

Assignment Code: DS-AG-005

Statistics Basics| Assignment

Instructions: Carefully read each question. Use Google Docs, Microsoft Word, or a similar tool to create a document where you type out each question along with its answer. Save the document as a PDF, and then upload it to the LMS. Please do not zip or archive the files before uploading them. Each question carries 20 marks.

Total Marks: 200

Question 1: What is the difference between descriptive statistics and inferential statistics? Explain with examples.

Answer:

- **Descriptive Statistics:** Summarizes and describes the main features of a dataset. It uses measures like mean, median, mode, standard deviation, and visualizations (histograms, boxplots).

Example: The average exam score of 100 students in a class.

- **Inferential Statistics:** Makes predictions or inferences about a population based on a sample of data. It uses hypothesis testing, confidence intervals, and regression.

Example: Using a sample of 100 voters to predict the outcome of an election for millions of people.

Question 2: What is sampling in statistics? Explain the differences between random and stratified sampling.

Answer:

- **Sampling:** The process of selecting a subset of individuals from a population to represent the whole.

- **Random Sampling:** Every individual has an equal chance of being selected.

Example: Picking 50 students randomly from a school of 1000.

- **Stratified Sampling:** Population is divided into subgroups (strata), and random samples are taken from each. Ensures representation from all groups.

Example: In a school, selecting students proportionally from each grade.



Question 3: Define mean, median, and mode. Explain why these measures of central tendency are important.

Answer:

- **Mean:** Arithmetic average.
- **Median:** Middle value when data is ordered.
- **Mode:** Most frequently occurring value.

Importance: They measure **central tendency**, helping summarize and compare datasets.

Question 4: Explain skewness and kurtosis. What does a positive skew imply about the data?

Answer:

- **Skewness:** Measure of asymmetry of a distribution.
 - Positive skew → tail on the right (e.g., income distribution).
 - Negative skew → tail on the left.
 - **Kurtosis:** Measure of "tailedness" (outliers).
 - High kurtosis → heavy tails.
 - Low kurtosis → light tails.
- Positive skew implies:** Most values are clustered on the left, but a few large values stretch the distribution to the right.

Question 5: Implement a Python program to compute the mean, median, and mode of a given list of numbers.

```
numbers = [12, 15, 12, 18, 19, 12, 20, 22, 19, 19, 24, 24, 24, 26, 28]
```

(Include your Python code and output in the code box below.)

Answer:

Paste your code and output inside the box below:

```
import statistics as stats

numbers = [12, 15, 12, 18, 19, 12, 20, 22, 19, 19, 24, 24, 24, 26, 28]

mean_val = stats.mean(numbers)
median_val = stats.median(numbers)
mode_val = stats.mode(numbers)

print("Mean:", mean_val)
print("Median:", median_val)
print("Mode:", mode_val)
```

Mean: 19.8
Median: 19
Mode: 12

Question 6: Compute the covariance and correlation coefficient between the following two datasets provided as lists in Python:

```
list_x = [10, 20, 30, 40, 50]  
list_y = [15, 25, 35, 45, 60]
```

(Include your Python code and output in the code box below.)

Answer:

Paste your code and output inside the box below:

```
import numpy as np  
  
list_x = [10, 20, 30, 40, 50]  
list_y = [15, 25, 35, 45, 60]  
  
cov_matrix = np.cov(list_x, list_y, bias=True)  
cov_xy = cov_matrix[0][1]  
corr_xy = np.corrcoef(list_x, list_y)[0][1]  
  
print("Covariance:", cov_xy)  
print("Correlation:", corr_xy)
```

Covariance: 200.0
Correlation: 0.989

Question 7: Write a Python script to draw a boxplot for the following numeric list and identify its outliers. Explain the result:

```
data = [12, 14, 14, 15, 18, 19, 19, 21, 22, 22, 23, 23, 24, 26, 29, 35]
```

