

Statistical Concepts: Q&A

Arnab Aich

Summary Statistics

Question 1

Question: What is the difference between the mean and median, and when would you use one over the other?

Question 1

Question: What is the difference between the mean and median, and when would you use one over the other?

Answer:

- The **mean** is the average of all data points.
- The **median** is the middle value when data are ordered.
- Use the **median** when the data is skewed or contains outliers, and the **mean** when the data is symmetrically distributed.

Question 2

Question: How do variance and standard deviation describe the spread of data?

Question 2

Question: How do variance and standard deviation describe the spread of data?

Answer:

- **Variance** measures how far data points are from the mean.
- **Standard deviation** is the square root of the variance and is in the same units as the data.

Z-Score

Question 1

Question: What does a z-score represent, and how is it calculated?

Question 1

Question: What does a z-score represent, and how is it calculated?

Answer:

A **z-score** shows how many standard deviations a data point is from the mean. It is calculated as:

$$Z = \frac{(X - \mu)}{\sigma}$$

Question 2

Question: Why are z-scores useful for comparing data across different datasets?

Question 2

Question: Why are z-scores useful for comparing data across different datasets?

Answer:

Z-scores standardize data, allowing you to compare values across datasets with different means and standard deviations.

Range

Question 1

Question: What is the range, and what does it tell you about a dataset?

Question 1

Question: What is the range, and what does it tell you about a dataset?

Answer:

The **range** is the difference between the maximum and minimum values in a dataset. It gives a simple measure of the spread, but it is sensitive to outliers.

Question 2

Question: When might the range be a misleading measure of spread?

Question 2

Question: When might the range be a misleading measure of spread?

Answer:

The range can be misleading when there are outliers, as it only considers the extreme values. A dataset with outliers will have a much larger range even if most data points are close together.

Quantiles and IQR

Question 1

Question: What are quantiles, and how do you interpret them?

Question 1

Question: What are quantiles, and how do you interpret them?

Answer:

Quantiles divide data into equal parts. Common quantiles include quartiles and percentiles. For example: - The 25th percentile (Q1) represents the value below which 25% of the data falls. - The 75th percentile (Q3) represents the value below which 75% of the data falls.

Question 2

Question: What is the interquartile range (IQR), and why is it a better measure of spread than the range?

Question 2

Question: What is the interquartile range (IQR), and why is it a better measure of spread than the range?

Answer:

The **IQR** is the range between the first quartile (Q_1) and the third quartile (Q_3). It measures the spread of the middle 50% of the data and is not affected by outliers, making it a better measure of spread than the range.

Histograms

Question 1

Question: What is the purpose of using bins in a histogram, and how does changing the number of bins affect the histogram?

Question 1

Question: What is the purpose of using bins in a histogram, and how does changing the number of bins affect the histogram?

Answer:

Bins group data into intervals, allowing visualization of the data's distribution. More bins provide more detail, but too many bins may make the histogram noisy. Too few bins can oversimplify the data.

Question 2

Question: How does a histogram differ from a bar plot?

Question 2

Question: How does a histogram differ from a bar plot?

Answer:

- **Histograms** are used for continuous numerical data. - **Bar plots** are used for categorical data.

Box Plots

Question 1

Question: What do the whiskers in a box plot represent, and how can you use a box plot to identify outliers?

Question 1

Question: What do the whiskers in a box plot represent, and how can you use a box plot to identify outliers?

Answer:

Whiskers show the spread of the data up to 1.5 times the interquartile range (IQR). Points outside this range are considered outliers and are shown individually.

Question 2

Question: How does a box plot show the skewness of a dataset?

Question 2

Question: How does a box plot show the skewness of a dataset?

Answer:

If the median is closer to one quartile, and the whiskers are of unequal length, the data is skewed in the direction of the longer whisker.

Pie Charts

Question 1

Question: When is it appropriate to use a pie chart over other types of plots?

Question 1

Question: When is it appropriate to use a pie chart over other types of plots?

Answer:

A pie chart is best used to show proportions of categories within a whole when there are only a few categories. It becomes less useful as the number of categories increases.

Question 2

Question: What are the limitations of using pie charts for data visualization?

Question 2

Question: What are the limitations of using pie charts for data visualization?

Answer:

Pie charts can be difficult to interpret when there are too many categories or when the differences between slices are small. Bar plots often make comparisons easier.