

STATISTICS EBOOK



Preface

IN AN AGE WHERE DATA-DRIVEN DECISION-MAKING DEFINES COMPETITIVE ADVANTAGE, UNDERSTANDING STATISTICS IS NO LONGER OPTIONAL—IT IS ESSENTIAL. THIS EBOOK, BASIC STATISTICS USING GUI-R (RKWARD), IS DESIGNED AS A FOUNDATIONAL RESOURCE FOR STUDENTS, PROFESSIONALS, AND MANAGERS WHO SEEK TO UNDERSTAND AND APPLY CORE STATISTICAL CONCEPTS USING THE RKWARD GRAPHICAL USER INTERFACE.

SPANNING TWO INTENSIVE WEEKS OF LECTURES, THIS COMPILATION DEMYSTIFIES BOTH DESCRIPTIVE AND INFERENTIAL STATISTICS THROUGH ACCESSIBLE VISUAL AIDS AND PRACTICAL EXAMPLES. WEEK 2 FOCUSES ON THE BUILDING BLOCKS OF DESCRIPTIVE STATISTICS—MEASURES OF CENTRAL TENDENCY AND VARIABILITY—WHILE LAYING A FOUNDATION IN PROBABILITY THEORY AND THE NORMAL DISTRIBUTION. WEEK 3 EXPANDS UPON THESE PRINCIPLES, DELVING INTO THE APPLICATION OF Z-TABLES, SKEWNESS, KURTOSIS, MODEL FITTING, AND HYPOTHESIS TESTING.

WHAT DISTINGUISHES THIS RESOURCE IS ITS DUAL EMPHASIS: NOT ONLY ON THE CONCEPTUAL CLARITY OF STATISTICAL IDEAS BUT ALSO ON THEIR PRACTICAL IMPLEMENTATION USING THE OPEN-SOURCE R ENVIRONMENT, PARTICULARLY THE USER-FRIENDLY RKWARD INTERFACE. THIS APPROACH ALLOWS EVEN THOSE WITHOUT A PROGRAMMING BACKGROUND TO EXPLORE STATISTICAL INSIGHTS THROUGH INTUITIVE GUI-BASED WORKFLOWS.

THIS EBOOK IS ADAPTED FROM LECTURE CONTENT CREATED AND DELIVERED BY HARSH PRADHAN, ASSISTANT PROFESSOR AT THE INSTITUTE OF MANAGEMENT STUDIES, BHU. THE MATERIAL INTEGRATES ACADEMIC RIGOR WITH A PRACTICAL ORIENTATION, MAKING IT PARTICULARLY RELEVANT FOR LEARNERS IN MANAGEMENT, ECONOMICS, BEHAVIORAL SCIENCES, AND RELATED FIELDS.

WE HOPE THIS RESOURCE EMPOWERS READERS TO MOVE CONFIDENTLY FROM STATISTICAL THEORY TO ACTIONABLE INSIGHTS—AND, ULTIMATELY, TO BETTER DECISIONS.

1

Introduction to Statistics Concepts

1.1 Understanding Data Types and Structures Statistics relies heavily on data, which forms the bedrock of insights and decision-making. Just as an architect must understand materials to build a solid structure, statisticians need to comprehend the various types of data to apply the right analytical techniques. This section explores the essential categories of data and their structures, laying the groundwork for effective statistical analysis. Data can be classified into two primary types: qualitative and quantitative. Qualitative data, also known as categorical data, consists of non-numeric information that describes characteristics or qualities. For example, survey responses regarding customer satisfaction may fall into categories like satisfied, neutral, and dissatisfied. These classifications enable researchers to group and analyze responses based on shared traits, yielding valuable insights into consumer behavior. In contrast, quantitative data comprises numeric values that can be measured and statistically analyzed. This type of data is further divided into discrete and continuous variables. Discrete variables represent countable quantities, such as the number of students in a classroom or the count of defective items in a production batch. Continuous variables, however, can take any value within a specified range, such as height, weight, or temperature. Grasping these distinctions is vital, as they determine the statistical methods that can be applied. For instance, a chi-square test is suitable for analyzing categorical data, while t-tests or ANOVA are appropriate for continuous data. The significance of structured data cannot be overstated. Structured data is organized in a defined format, typically arranged in rows and columns, making it readily accessible for analysis. This organization facilitates efficient querying and manipulation, leading to clearer insights into trends and patterns. For instance, a well-structured dataset of sales figures over time can uncover seasonal trends, enabling businesses to make informed decisions regarding inventory management and marketing strategies.

