



**ATLAS**  
SKILLTECH  
UNIVERSITY

Accredited with

**NAAAC**



Recognized by the  
University Grants Commission (UGC)  
under Section 2(f) of the UGC Act, 1956

COURSE NAME

**STATISTICS FOR BUSINESS**

COURSE CODE

**OL BBA BA 110**

**CREDITS: 3**



**ATLAS**  
SKILLTECH  
UNIVERSITY

Centre for Distance  
& Online Education



[www.atlasonline.edu.in](http://www.atlasonline.edu.in)





Accredited with

**NAAC**



Recognized by the  
University Grants Commission (UGC)  
under Section 2(f) of the UGC Act, 1956

COURSE NAME:

**STATISTICS FOR BUSINESS**

COURSE CODE:

**OL BBA BA 110**

**Credits: 3**



**Centre for Distance  
& Online Education**



[www.atlasonline.edu.in](http://www.atlasonline.edu.in)



## Content Review Committee

Members	Members
<b>Dr. Deepak Gupta</b> Director ATLAS Centre for Distance & Online Education (CDOE)	<b>Dr. Naresh Kaushik</b> Assistant Professor ATLAS Centre for Distance & Online Education (CDOE)
<b>Dr. Poonam Singh</b> Professor Member Secretary (Content Review Committee) ATLAS Centre for Distance & Online Education (CDOE)	<b>Dr. Pooja Grover</b> Associate Professor ATLAS Centre for Distance & Online Education (CDOE)
<b>Dr. Anand Kopare</b> Director: Centre for Internal Quality (CIQA) ATLAS Centre for Distance & Online Education (CDOE)	<b>Prof. Bineet Desai</b> Prof. of Practice ATLAS SkillTech University
<b>Dr. Shashikant Patil</b> Deputy Director (e-Learning and Technical) ATLAS Centre for Distance & Online Education (CDOE)	<b>Dr. Mandar Bhanushe</b> External Expert (University of Mumbai, ODL)
<b>Dr. Jyoti Mehndiratta Kappal</b> Program Coordinator: MBA ATLAS Centre for Distance & Online Education (CDOE)	<b>Dr. Kaial Chheda</b> Associate Professor ATLAS SkillTech University
<b>Dr. Vinod Nair</b> Program Coordinator: BBA ATLAS Centre for Distance & Online Education (CDOE)	<b>Dr. Simarieet Makkar</b> Associate Professor ATLAS SkillTech University

### Program Coordinator BBA:

**Dr. Vinod Nair**

Asst. Professor  
ATLAS Centre for Distance & Online Education (CDOE)

### Unit Preparation:

**Unit 1 –4****Dr. Mukul Bhatt**

Associate Professor  
ATLAS SkillTech University

**Unit 5 –9****Dr. Neha Karnik**

Associate Professor  
ATLAS SkillTech University

### Secretarial Assistance and Composed By:

Mr. Sarur Gaikwad / Mr. Prashant Nair / Mr. Dipesh More



## Detailed Syllabus

Block No.	Block Name	Unit No.	Unit Name
1	Data Organisation and Frequency Analysis	1	Introduction to Data
		2	Frequency Distributions
2	Probability Theory and Random Variables	3	Probability
		4	Random Variables
3	Statistical Measures and Distributions	5	Measures of Central Tendency
		6	Measures of Dispersion
		7	Probability Distributions
4	Correlation and Regression Analysis	8	Correlation
		9	Regression

Course Name: Statistics for Business

Course Code: OL BBA BA 110

Credits: 3

Teaching Scheme				Evaluation Scheme (100 Marks)	
Classroom (Online)	Session	Practical / Group Work	Tutorials	Internal Assessment (IA)	Term End Examination
9+1 = 10 Sessions		-	-	30% (30 Marks)	70% (70 Marks)
Assessment Pattern:	Internal			Term End Examination	
	Assessment I	Assessment II			
Marks	15	15		70	
Type	MCQ	MCQ		MCQ – 49 Marks, Descriptive questions – 21 Marks (7 Marks * 3 Questions)	

**Course Description:**

This course provides an introduction to the fundamental concepts and methods of Statistics for Business, covering data organization, frequency distributions, probability, random variables, measures of central tendency and dispersion, correlation, regression, and various probability distributions. Students will learn to collect, summarize, analyze, and interpret business data to support decision-making processes.

