**30 most commonly asked statistics questions for Data Analyst/Business Analyst interviews:**

**Basic Level**

1. **What is the difference between mean, median, and mode?**

**Answer** - Mean is the average of a dataset, calculated by adding all the numbers and dividing by the count.

Median is the middle value when the data is sorted in ascending or descending order.

Mode is the value that appears most frequently in a dataset.

2. **What are variance and standard deviation, and why are they important?**

**Answer**: Variance measures the spread of data points from the mean. A higher variance indicates greater spread.

Standard deviation is the square root of variance and gives the spread in the same units as the data.

Both are important because they quantify how spread out or clustered the data is, helping to understand data consistency and risk.

3. **What is skewness, and how does it affect data interpretation?**

**Answer**: Skewness measures the asymmetry of a data distribution. A positive skew means data is stretched to the right, and a negative skew means it is stretched to the left.

Skewness affects the interpretation of central tendency and the choice of statistical methods (e.g., mean is sensitive to skewness).

4. **What is a histogram, and how can it be used to visualize data distribution?**

**Answer**: A histogram is a graphical representation of data distribution, showing the frequency of data within specified ranges (bins).

It helps Identify the shape of the distribution (e.g., normal, skewed, bimodal).

### 5. What are percentiles and quartiles, and how do they help in data analysis?

**Answer**: Percentiles divide data into 100 equal parts. The 25th, 50th, and 75th percentiles are known as quartiles.

They help to understand data spread, identify outliers, and compare distributions.

### 6. What is the central limit theorem, and why is it relevant in data analysis?

**Answer**: The central limit theorem states that, for a sufficiently large sample size, the sampling distribution of the sample mean will be approximately normally distributed, regardless of the population's distribution.

It Is important because it allows for the use of normal distribution-based techniques even with non-normally distributed data.

### 7. What are descriptive vs inferential statistics? When would you use each?

**Answer**: Descriptive statistics summarize data through measures like mean, median, and standard deviation (used to describe and present data).

Inferential statistics make predictions or generalizations about a population based on sample data (used when making inferences).

### 8. What is a p-value, and how does it inform decision-making in hypothesis testing?

**Answer**: The p-value is the probability of observing the data, or something more extreme, under the null hypothesis.

A low p-value (usually $< 0.05$) suggests strong evidence against the null hypothesis, indicating that the result is statistically significant.

### 9. What are Type I and Type II errors, and how do they impact analysis?

**Answer**: Type I error occurs when the null hypothesis is incorrectly rejected (false positive).

Type II error occurs when the null hypothesis is incorrectly accepted (false negative).

Both errors impact the accuracy of decisions made from hypothesis testing.

10. **How do correlation and causation differ, and why is this important in analytics?**

**Answer**: Correlation indicates a relationship between two variables, but it does not imply causation.

Causation means one variable directly influences another.

Understanding the difference is critical to avoid making incorrect assumptions in business analysis.

**Intermediate Level**

11. **What is regression analysis, and how is it used in data analysis?**

**Answer**: Regression analysis models the relationship between a dependent variable and one or more independent variables to predict outcomes.

It Is used in forecasting, trend analysis, and risk assessment.

12. **What are residuals in regression, and what do they indicate?**

**Answer**: Residuals are the differences between the observed and predicted values in a regression model.

They indicate the model's accuracy, and patterns in residuals can reveal issues like heteroscedasticity or model misspecification.

13. **What is the difference between simple and multiple regression?**

**Answer**: Simple regression involves one independent variable to predict the dependent variable.

Multiple regression involves two or more independent variables.

14. **What is multicollinearity, and why is it a problem in regression analysis?**

**Answer**: Multicollinearity occurs when independent variables in a regression model are highly correlated with each other.

It makes it difficult to determine the individual effect of each variable and can distort coefficient estimates.

15. **Explain confidence intervals and how they are interpreted in data analysis.**

**Answer**: A confidence Interval is a range of values used to estimate the true value of a population parameter.

It gives a level of confidence (e.g., 95%) that the true value lies within that range.