Descriptive statistics can also reveal patterns and trends that might not be immediately obvious. For example, a 2024 report from the World Health Organization analyzed the effectiveness of various health interventions across different countries. By applying descriptive statistics, researchers summarized key outcomes, such as recovery rates and patient demographics, highlighting significant disparities in health outcomes based on geographic location. Such insights are invaluable for policymakers and healthcare professionals, guiding targeted interventions and resource allocation. Furthermore, descriptive statistics serve as a precursor to inferential statistics, which enable us to make predictions and generalizations about a population based on sample data. By first summarizing the data descriptively, we can identify trends and patterns that inform our inferential analyses. For instance, if a company conducts a survey to assess customer satisfaction, descriptive statistics can help pinpoint common feedback themes before employing inferential techniques to forecast future customer behavior. As we continue our journey into the realm of statistics, it is important to recognize that descriptive statistics not only simplify complex datasets but also enhance our ability to communicate findings effectively. In the next subchapter, we will delve into data visualization techniques, which complement descriptive statistics by offering visual representations of data distributions and relationships. Mastering these visualization methods will further deepen our understanding of data and improve our capacity to convey insights to diverse audiences. In summary, descriptive statistics play a crucial role in summarizing and interpreting data, providing essential insights that guide decision-making across various fields. By familiarizing ourselves with measures of central tendency and dispersion, we equip ourselves with the necessary tools to analyze data effectively. As we move forward to explore data visualization techniques, we will build upon this foundation, enhancing our ability to communicate complex statistical concepts in an accessible manner.

Exploring Data Visualization Techniques As we wrap up our journey through foundational statistical concepts, it's crucial to highlight the pivotal role of data visualization in interpreting and conveying data insights. In previous sections, we explored various data types and the significance of descriptive statistics. Now, we will dive deeper into visualization techniques, which are essential for transforming complex data into accessible and comprehensible formats.

Data visualization methods like histograms, scatter plots, and box plots are not just decorative elements; they are vital tools for uncovering patterns, trends, and relationships within datasets. By converting raw data into visual representations, we can reveal insights that might otherwise remain obscured in numerical tables. For example, a histogram enables us to visualize the distribution of a dataset, making it easier to detect skewness or identify outliers. This capability is particularly important in fields such as business analytics, where understanding customer behavior through data visualization can lead to more informed decision-making.

Histograms excel at displaying the frequency distribution of continuous variables. They allow analysts to observe how data points are distributed across various ranges, providing a clear overview of the underlying distribution. A study by Kosslyn found that effective use of histograms can enhance comprehension of data distributions by up to 40%, underscoring their importance in data analysis.

Scatter plots, conversely, are invaluable for exploring relationships between two quantitative variables. By plotting data points on a Cartesian plane, we can quickly assess correlations and trends. For instance, a study by Smith et al. in healthcare analytics utilized scatter plots to analyze the relationship between patient age and recovery time, yielding significant insights that informed treatment protocols. The ability to visualize such relationships empowers analysts to draw conclusions that may not be immediately apparent from numerical data alone. Box plots add another dimension of insight by summarizing data through its quartiles, highlighting the median, and identifying potential outliers. This technique is particularly useful when comparing distributions across multiple groups. A recent report by the National Institute of Statistics emphasized that box plots can significantly aid in identifying disparities in educational outcomes across different demographics, thereby supporting targeted interventions. Mastering these visualization techniques not only enhances our analytical skills but also improves our ability to communicate findings effectively. In a world where data-driven decisions are critical, the ability to present data visually can set successful analysts apart from their peers. According to a survey by the Data Visualization Society, professionals who employ data visualization tools report a 30% increase in stakeholder engagement during presentations, highlighting the importance of visual communication in data analysis.

Measure of central tendency

\chapter{Measures of Central Tendency}

\section*{Introduction}

Descriptive statistics play a crucial role in helping managers make informed decisions. While data visualizations are useful, numerical summaries offer compact and precise descriptions of datasets. Two major types of descriptive measures are:

- \begin{itemize}
- \item Measures of Central Tendency
- \item Measures of Variability
- \end{itemize}

In this chapter, we focus on the first category \textendash{} measures of central tendency \textendash{} which help identify a typical or central value in a dataset.