### Course Objectives:

1. To understand the basic concepts of data organization, types of data, and scales of measurement essential for business applications.
2. To construct and interpret various types of frequency distributions, including relative, percentage, and cumulative distributions, and their graphical representations.
3. To comprehend the principles and theorems of probability, including conditional probability and Bayes' Theorem, as a foundation for inferential statistics.
4. To apply the concepts of random variables, probability mass function (PMF), and cumulative distribution function (CDF).
5. To calculate and analyze measures of central tendency (mean, median, mode) and measures of dispersion (range, quartile, standard deviation) for data summarization.
6. To evaluate the relationship between two variables using correlation and regression techniques.

### Course Outcomes:

At the end of course, the students will be able to:

- CO1: Remember: Recall the key terms, definitions, and formulas related to different types of data, probability concepts, and basic statistical measures.
- CO2: Understand: Explain the significance of organizing data, constructing frequency distributions, and the theoretical underpinnings of probability and random variables in a business context.
- CO3: Apply: Compute and interpret the measures of central tendency (mean, median, mode) and dispersion (standard deviation, variance) for a given set of business data.
- CO4: Analyze: Differentiate between various concepts like types of correlation and regression and analyze the properties of different probability distributions (Binomial, Poisson, Normal).
- CO5: Evaluate: Judge the strength and direction of linear relationship between variables using Karl Pearson's and Spearman's coefficients of correlation and interpret the lines of regression.
- CO6: Create: Formulate business problems using statistical frameworks and synthesize the knowledge of probability and distributions to model real-world business scenarios.

Pedagogy: Online Class, Discussion Forum, Case Studies, Quiz etc

Textbook: Self Learning Material (SLM) From Atlas SkillTech University

Reference Book:

1. Newbold, P., Carlson, W. L., & Thorne, B. (2020). *Statistics for business and economics* (9th ed.). Pearson.
2. Keller, G. (2020). *Statistics for management and economics* (11th ed.). Cengage Learning.
3. Anderson, D. R., Sweeney, D. J., Williams, T. A., Camm, J. D., & Cochran, J. J. (2021). *Statistics for business and economics* (14th ed.). Cengage Learning.

Course Details:

Unit No.	Unit Description
1	Introduction to Data: Organisation of Data, Types of Data, Scales of Measurement, Statistical Series.
2	Frequency Distributions: Construction of Frequency Distributions, Relative and Percentage Frequency Distribution, Cumulative Frequency Distribution, Frequency Density, Bivariate Frequency Distribution, Graphical Representations of Frequency Distributions, Summary, Key Terms, Descriptive Questions
3	Probability: Introduction to Probability, Important Terms and Concepts, Definitions of Probability, Theorems on Probability, Conditional Probability, Multiplicative Theorem for Independent Events, Bayes' Theorem.
4	Random Variables: Introduction, Random Variables, Probability Mass Function (PMF), Cumulative Distribution Function (CDF), Two-Dimensional Discrete Random Variables.
5	Measures of Central Tendency: Measures of Average, Arithmetic Mean, Positional Averages: Median and Mode, Empirical Analysis of Central Tendency.
6	Measures of Dispersion: Objectives of Measuring Dispersion, Range, Quartile and Percentile Measures, Mean Deviation, Standard Deviation and Variance.
7	Probability Distributions: Introduction, Binomial Distribution, Poisson Distribution, Normal Distribution.
8	Correlation: Introduction, Types of Correlation, Scatter Diagram and Simple Graph, Karl Pearson's Coefficient of Correlation, Properties of Coefficient of Correlation, Spearman's Rank Correlation.

