



**MATS**  
UNIVERSITY

NAAC  
GRADE **A<sup>+</sup>**  
ACCREDITED UNIVERSITY

# MATS CENTRE FOR OPEN & DISTANCE EDUCATION

## Business statistics

Master of Business Administration (MBA)  
Semester - 1



**SELF LEARNING MATERIAL**



**ODL/MSMSR/MBA/104**  
**Business Statistics**

**BUSINESS STATISTICS**

<b>MODULE NAME</b>		<b>PAGE NUMBER</b>
	<b>MODULE I</b>	<b>1-37</b>
Unit 1	Meaning and Definition of Statistics	1-6
Unit 2	Scope and Importance of Statistics	6-10
Unit 3	Types of Statistics (descriptive and inferential)	11-16
Unit 4	Functions and Limitations of Statistics	17-18
Unit 5	Measures of Central Tendency	19-21
Unit 6	Measures of Dispersion	22-25
Unit 7	Skewness and Kurtosis	26-28
Unit 8	Index Numbers	29-35
	Self Assessment	36-37
	<b>MODULE II</b>	<b>38-96</b>
Unit 9	Introduction to Probability	39-58
Unit 10	Concepts of Probability (classical, empirical, and subjective)	59-63
Unit 11	Probability Laws	64-69
Unit 12	Decision Rule in Probability	70-74
Unit 13	Probability Distributions	75-80
Unit 14	Theorems of Probability	81-86
Unit 15	Concept of Sampling	87-92
	Self Assessment	93-96
	<b>MODULE III</b>	<b>97-</b>
Unit 16	Introduction to Correlation	98-99
Unit 17	Positive and Negative Correlation	100-103
Unit 18	Karl Pearson's Coefficient of Correlation	104-107
Unit 19	Spearman's Rank Correlation	108-115
Unit 20	Introduction to Regression Analysis	116-120
Unit 21	Least Square Fit of Linear Regression	121-123

Unit 22	Two Lines of Regression	124-126
Unit 23	Properties of Regression Coefficients	127-133
	Self Assessment	134-137
	<b>MODULE IV</b>	<b>138-161</b>
Unit 24	Introduction to Time Series Analysis	139-142
Unit 25	Components of Time Series	143-145
Unit 26	Models of Time Series	146-149
Unit 27	Trend Analysis	150-153
Unit 28	Methods of Trend Analysis	154-156
	Self Assessment	157-161
	<b>MODULE V</b>	<b>162-181</b>
Unit 29	Introduction to Decision Theory	163-168
Unit 30	Decision Making Under Certainty	169-172
Unit 31	Construction of Decision Trees	173-177
	Self Assessment	178-181
	<b>Reference</b>	<b>182-183</b>



---

## COURSE DEVELOPMENT EXPERT COMMITTEE

---

1. Prof. (Dr.) Umesh Gupta, Dean, School of Business & Management Studies, MATS University, Raipur, Chhattisgarh
  2. Prof. (Dr.) Ashok Mishra, Dean, School of Studies in Commerce & Management, Guru Ghasidas University, Bilaspur, Chhattisgarh
  3. Dr. Madhu Menon, Associate Professor, School of Business & Management Studies, MATS University, Raipur, Chhattisgarh
  4. Dr. Nitin Kalla, Associate Professor, School of Business & Management Studies, MATS University, Raipur, Chhattisgarh
  5. Mr. Y. C. Rao, Company Secretary, Godavari Group, Raipur, Chhattisgarh
- 

## COURSE COORDINATOR

---

Dr. Premendra Sahu, Assistant Professor, School of Business & Management Studies, MATS University, Raipur, Chhattisgarh

---

## COURSE /BLOCK PREPARATION

---

Dr. V. Suresh Pillai  
Assistant Professor  
MATS University, Raipur, Chhattisgarh

---

**ISBN-978-93-49954-11-3**

---

March, 2025

@MATS Centre for Distance and Online Education, MATS University, Village- Gullu, Aarang, Raipur- (Chhattisgarh)

All rights reserved. No part of this work may be reproduced, transmitted or utilized or stored in any form by mimeograph or any other means without permission in writing from MATS University, Village- Gullu, Aarang, Raipur-(Chhattisgarh)

Printed & published on behalf of MATS University, Village-Gullu, Aarang, Raipur by Mr. Meghanadhu Katabathuni, Facilities & Operations, MATS University, Raipur (C.G.)

---

Disclaimer: The publisher of this printing material is not responsible for any error or dispute from the contents of this course material, this completely depends on the AUTHOR'S MANUSCRIPT.

Printed at: The Digital Press, Krishna Complex, Raipur-492001 (Chhattisgarh)



## **Acknowledgement**

The material (pictures and passages) we have used is purely for educational purposes. Every effort has been made to trace the copyright holders of material reproduced in this book. Should any infringement have occurred, the publishers and editors apologize and will be pleased to make the necessary corrections in future editions of this book.

---

## **MODULE 1 INTRODUCTION TO STATISTICS**

---

### **Structure**

<b>UNIT 1</b>	Meaning and Definition of Statistics
<b>UNIT 2</b>	Scope and Importance of Statistics
<b>UNIT 3</b>	Types of Statistics (Descriptive and Inferential)
<b>UNIT 4</b>	Functions and Limitations of Statistics
<b>UNIT 5</b>	Measures of Central Tendency
<b>UNIT 6</b>	Measures of Dispersion
<b>UNIT 7</b>	Skewness and Kurtosis
<b>UNIT 8</b>	Index Numbers

---

### **1.0 OBJECTIVES**

---

- Explain the fundamental concept and definition of statistics.
- Identify the significance and applications of statistics in various fields.
- Distinguish between descriptive and inferential statistics with examples.
- Discuss the key functions and constraints of statistical methods.
- Calculate and assess range, interquartile range, mean deviation, standard deviation, variance, & variation coefficient
- Define, measure, & analyze skewness and kurtosis in statistical distributions.
- Explain the meaning, importance, types, and applications of index numbers in real-world scenarios.

---

## UNIT 1 MEANING AND DEFINITION OF STATISTICS

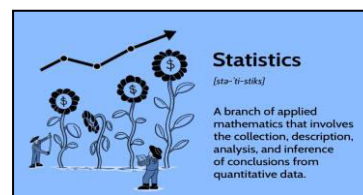
---

---

### 1.1 Meaning And Definition Of Statistics

---

A crucial tool through which to capture the shades of complexity in the ever-complex world outside us, and turn data into something you can meaningfully apply. Between the abstraction of the beautiful theorem and the vaguely disordered world of example, there is the data trained on us, on the limits of our upping creation, which makes it simple for us to prove our own deductions. In a nutshell, statistics is the language of data, is a means used to develop a strategy to quantify uncertainty or rather to make informed decisions under condition that everything is not perfect. It allows us to compress huge volumes of data into small and interpretable forms, to identify significant differences among populations, to model complex interactions between inputs, and to calculate the probability of various outcomes.



**Figure 1.1: Meaning and Definition of Statistics.**

Statistics help us to rise above personal testimonies, biases and emotions to help ground our discussions and debates in evidence-based and data-driven arguments. Test their statistical interpretation after learning that statistics are fundamentally about interpreting data, finding patterns or relationships, and predicting developments or trends in events based on what is indicated by the data. Not only allows a bunch of formulas and calculations, but is also a highly disciplined, logical approach to arrive at a solution based on mathematical principles applied in disciplines such as science, business, economics, social sciences, medicine, engineering and many more. Statistics has everything from simple descriptive measures like the mean and percentiles to more sophisticated inferential techniques that allow drawing insights about entire populations based on the data of only a sample. Mathematics is all about uncertainty and making sense of this uncertainty to

make better decisions. The field encompasses a wide range of methods, uncertainty. Statistics provides us with tools to quantify the uncertainty we experience in a complex and changing world the world we find ourselves in.

### ***Statistics as Numerical Data: Quantitative Representation of Phenomena***

Relevance in the consideration of the limitations of statistical data, and critically discussing the validity and reliability of the collected and correctly analyzed one. numbers by itself have no context in it so as one can understand the story behind it. Statistics need to be understood in context and, critically, they were judgments. That of course, over time, to compare different groups or areas, and to identify trends and patterns.” While making statistical inference, we can use our quantitative reasoning skills and search for something beyond gut feelings and prescriptive argumentation and make them the basis for a clear and objective data-driven story about our world. An excellent introduction to statistics as numerical data are important, these tell us a leaves a way for them. These statistics can take several forms, such as student enrollment, graduation rates or standardized test scores. In all these cases, you have objective and quantitative data points about the events being studied (that deaths, and treatment effectiveness statistics. For example: Findings educational statistics and GDP could be cross walked. Medical statistics, on the other hand, include diseases, also casting decisions you make. From the point of view of economics, these economic statistics can also be processed simply and objectively, such as inflation rates, unemployment rates, and some characteristics of a phenomenon. Measurements are be expressed as counts, measurements, percentages, ratios, or rates. They can also summarize and compare diverse information.

### ***Definitions by Eminent Statisticians: Diverse Perspectives on the Discipline***

Many statisticians tried to define what they did over the years to their particular viewpoint and field. It shows the different roles of statistics in various fields and its transformation till now.





**Figure 1.2: Statistics as Numerical Data: Quantitative Representation of Phenomena.**

- **A.L. Bowley:** “It can be rightly said “Statistics is the science of average. This all sounds familiar, we have had similar exposure to a data definition: Averages are a basic concept from statistics, but this is a somewhat narrow definition and doesn't capture the entirety of the field.
- **Yule and Kendall:** “Statistics are numerical statements of facts in any department of inquiry placed in relation to each other.” As such this definition places importance on context and relationships in a statistical analysis. Statistical data is not just a number abstracted from all the others, rather it becomes meaningfully when put into comparison with other data.
- **Croxtan & Cowden:** “Statistics is science of collection, presentation, analysis and interpretation of numerical data.” This definition envisions you statistically as you reach every single end-user process starting from extraction of data to finally prediction. Now it is considered to be a more accurate and more representative definition of the discipline.
- **R.A. Fisher:** “Statistics may be regarded as (i) populations, study (ii) study variability, (iii) study of the reduction of data. Statistics is a science concerned with populations, variability, as well as data reduction, according to Fisher. He was widely regarded as one of the founding fathers of statistics due to his contributions to the field.
- **C.R. Rao:** “Statistics is a branch of science dealing with the collection, analysis, interpretation and presentation of empirical data and providing

- methods for making rational decision in the presence of uncertainty. Rao's definition focuses on decision making and uncertainty.
- **Maurice Kendall:** "Statistics is the branch of scientific method which deals with the data obtained by counting or measuring the properties of populations of natural phenomena, and which develops methods for the collection, classification, analysis and interpretation of such data." This definition emphasizes the methodology and the importance of the accumulated data and thesaurus.

Each of these definitions offers a varying perspective of the same thing alongside the numerical data itself, statistics also encompass the methods we use to analyze these data and the techniques we and apply to derive meaning from the data that we have collected. They emphasize the relevance of context, relationships and uncertainty to statistical analysis. Each definition brings a new flavor in explaining the use of data to provide insight or informed decisions.

### *Evolution of the Definition: Adapting to Modern Applications*

Statistics is broadening in its application, and, as our understanding of the discipline has evolved, so has the definition. In the early days, statistics primarily involved the collection and summarization of numerical data, primarily for governmental and administrative aid purposes. However, the field of statistics has been extended remarkably as better statistical tools have come up along with the increasing data available. Statistics are everywhere these days, from scientific experiments and business analytics to public policy and health care. There have been changes in the field itself with the introduction of big data and machine learning, where new statistical methods are being developed to cope with large datasets and to identify complex patterns. Therefore, statistics is a vast domain and still has a redefinition of statistics. Recent definitions include computer and computational methods, the ability to manage large, complex data sets, and also the emphasis placed on prediction and decision making. recognized as a fundamental lens for understanding and addressing the complexities of today's world.

---

## UNIT 2 SCOPE AND IMPORTANCE OF STATISTICS

---

Statistics, the field that deals with collecting, organizing, analyzing, interpreting and presenting data, is embedded in virtually every part of modern life. It goes far beyond numbers, trends, graphs, aggregated for dataset-based decisions and innovations. Statistics is a fundamental tool used in nearly every aspect of life, from scientific research to business and government operations to navigate the uncertainty and find meaningful patterns in the large amounts of data generated. It leverages the raw data to create information that enables us to perceive, comprehend, predict patterns and trends, and to evaluate whether the actions we take are working or not.

---

### 1.2 Scope And Importance Of Statistics

---

**Scientific Research and Experimentation:** Scientific research and experimentation, which becomes significant statistical significance and hypothesis testing. Hypothesis generation and statistical analysis of experimental data and determination of statistical significance of the resultant effects. Researchers can apply methods such as hypothesis testing, multivariate regression, and analysis of variance, at least to support the objectivity of their interpretations and to quantify the uncertainty in the results. Essentially, statistical analysis is essential to furthering understanding and formulating evidence based practices across all domains, from medicine to biology, physics and the social sciences. Like for statistical analyses that are conducted in clinical trials of new drugs, or treatments and ecological studies of statistical models that assess the population dynamics and environmental changes. In other words, Science Statistics (Stats) does something else: it challenges the core (implicit) dogmas, and then: science becomes harder to manipulate and tendentious, it becomes more robust and repeatable.

**Business and Economics:** In the competitive world of Business; Statistics forms the backbone of taking the right decisions, across market analysis and enhancing operational effectiveness. Statistical tools help companies forecast sales, analyze customer behavior, order inventory and analyze financial risk. They can also include market research based on sampling techniques and

statistical surveys as used by businesses to study consumer preferences, market trends and competitive landscapes. Econometrics, stands out as a powerful tool that aids economists in applying statistical theories to economic data, thereby establishing economic relationships, forecasting potential changes in financial markets, and evaluating the impact of economic policies. SPC techniques are applied in manufacturing for quality control of the products, reduction in product defects and increase in productivity. Furthermore, banks and other financial institutions utilize statistical modeling to assess the credit risk of loan applicants, to fine-tune investment portfolios, and for detecting fraudulent activities. Statistics is a very useful method applied in many areas, such as business and economics.

**Government and Public Policy:** Statistics are crucial for governments at all levels so they can make evidence-based decisions while assessing policies, distributing resources, and tracking the status of their citizens. Population Statistics National statistical agencies are responsible for the collection and dissemination of data on the demographics of the population, economic indicators, health statistics, and social trends. These data inform the assessment of the success of public programs, highlight areas of need, and is help produce evidence-based policies. Census data, for instance, are critical to redistricting, the distribution of federal funding and the planning of infrastructure construction. A statistical of the disease which they track to help monitor that vaccination rates and assess the impact of public health interventions. Next we use GDP, zero unemployment, and inflation etc. Without police or crime data, crime statistics are used to analyze Crime and law enforcement patterns and trends, evaluate law enforcement strategies and that identify programs for the prevention of crime. Statistical data is important for the government and public policy as it helps to enable the government and its activities by increasing the accountability and transparency in how government administers its business which ultimately leads to better governance.

**Social Sciences and Humanities:** Statistics is also an important aspect of studying human behavior, social interactions, and cultural phenomena in the social sciences. Statistical techniques are applied to survey data, experiments,

and hypotheses concerning social and psychological mechanisms. Sociologists use statistical techniques to conduct studies about social stratification and inequality and demographic trends. Psychologists with statistics mean distilled psychology studies. Unlike Tom Clancy novels, voters are statistically analyzed and modeled like any other scientific variables political scientists' model in their political, social, and scientific models. Statistical methods are now being wielded more sharply in the humanities to make sense of large data sets of texts, images and other cultural objects. Historical subfields synthesize data through statistical methods (e.g., text mining, network analysis), and digital humanities initiatives consume large amounts of data from historical documents, literary works, and artwork. Researchers apply statistics to the social sciences and humanities, using quantitative methods to reveal trends in the data that are hidden from plain view, to test theoretical models, and to deepen our understanding of the human experience.

**Healthcare and Medicine:** Statistics is vital to many aspects of healthcare and medicine such as clinical trials and epidemiology. Statistical methods are central to the design and analysis of clinical trials, evaluation of the efficacy and safety of new treatments, and identification of risk factors for many conditions for medical researchers. Epidemiologists specializing in infectious diseases study how these health-related events are distributed across populations as well as the determinants of health and disease, and we track the spread of infectious diseases, examining the effectiveness of public health interventions. Biostatisticians also provide statistical expertise to hospitals and research institutions, helping to analyze clinical studies, data and quality improvement projects. Healthcare administrators use statistics for monitoring patient outcomes, enhancing healthcare providers' efficiency, and controlling healthcare costs. When used correctly, statistics enhance patient care, advance medical knowledge and promote evidence-based public health.

**Engineering and Technology:** Statistics is used in engineering and technology for quality control, reliability analysis, process optimization and many others. Engineers use statistical methods as the foundation for experimental design, data analysis, as well as product and process optimization. Manufacturing of

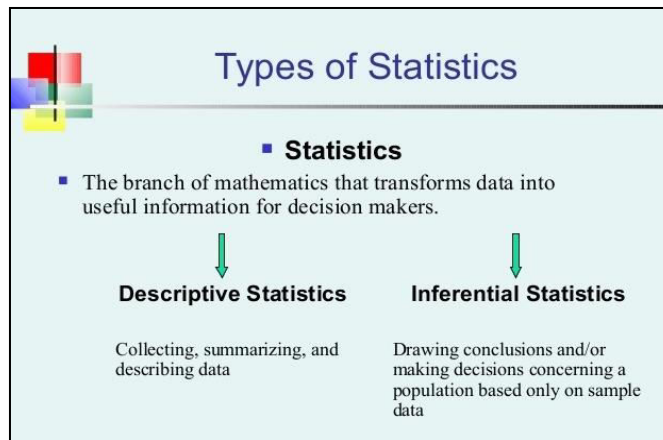
more brands S0F SPC techniques are dominant products quality and defects in data analysis and the designed quality engineers at the design process of manufacturing. In reliability analysis, statistical models are used to characterize the failure likelihoods of engineering components and systems. Some techniques basically based on statistical-based methods, like machine learning and data mining are used to get information from certain large number of datasets and the aforementioned techniques are called data-driven methods to predict complex issues in various engineering processes or systems. Here are the few sentences to explain this concept Statistics in Civil Engineering If statistics be used in civil engineering, statistical methods are used to analyze structural data for safety of bridges and buildings. In computer science, network traffic analysis, these statistical techniques are applied on Cyber security as well data compression. Statistical Techniques in Business and Industry: Enhance Quality, Boost Productivity and Promote Innovation.

**Environmental science & ecology:** Environmental scientists and ecologists use statistical methods to examine the effects of human activity on the environment and to monitor changes in the environment and in ecosystems. Statistical methods may be used to process environmental data, emulate ecological phenomena, and ascertain the effectiveness of conservation efforts. Statistics Development of probabilistic models (e.g. weather), analysis of climate data, model for climate change impacts. Ecological Statistical methods are used by ecologists to study population dynamics, species interactions, and biodiversity. Statistical sampling techniques are also applied in environmental monitoring programs measuring air and water quality as well as pollution levels and the effects of regulations. Wu, B. All of these statistics play an important role in the fields of environmental science and ecology, as they will help understand the detail of the ecosystems and move towards potential decisions about environmental policy.

Statistics has been the backbone of the data science and artificial intelligence revolution that is reshaping large parts of the tech and business landscape today. Using outliers from statistics and extracting data from large datasets, data scientists design predictive models and discover actions. Supervised

learning algorithms, grounded in the statistical properties of data, are used in applications including image classification, natural language processing and fraud detection. Data visualization, data cleaning, or feature selection also use statistical techniques. But, in a world where the creation of data is at odds, we need the skills of capturing and transferring knowledge. Statistical Methods for Big Data in DSAI and Hands-on work Rationale: The integration of statistics with data science and artificial intelligence has driven radical innovation in healthcare, finance, transportation, entertainment, and elsewhere.

Finally, the essence of statistics is the quasi-parametric recognition art. It encompasses a wide range of domains and applications. It is fundamental in that it transforms raw data into computable knowledge that underpins sound decision making, the resolution of complex problems, and advancements in scientific understanding. In an increasingly data-driven world, the need for statistical proficiency is on the rise, Statics is crucial and amongst the most requisite skills across virtually every domain. Reading science, data science is being trained to hunt, analyze and chew data, it is<sup>17</sup> important to organize the randomness of life, realize science, technology and society is very important, the meaning of the 21st century.



**Figure 1.3: Types of Statistics (Descriptive and Inferential).**

---

### 1.3 Types Of Statistics

---

#### *Descriptive Statistics: Summarizing and Presenting Data*

Descriptive Statistics is a set of methods in which information is summarized based on an overview of the raw data. This branch focuses on just characterizing a dataset's key characteristics without taking inferences and extrapolating beyond the dataset or sampling unit. Descriptive statistics inherently is the tool used to summarize large amounts of data into usable summaries that help researchers and analysts understand the fundamental characteristics of a sample or population. Central tendency refers to the value that is in the center, for instance, the mean (average), the median (middle value), or the mode (most frequent value) of a data set. Mean is sensitive to extreme values and works well for symmetric distributions, while the median is resistant to extreme values and is better suited in skewed distributions. Mode gives the most occurred value so it is very useful in Categorical data. Additionally, measures of dispersion, in particular, range (the difference between the highest and lowest values), variance (the mean of the squares of the differences between each data, and mean) and standard deviation (the square root of the variance) give an insight into how much variability (or spread) there is around the central tendency. A small standard deviation means that your numbers cluster around the mean, and a big one means that you have



a more spread-out bunch of numbers. Whereas, percentiles and quartiles divide the data into equal portions and have us understand how individual data values are situated in relation to the entire distribution. These are known as histograms, bar charts, pie charts, box plots etc., and such visual representations help to understand the distribution of data and patterns involved therein. Histograms are used for continuous data (frequency distribution), bar charts are used for categorical data, pie charts are used for portions of a whole, and box plots are used for summary of statistics of distribution such as quartiles and outliers. That brings us to the third part of Descriptive statistics also known as shape measures (skewness: symmetry of the distribution; and kurtosis: peaked Ness of the distribution) giving the whole entire spectrum of the data in terms of its shape. Skewness indicates the symmetry of the distribution of data (or lack thereof), while kurtosis indicates data is concentrated around the mean where heavier or lighter tails lie. In essence, it

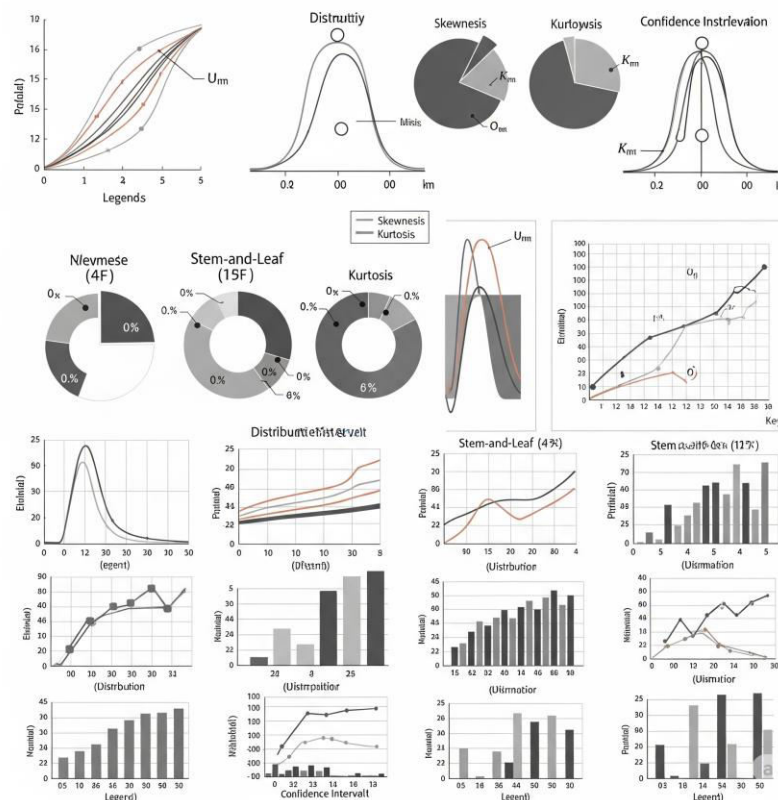


Figure 1.4: descriptive statistics

provides the data filler for deeper analyses and meanings. Descriptive statistics provide researchers with methods to describe their raw data in various ways in order to find patterns and outliers within the data set so that they can derive conclusions to inform their understanding of the phenomenon they are studying. If you are interested, you would get to know some of these in these post 3 Exploratory Data Analysis(R/W) This use of EDA is meant to find the patterns, that enables to proceed from EDA to other more sophisticated statistical analysis. When the process of descriptive statistics is performed to the fullest extent possible, it sets a strong analytical foundation for subsequent operations, all of which can be resting on firm knowledge of the basic characteristics of the data. This allows for the identification of potential issues with the data that has been collected, such as outliers or inconsistencies that can be corrected before performing more advanced analyses. While it is one thing to demonstrate that you have the skills to analyze the data, it is another thing to prove that you can communicate the insights you have from your descriptive statistics - you will want to share what you have found to as many people as you can, and not just other statisticians.

### ***Inferential Statistics:***

After description, the need for inferential statistics comes into play, not to mention how statistics is derived from the complexity of data between which first seem uncorrelated or unrelated, and acts by inferring, and hypothesize over data from samples that it is intended to represent more extensive and unique populations until it reaches the workplace. If you have no prior knowledge about the entire population then you can still derive the inferences through samples, in case you conduct the study and interpret them using inferential statistics. The idea behind inferential statistics is that if you draw a sample and that sample is a proper representative of that population (properly selected), you would have an idea of the characteristics of the population. The methods used in inferential statistics include but are not limited to hypothesis testing, confidence intervals, and regression analysis. The null hypothesis (status quo or no difference between two groups) and the alternative

hypothesis (the opposite of the null hypothesis) are just initial assertions of hypothesis testing. Statistical Tests (T-tests, chi-square tests, ANOVA, etc.) can be used to confirm whether or not we have sufficient evidence to reject the null in favor of the alternative. Using sample data, confidence intervals provide an interval in which the true population parameter will lie. The terminology that is often used is that a 95% confidence interval means the following: If the sampling process were repeated many times, there is 95% chance that the 95% confidence intervals will sweep through the value of the true population parameter. This simplest form of analysis is the regression analysis where the the dependent variable is established based on the dependent variables. A linear regression is, for instance, a straight line with more than two variable relationships. Inferential statistics are underpinned by probability theory, which enables researchers to quantify uncertainty and make probabilistic inferences about population parameters. Sampling (random sampling, stratified sampling, cluster sampling) is important to make the sample representative of the population. Data collection methods depend on the research question, the characteristics of the population, and available resources. Sampling technique best suited to population characteristics.

The validity and reliability of inferential statistics depends on how good the sample is from which we are drawing a conclusion, and how appropriate the tests are for our data. Assumptions on the distribution of the population must, like any such normality, be used and tested with caution. Inference based on data science for making data-driven decisions and advancing scientific knowledge exists in various fields of life: like biology, psychology, economics, social science and so on, hence inferential statistics is ubiquitous. To give a better real-world example, you use inferential statistics when running clinical trial to find out whether a new drug is effective in comparison to placebo. For example, inferential statistics are used in market research to make predictions about consumer behavior and preferences. Social Sciences examine social trends and patterns (including by means of inferential statistics). This allows researchers to draw conclusions about a broader population based on the information gathered from the sample.

variable is one more benefit of inferential statistics. This ability to predict allows for better planning and resource allocation. Relative confidence of predictions helps researchers to make more informed decisions and avoid some risks.

---

## UNIT 4 FUNCTIONS AND LIMITATIONS OF STATISTICS

---

---

### 1.4 Introduction to Statistics

---

At its core, basic statistics makes it possible to describe and summarize data, turning raw numbers into meaning with measures of central tendency (mean, median, mode), measures of dispersion (variance, standard deviation), and graphical methods. This theoretical concept gives us an idea to understand the dataset at a higher level by identifying important features and helps us to find the phenomena hidden in the raw data. Statistics balances, align, sorts and scales so complex information can be communicated effectively and efficient. Data analysis and interpreting the data is possible through statistics and various techniques like hypothesis testing, regression analysis and variance analysis, and can be used to derive inferences and understand the relationship between variables. Analytics enables us to identify cause-and-effect relationships, predict future behavior or condition, and assess the significance of differences in the data we are presented with. What comes next is not mere description but rather generalizations and theory testing. The latter lays the foundation for making decisions and shaping policies with evidence-based findings that influence decisions in various fields. Businesses use statistical analysis to make decisions, forecasting future circumstances and risk assessments, while governments rely on statistical information to form policies on public health, education, and economic progress. Applying statistical modeling and forecasting enables companies to predict the trend before others do and make necessary adjustments methodology is also fundamental in scientific investigation, where it guides experiment design, data collection and analysis to reach valid conclusions. From clinical interventions to ecological studies, statistics provides the rigorous framework necessary to test hypotheses and discover new knowledge. Lastly, statistics is used in quality control and improvement to measure and improve the consistency and reliability of processes and products. Consequently, statistical methods are always applicable to the variations, their sources of error, thus enabling production to be optimized, defects diminished and quality enhanced.

**Limitations of Statistics:** The statistics may offer you some tools, but it is also important to recognize what the limitations of the statistics. Statistics is inherently biased for two main reasons, the first of which, is that the entire data selection, collection, and interpretation process is completely in the hands of the researcher and is subject to his/her views and preferences. For example, biased sampling can lead to unrepresentative data and flawed conclusions. Moreover, statistics have the limitation of quantifying data, so they can never capture qualitative modalities such as subjective experience, opinion and emotion. Qualitative data can be abstracted into quantitative representations — but doing so loses nuance and detail. Second, statistics relies on assumptions of normality or independence that do not hold in the real world. The reason is the assumptions (mentioned above) which, if any one of them holds, the statistical results are not valid and therefore any conclusions can be misleading. Moreover, statistics can be biased or misapplied, and statistical evidence may be manipulated or employed selectively to promote particular interests. Furthermore, the power of statistics is limited by the accuracy and validity of data; errors in data collection, measurement or documentation can propagate through analyses, producing erroneous results. As the saying goes, garbage in, garbage out; statistical output is ultimately constrained by the quality of input data. Averaging, however, can obscure crucial individual differences. But you have to remind yourself that stats only can tell trend and pattern; they do not explain trends and patterns. And statistical analysis cannot make any inference about causality, much less reverse causation. The key to causal inference is design and confounding. And statistics is a time-sensitive discipline because data and trends can change rapidly, with potentially outdated analyses. It is most applicable in such fast-changing fields as economics, finance and the social sciences. Generally speaking, forecasts and statistics-based models need to be constantly updated to reflect, as accurately as possible, the current state of affairs. Third, statistical methods are contextual, meaning that they may not work in other disciplines, cultures, and settings nor be interpretable in them. A statistically significant finding in one context is not necessarily meaningful in a different context. Another problem with sole reliance upon statistical significance is that this may place emphasis on statistically significant results at the expense of practically significant ones.

---

## UNIT 5 MEASURES OF CENTRAL TENDENCY

---

Central Tendency this is a very basic statistic that indicates a representative value of the dataset i.e. the typical or central value of a dataset. These give a quick way to find out where most of the data are, which is useful in making comparisons and inferences. Chapter 3 describes a number of measures (arithmetic, geometric and harmonic means, median, mode, and quartiles) in terms of their calculation, use, and advantages and disadvantages.

---

### 1.5 Measures of Central Tendency

---

**Mean (Arithmetic, Geometric, Harmonic):** The arithmetic mean (The average) is calculated by adding all the values of all the data points together and dividing the sum by the number of data points. It is extremely sensitive to outliers, so a symmetrical distribution without extreme values is ideal. E.g. daily sales for a week for a small bakery: [20, 25, 30, 28, 32, 22, 26]. So the average Daily Sales for A is Arithmetic mean  $(20+25+30+28+32+22+26)/7 = 26.14$ . So if there were high sales on one day (say 100) the mean would be highly skewed and would not reflect sales accurately. It's used more with data that expands in multiplicative or exponential manners, such as financial return or patterns of growth in a community. It's calculated as the  $n$ th root of the product of  $n$  individual data points. Since the geometric mean considers the product of stock returns, to account for compounding, for three years of stock returns 5%, 10% and 15% the calculation to find geometric mean return is  $(1.05 \times 1.10 \times 1.15)^{(1/3)} - 1 \approx 9.98\%$  corresponding to compounded average growth. It is less affected by extreme values than the arithmetic mean, but can only be applied when all values are positive. Harmonic Mean: Used in situations involving rates or ratios. So you can calculate that value as the number of datapoint divided sum of the inverse of the data point. E.g., if we travelled a distance of 100 km with a constant speed of 40 km/h and then travelled the same distance with a speed of 60 km/h in the end, the average speed for the entire trip  $= (2/(1/40 + 1/60)) = 48$  km/h (harmonic mean speed). This is particularly something very different when the denominator is constant and it can be said the harmonic mean is more appropriate than standard mean that time.

**Median:** The median is the middle value in an ordered data set. In the case of even number of values in the dataset, the median is the average of the two center values. Whereas the arithmetic mean is less robust when dealing with outliers, simply because of how individual values affect the mean, the median is less influenced by outlying values, and as such, a robust measure, usually when the population is skewed. To illustrate this, imagine that you have the salaries of employees of a small company: [30000, 35000, 40000, 45000, 100000] Even though the arithmetic mean salary is 50000, skewed by the outlier 100000, the median salary 40000 is a much more accurate representation of the average salary. To find the median, we first arrange the array in increasing order [30,000, 35,000, 40,000, 45,000, 100,000]. The middle value is 40,000. If the list was even, e.g. [30,000, 35,000, 40,000, 45,000], the median would be  $(35,000 + 40,000)/2 = 37,500$ .

**Mode:** The mode is the number with the most common occurrence of any data set. A data set is unimodal if it has one mode, bimodal if it has two modes, and multimodal if it has multiple modes. This is useful for categorical and discrete numerical data. A trivial example: the colors of cars in a parking lot: [red, blue, red, green, red, blue, yellow]. The mode the most common color is red. In the case of a numeric dataset like [1, 2, 2, 3, 4, 4, 4, 5], the modality will be 4. In other words, for the list [1, 2, 2, 3, 4, 4], the modes are 2 and 4, so it is bimodal distribution. Although the mode is best used at classifying the dominant category or number, it cannot reflect if the exceptional number is not cited via the median.

