**1. What are the measures of central tendency, and how do they differ from each other?**

Measures of central tendency are statistical metrics used to identify the center point or typical value of a dataset. The main measures of central tendency are the mean, median, and mode. Each of these measures offers different insights into the data, and they can differ significantly depending on the nature and distribution of the data.

**Mean**

The mean, often referred to as the average, is calculated by summing all the values in a dataset and then dividing by the number of values.

**Median**

The median is the middle value of a dataset when it is ordered in ascending or descending order. If the dataset has an even number of values, the median is the average of the two middle numbers.

**Mode**

The mode is the value that occurs most frequently in a dataset. A dataset may have one mode, more than one mode, or no mode at all if no number repeats.

**Differences between Mean, Median, and Mode**

1. **Sensitivity to Outliers:**

   ✓ **Mean:** Highly sensitive to outliers. Extreme values can significantly skew the mean.
   ✓ **Median:** Not sensitive to outliers. It only depends on the middle values, making it a better measure when there are outliers.
   ✓ **Mode:** Not sensitive to outliers. It focuses on the most frequent value(s).

2. **Usefulness:**

   ✓ **Mean:** Useful when you need to consider every value in the dataset. It provides a measure of central tendency that includes all values.
   ✓ **Median:** Useful when you need to find the central value without the influence of outliers. It provides a more accurate reflection of a typical value in skewed distributions.
   ✓ **Mode:** Useful in categorical data to find the most common category. It can also indicate the most frequent value in numerical data.

3. **Data Type:**

   ✓ **Mean:** Requires interval or ratio data.
   ✓ **Median:** Can be used with ordinal, interval, or ratio data.
   ✓ **Mode:** Can be used with nominal, ordinal, interval, or ratio data.

**2. How do you interpret the standard deviation in the context of data variability?**

Standard deviation is a measure of the amount of variation or dispersion in a set of values. It quantifies how much the values in a dataset deviate from the mean (average) of the dataset. A low standard deviation indicates that the values are close to the mean, while a high standard deviation indicates that the values are spread out over a wider range.

**Interpretation of Standard Deviation**

1. **Low Standard Deviation:**

   ✓ Values are close to the mean.
   ✓ Indicates low variability within the dataset.
   ✓ Example: In a quality control scenario, if the standard deviation of the diameter of manufactured screws is low, it means most screws are close to the target diameter, indicating consistent manufacturing.

2. **High Standard Deviation:**

   ✓ Values are spread out from the mean.
   ✓ Indicates high variability within the dataset.
   ✓ Example: In test scores, a high standard deviation indicates a wide range of scores, suggesting that some students performed very well while others performed poorly.

## 3. What is a box plot, and what information can you extract from it?

A box plot, also known as a box-and-whisker plot, is a graphical representation of the distribution of a dataset that shows its central tendency, variability, and skewness. It provides a visual summary of the data through five key summary statistics: the minimum, first quartile (Q1), median (Q2), third quartile (Q3), and maximum.

**Information Extracted from a Box Plot**

1. **Central Tendency:** The median line inside the box shows the median value.
2. **Spread and Variability:** The length of the box (IQR) shows the variability of the middle 50% of the data.
3. **Skewness:** The position of the median within the box indicates skewness. If the median is closer to Q1, the data is right-skewed; if closer to Q3, the data is left-skewed.
4. **Range:** The whiskers show the range of the data, excluding outliers.
5. **Outliers:** Individual points plotted outside the whiskers represent outliers.

## 4. Explain the significance of the interquartile range (IQR) and how it is used to detect outliers.

The Interquartile Range (IQR) is a measure of statistical dispersion, which represents the range within which the middle 50% of the data values lie. It is calculated as the difference between the third quartile (Q3) and the first quartile (Q1).

**Significance of the IQR**

1. **Measure of Spread:** The IQR measures the spread of the central portion of the data. Since it excludes the extreme values (outliers), it provides a more robust measure of spread compared to the range.
2. **Identification of Outliers:** The IQR is used to detect outliers in the data. Outliers are data points that fall significantly outside the typical range of values.

**Example:**
Consider the dataset of house prices: 150, 200, 210, 220, 250, 260, 270, 300, 310, 1500

1. **Order the Data:** 150, 200, 210, 220, 250, 260, 270, 300, 310, 1500

2. **Calculate Q1 and Q3:**

    ✓ Q1 (25th percentile) = **210**
    ✓ Q3 (75th percentile) = **300**

3. **Calculate the IQR:**

    ✓ IQR=Q3−Q1=300−210=**90**

4. **Determine the Bounds:**

    ✓ Lower Bound = Q1−1.5×IQR=210−1.5×90=210−135=**75**
    ✓ Upper Bound = Q3+1.5×IQR=300+1.5×90=300+135=**435**

5. **Identify Outliers:**

    ✓ Any data points below 75 or above 435 are outliers.
    ✓ The value 1500 is above 435, so it is an outlier.

## 5. How Do Maximum Likelihood Estimators (MLE) Work

Maximum Likelihood Estimation (MLE) is a method used in statistics to estimate the parameters of a statistical model. The idea behind MLE is to find the parameter values that maximize the likelihood function, which measures how well the model explains the observed data.

1. **Define the Likelihood Function:** The likelihood function, $L(\theta)$, is defined as the probability of observing the given sample data as a function of the parameters of the model, $\theta$.
2. **Construct the Log-Likelihood:** Since the likelihood function often involves products of probabilities, it can be more convenient to work with the log-likelihood, $\ell(\theta)$, which is the natural logarithm of the likelihood function: $\ell(\theta)=\ln L(\theta)$. The log-likelihood is easier to differentiate and work with due to its additive properties.
3. **Maximize the Log-Likelihood:** Find the parameter values that maximize the log-likelihood function. This typically involves taking the derivative of the log-likelihood function with respect to the parameters, setting the derivatives to zero, and solving for the parameters.
4. **Solve for the Parameters:** The values of the parameters that maximize the log-likelihood function are the maximum likelihood estimates.