9	Regression: Introduction, Types of Regression, Methods of Studying Regression, Lines of Regression, Regression Coefficients, Properties of Lines of Regression (Linear Regression).
---	---

#### POCO Mapping

CO	PO 1	PO 2	PO 3	PO 4	PSO 1	PSO 2	PSO 3	PSO 4	PSO 5	PSO 6	PSO 7	PSO 8
CO 1	2	1	1	3	3	1	2	3	1	1	1	1
CO 2	2	2	2	3	3	1	2	3	1	1	1	1
CO 3	2	2	2	3	3	1	2	3	1	1	2	1
CO 4	2	2	2	3	3	1	2	3	1	1	2	1
CO 5	2	2	2	3	3	1	2	3	1	1	2	1
CO 6	2	3	3	3	3	1	2	3	1	1	2	1

## Unit 1: Introduction to Data

### Learning Objectives

1. Understand the importance and methods of data organisation, including the need for classifying, tabulating, and presenting data in a structured format.
2. Differentiate between the major types of data, such as qualitative vs. quantitative data, and primary vs. secondary data.
3. Identify and apply the four levels of measurement scales—nominal, ordinal, interval, and ratio—used in statistical analysis.
4. Construct and interpret various statistical series, including individual, discrete, and continuous series, for effective data representation.
5. Summarise data meaningfully using appropriate tools and techniques that highlight the central tendencies, dispersion, and distribution patterns.
6. Familiarise with key statistical terminology to enhance comprehension and communication of statistical findings.
7. Apply knowledge through descriptive questions and a real-life case study, reinforcing practical understanding and analytical skills.

### Content

- 1.0 Introductory Caselet
- 1.1 Organisation of Data
- 1.2 Types of Data
- 1.3 Scales of Measurement
- 1.4 Statistical Series
- 1.5 Summary
- 1.6 Key Terms
- 1.7 Descriptive Questions
- 1.8 References
- 1.9 Case Study

## 1.0 Introductory Caselet

### “Anita’s Academic Analytics: Turning Raw Scores into Smart Insights”

Anita, a data coordinator at a senior secondary school in Pune, was overwhelmed by the amount of academic data she had to process each term. From exam scores and attendance records to student feedback and extracurricular achievements, every department sent her spreadsheets filled with unorganised information.

The principal had tasked her with identifying academic trends and suggesting improvements across grades. However, Anita found it hard to make sense of the scattered raw data. There was no standard format, no consistent naming, and no clear method of presenting information. Even though the school had years of valuable data, it wasn’t being used effectively.

During a professional development seminar, Anita learned the basics of **statistical data organisation and analysis**. She discovered that not all data is the same—it can be **qualitative or quantitative**, and it must be measured on appropriate **scales**, such as **nominal, ordinal, interval, or ratio**. More importantly, she understood how to convert raw marks into **individual, discrete, and continuous statistical series**.

Returning to school, Anita first focused on the recent mid-term exam data. She **classified** it by grade, subject, and gender, and then **tabulated** the marks into frequency tables. She created **continuous series** to identify score ranges where most students fell and used **bar charts** to highlight subject-wise performance.

When she presented her findings at a staff meeting, the response was positive. Teachers could now see clear trends—for instance, most students scored between 60–70 in science but between 80–90 in language subjects. This helped them plan targeted revision sessions. Attendance records, previously just lists of dates, were turned into meaningful patterns using **time-series data**.

Anita’s structured approach transformed how academic data was viewed in the school. What was once just a pile of marks became a foundation for insight-driven teaching and decision-making.

#### **Critical Thinking Question:**

If you were in Anita’s position, what type of data would you organise first to make the biggest impact on academic planning: exam scores, attendance records, or extracurricular performance? Justify your answer based on the concepts of data types and statistical series.

## 1.1 Organisation of Data

When we collect data for any purpose—such as a survey, experiment, or research study—it usually comes in a raw or unprocessed form. This raw data is often scattered, disorganized, and difficult to interpret. Therefore, to make it useful, we must **organise** it in a meaningful and systematic way.