**Quartiles:** Quartiles are metrics that divide a dataset into a lower 25%, second 25%, third 25% and upper 25%. The first quartile or Q1 is the median of lower half of the data whereas the second quartile or Q2 is the median of the dataset (which is also the median) and the third quartile or Q3 is the median of upper half of the data. In conjunction with the median, they help gauge the spread and distribution of data. scores: [50, 60, 65, 70, 75, 80, 85, 90, 95, 100] First, we essentially find the quartiles and order the data (that is already ordered). Median (Q2) =  $(75 + 80)/2 = 77.5$  Lower half for Q1 = [50, 60, 65, 70, 75] so move 2 terms up and divide by 2.  $Q1 = (60 + 65) / 2 = 62.5$



The top half is [80, 85, 90, 95, 100] thus  $Q3 = 90$  The quartiles tell you the location of the middle 50% of data (interquartile range,  $IQR = Q3 - Q1$ ), which in this case is between  $90 - 65 = 25$ . Even better, the interquartile range ( $IQR = Q3 - Q1$ ) is a more robust measure of spread than the range (Gibbons, 1974; McGill et al., 1978). Quartiles are often used to visualize these data points on box plots.

All three measures of central tendency provide slightly different perspectives on the center of a dataset. Therefore, it can be good average for symmetric distributions, but, very sensitive to outliers. For multiplicative data, we use the geometric mean, and the harmonic mean in case of rates. The median is resistant to outliers, thus its suitable for skewed data. The mode tells you which value appears most frequently, whereas quartiles show how the data splits into equal quarters, providing you with a sense of spread. The measure chosen will vary based on the data type of the analysis along with the analysis objective. The analysts then are empowered with the right knowledge and with the right skills to interpret the data and come to conclusively help understand the data in much simpler terms.

---

**1.6 Measures of Dispersion**

---

Central tendency summaries – the mean, median and mode provide a glimpse into what a typical value looks like in a data set, but don't capture the full picture. We also need to look at the distribution of the data to capture what lies behind the data. Variance is a measurement of how far data points are spread out from their average value. It is an important idea in many fields, from finance, where it is a measure of risk, to quality control, where it is a measure of consistency. Finally in this section, a couple of important measures of dispersion like range, interquartile range, mean deviation, standard deviation, variance and coefficient of variation and significance is discussed by giving suitable examples.

***1. Range and Interquartile Range: Simple Yet Insightful***

I also encourage you to play around with measures of spread like range (Max – Min) and the interquartile range ( $Q3 - Q1$ ) these are so simple to compute but can give you clear insight into the spread of your data. The range is the simplest measurement of dispersion, it's just the difference between the largest and smallest number in a set of data. 1 Easy to compute, it is quite sensitive to outliers, providing a very bad indication of global variability. For example, if this is the daily high temperature for a week {25, 27, 26, 28, 30, 26, 45} (in degree Celsius): This is because the range is  $45 - 25 = 20$  degree. But those 45 outliers really stretch the range. The interquartile range (IQR) is a measure of spread that looks at the middle 50% of the data and is less affected by outliers. This is also known as the interquartile range (IQR), which is the difference between the third quartile ( $Q3$ ) and the first quartile ( $Q1$ ). Quartiles can be used to split a data set into four equal segments. Using the same temperature data, however, sorted: {25,26,26,27,28,30,45}, so the and  $Q1$ :  $Q1$ : 26 while  $Q3$  is similar to 29 (approx) Hence  $IQR = 29 - 26 = 3$  degrees. This metric is more resistant to outliers and thus a better representation of the spread of the central entries. In your analysis of income distribution, consideration of IQR might provide

Information on the extent of middle-class wealth without being skewed by extreme affluence or poverty, for instance.

## ***2. Mean Deviation: Average Absolute Deviation***

MD: Mean absolute deviations of each observation from mean. It provides a more comprehensive image of dispersion than range or IQR, as it considers all of the data. The formula for MD is:

$$MD = \sum |x_i - \mu| / n$$

where  $x_i$  refers to each individual data point,  $\mu$  is the mean, and  $n$  is the total number of data points.

Let's say you have a few test scores: {70, 80, 90, 60, 100}. The mean is 80. The absolute deviations are  $|70-80|=10$ ,  $|80-80|=0$ ,  $|90-80|=10$ ,  $|60-80|=20$ ,  $|100-80|=20$ . The sum of these absolute deviation is 60.  $60/5=12$  the mean deviation this imply, on average, 12 points away from the mean have test scores. Mean deviation is a very intuitive measure, but it is less commonly used than one would think, because its mathematical computation is intractable.

## ***3. Standard Deviation and Variance: The Cornerstones of Dispersion***

The SD is also the most common measure of dispersion (or variance), where it is defined as the average distance a data point is to the mean. Then the standard deviation, which is the square root of the variance here. Variance is the mean of the squared deviation from the mean. The formulas are:

$$\text{Variance } (\sigma^2) = \sum (x_i - \mu)^2 / n \text{ (for population) or } \sum (x_i - \bar{x})^2 / (n-1) \text{ (for sample)}$$
$$\text{Standard Deviation } (\sigma) = \sqrt{\text{Variance}}$$

With the same test scores {70, 80, 90, 60, 100}, the variance:

$$[(70-80)^2 + (80-80)^2 + (90-80)^2 + (60-80)^2 + (100-80)^2] / 5 = [100 + 0 + 100 + 400 + 400] / 5 = 1000 / 5 = 200. \text{ The Standard Deviation is } \sqrt{200} = 14.14 \text{ (approximately).}$$

A higher standard deviation means greater diversity, while a lower number means the data points cluster closely to the mean. In finance, greater standard deviation of stock returns mean greater risk. For example, in manufacturing, by showing the lower standard deviation of the product dimension indicates more uniform of the product dimension that leads to a higher product quality.

#### ***4. Coefficient of Variation: Relative Variability***

The CV is relative measure of dispersion expressed as a percentage. It's calculated as the ratio of the standard deviation to the mean:

$$CV = (\sigma / \mu) * 100\%$$

The CV deals with the variability of multiple datasets which could have varying units and very different means. Standard deviations, as a matter of convention, are completely irrelevant when comparing: e.g. natural comparisons, like the variability of stock prices (in dollars) and the variability of temperature (in degrees Celsius), are meaningless. However, the CV makes for a decent comparison.

Suppose two datasets have the following properties:

- Dataset A: Mean = 50, Standard Deviation = 10
- Dataset B: Mean = 200, Standard Deviation = 20

The standard deviation of Dataset B is higher, but the CVs are:

- $CV(A) = (10 / 50) * 100\% = 20\%$
- $CV(B) = (20 / 200) * 100\% = 10\%$

In Dataset A, we have more relative variability but less absolute variability (standard deviation). The CV is significant for finance and quality control, since it is needed to compare the relative risk nor process variation

#### ***Choosing the Right Measure for Insightful Analysis***

Without understanding the measure of dispersion, overall analysis about data remains incomplete. Although the range and IQR list all data points (none were included in this example), these options quickly summarize total spread and typical variability. Mean deviation measures the average absolute deviation, whereas the standard deviation and variance are the building blocks for measuring squared mean deviation. Finally, this property enables comparison of relative dispersion among different data sets by the coefficient of variation. Which measure is appropriate and in which case depend on the nature of the data and data context. And having an in-depth understanding of these metrics helps analysts to work with a deeper understanding of how much data can vary, making them lead to better decisions and right conclusions.

---

**1.7 Introduction To Statistics**

---

***1. Unveiling the Shape: Meaning and Interpretation of Skewness and Kurtosis***

The basic concepts of statistics are concerning the central tendency and variation of the data. These measures alone, however, often do not express enough about the underlying distribution. Skewness and kurtosis look deeper into the shape and symmetry of data sets. In layman terms, skewness tells you about the asymmetry of a distribution. A perfectly symmetric distribution (such as the bell-shaped normal distribution) has zero skewness. Having longer or fatter tail to the right denotes positive skewness: The mean of the distribution is higher than the median. This suggests that there are some very high values that are affecting the average. Conversely, in negative skewness (left skewness), the left side has a longer or thicker tail so that it has a mean lower than the median by extreme low-value.

Kurtosis, on the other hand, is analytics of tailenders or peaked Ness of a distribution. It measures how closely data points cluster around a mean and how heavy tails are. Leptokurtic: high kurtosis sharp peak heavy tails adding to the tail extremism Platykurtic distributions have low kurtosis and a lower peak with thinner tails and fewer extreme values. In particular, a normal distribution, the reference, has moderate kurtosis and is called mesokurtic. These properties of data are incredibly revealing in exposing the profound features of data to a level much deeper than basic characteristics of means and spread. Data on the risk side of the distribution tail, such as financial data that are influenced by extreme events, tend to have a high kurtosis. We will get normal distribution for data from stable process.

***2. Measuring Asymmetry: Delving into Measures of Skewness***

In order to measure the skewness, it needs to be quantified. An eternal method of measuring the skew, would be to use (D1) the first coefficient of skewness,

(Pearson), which is calculated between the mean vs mode. This metric is Computed as:  $(\text{Mean} - \text{Mode}) / \text{Standard Deviation}$  If the value is positive, it will have a positive skewness, if it is negative, it will have a negative skewness & if the value is very close to 0., then it is symmetric. However, this measure is sensitive to the mode that is not always reliably determined. Another popular measure is Pearson's second coefficient of skewness based on mean and median. Formula to calculate Skewness:  $3(\text{Mean} - \text{Median}) / \sigma$  (or)  $3(\text{Median} - \text{Mode}) / \sigma$  This is slightly better of a measure compared to the first, since the median is more robust against extreme values than the mode. The sign shows the direction of skewness and its absolute value, the force. A more subtle and routine technique uses the third moment of the distribution. This approach calculates the standardized third moment, resulting in a numerical score that reflects the degree of asymmetry. Typically, this is calculated through software. For example, we have a data set of scores for an exam and we used some statistical software to find out the p-values. A net +0.7 would suggest a “fairly positively skewed” distribution; that is, many of the scores are below the average, such that the higher scores “pull up” the mean. Where a slight negative skew would be -0.3 All this skewness is measure that give a little bit different insights into the nature of the data, it gives researchers and analytics to choose the kind that is better for them.

### ***3. Grasping the Tails: The Kurtosis Index and Its Significance***

Kurtosis, as mentioned earlier, describes the tailenders of a distribution. This property is measured with a number called the kurtosis index (kurtosis) Now, the above formula of kurtosis has the fourth moment of the distribution as its initial part, normalized to the degree that sets up for differences in scale. (Just know that the most common way software packages report this is as “excess kurtosis,” which is  $\text{kurtosis} - 3$ .) This is done so the normal distribution has excess kurtosis of 0 (the kurtosis of the normal distribution is 3).

- **Leptokurtic (positive excess kurtosis):** It has pointy peak and heavy tails (known as leptokurtic). This indicates that data points are clustered near the center, and a broader distribution of tail chances. Leptokurtic distributions are common in financial markets, particularly in stock returns, indicating that

- extreme positive or negative outcomes are more likely than what a normal distribution might imply. For example, a kurtosis index of 5 would imply a leptokurtic distribution while analyzing a data set of hourly stock price changes.
- **Platykurtic (negative excess kurtosis):** A platykurtic distribution has a lower, flatter peak with thinner tails (indicating more evenly dispersed data, such that extreme values are less likely). max is near to 1. A normal distribution ends at 3 std dev so this is more probably a special condition of Less squares or More squares condition where data is limited or controlled.
- **Mesokurtic (with excess kurtosis near-zero):** Mesokurtic distributions (for example, the normal distribution) have intermediate tails and a moderate peak. It is what you were trained to measure against.

From Kurtosis index you can have a hypothesis test of the tailedness of the distribution (how far it is from normality). Such data is vital for risk assessment, statistical modeling, and decision-making.

#### ***4. Practical Applications and Interpretive Nuances***

Skewness and kurtosis are not just themselves abstractions; they bear great practical meaning in various fields. Even in finance, most of these measures represent the risk of investments. Positive skewness in returns would mean that you have a higher chance of having higher returns while high kurtosis indicates a higher chance of lower returns or downside risk. Skewness in production, for example, may show bias in the manufacturing process, while kurtosis can show variation in the dimensions of parts produced. In the social sciences, such measures help facilitate understanding of how income, test scores and other such variables are distributed. Skewness and kurtosis make sense given context, however. In small-sized samples, the utility of skewness and kurtosis estimates can be questionable. Thus, confirm sample size, and use best practice. Additionally, histograms and box plots for data visualization act as Extra M/minimum to the numeric outcomes. In short, if Researchers understand skewness and kurtosis, they will get more insights and eventually will also make better and informed decisions.



---

## UNIT 8 INDEX NUMBERS

---

---

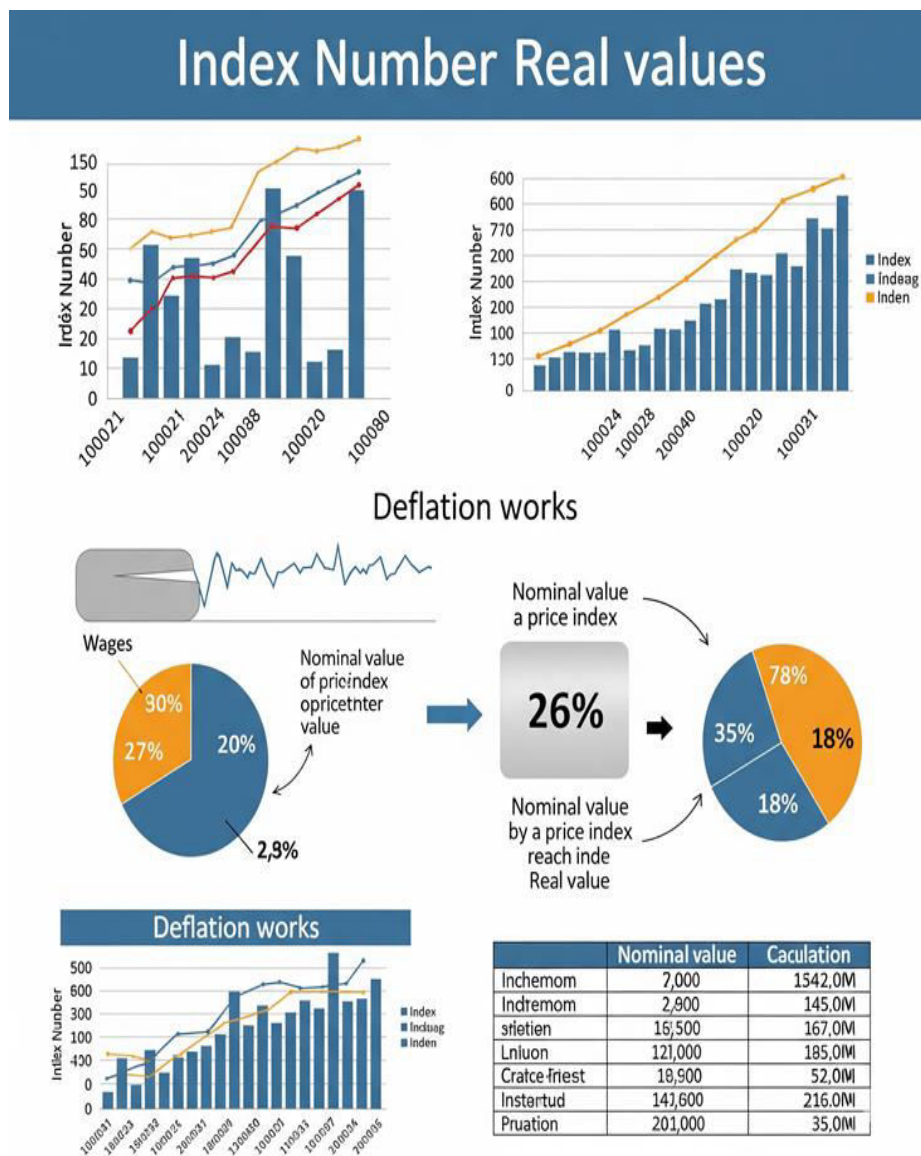
### 1.8 Index Numbers: Meaning And Importance

---

Index numbers are a very efficient statistical tools to measure the variation in a variable (that is a con Shared Attribute) or a set of related variables over time or from them in different locations. In short, they distill data into a number that communicates a lot with little explanation. Instead of working with raw data, which can be unwieldy, index numbers provide a comparative measure of change; an index number uses a base period or location as a reference point. The base is typically set at a value of 100 and the relative amounts are described in percentage terms in comparison to this base. The consumer price index (CPI) is an index number in which a number indicates how much prices have increased from a base period (100). They are important as they indicate trends and patterns not readily discernible otherwise. They are relied on by economists, policymakers, businesses, and researchers who seek to understand and analyze economic phenomena. Index numbers are a statistical measure that enables ongoing quantitative comparisons over time by recognizing that prices, outputs and other variables are always in flux. They assess the impact of economic policy, determine the cost of living, monitor inflation and guide business decisions.

Say, we want to compare the wheat production of a region over a decade. Rather than measuring in raw tonnage which would be misconstrued to larger variables such as area of land, Item of weather and many more, we may take index number. The concept is quite simple, we take a base year, we can say 2010 and index it at 100. This means here if Wheat Production in 2020 = 125 In 2010 we had a Wheat production of 100 and we observe 25% growth with compared figures of earlier year. Simplified it might be, but it makes for rapid, useful comparison. Index numbers also enable comparison across time and space. (For example, where you compare the CPI between countries to measure relative inflation differences. In business, they track sales performance, market share and productivity. Index numbers also aid in summarizing the changes, making informed decisions and strategic planning.) This reduction is not only useful for functions such as development,

entrepreneurship, and innovation (among many others), but also provides important insights due to them being compressed with the exploration of the resulting economic- and socio-historical vectors. But index numbers, you see, also allow you to deflate nominal values into real values. Nominal GDP may rise, but that rise may simply reflect inflation or it may reflect an increase in production. Real GDP measures the value of output produced in an economy while controlling it for inflation and using a price index to deflate the nominal GDP. Therefore, real GDP is adjusted for the price level in the economy.



**Figure 1.5: Index Number Real Values**

### ***Types of Index Numbers***

Broadly, index numbers can be classified on the basis of the variables measured, the methods of construction. Understanding these differences would help us select a relevant index for that use case.

1. **Price Index Numbers:** Price index numbers are most commonly used index numbers, as they measure changes in the general price level. The Consumer Price Index (CPI) is classic example, which seeks to measure average change over time in prices paid by urban consumers for market consumer basket goods & services. The WPI is a measure that tracks the prices of goods sold in bulk as well as in wholesale markets. Another inflation measure is the Producer Price Index (PPI), which looks at the average price increases domestic producers receive for their products.

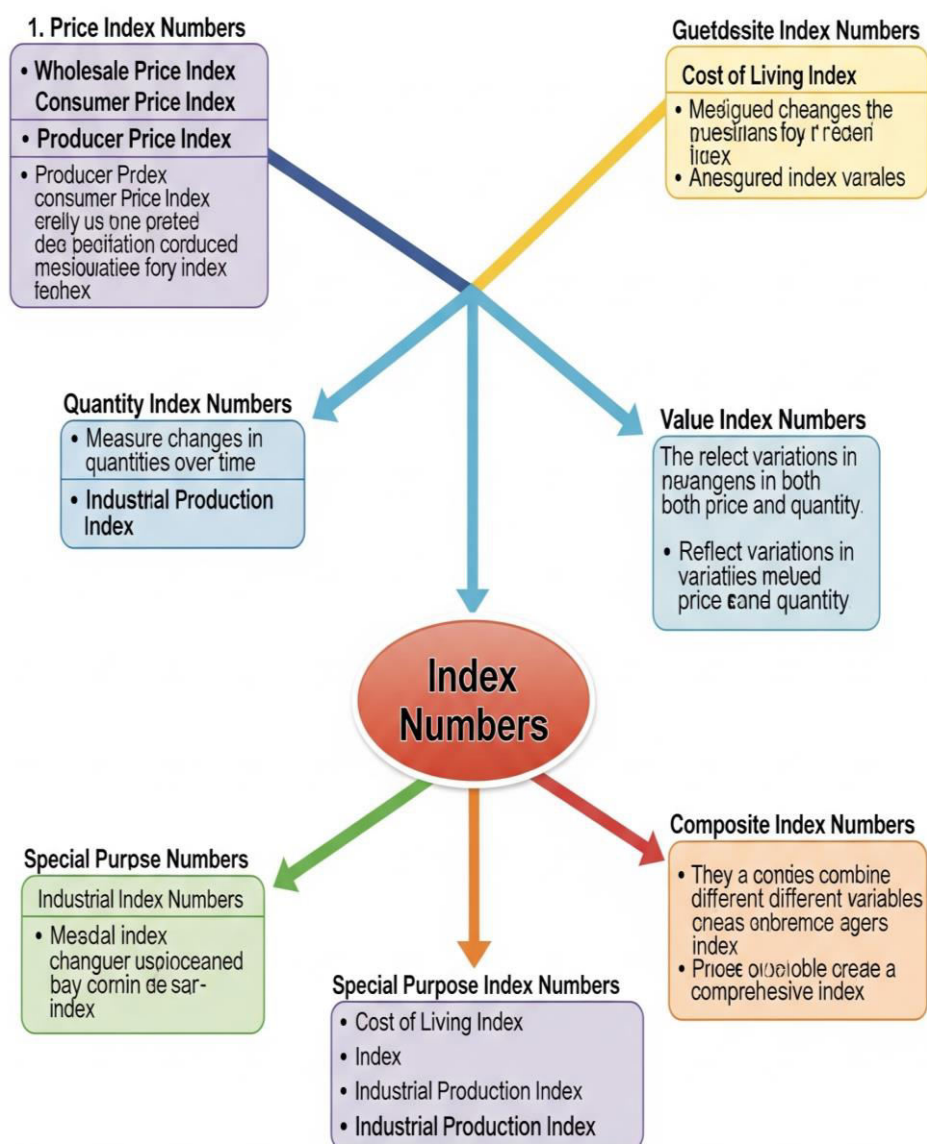
**Example:** the CPI for a country could demonstrate an increase from 100 to 110 from 2020 to 2023, which means that consumer prices rose by 10% over the course of three years.

2. **Quantity Index Numbers:** Volume/quantity of goods & services produced or consumed. To monitor this and arrive at a better assessment of the health of the industry, economists use a number of metrics, one available on a monthly basis Most importantly, the Index of Industrial Production (IIP), which measures growth in the physical volume of production across sectors in the economy

**Example:** If IIP goes up from 100 in one quarter to 105 in the next, it means that industrial output has expanded by 5%.

3. **Value Index Numbers:** Index numbers, which indicate the aggregate value of a variable determined by a combination of price and quantity. They combine both price and quantity movements.

**Example:** Value Higher prices and increased selling volume could lift value index retail sales.



**Figure 1.6: Types of Index Numbers**

The supplied table is a synthesis of Index Numbers, where the purpose and the way of construction is clearly arranged. The nucleus of this idea is "Index Numbers", which can be further classified into five principal varieties : Price Index Numbers, Quantity Index Numbers, Value Index Numbers, Special Purpose Index Numbers, and Composite Index Numbers.

4. **Special Purpose Index Numbers:** These are constructed to represent specific phenomena of change. An instance in this category are stock market indices, the S&P 500 being an example: this index tracks changes in stock prices; indices associated with agricultural production, exports, or imports also fall under this category.

○ **Example:** It would be similar to saying that the index of stock market grow by 15%, means the value of listed stocks increase exponentially.

5. **Composite Index Numbers:** Custom Email Manager You can configure a filter for your emails, and Custom Email Manager will wait them in your inbox all the same. For example, one possible composite economic index also would have production, employment and price indices.

○ **Example,** a number of individual indicators can be aggregated to create an index of economic sentiment, e.g. consumer confidence, business confidence and financial market indices.

Furthermore, index numbers can be constructed using different methods, such as:

- **Simple Aggregative Method:** This simply sums up prices/quantities of all items for a given period and compares to from the base period.
- **Weighted Aggregative Method:** Use this method, where you need to assign weight to each object based on their importance level. Indexing methods are commonly standardized using Laspeyres, Paasche and Fisher ideal index weights.
- **Average of Relatives Method:** For every item, we calculate adjusted price or quantity relatives (ratios) and average them.

Which index type to use, and how to build it would depend on the specific research question, as well as the properties of the data being analyzed.

### Uses of Index Numbers

Index numbers are used in many different fields, so they are an essential tool for analysis and decision-making.

1. **Economic Policy Formulation:** There are few notable applications of Index Numbers, they are listed as follows– Economic Policy Formulation Government and policy makers use the index numbers to keep track of the trends in economy and formulate the policies accordingly. The CPI, for instance, is a vital measure in gauging inflation, and adjusting monetary and fiscal policy. IIP assists to increase industrial growth and formulate plans for enhancing production.

- **Example:** A central bank may raise interest rates to curtail a rise in inflation based on CPI numbers.

2. **Business Decision-Making:** Companies use index numbers to identify sales, expenses, and productivity. They assist in predicting demand, pricing goods and making investment choices.

- **Example:** Using a sales index to detect seasonal trends and guide inventory adjustments.

3. **Wage and Salary Adjustments:** Many wage and salary agreements are linked to the CPI to ensure that workers' purchasing power is maintained in the face of inflation.

- **Example:** sales index to detect seasonal trends and guide inventory adjustments. Using a.

4. **International Comparisons:** It is often used in an index for wage and salary adjustments: Many of the agreements for wages and salaries are tied to the CPI to maintain the purchasing power of workers in the event of inflation.

- **Example:** in many labor contracts cost-of-living adjustments (COLAs) are based on changes in the CPI.

5. **Market Analysis:** In financial markets, stock market indices provide a snapshot of overall market performance and help investors make informed decisions.

- **Example:** A rise in the S&P 500 indicates an overall increase in the value of listed stocks, which can influence investment strategies.

- **Deflating Economic Data:** Inflation adjustment is done using index numbers so that nominal economic data reflect real changes. Nominal GDP, for example, can be deflated by a price index to get real gross domestic product.
  - **Example:** GDP growth is merely 2%. if nominal GDP has grown by 5% and CPI has gone up by 3%, the real
6. **Social Analysis:** They are also used in social analysis to measure a change in social indicators; for instance, poverty rates, health indicators, educational attainment, and health insurance--also referred to as an index number.
- **Example:** An index of human development may be constructed from life expectancy, education and income indices to gauge overall social progress.
7. **Forecasting:** Index numbers serve in time series analysis to discern trends and patterns, thereby facilitating the forecasting of future values.
- **Example:** In the IIP context, it is used to predict future industrial production levels through analysis of potential upcoming trends.

Last but not the least, index numbers are being powerful instruments for analyzing and interpreting economic and social statistics. This ability to take complex information and distil it down into a simple, stand raised form that can be absorbed and understood has made them a must have weapon in the arsenal of policy making, business decision, social analysis and forecasting.

---

## 1.9 SELF ASSESSMENT QUESTION

---

### 1.9.1 Multiple-Choice Questions (MCQs)

**1. What is the primary purpose of statistics?**

- a. To manipulate data randomly
- b. To collect, analyze, and interpret data
- c. To create unnecessary data
- d. To avoid decision-making

**2. Which of the following is an example of descriptive statistics?**

- a. Predicting next year's sales based on past data
- b. Calculating the average marks of students in a class
- c. Testing hypotheses about population parameters
- d. Drawing conclusions about a population from a sample

**3. Inferential statistics involves:**

- a. Summarizing data without making conclusions
- b. Drawing conclusions about a population from a sample
- c. Listing all observations in a table
- d. Measuring only qualitative data

**4. The measure of central tendency that is most affected by extreme values is:**

- a. Mean
- b. Median
- c. Mode
- d. Quartiles

**5. Which of the following correctly defines the median?**

- a. The most frequently occurring value in a dataset
- b. The middle value when data is arranged in ascending order
- c. The sum of all values divided by the total number of values
- d. The difference between the highest and lowest values

**6. Which of the following is true about quartiles?**

- a. They divide data into three equal parts
- b. They divide data into four equal parts
- c. They are always equal to the mean
- d. They are the same as percentiles



**7. Standard deviation measures:**

- a. The difference between the highest and lowest values
- b. The spread or dispersion of data around the mean
- c. The most frequently occurring value in a dataset
- d. The middle value of a dataset

**8. Component: Coefficient Variance (CV) Use the coefficient of variation (CV) to:**

- a. Assess the level of relative variability across solutions.
- b. Have its data range only.
- c. Determine the Most Common Value in a Data Set
- d. Find the average of a set of data.

**9. Skew a dataset is defined as:**

- a. The sharpest or the largest data distribution
- b. Degree and direction of distributional asymmetry in the data
- c. None of the above average for a dataset
- d. Extent: The difference of maximum and minimum values.

**10. What is the word for how pointy or flat a curve is?**

- a. Standard deviation
- b. Skewness
- c. Kurtosis
- d. Range

**1.9.2 Short Questions:**

- 1. What is statistics? Explain its scope.
- 2. Distinguish between descriptive and inferential statistics.
- 3. Explain The Mean, Median and Mode and Give at Least Three Illustrative Examples.
- 4. What are quartiles? Explain their significance.
- 5. Express the meaning of standard deviation and its significance.

**1.9.3 Long Questions:**

- 1. Discuss the utility of statistics and its limitations.
- 2. Explain the various central tendencies.
- 3. Distinguish between mean, median, and mode.
- 4. Explain the significance of the dispersion measures in the statistics.
- 5. Describe the importance of standard deviation and variance expressed from the data.

---

## **MODULE 2 PROBABILITY AND PROBABILITY DISTRIBUTIONS**

---

### Structure

<b>UNIT 9</b>	Introduction to Probability
<b>UNIT 10</b>	Concepts of Probability (Classical, Empirical, and Subjective)
<b>UNIT 11</b>	Probability Laws
<b>UNIT 12</b>	Decision Rule in Probability
<b>UNIT 13</b>	Probability Distributions
<b>UNIT 14</b>	Theorems of Probability
<b>UNIT 15</b>	Concept of Sampling

---

### **2.0 OBJECTIVES**

---

- Explain the concept and significance of probability in statistical analysis.
- Digest classical, empirical, and subjective probability.
- Apply the additive and multiplicative laws of probability to problem solving.
- Understand and employ probabilistic decision-making principles.
- Use basic results from probability theory in statistical calculation operations.
- Be aware of the applications and methods of sampling in statistics.

---

## UNIT 9 INTRODUCTION TO PROBABILITY

---

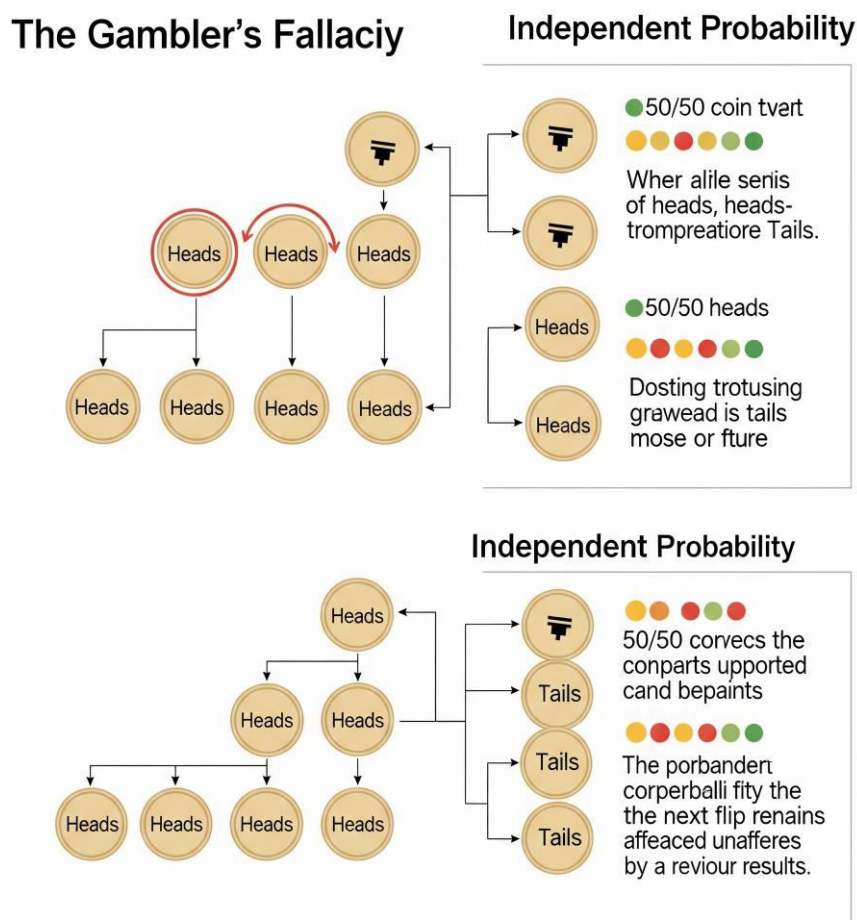
---

### 2.1 Introduction to Probability

---

At the core of everything we experience is probability, influencing our lives in the obvious and the subtle. It's an edifice that rules our paltry uncertainty, forming a vestigial lung we breathe in each day, a vast model of how to make sense of a world in which complete certainty is a rarity. At heart, probability is a measure of the extent to which different things could happen in situations of uncertainty. Consider the weather forecast and the fact that there is a 70% chance of rain, or when a doctor explains the percentage success rate of a medical procedure – these are probability in practice. Although most people think of probability only in terms of a game of dice, and cards, this concept applies to various other spheres of life from gambling to science, to medicine, insurance to financial industry, and right to the way we make decisions in our daily lives irrespective of all rational considerations. Probability has its origins in antiquity, and 16th-century Italian mathematicians Gerolamo Cardano and Pietro Cataldi are among the first to write of it. But during the 17th century, formal probability theory emerged in correspondence between French mathematicians Blaise Pascal and Pierre de Fermat, while working on gambling problems brought to them by a nobleman called the Chevalier de Méré. Their work introduced the concept of how to systematically compute probabilities of different outcome. From these simple origins grew probability theory, which over the centuries became an elegant part of mathematics with fundamental application in the real world. In everyday life we base a myriad of decisions, from the conscious to the automatic, on probability. When we look at the weather before deciding whether to take an umbrella, we're making a decision with probability. When we buy insurance, we're in effect paying to protect ourselves from rare but potentially catastrophic events. Our medical interventions are frequently formulaic, delivered based on statistical evidence of what works with masses of patients. Even a seemingly simple decision such as which route to take to work might entail a back-of-the-envelope calculation of which option is likely to experience less congestion. One of the more interesting things about probability is how it defies our intuition. Human intuition about chance

events is notoriously unreliable, leading to many common misconceptions and biases in our thinking. For example, after seeing five heads in a row when flipping a fair coin, many people intuitively feel that tails is "due" to appear next. This is the “gambler’s fallacy”: that if something happens more often than normal in one period, it will happen less than normal in the next period, or vice versa. When you boil it down, every single coin flip is it’s own event, so the odds of heads tails will always be 50/50, despite the events before. Learning to think like a probability can teach us to recognize and thwart these cognitive biases.



