



**Sravya Madipalli** ✓

# Basic Statistics for Data Science

Statistics is crucial for data science interviews. Whether you're analyzing trends, building machine learning models, or visualizing insights, a solid understanding of statistics is essential. It showcases your analytical mindset and problem-solving abilities.

Today, let's dive into below topics...

- Basic statistics concepts.
- Real-world applications to prepare for interviews.
- Interview questions, and ChatGPT prompts to deepen your understanding.

# Descriptive Statistics

## Measures of Central Tendency

- **Mean:** The arithmetic average of a dataset, commonly used for balanced datasets. It's a basic measure to summarize a dataset.
  - **Why it's important:** The mean gives an idea of the average value and is essential for data analysis tasks like calculating average sales or customer lifetime value.
  - **Application:** Useful in summarizing financial data, such as the average revenue per customer.
- **Median:** The middle value in an ordered dataset, less sensitive to outliers.
  - **Why it's important:** The median is critical when the dataset is skewed and the mean might give a misleading picture.
  - **Application:** Often used in real estate or salary distributions to provide a more accurate 'center' in skewed data.
- **Mode:** The most frequent value in a dataset, important for categorical data.
  - **Why it's important:** Helps identify the most common category or event in data.
  - **Application:** Used in customer preference analysis (e.g., most bought product).



# Descriptive Statistics

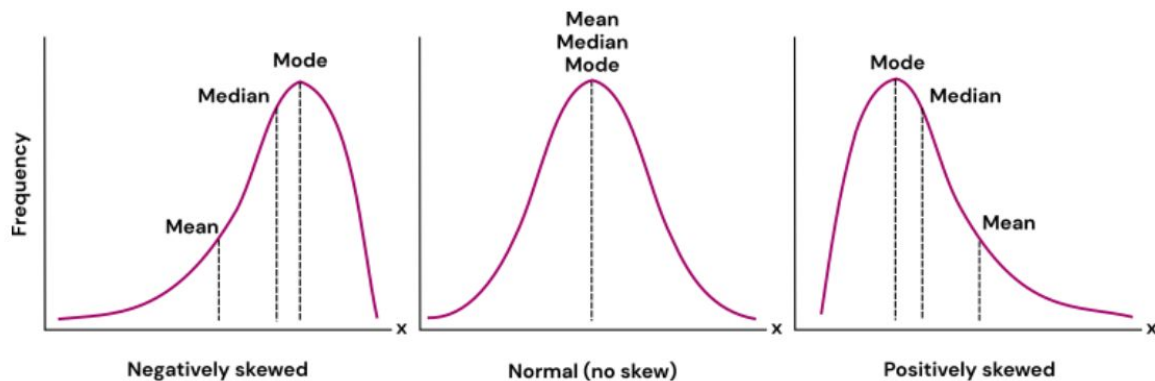
## Skewness and the Relationship Between Mean, Median, and Mode

- **Positive Skewness:** In a positively skewed distribution (right-skewed), the tail is longer on the right side. In such cases, the mean is usually greater than the median, and the mode is less than the median.
  - **Order:**  $\text{Mode} < \text{Median} < \text{Mean}$ .
  - **Why it's important:** If you observe this pattern, it tells you that the data is right-skewed, meaning there are some higher outliers pulling the mean up.
  - **Application:** This often happens in income distributions, where a few individuals have very high incomes compared to the rest of the population.
- **Negative Skewness:** In a negatively skewed distribution (left-skewed), the tail is longer on the left side. In such cases, the mean is usually less than the median, and the mode is greater than the median.
  - **Order:**  $\text{Mean} < \text{Median} < \text{Mode}$ .
  - **Why it's important:** A negative skew indicates that there are lower outliers pulling the mean down. This could indicate that most of the data is clustered toward the higher end.
  - **Application:** An example might be retirement age, where a few early retirees bring down the mean, while most people retire at older ages.

# Descriptive Statistics

## Skewness and the Relationship Between Mean, Median, and Mode

- **Symmetric Distribution:** In a perfectly symmetrical distribution (normal distribution), the mean, median, and mode are all equal.
  - **Order:** Mean = Median = Mode.
  - **Why it's important:** A symmetrical distribution indicates that the data is evenly distributed, which is often assumed in statistical modeling.



