

Name: Mohan Krishna Balaji

Reg No:192324221

Course Code: CSA0406

Course Name: Fundamental Data Science

11. **Scenario** : You are a data scientist working for a company that sells products online. You have been tasked with creating a simple plot to show the sales of a product over time.

Question:

1. Write code to create a simple line plot in Python using Matplotlib to predict sales happened in a month?
2. Write code to create a scatter plot in Python using Matplotlib to predict sales happened in a month?
3. Develop a Python program to create a bar plot of the monthly sales data.



12. Scenario: You are working on a data analysis project that involves analyzing the monthly temperature and rainfall data for a city. You have a dataset containing the monthly temperature and rainfall values for each month of a year. Your task is to develop a Python program that generates line plots and scatter plots to visualize the temperature and rainfall data.

Question:

1. Develop a Python program to create a line plot of the monthly temperature data.
2. Develop a Python program to create a scatter plot of the monthly rainfall data.

```
import matplotlib.pyplot as plt
import numpy as np

# Sample data for 12 months
months = ['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun', 'Jul', 'Aug', 'Sep', 'Oct', 'Nov', 'Dec']

# Example temperature (°C) and rainfall (mm) data
temperature = [5, 7, 12, 18, 23, 27, 30, 29, 25, 18, 10, 6] # Replace with actual temperature data
rainfall = [78, 60, 55, 42, 35, 20, 15, 18, 30, 50, 70, 85] # Replace with actual rainfall data

# Create figure and subplots
plt.figure(figsize=(14, 5))

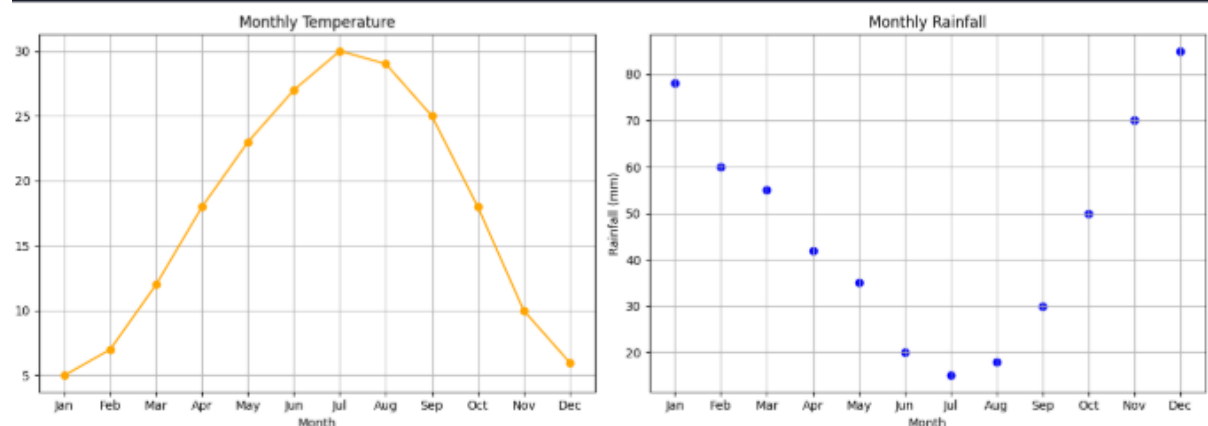
# 1. Line Plot - Monthly Temperature
plt.subplot(1, 2, 1)
plt.plot(months, temperature, marker='o', color='orange', linestyle='-')
plt.title("Monthly Temperature")
plt.xlabel("Month")
plt.ylabel("Temperature (°C)")
plt.grid(True)

# 2. Scatter Plot - Monthly Rainfall
plt.subplot(1, 2, 2)
plt.scatter(months, rainfall, color='blue')
plt.title("Monthly Rainfall")
plt.xlabel("Month")
plt.ylabel("Rainfall (mm)")
plt.grid(True)

# Layout adjustment and display
plt.tight_layout()
plt.show()
```

