

# Data Analysis Unit - 1

**Dr. Praveen Barapatre**

# Data Preparation and Cleaning

1. **Data Collection:** The initial step is to collect data relevant to your analysis. Ensuring you have the right and complete dataset is crucial.
2. **Data Review:** Review the data to understand what you have and whether it's suitable for your analysis.
3. **Identifying Errors:** Test the data to identify errors such as outliers, missing values, or inconsistencies.
4. **Data Cleansing:** Correct errors and clean the data. This may involve removing invalid data, imputing missing values, or correcting inaccuracies.

# Data Preparation and Cleaning

5. **Data Formatting:** Format the data in a way that is suitable for your analysis. This may involve converting dates to a standardized format and transforming string data into numerical values.
6. **Data Structuring:** Organize the data into a structured format, such as tables, to make analysis easier.
7. **Sharing Your Data:** If you intend to share your data with others, ensure it is secure, confidential, and ready for content control.

# Missing Data and Data Errors

## Missing Data:

- Missing data refers to observations or values that are absent from the dataset.
- Missing data can lead to biased estimates and reduce the power of statistical tests.
- Ignoring missing data or using inappropriate methods to handle it can distort results and lead to incorrect conclusions.
- Common approaches to handling missing data include deletion (removing observations with missing values), imputation (replacing missing values with estimated values), or using advanced techniques like multiple imputation.

# Missing Data and Data Errors

## Data Errors:

- Data errors encompass various inaccuracies or inconsistencies present in the dataset.
- Errors can arise due to data entry mistakes, measurement errors, or data processing issues.
- Data errors can distort summary statistics, correlations, and relationships between variables, leading to erroneous interpretations.
- Identifying and correcting data errors is crucial for ensuring the integrity and reliability of statistical analyses.
- Techniques for detecting and addressing data errors include outlier detection, data validation checks, and data cleaning procedures.

# Exploratory Data Analysis (EDA)

- Exploratory Data Analysis (EDA) is a critical initial step in statistical analysis that involves examining and summarizing data sets to understand their main characteristics, uncover patterns, and identify relationships among variables.

# Exploratory Data Analysis (EDA)

- 1. Summary Statistics:** Calculate and present descriptive statistics such as mean, median, mode, standard deviation, variance, minimum, maximum, and quartiles for each variable in the dataset. These statistics provide a basic understanding of the distribution and central tendencies of the data.
- 2. Data Visualization:** Create graphical representations of the data using histograms, box plots, scatter plots, bar plots, and heatmaps. Visualization helps in understanding the distribution, variability, and relationships between variables visually, making it easier to identify patterns and outliers.

# Exploratory Data Analysis (EDA)

- 3. Missing Data Analysis:** Identify and analyze missing values in the dataset. Understanding the extent and patterns of missing data is crucial for determining appropriate strategies for handling them, such as deletion, imputation, or modeling missingness.
- 4. Outlier Detection:** Detect outliers, which are data points that deviate significantly from the rest of the data. Outliers can affect the accuracy of statistical estimates and models. Techniques for outlier detection include graphical methods, such as box plots and scatter plots, as well as statistical methods like Z-score, Tukey's method, and clustering-based approaches.



# Exploratory Data Analysis (EDA)

5. **Exploring Relationships:** Analyze relationships between variables using correlation matrices, scatter plots, and pair plots. Understanding these relationships helps identify potential predictors or explanatory variables for further analysis.
6. **Data Transformation:** Explore transformations such as log transformations, square root transformations, or Box-Cox transformations to stabilize variance, reduce skewness, or make data more normally distributed, which can improve the performance of statistical models.

# Exploratory Data Analysis (EDA)

7. **Dimensionality Reduction:** Explore techniques like Principal Component Analysis (PCA) or t-distributed Stochastic Neighbor Embedding (t-SNE) to visualize high-dimensional data in lower dimensions and identify patterns or clusters.
8. **Interactive Exploration:** Utilize interactive tools and dashboards for exploring data dynamically, enabling users to interact with visualizations and filter data based on specific criteria.

# Exploratory Data Analysis (EDA)

By conducting EDA, analysts gain insights into the structure, quality, and characteristics of the data, informing subsequent analyses and modeling decisions in statistics. EDA is an iterative process that often guides researchers in formulating hypotheses and designing more targeted analyses.

# Statistical Error

- Statistical errors refer to mistakes or inaccuracies that occur during the process of statistical analysis, leading to incorrect conclusions or interpretations. These errors can occur at various stages of the statistical analysis, including data collection, data preparation, modeling, and interpretation of results.

# Statistical Error

1. **Sampling Errors:** These errors occur when the sample used in the analysis is not representative of the population from which it is drawn. Sampling errors can lead to biased estimates and incorrect inferences about the population parameters.
2. **Measurement Errors:** Measurement errors occur when there are inaccuracies in the measurement or recording of data. This can result from human error, faulty instruments, or inconsistencies in data collection procedures. Measurement errors can distort relationships between variables and lead to incorrect conclusions.