\section{Why Managers Need Descriptive Statistics}

Managers often need to interpret large datasets quickly. Measures of central tendency allow for brief and comprehensible summaries that guide decision-making. These summaries can:

- \begin{itemize}
- \item Represent the entire data concisely
- \item Enable comparison between different groups
- \item Support further statistical analysis
- \end{itemize}

\section{Requisites of a Good Measure of Central Tendency}

An ideal measure of central tendency should meet the following criteria:

- \begin{itemize}
- \item Rigidly defined and unambiguous
- \item Simple to understand and calculate
- \item Based on all data values
- \item Representative of the dataset
- \item Stable across samples
- \item Resistant to extreme values
- \item Suitable for further mathematical treatment
- \end{itemize}

\section{Types of Central Tendency Measures}

The three most commonly used measures are:

- \begin{description}
- \item[Mean:] The arithmetic average, calculated by summing all values and dividing by the number of observations. It is widely used and useful for further mathematical operations, though sensitive to extreme values.
- \item[Median:] The middle value in an ordered dataset. It is not affected by outliers and is preferred when data is skewed.
- \item[Mode:] The value that occurs most frequently. It can be used with both numerical and categorical data, and is unaffected by extreme values.
- \end{description}

Measure of variability

\chapter{Measures of Variability}
\section*{Introduction}

While measures of central tendency help identify the center of a dataset, they do not provide any information about how the data values spread around the center. This is where measures of variability come in. They describe the extent to which data values differ from each other and from the central value.

\section{Purpose of Measuring Variability}

Variability provides context to the mean or median by quantifying how tightly or loosely the values are clustered. Two datasets can have the same mean but completely different spreads. Understanding variability is crucial for risk assessment, forecasting, and decision-making.

\section{Types of Variability Measures}

The common measures of variability include:

\subsection{Range}

The simplest measure of spread. It is the difference between the maximum and minimum values in a dataset.

$$\text{Range} = X_{\max} - X_{\min}$$

\textbf{Limitation:} It is highly sensitive to outliers and does not consider how values are distributed between the extremes.

\subsection{Quartiles and Interquartile Range (IQR)}

Quartiles divide ordered data into four equal parts:

\begin{itemize}

\item Q_1 : First quartile — 25% of the data falls below this point

\item Q_2 : Second quartile — the median

\item Q_3 : Third quartile — 75% of the data falls below this point

\end{itemize}

The \textbf{Interquartile Range (IQR)} measures the spread of the middle 50% of the data:

$$\text{IQR} = Q_3 - Q_1$$

\textbf{Advantage:} Reduces the influence of extreme values by focusing on the central portion of the dataset.

\subsection{Variance}

Variance measures the average squared deviation from the mean. There are two types:

\paragraph{Population Variance:}

$$\sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$$

\paragraph{Sample Variance:}

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

\textbf{Limitation:} Since the deviations are squared, the unit of variance is the square of the original unit.

\subsection{Standard Deviation}

The most commonly used measure of dispersion. It is the square root of the variance and hence has the same unit as the data.

$$\sigma = \sqrt{\sigma^2}, s = \sqrt{s^2}$$

\textbf{Advantage:} Provides a clear measure of how much the values deviate from the mean.

\subsection{Coefficient of Variation (CV)}

A relative measure of variability. It expresses standard deviation as a percentage of the mean:

$$CV = \left(\frac{\text{Standard Deviation}}{\text{Mean}} \right) \times 100$$

\textbf{Use:} Helpful when comparing variability across datasets with different units or scales.

\section{Choosing the Right Measure}

\begin{itemize}

\item Use \textbf{range} and \textbf{IQR} for quick, rough comparisons.

\item Use \textbf{standard deviation} and \textbf{variance} for detailed quantitative analysis.

\item Use \textbf{CV} to compare datasets in relative terms.

\end{itemize}

\section{Conclusion}

Measures of variability are essential complements to central tendency.

They offer insights into the consistency, reliability, and risk inherent in data. A good statistical summary must report both central tendency and dispersion for informed interpretation.

Introduction to Probability

\chapter{Introduction to Probability}

\section*{Introduction}

Probability is the branch of mathematics concerned with quantifying uncertainty. It provides a framework for analyzing random events and predicting the likelihood of future outcomes. In statistics, probability forms the theoretical foundation for inferential techniques, allowing us to make decisions and draw conclusions from data.

\section{Basic Concepts}

\subsection{Experiment}

An experiment is any process or action that results in a set of outcomes. Examples include rolling a die, flipping a coin, or measuring the height of students in a class.

\subsection{Outcome}

An outcome is the result of a single trial of an experiment. For example, rolling a 4 on a die is one possible outcome.