**Figure 2.1: The Gambler's Fallacy**

The language of probability gives us precise ways to discuss uncertainty. We represent these probabilities as a value between 0 and 1 (or 0% and 100%). Zero probability event is impossible—the event can't happen, under no condition. The probability 1 is certainty, it's going to happen, guaranteed. All else is a matter of different sorts of likelihood. For example, a fair six-sided die has equal chance of  $1/6$  (0.167 or 16.7%) of landing on any of the numbers. This system of numbers provides a means to measure uncertainty and to compare different cases. We can classify the boiling of dice as well: so we can make sense of which events occur together, which don't, and if so, which happened first. Two independent events are two events such that if the first happens, it doesn't change the second's probability from whatever it was before we knew that the first event occurred — like independent coin flips. In contrast, dependent events have an influence on each other — such as when you draw cards from a deck without replacing them, and each successive draw affects the constituent cards that remain. If two events are mutually exclusive, then they both cannot happen at the same time — such as a die showing a 3 and a 4 in a single roll. In other words, complementary events are opposites — if one event doesn't happen, the other one must. These kinds of labels help us to choose the right rules when computing probabilities in complicated situations. Probabilities are spread out over the various possible outcomes according to a probability distribution. The easiest type of distribution is the uniform one, in which all results are equally probable, like for a fair die or coin. Unfortunately, many real world processes are not normally distributed. The normal distribution (or “bell curve”) is common in nature and social processes, from the sizes of people's weights and heights to errors in scientific measurements. Other widely used distributions include the binomial distribution (for cases with two possible outcomes, such as success or failure) and the Poisson distribution (for counting rare events on time or space). In probability, we often look for the probability of combined events. The addition rule enables us to determine the probability that one event or another event will occur. For mutually exclusive, together, we just add the probabilities of the individuals. For events that can occur simultaneously, we need to account for the overlap by subtracting the probability of both events occurring together. The multiplication rule helps us find the probability of two events both occurring.

For independent events, we multiply their individual probabilities. For dependent events, we multiply the probability of one event by the conditional probability of the second event given that the first has occurred. Conditional probability addresses how the likelihood of an event changes based on additional information. For example, the probability of a randomly selected person having a certain disease might be quite low. But if we know that subject possesses a particular symptom, then we might raise that probability considerably. We can write conditional probability as "the probability of A given B." It's a fundamental concept in probability, and it's used in many advanced probability ideas, such as Bayes's theorem, which provides a method for systematically updating probability estimates as new evidence or information comes to light. Bayes's theorem is one of the most influential and general ideas in probability. Named for the 18th-century English statistician and minister Thomas Bayes, this theorem offers a mathematical formula for updating beliefs when new evidence is received. It's especially useful when we are interested in knowing the probability of a cause given an observed effect. For example, if someone tests positive for a disease, Bayes' theorem can be used to compute the probability that the person actually has the disease, accounting for how accurate the test is and how common the disease is in the population. This approach can be applied in medical diagnosis, spam filtering, criminal investigation and machine learning classifiers. If the random process were repeated many times, the expected value is the mean (average) value of the random process. In probability theory, it is the product of all possible outcomes with their respective likelihoods and then summed.

In other words, in a game in which you win \$10 with probability 0.2 and lose \$2 with probability 0.8, you can expect to win  $(10 \times 0.2) + (-2 \times 0.8) = \$2 - \$1.60 = \$0.40$ . This implies you would average a gain of 40 cents per play over a large number of such plays. The notion of expected value is central in decision theory, insurance, gambling, investments, and a multitude of areas in which the long run is more important than individual outcomes. Statistics derives from probability theory and is concerned with the collection, analysis, interpretation, and presentation of data.

These statistical tricks give us the tools to make very strong inferences about an entire population on the basis of only samples, and a quantification of how uncertain we are in an estimate, then to test hypotheses, to say something about whether one variable causes another, we now have a way of thinking about those questions. Statistical analysis is indispensable in various fields, such as medical research, quality control in the production industry, development of public policy, and research in the social sciences. One of the general principles for relating theoretical model to reality is the law of large numbers. This is expressed that the average of the result grows close to the expected value as the number of the tasks is grown. For example, if you were to flip a fair coin only 10 times, you could get 7 heads and 3 tails — something far from the expected 50-50 split. Yet if you flip it 10,000 times, the percentage of heads will probably be a lot closer to 0.5. This is why casinos make a steady profit and are financially successful over time, since individual big wins can be averaged out to the casino's theoretical advantage. That randomness and unpredictable nature doesn't mean there is no pattern, not necessarily. In truth, random processes frequently show interesting and uniform phenomena when repeated for many iterations. Stochastic processes are used to model systems that advance in time according to an element of randomness. This is the case with stocks, the flow of particles in a fluid, or the spread of diseases within a population. These processes, however, may behave in nontrivial ways even though they are largely governed by probabilistic rules. Knowledge of these patterns enables scientists and analysts to model and predict systems that would seem at first to chaotic or unpredictable to analyze. Probability is essential for science when using the concept of "statistical significance". Scientists who carry out experiments must decide whether the results they observe reflect an actual effect — an experimental value, such as the speed of light — or is just the result of random chance. Tests of statistical significance allow us to estimate the likelihood of the observed data under the assumption of no genuine effect (the "null hypothesis"). If this probability is low enough (usually under 5 percent or 1 percent), scientists call the results statistically significant: The data suggest we're seeing something more than random chance at play.

This framework has been the foundation of the scientific method through the generations, though it is worth emphasizing that statistical significance does not equate to practical importance. Probability concepts are highly stressed in terms of risk assessment and management. Risk is such a thing for which we could express the probability of an eventual come upon of something bad and the magnitude of the bad thing was going to happen. In addition, insurers apply complex probability models to determine premiums that equilibrate infrequent large pay-outs against continuous income in the form of premiums. Engineers add in that margin of safety when they design systems. Risk assessments are employed by health care providers to identify those patients who are likely to be most in need of preventive interventions. Even personally, our intuitive sense of risk guides countless daily decisions, from how fast we should drive in various conditions to which investments might be appropriate for our retirement portfolios. Probabilistic theorem has developed extraordinarily since the advent of the computer, and computational techniques have crafted new horizons. Methods such as Monte Carlo simulation take random samples to approximate solutions to problems that are hard or impossible to solve analytically. For instance, a financial analyst could simulate thousands of potential future market realities to evaluate investment risks, or a physics researcher could use a random sample to estimate complicated multidimensional integrals. Most machine learning algorithms are based on probability theory, and learn without algorithms being specifically programmed to do so through statistical patterns in data used for decision or prediction making. These computational methods have transformed everything from climate modeling to AI. Probability is the hero of games of chance.” Card games and dice games and roulette and the lottery all provide those rules of probability. Knowing these rules isn’t a guarantee of winning (the house edge is designed to ensure that casinos never lose over the long term), but it allows players to make more informed decisions and avoid common misconceptions. One such example might be the basic strategy while playing blackjack is described by probability and can be used to lower the house edge. Poker is a game that is part probability and part psychology – the players need to assess the probabilities of different hands, as well as their opponents’ likely tactics. Even basic children’s games feature probability through dice or card draws.



In the early 20th century quantum mechanics brought probability to the very core of our understanding of physical reality. While classical physics is deterministic, quantum physics is fundamentally probabilistic. The well-known Schrödinger's wave equation doesn't describe the exact position or momentum of a particle, but a probability distribution — of where the particle might show up when you measure it. This probabilistic character of quantum systems is not a feature of our measuring tools or of our state of knowledge, but rather of the reality itself, at the level of the quantum world. This in itself was revolutionary: it overturned centuries of Determinism and still provokes questions among philosophers about the nature of reality. Genetic transmission occurs in a probabilistic manner so that probability theory is inherently important to both genetics and evolutionary biology. The laws discovered by Mendel are the ones that clarify how characteristics get from parents to the next generation and not simply by accident but in a predictable ratio. For a simple cross between two heterozygotes for a trait, each child would have a 25% chance of both of its inherited alleles being recessive for the recessive trait if the alleles are independent. Population genetics employs probabilistic models to follow the evolution of gene frequencies across generations as a result of forces such as natural selection, genetic drift, mutation and migration. Such models help to account for why the features of species are relatively fixed and why change over time occurs. Decision theory provides a formal structure for optimal decision under uncertainty, where probability and utility (a measure for the value or satisfaction) are combined. When one has to make a decision with uncertain consequences, according to the expected utility hypothesis, one ought to choose the option with the greatest expected utility--the sum of the utility of each possible outcome (albeit weighted by its probability). This model can explain a lot about how humans make choices, from decisions about money to choices about health. But in behavioral economics, we have research showing that people frequently do not adhere to this model of rational behavior, typically because of cognitive and emotional biases, or because their subjective assessments of probability don't match the actual probabilities. Information theory, established by Claude Shannon in the latter part of the 20 century, creates deep links between probability theory and the notion of information and entropy.

In this context, the information carried by message is a function of its unpredictability (rare messages carry more information than common ones). For example, getting a message that "the sun rose today" is at least almost useless because it is so likely and hardly surprising. On the other hand, learning "your lottery numbers actually won" contains a huge amount of information, precisely because it's so unlikely. These factors have applications in data compression, communication systems, cryptography and more recently, we have begun to understand its consequence in systems of biophysical interest such as neural networks, DNA, etc. It happens that probabilistic reasoning goes far beyond mathematics to affect in what mode we understand knowledge and certainty in everyday life. But hold. The project of the Bayesian philosophy of science is to set rational acceptance on the firm bases of probability theory applied to questions of what is known, shown, or believed at any given time, where rational belief people think should behave according to the laws they have come to recognize for everyday life. From this perspective, beliefs should be constantly revised as new evidence occurs, in accordance with Bayes' theorem. This is very different than the classical "yes or no" approach to knowledge and treats knowledge as a matter of degrees of belief and respective confidences. This probabilistic generation of knowledge fits nicely with the way science works, which is to draw tentative conclusions tempered by an openness to new evidence. There are many situations where Probability meets ethics and fairness. However, when resources or opportunities are allocated according to some probabilistic assessment, insurance premiums, loan applications or predictive policing, questions about fairness and discrimination can come into play. For example, pricing insurance on the basis of postal codes might discriminate indirectly against some demographic groups that are heavily represented in certain neighbourhoods. Likewise, machine learning algorithms which predict future outcomes using past data, may end up replicating the existing biases. These challenges have created increasing interest in "algorithmic fairness", creating techniques to ensure that, for example, a probabilistic decision system treats people fairly while still making statistically accurate predictions. There are some interesting facts about human cognition in there. Years of science reveal that people are prone to systematic errors when they reason about probability. We think that dramatic events (a

plane crash) are more likely to happen than they are, while events that are less dramatic but more likely to befall us (car crash) are less likely to happen than they are. We see patterns in truly random sequences and fail to appreciate the role of chance in many outcomes. As our thinking drives those we consult to frame probabilities (the same medical procedure described as having a “90% survival rate” is more attractive than the one with “10% mortality rate”) we become influenced by that framing. Knowing about these cognitive biases can help us to make better decisions in situations of risk and uncertainty. In the modern age, reading of probability has become even more important. Probability information about health risks, financial investments, weather forecasts, and election polls, to name a few, is constantly presented to the public. Misinterpreting probabilities can result in bad decisions with widespread ramifications. The interpreting badly of the results of medical screening can cause unnecessary anxiety or unwarranted courses of treatment. Similarly, failure to understand the margin of error in opinion polls can also produce confidence in the results of an election that may not exist. Better probability education could assist people in making more informed decisions about everything from personal health choices to policy preferences on complex societal issues. The idea of probability distributions generalizes to multivariate probability distributions, which cover situations in which multiple random variables are of interest at the same time. These joint distributions reflect not only the probability of specific outcomes, but also the degrees of association that may exist between variables. The correlation coefficient ranges from 1 to  $-1$ , indicates the strength and direction of a linear relationship between two variables, and 0 indicates no linear relationship. But correlation does not mean causation - this is a fallacy. And just because two variables Scaffidi discusses are correlated does not mean one is causing the other; it could be that both are affected by a third factor, or that the relationship is spurious. Appreciating these differences is important for correctly interpreting results of statistical analysis. Probability theory is still developing and new problems and applications are being addressed. One active area is the development of strategies for responding to extremely rare occurrences that, when they occur, can have huge effects — “black swans,” in the metaphor popularized by the finance expert Nassim Taleb.

Another frontier involves complex systems with many interacting components, where emergent behaviors can arise that are difficult to predict from individual elements. Yet another theme is the natural extensions of probability theory to describe structure and dynamics in networked systems (social networks, transportation systems, biological networks). These advances have only served to extend the range and relevance of probability. The probability theory is in the center of more and more complex and larger applied AI systems. Most machine learning algorithms employ probabilistic models to cope with the uncertainties in data and to predict. Language processing systems for natural languages use this probability to decide which sense a word has in a given sentence. Computer Vision systems score the likelihood that a potential object is what it has been trained to detect. Learning from reinforcement is guided by probability to strike a balance between exploring unknown strategies and exploiting established effective ones and is used to power systems that learn by trial and error. These are some of the most advanced and useful applications of the theory of probability in operation today. As the manner in which societies have perceived chance, randomness, and uncertainty has changed, so has probability theory. In past cultures it was common to attribute casual events to the Gods or to "Fate." The evolution of probability in diverse cultures has stimulated early interest in the study of probabilities. Classical period During the Renaissance, scientists such as Leonardo da Vinci sought to understand the mathematics of probability, but it was Stevin who put it on a firm theoretical basis. The 20th century brought transformative extensions through links with statistics, physics and computer science, among other areas. This evolution endures to this day, and with it probability has wormed its way deeper and deeper into ways we perceive and interact with our complicated world. Objective and subjective interpretations of probability present key dividing lines in philosophy. The frequentist interpretation identifies probability as the relative frequency of the event occurring in a large number of trials, conducted in the same or over similar circumstances, in the long run. This view views the probability as a real property in the world that operates irrespective of human knowledge or belief. The posterior Bayesian perspective, in contrast, views probability as a degree of belief, which can differ between people given what they know beforehand and how they interpret the evidence.

This subjective view permits us to make objective probability statements about one-off occurrences which can not be reproduced (e.g. The chance that it will rain tomorrow”). Both views have their merits and utility, and contemporary probability theory is inspired by elements of both traditions. Probability theory offers a set of important tools for reasoning under uncertainty, but it has some very real limitations and can be abused. Statistical measures can create a false sense of precision or certainty if their limitations aren't understood. Probability calculations are only as good as the assumptions and data that go into them. Bastard models are endearing when they perform well, but catastrophic when they do not. And even perfect probability knowledge does not dispense with value judgments, if we actually knew what the precise probability of different outcomes was, we'd still have to decide which outcome we want. These limitations emphasize the need for complementing probabilistic reasoning with critical thinking, domain knowledge, and ethical considerations when faced with crucial decisions. Finally, probability is one of the most potent intellectual tools available to humanity to make sense of, and macro-navigate, our lightning-strikingly uncertain world. Developed from a course for students of statistics and psychology, this book is relatively easy to read for anyone with high-school-level math. It includes a variety of problems with numerical answers. It allows us to interpret randomness, measure risk, update our beliefs in the face of evidence and make better decisions. At the same time, probability confronts us with the limitations of certainty and prediction. In a world in which we're constantly confronted with incomplete information and unknown outcomes, however, probability literacy provides the route toward a more rational, nuanced and effective engagement with life's essential vagaries. In embracing probabilistic thinking, we are not sacrificing certainty for uncertainty, but simply offloading some of the complexity into a framework better designed to deal with it. It seems to me that the yield is not complete and utter certainty (which may, in any case, be a mirage), but something just as valuable: a systematic way of navigating through the uncertainty that is essential to our personal and collective futures.

**Practical Applications of Probability in Daily Life:** Probability concepts permeate our everyday lives, often in ways we don't immediately recognize. Take weather forecasts, for instance, which we consult almost daily. When meteorologists predict a 30% chance of rain, they're indicating that, based on current atmospheric conditions, similar weather patterns have historically resulted in rainfall about 30% of the time. This likelihood information informs our practical decisions – to take an umbrella, rearrange outdoor plans, be ready for interruptions. The more we know about these probability statements, the better poised we are to make sense of them and to take measures without overreacting or underreacting to the forecast. Another field in which probability ideas are tangible is in the realm of personal finance. Capital allocation choices always come with an element of unknown associated with future return. Diversification, or spreading out investments among different types of assets, mitigates risk specifically because it's unlikely for all investment categories to perform poorly at the same time. Likewise, decisions with insurance are also a kind of intuitive probability reasoning. We buy insurance to guard against scenarios that are unlikely but potentially catastrophic, such as house fires or the diagnosis of a serious illness. The insurance firm charges premiums against the odds of these events and consumers agree to the protection depending on how much they care about the risks how much they are willing to pay. Even basic budgeting incorporates probability as we budget for variable items that vary and we cannot predict from month to month. Many healthcare decisions need to make judgements using probability (although often implicitly, not explicitly). The trade-off in deciding whether or not to undergo a screening test include our prior probability (pretest probability) of the condition in question, the sensitivity (the pretest positive probability) of the test to detect the condition if present, and the specificity (the pretest negative probability) of the test to determine that the tested person does not have the condition if the condition is truly absent. Understanding these probabilities can help patients and doctors make decisions about testing and treatment. And behaviour, such as the decision to drink or not, like diet, exercise and smoking, means assessing trade-offs between probabilities of health states and immediate benefits/ convenience. Although we do not perform these probability computations in a conscious manner, such

intuitive estimates underlie many health behaviors. Transport and travel planning are using probability in different shapes of forms. When we're deciding what time to leave for an important meeting, we naturally take into consideration the possibility of delays - adding some buffer time if we're on the road during rush hour, say, or when the weather is bad. Those GPS navigation apps that have always given estimated arrival times now show ranges of times, to account for uncertainty in conditions. Airlines overbook flights based on the expectation that some people won't show up, weighing the costs of occasionally having to pay for passengers they have to bump against the extra money they make by flying with fuller planes. The same holds true for connections between flights or trains, as when travelers with half a brain plan these transfers they factor in buffer time based on the likelihood of delays, knowing that tight connections raise the likelihood of missing a subsequent departure. Social life is full of probability calculations, even if we don't normally consciously think of it like that. When we read that someone commented on something, and that comment was sincere or sarcastic, we make a probability judgment based on the context, and possibly the tone, our knowledge of the person, and so on. Choosing whom to date and whom to lay are estimates of compatibility and long-term success derived from available information. In a professional environment, we could also strive to maintain connections with individuals most likely in the future to offer opportunities, or offer information. See even routine decisions about what you can and can't bring up in small talk amount to lightning fast assessments of what the other person will and won't tolerate. Consumers decisions often rely on judgements about probability. For consumers, the decision to buy an extended warranty comes down to how likely a product is to fail and the cost of the warranty. When we decide whether, say, to buy a name-brand product instead of a cheaper alternative we haven't tried, what we're often doing is making intuitive probability estimates both of quality and of how satisfied we'll be with the decision later. Deciding how much fresh food to buy is a matter of what you think the likelihood of eating it before it spoils. Purchasing decisions in online shopping involve assessments of the trustworthiness of merchants, the truthfulness of product descriptions and the chance of receiving timely delivery. These are not necessarily not

computations of formal probabilities, but they represent probabilistic reasoning. In reality, tasks around the house use probability in various applied forms. Homeowners have to determine which preventive maintenance steps are a good value — in part, based on the likelihood and expense of problems that could otherwise arise. For example, one's choice of frequency of gutter cleaning is responsive to this person's risk of water damage resulting from clogged gutters. And the same is true of these decisions on when to replace old appliances; it's a trade-off between the increasing chance of death and the cost of a new one.

But basic prudent acts in the home, like having extra light bulbs, batteries or pantry staples on hand, acknowledge some probability that a need will arise one day, even if you can't know for sure when you'll need them. Choices about education and careers require a variety of sophisticated probability judgements. When students pick a major or a course of study, they consider their relative "likelihood of success" in different fields, the availability of jobs in the future, and potential earnings. When it comes to deciding whether to switch jobs or careers, workers weigh the likelihood of positive outcomes against the risks they face in making a move. The choice to further your education or training is partly based on your estimate of the investment in your future in return for a higher-paying job or more personal fulfillment. Even if such estimates are never very precise, they are at least one example of probabilistic thinking about the uncertain future. Social media and knowledge sharing are based on probability judgments of accuracy, and relevance. At a time when we're all overloaded with information, and disinformation, those who read, view and listen to the media should always be questioning how good the source really is, and what the likelihood is that what it's presenting is correct. Multi-sourcing is an - if one independent source confirms, the likelihood of truth increases. Likewise, when we choose which news stories to open or which videos to watch, we are gambling very rapidly on the likelihood that this content will be most valuable or most entertaining as we make a lightning-fast probabilistic calculation from titles or previews and our past experience with similar content. Pleasant pastimes frequently involve challenges that are presented in a probabilistic context. Most board and card games have an element of luck,



where good strategies require a good assessment of probabilities. Fantasy sports participants choose players based in part on probability evaluations of future performance. Gardeners “zone plant,” using hardiness zones to determine which plants are likely to survive in different climates. Weather forecast determines what the outdoor-activity enthusiast does. Probability is also at play when we watch TV as we predict if we’re likely to like a new series before we hit play. These uses of probability thinking during leisure time enrich and are enjoyable. There are about a gazillion probability judgments in cooking and cooking-like activities. Similarly, experienced cooks have an almost intuitive sense for how likely it is that certain techniques will achieve their desired results. At the time of meal planning, it is difficult to estimate if there will be enough time and energy to execute a planned meal on a certain day.

Good food storage is a judgment of the likelihood of needing something versus the risk of it going to waste. Probability enters recipe following, too, as cooks manipulate technique in light of the likely behavior of their specific ingredients and equipment. These problem solving situations with food emphasize the ubiquitous nature of probability thinking in ordinary life. Probability is used in energy use and conservation. Thermostat setting decisions trade-off comfort against energy cost, and programmable thermostats can be used to have different settings depending on the likelihood of occupation. With investments in energy-saving appliances or home-strengthening upgrades, it's a matter of gauging whether you'll save enough over time to make it pay. If nothing else, even little behaviors, like switching off lights when you leave rooms, convey a probabilistic computation of the odds of return in a short term. With worries about the climate on the rise, more and more consumers are taking personal responsibility for their energy choices — from the cars they drive to the light bulbs in their lamps. Being a parent is a constant risk assessment of child safety, development and well-being. The trick for parents is balancing the fact that it's very unlikely their child will be seriously hurt running around at the playground with the developmental value of letting the child take measured risks and experience some independence. The judgments confronted also include when children are judged capable of

new privileges or responsibilities, and these are probabilistic judgments also. There are times when even simple decisions, such as how much food to cook or at what point during the day to set out for a day at the beach, depend to at least some degree on conjectures based on past experience of the likelihoods of various outcomes. This is also an inherent part of good parenting – adjusting these probability estimates as children age and acquire new skills.

Evidence of probability thinking in daily life are time management strategies. When making to-do lists or schedules, we already take into account the likelihood of finishing things according to schedule. Decisions about what to do first are often made not merely as a matter of importance, but as a function of how bad things will get if a task is held off. Padding time between patients recognizes the likelihood that things don't go as planned. Some even involve choices about when to multitask and when to instead attend to one activity, with a consideration of the likelihood of errors or inefficiency when attention is fragmented. For all of this to work, and to use inferential perspective, one would need to make good judgements of the probabilities of both how long tasks will take, and how likely they are to be completed. Depending on many things to which they can't subscribe to probability, and which, if they could, would result with deterrence-which are to say, lives. Speed limits are established in part based on the likelihood and severity of accidents at different speeds. Defensive driving strategies aim to lower the risk of such collisions by properly educating and understanding the dangers associated with driving. The "three second rule" to maintain distance from the vehicle in front makes driving safer, and takes into account the fact that vehicles in front might suddenly stop. Probability is even used in the design of highway systems, as can be seen in such features as merge lanes, traffic circles, and signal timings to reduce the probability of collisions. Local routing decisions, times of departure and arrival all seem to be attempts to compromise between the time that we spend travelling and the odds that we're going to have a crash. Totting up includes delicate probability judgements about what the recipient would like and how they would react. People who are good at giving gifts tend to be good at predicting the likelihood that an individual will like the particular thing one buys. Gift receipts are recognition that judges these

questions and allows them to be revisited if the guess work reflected by them is invalidated in reality. Price ranges on gifts are generally based on a measuring of the importance of the relationship in compromise with the likelihood that items falling within a particular price range may be found. Even choices of when to give a gift card versus a specific item are probability judgments about what the recipient wants, and what the giver knows about what the recipient wants.

## 1. Foundations: Defining Probability and its Core Concepts

At its most fundamental, probability is measure of how likely an event is to occur. This framework allows measuring uncertainty and decision making in the presence of randomness. It's, in a way, a mathematically distilled knowing number that tells you how likely it is that something will happen, which is to say somewhere between 0 (impossible) and 1 (certain). Probability is involved in all things in our lives, such as predicting the weather, diagnosing a person's disease, and even the winning score of games and the closing price of the stock market. To talk about probability, we first need to establish some fundamental concepts. An experiment is simply a method or action that produces an observable outcome. The collection of all possible outcomes of an experiment is called sample space & is usually denoted as  $S$ . An event is subset of sample space that describes a single outcome or outcomes. collection to give an example, consider flipping of a coin. The sample space is  $\{\text{Heads}, \text{Tails}\}$ . For example, this second event "getting heads" is defined as the set,  $\{\text{Heads}\}$ .  $P(A)$  = fraction of favorable number outcomes divided by the total number of possible outcomes when all things are equally likely. In case,  $P(A) = n(A) / n(S)$ ; where  $n(A)$  is number of events in event -A, &  $n(S)$  is number events number in sample space  $S$ . That is classical definition of probability which assumes that all possible outcomes an experiment have same chance of happening, regardless of how likely they are to occur. Instead, we use the empirical definition of probability (or relative frequency approach) in situations where probabilities of outcomes are not equal. This is like establishing probability of an event based on empirical data. Thus, the empirical probability, according to the empirical definition of probability is

given as: If an experiment is repeated 'n' times & event 'A' occurred 'm' times, then empirical probability of A is approximated as  $P(A) = m/n$ . As 'n' becomes large, the empirical probability of A converges to the true probability. To demonstrate this, let us take the example of rolling a fair 6-sided die. Classical probability is given by ratio of number favorable outcomes to the total number of possible outcomes, such as the statement, (number of favorable outcomes/rolling 4)/(number outcomes/1, 2, 3, 4, 5, 6)  $\Rightarrow$  number of favorable outcomes = 1, number of outcomes = 6 and, thus, probability of rolling a '4' =  $1/6$ ., when we roll the die 100 times and get '4' 18 times then the empirical probability =  $18/100 = 0.18$ , and it is pretty close to classical probability  $1/6$  ( $\approx 0.1667$ ). So, we raise the number of rolls say to 1000, and observe 1000 rolls. We wanted the empirical probability to be closer to  $1/6$ . The code simulates this process. You have now mastered conditions and loops now let's write a code, that simulates 1000 rolls of a die and tells you the empirical probability of an even number being rolled. And sure enough, the result of the run (for example 0.505) is quite close to the theoretical probability of the outcome of 0.5 (i.e., three even digits of six possible outcomes). This illustrates that classical probability can approximate empirical probability, with many high numbers of trials.

**2. Conditional Probability and Independence:** In numerous real-world scenarios, events are interconnected rather than isolated. Conditional probability refers to the likelihood of an event (A) occurring, contingent upon the occurrence of another event (B). This allows us to modify our predicted odds as new information emerges.  $P(A|B)$  denotes the conditional probability of event 'A' occurring provided that event 'B' has transpired.  $P(A|B) = P(A \cap B) / P(B)$  For instance, drawing two cards from a regular deck of 52 cards without replacement exemplifies a straightforward scenario. The probability that the second card is a king, given that the first card was a king, is defined as follows: let A represent "the second card is a king" and B represent "the first card is a king". In the first scenario, there are 4 kings in a deck of 52 cards, hence  $P(B) = 4/52$ . Assume we select a king. Among the remaining 51 cards, only 3 are kings. Thus,  $P(A|B) = 3/51$ . The latter refers to the preceding event and provides a general indication of how the likelihood of an event alters with

the occurrence of prior events. Conversely, independent events are occurrences whose consequences do not affect one another. Events A and B are considered independent if  $P(A|B) = P(A)$  or, equivalently,  $P(B|A) = P(B)$ . Mathematically,  $P(A \cap B) = P(A) * P(B)$ . We can commence by flipping a coin twice. The outcome of the initial flip does not influence the outcome of the subsequent flip. The critical inquiry is the result of the second flip, which is entirely independent of the first flip's conclusion, whether heads or tails, despite the game's total being  $1/2$ . Let A represent the event of obtaining heads on the first flip, and let B denote the occurrence of obtaining heads on the second flip. Therefore,  $P(A \cap B) = P(A) * P(B) = (1/2) * (1/2)$ . The likelihood of achieving heads on both flips is  $(1/2) \times (1/2) = 1/4$ . The law of total probability asserts that if the occurrences  $B_1, B_2, \dots, B_n$  constitute a partition of S (being mutually exclusive and collectively exhaustive), then for each event A, the equation  $P(A) = P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + \dots + P(A|B_n)P(B_n)$  is valid. This will assist us in deconstructing the problem into smaller components. To illustrate, consider a factory that has two machines, M1 and M2, that make light bulbs. Let the machines be M1, M2, M3. Machine M1 makes 60% of the bulbs it produces, which has a 3% fault rate. Machine M2 makes 40% of the bulbs, 5% of which are defective. If a light bulb is selected at random, what is the chance that it will be defective? We will let A be the event that you get a faulty bulb. We are given  $P(M1) = 0.6$ ,  $P(M2) = 0.4$ ,  $P(A|M1) = 0.03$ , and  $P(A|M2) = 0.05$ . Using law of total probability:  $P(A) = (0.03 * 0.6) + (0.05 * 0.4) = 0.018 + 0.02 = 0.038$ . Therefore, the probability of a randomly drawn bulb being defective is 0.038 or 3.8%.

### ***3. Random Variables and Probability Distributions: Modeling Random Phenomena***

We introduce random variable to formalize the Manera of handling and analyzing random phenomena. A random variable is set of values whose values are the numerical outcomes of stochastic event. It is gotten on: sample space real numbers. Random variable is either discrete or continuous. This term typically refers to a countably infinite random variable with values that

might include, for example, the number of heads flipped after tossing coin  $n$  times, or the number of bits of a broken part produced by a machine. In the case of continuous random variable, it can take infinitely many values in certain range ( $x$  (e.g., height of a person, temperature of a room, etc.)). Each random variable is associated with probability distribution that describes likelihoods of its possible values. In common, the chance distribution for discrete random variable is defined via a chance mass serve as (PMF), as many probabilities assigned to every potential value. Take the simple example of flipping fair coin three times. Let us say that number of heads, say  $X$ , is random variable. As a result,  $X$  can take on values 0, 1, 2, 3. The random variable  $X$  has probability mass function (PMF):  $P(X=0) = 1/8$ ;  $P(X=1) = 3/8$ ;  $P(X=2) = 3/8$ ;  $P(X=3) = 1/8$ . In case of a continuous random variable, the probability distribution is defined by a probability density function (PDF) which describes relative likelihood of the random variable taking on a given value. Between two points under the PDF curve lies the probability that our random variable belongs to that interval. It represents one of the most widely used continuous probability distributions, commonly known as the normal (or Gaussian) distribution and represented with statistics favorable curve. Normal distribution is commonly used to approximate certain distributions; for example, weight, height, and exam scores.  $E(X)$ : Expected Value of a Random Variable Expectation or mean of random variable  $E(X)$  represents expected value of random variable, which we can define as a variable that takes on random value according to some probability distribution.

---

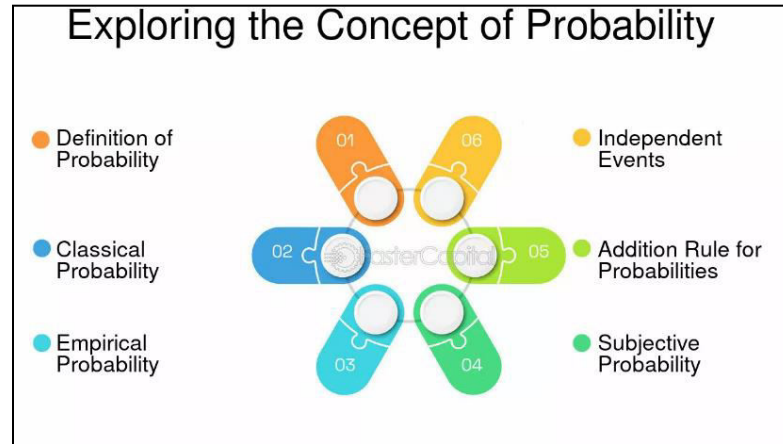
## UNIT 10 CONCEPTS OF PROBABILITY (CLASSICAL, EMPIRICAL, AND SUBJECTIVE)

---

---

### 2.2 Introduction To Statistics

---



**Figure 2.2: Concepts of Probability**

#### ***1. Classical Probability: The Realm of Equally Likely Outcomes***

Classical probability, also known as a priori probability, is founded on basis of equal likelihood of all outcomes of an experiment. This works only in very specific situations such as coin tosses, dice rolls, card draws. The definition states that probability of an event (A) is number ratio of positive outcomes ( $n(A)$ ) to total number of possible outcomes ( $n(S)$ )

Mathematically, this is represented as:

$$P(A) = n(A) / n(S)$$

Classical probability works because of its simplicity, its logical foundations. However, its limitations should be appreciated. It depends on our perfect fairness and symmetry, neither of which necessarily exists in the real world.

#### **Numerical Example 1: Rolling a Fair Die**

Consider standard six-sided die. What is probability of rolling an even number?

- **Total Possible Outcomes (S):**  $\{1, 2, 3, 4, 5, 6\} \Rightarrow n(S) = 6$
- **Favorable Outcomes (A):**  $\{2, 4, 6\} \Rightarrow n(A) = 3$
- **Probability of Rolling an Even Number:**  $P(A) = 3 / 6 = 1/2$  or 0.5 or 50%

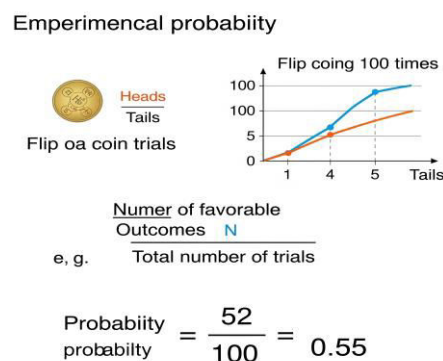
### Numerical Example 2: Drawing a Card

What is probability of drawing an Ace from a standard deck of 52 playing cards?