# Descriptive Statistics

## Measures of Dispersion

- **Range:** The difference between the highest and lowest values in a dataset.
  - **Why it's important:** Provides a quick sense of how spread out the data is.
  - **Application:** Common in reporting temperature variations or stock price ranges.
- **Variance:** Measures how much the data points differ from the mean.
  - **Why it's important:** Variance helps assess the variability in data, which is crucial in fields like finance and machine learning.
  - **Application:** Used in risk assessment to understand fluctuations in stock prices.
- **Standard Deviation:** The square root of the variance, gives spread in original units.
  - **Why it's important:** Easier to interpret than variance and widely used to gauge data consistency.
  - **Application:** Used to measure market volatility in finance or customer behavior in A/B tests.



# Descriptive Statistics

## Interview Questions for Descriptive Statistics

1. When would you use median instead of mean to summarize data?
2. How does standard deviation help in evaluating data consistency?
3. What are the limitations of using the range as a measure of dispersion?
4. Why might mode be useful in analyzing categorical data, and can there be more than one mode?
5. What is the relationship between mean, median, and mode in positively skewed and negatively skewed data, and why does this matter?

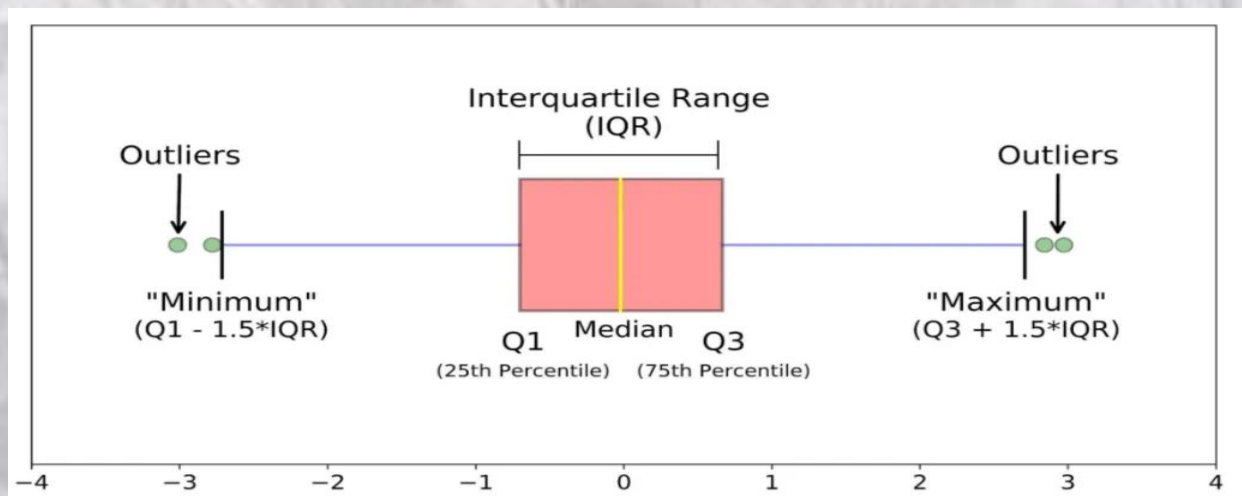
## ChatGPT Prompt for Descriptive Statistics

*"Explain how measures of central tendency (mean, median, mode) and measures of dispersion (range, variance, standard deviation) work together to give you a complete picture of a dataset. Provide an example dataset and show how each measure helps in understanding the data."*

# Quartiles, Interquartile Range (IQR), and Percentiles

## Quartiles and Interquartile Range (IQR)

- **Quartiles:** Divide the data into four equal parts. The first quartile (Q1) is the 25th percentile, and the third quartile (Q3) is the 75th percentile.
  - **Why it's important:** Provides insight into the distribution of the dataset, especially in skewed data.
  - **Application:** Used in box plots to identify the spread of the middle 50% of data.
- **IQR:** The difference between Q1 and Q3, helps identify outliers.
  - **Why it's important:** Provides a robust measure of spread. IQR is preferred when outliers are present, as it ignores the extremes in the dataset.
  - **Application:** IQR is widely used in outlier detection and in visualizing data using box plots.





# Quartiles, Interquartile Range (IQR), and Percentiles

## Percentiles

- **Percentiles:** Show the relative position of a data point within a dataset. For example, the 90th percentile means 90% of data points are below this value.
  - **Why it's important:** Helps compare individual data points to the rest of the dataset. Percentiles are commonly used in standardized testing.
  - **Application:** Used to analyze distributions in test scores, salary ranges, or population studies to understand relative standings.

# Quartiles, Interquartile Range (IQR), and Percentiles

## Interview Questions for Quartiles, IQR, and Percentiles

1. How can the IQR help in identifying outliers?
2. Why are percentiles useful when comparing individual performances or outcomes?
3. How does the interquartile range differ from the range as a measure of data spread?

