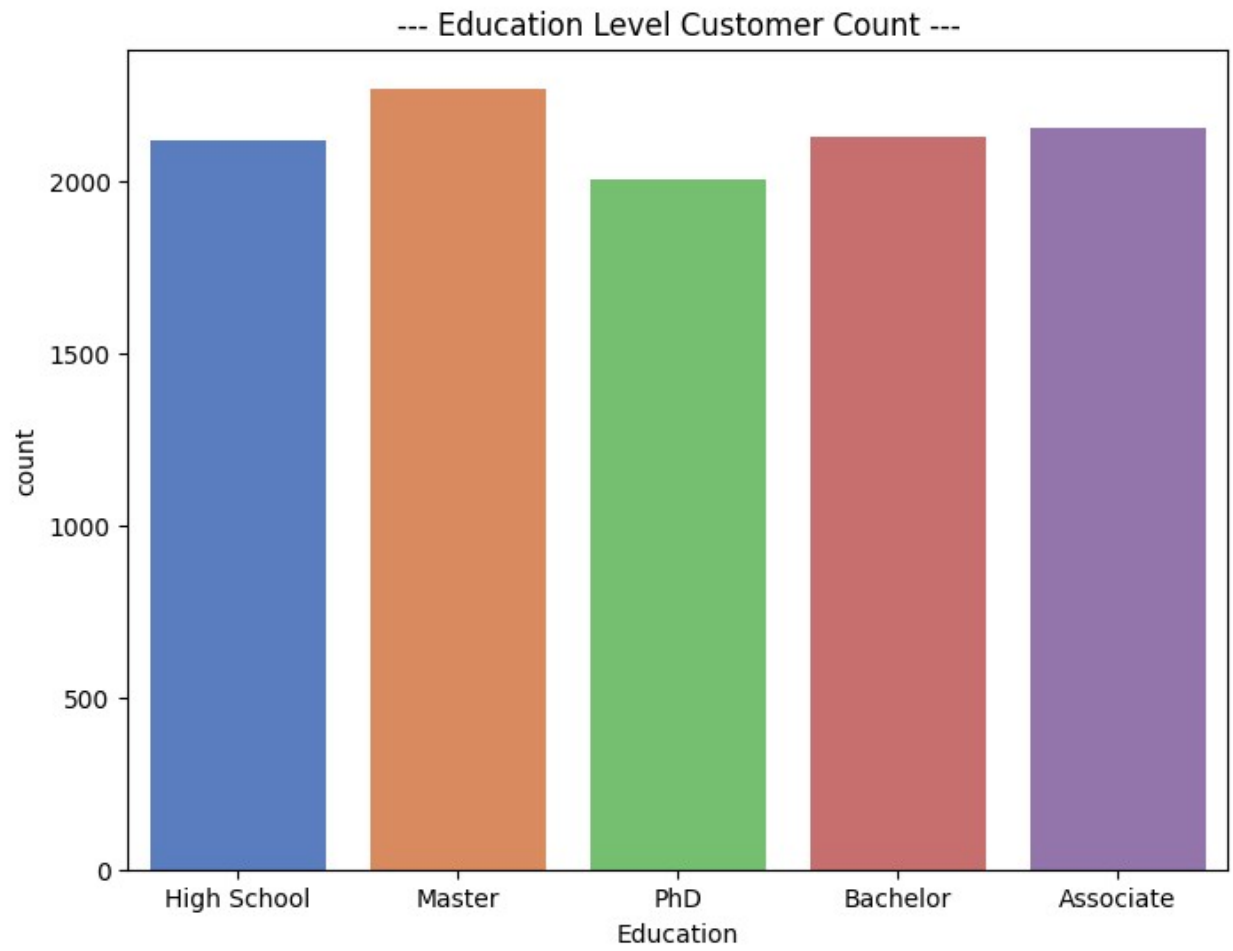
- Create a bar chart for Gender, Education, State

