

Basic Statistics:

Mean

What it is: The mean, also known as the average, is simply the sum of all the values in the data set divided by the number of values.

How to find it:

1. Add all the values in your data set.
2. Divide the sum by the total number of values.

Example: Let's say you have the following test scores: 78, 85, 92, 80, 75.

1. Add them up: $78 + 85 + 92 + 80 + 75 = 410$.
2. Divide by the number of scores (5): $410 / 5 = 82$.
3. The mean score is 82.

Median

What it is: The median is the "middle" value when the data is ordered from least to greatest.

How to find it:

1. Arrange the data in ascending or descending order.
2. If you have an odd number of values, the median is the middle value.
3. If you have an even number of values, the median is the average of the two middle values.

Example: Use the same test scores: 75, 78, 80, 85, 92 (already ordered).

1. Since we have an odd number (5), the middle value is the median.
2. The middle value is 80.
3. The median score is 80.

Mode

What it is: The mode is the most frequent value in the data set.

How to find it:

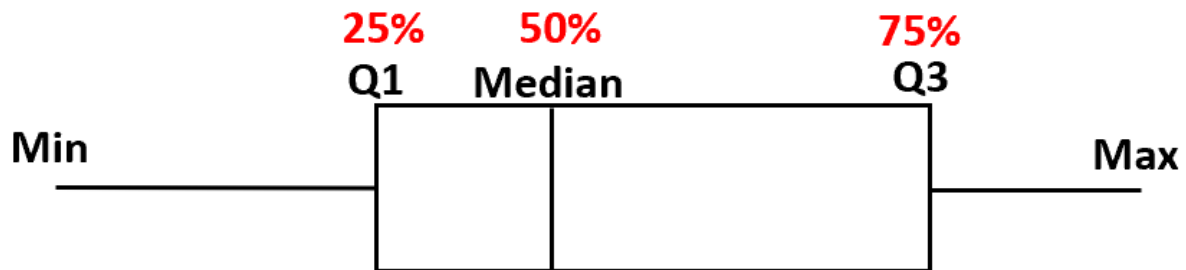
Examine the data and see which value appears most often.

Example: Use the same test scores: 75, 78, 80, 85, 92.

No value appears more than once. In this case, there is no mode.

Important points to remember:

- The mean is sensitive to outliers (extreme values) in the data set. The median is generally less affected by outliers.
- The mode may not exist if no value repeats, or there may be multiple modes if there are two or more values that tie for the most frequent.



A box and whisker plot, also called a box plot, is a way to visually represent the distribution of a set of data using quartiles. Here's how it works with Q1, Q2, and Q3:

- **Quartiles:** These are specific points that divide the data into four equal quarters.
 - Q1 (First Quartile): Represents the value where 25% of the data falls below it and 75% falls above.
 - Q2 (Second Quartile): This is the median, the middle value when the data is ordered from least to greatest. Here, 50% of the data falls below and 50% falls above.
 - Q3 (Third Quartile): Represents the value where 75% of the data falls below it and 25% falls above.
- **The Box:** The box in the plot depicts the interquartile range (IQR), which is the spread of the middle 50% of the data. The box stretches from Q1 (bottom) to Q3 (top).
- **Whiskers:** These lines extend from the box outwards. They typically reach the farthest data points within 1.5 times the IQR from Q1 and Q3. Data points beyond these whiskers are considered potential outliers.

Key points to remember:

- The median (Q2) is visualized by a line segment inside the box, dividing it in half (sometimes). If the distribution is symmetrical, the median will be exactly in the center of the box.
- The box plot allows you to see:
 - The center of the data (median)
 - The spread of the middle 50% (IQR)

- The presence of potential outliers (data points far from the box)

Comparing Data:

Bivariate data refers to information collected on two variables for each element or observation. Each element essentially has a paired value, allowing you to investigate how these values change in relation to each other. Imagine a dataset with two columns, representing the two variables you're interested in. Each row represents an element with corresponding values for both variables.

Understanding Relationships

The key aspect of bivariate data is exploring the relationship between the two variables. This can involve:

- **Direction:** Do the values of one variable tend to increase or decrease as the values of the other variable change?
- **Strength:** How strong is the observed change? Is it a subtle shift or a dramatic difference?
- **Linearity:** Does the change follow a straight-line pattern, or is it more complex?

Comparison Techniques

Several techniques help analyze and compare variables in bivariate data:

- **Scatter Plots:** This visual representation plots each element as a point on a graph, with one variable on each axis. By observing the distribution of these points, you can identify patterns and potential relationships between the variables.
- **Correlation Coefficient:** This statistical measure quantifies the strength and direction of the linear relationship between the two variables. Values closer to +1 indicate a strong positive correlation (both variables tend to increase together), while values closer to -1 indicate a strong negative correlation (as one increases, the other decreases). A value close to 0 suggests little to no linear relationship.

Important to Consider

- **Correlation vs. Causation:** A correlation between variables doesn't necessarily imply causation (one causing the other). There might be a third, unseen factor influencing both variables.
- **Data Type:** The chosen analysis techniques may depend on the type of data your variables represent (numerical, categorical, etc.).

Cumulative Frequency Explained with Graph and Table

Cumulative frequency refers to the total number of observations in a data set that fall at or below a specific value. It's a way to understand how the data accumulates as you move through the range of values.

Visualizing Cumulative Frequency with a Graph

The graph you see above is a cumulative frequency graph. Let's break down its elements:

- **Horizontal Axis (X-axis):** This axis represents the different values (or intervals) in your data set. In this example, it might be test scores ranging from 20 to 100.
- **Vertical Axis (Y-axis):** This axis shows the cumulative frequency. It starts at 0 and increases as you move up the graph.
- **Line:** The line on the graph connects the data points, representing the cumulative frequency for each value (or interval).

How to Read the Graph:

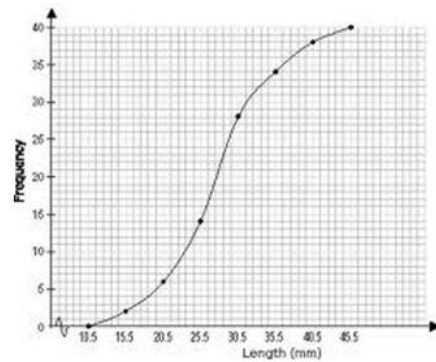
Imagine you want to find the cumulative frequency for a value of 70 on the X-axis. Trace a vertical line from 70 on the X-axis until it meets the line graph. Then, follow a horizontal line from that point on the graph to the Y-axis. The value you read on the Y-axis is the cumulative frequency. In this example, let's say it's 22. This means there are 22 or fewer observations in the data set that have a score of 70 or below.

Cumulative frequency

Cumulative frequency table

Class Limits	Frequency	Cumulative Frequency
5-10	1	1
10-15	2	3
15-20	4	7
20-25	0	7
25-30	3	10
30-35	5	15
35-40	6	21

Cumulative frequency graph



dr

- Frequency: This column shows how many observations fall within each value range (interval).
- Cumulative Frequency: This column shows the total number of observations that fall at or below the upper limit of each interval. It's calculated by adding the frequency of the current interval to the cumulative frequency of the previous interval(s).