- **Total Possible Outcomes (S):** 52 cards  $\Rightarrow n(S) = 52$
- **Favorable Outcomes (A):** 4 Aces  $\Rightarrow n(A) = 4$
- **Probability of Drawing an Ace:**  $P(A) = 4 / 52 = 1 / 13$

**Explanation extension:** When we are learning these terms there is other one term that we have to understand that is SAMPLE SPACE. In probability theory, sample space is set of all possible outcomes in a stochastic experiment. So, in the dice problem above, the sample space would be  $\{1, 2, 3, 4, 5, 6\}$ . So, the sum of all possibilities in the sample space must be equal to 1. A die with six faces stands a  $1/6$  chance of falling on any one of the facsimiles on its six sides. For example, by adding  $1/6$  6 times, you obtain 1. One may next consider the case of classical probability. Classical probability has a nice property when it comes to those things where we would expect true random outcomes, like many games of chance.

### 2. Empirical Probability: Learning from Observations



**Figure 2.3: Empirical probability**



**Empirical probability:** It is based on observed data and previous experience; also known as relative frequency probability. It is about how likely an event is based on how often it appeared in trials.

The formula for empirical probability is:

$$P(A) = \text{Number of times event A occurs} / \text{Total number of trials}$$

This is convenient for instances where an application of classical probability cannot be applied due to the fact that there isn't an equally likely outcome. Such as predict weather patterns, predicting failure rate from manufactured products, analyzing customer behavior etc.

### **Numerical Example 3: Coin Toss Experiment**

Assume you flip a coin 100 times & record 53 heads. What is empirical chance of obtaining heads?

- **Number of Times Heads Occur:** 53
- **Total Number of Trials:** 100
- **Empirical Probability of Heads:**  $P(\text{Heads}) = 53 / 100 = 0.53$  or 53%

### **Numerical Example 4: Manufacturing Defects**

A factory produces 10,000 units of certain product. Upon inspection, 250 units are found to be defective. What is empirical probability of a product being defective?

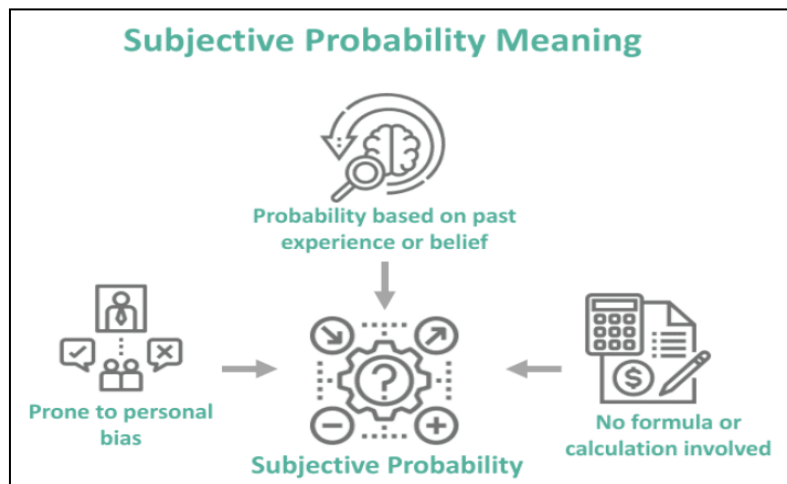
- **Number of Defective Units:** 250
- **Total Number of Units Produced:** 10,000
- **Empirical Probability of Defect:**  $P(\text{Defect}) = 250 / 10,000 = 0.025$  or 2.5%

**Explanation extension:** This is a very useful method to analyze the outcomes of events for which equal probability of all outcomes is not possible and classical probability is not applicable. Example: Weather pattern prediction,

Failure rate prediction of manufactured products, Customer behavior analysis etc.

Probability  
and  
Probability  
Distributions

### 3. Subjective Probability: The Role of Personal Beliefs



**Figure 2.4: Subjective Probability: The Role of Personal Beliefs.**

And this is especially true for rare or unprecedented events for which objective data are scarce or nonexistent. Subjective probability is estimating the probability of something based on how people feel and what they know. It is and often in context such as predicting the success of a new business venture or the outcome of a political election, or the likelihood of a rare medical condition.

#### Numerical Example 5: Startup Success

An entrepreneur thinks that their startup will be successful 70% of the time due to their market research, experience, and instinct. This is a subjective probability assessment.

- $P(\text{Startup Success}) = 0.70$  or 70%

**Numerical Example 6: Medical Diagnosis** A doctor decides that there is a 10% chance, based on a patient's symptoms, medical history, and how

common the disease is, that the patient has a rare disease. Note that this is a subjective probability estimate.

- $P(\text{Rare Disease}) = 0.10$  or 10%

**Explanation extension:** Of the three, subjective probability is the most poorly defined (and therefore the most contentious), because it is so dependent on individual bias. Two very different people who have access to different information might determine very different levels of probability for the exact same event, and be correct. Thus, we often use subjective probability, when objective facts cannot be established. Though individual opinions vary, they remain helpful in risk assessment, and decision making. We, in a lot of different professions, rely on experience, and judgement to make decisions about likely outcomes.

#### ***4. Interplay and Applications: Blending the Approaches***

**Conditional Probability** When discussing the different types of probabilities, it is worth mentioning that in many ordinary life situations classical, empirical and subjective probabilities are used simultaneously. For instance, suppose an insurance company wants to calculate risk of its clients to get in a car accident: It could use classical probability example to measure the probability of accidents, use empirical probability to assess historical claim data and use subjective probability to accounts for individuals risk profile.

Requiring knowledge about and application of these perspectives of probability is critical to making informed choices in many domains, including:

- Finance: Pricing financial instruments, evaluating investment risks.
- Medicine: Disease diagnosis, treatment efficacy assessment.
- Engineering: Studying systems reliability, safety development.
- Business: Sales prediction, marketing campaign optimization.
- Science: Statistical analyses, interpreting experimental results

However, do you know what is powerful tool that allows you to better deal with uncertainty and make sound judgment in a dynamic world by mastering the concepts of classical, empirical, and subjective probability? It is one of the basic corner stones of statistical analysis, and its principals are useful in our daily life.

---

## UNIT 11 PROBABILITY LAWS

---

---

### 2.3 PROBABILITY LAWS

---

Probability Laws: Navigating the Realm of Chance

**1. The Additive Law:** The Additive Law The additive law of probability is critical to calculating the probability of one event or another event. This theorem applies significantly to the cases of mutually exclusive events and non-mutually exclusive events. Mutually exclusive events cannot occur simultaneously, while nonmutually exclusive events can. Disjoint Events (or mutually exclusive events) If A and B are two events which cannot happen at the same time  $P(A \text{ or } B) = P(A) + P(B)$ . Mathematically, we interpret this as:  
 $P(A \text{ or } B) = P(A) + P(B)$

This fits with what we'd expect to happen according to common sense. In cases where two events cannot both happen at the same time, the probability of either occurrence is just the sum of their probability as separate events.

Illustrative Example: Utilize a standard six-sided die. Let event A denote the occurrence of rolling a 2, and let event B denote the occurrence of rolling a 5. The occurrences are mutually incompatible, as it is impossible to roll both a 2 and a 5 simultaneously in a single throw.

$P(A) = 1/6$  (probability of rolling a two)  $P(B) = 1/6$  (probability of rolling a five)

Applying the additive law:  $P(A \text{ or } B) = P(2 \text{ or } 5) = P(2) + P(5) = 1/6 + 1/6 = 2/6 = 1/3$

Consequently, the likelihood of rolling either a 2 or a 5 is  $1/3$ .

When events are not mutually exclusive, meaning they can occur simultaneously, the addition law must be adjusted. NOTICE Due to instances where both events occur, it is necessary to eliminate them to avoid double counting. The equation is expressed as:  $P(A \text{ or } B) = P(A) + P(B) - P(A \& B)$ ,  $P(A \cap B)$  denotes the intersection of occurrences A and B, representing the probability that both events occur simultaneously.

### Numerical Example:

Imagine you are drawing a card from a normal 52-card deck. Let A be event of drawing heart, & B be event of drawing king. [Because one can draw the king of hearts.

- $P(A) = 13/52 = 1/4$  (probability of drawing heart)
- $P(B) = 4/52 = 1/13$  (probability of drawing king)
- $P(A \text{ and } B) = 1/52$  (probability of drawing king of hearts)

Using the additive law for non-mutually exclusive events:

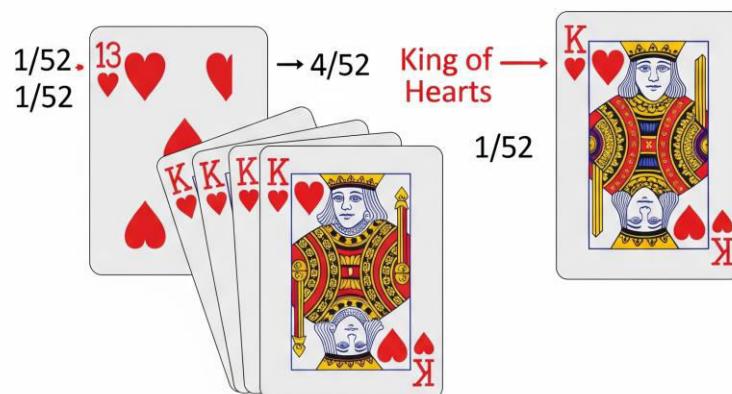
$$P(A \text{ or } B) = P(\text{heart or king}) = P(\text{heart}) + P(\text{king}) - P(\text{heart and king})$$

$$P(A \text{ or } B) = 1/4 + 1/13 - 1/52 = 13/52 + 4/52 - 1/52 = 16/52 = 4/13$$

So, the chance of drawing a heart or a king =  $4 / 13$

The additive law is indispensable from figuring out the chances of winning a lottery to assessing the odds of contracting a disease. It helps us to create scenarios and calculating the possibility of joint events happen that than the foundation of our informed decisions.

## Additive Law of Probability



$$P(\text{Heart or King}) = P(\text{Heart}) + P(\text{King}) - P(\text{Heart and King})$$

$$(13/52) + (4/52) - (1/52) = 16/52$$



Figure 2.5: additive law of probability

## 2. The Multiplicative Law: Determining the Probability of "Both/And" Events

The multiplicative law of probability concerns probability of simultaneous occurrence of two or more events. This is especially important when calculating independent and dependent events. Dependent Events: An event that has the property that the prediction of one event affect another event.

### 2.1. Independent Events:

For independent events that involve A & B, then chances for both the events to happen will be simply the multiplication of probabilities of A & B. Mathematically, this is expressed as:

$$P(A \& B) = P(A) * P(B)$$

The idea is that =total probability of a joint event is product of probabilities of its component events which occur independently of each other.

### Numerical Example:

Example 1: Tossing a fair coin twice Let A be the event that we get heads on first flip, & B be event that we get heads on second flip. The result of one flip does not affect the next; these events are independent.

- $P(A) = 1/2$  (probability of heads on first flip)

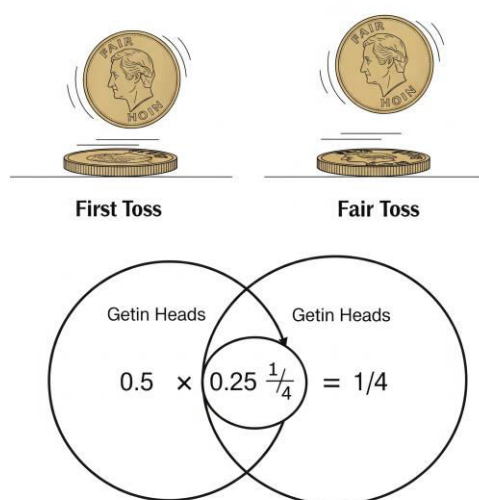


Figure 2.6: Multiplicative Law of Probability

- $P(B) = 1/2$  (probability of heads on second flip)

Using the multiplicative law:

$$P(A \& B) = P(\text{heads} \& \text{heads}) = P(\text{heads}) * P(\text{heads}) = 1/2 * 1/2 = 1/4$$

Therefore, probability of getting heads on both flips is 1/4.

### 2.2. *Dependent Events:*

For dependent events, where one event has an impact on probability of other.

The multiplicative law is based on conditional probability  $P(B|A)$ , The probability of event B occurring, given that event A has already happened.

The equation is expressed as::

$$P(A \text{ and } B) = P(A) * P(B|A)$$

This formulation accounts for dependency among the events, adjusting the likelihood of the second event given the first.

### **Numerical Example:**

Let us think about drawing two cards from a 52-card deck without replacement. Let event A be that we draw a king on the first draw, and event B be that we draw a queen on the second draw. But they are dependent events, because the result of your first draw directly (albeit indirectly) determines the contents of the rest of the deck.

- $P(A) = 4/52 = 1/13$  (likelihood of selecting a king on initial draw)
- $P(B|A) = 4/51$  (the likelihood of drawing queen on second draw, contingent upon a king being drawn first)

Using the multiplicative law for dependent events:

$$P(A \& B) = P(\text{king} \& \text{queen}) = P(\text{king}) * P(\text{queen} | \text{king})$$
$$P(A \text{ and } B) = 1/13 * 4/51 = 4/663$$

So the probability of drawing, without replacement, a king followed by a queen would be 4/663.

One of the most important laws in standalone form is known as the law of multiplication, it is applied in many of the science fields like genetics,

finance, engineering, etc. It allows us to deduce probabilities of complicated events by breaking them up into simpler, subsequent stages. Understanding whether events are dependent or independent is essential to wisdom of the appropriate implementation of this law.

### ***3. Integrating Additive and Multiplicative Laws: Real-World Applications***

They are not exclusive laws and most of the time you use them in conjunction to solve a problem on complex probability solving. There are typically two halves of real-world cases “either/or” and “both/and” “conditions” that should be reconciled.

#### **Example: Quality Control**

Let us consider an example of such a situation we have a manufacturing process where two machines, M1 and M2 produce items: Machine M1 occupies 60% of the product and have defect rate = 2% Machine M2 occupies 40% of the product and have defect rate = 3%.

We are interested in getting the probability for randomly chosen item being defective.

Let:

- A = item produced by M1
- B = item produced by M2
- D = item is defective

We have:

- $P(A) = 0.60$
- $P(B) = 0.40$
- $P(D|A) = 0.02$  (probability of defective given item from M1)
- $P(D|B) = 0.03$  (probability of defective given item from M2)

We need to find  $P(D)$ . We can use law of total probability, which combines the additive and multiplicative laws:

$$\begin{aligned} P(D) &= P(D \text{ and } A) + P(D \text{ and } B) \\ P(D) &= P(A) * P(D|A) + P(B) * P(D|B) \\ P(D) &= (0.60 * 0.02) + (0.40 * 0.03) \\ P(D) &= 0.012 + 0.012 \\ P(D) &= 0.024 \end{aligned}$$

Therefore, the probability that a randomly selected item is defective is 0.024 or 2.4%.



Business  
Statistics

This will give a you an example of how the additive and multiplicative laws come together. By knowing and understanding these basic laws that will allow us to record and analyze uncertainty and make smart decisions. More specifically these probability laws underlie complex probabilistic models and statistical analyses that are employed to better understand the inherent randomness in the world around us.

---

## UNIT 12 DECISION RULE IN PROBABILITY

---

Probability  
and  
Probability  
Distributions

---

### 2.4 Decision Rule In Probability

---

Deciding under uncertainty is a fact of human existence. Whether it is a doctor diagnosing a patient, a financial analyst predicting prices and future market trends, or a weather forecaster estimating the chance of rain, having to decide (for those responsible for the decision) the right option out of a limited (or vague) amount of information is a fundamental task. In order to measure and handle this uncertainty, we turn to math: probability. In effect, a decision rule is a rule-based assumption used to make a decision based on probability of the occurrence of certain events. It bridges subjective probabilities with tangible actions; less-than probabilities translate into objective choices. Probabilistic reasoning in fact giving numeric values, of probability, to what is to happen. These probabilities provide an idea on the basis of available information or based upon previous experiences or deduction. As an example, flipping fair coin, we would say that event heads have a probability (0.5 or 50%) and the event tails (0.5). Reality is not always so convenient. That means there are frequently situations where probabilities are unknown, or they vary with new information. And then enter decision rules and the mechanistic way of making decisions even when faced with ambiguity.

A decision rule usually involves four components: (a) a description of the possible states of the world, (b) a description of a probability distribution over those states, (c) a set of possible actions, and (d) a description of a criterion for selecting the preferred action (decision rule). This criterion is usually expressed in terms of minimizing expected loss or maximizing expected utility. The Expected utility is an assessment of how attractive a certain act is, and it can be how likely its sorted outcomes will appear, and the worth of those outcomes. Expected loss, on other hand, serves as an indicator of how much downside risk we are taking on by taking an action. So, let's consider a simple example: A retailer needs to decide how many units of a perishable product to order. What they have to sell is unknown and excess product at the end of the day must be thrown out. The retailer can use historical transaction data to predict the probability of various demand levels. For example, they

would consider a 30% probability of low demand, a 50% probability of medium demand, and a 20% probability of high demand. They can then compute the expected profit for different stocking levels and choose one that yields highest expected profit. This is how you can implement a decision rule in real world.

And decision rules use thresholds (or some cut-off point). For example, a test for a medical condition might have a threshold probability over which a positive test would be clinically significant. If the probability exceeds this threshold, the doctor might recommend further testing or treatment. This rule is a decision criterion that minimize false positive risk (treat a non-sick patient) against false negative risk (miss a diagnosis). Choosing this threshold is critical because anything in context and relative costs of errors matter.

## ***2. Building Robust Decision Rules: Expected Value, Bayesian Inference, and Risk Assessment***

Sound decision-making requires sound knowledge of probability theory and statistical methods. Beneath it all, one revolves around expected value. For every possible value of  $X$ , one multiplies it by the probability of  $X$  being that value, and then they sum all the products to compute the expected value of  $X$ . It calculates the average outcome of a random event over long period of time. Consider, for example, a lottery ticket that costs \$1 and has a 1% chance of paying off \$100. It will have an expected value of  $(0.01 * \$100) + (0.99 * -\$1) = \$1 - \$0.99 = \$0.01$ . That is to say, for the average person who buys lots of tickets, they'll lose \$.99 for every ticket they buy. Sure, some hypothetical someone comes out on top and wins, but in terms of expected value, the long-term picture is bleak.

Bayesian inference is another strong way to use to create decision rules. It gives us the ability to update our beliefs about the likelihood of events based on new information. This is particularly useful for fields with knowledge that is constantly changing. So, for example, a self-driving car might have initial beliefs about how likely a person will cross the same street and it could use information collected from sensors to adjust those beliefs using something

like Bayesian inference. For demonstration purpose let us take a numeric example. Consider case of a diagnostic test for a rare disease. The test is 95 percent sensitive (correctly identifies 95 percent of people with the disease) and 90 percent specific (correctly identifies 90 percent of people without the disease). The disease affects 1% global population. If person tests positive, how likely is it that they actually have the disease?

Employing Bayes' theorem, we may get the posterior probability:

- Prior probability of having disease ( $P(D)$ ) = 0.01
- Prior probability of not having disease ( $P(\neg D)$ ) = 0.99
- Probability of a positive test given having disease ( $P(+|D)$ ) = 0.95
- Probability of positive test given not having disease ( $P(+|\neg D)$ ) = 0.10

The posterior probability of having disease given positive test ( $P(D|+)$ ) is:

$$P(D|+) = [P(+|D) * P(D)] / [P(+|D) * P(D) + P(+|\neg D) * P(\neg D)]$$

$$P(D|+) = (0.95 * 0.01) / (0.95 * 0.01 + 0.10 * 0.99)$$

$$P(D|+) = 0.0095 / (0.0095 + 0.099)$$

$$P(D|+) = 0.0095 / 0.1085$$

$$P(D|+) \approx 0.0876$$

And this means that even if you get positive test result, probability that you actually have disease is roughly 8.76%. This highlights the delicate balance between prior probabilities and test characteristics that must be struck when considering test results. Decision rule development is really a risk assessment process. This involves the process of identifying potential risks, assessing the probability and consequences of those risks, and developing strategies to mitigate those risks. This can be done using one of many popular methods used for risk assessment, such as sensitivity analysis, scenario analysis, or decision tree analysis. Sensitivity analysis examines how variation in the input of a decision rule impacts its overall output. In fact, scenario analysis enables

to scope out different scenarios while decision tree analysis provides a diagrammatic aid displaying the different pathways taken to arrive at a decision along with the probability and the payoff associated with each. Such techniques add more stability and caution to the decision rules.

### ***3. Implementing and Evaluating Decision Rules: Practical Considerations and Ethical Implications***

You do not train on data past said date, so you have real business decisions to make to train the rules that matter. Use of poor-quality data can never be fixed by even well-trained algorithms, and in the absence of accurate and complete data, poor decisions are bound to be made. Some decision rules are computationally hard and require specialized algorithms and software. Furthermore, human judgment is often critical in the interpretation of probabilistic information and final decision-making. Finally, the first instance in finance trading we can identify are algorithmic trading systems that are systems of decision rules that are programmed with the ability to automatically execute trades based on market data and parameters fed in ahead of time. Propelled by large datasets and sophisticated algorithms free of human bias, these systems can sniff out profitable trading opportunities. However, these systems still need an overseer, in the form of human traders, to be able to monitor their performance and make adjustments when required.

Performance assessment of decision rules is a fundamental issue for the reports of violence. Methodologies like backtesting, simulation and empirical experimentation are used to make this possible. Backtesting means applying a decision rule to past data to check how well it would have performed. That in any case simulation is a way of literally modeling a system and then using that model you wrote to enter all kinds of various decision rules into the model you just wrote. We call this approach real-world experimentation: an effort to implement a decision rule, under controlled circumstances, in the real world, and measure its impact. The creation and application of decision rules also raises ethical dilemmas. Some of the decision rules devoured by AIs could have pernicious consequences or could entrench biases already present in society. For instance, decision rules that are

implemented in criminal justice systems can have unequal impacts on subpopulations. Decision rules have to be fair, transparent and ethical.

Probability  
and  
Probability  
Distributions

Furthermore, the increasing use of artificial intelligence (AI) and machine learning in decision making raises new ethical concerns. And while that is true, an AI algorithm can learn incredibly complex patterns from data; it can just as easily learn to amplify existing biases present in the data set. The challenge for us is to ensure that algorithmic decision systems are fair, transparent and explainable. The takeaway: decision rules are a big-picture approach for dealing with uncertainty and making low-regret choices. What differentiates us is the ability to derive valid decision rules to optimize these outcomes through the use of probabilistic reasoning, statistical methodologies, and ethical constraints. As far as new trends in data science and AI are concerned, decision rules will be an evolving pun.

---

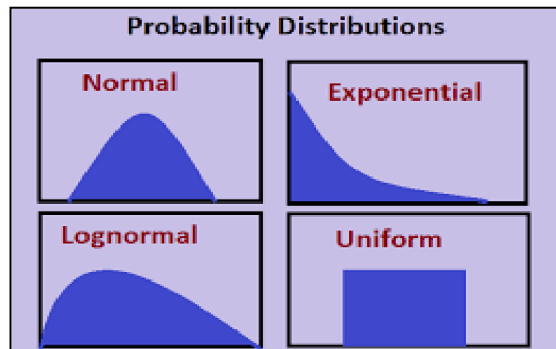
## UNIT 13 PROBABILITY DISTRIBUTIONS

---

---

### 2.5 PROBABILITY DISTRIBUTIONS

---



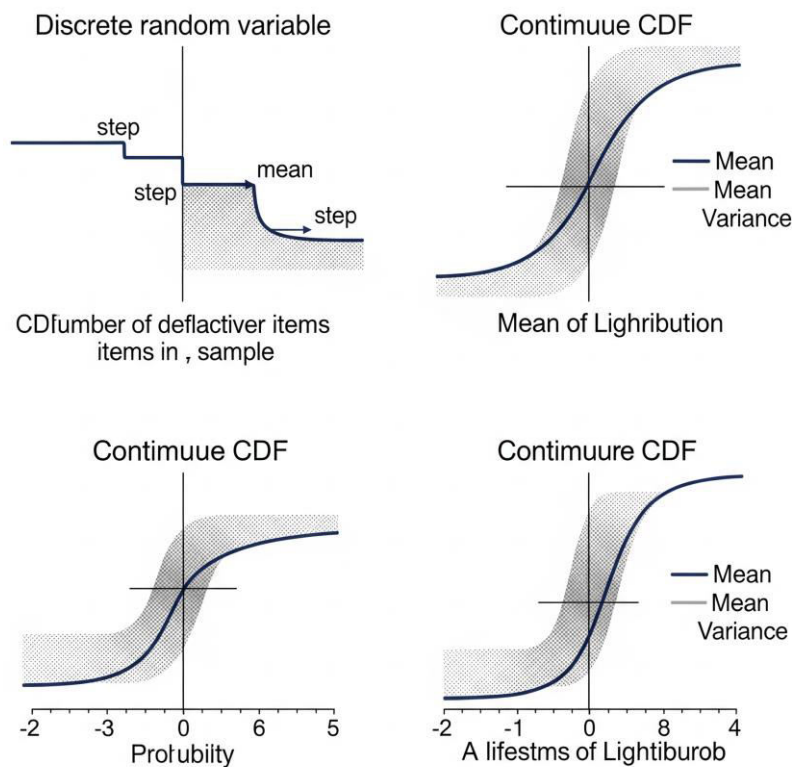
**Figure 2.7: Probability Distributions.**

#### *1. The Foundation: Understanding Probability Distributions*

Probability distributions form the bedrock for statistical inference and predictive modeling. They offer a mathematical structure for characterizing various probability outcomes in stochastic event. Every possible outcome of a random variable has probability mass assigned to it by probability distribution. The occurrence of random phenomena is an event whose fate is absolutely impossible to predict, yet this concept, albeit a little confusing, corresponds to the mathematical field of random variable, which is a variable amount that varies in accordance with the outcome of the real event. There are two types: discrete & continuous random variables. In contrast, discrete random variables have finite or countably infinite domain different values (e.g., the number of heads of coin tosses, the number of defects). The simple answer is that we are ultimately trying to get a better understanding of the uncertainty, and nothing captures the uncertainty better than the probability distribution. Instead of simply stating this event might happen, we can provide a pros and cons of it happening. This enables us to take action and make predictions based on likelihood of different outcomes. PMF indicates probability corresponding to every actual value of PMF. Discrete Stochastic Variables For continuous random variables, PDF (probability density function) describes probability distribution of the continuous random variable

and indicates relative probability that that random variable will equal a given true value. Knowing that CDF is found through integration of probability density function.

One of major tools is cumulative distribution function (CDF). It represents probability that a random variable is no greater than some specified value. The cumulative distribution function (CDF) generalizes to both discrete & continuous random variables. This is useful because predictive distributions only make sense if you understand what every type of parameter represents, so having a mental map of how they act and influence predictions will allow you to more easily navigate their practical functioning. The mean, or expected value  $E(X)$  or  $\mu$ , measures average value of the random variable, and the variance  $\sigma^2 = \text{Var}(X)$  measures the spread of values around that mean. As such,



**Figure 2.8: Cumulative Distribution Function (CDF)**



these properties offer a complete picture of the distribution's shape and where it lies.

## ***2. Discrete Distributions: Binomial and Poisson***

### *2.1 Binomial Distribution: The Probability of Successes*

The Binomial probability distribution is type of probability distribution that describes number of successes in fixed experimental number trials. A Bernoulli trial is a stochastic experiment (such as flipping a coin) that results in a binary outcome, with each possible outcome being assigned either the label of success or failure. These experiments are independent: The outcome of one trial does not influence outcomes of any other experiment. The fastest method is to take advantage of the Bernoulli distribution, which reflects a constant probability of success ( $p$ ) on every trial. There are two key components of the binomial distribution, number of trials,  $n$ , & success probability,  $p$ .

The PMF of binomial distribution can be written as:

The probability formula they provided is probability mass function (PMF) of binomial distribution:

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

Rewriting it with factorial notation:

$$P(X = k) = \frac{n!}{k!(n - k)!} p^k (1 - p)^{n-k}$$

where:

- $X$  denotes random variable that signifies quantity of successes.
- $k$  represents quantity of successes (0, 1, 2, ...,  $n$ )

- The binomial coefficient, denoted as  $\binom{n}{k}$ , signifies number of methods to select  $k$  successes from  $n$  trials. The calculation is expressed as  $n! / (k! * (n - k)!)$ .
- $p$  denotes probability of success in one trial.
- $(1 - p)$  represents likelihood of failure in singular trial.

### Numerical Example:

Consider fair coin being tossed ten times. What is likelihood of obtaining precisely 6 heads?

- $n = 10$  number of trials)
- $k = 6$  (quantity of successes)
- $p = 0.5$  (probability of getting head)

$$P(X = 6) = \binom{10}{6} * (0.5)^6 * (0.5)^4 \quad P(X = 6) = (10! / (6! * 4!)) * (0.5)^{10}$$

$$P(X = 6) = 210 * 0.0009765625 \quad P(X = 6) \approx 0.2051$$

The likelihood of obtaining precisely 6 heads in 10 throws is roughly 0.2051.

The mean (expected value) of binomial distribution is expressed as:

The equation:

$$\mu = n \cdot p$$

The formula for variance in a binomial distribution is:

$$\sigma^2 = np(1 - p)$$

### 2.2 Poisson Distribution: The Probability of Rare Events

Its proof is beyond the scope of the present discussion; in a few instances, some authors employ some distributions, for example Poisson. The Poisson distribution is used to model events that are rare in nature.

For the Poisson distribution, there is one parameter that we need to consider,  $\lambda$  (lambda), or average number of occurrences in given interval.

So, the probability mass function (PMF) of the Poisson distribution is given by:

where:

- $X$  denotes random variable that signifies quantity of occurrences.
- $k$  is number of events (0, 1, 2, ...).
- $\lambda$  is average number of events in given interval.
- $e$  is
- base of natural logarithm (approximately 2.71828).

### **Numerical Example:**

For example, if call center receives an average of 5 calls/min.  $\lambda = 5$  (average number of calls per minute)

- $k = 3$  (number of calls)

$$P(X = 3) = \frac{e^{-5} * 5^3}{3!} \quad P(X = 3) = \frac{0.006737947 * 125}{6} \quad P(X = 3) \approx 0.1404$$

Hence, The probability of getting exactly 3 calls in minute is approximately 0.1404.

The mean & variance of Poisson distribution are both equivalent to  $\lambda$ :

$$\mu = \lambda \quad \sigma^2 = \lambda$$

### **3. Continuous Distributions: Normal Distribution**

#### *3.1 Normal Distribution: The Bell Curve*

Normal Distribution Also known as a Gaussian distribution, it is continuous probability distribution that is symmetric about its mean, giving it a bell-

shaped appearance. This makes normal distribution one of most important distributions in statistics because many natural phenomena and empirical data are often normally distributed. It is defined by two parameters, average ( $\mu$ ) & standard deviation ( $\sigma$ ). The mean gives center of distribution and standard deviation gives distribution.

The normal distribution is defined by its probability density function:

$$f(x) = (1 / (\sigma * \sqrt{2\pi})) * e^{-(x - \mu)^2 / (2\sigma^2)}$$

where:

- $x$  is random variable.
- $\mu$  is mean.
- $\sigma$  is standard deviation.
- $\pi$  is approximately 3.14159.
- $e$  is base of natural logarithm (approximately 2.71828).

### Numerical Example:

Let's say heights of the adult males in particular community are normally distributed with average = 175 cm & standard deviation = 8 cm. Finally, we can standardize the value 190 cm using z-score formula:

First, we need to standardize the value 190 cm using z-score formula:

$$z = (x - \mu) / \sigma \quad z = (190 - 175) / 8 \quad z = 15 / 8 \quad z = 1.875$$

Then we want  $P(Z > 1.875)$ , with  $Z$  a standard normal random variable with mean 0 & standard deviation 1. So, by looking at the regular normal distribution table or calculator, we see that:

$$P(Z > 1.875) \approx 0.0304$$

Therefore, probability that randomly selected male is taller than 190 cm is approximately 0.0304.

## UNIT 14 THEOREMS OF PROBABILITY

### 2.6 Foundations of Probability: Theorems and Applications

#### 1. The Fundamental Principles: Defining Probability and Basic Theorems

The Central Limit Theorem states that sampling distribution of mean tends to be normal, no matter what initial sample distribution looks like, as sample size gets sufficiently large. This theorem underlies many themes of statistical procedures hypothesis testing, estimation of confidence intervals, etc.

#### 1.1 Defining Probability:

- Probability is represented as a numerical value ranging from 0 to 1, inclusive. A probability of 0 signifies that an event is impossible, whereas probability of 1 denotes that an event is certain.
- The probability of an occurrence A, represented as  $P(A)$ , is mathematically defined inside sample space (S) that encompasses all possible outcomes.:
- $P(A) = \text{Number of good results in A divided by total number of outcomes in S}$
- It is important to understand that sample space must contain all possible outcomes.

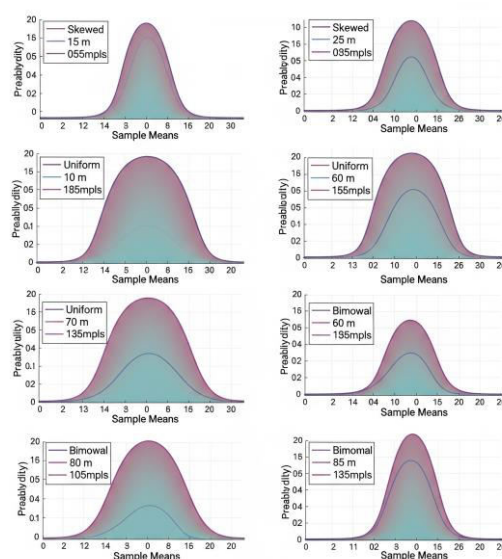


Figure 3.1 Central Limit Theorem (CLT).

### Numerical Example:

- Examine an equitable six-faced dice. The sample space is  $S = \{1, 2, 3, 4, 5, 6\}$ .
- The event of rolling an even number is  $A = \{2, 4, 6\}$ .
- Therefore,  $P(A) = 3/6 = 1/2$ .

### 1.2 Basic Theorems:

- Theorem 1: Probability of the Impossible:
  - § If something cannot happen, the probability is 0.
  - §  $P(\emptyset) = 0$ , with  $\emptyset$  being the empty set.
- Theorem 2 The Probability of a Particular Event:
  - § If an event is certain to happen then its probability is one.
  - §  $P(S) = 1$  (Here, S is sample space).
- Theorem 3: The complement rule:
  - § The probability of an event NOT occurring is 1 minus the probability that the event does occur.
  - §  $P(A') = 1 - P(A)$  where  $A'$  is the complement of event A.