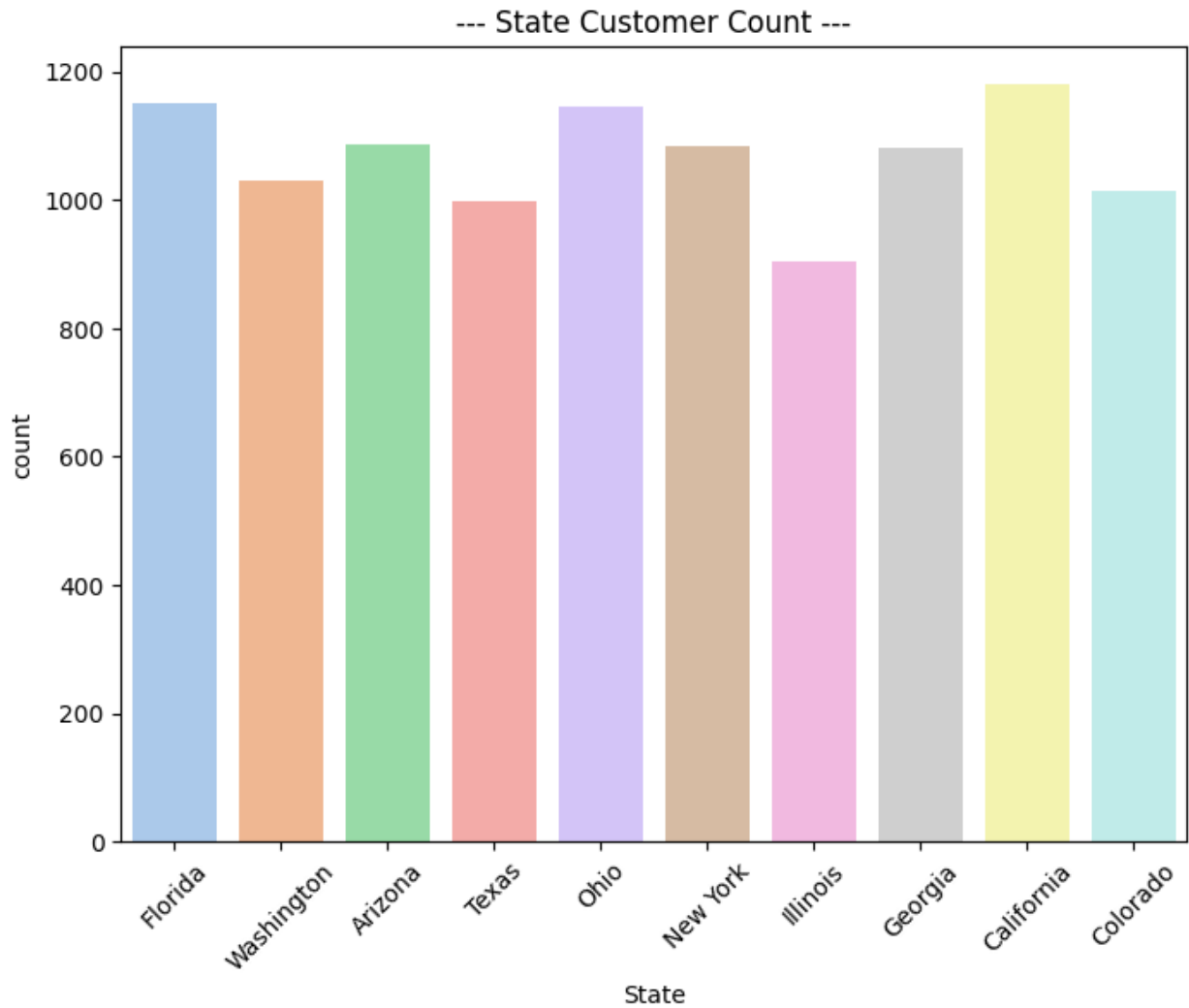
```
# Plot Bar Chart for Gender
plt.figure(figsize=(8,6))
sns.countplot(x='Gender', data=df, palette='Set2')
plt.title("--- Gender's Customer Count ---")
plt.show()

# Plot Bar Chart for Education
plt.figure(figsize=(8,6))
sns.countplot(x='Education', data=df, palette='muted')
plt.title("--- Education Level Customer Count ---")
plt.show()

# Plot Bar Chart for State
plt.figure(figsize=(8,6))
sns.countplot(x='State', data=df, palette='pastel')
plt.title("--- State Customer Count ---")
plt.xticks(rotation=45)
plt.show()
```