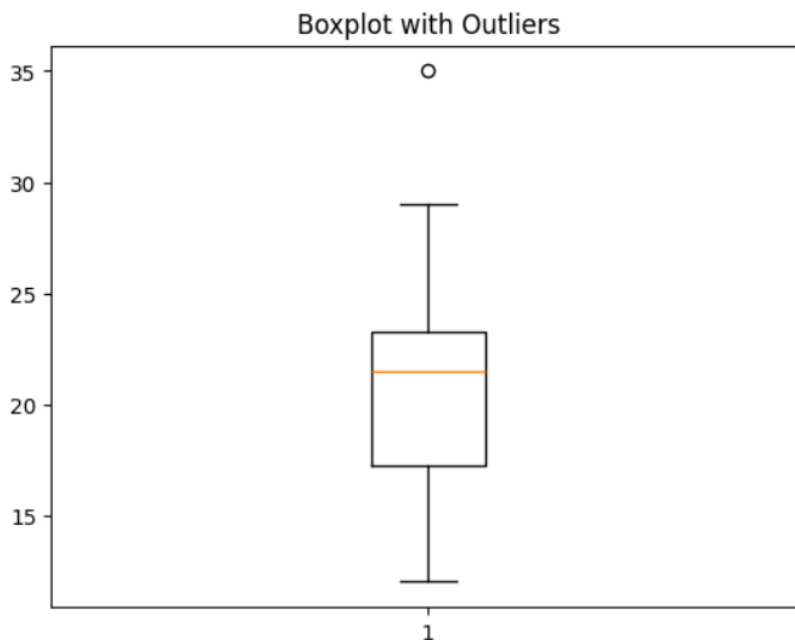
(Include your Python code and output in the code box below.)

Answer:

```
import matplotlib.pyplot as plt
```

```
data = [12, 14, 14, 15, 18, 19, 19, 21, 22, 22, 23, 23, 24, 26, 29, 35]
```

```
plt.boxplot(data)  
plt.title("Boxplot with Outliers")  
plt.show()
```



Question 8: You are working as a data analyst in an e-commerce company. The marketing team wants to know if there is a relationship between advertising spend and daily sales.

- Explain how you would use covariance and correlation to explore this relationship.
- Write Python code to compute the correlation between the two lists:

```
advertising_spend = [200, 250, 300, 400, 500]
```

```
daily_sales = [2200, 2450, 2750, 3200, 4000]
```

(Include your Python code and output in the code box below.)

Answer:

To explore the relationship between **advertising spend** and **daily sales**, we use:

- **Covariance:** Measures whether two variables move together.
 - Positive covariance → as advertising spend increases, sales also increase.
 - Negative covariance → as advertising spend increases, sales decrease.
 - But covariance does not indicate the strength of the relationship.
- **Correlation:** Standardizes covariance to a value between **-1 and +1**.
 - **+1** → perfect positive relationship
 - **0** → no relationship
 - **-1** → perfect negative relationship
 - Correlation is better for understanding **strength and direction**.

In this case, a **high positive correlation** would mean that more advertising spend leads to more sales.

```
import numpy as np

advertising_spend = [200, 250, 300, 400, 500]
daily_sales = [2200, 2450, 2750, 3200, 4000]

# Compute covariance
cov_matrix = np.cov(advertising_spend, daily_sales, bias=True)
cov_xy = cov_matrix[0][1]

# Compute correlation
corr_xy = np.corrcoef(advertising_spend, daily_sales)[0][1]

print("Covariance:", cov_xy)
print("Correlation:", corr_xy)
```

```
Covariance: 87500.0
Correlation: 0.9966158955271536
```

Question 9: Your team has collected customer satisfaction survey data on a scale of 1-10 and wants to understand its distribution before launching a new product.

- Explain which summary statistics and visualizations (e.g. mean, standard deviation, histogram) you'd use.
- Write Python code to create a histogram using Matplotlib for the survey data:

```
survey_scores = [7, 8, 5, 9, 6, 7, 8, 9, 10, 4, 7, 6, 9, 8, 7]
```

(Include your Python code and output in the code box below.)

Answer:

To understand the distribution of customer satisfaction survey scores (1–10 scale), we should use:

- **Summary Statistics:**
 - **Mean** → average satisfaction level.
 - **Median** → middle score, less affected by outliers.
 - **Mode** → most common score.
 - **Standard Deviation (SD)** → spread of scores; higher SD means more variation in customer opinions.
- **Visualizations:**
 - **Histogram** → shows frequency distribution of scores (how many customers gave each rating).
 - **Boxplot** (optional) → shows spread and potential outliers.

These help determine if customers are generally satisfied (scores clustered at high values) or opinions are mixed.

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

survey_scores = [7, 8, 5, 9, 6, 7, 8, 9, 10, 4, 7, 6, 9, 8, 7]

# Summary statistics
mean_val = np.mean(survey_scores)
median_val = np.median(survey_scores)
mode_val = stats.mode(survey_scores)
std_dev = np.std(survey_scores)

print("Mean:", mean_val)
print("Median:", median_val)
print("Mode:", mode_val)
print("Standard Deviation:", std_dev)

# Histogram
plt.hist(survey_scores, bins=6, edgecolor="black")
```

```
Mean: 7.333333333333333
Median: 7.0
Mode: 7
Standard Deviation: 1.577621275493231
```