Click Run or press **shift + ENTER** to run code

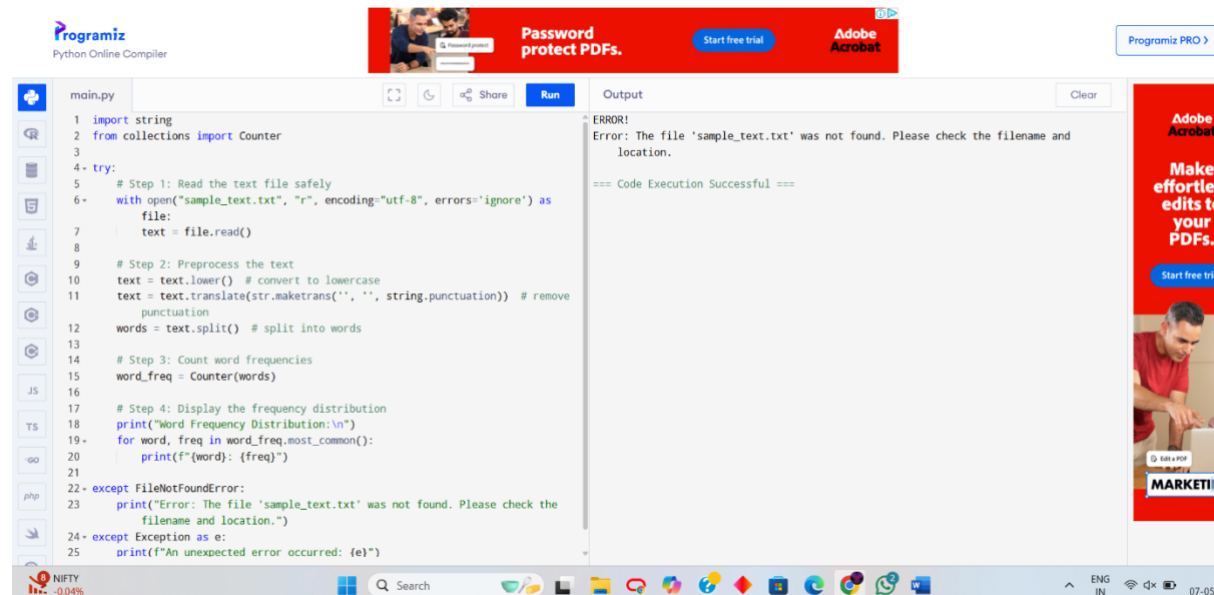
Enable code completions



13. **Scenario:** You are working on a text analysis project and need to determine the frequency distribution of words in a given text document. You have a text document named "sample_text.txt"

containing a paragraph of text. Your task is to develop a Python program that reads the text document, processes the text, and generates a frequency distribution of the words.

Question: How would you develop a Python program to calculate the frequency distribution of words in a text document?