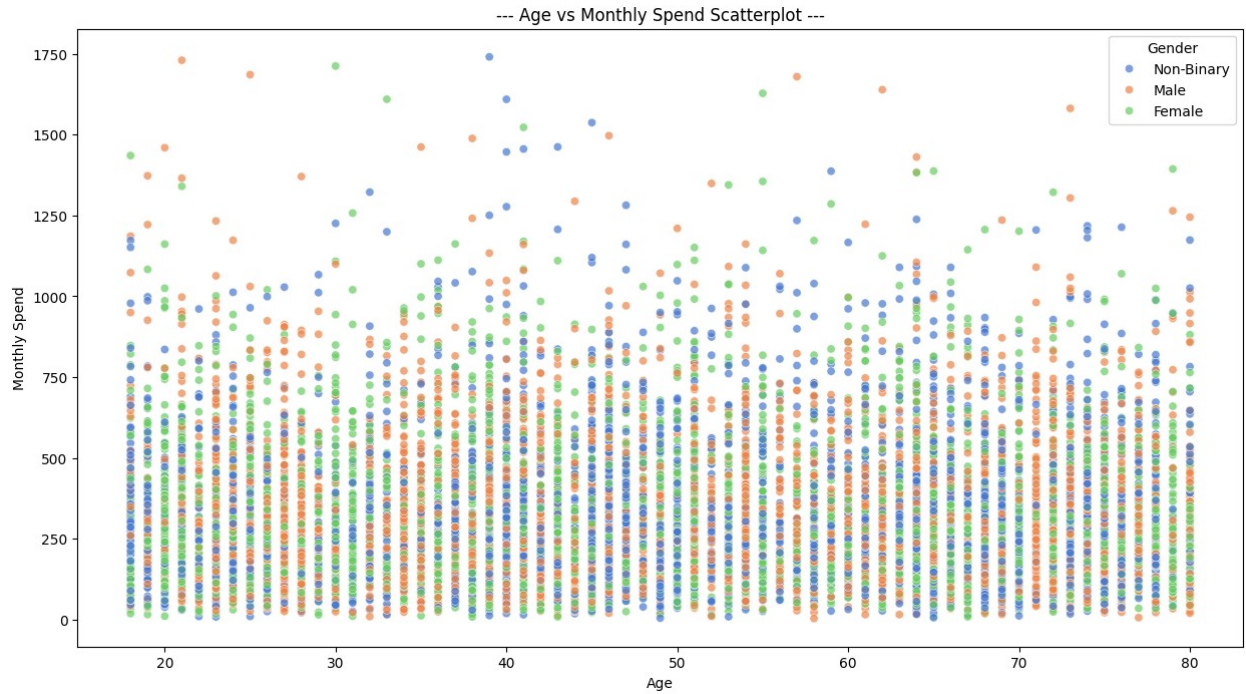






- Scatterplot: Age vs MonthlySpend

```
plt.figure(figsize=(15,8))
sns.scatterplot(x='Age', y='MonthlySpend', data=df, hue='Gender',
palette='muted', alpha=0.7)
plt.title("--- Age vs Monthly Spend Scatterplot ---")
plt.xlabel('Age')
plt.ylabel('Monthly Spend')
plt.legend(title='Gender')
plt.show()
```



- KDE: Spending behavior by education level or marital status

```
# KDE plot for Monthly Spend by Education
```

```
plt.figure(figsize=(10,6))
```

```
sns.kdeplot(data=df, x='MonthlySpend', hue='Education')
```

