

# Data Analysis Unit - 1

**Dr. Praveen Barapatre**

# Statistical methods

- **Statistical methods** are techniques used to analyze and interpret data in order to make inferences or decisions based on probability and mathematical principles. These methods are widely used in various fields including science, engineering, business, medicine, social sciences, and more.

# Statistical methods

- 1. Descriptive Statistics:** Descriptive statistics are used to summarize and describe the main features of a dataset. Measures such as mean, median, mode, variance, and standard deviation are often used to describe the central tendency, variability, and distribution of data.
- 2. Inferential Statistics:** Inferential statistics involve making inferences or predictions about a population based on a sample of data. This includes hypothesis testing, confidence intervals, and regression analysis.

# Statistical methods

- 3. Hypothesis Testing:** Hypothesis testing is a method used to determine whether there is enough evidence in a sample of data to infer that a certain condition or hypothesis about a population is true. Common tests include t-tests, chi-square tests, ANOVA (Analysis of Variance), and z-tests.
- 4. Regression Analysis:** Regression analysis is used to investigate the relationship between two or more variables. It helps to understand how changes in one variable are associated with changes in another variable.

# Statistical methods

5. **ANOVA (Analysis of Variance):** ANOVA is a statistical method used to compare means of two or more groups to determine if there are statistically significant differences between them.
6. **Correlation Analysis:** Correlation analysis examines the strength and direction of the relationship between two or more variables. The Pearson correlation coefficient is commonly used to measure the strength of linear relationships.
7. **Non-parametric Methods:** Non-parametric methods are used when data does not meet the assumptions of parametric statistics. These methods do not require the data to be normally distributed and are based on ranks or frequencies rather than specific numerical values.

# Statistical methods

8. **Time Series Analysis:** Time series analysis is used to analyze data collected over time to identify patterns, trends, and seasonal variations.
9. **Experimental Design:** Experimental design involves planning and conducting experiments in order to optimize the efficiency and effectiveness of data collection, ensuring that valid conclusions can be drawn from the results.
10. **Bayesian Methods:** Bayesian methods involve using probability to represent uncertainty in statistical inference. They incorporate prior knowledge or beliefs about the parameters being estimated, updating these beliefs as new data becomes available.

# Misuse of Statistical Methods

**Misuse of statistical methods** occurs when these techniques are applied incorrectly, leading to misleading or incorrect conclusions.

- 1. Inappropriate Sample Size:** Drawing conclusions from a sample size that is too small can lead to unreliable results. Conversely, using an excessively large sample size may waste resources without providing additional meaningful insights.
- 2. Cherry-Picking Data:** Selectively choosing data that supports a particular hypothesis while ignoring contradictory evidence can bias results and lead to erroneous conclusions.

# Misuse of Statistical Methods

- 3. Data Dredging or P-Hacking:** Repeatedly analyzing data until a significant result is found, without adjusting for multiple comparisons, can lead to false positives. This practice undermines the integrity of statistical inference.
- 4. Misinterpretation of Correlation as Causation:** Assuming that a correlation between two variables implies a causal relationship without considering other factors or conducting further research can lead to incorrect conclusions.



# Misuse of Statistical Methods

- 5. **Overfitting:** Fitting a statistical model too closely to a particular dataset, including noise or random fluctuations, can result in poor generalization to new data and misleading interpretations.
- 6. **Confounding Variables:** Failing to account for confounding variables—factors that influence both the independent and dependent variables—can lead to spurious correlations and incorrect conclusions about causal relationships.
- 7. **Publication Bias:** Journals may preferentially publish studies with significant results, leading to an overrepresentation of positive findings in the literature and distorting the overall body of evidence.

# Misuse of Statistical Methods

- 8. Misleading Visualizations:** Presenting data using misleading graphs or charts, such as using inappropriate scales or omitting relevant information, can distort the interpretation of results.
- 9. Inadequate Statistical Reporting:** Failing to report essential details about the statistical methods used, including assumptions, sample sizes, and measures of uncertainty, can make it difficult for others to evaluate the validity of the analysis.
- 10. Misleading Interpretation of Statistical Significance:** Overemphasis on statistical significance without considering effect size, practical significance, or contextual factors can lead to exaggerated claims or unwarranted conclusions.

# Misuse of Statistical Methods

- To mitigate the risk of misuse, researchers and practitioners should adhere to best practices in statistical analysis, including transparent reporting, careful consideration of assumptions and limitations, and robust validation of results through replication and independent verification. Additionally, interdisciplinary collaboration and peer review can help identify and address potential sources of bias or error in statistical analyses.

# Sampling & Sampling Size

- Sampling is the process of selecting a subset of individuals or items from a larger population to estimate characteristics of the whole population. In statistical methods, sampling is crucial because it's often impractical or impossible to collect data from an entire population. Therefore, researchers use samples to make inferences about population parameters.
- Sampling Size: Sampling size refers to the number of individuals or items selected for inclusion in the sample. The size of the sample is a critical consideration in statistical analysis, as it affects the precision and reliability of the estimates derived from the sample.

# Sampling & Sampling Size

- 1. Representativeness:** A sample must be representative of the population from which it is drawn to ensure that inferences made from the sample can be generalized to the population as a whole. A larger sample size generally improves the representativeness of the sample.
- 2. Precision and Accuracy:** Larger sample sizes tend to produce more precise estimates of population parameters. As the sample size increases, the variability or margin of error in the estimates decreases, leading to more accurate results.

# Sampling & Sampling Size

- 3. Statistical Power:** Statistical power refers to the ability of a study to detect true effects or differences when they exist. Larger sample sizes increase the statistical power of a study, making it more likely to detect real effects or relationships.
- 4. Cost and Feasibility:** While larger sample sizes offer benefits in terms of precision and power, they also tend to be more expensive and time-consuming to collect and analyze. Researchers must balance the desire for larger samples with practical constraints such as budget, time, and available resources.

# Sampling & Sampling Size

- 5. Effect Size and Research Goals:** The effect size, or the magnitude of the difference or relationship being studied, can influence the optimal sample size. Studies aiming to detect small effect sizes may require larger samples to achieve adequate statistical power.
- 6. Sampling Method:** The method of sampling (e.g., random sampling, stratified sampling, cluster sampling) can also influence the required sample size. Certain sampling methods may require larger sample sizes to achieve the same level of precision as others.

# Sampling & Sampling Size

- 7. Population Variability:** Higher variability within the population typically requires larger sample sizes to achieve the same level of precision compared to populations with lower variability.
- 8. Confidence Level and Margin of Error:** The desired level of confidence (e.g., 95%, 99%) and acceptable margin of error also influence the determination of sample size. Higher confidence levels and smaller margins of error generally require larger sample sizes.