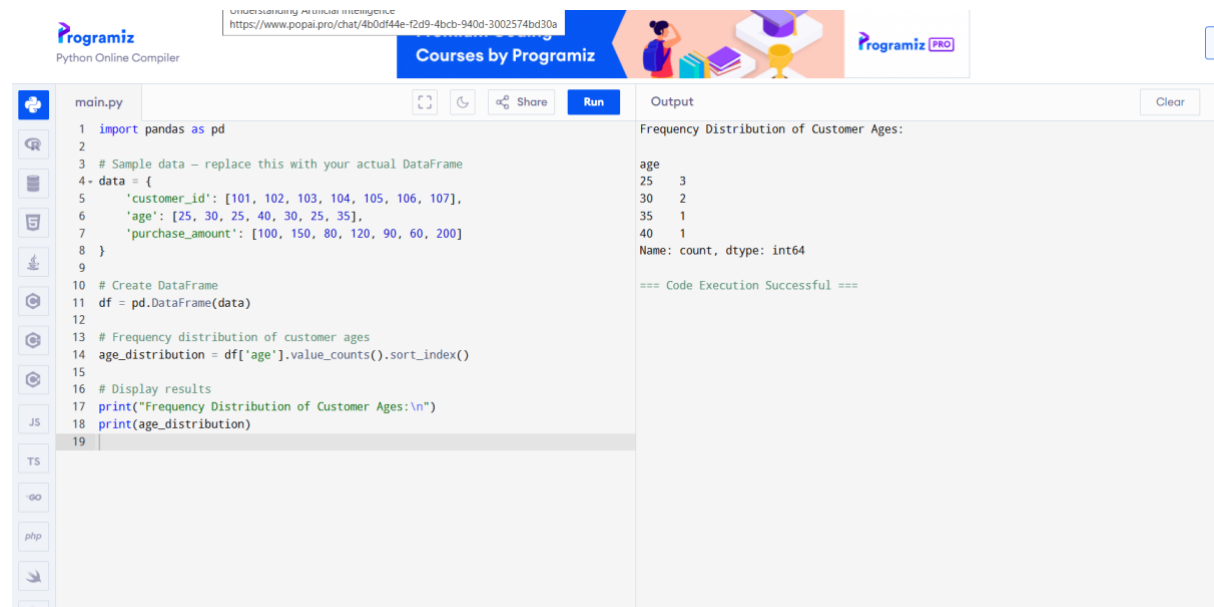
The screenshot shows the Programiz Python Online Compiler interface. The code in `main.py` is as follows:

```
1 import string
2 from collections import Counter
3
4- try:
5     # Step 1: Read the text file safely
6     with open("sample_text.txt", "r", encoding="utf-8", errors='ignore') as file:
7         text = file.read()
8
9     # Step 2: Preprocess the text
10    text = text.lower() # convert to lowercase
11    text = text.translate(str.maketrans('', '', string.punctuation)) # remove punctuation
12    words = text.split() # split into words
13
14    # Step 3: Count word frequencies
15    word_freq = Counter(words)
16
17    # Step 4: Display the frequency distribution
18    print("Word Frequency Distribution:\n")
19    for word, freq in word_freq.most_common():
20        print(f"{word}: {freq}")
21
22- except FileNotFoundError:
23    print("Error: The file 'sample_text.txt' was not found. Please check the filename and location.")
24- except Exception as e:
25    print(f"An unexpected error occurred: {e}")
```

The output window shows an error message: "ERROR! Error: The file 'sample_text.txt' was not found. Please check the filename and location." Below the error message, it says "=== Code Execution Successful ===".

14. **Scenario:** You are a data analyst working for a company that sells products online. You have been tasked with analyzing the sales data for the past month. The data is stored in a Pandas data frame.

Question: Develop a code in python to find the frequency distribution of the ages of the customers who have made a purchase in the past month.



The screenshot shows the Programiz Python Online Compiler interface. The code in `main.py` is as follows:

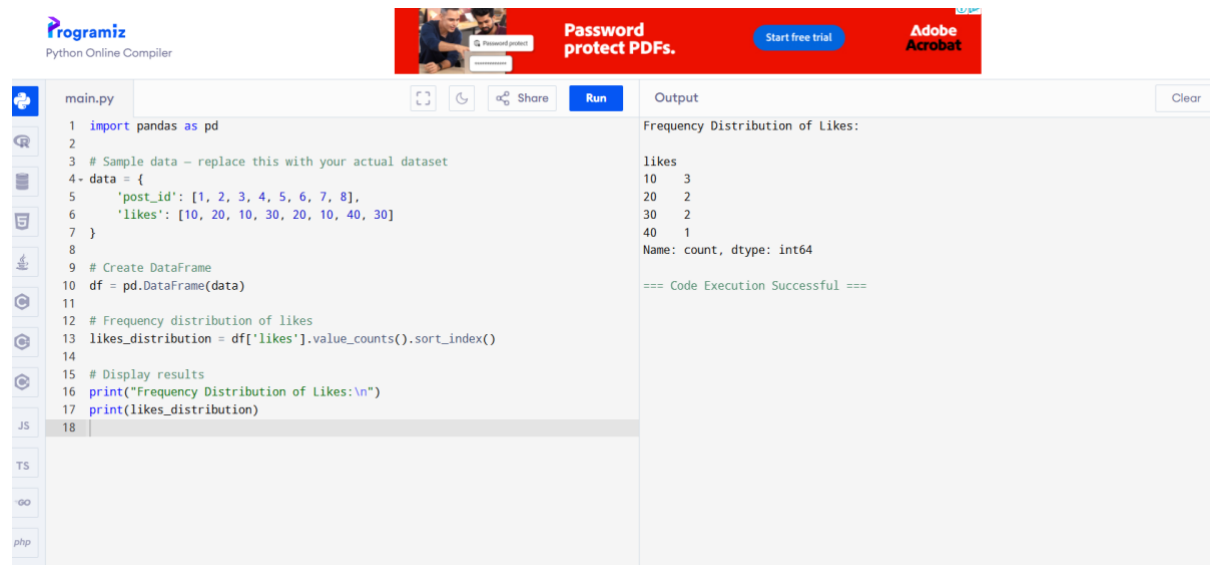
```
1 import pandas as pd
2
3 # Sample data - replace this with your actual DataFrame
4- data = {
5     'customer_id': [101, 102, 103, 104, 105, 106, 107],
6     'age': [25, 30, 25, 40, 30, 25, 35],
7     'purchase_amount': [100, 150, 80, 120, 90, 60, 200]
8 }
9
10 # Create DataFrame
11 df = pd.DataFrame(data)
12
13 # Frequency distribution of customer ages
14 age_distribution = df['age'].value_counts().sort_index()
15
16 # Display results
17 print("Frequency Distribution of Customer Ages:\n")
18 print(age_distribution)
```

The output window shows the frequency distribution of customer ages:

```
Frequency Distribution of Customer Ages:
age
25    3
30    2
35    1
40    1
Name: count, dtype: int64
=== Code Execution Successful ===
```

15. **Scenario:** You are a data analyst working for a social media platform. As part of your analysis, you have a dataset containing user interaction data, including the number of likes received by each post. Your task is to develop a Python program that calculates the frequency distribution of likes among the posts.

Question: Develop a Python program to calculate the frequency distribution of likes among the posts?



The screenshot shows the Programiz Python Online Compiler interface. The code in the editor is as follows:

```

1 import pandas as pd
2
3 # Sample data - replace this with your actual dataset
4 data = {
5     'post_id': [1, 2, 3, 4, 5, 6, 7, 8],
6     'likes': [10, 20, 10, 30, 20, 10, 40, 30]
7 }
8
9 # Create DataFrame
10 df = pd.DataFrame(data)
11
12 # Frequency distribution of likes
13 likes_distribution = df['likes'].value_counts().sort_index()
14
15 # Display results
16 print("Frequency Distribution of Likes:\n")
17 print(likes_distribution)