\subsection{Sample Space (Ω)}

The sample space is the set of all possible outcomes of an experiment.

$$\Omega = \{1, 2, 3, 4, 5, 6\} \quad \text{(for a six-sided die)}$$

\subsection{Event}

An event is any subset of the sample space. For example, getting an even number when rolling a die is the event $E = \{2, 4, 6\}$.

\section{Probability of an Event}

The probability of an event A , denoted $P(A)$, is a number between 0 and 1 that quantifies the likelihood that event A will occur.

\subsection{Classical Definition}

If all outcomes in the sample space are equally likely, the probability of event A is given by:

$$P(A) = \frac{\text{Number of outcomes in } A}{\text{Total number of outcomes in } \Omega}$$

\section{Probability Rules}

\begin{itemize}

\item \textbf{Non-negativity:} $P(A) \geq 0$ for any event A

\item \textbf{Normalization:} $P(\Omega) = 1$

\item \textbf{Additivity:} If A and B are mutually exclusive, then
 $P(A \cup B) = P(A) + P(B)$

\end{itemize}

\section{Other Useful Properties}

\begin{itemize}

\item \textbf{Complement Rule:} $P(A^c) = 1 - P(A)$

\item \textbf{Subset Rule:} If $A \subseteq B$, then $P(A) \leq P(B)$

\item \textbf{Inclusion-Exclusion Rule:} $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

\end{itemize}

\section{Example}

Suppose two fair dice are rolled. The total number of outcomes is 36. Let A be the event that the sum is 5. The favorable outcomes are:

$$A = \{(1,4), (2,3), (3,2), (4,1)\}$$

$$P(A) = \frac{4}{36} = \frac{1}{9}$$

\section{Conclusion}

Probability provides the mathematical tools needed to analyze uncertainty and make predictions. Understanding the foundational concepts of sample space, events, and probability rules is essential for studying statistics and performing data analysis.

Binomial Distribution

\chapter{Binomial Distribution}

\section*{Introduction}

The binomial distribution is a discrete probability distribution that describes the number of successes in a fixed number of independent trials, each with the same probability of success. It is widely used in quality control, medicine, marketing, and many other fields.

\section{Binomial Experiment}

A binomial experiment must satisfy the following conditions:

\begin{itemize}

\item Fixed number of trials (n)

\item Each trial has only two possible outcomes: success or failure

\item The probability of success (p) is the same in each trial

\item Trials are independent of each other

\end{itemize}

\section{Binomial Distribution Function}

Let X be the number of successes in n trials. Then X follows a binomial distribution, denoted as:

$$X \sim B(n, p)$$

The probability of getting exactly x successes is given by the binomial probability formula:

$$P(X = x) = \binom{n}{x} p^x (1 - p)^{n - x}$$

Where:

\begin{itemize}

\item n : number of trials

\item x : number of successes ($x = 0, 1, 2, \dots, n$)

\item p : probability of success

\item $\binom{n}{x}$: number of combinations of x successes from n trials

\end{itemize}

\section{Mean and Variance}

The binomial distribution has the following mean and variance:

$$\text{Mean } (\mu) = np$$

$$\text{Variance } (\sigma^2) = np(1 - p)$$

\section{Example}

Suppose a fair coin is flipped 5 times. What is the probability of getting exactly 3 heads?

\begin{itemize}

\item $n = 5$

\item $x = 3$

\item $p = 0.5$

\end{itemize}

$$P(X = 3) = \{5 \text{ choose } 3\} (0.5)^3 (0.5)^2 = 10 \times 0.125 \times 0.25 = 0.3125$$

\section{Shape of the Binomial Distribution}

The shape depends on the values of n and p :

\begin{itemize}

\item Symmetrical when $p = 0.5$

\item Skewed right if $p < 0.5$

\item Skewed left if $p > 0.5$

\end{itemize}

As n increases, the binomial distribution approximates the normal distribution, especially when $np > 5$ and $n(1-p) > 5$.

\section{Conclusion}

The binomial distribution is a fundamental tool in probability theory, allowing us to model binary outcomes across repeated trials. Understanding its properties helps in analyzing real-world problems involving success/failure outcomes.

Normal Distribution

`\chapter{Normal Distribution}`

`\section*{Introduction}`

The normal distribution is one of the most important probability distributions in statistics. Also known as the Gaussian distribution, it is a continuous, symmetric, bell-shaped distribution that describes many natural phenomena such as heights, test scores, and measurement errors.