**Organisation of data** means arranging it so that patterns, relationships, or insights become visible. This involves **classifying, tabulating, and presenting** the data in forms that are easy to read and analyze, such as tables, graphs, and charts.

This process is the **foundation of statistical analysis**, which helps in drawing conclusions, making decisions, and communicating findings effectively.

### 1.1.1 Meaning and Importance of Data Organisation

#### Meaning:

Data organisation refers to a set of methods used to **transform raw data into structured formats**. When we receive data from surveys, observations, or experiments, it's often not ready to use. We need to sort it, group it, and represent it in ways that make sense.

For instance, if a teacher collects marks from 200 students across five subjects, the raw list will be lengthy and hard to read. But if the marks are organised by student, subject, and class average, patterns start to emerge.

#### Importance of Organising Data:

1. **Simplifies Complex Data:**

Raw data is often too large and unstructured. Organising it helps in making it readable and manageable.

2. **Facilitates Analysis:**

Well-organised data is easier to compare, contrast, and analyze. You can quickly identify highs, lows, and averages.

3. **Reduces Errors:**

Systematic organisation minimizes duplication, confusion, and misinterpretation of data.

4. **Saves Time:**

Once data is structured, it becomes faster to interpret, which speeds up decision-making.

5. **Helps in Presentation:**

Whether for academic, business, or policy purposes, organised data can be presented to others in a professional manner.

## 6. Supports Informed Decision-Making:

In business, government, and research, decisions are often based on data. Good organisation ensures those decisions are well-informed.

### 1.1.2 Methods of Data Organisation

There are several methods used to organise data, each depending on the nature of the data and the purpose of analysis. The main methods include:

#### 1. Classification:

This is the process of sorting data into categories based on common characteristics. For example, grouping students by grade level or products by category.

#### 2. Tabulation:

This refers to arranging data in a table format with rows and columns. It makes comparison and summarisation easier.

#### 3. Data Presentation (Graphs and Charts):

Visual representation helps in quickly grasping the meaning of data. Graphs, diagrams, and charts make it easier to interpret large sets of information.

Each of these methods plays a role in turning raw data into meaningful insights.

### 1.1.3 Classification of Data

**Classification** means arranging data into **groups or categories** that share similar characteristics. This helps to identify relationships and trends. There are four common ways to classify data:

#### 1. Chronological Classification:

Data is classified **according to time**—hours, days, months, years, decades.

*Example:* Tracking the number of cars sold in a company from 2015 to 2025.

This type of classification helps show **trends over time**.

#### 2. Geographical Classification:

Data is arranged based on **location**—country, state, city, or region.

*Example:* Comparing literacy rates across different states of India.

This is useful for understanding **regional differences** and planning area-specific policies.

#### 3. Qualitative Classification:

Data is grouped based on **non-numeric attributes** such as gender, religion, or occupation.

*Example:* Classifying employees as male or female, or by job type (teacher, engineer, doctor).

This is commonly used in **demographic and social studies**.

#### **4. Quantitative Classification:**

Data is grouped based on **numerical values**.

*Example:* Income groups (below ₹10,000, ₹10,001–₹20,000, above ₹20,000).

This method is useful when dealing with **measurable variables** like age, income, weight, etc.

#### **Why Classification is Important:**

- It simplifies complex data.
- Enables targeted analysis (e.g., by age group, region, income level).
- Helps compare similar groups or identify outliers.

#### **1.1.4 Tabulation of Data**

**Tabulation** is the process of presenting data in a structured table format. This makes the data more organized, systematic, and easy to understand.

#### **Structure of a Table**

A properly constructed table usually includes:

- **Table number** (for reference, e.g., Table 1.1)
- **Title** (describes what the table shows)
- **Row and Column Headings** (clearly labeled categories)
- **Body** (the main data entries)
- **Footnotes** (to explain special entries or units)
- **Source** (if data is taken from a publication or survey)

#### **Types of Tables**

1. **Simple Table** – Shows data related to one characteristic (e.g., population by year).
2. **Complex or Multi-way Table** – Involves two or more characteristics (e.g., population by gender and age group).