```

The output on the right shows the frequency distribution of likes:

```

Frequency Distribution of Likes:

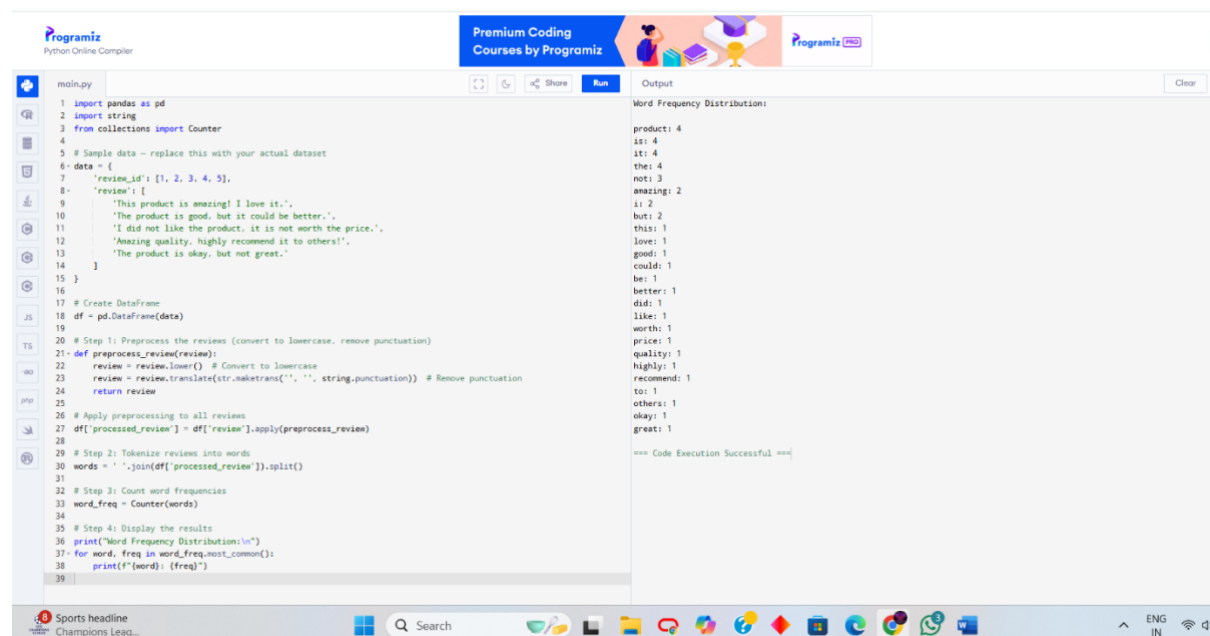
likes
10    3
20    2
30    2
40    1
Name: count, dtype: int64

```

Below the output, it says "=== Code Execution Successful ===".

16. Scenario: You are working on a project that involves analyzing customer reviews for a product. You have a dataset containing customer reviews, and your task is to develop a Python program that calculates the frequency distribution of words in the reviews.

Question: Develop a Python program to calculate the frequency distribution of words in the customer reviews dataset?



The screenshot shows the Programiz Python Online Compiler interface. The code in the editor is as follows:

```

1 import pandas as pd
2 import string
3 from collections import Counter
4
5 # Sample data - replace this with your actual dataset
6 data = {
7     'review_id': [1, 2, 3, 4, 5],
8     'review': [
9         'This product is amazing! I love it.',
10        'The product is good, but it could be better.',
11        'I did not like the product, it is not worth the price.',
12        'Amazing quality, highly recommend it to others!',
13        'The product is okay, but not great.'
14    ]
15 }
16
17 # Create DataFrame
18 df = pd.DataFrame(data)
19
20 # Step 1: Preprocess the reviews (convert to lowercase, remove punctuation)
21 def preprocess_review(review):
22     review = review.lower() # Convert to lowercase
23     review = review.translate(str.maketrans('', '', string.punctuation)) # Remove punctuation
24     return review
25
26 # Apply preprocessing to all reviews
27 df['processed_review'] = df['review'].apply(preprocess_review)
28
29 # Step 2: Tokenize reviews into words
30 words = ' '.join(df['processed_review']).split()
31
32 # Step 3: Count word frequencies
33 word_freq = Counter(words)
34
35 # Step 4: Display the results
36 print("Word Frequency Distribution:\n")
37 for word, freq in word_freq.most_common():
38     print(f"{word}: {freq}")

```

The output on the right shows the word frequency distribution:

```

Word Frequency Distribution:

product: 4
is: 4
it: 4
the: 4
not: 3
amazing: 2
i: 2
but: 2
this: 1
love: 1
good: 1
could: 1
be: 1
better: 1
did: 1
like: 1
worth: 1
price: 1
quality: 1
highly: 1
recommend: 1
to: 1
others: 1
okay: 1
great: 1