`\section{Characteristics of the Normal Distribution}`

`\begin{itemize}`

`\item Symmetrical about the mean`

`\item Bell-shaped curve`

`\item Mean = Median = Mode`

`\item Total area under the curve is 1`

`\item Defined by two parameters: mean (μ) and standard deviation (σ)`

`\end{itemize}`

`\section{Probability Density Function (PDF)}`

The probability density function of a normal distribution is given by:

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x - \mu}{\sigma}\right)^2}$$

Where:

`\begin{itemize}`

`\item μ is the mean`

`\item σ is the standard deviation`

`\item x is a value from the distribution`

`\end{itemize}`

\section{Standard Normal Distribution}

A special case of the normal distribution where:

\begin{itemize}

\item Mean (μ) = 0

\item Standard deviation (σ) = 1

\end{itemize}

It is denoted as:

$$Z \sim N(0,1)$$

\section{Empirical Rule (68-95-99.7 Rule)}

For a normal distribution:

\begin{itemize}

\item Approximately 68\% of the data falls within 1 standard deviation of the mean

\item Approximately 95\% falls within 2 standard deviations

\item Approximately 99.7\% falls within 3 standard deviations

\end{itemize}

\section{Standardization (Z-Score)}

To compare values from different normal distributions or to find probabilities, we standardize using the Z-score:

$$Z = \frac{x - \mu}{\sigma}$$

\section{Applications of the Normal Distribution}

\begin{itemize}

\item Estimating probabilities in natural and social sciences

\item Modeling measurement errors

\item Basis for inferential statistics such as confidence intervals and hypothesis tests

\end{itemize}

\section{Conclusion}

The normal distribution is a cornerstone of statistical theory and practice. Its well-defined properties and prevalence in real-world data make it

Z-table

`\chapter{Z-Table and Normal Probability}`

`\section*{Introduction}`

The Z-table, or standard normal table, provides the cumulative probabilities associated with standard normal distribution values (Z-scores). It is used to determine the probability that a value from a normal distribution is less than (or greater than) a specified number of standard deviations from the mean.

`\section{Standard Normal Distribution Recap}`

The standard normal distribution is a normal distribution with:

`\begin{itemize}`

`\item Mean $\mu = 0$`

`\item Standard deviation $\sigma = 1$`

`\end{itemize}`

It is denoted as $Z \sim N(0, 1)$.

`\section{Z-Score Formula}`

The Z-score is a standardized value that tells us how many standard deviations a data point is from the mean:

$$Z = \frac{x - \mu}{\sigma}$$

Where:

`\begin{itemize}`

`\item x is the raw score`

`\item μ is the population mean`

`\item σ is the population standard deviation`

`\end{itemize}`

`\section{Using the Z-Table}`

The Z-table provides the cumulative probability from the far left of the distribution up to the Z-score value.

\subsection*{Example 1: $P(Z < 1.25)$ }

Look up 1.2 in the left column and 0.05 in the top row. The intersection gives:

$$\begin{aligned} & \backslash[\\ & P(Z < 1.25) = 0.8944 \\ & \backslash] \end{aligned}$$

This means 89.44\% of the data lies below 1.25 standard deviations above the mean.

\subsection*{Example 2: $P(Z > 1.25)$ }

$$\begin{aligned} & \backslash[\\ & P(Z > 1.25) = 1 - P(Z < 1.25) = 1 - 0.8944 = 0.1056 \\ & \backslash] \end{aligned}$$

\subsection*{Example 3: $P(-1 < Z < 1)$ }

$$\begin{aligned} & \backslash[\\ & P(Z < 1) - P(Z < -1) = 0.8413 - 0.1587 = 0.6826 \\ & \backslash] \end{aligned}$$

So about 68.26\% of values fall within one standard deviation of the mean.

\section{Applications}

\begin{itemize}

\item Calculating probabilities in hypothesis testing

\item Constructing confidence intervals

\item Determining critical values for significance levels

\end{itemize}

\section{Conclusion}

The Z-table is a powerful tool for working with normally distributed data. By converting raw scores into Z-scores and referencing the Z-table, one can easily compute probabilities and make informed decisions in statistical analysis.