## ChatGPT Prompt for Quartiles, IQR, and Percentiles

*"Can you explain how quartiles and percentiles divide a dataset? Give an example of how quartiles and IQR are used in outlier detection, and percentiles are used in reporting test scores."*

# Data Types and Scales of Measurement

## Nominal, Ordinal, Interval, and Ratio

- **Nominal:** Data that categorizes without any order (e.g., colors, gender).
  - **Why it's important:** Nominal data is essential in categorical analysis where order does not matter.
  - **Application:** Commonly used in customer segmentation and demographics.
- **Ordinal:** Data with a meaningful order but no consistent interval between values (e.g., rankings).
  - **Why it's important:** Useful in surveys where items are rated or ranked.
  - **Application:** Used in satisfaction surveys, where respondents rank services or products.
- **Interval:** Data with meaningful differences between values but no true zero point (e.g., temperature).
  - **Why it's important:** Interval data allows for the measurement of differences but lacks the concept of an absolute zero.
  - **Application:** Used in education to compare test scores, where zero does not indicate a complete absence of knowledge.



# Data Types and Scales of Measurement

## Nominal, Ordinal, Interval, and Ratio

- **Ratio:** Data with meaningful differences and a true zero (e.g., height, weight).
  - **Why it's important:** Ratio data allows for the full range of mathematical operations, making it highly versatile in data analysis.
  - **Application:** Common in finance, where metrics like revenue or expenses are analyzed.

# Data Types and Scales of Measurement

## Interview Questions for Data Types and Scales of Measurement

1. What is the difference between ordinal and nominal data?
2. When would you use ratio data over interval data in an analysis?
3. Why is it important to distinguish between different scales of measurement?

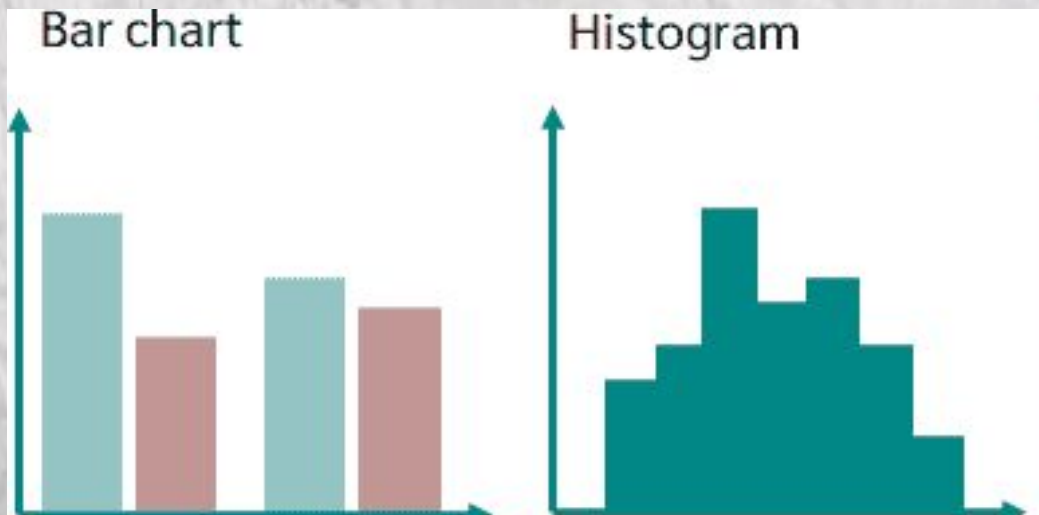
## ChatGPT Prompt for Data Types and Scales of Measurement

*"Explain the four scales of measurement (nominal, ordinal, interval, and ratio) with examples. How do these scales determine the type of statistical analysis you can perform?"*

# Data Visualization

## Histograms, Bar Charts, Box Plots, and Scatter Plots

- **Histograms:** Visualize the distribution of continuous data by grouping values into bins.
  - **Why it's important:** Helps in understanding the frequency distribution of data points.
  - **Application:** Common in displaying exam score distributions or sales volumes.
- **Bar Charts:** Compare different categories by showing the frequency or magnitude of each category.
  - **Why it's important:** Ideal for categorical data and comparing different groups.
  - **Application:** Used in comparing sales across regions or customer segments.