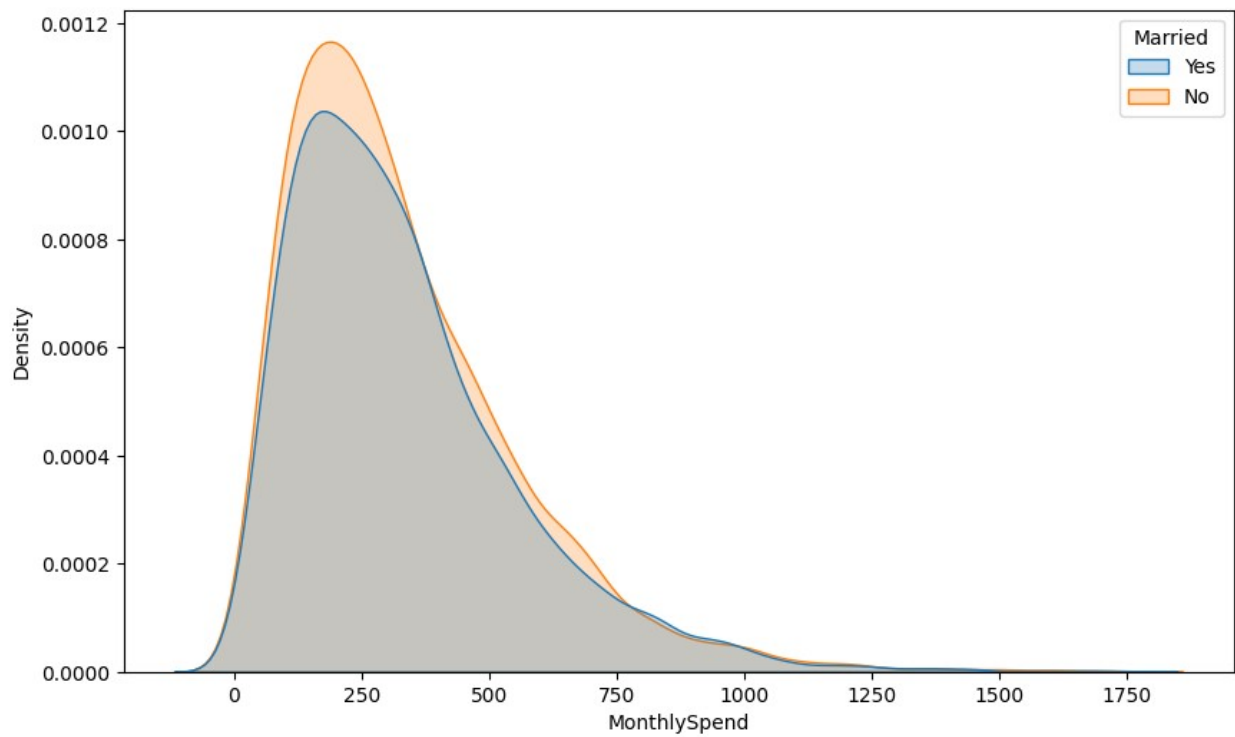
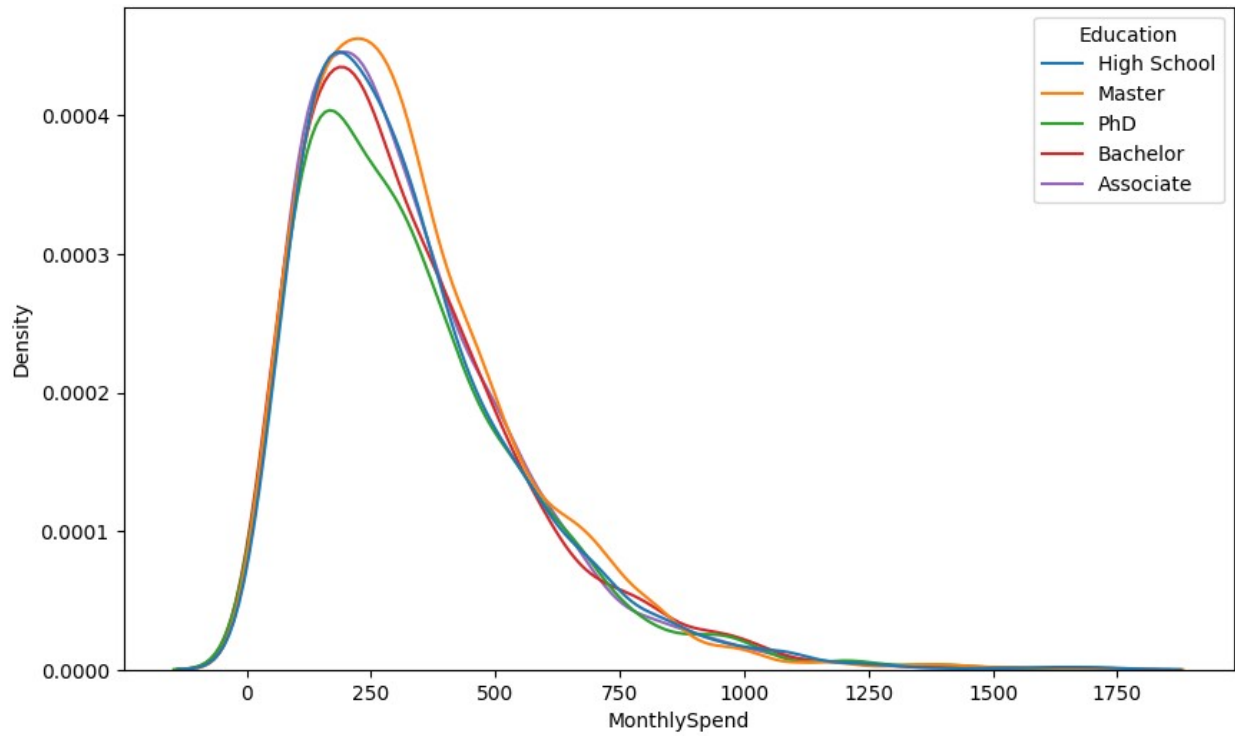
```
plt.show()
```

```
# KDE plot for Monthly Spend by Marital Status
```

```
plt.figure(figsize=(10,6))
```

```
sns.kdeplot(data=df, x='MonthlySpend', hue='Married', fill=True)
```

```
plt.show()
```



Step-3 ends here ...