### Numerical Example:

- § For the die above; probability of not obtaining an even number ( $A'$ ) is :
  - §  $P(A') = 1 - P(A) = 1 - 1/2 = 1/2$ .
- Theorem 4 Probability Range:
  - § For any event A,  $0 \leq P(A) \leq 1$ . This implies that risk probabilities will be set in between this range.

### 2. The Addition Theorem: Combining Probabilities

The addition theorem is essential for determining probability of occurrence of either event.

## 2.1 Mutually Exclusive Events:

- Two occurrences are mutually exclusive if they cannot happen at same time.
- If A & B are mutually exclusive, then  $P(A \cap B) = 0$ , where  $\cap$  denotes the intersection of events..

### Addition Theorem for Mutually Exclusive Events:

- $P(A \cup B) = P(A) + P(B)$ , where  $\cup$  denotes union of events.

### Numerical Example:

- Contemplate selecting one card from a regular 52-card deck.
- Let A be the event of drawing heart, and B be event of drawing spade.
- These events are mutually exclusive.
- $P(A) = 13/52 = 1/4$ , and  $P(B) = 13/52 = 1/4$ .
- The probability of drawing a heart or a spade is:
- $P(A \cup B) = 1/4 + 1/4 = 1/2$ .

## 2.2 Non-Mutually Exclusive Events:

- Two events are non-mutually exclusive if they can occur simultaneously.

### Addition Theorem for Non-Mutually Exclusive Events:

- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

### Numerical Example:

- Consider drawing a single card from a standard 52-card deck.
- Let A be event of drawing a king, and B be the event of drawing a heart.
- These events are not mutually exclusive because you can draw the king of hearts.
- $P(A) = 4/52 = 1/13$ ,  $P(B) = 13/52 = 1/4$ , and  $P(A \cap B) = 1/52$ .
- The probability of drawing king or a heart is

- $P(A \cup B) = 1/13 + 1/4 - 1/52 = (4 + 13 - 1)/52 = 16/52 = 4/13.$

### 3. The Multiplication Theorem: Independent and Dependent Events

The multiplication theorem facilitates the computation of probability of simultaneous occurrence of two or more occurrences. It distinguishes between independent and dependent occurrences.

**3.1 Independent Events:** Two occurrences are independent if occurrence of one event does not influence occurrence of other.

#### **Multiplication Theorem for Independent Events:**

- $P(A \cap B) = P(A) * P(B)$

#### **Numerical Example:**

- Consider flipping a fair coin twice
- Let A denote event of obtaining heads on initial flip, & B denote event of obtaining heads on the subsequent flip.
- These occurrences are autonomous.
- $P(A) = 1/2$  &  $P(B) = 1/2.$
- The likelihood of obtaining heads on both flips is:
- $P(A \cap B) = (1/2) \times (1/2) = 1/4.$

### 3.2 Dependent Events and Conditional Probability:

Two events are dependent if the occurrence of one affects cause the occurrence of the other.

#### **Conditional Probability:**

- The probability of event A happening is expressed as  $P(A)$ § The probability of event B happening, if A has already occurred is known as  $P(B|A).$
- $P(B|A) = P(A \cap B) / P(A)$  if  $P(A) > 0$



- **Multiplication Theorem for Dependent Events:**

- $P(A \cap B) = P(A) * P(B|A)$
- **Numerical Example:**
  - Consider selecting two cards from normal 52-card deck without replacement.
  - Let A represent event of drawing a king on initial draw, and B denote event of drawing a king on subsequent draw.
  - These four actions are interrelated.
  - $P(A) = 4/52 = 1/13$ .
  - If a King is drawn on the first draw, there are 3 more Kings remaining in the other 51 cards.
  - $P(B|A) = 3/51 = 1/17$ .
  - The chance at picking out two kings is:
  - $P(A \cap B) = (1/13)*(1/17) = 1/221$ .

#### *4. Advanced Theorems and Applications*

Beyond the fundamental principles, Probability theory encompasses sophisticated theorems that are crucial for addressing intricate situations and practical applications.

#### **4.1 Bayes' Theorem:**

- Bayes' Theorem delineates likelihood of an event, contingent upon prior knowledge of conditions potentially associated with the event.
- It is given by:  $P(A|B) = [P(B|A) * P(A)] / P(B)$

Where:

- $P(A|B)$  is posterior probability of event A occurring, contingent upon truth of event B.
- $P(B|A)$  represents the probability of event B occurring contingent upon the truth of event A.
- $P(A)$  denotes prior probability of event A

- $P(B)$  denotes prior probability of event B.

### Numerical Example:

- A medical test has a 95% accuracy rate. 1% of population has the disease. If person tests positive, what is probability they have disease?
- Let  $D$  = having disease, &  $+$  = testing positive.
- $P(D) = 0.01$ ,  $P(+|D) = 0.95$ ,  $P(+|D') = 0.05$ .
- $P(+) = P(+|D) * P(D) + P(+|D') * P(D') = 0.95 * 0.01 + 0.05 * 0.99 = 0.059$ .
- $P(D|+) = (0.95 * 0.01) / 0.059 = 0.161$  (approximately). Therefore, even though test is 95% accurate, because occurrence of the disease is so rare, there is only a 16.1% chance the person has the disease if they test positive.

### 4.2. Law of Total Probability:

- This theorem provides a way to calculate the probability of an event that can happen in more than one way.
- If  $A_1, A_2, \dots, A_n$  is mutually exclusive & exhaustive &  $B$  is an event, then:
- $P(B) = P(B|A_1) P(A_1) + P(B|A_2) P(A_2) + \dots + P(B|A_n) P(A_n)$

### 4.3. Applications:

These theorems are vital in numerous fields:

- **Statistics:** Hypothesis testing, confidence intervals.
- **Finance:** Risk assessment, portfolio management.
- **Medicine:** Diagnostic testing, epidemiological studies.
- **Computer science:** Machine learning, artificial intelligence.

By mastering these fundamental and advanced theorems, one gains the ability to navigate the complex world of probability and apply its principles effectively to solve a wide range of real-world problems.

---

## UNIT 15 CONCEPT OF SAMPLING

---

---

### 2.7 CONCEPT OF SAMPLING

---

#### *1. Unveiling the Need for Sampling: From Vast Populations to Manageable Insights*

Some populations (like the entire country of China, for example) are simply too large or too complex to be able to study head to toe allowing researchers and analysts to cherry pick a smaller, manageable sample to draw conclusions. Imagine trying to parse the sentiment of every citizen in a country, catalog the quality of every good coming off production line or model growth of every tree in giant forest. Such efforts would be far too time consuming and costly not to mention, logistically impossible. This is where the concept of sampling comes into play.) Sampling is the technique of assessing a part or sample of a bigger population to represent the features of the whole population.

So rather than trying to take on the entire population, we are dealing with a few, more tractable entities, to extrapolate from them to the larger whole. The reasoning goes that as long as a sample is representative of population, we can get useful information without needing to look at every single case. Not only is sampling practical, it is also efficient. Focusing our attention on a single sample allows us to conserve a great deal of resources: time, money, people. Mind that this timeliness is critical in disciplines like market research, where time-to-insight is crucial for business decisions. So, for instance, a company launching a new product might create a test event featuring a select audience of target customers to gauge interest in the product before committing to a full production run. Similarly, in the medical domain the clinical trials most often refers to a sequence of testing new pharmaceutical or treatment on a subset of patients in order to validate efficacy and safety before large scale deployment in patient population. Generalizability, the ability to apply knowledge derived from a sample to all of (or some relevant portion of) a population, is the cornerstone of scientific discovery and the evidence-based policymaking that drives much of the contemporary world.

The effectiveness of sampling, however, depends upon how representative the sample is. Assuming sample is representative of population findings will be valid, but if it turns out to be a biased sample, the resulting conclusions will be incorrect. Sampling aims to eliminate bias by making sure that sample reflects diversity and community. Characteristics This means being intentional about how the sample is drawn, how many people to sample, and what potential sources of error exist. But numerous sampling methods have been developed, each with distinct advantages and disadvantages. The selection process can also be different based on the requirements of research, characteristics of the population studied, and available resources at play. Thus, a proper sampling strategy is vital in order to verify the research results

### **Numerical Example:**

For example, a producer produces 100K lamps a day. They want to estimate the percentage of defective bulbs. There are 100,000 of them, so testing all of them isn't feasible. Instead they go with a sample. They choose a random sample of 1,000 bulbs. They are tested, and 20 of them are found to be faulty. What does this mean at this level: This means that the sample defect rate was 2% (20/1000) From this sample data, they can extrapolate that 2 percent of the overall batch of 100,000 bulbs is probably defective and that 2,000 bulbs are likely faulty. This conclusion is not the best, but rather a good approximation based on the sample.

## ***2. Navigating the Sampling Landscape: Types of Sampling Techniques***

Selecting a suitable sampling method is one of the factors that is critical in the research process since it affects the sample's representativeness and the research results' generalizability. Broadly, the two sampling techniques can be defined as Probability sampling: The method of sample selection gives each member of population known, non-zero chance of being chosen. This allows for sample representation and enables the population's statistical conclusions.

What you have is random sampling, where it is done randomly, this represents something roughly along the lines of "with probability," so no bias should be around here. However, in non-probability sampling there is no point or indicator, and some bias is introduced into the sample.

### **Probability Sampling Techniques:**

- **Simple Random Sampling:** This is the simplest form of probability sampling, wherein each individual in population has the same chance of being chosen. It's kind of like drawing names from a hat. While the technique is straightforward, it is difficult to apply at scale, particularly where populations are geographically separated.
- **Systematic Sampling:** It refers to selecting every  $n$ th member of population (here  $n$  is fixed sampling interval). For example, in case of a population size of 1,000 and sample you want to get of 100, your sampling interval will be:  $1,000 / 100 = 10$ , every 10th member will be selected. While this is very efficient, it can introduce bias if there is some hidden pattern in population.
- **Stratified Sampling:** This technique segments a population into strata or subgroups according to specific characteristics (such as age, gender, or income). A basic random sample is subsequently extracted from each stratum in a manner that ensures the proportions of these traits in the sample mirror those seen in the population. This is especially beneficial when engaging with varied communities.
- **Cluster Sampling:** In stratified sampling, the population is segmented into clusters, such as geographical regions or educational institutions, from which random clusters are then chosen. All units inside the designated clusters are incorporated in the sample.
- **Multi-stage Sampling:** This technique combines multiple sampling methods (eg, stratified, cluster), to create a sample that is both more efficient and representative. For instance, a researcher may want to first stratify the population by region of the country, and then randomly select clusters from

- within each region, and then take a simple random sample from clusters samples.

### **Non-Probability Sampling Techniques:**

- **Convenience Sampling:** Where samples are selected within the reach of the researcher, and are easy to access. An example might be a researcher interviewing people walking by on a street corner. Cheap and easy to implement; however, method has bias issues **Judgment sampling:** A process of collecting samples in an image while the researcher pulls from their expertise or skill of the material. In one, a marketing manager selects a sample of customers whom she believes accurately represents her target market. This is helpful when certain knowledge is required, but this leads to bias if the researcher's judgement was wrong (quantitative).
- **Quota Sampling:** In this method of sampling, a sample is selected according to a specific quota for certain types of characteristics such as sex or age group, education level, etc. That could be, for instance, a researcher who wants to interview an equal number of men and women. This is similar to stratified sampling, except that, you do not have to do the random selection here.
- **Snowball Sampling:** This sampling technique is applied in cases of some hard-to-access populations like drug users, or homeless individuals. It is simply identifying small group of people in population and asking them to refer more. This method is useful for obtaining samples from hidden populations, however, could introduce bias in the outcome if the first group of individuals was not truly representative of population.

### **Numerical Example:**

A university wants to understand how students feel about the services on campus. So they will perform stratified sampling. There are four strata in the student population: freshman, sophomore, junior, and senior. The university ensures that the sample is proportionally representative of each class. Alternatively, if the university's population consists of 25% each of the classes, Freshman, Sophomore, Junior, Senior, then a sample of 400 would yield 100

Freshman, 100 Sophomores, and so on. Doing so will ensure classes are not being misrepresented.

### ***3. Sizing Up the Sample: Determining the Right Sample Size***

The size of the sample it generates in a sampling process is one of the major components of sampling. If sample is small enough, it may misrepresent population, resulting in erroneous results. Or too large a sample size an unnecessary drain of time & money.

#### **Factors Affecting Sample Size:**

- **Population Size:** Larger populations require larger samples to be representative. But it's not a straight line between the two. Once a population reaches a certain size, increasing the sample size provides diminishing returns.
- **Precision:** The margin of error expresses precision, the range within which responses from the sample are presumed to reflect values in the population. Smaller margin of error requires a larger sample size.
- **Variability of the Characteristics Being Investigated:** Larger sample sizes are needed to detect substantial variation in the characteristics under scrutiny. In an opinion neutral about any topic an extremely large sample size is needed in order to identify difference.
- **Confidence level:** This is the degree of certainty that the sample outcome falls within the margin of error. A more confident level needs bigger sample size. Most common confidence levels are 95% and 99%.
- **Sample Size Formulas:**

Depending on type of data being collected & desired level of precision, several different formulas may be used to determine an appropriate sample size. The formula for sample size related to proportion is:

The formula you provided is:

$$n = \frac{Z^2 \cdot p \cdot (1 - p)}{E^2}$$

Where:

- n is sample size
- Z is Z-score corresponding to desired confidence level
- p is the estimated population proportion
- E is desired margin of error

To estimate number of voters supporting a specific candidate with 95% confidence level & a 3% margin of error, assuming a population proportion of 50%, the required sample size is:

$$n = (1.96^2 \cdot 0.5 \cdot 0.5) / 0.03^2 = 1067.11$$

Therefore, the researcher would need a sample size of approximately 1,0



---

## 2.8 SELF ASSESSMENT QUESTION

---

### 2.8.1 Multiple-Choice Questions (MCQs)

**1. What is the probability of an impossible event?**

- a. 1
- b. 0.5
- c. 0
- d. 100%

**2. Which of the following is a type of probability based on historical data?**

- a. Theoretical probability
- b. Experimental probability
- c. Subjective probability
- d. Axiomatic probability

**3. It is a characteristic of the additive law of probability that if two events are mutually exclusive, then the probability of either of them occurring is the sum of their probabilities. What formula signifies this law?**

- a.  $P(A \cap B) = P(A) + P(B)$
- b.  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- c.  $P(A \cup B) = P(A) + P(B)$
- d.  $P(A | B) = P(A) / P(B)$

**4. What probability distribution is used if an experiment results in exactly two potential outcomes (success and failure)?**

- a. Poisson distribution
- b. Binomial distribution
- c. Normal distribution
- d. Exponential distribution

**5. What percentage of data are within one standard deviation in a normal distribution?**

- a. Fifty percent
- b. Sixty-eight percent
- c. Seventy-five percent
- d. Ninety-five percent

**6. What is 1 feature of the Poisson distribution?**

- a. It is used when the data is continuous.
- b. It used for occasional events in a period of time.
- c. It is possible only in case of normal distribution.
- d. It follows a binomial distribution.

**7. Given that  $P(A) = 0.6$  and  $P(B) = 0.3$ , and that occurrences A and B are independent, what is  $P(A \cap B)$ ?**

- a. 0.9
- b. 0.18
- c. 0.3
- d. 0.6

**8. Which of the following best defines the decision rule in probability?**

- a. A rule that helps to choose between two probabilities
- b. A rule to determine whether to reject or accept a null hypothesis
- c. A method to calculate expected values
- d. A formula for binomial probability

**9. The sum of probabilities of all possible outcomes in a sample space must be:**

- a. 1
- b. 0
- c. Between 0 and 1
- d. Greater than 1

**10. What is the primary assumption of the composer of the binomial?**

- a. As many attempts as you like
- b. Variable chance of success
- c. Fixed number of trials and independent events.
- d. Probability Distribution

**11. The theorem which describes the probability of some other incident occurring when another event has already occurred is represented by  $P(A|B) = P(A \cap B) / P(B)$ ?**

- a. Total Law of Probability
- b. Bayes's Theorem
- c. Probability given that something happens (botherhead 2.3-7)
- d. Multiplication Rule

**12. What is significance of sampling in probability?**

- a. It complicates the study.
- b. It also assists in investigating large numbers of populations by small ones.
- c. It gives the results with 100% accuracy.
- d. That is it removes all doubt.

**13. Which of the following distributions is continuous?**

- a. The binomial distribution
- b. Poisson distribution
- c. Gaussian distribution
- d. Hypergeometric distribution

**14. What is an application of the Poisson distribution in real life?**

- a. Pass student in an examination
- b. Number of calls received in a call center every hour
- c. Heights of the students in a class are given by:
- d. Monthly sales of a product.

**15. In probability, one event that has no impact on another event is:**

- a. Dependent Event
- b. Dependent event if not independent event
- c. Conditional event.
- d. None of the above

### **2.8.2 Short Questions**

1. Define probability and its significance
2. Explain the additive and multiplicative laws of probability.
3. What is the decision rule in probability?
4. Define binomial distribution and its properties.
5. What are the characteristics of a normal distribution?
6. Explain the Poisson distribution and its applications.
7. What are the basic theorems of probability?
8. What is the role sampling in probability?
9. What role do probability distributions play in data analysis?

### **2.8.3 Long Questions**

1. Describe the various types of probability with examples.
2. Explain the addition and multiplication laws of probability with examples.
3. Describe the properties of binomial, Poisson, and normal distributions.
4. State the real-life applicability of the theorems of probability.
5. How is probability being used in decision making in business?
6. Discuss the idea of sampling and its dirnlication to statistics.
7. Explain the decision rule in probability and its significance.
8. Compare and contrast binomial and normal distributions.

---

## **MODULE 3 CORRELATION AND REGRESSION ANALYSIS**

---

### **Structure**

**UNIT 16** Introduction to Correlation

**UNIT 17** Positive and Negative Correlation

**UNIT 18** Karl Pearson's Coefficient of Correlation

**UNIT 19** Spearman's Rank Correlation

**UNIT 20** Introduction to Regression Analysis

**UNIT 21** Least Square Fit of a Linear Regression

**UNIT 22** Two Lines of Regression

**UNIT 23** Properties of Regression Coefficients

---

### **3.0 OBJECTIVES**

---

- Describe the meaning and importance of correlation in statistical analysis.
- Determine & explain the direction & strength of relationships among variables.
- Calculate and interpret linear correlation by the Pearson's method.
- Compute and interpret the rank correlation coefficient of non-parametric data.
- Use linear regression and R-Square implementation with Least Square Method to fit a line to data and calculate your square of your fit another straight and another data group.
- Interpretations of the regression lines equation for the two variables  
Understand and interpret the equations of regression lines of two variables.
- Identify and discuss key properties and implications of regression coefficients.

---

## UNIT 16 INTRODUCTION TO CORRELATION

---

Correlation  
And  
Regression  
Analysis

---

### 3.1 Introduction To Correlation

---

#### *1. Unveiling the Relationship: The Essence of Correlation*

Correlation is statistical concept that quantifies degree of association between two variables. It allows us to determine whether alterations in one variable are associated with modifications in another. The association does not imply causation, but shows correlation & dependencies that can be of great value in other areas.

- **Defining Correlation:**

- Correlation analysis investigates the degree and direction of a linear relationship between two quantitative variables.
- We use it to make decisions, such as: “As A increases, does B go up, down, or stay the same.”

- **The Significance of Correlation:**

- Correlation is a bedrock of data analysis, research, & decision making.
- In science, it can help to establish possible links between observations.
- It is useful in many business for understanding customers preferences and market orientations.
- In finance, it measures the correlation between asset prices.

- **Correlation vs. Causation:**

- It is essential to note that correlation does not imply causality. The correlation between two variables does not imply causation.
- There might be a third, unobserved variable influencing both, or the relationship could be coincidental.
- An investigation may reveal a correlation between ice cream sales & crime rates. Nonetheless, it seems more probable that elevated temperatures augment both ice cream sales & crime rates.

#### *2. Measuring the Strength and Direction: Correlation Coefficients*

Correlation coefficients yield a numerical value indicating degree & direction of linear association between two variables. Pearson's r is most often utilized coefficient.

### **Pearson's Correlation Coefficient (r):**

- Pearson's r quantifies linear correlation between two variables.
- It ranges from -1 to +1:
  - +1 signifies an impeccable positive association.
  - -1 signifies an ideal negative correlation.
  - 0 indicates no linear correlation.
- **Understanding the Values:**
  - Values approaching +1 or -1 signify a robust association.
  - Values approaching 0 signify a weak or nonexistent association.
  - Example values.
  - $r = 0.9$ : Strong positive correlation.
  - $r = -0.7$ : Strong negative correlation.
  - $r = 0.1$ : Weak positive correlation.
  - $r = -0.2$ : weak negative correlation.
  - $r = 0$ : no correlation.
- **Calculating Pearson's r:**
- Pearson's r formula incorporates the covariance of the two variables along with their standard deviations.

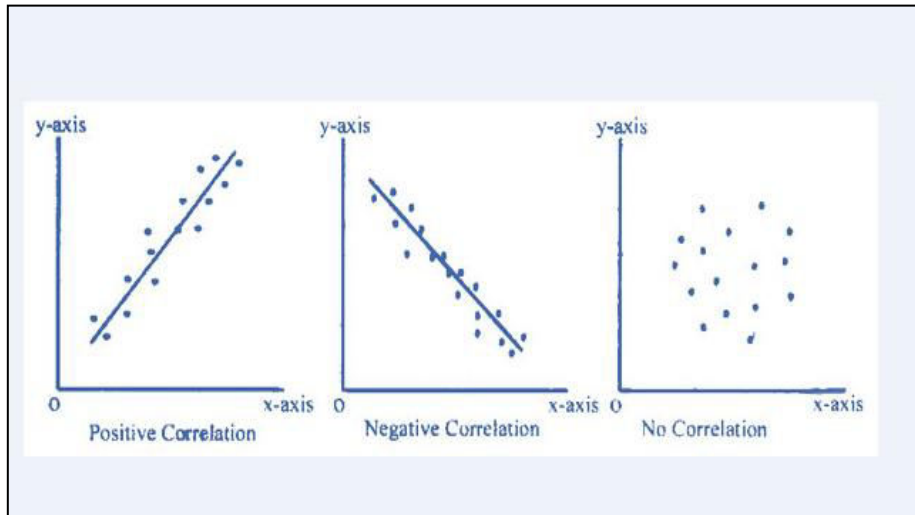
### **Formula:**

- $$r = [\Sigma(x - \bar{x})(y - \bar{y})] / [\sqrt{\Sigma(x - \bar{x})^2} * \sqrt{\Sigma(y - \bar{y})^2}]$$
- Where:
- x and y are the variable values.
- $\bar{x}$  and  $\bar{y}$  are the means of x and y.
- $\Sigma$  denotes the sum.

---

**3.2 Positive And Negative Correlation**

---



**Figure 3.2: Positive and Negative Correlation.**

***1. Understanding Correlation: The Foundation of Relationships***

- Correlation is statistical metric that quantifies degree to which two variables fluctuate in relation to one another. This is a key notion in data analysis that enables the identification of patterns and correlations within datasets.
- It is essential to recognize that correlation does not signify causality. A relationship between two variables does not mean a cause-and-effect-relationship between them. There may be other factors at play that Nino is influencing.
- We will delve into how correlation is calculated, how it is read, and what cannot be told from correlation.
- Correlations are scored from -1 to +1.
- A value of +1 corresponds to perfect positive correlation.
- A -1 value represents perfect negative correlation.



- A value of 0 indicates no correlation.

### **Visualizing Correlation: Scatter Plots:**

Scatter plots are essential tools for depicting the relationship between two variables. Each point on the graph represents a pair of values, with one variable shown on the x-axis and the other on the y-axis.

- By examining the configuration of the points, we may ascertain the intensity and direction of the link.
- A trend of points ascending from left to right signifies a favorable association.
- A decreasing trend of points from left to right signifies a negative association.
- Randomly spread points indicate minimal or no association.

### **The Correlation Coefficient:**

The correlation coefficient, represented as "r," measures the degree and direction of the linear relationship between two variables.

Pearson's correlation coefficient is the primary type of correlation coefficient, evaluating the linear relationship between two continuous variables. Comprehending the magnitude of association.

- Values approaching +1 or -1 signify a robust association.
- Values approaching 0 signify a weak or nonexistent association.

For instance:

- $r = 0.9$ : Indicating a robust positive association
- $r = -0.7$ : Indicating a strong negative connection
- $r = 0.1$ : indicates a weak positive connection.

### **Numerical example of calculating Correlation:**

- To show a basic example, we will use a small dataset.
- Lets say we have the following data of study hours and exam scores.
- Study Hours(x): 1, 2, 3, 4, 5.
- Exam Scores(y): 50, 60, 65, 80, 90.
- We can then calculate Pearson correlation coefficient. This involves finding mean of x & y, standard deviation of x & y, & covariance of x & y.
- After the calculations, we would find a high positive correlation. This means that as study hours increase, exam scores also increase.
- Explaining the formula of Pearsons correlation is very technical, therefore it is more important to explain the meaning of the resulting number.

## ***2. Positive Correlation: When Variables Move Together***

- **Definition and Characteristics:**
  - A positive correlation transpires when two variables simultaneously grow or decrease. In other words, an increase in one variable correlates with an increase in other variable, whereas a reduction in one variable correlates with decrease in other variable.
  - This relationship is represented by a positive correlation coefficient.
  - Examples of positive correlation are abundant in various fields.
- **Real-World Examples:**
  - **Height and Weight:** Generally, taller people tend to weigh more, demonstrating a positive correlation.
  - **Study Time and Exam Scores:** As study duration grows, examination scores often enhance.
  - **Advertising Spending and Sales:** Increased advertising spending often leads to increased sales.
  - **Temperature and Ice Cream Sales:** As the temperature rises, the sales of ice cream tend to increase.
  - **Exercise and Calorie Expenditure:** The more someone exercises the more calories they will burn.

- **Numerical Example:**

Let us examine the correlation between weekly exercise duration and caloric expenditure.

**Data:**

- Hours of Exercise (x): 1, 2, 3, 4, 5
- Calories Burned (y): 200, 400, 600, 800, 1000
- In this example, as number of hours spent exercising increases, number of calories burned also increases proportionally. This is a clear illustration of positive correlation.
- If we were to plot this data on a scatter plot, the points would form an upward sloping line.
- If we calculated the Pearsons Correlation coefficient, the result would be a number very close to 1.

---

## UNIT 18 KARL PEARSON'S COEFFICIENT OF CORRELATION

---

Correlation  
And  
Regression  
Analysis

---

### 3.3 Karl Pearson's Coefficient Of Correlation

---

Karl Pearson's correlation coefficient 'r' is a statistic that quantifies linear correlation between two continuous variables. It quantitatively assesses extent to which a linear equation can represent the relationship between those variables. The coefficient resides within the interval of -1 to +1, where:

- +1 signifies perfect positive linear correlation, indicating that when one variable rises by 2, other also increases proportionally by 2, with all points aligning precisely on a straight line with positive slope.
- -A correlation of -1 indicates perfect negative linear relationship, wherein an increase in one variable corresponds to a drop in other, with all data points aligning precisely along a straight line with negative slope.
- 0 means no linear correlation, so no straight-line relationship between variables. This doesn't necessarily mean there is no relationship, it may be non-linear relationship.
- A value between -1 & +1 signifies varying degrees of linear correlation. The value between +1 and -1 quantifies linear relationship strength. The closer the value is to 0, weaker linear relationship is.

This is determined by ratio of covariance of two variables to the product of their standard deviations. Covariance measures the degree to which two random variables co-vary, whereas standard deviation quantifies extent to which values of each variable diverge from the mean. Karl Pearsons Coefficient of Correlation Formula:

$$r = \text{Cov}(X, Y) / (\sigma X * \sigma Y)$$

Where:

- r is Pearson correlation coefficient.
- Cov (X, Y) is covariance between variables X & Y.

- $\sigma_X$  is standard deviation of variable X.
- $\sigma_Y$  is standard deviation of variable Y.

Alternatively, using raw scores, the formula can be expressed as:

$$r = \frac{[n(\sum XY) - (\sum X)(\sum Y)]}{\sqrt{\{[n(\sum X^2) - (\sum X)^2][n(\sum Y^2) - (\sum Y)^2]\}}}$$

Where:

- n is number of data pairs.
- $\sum XY$  is sum of products of paired scores.
- $\sum X$  is sum of X scores.
- $\sum Y$  is the sum of Y scores.
- $\sum X^2$  is sum of squared X scores.
- $\sum Y^2$  is sum of squared Y scores.

### Numerical Example:

Now let us consider a numerical example, calculating Karl Pearson's correlation coefficient. Let us say we have the following dataset for the Study hours (X) & Test scores (Y) of 6 students:

Student	Study Hours (X)	Test Scores (Y)
1	2	50
2	3	60.0
3	4	65
4	5	75
5	6	80
6	7	90

To calculate 'r', we need to compute Following:

1. **Calculate  $\sum X$ ,  $\sum Y$ ,  $\sum XY$ ,  $\sum X^2$ , and  $\sum Y^2$ :**
  - $\sum X = 2 + 3 + 4 + 5 + 6 + 7 = 27$

- $\sum Y = 50 + 60 + 65 + 75 + 80 + 90 = 410$
- $\sum XY = (2*50) + (3*60) + (4*65) + (5*75) + (6*80) + (7*90) = 1940$
- $\sum X^2 = 2^2 + 3^2 + 4^2 + 5^2 + 6^2 + 7^2 = 159$
- $\sum Y^2 = 50^2 + 60^2 + 65^2 + 75^2 + 80^2 + 90^2 = 28850$

2. **Plug the values into the formula:**

$$r = [6(1940) - (27)(410)] / \sqrt{\{6(159) - (27)^2 - (410)^2\}}$$

$$r = [11640 - 11070] / \sqrt{\{954 - 729\}}$$

$$r = 570 / \sqrt{\{(225)(5000)\}}$$

$$r = 570 / \sqrt{1125000}$$

$$r = 570 / 1060.66$$

$$r \approx 0.537$$

Therefore, the Karl Pearson's coefficient of correlation between study hours and test scores is approximately 0.537. It can be observed that this is a positive linear relationship. Test scores rise as time spent studying rises, but the connection is slightly less than perfectly linear.

### Interpretation and Significance

Correlation coefficient should be interpreted taking into account its magnitude and sign.

- **Magnitude:** The absolute value of 'r' signifies intensity of linear correlation.
  - $|r| \geq 0.8$ : Strong correlation
  - $0.5 \leq |r| < 0.8$ : Moderate correlation
  - $0.2 \leq |r| < 0.5$ : Weak correlation
  - $|r| < 0.2$ : Very weak or no correlation
- **Direction:** The sign of 'r' indicates direction of linear relationship.
  - Positive 'r': Positive linear correlation (variables increase together).
  - Negative 'r': Negative linear correlation (variables move in opposite directions).

It is essential to recognize that correlation does not imply causality. The more they are positively correlated does not mean that if it happens A B it microsoftm means that it isAB. There could be other variables affecting both, or this relation might be spurious.

---

## UNIT 19 SPEARMAN'S RANK CORRELATION

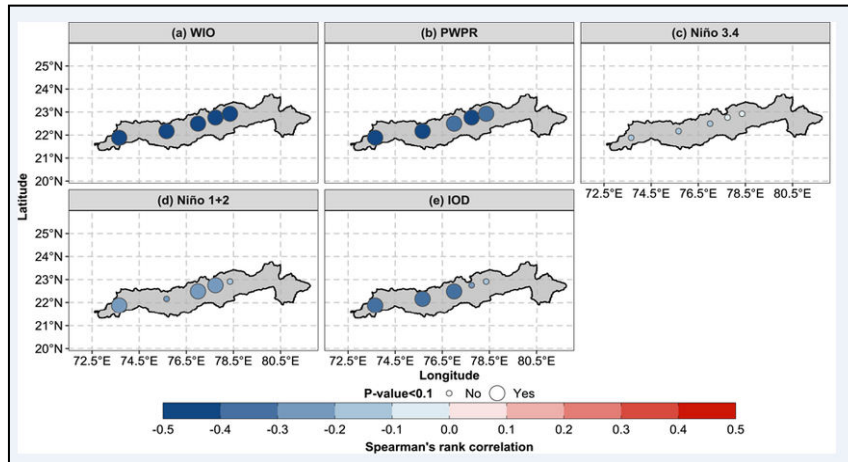
---

Correlation  
And  
Regression  
Analysis

---

### 3.4 Spearman's Rank Correlation

---



**Figure 3.3: Spearman's Rank Correlation Coefficient.**

#### *Understanding Non-Parametric Correlation*

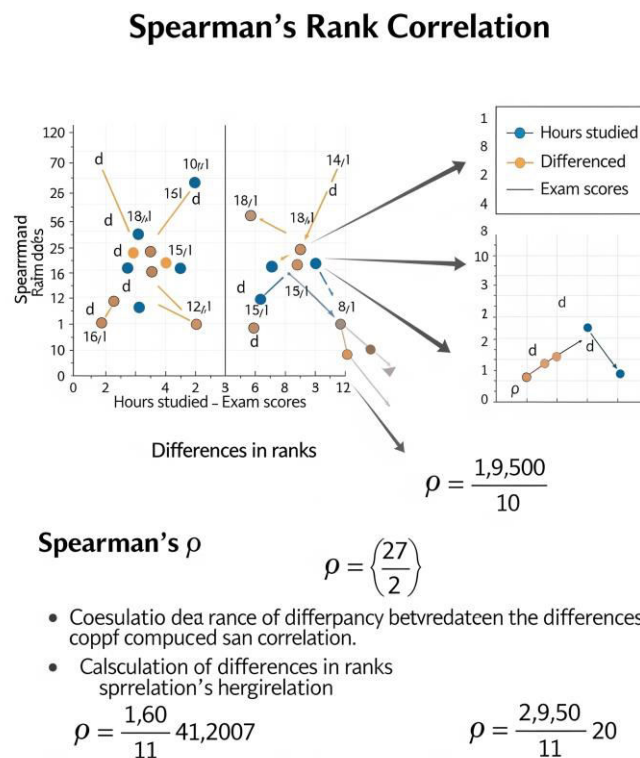
Spearman's Rank Correlation ( $\rho$ ) serves as non-parametric alternative to Pearson's correlation coefficient. Pearson's correlation is confined to linear associations among continuous variables, Spearman's correlation analyzes monotonic relationships between ranked data, where outliers and non-normally distributed data will not affect results significantly. Basically, it describes how well the relationship between two variables can be explained through monotonic functions: If one variable goes up, the other one will also go up (or down) but that does not have to be on a constant rate. Hence, Spearman rank correlation is especially valuable when dealing with ordinal data, such as survey Likert-scale responses, or when data is continuous but violates the assumptions of normality that are necessary for a valid Pearson's correlation. To be even more specific, heart of Spearman's correlation is converting the raw data to ranks and then finding a correlation coefficient on these ranks. This method works because it removes the influence of extreme values and considers the relative ranks of the data points we have, so we can get a true measure of the association regardless of the skewness in the distribution or outliers. Since you are concerned only with ranks instead of



data points, Spearman's correlation focuses on the trend of how two variables vary with respect to each other, regardless of the exact numerical distances between them. Due to its applicability to diverse datasets, it serves as a potent instrument in disciplines such as social sciences, psychology, and market research, where data seldom adhere to normal distribution. The coefficient,  $\rho$ , which varies from -1 to +1, indicates the presence of a statistical relationship between the data, whether positive or negative. +1 signifies a perfect positive monotonic relationship, -1 denotes a perfect negative monotonic relationship, & 0 represents the absence of a monotonic relationship. The intensity of the association is shown by the size of the coefficient, while its direction is denoted by the sign.