16. **What are the different types of sampling methods (random, stratified, cluster)?**

**Answer**: Random sampling: Every individual has an equal chance of being selected.

Stratified sampling: The population is divided into subgroups, and random samples are taken from each subgroup.

Cluster sampling: The population is divided into clusters, and entire clusters are randomly selected.

17. **What is a normal distribution, and why is it commonly used in statistics?**

**Answer**: A normal distribution is a symmetric, bell-shaped distribution characterized by its mean and standard deviation.

It is commonly used because many statistical methods assume normality, and many natural phenomena follow a normal distribution.

18. **What is the null hypothesis and alternative hypothesis, and how do they guide hypothesis testing?**

**Answer**: The null hypothesis suggests no effect or relationship, while the alternative hypothesis suggests that an effect or relationship exists.

They guide hypothesis testing by providing a framework for statistical testing (usually aiming to reject the null).

19. **What is a t-test, and when would you use it in a business context?**

**Answer**: A t-test compares the means of two groups to determine if they are significantly different.

It is used in business to compare two groups, such as testing the effect of a new marketing strategy on sales.

20. **What is the difference between parametric and non-parametric tests?**

**Answer**: Parametric tests assume data follows a known distribution (e.g., normal distribution) and are used when data meets these assumptions.

Non-parametric tests do not assume any specific distribution and are used when data does not meet parametric test assumptions.

21. **What is an ANOVA test, and how is it used to compare datasets?**

**Answer**: ANOVA (Analysis of Variance) tests for significant differences between the means of three or more groups.

It Is used when comparing multiple groups to determine if any significant differences exist.

**Advanced Level**

22. **What is regularization (e.g., Ridge and Lasso regression) in the context of preventing overfitting?**

**Answer**: Regularization adds a penalty term to the loss function to constrain model complexity and prevent overfitting.

Ridge regression adds a penalty proportional to the square of coefficients (L2 regularization).

Lasso regression adds a penalty proportional to the absolute value of coefficients (L1 regularization), which can also perform feature selection.

### 23. What is cross-validation, and why is it useful for assessing model performance?

**Answer**: Cross-validation involves splitting data into multiple subsets and training/testing the model on each subset to assess its performance.

It helps to reduce overfitting and gives a more reliable estimate of model performance.

### 24. How would you handle missing data in a dataset?

Missing data can be handled by:

Imputation: Filling missing values using mean, median, mode, or predictive models.

Deletion: Removing rows or columns with missing data.

Using models that handle missing values directly.

### 25. What is the difference between joint probability, marginal probability, and conditional probability?

**Answer**: Joint probability is the probability of two events happening simultaneously.

Marginal probability is the probability of an event occurring, ignoring the effect of other variables.

Conditional probability is the probability of an event occurring given that another event has already occurred.

### 26. What is A/B testing, and how is it applied in data-driven decision-making?

**Answer**: A/B testing is an experimental approach to compare two versions (A and B) of a product or service to determine which performs better.

It Is used to optimize user experience, marketing strategies, and product features.

### 27. How would you explain seasonality in time series data, and why is it important?

**Answer**: Seasonality refers to regular, predictable patterns in time series data (e.g., higher sales in the holiday season).

It is important because it helps to forecast and make informed decisions based on recurring trends.

### 28. How do you detect and handle outliers in a dataset?

**Answer**: Detection: Outliers can be detected using statistical methods (e.g., IQR, Z-scores) or visualization (e.g., box plots).

Handling: Outliers can be removed, transformed, or capped based on business requirements and data context.

### 29. What is the difference between a box plot and a histogram, and when would you use each?

**Answer**: A box plot displays the median, quartiles, and outliers, useful for identifying the spread and outliers in the data.

A histogram show's the frequency distribution of continuous data, useful for understanding data distribution and identifying skewness.

### 30. How do you interpret regression coefficients in a linear regression model?

**Answer**: Regression coefficients represent the change in the dependent variable for each one-unit change in an independent variable, holding other variables constant.

Positive coefficients indicate a direct relationship, while negative coefficients indicate an inverse relationship.