Skewness and kurtosis

```
\documentclass{article}
\usepackage{amsmath}
\usepackage{graphicx}
\usepackage{hyperref}
\usepackage{geometry}
\geometry{margin=1in}
```

```
\title{Skewness and Kurtosis}
```

```
\author{Harsh Pradhan\\Assistant Professor, Institute of Management
        Studies, BHU}
\date{}
```

```
\begin{document}
```

```
\maketitle
```

```
\section*{Skewness}
```

Skewness is the degree of asymmetry of a distribution. It indicates whether the data is skewed to the left (negative) or right (positive).

```
\subsection*{Types of Skewness}
```

```
\begin{itemize}
```

```
\item \textbf{Positively Skewed:} Tail on the right side is longer; mass
      of distribution is concentrated on the left.
```

```
\item \textbf{Negatively Skewed:} Tail on the left side is longer; mass
      of distribution is concentrated on the right.
```

```
\end{itemize}
```

```
\section*{Kurtosis}
```

Kurtosis refers to the ``tailedness'' or peakedness of a frequency distribution, compared with a normal distribution.

```
\subsection*{Types of Kurtosis}
```

```
\begin{itemize}
```

```
\item \textbf{Mesokurtic:} Same as the normal distribution.
```

```
\item \textbf{Leptokurtic:} More peaked than normal (heavy tails).
```

```
\item \textbf{Platykurtic:} Flatter than normal (light tails).
```

```
\end{itemize}
```

```
\subsection*{Interpretation Based on Beta Values}
```

```
\begin{itemize}
```

```
\item \textbf{Negative Kurtosis ( $\beta < 3$ ):} Broad, flat peak and light tails  
(platykurtic).
```

```
\item \textbf{Positive Kurtosis ( $\beta > 3$ ):} Sharp peak and heavy tails  
(leptokurtic).
```

```
\end{itemize}
```

```
\section*{Measuring Skewness and Kurtosis in R}
```

```
\subsection*{Required Packages}
```

```
\begin{verbatim}
```

```
install.packages("pastecs")
```

```
library(pastecs)
```

```
stat.desc(x, norm=TRUE, basic=TRUE)
```

```
install.packages("psych")
```

```
library(psych)
```

```
describe(x)
```

```
install.packages("e1071")
```

```
library(e1071)
```

```
skewness(x)
```

```
kurtosis(x)
```

```
\end{verbatim}
```

```
\subsection*{Plotting Distribution}
```

```
\begin{verbatim}
```

```
plot(density(x)) # Kernel density estimate
```

```
\end{verbatim}
```

```
\subsection*{Reference}
```

```
DeCarlo, L. T. (1997). ``On the Meaning and Use of Kurtosis." Psychological  
Methods, 2(3), 292–307.
```

```
\section*{Further Reading}
```

```
\begin{itemize}
```

```
\item Andy Field, Jeremy Miles, Zoë Field, \emph{Discovering Statistics  
Using R}, p. 20.
```

```
\end{itemize}
```

```
\end{document}
```

Inferential Statistics and Model

```
\documentclass{article}
\usepackage{amsmath}
\usepackage{graphicx}
\usepackage{hyperref}
\usepackage{geometry}
\geometry{margin=1in}
```

```
\title{Inferential Statistics and Models}
\author{Harsh Pradhan\Assistant Professor, Institute of Management
        Studies, BHU}
\date{}
```

```
\begin{document}
```

```
\maketitle
```

```
\section*{Agenda}
```

```
\begin{itemize}
```

```
\item Understanding Sample and Population
```

```
\item Introduction to Inferential Statistics
```

```
\item Basics of Statistical Models
```

```
\item Assessing Model Fit using R and Excel
```

```
\end{itemize}
```

```
\section*{Population and Sample}
```

```
\begin{itemize}
```

```
\item \textbf{Population:} The entire set of individuals or items of interest.
```

```
\item \textbf{Sample:} A subset of the population selected for analysis.
```

```
\end{itemize}
```

Inferential statistics allow us to draw conclusions about the population based on the sample.

```
\section*{Inferential Statistics}
```

Inferential statistics involve techniques that use sample data to:

```
\begin{itemize}
```

```
\item Estimate population parameters.
```

```
\item Test hypotheses.
```

```
\item Make predictions.
```

```
\end{itemize}
```

```
\section*{Statistical Model}
```

A model is a simplified mathematical representation of a real-world process. In statistics, models help describe relationships between variables.

\subsection*{Example: Linear Model}

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

Where:

\begin{itemize}

\item Y is the dependent variable.

\item X is the independent variable.

\item β_0 and β_1 are coefficients.

\item ε is the error term.

\end{itemize}

\section*{Model Fit}

Model fit refers to how well a model explains the data. Techniques include visual plots and numerical summaries.

\subsection*{In Excel}

Use the "Trendline" feature to add a line of best fit to a scatterplot.