### ***Calculating and Interpreting Spearman's Rank Correlation: A Step-by-Step Guide with Numerical Examples***

collected from five students their scores on that exam. Data were with an example to understand the process. For instance, consider examining the impact of the number of hours students dedicate to preparing for an impending



**Figure 3.4: Spearman's Rank Correlation.**

examination on correlation, which is computed by: Now let us go through the steps Spearman's Rank:

Student	Hours Studied (X)	Exam Score (Y)
A	10	20
B	15	25
C	8	18
D	20	35
E	12	22

### Step 1: Rank the Data

First, we rank the values of X & Y separately in ascending order. If there are ties, we assign the average rank to the tied values.

Student	Hours Studied (X)	Rank of X (Rx)	Exam Score (Y)	Rank of Y (Ry)
A	10	2	20	2
B	15	4	25	4
C	8	1	18	1
D	20	5	35	5
E	12	3	22	3

### Step 2: Calculate the Differences in Ranks (d)

Next, we calculate the difference (d) between ranks of each pair of observations ( $R_x - R_y$ ).

Student	Rx	Ry	d ( $R_x - R_y$ )
A	2	2	0
B	4	4	0
C	1	1	0
D	5	5	0
E	3	3	0

### Step 3: Square the Differences ( $d^2$ )

We then square the differences ( $d^2$ ) to eliminate negative values.

Student	d	$d^2$
A	.00	0
B	.00	0
C	.00	0
D	.00	0
E	.00	0

### Step 4: Sum Squared Differences ( $\Sigma d^2$ )

We sum the squared differences ( $\Sigma d^2$ ). In our example,  $\Sigma d^2 = 0 + 0 + 0 + 0 + 0 = 0$ .

### Step 5: Apply Spearman's Rank Correlation Formula

The formula for Spearman's Rank Correlation is:

$$\rho = 1 - (6\Sigma d^2) / (n(n^2 - 1))$$

Where:

- $\rho$  is Spearman's Rank Correlation coefficient.
- $\Sigma d^2$  is sum of squared differences in ranks.
- $n$  is number of data pairs.

In our example,  $n = 5$ , and  $\Sigma d^2 = 0$ . Plugging these values into formula:

$$\rho = 1 - (6 * 0) / (5(5^2 - 1)) \quad \rho = 1 - 0 / (5 * 24) \quad \rho = 1 - 0 \quad \rho = 1$$

This result indicates a perfect positive monotonic relationship between number of hours studied & exam scores.

### A More Complex Example with Ties

Let's consider another example with ties in the data:

Student	Study Time (X)	Exam Performance (Y)
F	12	75
G	15	80
H	10	70
I	15	80
J	18	90

### Step 1: Rank the Data with Ties

For X: 10, 12, 15, 15, 18. The ranks are 1, 2, 3.5, 3.5, 5 (15 is tied, so we take the average of 3 and 4). For Y: 70, 75, 80, 80, 90. The ranks are 1, 2, 3.5, 3.5, 5 (80 is tied, so we take the average of 3 and 4).

Student	X	R <sub>x</sub>	Y	R <sub>y</sub>
F	12	2	75	2
G	15	3.5	80	3.5
H	10	1	70	1
I	15	3.5	80	3.5
J	18	5	90	5

### Step 2: Calculate Differences (d)

Student	R <sub>x</sub>	R <sub>y</sub>	d
F	2	2	0
G	3.5	3.5	0
H	1	1	0
I	3.5	3.5	0
J	5	5	0

### Step 3: Square the Differences (d<sup>2</sup>)

Student	d	d <sup>2</sup>
F	0.0	0
G	0.0	0
H	0.0	0
I	0.0	0
J	0	0

**Step 4: Sum the Squared Differences ( $\Sigma d^2$ )**

$$\Sigma d^2 = 0$$

**Step 5: Apply the Formula**

$$\rho = 1 - (6 * 0) / (5(5^2 - 1)) \rho = 1$$

Again, we get perfect positive correlation.

Let's consider a different set of data that creates a result that is not 1.

Student	Study Time (X)	Exam Performance (Y)
K	10	90
L	12	80
M	15	75
N	18	70
O	20	60

**Step 1: Rank the Data**

Student	X	R <sub>x</sub>	Y	R <sub>y</sub>
K	10	1	90	5
L	12	2	80	4
M	15	3	75	3
N	18	4	70	2
O	20	5	60	1

**Step 2: Calculate Differences (d)**

Student	R <sub>x</sub>	R <sub>y</sub>	d
K	1	5	-4

Correlation  
And  
Regression  
Analysis

### Step 3: Square the Differences (d<sup>2</sup>)

Student	d	d <sup>2</sup>
K	-4	16
L	-2	4
M	0	0
N	2	4
O	4	16

### Step 4: Sum the Squared Differences (Σd<sup>2</sup>)

$$\Sigma d^2 = 16 + 4 + 0 + 4 + 16 = 40$$

### Step 5: Apply the Formula

$$\rho = 1 - (6 * 40) / (5(5^2 - 1)) \quad \rho = 1 - (240) / (5 * 24) \quad \rho = 1 - 240 / 120 \quad \rho = 1 - 2 \quad \rho = -1$$

In this case, we have a perfect negative correlation.

Now, let's consider a scenario with less perfect correlation.

Student	Study Time (X)	Exam Performance (Y)
P	10	75
Q	12	80
R	15	70
S	18	85
T	20	65

### Step 1: Rank the Data

Student	X	R <sub>x</sub>	Y	R <sub>y</sub>
P	10	1	75	3
Q	12	2	80	4
R	15	3	70	2
S	18	4	85	5
T	20	5	65	1

### Step 2: Calculate Differences (d)

Student	R <sub>x</sub>	R <sub>y</sub>	d
P	1	3	-2
Q	2	4	-2
R	3	2	1
S	4	5	-1
T	5	1	4

### Step 3: Square the Differences (d<sup>2</sup>)

Student	d	d <sup>2</sup>
P	-2	4
Q	-2	4
R	1	1
S	-1	1
T	4	16

### Step 4: Sum the Squared Differences (Σd<sup>2</sup>)

$$\Sigma d^2 = 4 + 4 + 1 + 1 + 16 = 26$$

### Step 5: Apply the Formula

$$\rho = 1 - (6 * 26) / (5(5^2 - 1)) \quad \rho = 1 - (156) / (5 * 24) \quad \rho = 1 - 156 / 120 \quad \rho = 1 - 1.3$$

$$\rho = -0.3$$

In this case, we have a moderate negative correlation.

Correlation  
And  
Regression  
Analysis

### Interpreting the Results

- **ρ = +1:** Ideal positive monotonic correlation. As one variable escalates, the other concomitantly escalates consistently.
- **ρ = -1:** Ideal negative monotonic correlation. As one variable escalates, the other invariably diminishes.
- **ρ = 0:** No monotonic correlation. The variables are not related in a consistent increasing or decreasing manner.
- **Values between -1 and +1:** Indicate varying degrees of correlation. The proximity of value to +1 or -1 indicates a higher association. A correlation closer to 0 indicates a weaker relationship.

---

## UNIT 20 INTRODUCTION TO REGRESSION ANALYSIS

---

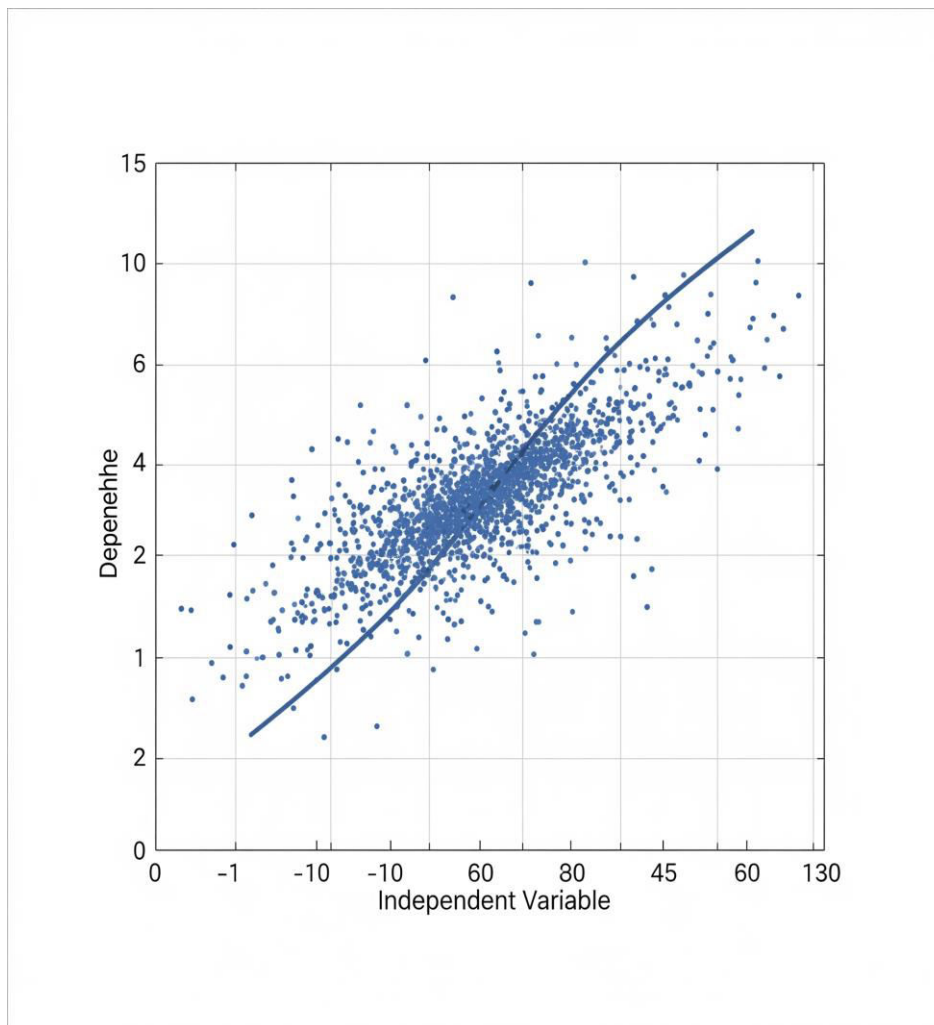
Correlation  
And  
Regression  
Analysis

---

### 3.5 Introduction To Regression Analysis

---

It can either be simple or multiple depending upon the number of which they have to relate. The primary aim is to understand the correlation between changes in independent factors and changes in dependent variable. In summary, regression seeks to establish a line or curve that accurately represents relationship between variables, enabling use of independent variable



**Figure 3.5: foundational concepts and purpose of regression analysis.**



values to forecast the dependent variable's value. Regression's capacity to forecast future outcomes from historical data renders it one of the most essential statistical models now employed, with applications across diverse domains such as economics, finance, social sciences, and engineering. This will allow researchers to detect and study these interactions, quantify their strength and direction, and so predict and generalize results. Linear regression is the fundamental form of regression that assumes a linear relationship between variables, although polynomial regression and multiple regression can accommodate non-linear correlations and numerous predictors. Regression provides methods to evaluate model's goodness of fit, indicating its explanatory power about the data, and to analyze the statistical significance of the predictors.; and flag potential outliers or influential data points. Well, it is essential since regression analysis offers a mechanism that helps understand and qualify relationships, including how variables influence each other.

***Building and Interpreting a Linear Regression Model: A Step-by-Step Numerical Example***

We will use a numerical example to demonstrate how to build and interpret a simple linear regression model. Let's say we wish to study the effect of number of hours students' study for an exam (independent variable, X) on their score in exam (dependent variable, Y). Data we collected from six students:

Student	Hours Studied (X)	Exam Score (Y)
A	2.0	55
B	3.0	60
C	4.0	68
D	5.0	72
E	6.0	78
F	7	85

**Step 1: Calculate Mean of X and Y**

First, we calculate mean of X (denoted as  $\bar{X}$ ) & mean of Y (denoted as  $\bar{Y}$ ).

$$\bar{X} = (2 + 3 + 4 + 5 + 6 + 7) / 6 = 27 / 6 = 4.5 \quad \bar{Y} = (55 + 60 + 68 + 72 + 78 + 85) / 6 = 418 / 6 = 69.67$$

### Step 2: Calculate Deviations from Mean

Next, we calculate deviations of each X value from  $\bar{X}$  ( $x = X - \bar{X}$ ) & deviations of each Y value from  $\bar{Y}$  ( $y = Y - \bar{Y}$ ).

Student	X	Y	x (X - $\bar{X}$ )	y (Y - $\bar{Y}$ )
A	2	55	-2.5	-14.67
B	3	60	-1.5	-9.67
C	4	68	-0.5	-1.67
D	5	72	0.5	2.33
E	6	78	1.5	8.33
F	7	85	2.5	15.33

### Step 3: Calculate the Products of Deviations (xy) and Squared Deviations ( $x^2$ )

We then calculate the product of deviations (xy) and squared deviations of X ( $x^2$ ).

Student	x	y	xy (x * y)	$x^2$ (x * x)
A	-2.5	-14.67	36.675	6.25
B	-1.5	-9.67	14.505	2.25
C	-0.5	-1.67	0.835	0.25
D	0.5	2.33	1.165	0.25
E	1.5	8.33	12.495	2.25
F	2.5	15.33	38.325	6.25

### Step 4: Calculate Sums of xy and $x^2$

We calculate sums of xy ( $\sum xy$ ) and  $x^2$  ( $\sum x^2$ ).

$$\Sigma xy = 36.675 + 14.505 + 0.835 + 1.165 + 12.495 + 38.325 = 104 \quad \Sigma x^2 = 6.25 + 2.25 + 0.25 + 0.25 + 2.25 + 6.25 = 17.5$$

### Step 5: Calculate Slope (b) and Intercept (a)

The slope (b) of regression line is calculated as:

$$b = \Sigma xy / \Sigma x^2 = 104 / 17.5 = 5.94 \text{ (approximately)}$$

The intercept (a) is calculated as:

$$a = \bar{Y} - b\bar{X} = 69.67 - (5.94 * 4.5) = 69.67 - 26.73 = 42.94 \text{ (approximately)}$$

### Step 6: Write Regression Equation

The regression equation is:

$$\hat{Y} = a + bX$$

Where:

- $\hat{Y}$  is predicted value of Y.
- a is intercept.
- b is slope.
- X is independent variable.

In our example, regression equation is:

$$\hat{Y} = 42.94 + 5.94X$$

### Step 7: Interpret the Results

- **Slope (b):** The slope of 5.94 signifies that for each additional hour studied, exam score is anticipated to rise by an average of 5.94 points.
- **Intercept (a):** The intercept (42.94) is the estimated exam score when the number of hours studied is zero. However, in this case, this may not be meaningful as one does not read for zero hours.

- **Regression Equation:** The equation  $\hat{Y} = 42.94 + 5.94X$  can be used to predict exam scores for different study times. For example, if a student studies for 8 hours, the predicted exam score would be:  $\hat{Y} = 42.94 + (5.94 * 8) = 42.94 + 47.52 = 90.46$ .

Correlation  
And  
Regression  
Analysis

**Step 8: Assess the Goodness of Fit (R-squared):** R-squared ( $R^2$ ) quantifies proportion of variance in dependent variable that can be anticipated from independent variable. It varies from 0 to 1, with 1 signifying an ideal fit.

To calculate R-squared, we need to find sum of squares regression (SSR) and total sum of squares (SST).

$$SSR = \sum(\hat{Y} - \bar{Y})^2 \quad SST = \sum(Y - \bar{Y})^2$$

Then,  $R^2 = SSR / SST$

Using statistical software or calculators, we can determine the R-squared value for this example. A high R-squared value indicates that model fits the data well.

### **Step 9: Test the Significance of the Regression Coefficients**

Then, we can conduct hypothesis tests to check if the slope and the intercept are statistically significant. Therefore, computing t-statistics and p-values. Reject null hypothesis if p-values are below significance level (e.g., 0.05), indicating that coefficients are significant.

**Step 10: Analyze Residuals:** These residuals are the differences of actual Y and predicted  $\hat{Y}$ . Residuals analysis also assists in detecting outliers, non-linearities, and assumption violations. To check for patterns we can plot residuals against predicted values or independent variables.

**Multiple Regression:** With more than one independent variable involved, we conduct multiple regression. The process is similar, but the math gets trickier. Multiple regression analysis is typically undertaken using statistical software.

---

## UNIT 21 LEAST SQUARE FIT OF LINEAR REGRESSION

---

---

### 3.6 Least Square Fit Of Linear Regression

---

Linear regression is arguably most elementary statistical technique for modeling relationship between two variables: an independent variable (predictor) & dependent variable (target). We are doing linear regression to identify the line that optimally fits this data in terms of least squares. The predominant approach for doing this is "least squares fit" method. It aims to minimize squared sum of the discrepancies between the observed values of the dependent variable and the values predicted by linear function. These discrepancies, termed residuals, represent the errors between the model and the actual data points. This would reduce the total error: the aggregate of all squared projected errors throughout the dataset to identify the line that most accurately represents the linear connection, offering a valuable framework for analyzing or predicting trends. This foundational technique is employed across various fields, including economics, finance, engineering, & social sciences, enabling analysis & prediction of linear relationships. The derived linear equation, typically expressed as  $y = mx + b$  (where  $m$  represents slope &  $b$  denotes the y-intercept), offers a straightforward and efficient method for analyzing relationships and generating data predictions. The slope ( $m$ ) indicates variation in the dependent variable for each unit change in the independent variable, whereas y-intercept ( $b$ ) denotes value of dependent variable when independent variable is zero. We choose the least squares method because it is an optimal and unique solution and makes sure that the resulting line is the best linearization of the data. It is also mathematically tractable, familiar formulas for slope and intercept can be derived, making it feasible to do the math's manually and not just place the formula on the computational side.

***Calculating the Least Squares Line: A Step-by-Step Numerical Example***

Now, we will demonstrate how to find the least squares fit using an example with numbers. Let's say we're trying to figure out the relationship between how many hours students' study (x) & their exam scores (y): We collect following data:

Hours Studied (x)	Exam Score (y)
1	2
2	4
3	5
4	4
5	7

**Step 1: Calculate the Sums**

We first calculate the sums of x, y,  $x^2$ , and xy:

- $\Sigma x = 1 + 2 + 3 + 4 + 5 = 15$
- $\Sigma y = 2 + 4 + 5 + 4 + 7 = 22$
- $\Sigma x^2 = 1^2 + 2^2 + 3^2 + 4^2 + 5^2 = 1 + 4 + 9 + 16 + 25 = 55$
- $\Sigma xy = (1 * 2) + (2 * 4) + (3 * 5) + (4 * 4) + (5 * 7) = 2 + 8 + 15 + 16 + 35 = 76$

**Step 2: Calculate Number of Data Points (n)**

In this case,  $n = 5$ .

**Step 3: Calculate Slope (m)**

The formula for the slope (m) is:  $m = (n\Sigma xy - \Sigma x \Sigma y) / (n\Sigma x^2 - (\Sigma x)^2)$

Plugging in the values:  $m = (5 * 76 - 15 * 22) / (5 * 55 - 15^2)$   $m = (380 - 330) / (275 - 225)$   $m = 50 / 50$   $m = 1$

**Step 4: Calculate the Y-Intercept (b)**

The formula for the y-intercept (b) is:

$$b = (\Sigma y - m\Sigma x) / n$$

Plugging in the values:

$$b = (22 - 1 * 15) / 5 \quad b = (22 - 15) / 5 \quad b = 7 / 5 \quad b = 1.4$$

### Step 5: Write the Linear Equation

The equation of least squares line is:  $y = mx + b$   $y = 1x + 1.4$   $y = x + 1.4$

**Interpretation:** Our slope (where  $m = 1$ ) means that for every extra hour studied, exam score is 1 point more. The y-intercept ( $b=1.4$ ): this is the predicted amount of exam score when the student spends 0 hours studying

**Assessing the Fit:** The coefficient of determination ( $R^2$ ) can be computed to assess adequacy of line's fit to the data.  $R^2$  multiplied by 100 yields the percentage of variance in  $y$  that is accounted for by  $x$ . NOTE: A higher  $R^2$  means the regression fits data better.

### Calculating $R^2$

1. Calculate mean of  $y$  ( $\bar{y}$ ):  $\bar{y} = \Sigma y / n = 22 / 5 = 4.4$
2. Calculate total sum of squares (SST):  $SST = \Sigma(y - \bar{y})^2$
3. Calculate the regression sum of squares (SSR):  $SSR = \Sigma(\hat{y} - \bar{y})^2$  (where  $\hat{y}$  is the predicted  $y$ )
4.  $R^2 = SSR / SST$

By computing these sums and applying the formula, we can determine the  $R^2$  value and assess the goodness of fit of the linear regression model.

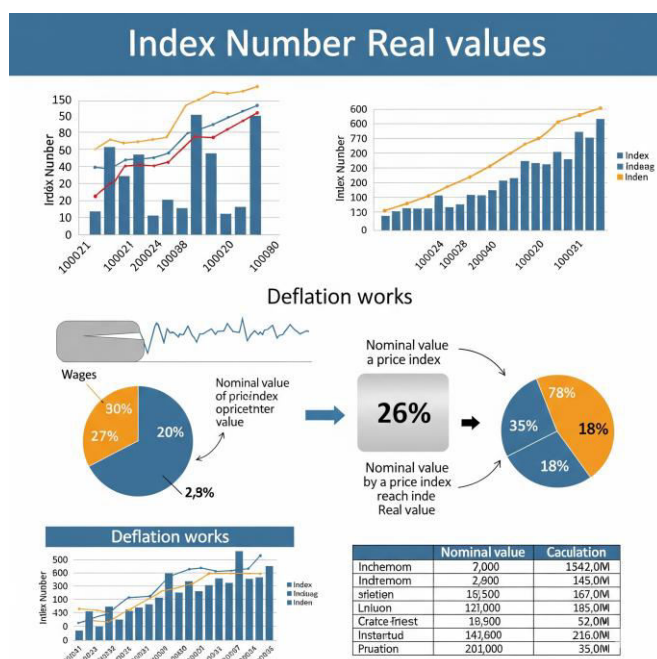
**Applications and Importance:** Least squares linear regression, used throughout many areas. In economics, it can model the correlation between GDP and unemployment. In finance, it can forecast stock prices from the market indicators. In engineering, it can study correlation between input and output variables in a system. In the world of social sciences, it can concisely describe the relationship between educational attainment and income.

## UNIT 22 TWO LINES OF REGRESSION

Correlation  
and  
Regression  
Analysis

### 3.7 Understanding Regression and its Dual Nature

Regression analysis is statistical technique used to model and examine relationship between two or more variables. For two variables, it aims to determine a line that optimally fits data points on a scatter plot, enabling prediction of one variable's value based on other variable's value. The concept of "best fit" can be understood in two distinct manners, resulting in two regression lines: the Y on X regression line ( $Y = a + bX$ ) and X on Y regression line ( $X = c + dY$ ). The regression line of Y on X is utilized to forecast the values of Y based on value of X, with X being the independent variable (predictor) and Y dependent variable (response). The regression line of X on Y is utilized to forecast the values of X based on the values of Y, with Y designated as the independent variable and X as the dependent variable. These two lines illustrate differing viewpoints of the same relationship, with the slope and intercept defining the nature and degree of that association. The mean for both variables is the intersection point of these two lines. Having an understanding of the context of the data and where you want to predict is



**Figure 3.6 Index Number Real Values**



important to identify which regression line to use. Overlap of data on those lines suggests the accuracy level of prediction.

### ***Calculating and Interpreting Two Lines of Regression: A Practical Approach with Numerical Examples***

Now I want to give you a numerical example to demonstrate the computation and meaning of two lines of regression. Let us assume we want to study relationship between number of hours students' study (X), & their exam scores (Y). We gather data from 5 students.:

Student	Hours Studied (X)	Exam Score (Y)
A	2	50
B	4	60
C	6	70
D	8	80
E	10	90

#### **1. Calculate Means of X & Y:**

- Mean of X ( $\bar{X}$ ) =  $(2 + 4 + 6 + 8 + 10) / 5 = 30 / 5 = 6$
- Mean of Y ( $\bar{Y}$ ) =  $(50 + 60 + 70 + 80 + 90) / 5 = 350 / 5 = 70$

#### **2. Calculate the Sum of Squares and Cross-Products:**

- $\Sigma(X - \bar{X})^2 = (2-6)^2 + (4-6)^2 + (6-6)^2 + (8-6)^2 + (10-6)^2 = 16 + 4 + 0 + 4 + 16 = 40$
- $\Sigma(Y - \bar{Y})^2 = (50-70)^2 + (60-70)^2 + (70-70)^2 + (80-70)^2 + (90-70)^2 = 400 + 100 + 0 + 100 + 400 = 1000$
- $\Sigma(X - \bar{X})(Y - \bar{Y}) = (2-6)(50-70) + (4-6)(60-70) + (6-6)(70-70) + (8-6)(80-70) + (10-6)(90-70) = 80 + 20 + 0 + 20 + 80 = 200$

#### **3. Calculate the Regression Coefficients:**

- **Regression Coefficient of Y on X (b):**  $b = \Sigma(X - \bar{X})(Y - \bar{Y}) / \Sigma(X - \bar{X})^2 = 200 / 40 = 5$
- **Regression Coefficient of X on Y (d):**  $d = \Sigma(X - \bar{X})(Y - \bar{Y}) / \Sigma(Y - \bar{Y})^2 = 200 / 1000 = 0.2$

#### 4. Calculate the Intercepts:

- **Intercept of Y on X (a):**  $a = \bar{Y} - b\bar{X} = 70 - (5 * 6) = 70 - 30 = 40$
- **Intercept of X on Y (c):**  $c = \bar{X} - d\bar{Y} = 6 - (0.2 * 70) = 6 - 14 = -8$

#### 5. Write the Regression Equations:

- **Regression Line of Y on X:**  $Y = a + bX = 40 + 5X$
- **Regression Line of X on Y:**  $X = c + dY = -8 + 0.2Y$

#### Interpretation:

- **Y on X ( $Y = 40 + 5X$ ):** For every one-hour increase in study time (X), exam score (Y) is predicted to increase by 5 points. The intercept, 40, represents predicted exam score when no hours are studied, though this may not be practically meaningful.
- **X on Y ( $X = -8 + 0.2Y$ ):** For every one-point increase in exam score (Y), the study time (X) is predicted to increase by 0.2 hours. The intercept, -8, represents the predicted study time when the exam score is zero, which is also not practically meaningful.

#### Using the Equations for Prediction:

- If a student studies for 7 hours ( $X = 7$ ), the predicted exam score (Y) is:  $Y = 40 + (5 * 7) = 40 + 35 = 75$ .
- If a student scores 85 on the exam ( $Y = 85$ ), the predicted study time (X) is:  $X = -8 + (0.2 * 85) = -8 + 17 = 9$  hours.

#### Important Notes:

- The regression lines should intersect at the mean values ( $\bar{X}$ ,  $\bar{Y}$ ), which in our example is (6, 70).
- The coefficients (b and d) represent the extent of change in dependent variable corresponding to a unit change in independent variable.

## UNIT 23 PROPERTIES OF REGRESSION COEFFICIENTS

### 3.8 Understanding the Foundation of Regression Coefficients

Regression analysis, a prevalent activity in statistical modeling, seeks to ascertain the response of a dependent variable (Y) to variations in one or more independent variables (X). This approach centers on regression coefficients, which indicate the amount and direction of the influence of each independent variable on the dependent variable. In a fundamental linear regression model ( $Y = \beta_0 + \beta_1 X + \epsilon$ ), the coefficients denote the Y-intercept ( $\beta_0$ , the value of Y when X equals zero) and the slope ( $\beta_1$ , the variation in Y for each unit increment in X). In these cases, the least squares method is utilized to ascertain the coefficient values that minimize the sum of squared residuals between the observed Y and the predicted values  $\hat{Y}$ . The characteristics of these coefficient values, such as unbiasedness, consistency, and efficiency, are essential for the reliability and validity of the regression model. Understanding these qualities enables researchers to make informed decisions about model selection, interpretation, and inferential implications. The coefficients are random variables which are calculated from

#### Statistical Model with Regression Coefficients

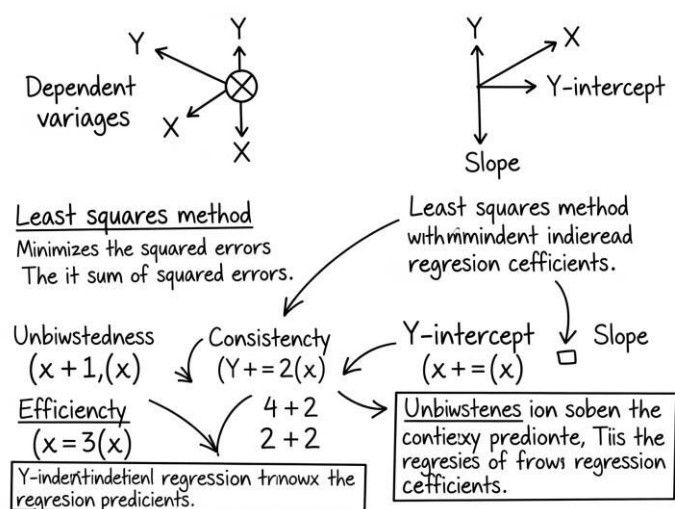


Figure 3.7 Statistical Model With Regression Coefficients

sample data, and their distributions are necessary for hypothesis testing and confidence interval construction. They are subject to the assumptions of the linear regression model (e.g., linearity, independence, homoscedasticity, normality of errors). If these assumptions are violated, the estimates may become biased or inefficient, which can affect the accuracy and generalizability of the regression outcomes.

### ***Key Properties and Numerical Illustration: Deconstructing the Behavior of $\beta_0$ and $\beta_1$***

Regression coefficients have several important properties that make them reliable and useful in statistical inference. The ordinary least squares estimators of regression coefficients are unbiased when the classical linear regression model (CLRM) conditions hold. This indicates that, on average, the predicted coefficients will correspond to the genuine population coefficients. Secondly, they exhibit consistency, indicating that as sample size rises, calculated coefficients converge to true population values. Third, they are efficient, i.e. OLS estimators have minimum variance among every linear unbiased estimator. Fourth, OLS estimators follow a normal distribution which aids in hypothesis testing and creating confidence intervals. The covariance between the estimated coefficients reveals the degree of interdependence among them. Now, let us proceed with a numerical example to put together these properties. Example: In correlation analysis, we may want to study the relation between no. of hours of study (X) and the unsigned exam scores (Y) for a group of students. We collect following data:

Student	Hours Studied (X)	Exam Score (Y)
A	2.0	60
B	3.0	70
C	4.0	80
D	5.0	90
E	6.0	100

We want to estimate simple linear regression model:  $Y = \beta_0 + \beta_1 X + \varepsilon$ .

### ***1. Calculating Regression Coefficients:***

We can calculate the regression coefficients using the following formulas:

$$\beta_1 = \Sigma[(X_i - \bar{X})(Y_i - \bar{Y})] / \Sigma(X_i - \bar{X})^2 \quad \beta_0 = \bar{Y} - \beta_1 \bar{X}$$

Where:

- $\bar{X}$  is mean of X.
- $\bar{Y}$  is mean of Y.

$$\bar{X} = (2 + 3 + 4 + 5 + 6) / 5 = 4 \quad \bar{Y} = (60 + 70 + 80 + 90 + 100) / 5 = 80$$

Now, we calculate the necessary sums:

$$\begin{aligned} \Sigma [(X_i - \bar{X})(Y_i - \bar{Y})] &= (-2)(-20) + (-1)(-10) + (0)(0) + (1)(10) + (2)(20) = 40 + \\ &10 + 0 + 10 + 40 = 100 \quad \Sigma(X_i - \bar{X})^2 = (-2)^2 + (-1)^2 + (0)^2 + (1)^2 + (2)^2 = 4 + 1 + 0 \\ &+ 1 + 4 = 10 \end{aligned}$$

$$\beta_1 = 100 / 10 = 10 \quad \beta_0 = 80 - 10 * 4 = 80 - 40 = 40$$

Therefore, the estimated regression equation is:  $Y = 40 + 10X$ .

**2. Unbiasedness:** In repeated sampling, the mean of the predicted  $\beta_1$  values would converge to the true population  $\beta_1$ . If we were to replicate sampling and estimating procedure multiple times, average of the  $\beta_1$  values would be close to 10.

**3. Consistency:** As the sample size increases, the estimated  $\beta_1$  and  $\beta_0$  values become closer to the true population values. If we collected data from a larger group of students, the estimated coefficients would be more accurate.

**4. Efficiency:** Among all linear unbiased estimators, the Ordinary Least Squares (OLS) estimators exhibit the minimal variation. This indicates that the predicted coefficients are the most accurate.

**5. Normality:** Under the CLRM assumptions, the estimated coefficients are normally distributed. This allows us to perform hypothesis tests and construct

confidence intervals. For instance, we can test null hypothesis that  $\beta_1 = 0$  (no relationship between hours studied & exam scores) using a t-test.

**6. Covariance:** The covariance between  $\beta_0$  &  $\beta_1$  indicates how they vary together. A negative covariance suggests that as  $\beta_1$  increases,  $\beta_0$  tends to decrease, and vice versa. This is often observed in regression models.

**7. Variance of the Coefficients:** The variances of regression coefficients are crucial for assessing reliability of the estimates. They are calculated as follows:

$$\text{Var}(\beta_1) = \sigma^2 / \sum (X_i - \bar{X})^2 \quad \text{Var}(\beta_0) = \sigma^2 [1/n + \bar{X}^2 / \sum (X_i - \bar{X})^2]$$

Where  $\sigma^2$  is variance of error terms. The standard errors of coefficients are square roots of these variances.