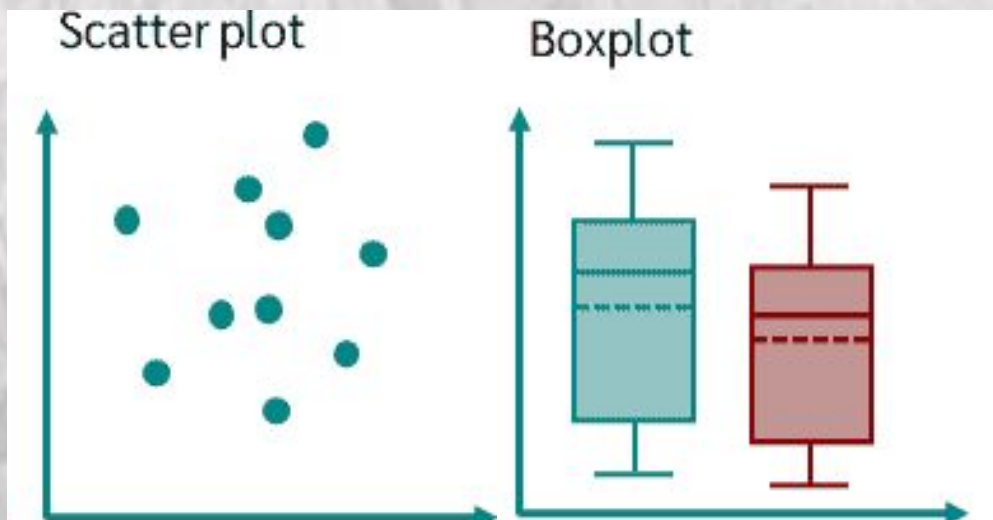




# Data Visualization

## Histograms, Bar Charts, Box Plots, and Scatter Plots

- **Box Plots:** Summarize data distribution through five-number summaries (minimum, Q1, median, Q3, maximum) and identify outliers.
  - **Why it's important:** Provides a clear view of data spread and highlights outliers.
  - **Application:** Frequently used in financial data analysis to compare stock price distributions.
- **Scatter Plots:** Show the relationship between two continuous variables, often used to identify correlations.
  - **Why it's important:** Critical in regression analysis and identifying trends or relationships.
  - **Application:** Common in exploring the relationship between advertising spend and sales revenue.



# Data Visualization

## Interview Questions for Data Visualization

1. When is a histogram more appropriate than a bar chart?
2. How can you use a scatter plot to identify relationships between two variables?
3. What insights can you gain from a box plot that other charts may not show?

## ChatGPT Prompt for Data Visualization

*"Explain how different types of data visualizations (histograms, bar charts, box plots, and scatter plots) are used to communicate insights. Provide an example dataset and show how each visualization reveals different aspects of the data."*

# Data Collection and Sampling

## Random, Stratified, and Systematic Sampling

- **Random Sampling:** Ensures every individual has an equal chance of being selected.
  - **Why it's important:** Reduces bias and increases the representativeness of the sample.
  - **Application:** Used in surveys and experiments to gather unbiased data.
- **Stratified Sampling:** Divides the population into subgroups (strata) and selects samples from each.
  - **Why it's important:** Ensures representation of all key subgroups in the data.
  - **Application:** Commonly used in political polling to ensure representation from different demographic groups.
- **Systematic Sampling:** Selects every  $n$ th individual from a population.
  - **Why it's important:** Simple and efficient for large populations when a complete list is available.
  - **Application:** Often used in quality control, where every  $n$ th product is inspected.



# Data Distributions

## Normal Distribution, Skewness, and Kurtosis

- **Normal Distribution:** A bell-shaped curve where most data points are centered around the mean.
  - **Why it's important:** Many statistical methods assume normal distribution. It's a baseline for hypothesis testing and regression analysis.
  - **Application:** Common in analyzing heights, IQ scores, or exam scores where data tends to follow a normal pattern.
- **Skewness:** Measures the asymmetry of a data distribution.
  - **Why it's important:** Skewness helps understand whether the data is biased towards higher or lower values.
  - **Application:** Useful in finance, where skewed data can indicate risk or unusual market behavior.

# Data Distributions

## Normal Distribution, Skewness, and Kurtosis

- **Kurtosis:** Describes the "tailedness" of a distribution.
  - **Why it's important:** Understanding the kurtosis of a dataset helps in identifying the presence of outliers and the general shape of the data distribution. This insight is critical for fields like finance and risk management, where extreme values can have significant impacts. Knowing whether a dataset is leptokurtic or platykurtic allows you to anticipate whether extreme events (outliers) are likely or rare, helping to adjust strategies accordingly.
  - **Application:** In investment analysis, a leptokurtic distribution might indicate higher risk because of the frequent occurrence of extreme returns. In contrast, a platykurtic distribution could suggest a more stable and predictable investment with fewer extreme outcomes.