### Benefits of Tabulation

- Provides a clear and concise summary of data.
- Helps in quick comparison between categories.
- Saves time and space while handling large data sets.
- Enables easy extraction of key figures.

### Example Table

**Table 1.1: Population of a Country by Year and Gender (in millions)**

Year	Male Population	Female Population	Total Population
2010	620	590	1,210
2015	670	635	1,305
2020	700	650	1,350
2025	740	690	1,430

### “Activity : Construction of a Simple Table”

Collect any set of information from your daily life, such as your class attendance for a week, monthly household expenses, marks obtained in different subjects, or the number of books read in different months. Organize this data into a properly structured simple table that includes a table number, title, row and column headings, body (data), and source (if applicable). For example, a student’s monthly expenses can be shown as follows:

*Table 1.2: Monthly Expenses of a Student (in INR)*

Item	Amount (₹)
Food	3,500
Transport	1,200
Stationery	800
Internet/Phone	1,000
Miscellaneous	1,500
<b>Total</b>	<b>8,000</b>

### 1.1.5 Presentation of Data: Tables, Charts, and Diagrams

Once data is tabulated, it can be presented visually to improve understanding. Visual presentation is particularly helpful when the audience is not familiar with statistical analysis, as it allows trends, proportions, and comparisons to be observed quickly and clearly.

#### 1. Tables

- We’ve already discussed tabulation.
- Tables are the foundation for all further presentation and provide the raw data that can later be represented visually.

#### Example Table

Table 1.3: Quarterly Sales of a Company (in ₹ Lakhs)

Quarter	Product A	Product B	Product C	Total Sales
Q1	120	90	60	270
Q2	140	110	80	330
Q3	160	130	100	390
Q4	180	150	120	450

Source: Hypothetical Data

#### 2. Charts

Charts are graphic representations that show relationships, proportions, or trends in data.

- **Bar Chart:** Uses bars to show the frequency or value of different items. Best for comparing categories.

Example (Sales of Products in Q4)



- **Pie Chart:** A circular chart divided into slices, each representing a proportion of the whole. Useful for showing percentage distribution.

Example (Share of Q4 Sales)

Product A: 40%

Product B: 33%

Product C: 27%

(Visual: Imagine a circle divided into 3 slices with above proportions)

- **Line Chart:** Shows data points connected by a line. Best for showing trends over time.

*Example (Total Sales Trend across Quarters)*

Q1 ● — Q2 ● — Q3 ● — Q4 ●  
270 330 390 450

### 3. Diagrams

Used for more specific or complex types of data.

- **Histogram:** A type of bar chart used for continuous data, showing frequency distribution.
- **Frequency Polygon:** A line graph that shows frequencies by connecting midpoints of histogram bars.
- **Pictogram or Pictograph:** Uses images or icons to represent data values, making data easy to understand for children or the general public.