# Statistical Error

3. **Type I Error (False Positive):** Type I error occurs when the null hypothesis is incorrectly rejected when it is actually true. In other words, it is the probability of concluding that there is an effect or relationship in the data when there is none. The significance level ( $\alpha$ ) of a statistical test determines the probability of making a Type I error.
4. **Type II Error (False Negative):** Type II error occurs when the null hypothesis is incorrectly accepted when it is actually false. It is the probability of failing to detect an effect or relationship in the data when it exists. The power of a statistical test is the probability of correctly rejecting a false null hypothesis and avoiding a Type II error.

# Statistical Error

5. **Confounding Variables:** Confounding occurs when an extraneous variable is associated with both the independent and dependent variables, leading to a spurious relationship between them. Failure to account for confounding variables can result in incorrect conclusions about the true relationship between variables.
6. **Overfitting:** Overfitting occurs when a statistical model captures noise or random fluctuations in the data instead of the underlying patterns. Overfitted models may perform well on the training data but generalize poorly to new data, leading to unreliable predictions.

# Statistical Error

7. **Selection Bias:** Selection bias occurs when certain groups or observations are systematically excluded or overrepresented in the sample, leading to biased estimates of population parameters.
8. **Misinterpretation of Results:** Errors can also arise from misinterpretation of statistical results, such as misreading p-values, misinterpreting confidence intervals, or attributing causality to correlation.



# Statistical Error

- It's important for researchers and analysts to be aware of these potential sources of error and take appropriate steps to minimize their impact on the validity and reliability of statistical analyses.
- This includes careful planning of study design, rigorous data collection and validation procedures, appropriate choice and execution of statistical methods, and cautious interpretation of results. Additionally, peer review and replication of findings can help identify and correct errors in statistical analyses.

# Statistical Modeling

Statistical modeling involves using mathematical and statistical techniques to analyze data, understand relationships between variables, and make predictions or inferences.

- 1. Problem Formulation:** Clearly define the research question or problem that the statistical model aims to address.
- 2. Data Collection:** Gather relevant data that are suitable for addressing the research question.
- 3. Exploratory Data Analysis (EDA):** Explore the data using descriptive statistics and data visualization techniques to understand its characteristics and identify patterns.

# Statistical Modeling

4. **Model Selection:** Choose an appropriate statistical model based on the nature of the data and the research question. Common models include linear regression, logistic regression, time series models, and machine learning algorithms.
5. **Model Development:** Develop the statistical model by specifying the mathematical relationship between the variables and estimating model parameters using statistical techniques such as maximum likelihood estimation or Bayesian inference.
6. **Model Evaluation:** Assess the performance of the model using various techniques, such as goodness-of-fit tests, diagnostic plots, and cross-validation, to ensure that it adequately captures the underlying patterns in the data.

# Statistical Modeling

7. **Model Interpretation:** Interpret the results of the statistical model to gain insights into the relationships between variables and the factors influencing the outcome of interest.
8. **Prediction or Inference:** Use the fitted statistical model to make predictions about future observations or draw conclusions about the population based on the sample data.
9. **Validation and Sensitivity Analysis:** Validate the model's performance on independent data and conduct sensitivity analysis to assess the robustness of the results to variations in model assumptions or parameter values.
10. **Communication of Results:** Present the findings of the statistical modeling analysis in a clear and understandable manner, emphasizing key insights and implications for decision-making.

# Statistical Modeling

Overall, statistical modeling is a systematic approach to analyzing data and extracting meaningful information to support decision-making in various fields, including science, engineering, business, and social sciences. It combines mathematical theory, statistical methods, computational techniques, and domain knowledge to develop models that accurately represent real-world phenomena.

# Computational Statistics

- Computational statistics utilizes numerical methods, algorithms, and statistical techniques to efficiently analyze data and solve complex statistical problems.
- It involves simulation, algorithm design, optimization, and the use of statistical software to handle large datasets, perform data analysis, and make data-driven decisions.
- Computational statistics plays a crucial role in modern data analysis, scientific research, and decision-making by providing tools and methods to extract meaningful insights from data efficiently.

# **Computational Statistics Includes**

- 1.Numerical Methods**
- 2.Simulation**
- 3.Algorithm Design**
- 4.Computational Complexity**
- 5.Statistical Software**
- 6.High-Performance Computing**
- 7.Optimization**
- 8.Validation and Verification**

# Definitions

**Inference:** In statistics, inference refers to the process of drawing conclusions or making predictions about a population based on sample data. It involves using statistical methods to generalize findings from a sample to a larger population, while acknowledging uncertainty.

**Bias:** Bias refers to systematic errors or inaccuracies in data collection, analysis, or interpretation that lead to deviations from the true values or relationships in the population. Bias can arise from various sources, such as sampling methods, measurement errors, or confounding variables, and can distort the results of statistical analyses.



# Definitions

**Confounding:** Confounding occurs when an extraneous variable (confounder) is associated with both the independent variable and the outcome variable in a study, leading to a spurious or misleading association between them. Confounding variables can obscure the true relationship between the variables of interest, making it difficult to draw valid conclusions.