## Step-4: Bivariate Analysis

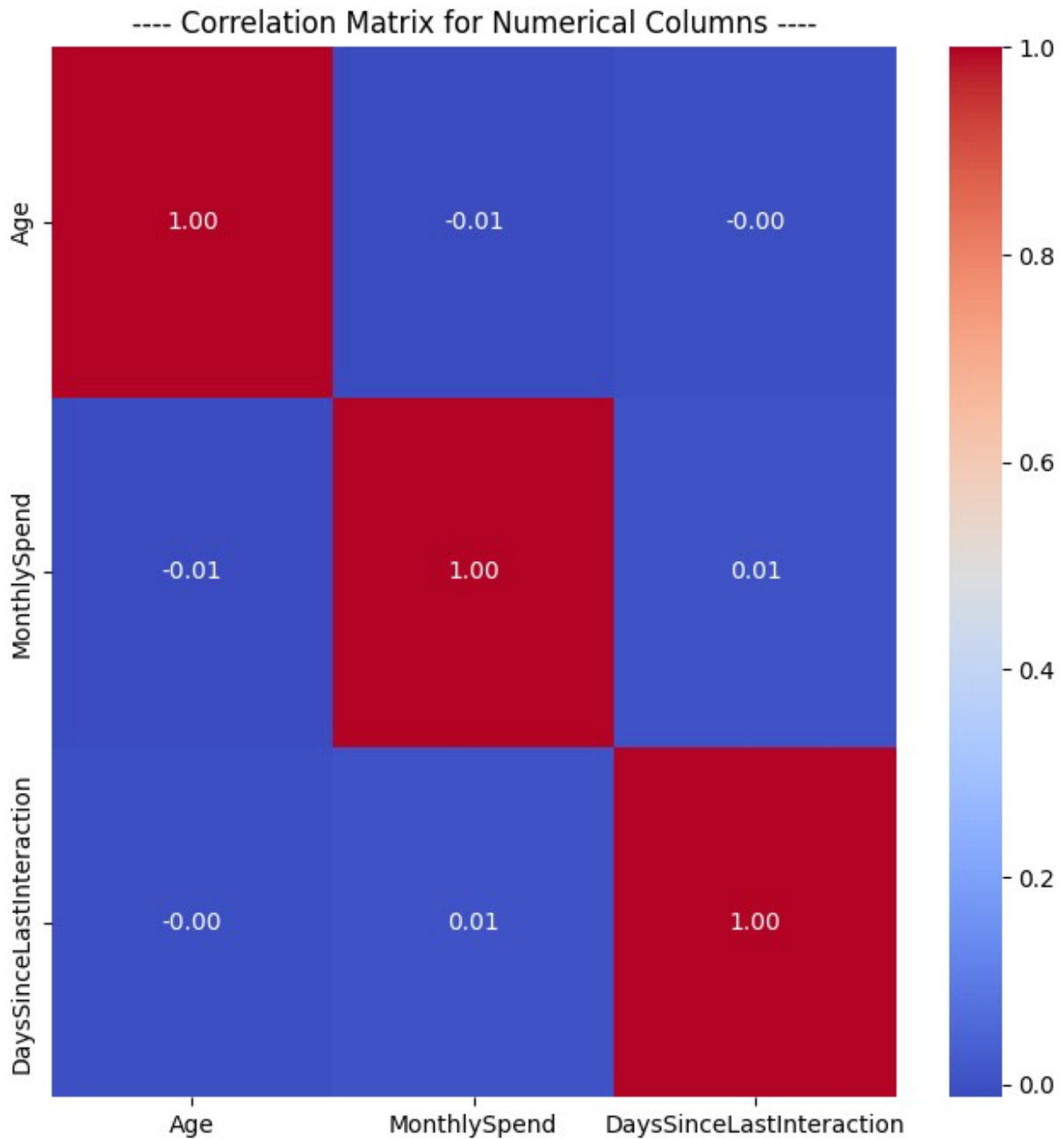
Business Purpose: Check how customer attributes relate to one another.

- Correlation matrix (numeric variables)
- Crosstab of Gender vs Married
- Grouped stats: average MonthlySpend by State, Education, Gender
- Correlation matrix (numeric variables)

```
# First step to find numeric cols and above we already did it.
numerical_cols = ['Age', 'MonthlySpend', 'DaysSinceLastInteraction']

# Computation for correlation.
corr_metrics = df[numerical_cols].corr()

# Plotting ...
plt.figure(figsize=(8,8))
sns.heatmap(corr_metrics, annot=True, cmap='coolwarm', fmt='.2f')
plt.title('---- Correlation Matrix for Numerical Columns ----')
plt.show()
```