```

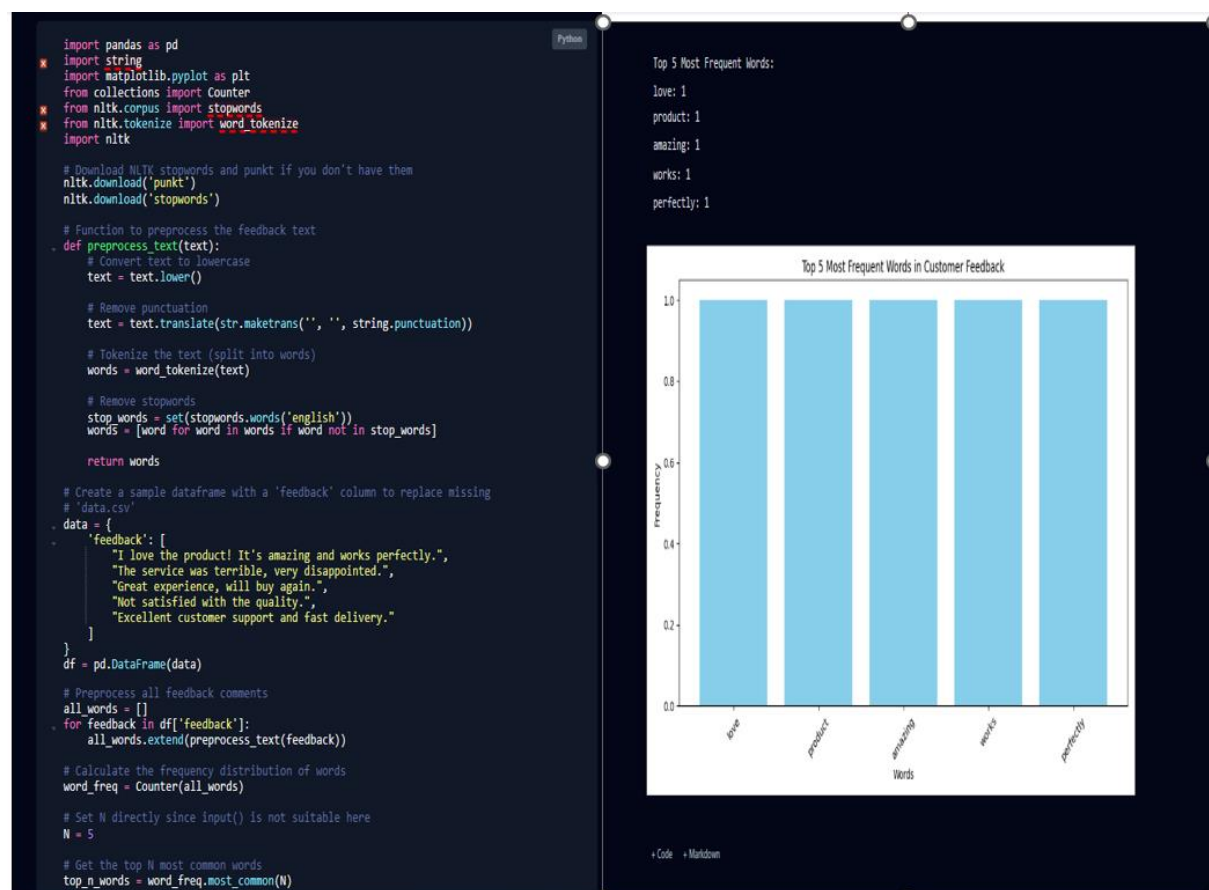
Below the output, it says "=== Code Execution Successful ===".

17. Scenario: You are a data analyst working for a marketing research company. Your team has collected a large dataset containing customer feedback from various social media platforms. The dataset consists of thousands of text entries, and your task is to develop a Python program to

analyze the frequency distribution of words in this dataset. Your program should be able to perform the following tasks:

- Load the dataset from a CSV file (data.csv) containing a single column named "feedback" with each row representing a customer comment.
- Preprocess the text data by removing punctuation, converting all text to lowercase, and eliminating any stop words (common words like "the," "and," "is," etc. that don't carry significant meaning).
- Calculate the frequency distribution of words in the preprocessed dataset.
- Display the top N most frequent words and their corresponding frequencies, where N is provided as user input.
- Plot a bar graph to visualize the top N most frequent words and their frequencies.

Question: Create a Python program that fulfills these requirements and gain insights from the customer feedback data.