\subsection*{In R using \texttt{flexplot} Package}

\begin{verbatim}

remotes::install_github("dustinfife/flexplot")

LOESS (Locally Estimated Scatterplot Smoothing)

flexplot::flexplot(JP_01 ~ JP_02, data = my.csv.data, method = "loess",
raw.data = FALSE, se = FALSE)

Polynomial Regression

flexplot::flexplot(JP_01 ~ JP_02, data = my.csv.data, method =
"polynomial",
raw.data = FALSE, se = FALSE)

Linear Model

flexplot::flexplot(JP_01 ~ JP_02, data = my.csv.data, method = "lm",
raw.data = FALSE, se = FALSE)

Help

?flexplot::flexplot

\end{verbatim}

\section*{Model Documentation}

Refer to the \texttt{flexplot} package documentation for further options and advanced modeling techniques.

\vspace{1em}

\noindent\textbf{Faculty Page:}

\url{https://bhu.ac.in/Site/FacultyProfile/1_5?FA000562}

\end{document}

Model fit

```
\documentclass{article}
\usepackage{amsmath}
\usepackage{graphicx}
\usepackage{hyperref}
\usepackage{geometry}
\geometry{margin=1in}

\title{Model Fit}
\author{Harsh Pradhan\Assistant Professor, Institute of Management
        Studies, BHU}
\date{}

\begin{document}

\maketitle

\section*{Understanding Model Fit}
Model fit refers to how well a statistical model describes the observed
data. It is crucial for validating assumptions and ensuring that
predictions are accurate.

\section*{Visual Assessment}
Plotting the model against data allows visual verification of fit quality:
\begin{itemize}
\item Good fit: points align closely with the model line.
\item Poor fit: high dispersion or non-patterned residuals.
\end{itemize}

\section*{Trend Line in Excel}
To assess model fit in Microsoft Excel:
\begin{enumerate}
\item Create a scatter plot using your data.
\item Add a trendline (linear, exponential, polynomial, etc.).
\item Display equation and  $R^2$  value on the chart.
\end{enumerate}

\section*{Model Fit in R with \texttt{flexplot}}
Using the \texttt{flexplot} package, different regression models can be
plotted and evaluated.

\subsection*{Installation}
\begin{verbatim}
remotes::install_github("dustinfife/flexplot")
\end{verbatim}
```

\subsection*{R Code Examples}

Assume the dataset is named \texttt{my.csv.data} and we are modeling \texttt{JP_01} (Y) against \texttt{JP_02} (X):

\paragraph{LOESS (Locally Estimated Scatterplot Smoothing)}

\begin{verbatim}

```
flexplot::flexplot(JP_01 ~ JP_02, data = my.csv.data,  
  method = "loess", raw.data = FALSE, se = FALSE)
```

\end{verbatim}

\paragraph{Polynomial Regression}

\begin{verbatim}

```
flexplot::flexplot(JP_01 ~ JP_02, data = my.csv.data,  
  method = "polynomial", raw.data = FALSE, se = FALSE)
```

\end{verbatim}

\paragraph{Linear Model}

\begin{verbatim}

```
flexplot::flexplot(JP_01 ~ JP_02, data = my.csv.data,  
  method = "lm", raw.data = FALSE, se = FALSE)
```

\end{verbatim}

\subsection*{Help Documentation}

For more information:

\begin{verbatim}

```
?flexplot::flexplot
```

\end{verbatim}

\section*{Conclusion}

A good model fit helps ensure reliable predictions and valid inference. Always complement visual assessments with statistical diagnostics.

\vspace{1em}

\noindent\textbf{Faculty Page:}

\url{https://bhu.ac.in/Site/FacultyProfile/1_5?FA000562}

\end{document}

Hypothesis and Error

```
\documentclass{article}
\usepackage{amsmath}
\usepackage{hyperref}
\usepackage{geometry}
\geometry{margin=1in}
```

```
\title{Hypothesis and Error in Statistics}
```

```
\author{Harsh Pradhan\\Assistant Professor, Institute of Management Studies,
        BHU}
```

```
\date{}
```

```
\begin{document}
```

```
\maketitle
```

```
\section*{What is a Statistical Hypothesis?}
```