**8. R-squared and Adjusted R-squared:** Understanding R-squared and Adjusted R-squared in Statistical Modeling

R-squared ( $R^2$ ) is one of the most widely used metrics for evaluating the goodness-of-fit of statistical models, particularly in regression analysis. At its core,  $R^2$  represents the proportion of variance in the dependent variable that is explained by the independent variable(s) in the model. This metric provides analysts with a straightforward interpretation: an  $R^2$  value of 0.75 indicates that approximately 75% of the variability in the outcome can be explained by the predictor variables included in the model.

However,  $R^2$  has a fundamental limitation that necessitates caution in its application and interpretation. By mathematical construction, the  $R^2$  value will always increase or, at minimum, remain unchanged when additional independent variables are introduced to the model, regardless of whether these new variables genuinely contribute meaningful explanatory power. This property creates a problematic incentive in model building, as it can lead analysts to artificially inflate their models with superfluous variables merely to achieve a higher  $R^2$  value, potentially resulting in overfitting and reduced model generalizability.

This inherent limitation of  $R^2$  led to the development of adjusted R-squared, which incorporates a penalty for each additional predictor variable added to the model. Unlike standard  $R^2$ , adjusted R-squared increases only if the new variable improves the model more than would be expected by chance alone. In some cases, adjusted R-squared can decrease when irrelevant variables are added, providing a more reliable indicator of model quality and a safeguard against unnecessarily complex models. When applying these concepts to practical data analysis, calculating both  $R^2$  and adjusted R-squared offers valuable insights about model performance. The  $R^2$  value provides a straightforward indication of how well the model captures the variance in the dependent variable, while adjusted R-squared serves as a check against overfitting by balancing explanatory power against model complexity. Together, these metrics form an essential part of the model evaluation toolkit, although they should be interpreted alongside other diagnostic measures such as residual analysis, hypothesis tests, and information criteria for a comprehensive assessment of model adequacy.

### **9. Hypothesis Testing:**

We can conduct t-tests to ascertain statistical significance of regression coefficients. For instance, we can assess if  $\beta_1$  is statistically distinct from zero.

T-tests in hypothesis testing are essential in regression research, offering a rigorous statistical framework to ascertain whether the patterns identified in our data likely represent true linkages in the larger population or are simply due to sampling variability. In regression analysis, we derive coefficient estimates (such as  $\beta_1$ ) that quantify the associations between independent variables and the dependent variable. Nevertheless, these estimates are prone to sampling error, necessitating a methodical approach to assess their trustworthiness. The t-test for regression coefficients fulfills this requirement by enabling us to evaluate whether a coefficient significantly differs from zero. A non-zero coefficient indicates that the associated independent variable significantly influences the dependent variable, while a coefficient indistinguishable from zero signals that the variable may lack substantial explanatory power in the model.

The procedure commences with the formulation of null and alternative hypotheses. The null hypothesis ( $H_0$ ) posits that the coefficient is zero ( $H_0: \beta_1 = 0$ ), indicating an absence of correlation between the independent variable and the dependent variable. The alternative hypothesis ( $H_1$ ) posits that the coefficient is not equal to zero ( $H_1: \beta_1 \neq 0$ ), signifying the presence of a significant link. To conduct the test, we compute a t-statistic by dividing the estimated coefficient by its standard error:  $t = \beta_1 / SE(\beta_1)$ . The t-statistic quantifies the number of standard errors the calculated coefficient deviates from zero. The greater the absolute value of the t-statistic, the more compelling the evidence against the null hypothesis.

We then compare this t-statistic to critical values from the t-distribution with the appropriate degrees of freedom (typically  $n-k-1$ , where  $n$  is the sample size and  $k$  is the number of independent variables). Alternatively, we can calculate the p-value, which represents the probability of observing a t-statistic as extreme as ours if the null hypothesis were true. A small p-value (typically below 0.05) suggests that it's unlikely to observe our results by chance alone if no relationship exists, leading us to reject the null hypothesis. In business applications, these tests help determine which variables significantly influence outcomes of interest. For example, a marketing team might analyze whether advertising expenditure significantly affects sales, or a financial analyst might assess whether certain economic indicators reliably predict stock returns. By applying hypothesis testing to regression coefficients, business professionals can make data-driven decisions with quantifiable levels of confidence, distinguishing between meaningful factors and statistical noise. While hypothesis testing provides valuable insights, it's important to interpret results in context, considering practical significance alongside statistical significance, particularly when working with large sample sizes where even small effects may appear statistically significant. Additionally, multiple hypothesis testing requires appropriate adjustments to control error rates across the entire set of tests.



## 10. Confidence Intervals:

Confidence intervals provide range of plausible values for regression coefficients. They are calculated as:

$$\beta_1 \pm t(\alpha/2, n-2) * SE(\beta_1) \quad \beta_0 \pm t(\alpha/2, n-2) * SE(\beta_0)$$

Where  $t(\alpha/2, n-2)$  is critical value from t-distribution with  $n-2$  degrees of freedom.

In this post, we will cover some essential properties of regression coefficients and what they can tell you about the relationships between variables in your data. These properties are crucial to the validity and utility of regression analysis in various disciplines.

---

### 3.9 SELF ASSESSMENT QUESTION

---

Correlation  
And  
Regression  
Analysis

#### 3.9.1 Multiple-Choice Questions (MCQs)

**1. What does correlation measure?**

- a. The difference between two variables
- b. The strength and direction of the relationship between two variables
- c. The causation between two variables
- d. The average value of two variables

**2. Which of the following correlation values indicates the strongest relationship?**

- a. -0.85
- b. 0.65
- c. 0.25
- d. -0.20

**3. What does a positive correlation indicate?**

- a. One variable increase while the other decreases
- b. Both variables increase or decrease together
- c. There is no relationship between variables
- d. One variable remains constant while the other increases

**4. Which method is commonly used to measure correlation?**

- a. Standard deviation
- b. Karl Pearson's Coefficient of Correlation
- c. Moving average method
- d. Chi-square test

**5. What is the range of Karl Pearson's correlation coefficient?**

- a. -2 to 2
- b. 0 to 1
- c. -1 to 1
- d.  $-\infty$  to  $\infty$

**6. Which type of correlation does Spearman's Rank Correlation measure?**

- a. Linear correlation
- b. Non-linear correlation
- c. Rank-based correlation
- d. None of the above

**7. Which of the following is a key difference between correlation and regression?**

- a. Correlation measures dependence, while regression measures association
- b. Correlation does not imply causation, whereas regression does
- c. Correlation only describes the relationship, while regression predicts one variable based on another
- d. Correlation requires more data points than regression

**8. What does the regression equation  $Y = a + bX$  represent?**

- a. A correlation equation
- b. The relationship between independent and dependent variables
- c. The calculation of mean and median
- d. A probability distribution function

**9. What are the two lines of regression called?**

- a. Regression of X on Y and Regression of Y on X
- b. Simple regression and Multiple regression
- c. Karl Pearson's regression and Spearman's regression
- d. Linear regression and Non-linear regression

**10. What does the Least Squares Method in regression do?**

- a. It finds the median of the dataset

- b. It minimizes the sum of squared differences between observed and predicted values
  - c. It maximizes the correlation coefficient
  - d. It eliminates all errors in data
- Correlation  
And  
Regression  
Analysis

**11. Which of the following is a property of regression coefficients?**

- a. They are always greater than 1
- b. They are independent of measurement units
- c. They remain constant for all datasets
- d. They indicate the change in the dependent variable for a unit change in the independent variable

**12. Which of the following is NOT an application of regression analysis?**

- a. Predicting stock prices
- b. Finding relationships between economic indicators
- c. Calculating the mean of a dataset
- d. Forecasting business trends

**13. What is the main advantage of using regression analysis?**

- a. It helps in establishing cause and effect relationships
- b. It calculates averages quickly
- c. It eliminates errors in statistical data
- d. It ensures that correlation is always equal to one

**14. Which type of regression is used when there are multiple independent variables?**

- a. Simple linear regression
- b. Multiple regression
- c. Rank regression
- d. Exponential regression

**15. In financial forecasting, regression analysis is used to predict:**

- a. Historical stock prices
- b. Future trends based on past data
- c. Fixed values of assets
- d. The probability of an event occurring

**3.9.2 Short Questions:**

- 1. Define correlation and explain its importance.
- 2. What is the difference between positive and negative correlation?
- 3. Explain Karl Pearson's Coefficient of Correlation.
- 4. What is Spearman's Rank Correlation?
- 5. Define regression and its significance.
- 6. What are the two lines of regression?
- 7. Explain the least square method in regression.
- 8. What are the properties of regression coefficients?
- 9. How does correlation differ from regression?
- 10. What are the applications of regression analysis in business?

**3.9.3 Long Questions:**

- 1. Explain correlation analysis and its significance.
- 2. Discuss the difference between Pearson and Spearman correlation.
- 3. Explain the regression analysis with examples.
- 4. Describe the least square method and its application in regression.
- 5. What are the properties of regression coefficients?
- 6. Explain how correlation and regression are used in real-world scenarios.
- 7. Compare Karl Pearson's and Spearman's correlation methods.
- 8. What are the advantages and limitations of regression analysis?
- 9. How does correlation help in predictive analytics?
- 10. Discuss the role of regression in financial forecasting.

---

## **MODULE 4 TIME SERIES ANALYSIS**

---

### **Structure**

**UNIT 24** Introduction to Time Series Analysis

**UNIT 25** Components of Time Series

**UNIT 26** Models of Time Series

**UNIT 27** Trend Analysis

**UNIT 28** Methods of Trend Analysis

---

### **4.0 OBJECTIVES**

---

- Explain the concept, significance, and applications of time series analysis.
- Recognize and describe the different components of time series.
- Explain and compare additive, multiplicative, & mixed models of time series.
- Understand concept of trend analysis and its importance in forecasting.
- Explain and implement free hand curve, semi-averages, moving averages, and least square methods for trend estimation.

---

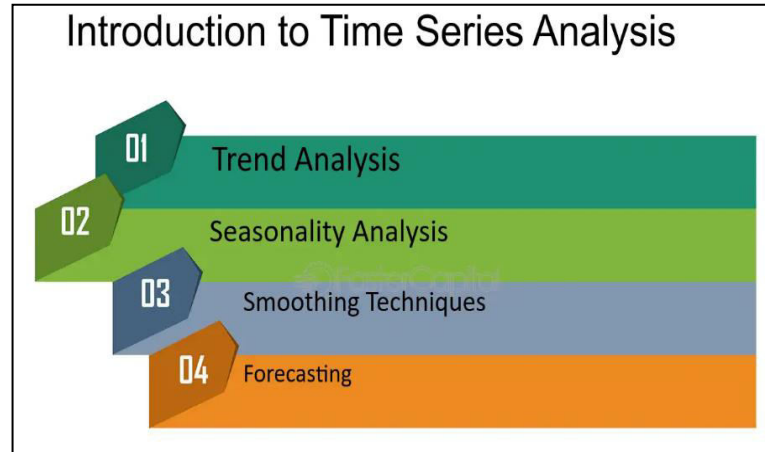
## UNIT 24 INTRODUCTION TO TIME SERIES ANALYSIS

---

---

### 4.1 INTRODUCTION TO TIME SERIES ANALYSIS

---



**Figure 4.1: Introduction to Time Series Analysis.**

Time series analysis is study of data points collected, or recorded, at specific time intervals and allows you to analyze the data point readings over time to better understand what happens in the future based on previously determined values. In contrast to cross-sectional data, which reflects a snapshot of observations at a given point in time, time-series data exposed trends, seasonality, and cyclical behavior that are endemic to temporal sequences. Such analysis is vital in many fields ranging from economics (predicting stock prices or inflation) to environmental science (weather and climate patterns) to even signal processing (understanding the variation in audio waves). A time series is a type of dependent data; for any point in time, the value will usually depend on the previous value. For a better analysis of time series, we usually decompose it into a few components: a trend (long-term movement), seasonality (repeated patterns with a fixed time interval), cyclical component (long-term variance), and random or irregular components (unpredictable noise). By comprehending these factors, we can simulate the fundamental mechanisms and generate educated forecasts. For example: retail sales may show a yearly trend of increase, seasonal peaks around holidays, and outlier drop/ups due to unexpected occurrences.

***Numerical Example: Analyzing Monthly Sales Data***

Let's illustrate time series analysis with a simple numerical example. Suppose we have monthly sales data for a small bookstore over a year:

Month	Sales (Units)
Jan	120
Feb	130
Mar	150
Apr	160
May	170
Jun	180
Jul	190
Aug	200
Sep	180
Oct	160
Nov	220
Dec	250

**1. Visualizing the Time Series:**

The first task is to plot data, specifically time series with months for x-axis and sales for the y axis. This image shows a positive line, indicating sales are better throughout the year. You also see a peak of sales in November and December, which suggests some seasonality due to holiday shopping.

**2. Identifying Trend:**

To identify the trend, we can use a moving average. A 3-month moving average smooths out short-term fluctuations & highlights the longer-term trend. For example, the moving average for March is  $(120 + 130 + 150) / 3 = 133.33$ .



Month	Sales (Units)	3-Month Moving Average
Jan	120	-
Feb	130	-
Mar	150	133.33
Apr	160	146.67
May	170	160
Jun	180	170
Jul	190	183.33
Aug	200	190
Sep	180	193.33
Oct	160	180
Nov	220	200
Dec	250	210

### 3. Detecting Seasonality:

Seasonal indices can be calculated in order to detect seasonality. For ease of calculation, let's examine the December spike. We will take the average sales across all months and compare the sales for December to this average. Monthly Average Sales:

$$(120+130+150+160+170+180+190+200+180+160+220+250)/12=184.17$$

December seasonality index =  $250/184.17 = 1.36$  This means that sales in December is about 36% more than monthly average sales.

### 4. Simple Forecasting:

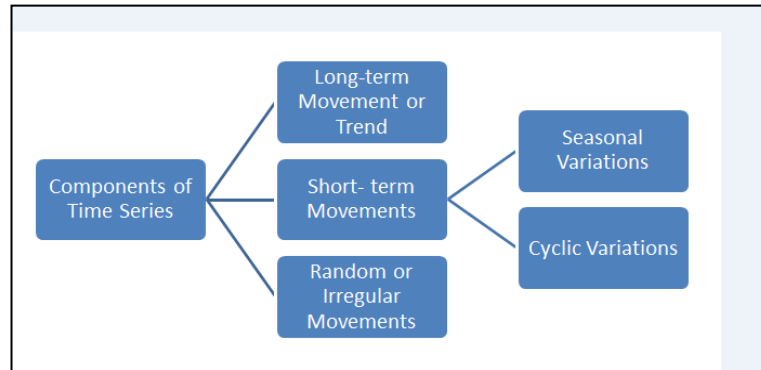
We can compute a naive forecast using trend and seasonality. Using seasonal adjustment, extrapolate up to January of the following year assuming those trends hold. But for convenience, we may also take the average of the last few months moving average, and consider slight uptrend.

**Further Analysis:** Applications for more broad-spectrum techniques such as ARIMA models, exponential smoothing, decomposition methods, can also be used for more clarified forecasting here. These are adjusted for autocorrelation, the correlation of values at different time points. This is a simple example on how time series analysis works. Analyzing load data, we train time series models to make predictions in production systems.

Time  
Series  
Analysis

## UNIT 25 COMPONENTS OF TIME SERIES

### 4.2 Unraveling Dynamics of Time-Dependent Data



**Figure 4.2: Components Concerning Time Series.**

Feature Engineering for Time series data Time series data is kind of data that is used in time series analysis which is an important analytical method that used to analyze time series data to extract interesting statistics and other characteristics data. Seemingly, this data sets are collected over time, and are coming in at regular intervals, and such data usually has complex patterns about them which can be broken down into several components. Understanding and separation of these elements are necessary for proper prognostication & rationalization of the business decision. Trend, seasonality, cyclical variations, & irregular fluctuations that are four main components of any time series. The trend refers to long-term movement of data, whether up or down, over several months or years. Seasonality is the repetitive patterns that happen on a shorter time span, like daily, weekly, monthly or yearly. Cyclical variations are long-run oscillations of indefinite frequency associated with business cycles or economic conditions. Finally, uneven oscillations (or random noise) are variations that cannot be attributed to any of the other components; they are unfurling in a random manner. Extracting these components from a time series provides us with useful information about the main mechanisms that drive the time series, helps generate better predictions, and helps develop a clearer picture of the underlying process that generates the observed results.

### *Numerical Example: Decomposing Sales Data*

Consider a company's quarterly sales data for three years (12 quarters). Let's illustrate how these components might manifest and how we can conceptualize their impact.

Quarter	Year 1	Year 2	Year 3
Q1	110	130	155
Q2	120	145	170
Q3	105	125	150
Q4	135	160	190

**1. Trend:** Note that total sales figures are increasing over the three years. This also means that the trend is a positive one. Thus, if we plot the quarterly sales, we can see the general upward slope. From week to week, it can look like a mountain range so using a simple moving average to smooth the bumps out and show the general trend helps. For example, a four-quarter moving average would smooth sales over four successive quarters, uncovering the underlying upward trend.

**2. Seasonality:** Note that Q4 always has the highest sales, while Q3 has the lowest. “Such seasonal patterns may be driven by holiday shopping-related events in Q4. We can discuss seasonal indices to quantify this seasonality. We can compute the average sales for that quarter across years and then divide it by the overall average sales. This measures the amount that seasonal effects cause an individual quarter to vary from the overall mean.

- Average Q1:  $(110+130+155)/3 = 131.67$
- Average Q2:  $(120+145+170)/3 = 145$
- Average Q3:  $(105+125+150)/3 = 126.67$
- Average Q4:  $(135+160+190)/3 = 161.67$
- Overall Average:

$$(110+120+105+135+130+145+125+160+155+170+150+190)/12 = 143.33$$

- Seasonal index for Q1:  $131.67/143.33 = 0.92$

Business Statistics	•	Seasonal index for Q2: $145/143.33 = 1.01$
	•	Seasonal index for Q3: $126.67/143.33 = 0.88$
	•	Seasonal index for Q4: $161.67/143.33 = 1.13$

These indices show Q4 sales are about 13% higher than average due to seasonality, and Q3 sales about 12% lower.

**3. Cyclical Variations:** Were this company to exist in a cyclical industry, we might witness longer-term swings beyond seasonal trends. Sales might drop off over a few years and then recover behind a broader economic downturn, for instance. Spotting cyclical fluctuations typically needs longer time series data and advanced statistical methods.

**3. Irregular Fluctuation:** After removing trend, seasonality, and cyclical variations from the data, there will be still be random variations. These may be because something unexpected happened, like a shift in consumer behavior, the unexpected success of a marketing campaign, or a supply chain problem. These variations are non-deterministic and are usually described as a random noise.

By identifying and separating these components we are able to create more accurate forecasting models. We can time-shift the data by dividing the actual sales by the seasonal indices to separate out what underlying trend is actually there. It can capture the longer-term trend as well as the repeating seasonal patterns for a better prediction of future sales.

---

## UNIT 26 MODEL OF TIME SERIES

---

---

### 4.3 MODEL OF TIME SERIES

---

Time series data which is a sequence of observations recorded over a period of time usually show complex patterns that can hide underlying trends or seasonal fluctuations. In short, we can use different techniques to decompose time series into its elements to then analyze and forecast it. These elements often consist of a trend component (long-term trend), a seasonal component (repeatable fluctuations), a cyclic component (long-term disturbances), and a residual or irregular component (random noise in general). Additive, multiplicative, and mixed models are among the common decomposition models that help determine the models as per how the components interact. The selection of model is depending on data as well as the different relationships among its constituent components. All components are assumed to be independent and additively contribute to the final outcome in the additive model. A multiplicative model multiplies the components together with dependent effects. A mixed model is a combination of both approaches, which provides a better representation for more complicated time series. This analysis offers crucial insights into the underlying dynamics, allowing businesses and researchers to be equipped with data-driven decisions and predictions based on past behavior and trends these become apparent.

#### *Additive and Multiplicative Models: Contrasting Approaches*

This algebraic equation of additive time series model for  $Y_t$  which is the value/time series is the sum or addition of Trend ( $T_t$ ), Seasonal ( $S_t$ ), Cyclical ( $C_t$ ), and Irregular ( $I_t$ ). This is ideal for seasonality when the absolute size of the seasonal variations are similar, over time, independent of the trend level. For examples, suppose monthly ice cream sales, increase or decrease by a fixed amount every year regardless of the total sales trend. This would indicate that the additive model would be appropriate.

**Multiplicative Model:** This model assumes that time series is result of components multiply together to give the time series  $Y_t = \text{Trend } (T_t) *$

Seasonal (St) \* Cyclical (Ct) \* Irregular (It). This model is suitable when amplitude of the seasonal variation's changes in proportion with trend level. For instance, multiplicative model would be more suitable if the monthly sales of a luxury product go through a more pronounced seasonal variability when sales are high and a more moderate seasonal variability when sales are low.

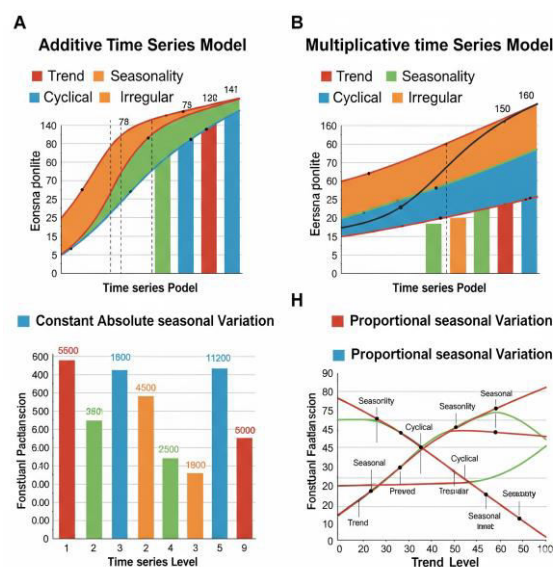
### *Numerical Example: Comparing Additive and Multiplicative Models*

Let's illustrate these models with a numerical example. Suppose we have quarterly sales data for a product over two years:

Quarter	Year 1 Sales	Year 2 Sales
Q1	110	121
Q2	120	132
Q3	130	143
Q4	140	154

#### 1. Trend Component:

First, we calculate the trend using a moving average. For simplicity, we'll use a 4-quarter moving average.



**Figure 4.3: Additive and Multiplicative Models:**

- $(110+120+130+140)/4 = 125$  Year 2:
- $(121+132+143+154)/4 = 137.5$

## 2. Seasonal Component (Additive Model):

To estimate the seasonal component for additive model, we calculate average deviation from the trend for each quarter.

- Q1:  $(110-125) + (121-137.5)/2 = -15.5$
- Q2:  $(120-125) + (132-137.5)/2 = -5.5$
- Q3:  $(130-125) + (143-137.5)/2 = 5.5$
- Q4:  $(140-125) + (154-137.5)/2 = 15.5$

## 3. Seasonal Component (Multiplicative Model):

For the multiplicative model, we calculate average ratio of actual sales to trend for each quarter.

- Q1:  $(110/125) + (121/137.5)/2 = 0.88 + 0.88/2 = 0.88$
- Q2:  $(120/125) + (132/137.5)/2 = 0.96 + 0.96/2 = 0.96$
- Q3:  $(130/125) + (143/137.5)/2 = 1.04 + 1.04/2 = 1.04$
- Q4:  $(140/125) + (154/137.5)/2 = 1.12 + 1.12/2 = 1.12$

## 4. Decomposed Values:

- **Additive Model:**
  - Year 1 Q1:  $125 - 15.5 = 109.5$
  - Year 1 Q2:  $125 - 5.5 = 119.5$
  - Year 1 Q3:  $125 + 5.5 = 130.5$
  - Year 1 Q4:  $125 + 15.5 = 140.5$
  - Year 2 Q1:  $137.5 - 15.5 = 122$
  - Year 2 Q2:  $137.5 - 5.5 = 132$
  - Year 2 Q3:  $137.5 + 5.5 = 143$
  - Year 2 Q4:  $137.5 + 15.5 = 153$



- **Multiplicative Model:**

- Year 1 Q1:  $125 * 0.88 = 110$
- Year 1 Q2:  $125 * 0.96 = 120$
- Year 1 Q3:  $125 * 1.04 = 130$
- Year 1 Q4:  $125 * 1.12 = 140$
- Year 2 Q1:  $137.5 * 0.88 = 121$
- Year 2 Q2:  $137.5 * 0.96 = 132$
- Year 2 Q3:  $137.5 * 1.04 = 143$
- Year 2 Q4:  $137.5 * 1.12 = 154$

In this simplified example, the multiplicative model exactly reproduces the original data, suggesting it is a better fit. However, real-world data is rarely this perfect.

### ***Mixed Model and Model Selection***

The mixed model is a combination of both the additive model and multiplicative model, and implementations of this model can be more complex than both components. For instance, it could assume that trend and cyclical components are additive, but seasonal and irregular ones are multiplicative. A log additive model is beneficial in cases where the data has both additive and multiplicative components. A mixed model can be articulated in several forms' contingent upon its intended application. For example,  $Y_t = T_t + S_t I_t$ . This involves examining features of the time series to identify trending behavior or seasonal patterns within it. An initial impression can be obtained through visual inspection of the time series plot. Seasonal fluctuations can be constant or can be proportional to the trend statistical tests like the F-test for homogeneity of variance can be performed in order to decide. Also, the analysis of the next residuals (the difference between the real and decomposed values) can inform us about the model chosen. If the residuals form a random pattern then model is said to be a good fit. Looking at the residuals should all be random and independent of the fitted values, if they are systematic, including being auto correlated or heteroscedastic, we need to adjust the models.

---

## UNIT 27 TREND ANALYSIS

---

Time  
Series  
Analysis

---

### 4.4 Introduction to Statistics

---

I still consider myself a newbie in this domain, but I like to know about Trend Analysis which is a statistical analysis made over time series data to identify patterns and direction. So it looks at data that gets collected regardless, at regular intervals, like daily figures on sales, monthly reports on web visitors, or annual statistics on economic metrics, so that it can analyze the trends they form and project the likely values they will have at future points. While descriptive statistics provide a summary of data at a specific moment in time, trend analysis looks at change in data over time to identify long-term trends, seasonal variations and cyclical shifts. Accurate forecasting is necessary for decision-making in many domains, ranging from business forecasts and financial planning to scientific research and social policy formulation. Through data analysis and the identification of trends, organizations can foresee challenges and opportunities on the horizon, optimize resource allocation, and implement proactive measures. A retailer, for instance, may use trend analysis to anticipate seasonal demand for goods, a financial analyst could use it to project stock prices, or a public health official may use it to monitor the spread of a disease. Time series analysis is essentially about breaking down the time-series data and separating the trend, seasonality, cycles, and noise. This allows us to decompose the time series into various components as we already see, where one often cares about the trend, which is the long-term movement in the data after removing the effects of other component. The trend (meaning up, down, or flat) tells you whether we are growing, declining, or stable. Different techniques like moving averages, linear regression, and exponential smoothing are used to model and forecast none of which have a monopoly on strengths or weaknesses.

### ***Methods and Numerical Example: Linear Trend Analysis***

Linear trend analysis is one of the easiest and popular methods for trend analysis where its assumption is the data is following a linear pattern in time.

**Linear Regression:** This method involves fitting straight line to time series by linear regression, utilizing time as independent variable & observed values as dependent variable. The equation of line is expressed as  $y = a + bx$ , where  $a$  represents y-intercept &  $b$  denotes slope. The ' $b$ ' represents the slope of the linear trend, indicating rate of change, whereas ' $a$ ' (the intercept) denotes the initial value. To have further insight, let us do a numerical example. Let us examine the subsequent sales statistics of the company over a five-year period.:

<b>Year (X)</b>	<b>Sales (Y) (in thousands)</b>
1	10
2	12
3	15
4	18
5	20

To perform linear trend analysis, we first need to assign numerical values to the years. We can simply use the year number (1, 2, 3, 4, 5) as the independent variable. Next, we calculate the necessary sums:

- $\Sigma X = 1 + 2 + 3 + 4 + 5 = 15$
- $\Sigma Y = 10 + 12 + 15 + 18 + 20 = 75$
- $\Sigma X^2 = 1^2 + 2^2 + 3^2 + 4^2 + 5^2 = 55$
- $\Sigma XY = (1 * 10) + (2 * 12) + (3 * 15) + (4 * 18) + (5 * 20) = 249$
- $n = 5$  (number of data points)

Now, we can calculate slope ' $b$ ' and the intercept ' $a$ ' using the following formulas:

- $b = (n\Sigma XY - \Sigma X\Sigma Y) / (n\Sigma X^2 - (\Sigma X)^2)$
- $a = (\Sigma Y - b\Sigma X) / n$

Plugging in the values:

- $b = (5 * 249 - 15 * 75) / (5 * 55 - 15^2) = (1245 - 1125) / (275 - 225) = 120 / 50 = 2.4$
- $a = (75 - 2.4 * 15) / 5 = (75 - 36) / 5 = 39 / 5 = 7.8$

Therefore, the linear trend equation is  $y = 7.8 + 2.4x$ . This equation indicates that the company's sales are increasing by 2.4 thousand units per year, with a starting point of 7.8 thousand units. To forecast sales for the next year (Year 6), we can plug in  $x = 6$ :

- $y = 7.8 + 2.4 * 6 = 7.8 + 14.4 = 22.2$

Thus, the forecasted sales for Year 6 are 22.2 thousand units. This method provides a simple and effective way to estimate and project linear trends, but it's important to note that it assumes a constant rate of change, which may not always hold true in real-world scenarios.

### ***Beyond Linearity: Advanced Trend Analysis Techniques***

linear trend is a great fit for simple datasets, most time series in the real world exhibit more complex trends. These complexities require advanced techniques to capture them. For example, moving averages smooth out short-term fluctuations by averaging data points over specified period. By averaging, we mitigate random noise and may spot hidden trends. Where exponential smoothing applies exponentially decreasing weights to past observations, focusing more on recent observations. This method is especially effective at predicting time series that has trends and seasonality. Statistical Methods for Logistic Regression Seasonal Decomposition Seasonal decomposition is an effective technique employed to disaggregate time series into its constituent components: trend, seasonal, & residual elements. This allows analysts to examine each individual segment without deciphering concealed meanings in the data. As an example, a retailer can use seasonal decomposition to analyze sales data and determine the seasonal peaks and troughs. techniques such as spectral analysis and wavelet analysis can also be applied to cyclical fluctuations that are essentially long-term variations of the trend. Such

techniques enable the classification of periodic patterns and project future cycles. Apart from these classical methods, various machine learning techniques like ARIMA (Autoregressive Integrated Moving Average) and neural networks are also being used for trend analysis. Such ARIMA models tend to capture the autocorrelation and moving average components while neural networks are able to learn complex non-linearities. These advanced methods offer more precise forecasts and insights, particularly for intricate and fluctuating time series. They do, however, also need more computational resources and expertise. Assessing trend analysis accuracy is key to making accurate predictions. Different metrics, like mean absolute error (MAE), root mean square error (RMSE), and mean absolute percentage error (MAPE), are commonly used to measure the discrepancy between forecasted and actual values. Introduction In order to decide what trend analysis method to use, analysts can compare the performance of various trend analysis methods. The data characteristics, accuracy requirements and resource availability will determine the appropriate trend analysis method. Time-series analysis is a powerful tool that can be used to extract insights from a wide range of data sources, and by understanding different methods available, analysts can better leverage these techniques to inform their decision-making process.

---

## UNIT 28 METHODS OF TREND ANALYSIS

---

Time  
Series  
Analysis

---

### 4.5 The Significance of Trend Analysis

---

Trend analysis is an important statistical approach that is used to analyse the pattern and direction of time series data. Analyzing trends involves discerning patterns and trends in values recorded over time, usually during regular intervals. This is critical across different fields, including economics and finance, environmental science, and marketing. By identifying long-term movements, cyclical variations and seasonal fluctuations, businesses can forecast sales, governments can plan infrastructure and researchers can gain an understanding of changing phenomena. Trend analysis allows us to identify the signal from the noise the basic trend that a dataset is following and predict where it might head in the future. This data however is crucial for the comprehension of the past, present and possible future of datasets, it is inevitable. There are multiple approaches to accomplish this, which vary in benefits and constraints, and are more or less suitable for various data types and analytical requirements.

#### ***1. Free Hand Curve: A Visual Approach to Trend Identification***

Logically, the easiest and subjective method of trend analysis is the freehand curve method. These involve plotting time-series data and drawing a graph by hand, a smooth curve which best fits the general trend. This quick and simple method requires no complex calculations, suitable for a preliminary overview or with small datasets. But its system is subjective, so different analysts might draw different curves and thus get different results. As an example, take the yearly sales figures of a small book shop for 5 years: [20, 25, 30, 35, 40]. If we plot these points and fit a line that tends to follow the upward direction, we can get a rough idea of the trend in sales. Although it is useful for a preliminary overview, it is not precise and objective as more sophisticated ways. It is most useful for a rapid first pass at the data, most specifically when a back-of-the envelope sense of the trend is all that is required.

#### ***2. Semi-Averages Method: Simplifying Trend Calculation***

The semi-averages method tries to add more objectivity into trend analysis, for each half, you need to calculate the average value immediately. Averages are computed and then plotted at the midpoint of their respective time periods, with a straight line drawn between them. This line shows the trajectory. For example, you may have ten years' worth of sales data: [10, 12, 15, 18, 20, 22, 25, 28, 30, 32]. Splitting it like this leads us to [10, 12, 15, 18, 20] and [22, 25, 28, 30, 32]. They're averaging 15 and 27.4, respectively. Plotting these averages at the midpoints of their halves and drawing a connecting line gives a trend line. This method is easy and straightforward and also less subjective in comparison with a custom freehand curve. Yet, it assumes a linear behavior and it may not eventually reflect more complex behavior. It is handy when you need a fast, less subjective approximation of a linear trend.

### ***3. Moving Averages Method: Smoothing Out Fluctuations***

The moving averages method is another highly popular method, which allows smoothing out the noise/volatility in the data and highlight the general direction in a long-term. The employed technique is moving average, which computes average value of specified number of successive data points. That average is then displayed at the halfway point of the period that the average covers. The more number of data points you take for the average, smoother will be the trend line. For instance, for the sales data [10, 12, 15, 18, 20, 22, 25, 28, 30, 32], we compute three-year moving averages like  $(10+12+15)/3 = 12.33$ ,  $(12+15+18)/3 = 15$ , etc. Plotting these averages shows a smoother trend line than the raw data. Moving averages method is the most common technique used to smooth the data as it effectively smooths with time ahead and helps to identify the long-term trend by reducing the impact of random variation. However, it may lag actual data especially during periods of rapid change and does not correspond to trends for the beginning or end of the time series. Choosing the moving average period is important and should be based on characteristics of the data & desired amount of smoothing.