**Hypothesis Testing:** Hypothesis testing is a statistical method used to make decisions or draw inferences about a population parameter based on sample data. It involves formulating a null hypothesis ( $H_0$ ) and an alternative hypothesis ( $H_a$ ), collecting data, calculating a test statistic, and comparing it to a critical value or p-value to determine whether there is sufficient evidence to reject the null hypothesis in favor of the alternative hypothesis. The outcome of a hypothesis test helps researchers assess the strength of evidence for or against a particular hypothesis.

# Definitions

**Confidence Interval** : confidence interval is a range of values around an estimated parameter, such as a mean or proportion, calculated from sample data. It provides a measure of uncertainty about the true value of the parameter and is expressed with a specified level of confidence, typically 95% or 99%. The interval is interpreted as "We are [confidence level]% confident that the true parameter lies within this range." Confidence intervals help researchers assess the precision of their estimates and make inferences about population parameters based on sample data.

# Statistical significance

- Statistical significance refers to the likelihood that an observed difference or relationship in data is not due to chance variation. In statistical hypothesis testing, statistical significance is assessed by comparing the observed data to a null hypothesis ( $H_0$ ), which typically states that there is no effect, difference, or relationship in the population.
- The concept of statistical significance is often associated with p-values. A p-value is a measure of the strength of evidence against the null hypothesis. It represents the probability of observing the data, or more extreme data, if the null hypothesis were true. A smaller p-value indicates stronger evidence against the null hypothesis and suggests that the observed effect or relationship is unlikely to be due to random chance alone.

# Statistical significance

- Generally, if the p-value is less than a pre-determined significance level (often denoted as  $\alpha$ , commonly set to 0.05), the null hypothesis is rejected, and the result is considered statistically significant. This means that the observed difference or relationship is unlikely to have occurred by chance alone, and there is evidence to support the alternative hypothesis ( $H_a$ ).
- It's important to note that statistical significance does not necessarily imply practical significance or importance. Even if a result is statistically significant, its effect size and relevance in real-world contexts should also be considered when interpreting the findings. Additionally, statistical significance depends on factors such as sample size, effect size, and variability in the data, so it's essential to interpret significance in the context of the specific study and its limitations.

# Power

1. Power refers to the probability of correctly rejecting a false null hypothesis in a hypothesis test. In other words, it is the probability of detecting a true effect or relationship when it exists.
2. A statistical test with high power is more likely to detect an effect or difference if it truly exists in the population, whereas a test with low power is less sensitive and may fail to detect real effects.
3. Power depends on various factors, including the sample size, effect size, significance level ( $\alpha$ ), and variability in the data. Increasing the sample size or effect size typically increases the power of a statistical test.

# Robustness

1. Robustness refers to the ability of a statistical method or model to provide valid results even when its assumptions are not perfectly met or when there are deviations from ideal conditions.
2. A robust statistical method is less sensitive to violations of assumptions or outliers in the data and tends to produce reliable results across different conditions.
3. Robustness is particularly important in practical applications where data may not always adhere to the assumptions of traditional statistical methods. Robust methods can provide more accurate and stable estimates in such cases.

# Degrees of freedom

- Degrees of freedom (df) is a concept used in statistics to describe the number of independent pieces of information available in a dataset or the number of values in a sample that are free to vary. The concept of degrees of freedom is relevant in various statistical calculations, including hypothesis testing, estimation of parameters, and model fitting.
- In general, the degrees of freedom can be defined differently depending on the context:

# Degrees of freedom

- **Sample Size Minus One:** In the context of estimating a population parameter from sample data, the degrees of freedom is typically equal to the sample size minus one ( $n - 1$ ). This adjustment accounts for the fact that when estimating parameters from sample data, we lose one degree of freedom because the sample mean is used as an estimate of the population mean.
- **Model Complexity:** In the context of regression analysis or analysis of variance (ANOVA), degrees of freedom represent the number of independent pieces of information available to estimate model parameters. For example, in simple linear regression with one predictor variable, the degrees of freedom for the regression model would be the total sample size minus two ( $n - 2$ ): one degree of freedom is used for estimating the slope coefficient, and another for the intercept.



# Degrees of freedom

- **Residuals:** Degrees of freedom can also refer to the number of independent observations remaining after fitting a model or conducting a statistical test. In hypothesis testing, the degrees of freedom associated with the test statistic reflect the number of observations that are free to vary given the constraints imposed by the null hypothesis.
- Understanding degrees of freedom is crucial for correctly interpreting statistical results and selecting appropriate statistical methods. It helps ensure that statistical analyses are conducted with the appropriate level of precision and accuracy, accounting for the variability in the data and the constraints imposed by the statistical model or hypothesis being tested.

# Non parametric analysis

- Non-parametric analysis refers to a set of statistical methods that don't rely on assumptions about the underlying distribution of the data. Instead of assuming a specific distribution, these methods focus on ranking or ordering data.
- They are useful when data don't meet the requirements of parametric tests or when dealing with categorical or ordinal data. Non-parametric tests are less sensitive to outliers and provide robust alternatives for hypothesis testing and data analysis without strict distributional assumptions.
- Examples include the Mann-Whitney U test, Kruskal-Wallis test, and Spearman's rank correlation coefficient.