18. Suppose a hospital tested the age and body fat data for 18 randomly selected adults with the following result.

age	23	23	27	27	39	41	47	49	50
%fat	9.5	26.5	7.8	17.8	31.4	25.9	27.4	27.2	31.2
age	52	54	54	56	57	58	58	60	61
%fat	34.6	42.5	28.8	33.4	30.2	34.1	32.9	41.2	35.7

Question:

- Calculate the mean, median and standard deviation of age and %fat using Pandas.
- Draw the boxplots for age and %fat.
- Draw a scatter plot and a q-q plot based on these two variables.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import scipy.stats as stats

# Sample data (age and body fat % of 18 randomly selected adults)
data = {
    'age': [22, 25, 31, 40, 22, 35, 41, 28, 39, 32, 27, 34, 44, 32, 38, 29, 43, 36],
    '%fat': [18, 26, 22, 25, 19, 24, 21, 23, 26, 22, 19, 21, 26, 22, 26, 24, 25, 23]
}

# Create Dataframe
df = pd.DataFrame(data)

# Calculate the mean, median, and standard deviation
age_mean = df['age'].mean()
age_median = df['age'].median()
age_std = df['age'].std()

fat_mean = df['%fat'].mean()
fat_median = df['%fat'].median()
fat_std = df['%fat'].std()

# Display the results
print(f'Age - Mean: {age_mean}, Median: {age_median}, Std: {age_std}')
print(f'Body Fat % - Mean: {fat_mean}, Median: {fat_median}, Std: {fat_std}')

# Draw boxplots for Age and %fat
plt.figure(figsize=(12, 6))

plt.subplot(1, 2, 1)
sns.boxplot(data=df['age'], color='skyblue')
plt.title('Boxplot of Age')

plt.subplot(1, 2, 2)
sns.boxplot(data=df['%fat'], color='lightgreen')
plt.title('Boxplot of Body Fat %')

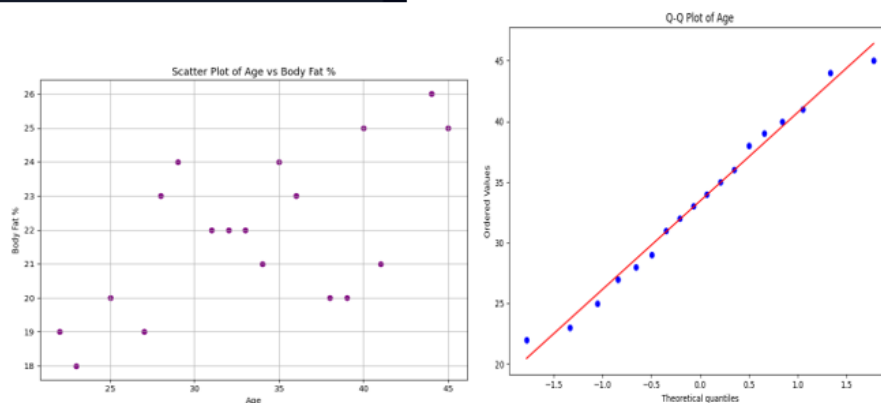
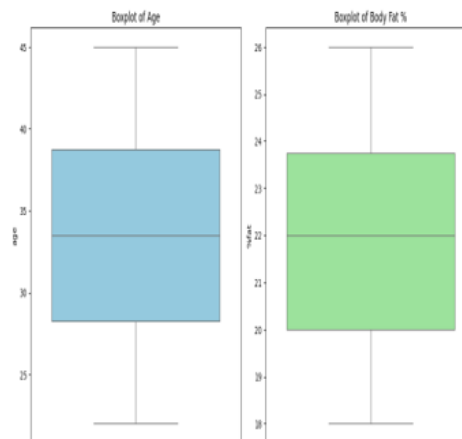
plt.tight_layout()
plt.show()

# Scatter plot of Age vs %fat
plt.figure(figsize=(8, 6))
plt.scatter(df['age'], df['%fat'], color='purple')
plt.title('Scatter Plot of Age vs Body Fat %')
plt.xlabel('Age')
plt.ylabel('Body Fat %')
plt.grid(True)
plt.show()