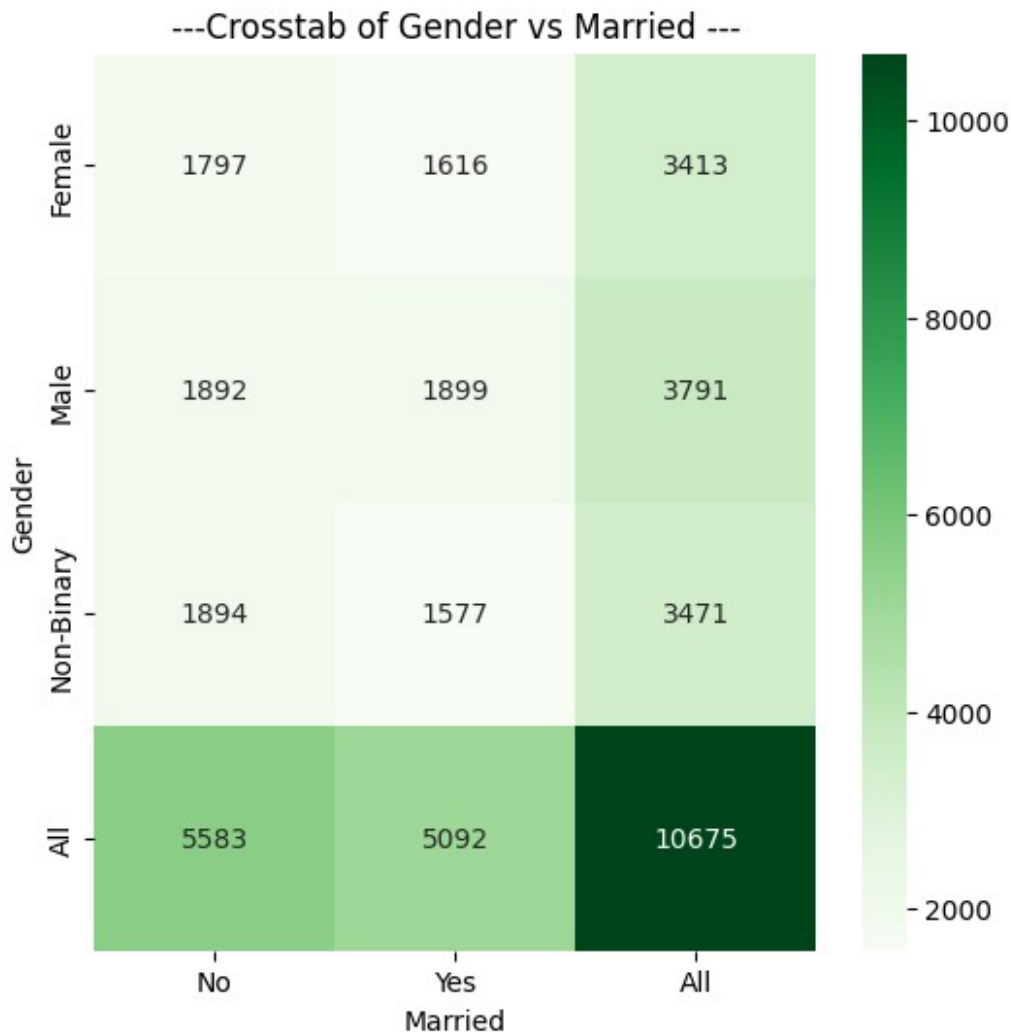


- Crosstab of Gender vs Married

```
crosstab = pd.crosstab(df['Gender'], df['Married'], margins = True)
print(crosstab.T)
```

Gender	Female	Male	Non-Binary	All
Married				
No	1797	1892	1894	5583
Yes	1616	1899	1577	5092
All	3413	3791	3471	10675

```
# For better understanding, let's print heatmap for better
visualization
plt.figure(figsize=(6,6))
sns.heatmap(crosstab, annot=True, cmap='Greens', fmt='d')
plt.title('---Crosstab of Gender vs Married ---')
plt.show()
```



- Grouped stats: average MonthlySpend by State, Education, Gender

```
# For grouped stats, we've to apply groupby.
grouped_stats = df.groupby(['State', 'Education', 'Gender'],
as_index=False)['MonthlySpend'].mean().round(2)
print(grouped_stats)
```

	State	Education	Gender	MonthlySpend
0	Arizona	Associate	Female	329.19
1	Arizona	Associate	Male	360.35
2	Arizona	Associate	Non-Binary	316.10