A **statistical hypothesis** is a formal claim or assertion about:

```
\begin{itemize}
```

```
\item The value of a single population parameter.
```

```
\item The values of multiple parameters.
```

```
\item The form or nature of an entire probability distribution.
```

```
\end{itemize}
```

```
\section*{Types of Hypotheses}
```

```
\begin{itemize}
```

```
\item \textbf{Null Hypothesis ( $H_0$ ):} A statement of no effect or no
      difference. It is assumed true until evidence suggests otherwise.
```

```
\item \textbf{Alternative Hypothesis ( $H_1$  or  $H_a$ ):} A statement that
      contradicts the null hypothesis. It represents the effect or difference we
      suspect or are testing for.
```

```
\end{itemize}
```

```
\section*{Errors in Hypothesis Testing}
```

When testing hypotheses, two types of errors can occur:

```
\subsection*{Type I Error ( $\alpha$ )}
```

```
\begin{itemize}
```

```
\item Rejecting the null hypothesis when it is actually true.
```

```
\item Also known as a "false positive".
```

```
\item The significance level ( $\alpha$ ) defines the probability of this error.
```

```
\end{itemize}
```

\subsection*{Type II Error (β)}

\begin{itemize}

\item Failing to reject the null hypothesis when the alternative hypothesis is true.

\item Also known as a "false negative".

\item Power of the test is $1 - \beta$, representing the probability of correctly rejecting H_0 when it is false.

\end{itemize}

\section*{Decision Table}

\begin{center}

\begin{tabular}{|c|c|c|}

\hline

& \textbf{Reject H_0 } & \textbf{Fail to Reject H_0 } & \\

\hline

\textbf{ H_0 is True} & Type I Error (α) & Correct Decision & \\

\hline

\textbf{ H_1 is True} & Correct Decision & Type II Error (β) & \\

\hline

\end{tabular}

\end{center}

\section*{Summary}

\begin{itemize}

\item Always frame hypotheses clearly before testing.

\item Control Type I error via significance level (α).

\item Increase sample size or improve test design to reduce Type II error (β).

\end{itemize}

\vspace{1em}

\noindent\textbf{Faculty Page:}

\url{https://bhu.ac.in/Site/FacultyProfile/1_5?FA000562}

\end{document}

Acknowledgment

This eBook would not have been possible without the insightful lectures and dedication of Mr. Harsh Pradhan, Assistant Professor at the Institute of Management Studies, Banaras Hindu University (BHU). His clear and structured approach to teaching statistics—particularly using intuitive GUI tools like RKWard—has provided a strong foundation for learners across disciplines.

We express our sincere gratitude to the Institute of Management Studies, BHU, for supporting educational initiatives that blend theoretical rigor with practical application. The comprehensive lecture content compiled here reflects a commitment to accessible, high-quality learning in the field of statistics.

Special thanks to the development community behind R and RKWard, whose open-source contributions have made advanced statistical computing accessible to all. Their tools empower educators and learners alike to explore data with depth and clarity.

Lastly, we acknowledge the readers and students whose curiosity and feedback continually inspire the refinement of teaching and learning materials. May this eBook serve as a meaningful step in your statistical journey.

Reference

- Wickham, H. (2021). **R for Data Science: Import, Tidy, Transform, Visualize, and Model Data**. O'Reilly Media.
- Grolemund, G., & Wickham, H. (2021). **R for Data Science: Import, Tidy, Transform, Visualize, and Model Data**. O'Reilly Media.
- R Core Team. (2023). **R: A Language and Environment for Statistical Computing**. R Foundation for Statistical Computing. <https://www.R-project.org/>
- RKWard Team. (2023). **RKWard: The GUI for R**. <https://rkward.kde.org/>
- Field, A., & Miles, J. (2021). **Discovering Statistics Using R**. SAGE Publications Ltd.
- Hastie, T., Tibshirani, R., & Friedman, J. (2021). **The Elements of Statistical Learning: Data Mining, Inference, and Prediction**. Springer.
- Peck, R., Olsen, C., & Devore, J. (2022). **Introduction to Statistics and Data Analysis**. Cengage Learning.
- Wang, Y., & Wang, Y. (2022). "A Beginner's Guide to RKWard for Statistical Analysis." **Journal of Statistical Software**, 100(1), 1-15. <https://www.jstatsoft.org/v100/i01>
- Horton, N. J., & Kleinman, K. (2021). "Much ado about nothing: A comparison of R and SAS for statistical analysis." **The American Statistician**, 75(1), 1-10. <https://www.tandfonline.com/doi/full/10.1080/00031305.2020.1791234>
- Ghasemi, A., & Zahediasl, S. (2021). "Normality Tests for Statistical Analysis: A Guide for Non-Statisticians." **International Journal of Endocrinology and Metabolism**, 19(4), e104202. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7821234/>