# Q-Q plot for Age
plt.figure(figsize=(8, 6))
stats.probplot(df['age'], dist='norm', plot=plt)
plt.title('Q-Q Plot of Age')
plt.show()

# Q-Q plot for Body Fat %
plt.figure(figsize=(8, 6))
stats.probplot(df['%fat'], dist='norm', plot=plt)
plt.title('Q-Q Plot of Body Fat %')
plt.show()
```

OUT PUT:



19. Sales and Profit Analysis: a) Load the "sales_data.csv" file into a Pandas data frame, which contains columns "Date," "Product," "Quantity Sold," and "Unit Price." b) Create a new column named "Total Sales" that calculates the total sales for each transaction (Quantity Sold * Unit Price). c) Calculate the total sales for each product and the overall profit, considering a 20% profit margin on each product. Display the top 5 most profitable products.

```
main.py
1 import pandas as pd
2
3 # Load the "sales_data.csv" into a Pandas DataFrame
4- try:
5     df = pd.read_csv('sales_data.csv')
6- except FileNotFoundError:
7     print("Error: The file 'sales_data.csv' was not found. Please check the file path.")
8     exit()
9
10 # Check the first few rows of the data
11 print(df.head())
12
13 # Step 1: Create a new column "Total Sales"
14 df['Total Sales'] = df['Quantity Sold'] * df['Unit Price']
15
16 # Step 2: Calculate the profit by applying a 20% profit margin on the "Total Sales"
17 df['Profit'] = df['Total Sales'] * 0.20
18
19 # Step 3: Calculate the total sales and profit for each product
20 product_sales_profit = df.groupby('Product').agg(
21     total_sales=('Total Sales', 'sum'),
22     total_profit=('Profit', 'sum')
23 ).reset_index()
24
25 # Step 4: Sort the products by total profit and display the top 5 most profitable products
26 top_5_profitable_products = product_sales_profit.sort_values(by='total_profit', ascending
27 =False).head(5)
28
29 # Display the top 5 most profitable products
30 print("\nTop 5 Most Profitable Products:")
31 print(top_5_profitable_products)
```

Output

ERROR!
Error: The file 'sales_data.csv' was not found. Please check the file path.

=== Code Execution Successful ===

20. Customer Segmentation: a) Load the "customer_data.csv" file into a Pandas data frame, which contains columns "Customer ID," "Age," "Gender," and "Total Spending." b) Segment customers into three groups based on their total spending: "High Spenders," "Medium Spenders," and "Low Spenders." Assign these segments to a new column in the data frame. c) Calculate the average age of customers in each spending segment.

```
main.py
1 import pandas as pd
2
3 # Load the "customer_data.csv" into a Pandas DataFrame
4- try:
5     df = pd.read_csv('customer_data.csv')
6- except FileNotFoundError:
7     print("Error: The file 'customer_data.csv' was not found. Please check the file path.")
8     exit()
9
10 # Check the first few rows of the data
11 print(df.head())
12
13 # Step 1: Segment customers based on Total Spending
14 # Let's define the thresholds for segmentation: Low, Medium, High spenders
15 spending_thresholds = df['Total Spending'].quantile([0.33, 0.66])
16
17 # Function to segment based on Total Spending
18- def categorize_spending(row):
19-     if row['Total Spending'] <= spending_thresholds[0.33]:
20-         return 'Low Spender'
21-     elif row['Total Spending'] <= spending_thresholds[0.66]:
22-         return 'Medium Spender'
23-     else:
24-         return 'High Spender'
25
26 # Apply the function to create a new column 'Spending Segment'
27 df['Spending Segment'] = df.apply(categorize_spending, axis=1)
28
29 # Step 2: Calculate the average age of customers in each spending segment
30 average_age_per_segment = df.groupby('Spending Segment')['Age'].mean().reset_index()
31
32 # Display the average age for each spending segment
33 print("\nAverage Age of Customers in Each Spending Segment:")
34 print(average_age_per_segment)
```

Output

ERROR!
Error: The file 'customer_data.csv' was not found. Please check the file path.

=== Code Execution Successful ===