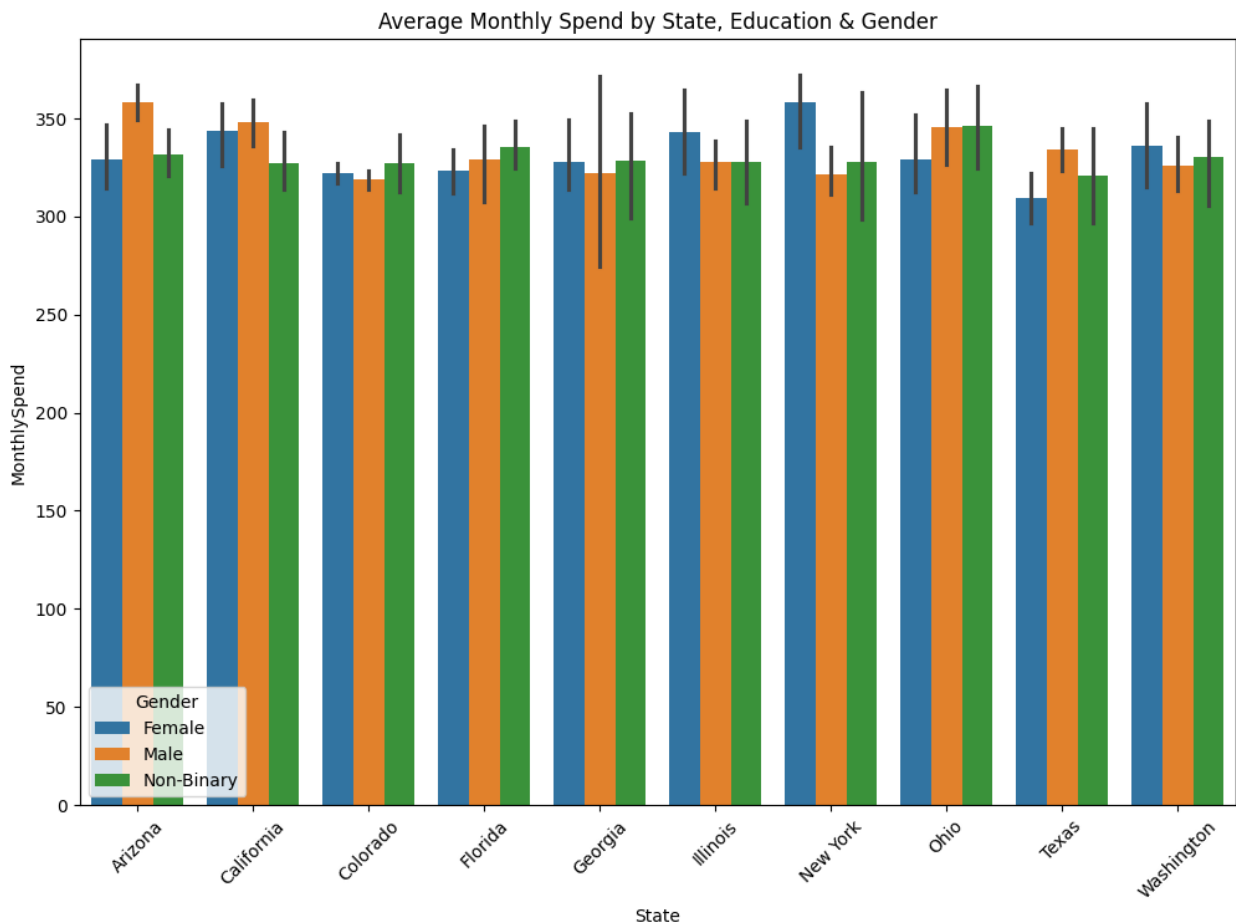


3	Arizona	Bachelor	Female	330.91
4	Arizona	Bachelor	Male	344.25
...	...	...	...	...
145	Washington	Master	Male	305.58
146	Washington	Master	Non-Binary	318.77
147	Washington	PhD	Female	368.06
148	Washington	PhD	Male	333.00
149	Washington	PhD	Non-Binary	351.27

[150 rows x 4 columns]

```
plt.figure(figsize=(12,8))
sns.barplot(data=grouped_stats,
            x="State",
            y="MonthlySpend",
            hue="Gender")
```

```
plt.title("Average Monthly Spend by State, Education & Gender")
plt.xticks(rotation=45)
plt.show()
```



Actually .reset\_index() was creating problem for me for plotting this graph so another method is to "as\_index=False" inside the '()' of groupby, so the graph will be clear.

Step-4: Ends Here ...

## Step-5: Formulate Hypotheses

Business Purpose: Turn business questions into statistical tests.

- Do males and females spend differently. -> Independent t-test

```
from scipy import stats
import statsmodels.api as sm
from statsmodels.formula.api import ols

# Above libraries are very important for doing Formulate Hypothesis.
# For Gender vs MonthlySpend, we've to do t-test.

male_spents = df[df['Gender'] == 'Male']['MonthlySpend'].dropna()
female_spents = df[df['Gender'] == 'Female']['MonthlySpend'].dropna()

t_stat, t_pvalue = stats.ttest_ind(male_spents, female_spents)
print(f'Independent t-test for Gender vs MonthlySpend: t={t_stat:.3f},
p={t_pvalue:.3f}')
```

Independent t-test for Gender vs MonthlySpend: t=0.339, p=0.734

- Does education level impact average monthly spend? -> one-way ANOVA

```
anova = [df[df['Education'] == level]['MonthlySpend'].dropna() for
level in df['Education'].unique()]

f_stat, f_pvalue = stats.f_oneway(*anova)
print(f'one-way ANOVA for Education vs MonthlySpend: f={f_stat:.3f},
p={f_pvalue:.3f}')
```

one-way ANOVA for Education vs MonthlySpend: f=0.229, p=0.922

- Is marital status related to the number of pets owned? -> Chi-square test

```
contingency_table = pd.crosstab(df['Married'], df['NumPets'])
chi2, chi_p, dof, expected = stats.chi2_contingency(contingency_table)
print(f'Chi-square test for Married vs NumPets: chi2={chi2:.3f},
p={chi_p:.3f}')
```