#### ***4. Least Square Method: Precise Trend Line Fitting***

Time  
Series  
Analysis

The least squares method is statistical technique that determines the optimal straight line by reducing total of squared deviations between observed data points & line. Its accuracy based solely on math's, unlike always subjective based judgments. Trend-related equations are typically expressed as:  $y = a + bx$ , where  $y$  represents predicted value,  $x$  denotes time period,  $a$  signifies the  $y$ -intercept, and  $b$  indicates the slope. Let us examine data set [5, 8, 10, 12, 15] as an example. The slope & intercept of optimal line can be determined using the least squares approach. The slope signifies the pace of variation. While the intercept refers to the starting value. A method often used for forecasting, trend analysis, particularly when it is assumed that there is a linear trend; the method is quite accurate. Because it is often computationally expensive and may not perform well with nonlinear trends. If accuracy and objectivity are paramount, as is the case with most statistical applications, use the least squares method that produces a trend line with the strongest statistical characteristics. The least squares method is a widely used statistical technique for determining the optimal straight line that best fits a given set of data points. It is primarily employed in regression analysis and trend forecasting to establish a mathematical relationship between dependent and independent variables. By minimizing the sum of the squared deviations between observed data points and the fitted line, the least squares method ensures an optimal representation of the data trend.

Unlike subjective judgment-based methods, which may introduce bias or inconsistency, the least squares method relies purely on mathematical principles. This makes it a preferred approach for analysts and researchers who seek objective and statistically robust models for decision-making.



---

## 4.6 SELF ASSESSMENT QUESTION

---

### 4.6.1 Multiple Choice Questions (MCQs)

**1. What is Time Series Analysis?**

- A) The study of historical data to identify patterns over time
- B) The process of calculating averages of unrelated data
- C) A method used only for financial forecasting
- D) A technique to collect survey data randomly

**2. Which of the following is NOT a component of time series?**

- A) Trend
- B) Seasonality
- C) Random Variations
- D) Hypothesis Testing

**3. In an additive time series model, how are the components combined?**

- A) Multiplication
- B) Subtraction
- C) Addition
- D) Division

**4. Which of the following is an example of a multiplicative time series model?**

- A)  $Y = T + S + C + R$   $Y = T + S + C + R$
- B)  $Y = T \times S \times C \times R$   $Y = T \times S \times C \times R$
- C)  $Y = (T + S) \times C$   $Y = (T + S) \times C$
- D)  $Y = T - S - C - R$   $Y = T - S - C - R$

**5. What does the Free-Hand Curve method help in identifying?**

- A) Cyclical variations
- B) Trend component
- C) Seasonal variations
- D) Residual error

**6. What is the Semi-Averages method used for?**

- A) To calculate moving averages
- B) To split data into two equal parts and find trends
- C) To analyze cyclical variations
- D) To measure seasonal effects

**7. In the Moving Average method, what happens when the window size increases?**

- A) The trend line becomes smoother
- B) The fluctuations increase
- C) The seasonal variations become more prominent
- D) The analysis becomes less reliable

**8. The Least Squares Method is primarily used to:**

- A) Find the relationship between two independent variables
- B) Fit a trend line to historical data
- C) Remove seasonal fluctuations
- D) Analyze random variations

**9. Which of the following is a major application of time series analysis?**

- A) Medical research
- B) Forecasting future sales
- C) Analyzing survey responses
- D) Predicting election results

**10. Why is Time Series Analysis important in forecasting?**

- A) It identifies trends and patterns in historical data
- B) It eliminates all fluctuations in data
- C) It removes randomness from financial markets
- D) It guarantees accurate future predictions

**11. What is the primary objective of Trend Analysis in time series?**

- A) Identifying long-term movement in data
- B) Removing seasonal fluctuations
- C) Adjusting cyclical variations
- D) Predicting short-term random changes

**12. Which of the following is NOT a trend analysis technique?**

- A) Free-Hand Curve Method
- B) Semi-Averages Method
- C) Regression Analysis
- D) Monte Carlo Simulation

**13. In which sector is Time Series Analysis widely used?**

- A) Financial markets
- B) Meteorology
- C) Sales forecasting
- D) All of the above

**14. How does Time Series Analysis help in stock market predictions?**

Time  
Series  
Analysis

- A) By ensuring future stock prices
- B) By identifying historical patterns and trends
- C) By eliminating market risks
- D) By removing external economic factors

**15. What is a common challenge in Time Series Forecasting?**

- E) Data is always accurate
- F) Market trends remain constant
- G) Presence of random variations and external factors
- H) Lack of statistical models

**4.6.2 Short Questions:**

1. What is time series analysis?
2. Explain the different components of a time series.
3. What is the difference between additive and multiplicative models?
4. Describe the free-hand curve method for trend analysis.
5. What are semi-averages in time series analysis?
6. How is the least square method used in trend analysis?
7. What are the applications of time series analysis?
8. How does time series analysis help in forecasting?
9. What is the importance of trend analysis?

**4.6.3 Long Questions:**

1. Explain time series analysis and its significance.
2. Describe the different models used in time series analysis.
3. Discuss the various methods of trend analysis with examples.
4. Explain the least square method and its application in time series.
5. What are the advantages of using moving averages in trend analysis?

6. How does time series analysis help in business forecasting?
7. Compare the different trend analysis techniques.
8. Discuss the impact of time series analysis on financial decision-making.
9. Explain the role of trend analysis in stock market predictions.
10. What are the challenges in time series forecasting?

---

## **MODULE 5 DECISION THEORY**

---

### **Structure**

<b>UNIT 29</b>	Introduction to Decision Theory
<b>UNIT 30</b>	Decision Making Under Certainty
<b>UNIT 31</b>	Construction of Decision Trees

---

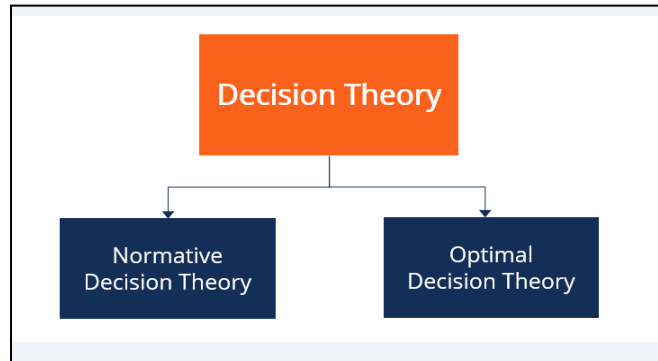
### **5.0 OBJECTIVES**

---

- Explain the concept, significance, and applications of decision theory in problem-solving.
- Understand and apply decision-making principles in situations with known outcomes.
- Develop decision trees to visualize and evaluate different decision-making scenarios.

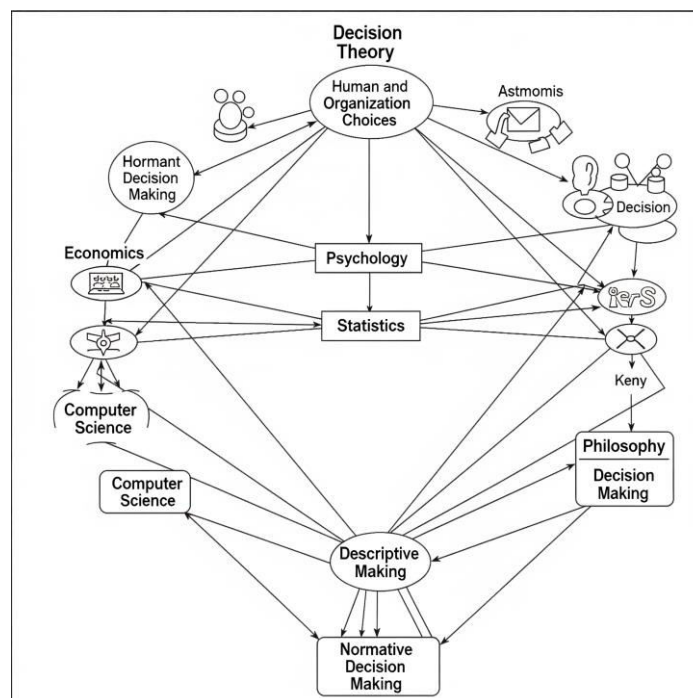
## UNIT 29 INTRODUCTION TO DECISION THEORY

### 5.1 Defining Decision Theory and Its Relevance



**Figure 5.1: Decision Theory**

At a basic level decision theory is the study of how humans and organizations make choices. It's an interdisciplinary field, pulling from economics, psychology, statistics, philosophy, computer science and others. It attempts to understand the processes that underlie decision-making, which can include both descriptive (how people actually decide) and normative (how people



**Figure 5.2: Relevance of Decision Theory**

should decide). We start with the innate complexity of choice. We are constantly faced with decisions in life, from the mundane and quotidian (what to eat for breakfast) to the profound and life-changing (career choices, investments, etc.). This article is suggested by Decision theory. The basic idea is that decisions are made in face of uncertainty. We seldom know enough about the consequences of our decisions. You may not know everything you need to know to make predictions, or events may defy predictions, or other beings may choose actions that create uncertainty in the future, even with optimal knowledge. Decision theory uses elements like probabilities, utilities, and risk to avoid becoming mired in uncertainty. Utilities are used to reflect the expected value or satisfaction coming from particular scenarios, while probabilities show how likely those scenarios are. Risk, in turn, represents the possibility of downside.

We need to distinguish descriptive from normative decision theory. Influencing behavioral decision making, the descriptive decision theory analyses how people do make decisions and commonly describe biases and irrationalities. For instance, the field of behavioral economics has revealed psychological phenomena, such as loss aversion, which is the fact that we feel more pain from losses than we derive pleasure from equivalent gains. In contrast, normative decision theory sets out how one should decide to achieve the most preferred outcomes, generally in a rational manner. This method is based on principles such as expected utility maximization, which takes into account the potential results of each decision and balances them in accordance to their probabilities and utilities. Decision theory is not just an ivory tower exercise; it is more than a lot of theorems stated without proof; it has strong real life implications. In business, it guides strategic planning, investment decisions and risk management. In medicine, it informs treatment decisions and public health policies. In A.I. it forms the foundation for the creation of intelligent agents capable of makes autonomous decisions. It can guide us to better decision-making in day-to-day scenarios.

**Key Concepts to Introduce and Elaborate:**



- **Decision-Making Process:** The steps taken in decision-making process such as recognizing issues, gathering information, developing options, evaluating options, and making a decision and reviewing it.
- **Rationality:** The idea of making economic decisions that are aligned with your preferences and values. Embrace the imperfection of rationality and understand bounded rationality.
- **Uncertainty and Risk:** Understanding the difference between uncertainty (when the outcomes are not known) and risk (when the probabilities of outcomes are known). How could we collaborate to identify types of risk (financial, operational, etc.)
- **Chance:** The probability of an eventuates happening. Introduce subjective probability and objective probability.
- **Utility:** The subjective value or satisfaction associated with an outcome. Paraphrase
- **Expected Value and Expected Utility:** Teaching how to compute expected value (average outcome) and expected utility (average satisfaction).
- **Decision Trees:** A visual decision-making process used to examine possible outcomes.
- **Real Life Examples:** Give examples of how decision theory is used in business, finance, medicine, public policy, etc.
- **Cognitive Biases:** You can explain cognitive biases and how they impact your decision making. For example availability heuristic, anchoring bias, confirmation bias.

## ***II. Navigating the Unknown: Tools and Frameworks in Decision Theory***

Now that we have established foundational knowledge, we can go into some of the core tools and frameworks possessed by the field of decision theory that we can leverage to analyze and improve decision process. This is where you apply your theoretical learnings in practice. Perhaps, one of the simplest foundational tools is a decision matrix, where you line up potential choices, their possible outcomes, and the relative utilities or payoffs. It facilitates an structured comparison of options' A company

deciding whether or not to launch a new product could, for instance, build a decision matrix that lays out the potential outcomes (success, moderate success, failure) against the profits or losses for each scenario.

Bayesian decision theory updating and the test outcome. Sequential decision-making, where decisions are made based on an evolving body of information, is a key application for Bayesian beneficial especially when not all the information is available or well-defined. Examples would include like how a physician diagnosing a patient would use Bayesian reasoning to revise probability of a disease based on the symptoms that the patient presents with evidence to update probabilities. This method is a complementary powerful framework, incorporating previous beliefs and new and prisoner's dilemma help in understanding how individuals and organizations behave in strategic situations. interactions (e.g., in auctions, negotiations, or competitive markets). Concepts from game theory such as the Nash equilibrium circumstances with multiple decision-makers that might have conflicted or aligned goals. It studies strategic Just as decision theory studies choice under uncertainty, game theory generalizes it to MCDA methods allow for prioritization and weighting of these objectives. involves the location of a factory, where you decide based on cost, environment, and nearness to customers, etc. Tools such as conflicting objectives. An instance Simultaneously, multi-criteria decision analysis (MCDA) addresses decision-making involving many criteria random sampling analysis studies the effect of varying inputs on outputs, whereas scenario planning investigates possible future scenarios and their consequences. Monte Carlo simulation models the probability of different outcomes with in business refers to the variability of future outcomes, and methods while quantifying and managing risk include sensitivity analysis, scenario planning, and Monte Carlo simulation, etc. For Looking Back — Sensitivity Analysis and Scenario Planning: Sensitivity in decision theory. Risk analysis is a fundamental discipline.

### **Key Concepts to Introduce and Elaborate:**

- **Decision Matrices:** Constructing and interpreting decision matrices

- **Bayesian Decision Theory:** Bayes' theorem, prior and posterior probabilities, belief updating.
- **Game Theory:** Nash equilibrium, prisoner's dilemma, strategic interactions
- **Multi-Criteria Decision Analysis (MCDA):** Weighting of criteria, scoring of alternatives, ranking approaches.
- **Risk Analysis:** Sensitivity, scenario, Monte Carlo.
- **Value of Information** - What is the cost of obtaining further information?
- **Information Systems:** Role of technology in decision support.
- **Real World Examples:** Instances of accurate techniques in respective fields.

### ***III. The Human Element: Behavioral Insights and Ethical Considerations***

and psychology that people frequently diverge from rationality, often as a result of cognitive biases, emotions and social influences. bases decisions on cold calculations and rational choices, but we must remember the humanity behind it all. It has been shown by behavioral economics Normative decision theory to combat them. on the first information given), and loss aversion (the tendency to prefer avoiding losses over acquiring equivalent gains). By being aware of these biases, we can make better choices and design interventions The field of behavioral decision theory delves into the nature of these deviations, examining conceptual occurrences such as framing effects (the impact of how a decision is framed on the decision), anchoring bias (the tendency for an individual to rely too heavily.

Emotions drive many of the decisions we make. These feelings of fear, regret, excitement, can affect our choices; sometimes in even irrational ways. Decision theory asset us to understand and navigate these emotional ensnarement's. Social bonds also affect our choices. Meaning, we are affected by what other people think of us and do, as well as what others say is right or wrong. Decision theory can help us make sense of how

these social influences impact our decisions. Ethical considerations are paramount in decision-making. Any decision we make has the potential to affect either others or society greatly, and as such, we need to also therefore be wary of the ethics of our decisions. For example, the principles and values that should dictate the choices we make can be framed using decision theory.

Also Important are Long-term vs. Short-term decisions. Most decisions are made on the basis of immediate gratification; however, the best decision may be the one that'll give the best outcome in the long run. Consideration of decision theory allows us to narrow down a preferred long term action.

---

## UNIT 30 DECISION MAKING UNDER CERTAINTY

---

---

### 5.2 Decision Making Under Certainty

---



**Figure 5.3: Decision-Making under Conditions of Certainty**

Decision theory, a cornerstone of rational choice, provides a framework for understanding and analyzing how individuals and organizations make choices in the face of uncertainty. It is a deep dive into the ways that we assess choices, consider the potential consequences, and finally make a decision that is consistent with our objectives. Basically, decision theory is the systematic study of decision-making, making choices that maximize the expected payoff and minimize the expected loss. It is a trans-disciplinary field that spans economics, psychology, statistics, philosophy, artificial intelligence, management, etc. The written word is the most efficient route for conveying a structured framework down to addressing complex matters, whether components of everyday living, enhancing strategic objectives or critical planning decisions. Both decision theory and HJB theory are not based on the idea that decision making is a random occurrence, but that we are deliberate in our choices given our beliefs, preferences, and available data. It can help codify these influences so that we can construct models to predict and prescribe the best choices. Decision theory starts with some basics: Alternatives, outcomes, probabilities and utilities. Alternatives are the actions or decisions that the decision-maker can take or make, each with different outcomes. Outcomes are the results of these events and can be known outcomes or unknown outcomes. Probabilities measure how likely each

outcome is to happen, capturing the decision-maker's beliefs about how the world works.

Decision  
Theory

Utilities are, instead, the subjective value or desirability of each outcome and therefore embody the preferences of the decision-maker. Decision-making can be roughly defined as the process of selecting the alternative that maximizes expected benefit, influenced by many factors. This entails calculating the weighted average of the utilities of all potential outcomes, with the weights corresponding to the probability of those possibilities. Decision theory distinguishes between decisions made under certainty, risk, and uncertainty. **Decision-Making under Conditions of Certainty** Decision-making under certainty pertains to scenarios where the outcomes of all alternatives are unequivocally known. While this is a rather basic situation, it serves as a foundation for more complex cases. Decision-making under risk refers to circumstances where the outcome is uncertain, but the probabilities of outcomes are known or can be estimated. This is the most basic situation covered in decision theory, where on the basis of expected utility a concrete conclusion is drawn. How to make decision under uncertain -- the situations where the results are not guaranteed, and the redundancies of these results are nothing but guess or estimation that may or may not work. This becomes quite a task since expected utility calculations cannot be applied normally. A number of different approaches have been devised for this, including subjective probabilities, robust decision-making, and ambiguity aversion. The first examines deductive normative approaches, while the second explores a variety of both normative and descriptive approaches. Normative decision theory is an attempt to tell rational people how to make decisions according to rules of logic and axioms. It sets up an ideal standard of decision making, thereby giving a yardstick to measure reality against. In contrast, descriptive decision theory fares an attempt to characterize the way people really make decisions, often admitting that human behavior is irrational. It integrates psychological elements, including cognitive biases and heuristics, to understand where such deviations arise. We are all taught the great key concepts of decision theory the principle of dominance, which is when rational decision-makers will always choose the option that is best in all states

of the world, and so on. They can be used to describe very different preferences of decision-making: the transitivity axiom states that if a decision-maker prefers alternative A over alternative B and B over alternative C, then A must be preferred over C too, whereas the independence axiom states that preference between A and B must not change if a third alternative, not relevant to the choice, is included. These principles underlie rational choice theory, which posits that rational beings make consistent and coherent choices.

Decision trees and influence diagrams are two important tools used to help people understand decision problems and to analyze complex scenarios in decision theory. They can help us understand decision trees, which are graphical representations of the decision situation, explaining the sequence of decisions, chance events, and the resulting outcomes. They are especially useful for sequential decision problems where the outcome of one decision impacts future decisions. Other than Influence diagrams highlight the relationships among the variables, decisions, and outcomes, showing the dependencies and the flow of information, They are useful for the study of complex systems with multiple causes interacting. Game theory, a closely related field, generalizes decision theory to cases with multiple decision-makers with conflicting or aligned interests. It studies strategic interactions, where the payoff of one decision maker's action depends on the actions of others. Game theory explains competitive and cooperative behavior, with applications in fields from economics and political science to evolutionary biology. Behavioral decision theory takes insights from psychology to explain how cognitive. It recognizes that human decision making may not always be rational in the sense of expected utility theory. Such biases include framing effects -- when the way a problem is presented makes a difference to the choices made; anchoring effects -- as when the first piece of information received biases subsequent judgments; and availability heuristics, when information that comes to mind easily is overweighted.

These perceptual and cognitive biases can introduce or exacerbate systematic errors in how we make important decisions; and so, they are in danger of

being misunderstood or misapplied, highlighting the need for a thorough understanding of the sources and influences of these sugars. Decision theory also investigates the phenomena of risk aversion, where individuals prefer known risks over unknown risks, given the same expected value. Individual preferences, cultural factors, and situational context affect people's risk attitudes. Another area of focus is making decisions under ambiguity, where probabilities are unknown or uncertain. Ambiguous Aversion: Likely to avoid from options with unknown probabilities even when the expected utility is likely the same as options that have known probabilities. Robust decision making is concerned with making decisions under deep uncertainty; where the probabilities of the outcomes are poorly understood. It means creating strategies that will prove robust to a broad range of potential futures, instead of aiming for accurate predictions. This obviously include new advancements and ideas from various fields. It offers an empowering platform for understanding and improving decision-making across a diverse scope of frameworks. These are the key to better decisions leading to improved outcomes, whether they be individual or organizational. Except that an instructional process that is prescriptive (top-down rules) does not allow for any abductive reasoning about shared context between multiple disciplines. We live in a time of uncertainty and complexity, and in such an environment, decision theory can serve as an important guide for how we approach the challenges, challenges will face us, and opportunities ahead, that we need rational and effective decision-making.



---

## UNIT 31 CONSTRUCTION OF DECISION TREES

---

---

### 5.3 Construction of Decision Trees

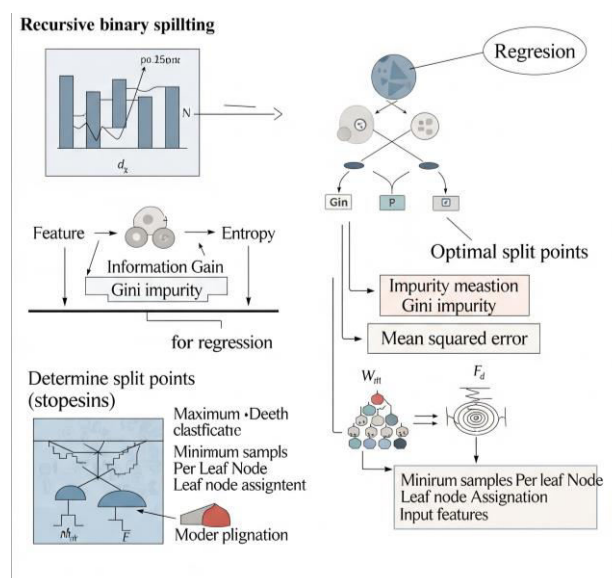
---

**Core Principles:** One powerful tool within business analytics is decision trees, a visual representation of the decision-making process, including potential outcomes, probabilities, and costs associated with each choice. They are constructed based on a recursive partitioning scheme, where the data is divided according to values of attributes that maximize information gain or minimize some measure of impurity. This begins from a root node that contains the entire data set and divides into internal nodes, which represent decision points around a specific attribute. The leaf nodes, which are the terminal points, represent the final outcomes, classified according to their respective categories or numerical values.

For this reason, the primary objective is an accurate model to predict outcomes in addition with interpretability so that the business could comprehend why decisions are made. One of the common algorithms for this construction is called decision trees, which utilize the chosen splitting criteria, such as Gini impurity or entropy for categorical variables and variance reduction for numerical variables, to choose what attributes at each node provides the most information. The basic idea behind pruning is an application of techniques, such as cost-complexity pruning, to reduce the complexity of the model and help ensure the model does not overfit to the train dataset, and does well on previously unseen data. The structure of the tree is constructed in an iterative manner, where all possible splits are evaluated, and the one that separates the outcomes best is selected, and this is done until some stopping condition is reached, such as minimum number of samples in a leaf node or maximum depth of the tree. This yields what we call the decision tree: a clear, hierarchical decision space allowing the organization to see the risk versus reward of each of the decisions.

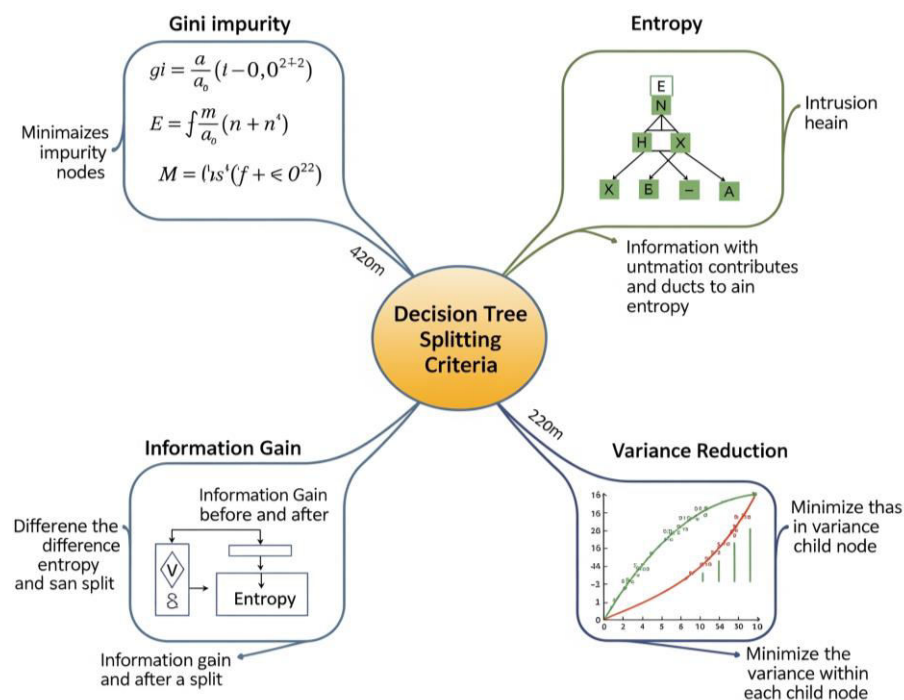
**Data Processing and Preprocessing:** Before any decision tree is made, it is essential to start from quality input data. Until the construction, the data should be cleaned and preprocessed carefully. This includes dealing with missing values through imputation or removal, addressing outliers which can

skew the model and transforming variables when required. Feature engineering is key, where you can create new features based on the existing ones to improve predictive power. Data cleaning is a form of organization in its purest form, ensuring consistency and accuracy while eliminating duplicates and errors. Depending on the data preprocessing that one applies, categorical variables are transformed into numeric values (like one-hot encoding or label encoding) to make it easier to work with them. Dimensionality reduction methods, such as feature selection, can help you focus on most relevant features to improve model's performance. The datasets can be divided into training and testing datasets, training datasets are used to construct the tree while testing datasets are used for testing the constructed tree statistically. This split allows to cover on model generalization to unseen data and prevent overfitting. You assess the distribution of classes within the dataset, and you may employ techniques like oversampling or under sampling to balance imbalanced datasets, ensuring that all classes are adequately represented in the model. If there are a number of different numerical features that are on very different scales, data normalization or scaling may be necessary since some of the splitting criteria can be sensitive to feature magnitude. The Preprocessing phase is an iterative one and might need to be adjusted as you try to fit and test your model.



**Figure 5.4: Decision Splitting**

**Selection of Splitting Criteria:** This is one of most important aspects of a decision tree. For categorical type target variables; Gini impurity and entropy are widely adopted. Gini impurity estimates the likelihood of mislabeling a randomly chosen item if it is randomly labeled according to the distribution of labels in the subset. Gini impurity: a lower value means a more homogeneous subset. Entropic, on the other hand, measures the unruliness or randomness in a fraction. This change in entropy (less entropy value) when we split on a specific attribute is termed as information gain; information gain is derived from entropy. The maximum information gain is chosen as the splitting criterion. Variance reduction is commonly used for numerical target variables. This is based on the variance reduction when dividing the node according to a certain attribute.



**Figure 5.5: Decision Tree Splitting Criteria**

As in decision trees, the attribute resulting in maximum reduction of variance is chosen. Split can also be assessed using other criteria, for example chi-square test. The decision

between splitting criteria depends on the nature of dataset & particular aspect of analysis that one is interested in. Gini impurity, for example, is computationally faster than entropy, and therefore well suited for handling very large datasets. Choosing the splitting criterion is a canonical step in the construction process that significantly affects the capability of the tree to accurately classify or predict outcomes. For each potential split, the criteria are calculated and the split that creates the maximum of the selected criterion is used.

**Tree Growth and Pruning:** the growth of the tree is similar to building the database recursive partitioning the data, and stops when a criterion is met. Some common criteria include minimum leaf node sample, maximum tree depth, maximum number of leaf nodes. Because decision trees are prone to overfitting the training data when pruning is not applied, this often results in weak model performance in terms of generalizing to unseen data. Techniques used to prune trees to avoid overfitting. One popular approach for pruning Decision Trees is cost-complexity pruning, also referred to as weakest link pruning. It introduces a complexity parameter -  $\alpha$  - that governs the balance between accuracy and size of the tree. The algorithm pruning begins by cutting off the weakest link, that is, the node that provides the least amount of error reduction, and continues until the desired pruning level. The value of  $\alpha$  is typically optimized through cross-validation to strike a balance between bias and variance. There are also other pruning methods like Lower Error Pruning and Pessimistic Error Pruning which help trim the tree by removing those nodes that do not yield significant improvement. The second tree is simpler and even more interpretable than the first tree, thus it will be easier to understand and keep in mind while applying it to business decisions.

**Evaluation and Interpretation:** To provide the optimal information for the system, proper data running strategies should be in place. Evaluation Metrics are based on types of target variable. Common metrics for categorical variables include: accuracy, Precision, recall, F1-score, and area under receiver operating characteristic curve (AUC). Accuracy quantifies the proportion of correctly classified cases. Precision is ratio of accurately

anticipated positive instances to the total expected positive instances. Recall: The proportion of True Positives to Total Positives. Precision and recall are derived from F1-score, which represents harmonic mean of both metrics. AUC represents a comprehensive measure of performance across all potential classification criteria. When predicting numerical target variables, metrics such as mean squared error (MSE), root mean squared error (RMSE), or mean absolute error (MAE) can be employed to assess the model's predictive capability. This can aid in comprehending the outcomes by tracing the paths from the root node to each leaf node, which describes decision rules and distribution of outcomes across leaf nodes. You can evaluate feature importance by checking how often a feature is used to split a node and how much impurity or variance is reduced due to a feature. Graphviz, for example, can be a straightforward way to visualize a tree, as can the plot tree function from the scikit-learn. The resulting decision tree offers visual representation of the data, highlighting the factors that contribute to different outcomes. Train data until the decision tree is retrieving better results Evaluation & Interpretation: This stage confirms if the decision tree is accurate and usable, so that it can provide useful insights about business decisions.

**Applications in Business:** In marketing, they are often applied for purposes like customer segmentation, target audience identification, and forecasting customer turnover. In the field of finance, they can be used for credit risk assessment, fraud detection, and portfolio management. In business operations, they can be employed for streamlining the supply chain, tracking inventory levels, and maintaining quality control. In HR, they can apply to employee performance evaluations, hiring, and training. They are also used in decision support systems where the algorithm recommends a best decision for a complicated decision-making scenario involving multiple attributes. In health care, they are used for diagnosis for disease diagnosis, treatment, and assessment of patient risk. Decision trees are interpretable which makes them very useful especially when you need to understand how decisions are made.

---

## 5.4 SELF ASSESSMENT QUESTION

---

Decision  
Theory

### 5.4.1 Multiple-Choice Questions (MCQs)

**1. What is Decision Theory primarily concerned with?**

- a. Probability calculations
- b. Making optimal choices under uncertainty
- c. Financial accounting
- d. Manufacturing processes

**2. Which of the following is NOT a type of decision-making environment?**

- a. Decision-making under certainty
- b. Decision-making under uncertainty
- c. Decision-making under dictatorship
- d. Decision-making under risk

**3. Which decision-making condition involves complete knowledge of outcomes?**

- a. Uncertainty
- b. Risk
- c. Certainty
- d. Probability-based decision-making

**4. A decision tree is mainly used for:**

- a. Predicting financial losses
- b. Evaluating decision alternatives systematically
- c. Conducting experiments
- d. Measuring economic growth

**5. Which component is NOT part of a decision tree?**

- a. Decision nodes
- b. Probability nodes
- c. Regression equations
- d. Outcome nodes

**6. Which of the following represents a decision-making technique that evaluates multiple possible outcomes?**

- a. Decision tree
- b. Pie chart
- c. Histogram
- d. Time series analysis

**7. What does "Maximin" strategy imply in decision-making?**

- a. Choosing the alternative with the best worst-case scenario
- b. Maximizing profits at any cost
- c. Ignoring uncertainties
- d. Selecting random alternatives

**8. In decision-making under risk, probabilities of outcomes are:**

- a. Unknown
- b. Known
- c. Assumed to be equal
- d. Ignored

**9. What is the purpose of Expected Monetary Value (EMV) in decision-making?**

- a. To determine the worst possible outcome
- b. To calculate the most likely profit or loss
- c. To eliminate uncertainty
- d. To ignore risks

**10. Which of the following is NOT a component of decision theory?**

- a. Alternatives
- b. Outcomes
- c. psychological factors
- d. Payoffs

**11. What is a key advantage of using decision trees?**

- a. They eliminate risk
- b. They provide a structured and visual representation of choices
- c. They guarantee maximum profit
- d. They are only useful for large businesses

**12. Bayesian decision theory is based on:**

- a. Subjective opinions
- b. Probability and statistics
- c. Random selection
- d. Maximizing losses

**13. The Hurwicz criterion is used when decision-makers:**

- a. Are highly risk-averse
- b. Are optimistic or pessimistic about outcomes
- c. Have complete certainty
- d. Use decision trees only

**14. Which tool is commonly used for decision-making under uncertainty?**

- a. Probability distributions
- b. Regression analysis
- c. SWOT analysis
- d. Demand forecasting



**15. Which of the following best describes a "Payoff Matrix"?**

- a. A mathematical tool showing possible outcomes for each decision alternative
- b. A graphical representation of financial trends
- c. A type of accounting statement
- d. A time-series model

**5.4.2 Short Questions**

- 1. What is decision theory?
- 2. Explain decision-making under certainty.
- 3. What are decision trees in statistics?
- 4. How does decision theory impact business decisions?
- 5. What are the advantages of decision trees?

**5.4.3 Long Questions:**

- 1. Explain the process of decision-making in uncertainty.
- 2. Discuss the importance of decision trees in business strategy.

# **MATS UNIVERSITY**

**MATS CENTRE FOR DISTANCE AND ONLINE EDUCATION**

**UNIVERSITY CAMPUS:** Aarang Kharora Highway, Aarang, Raipur, CG, 493 441

**RAIPUR CAMPUS:** MATS Tower, Pandri, Raipur, CG, 492 002

**T : 0771 4078994, 95, 96, 98 Toll Free ODL MODE : 81520 79999, 81520 29999**

**Website:** [www.matsodl.com](http://www.matsodl.com)