# Data Distributions

## Types of Kurtosis:

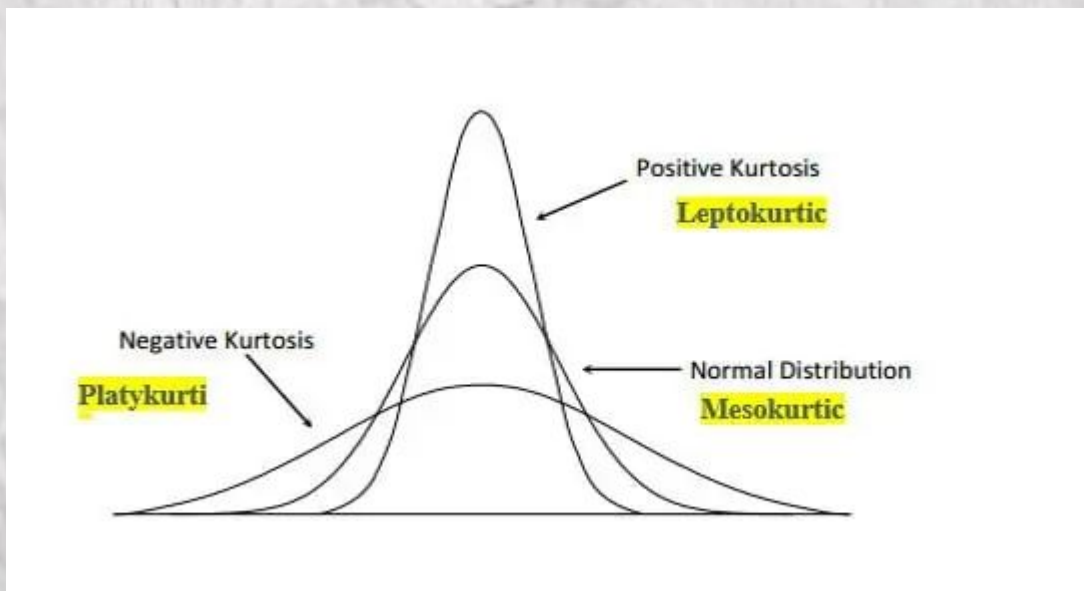
- **Leptokurtic:** A distribution with positive kurtosis (greater than 3) is called **leptokurtic**. These distributions have heavier tails and a sharper peak around the mean, indicating more outliers.
  - **Why it's important:** Leptokurtic distributions signal that extreme values (outliers) are more frequent than in a normal distribution. This can impact models where outliers skew results.
  - **Application:** Often seen in financial markets, where extreme changes in stock prices or returns are more common than a normal distribution would suggest. A leptokurtic distribution warns of potential market volatility.
- **Platykurtic:** A distribution with negative kurtosis (less than 3) is called **platykurtic**. These distributions have lighter tails and a flatter peak around the mean, indicating fewer outliers and less extreme variability.
  - **Why it's important:** Platykurtic distributions suggest that data points are more evenly distributed and extreme outliers are less likely.
  - **Application:** Used in quality control processes where variability needs to be minimized, such as in manufacturing. A platykurtic distribution indicates that most products meet the standard quality, with fewer extreme deviations.



# Data Distributions

## Types of Kurtosis:

- **Mesokurtic:** A distribution with kurtosis close to 3 (the same as a normal distribution) is called **mesokurtic**. These distributions have moderate tails and follow the pattern of a normal distribution.
  - **Why it's important:** Mesokurtic distributions are often assumed in many statistical models (like regression or hypothesis testing) and serve as a baseline for comparison.
  - **Application:** Common in fields like social sciences, where normal distribution assumptions often hold for variables like test scores or IQ.



# Data Distributions

## Interview Questions for Data Distributions

1. Why is normal distribution important in statistical analysis?
2. What does skewness tell you about the nature of your data?
3. How does kurtosis help in identifying outliers in a dataset?
4. What does kurtosis tell you about the distribution of your data?
5. How does leptokurtic distribution differ from platykurtic, and in what scenarios would you expect to find each?
6. Why is kurtosis important in understanding the risks of outliers in data?

## ChatGPT Prompt for Data Distributions

*"Describe the importance of understanding data distributions, including normal distribution, skewness, and kurtosis. Provide an example dataset and explain how each of these concepts applies to analyzing the data."*



**Sravya Madipalli** ✓

## Was this Helpful?



Save it



Follow Me



Repost and Share it  
with your friends