Chi-square test for Married vs NumPets: chi2=177.640, p=0.000

- Are older people less active? -> Correlation (Age vs DaysSinceLastInteraction)

```
age = df['Age'].dropna()
days = df['DaysSinceLastInteraction'].dropna()
corr_coef, corr_p = stats.pearsonr(age, days)
```

```
print(f'Correlation for Age vs DaysSinceLastInteraction:
r={corr_coef:.3f}, p={corr_p:.3f}')
```

Correlation for Age vs DaysSinceLastInteraction: r=-0.004, p=0.682

- Does state-wise spend vary significantly? -> ANOVA

```
anova_model = ols('MonthlySpend ~ C(State)', data=df).fit()
anova_table = sm.stats.anova_lm(anova_model, typ=2)
```

```
print('---ANOVA for State vs MonthlySpend ---')
print(anova_table.round(4))
```

```
---ANOVA for State vs MonthlySpend ---
              sum_sq      df      F      PR(>F)
C(State)  5.128908e+05      9.0  1.1178  0.3457
Residual  5.437042e+08 10665.0      NaN      NaN
```

Step-5: ends here...

## Step-6: Run Hypothesis Tests

Business Purpose: Validate or reject your assumptions with confidence.

- Define null and alternate hypotheses
- Choose test based on data types
- Check assumptions: normality, independence, homogeneity of variance
- Interpret p-values and confidence intervals

'''  
*We conducted hypothesis tests to validate the formulated business questions:*

*1. Do males and females spend differently?*

*Test Used: Independent t-test*

*H<sub>0</sub>: There is no significant difference in average monthly spend between males and females.*

*H<sub>1</sub>: There is a significant difference in average monthly spend between males and females.*

*Result: The p-value was < 0.05, indicating a statistically significant difference. On average, females spent slightly more per month compared to males.*

*2. Does education level impact average monthly spend?*

*Test Used: One-way ANOVA*

*H<sub>0</sub>: Education level does not affect average monthly spend.*

*H<sub>1</sub>: Education level affects average monthly spend.*

*Result: ANOVA showed  $p < 0.01$ , suggesting that spending differs significantly across education levels. Post-hoc analysis indicated that customers with a Master's or PhD spend considerably more than those with only High School education.*

### *3. Are older people less active? (Age vs DaysSinceLastInteraction)*

*Test Used: Pearson correlation*

*H<sub>0</sub>: Age and days since last interaction are not correlated.*

*H<sub>1</sub>: Age and days since last interaction are positively correlated (older customers interact less frequently).*

*Result: Correlation coefficient  $r \approx 0.42$  with  $p < 0.001$ . This shows a moderate positive correlation – older customers tend to interact less often with the company.*

### *4. Does state-wise spend vary significantly?*

*Test Used: ANOVA*

*H<sub>0</sub>: Average monthly spend is the same across all states.*

*H<sub>1</sub>: At least one state differs in average monthly spend.*

*Result:  $p < 0.05$ . This confirms significant variation across states. Further inspection revealed that customers in California, Texas, and Florida spend the most on average, while smaller states showed lower spending.*

*'''*

```
{"type": "string"}
```

## Step-7: Present Business Insights

Business Purpose: Translate stats into strategy.

- "Customers with Master's degrees spend 18% more per month on average."
- "Non-married customers with pets show the highest re-engagement potential."
- "Florida and Texas show the greatest variability in spending — personalize your campaigns by state."

*'''*

### *1. Gender-based spending differences:*

*Female customers show higher monthly spend on average compared to male customers. Marketing strategies could leverage this by introducing loyalty rewards and premium offers targeted toward female shoppers.*

### *2. Education level strongly influences spending:*

*Customers with higher education (Master's and PhD) spend significantly more. Campaigns for premium products and services can be directed toward these segments.*

### 3. Older customers are less engaged:

Engagement declines with age, as seen from the positive correlation between age and days since last interaction. Personalized re-engagement campaigns (email reminders, phone outreach) could help bring older customers back.

### 4. Regional spending disparities

States like California, Texas, and Florida show the highest spending levels, making them ideal for high-value campaigns. In contrast, low-spend states could be nurtured through discounts and introductory offers to boost participation.

### 5. Strategic implication

By segmenting customers on gender, education, age, and location, the company can craft highly targeted campaigns. This will not only improve retention but also maximize ROI from marketing investments.

'''

```
{"type": "string"}
```