**Example Pictogram (Books Read by Students in a Month)**

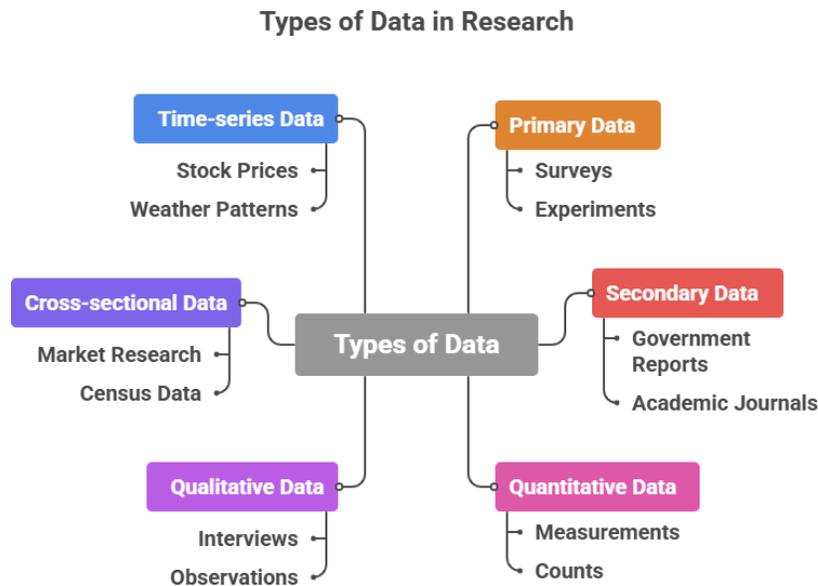
☐ = 2 books

- Student A: ☐☐☐☐ (8 books)
- Student B: ☐☐☐ (6 books)
- Student C: ☐☐☐☐☐☐ (10 books)

### Advantages of Visual Presentation

- Makes large amounts of data easier to understand.
- Highlights trends and comparisons quickly.
- More engaging and accessible for presentations and reports.

## 1.2 Types of Data



**Fig.1.1. Types of Data**

In statistics, **data** refers to the collection of facts, figures, or information that can be used for analysis. Understanding different **types of data** is essential because the methods of collection, organisation, and analysis depend on the nature of the data.

Data can be classified in different ways based on **how it is collected**, **what it describes**, and **when it was observed**.

### 1.2.1 Primary Data

**Definition:**

Primary data is the data that is **collected first-hand by the researcher** for a specific purpose or study. It is original and collected directly from the source.

**Examples:**

- A researcher conducting a survey to find out students’ study habits.
- A scientist performing an experiment and recording observations.
- A company collecting customer feedback through questionnaires.

### **Methods of Collecting Primary Data:**

- **Surveys and Questionnaires**
- **Interviews (face-to-face or telephonic)**
- **Experiments**
- **Observations**
- **Focus Groups**

### **Advantages:**

- Highly **relevant** to the research purpose.
- **Up-to-date** and **accurate** when collected properly.
- Researcher has **full control** over data collection.

### **Disadvantages:**

- Time-consuming and expensive.
- Requires planning and careful design.

## **1.2.2 Secondary Data**

### **Definition:**

Secondary data is the data that has **already been collected and published** by someone else for a different purpose. It is reused for a new analysis or study.

### **Examples:**

- Census data published by the government.
- Reports from the World Bank or WHO.
- Information from books, journals, newspapers, and websites.
- Company annual reports.

### **Sources of Secondary Data:**

- **Government publications**
- **International organisations (UN, IMF, WHO)**
- **Research papers and theses**
- **Online databases and statistics portals**
- **Business and market reports**

**Advantages:**

- Easily available and **less expensive**.
- **Timesaving** as it is already collected.
- Useful for conducting **preliminary research**.

**Disadvantages:**

- May be **outdated or irrelevant** to current study.
- Data accuracy and reliability may be **uncertain**.
- The researcher has **no control** over how it was collected.

**1.2.3 Quantitative and Qualitative Data**

Data can also be classified based on its **nature** or **characteristics**:

**1. Quantitative Data:**

This is **numerical data**, meaning it can be measured and expressed in numbers. It deals with **quantities** and can be used for mathematical calculations and statistical analysis.

**Types of Quantitative Data:**

- **Discrete data:** Whole numbers (e.g., number of children in a family).
- **Continuous data:** Any value within a range (e.g., height, weight, temperature).

**Examples:**

- Age of people in years
- Monthly income of individuals
- Number of books in a library

**2. Qualitative Data:**

This is **descriptive data**, representing characteristics, qualities, or categories. It cannot be measured numerically but can be classified or labeled.

**Examples:**

- Gender (male, female, other)
- Nationality (Indian, Chinese, Brazilian)
- Eye color (blue, brown, green)
- Type of job (teacher, engineer, doctor)

**Key Differences:**

Feature	Quantitative Data	Qualitative Data
Nature	Numerical	Descriptive
Measurement	Measurable	Not measurable
Examples	Height, Age, Salary	Gender, Religion, Occupation
Analysis	Statistical methods	Classification, frequency

### 1.2.4 Cross-sectional vs. Time-series Data

This classification is based on time—whether data is collected at **one point in time** or **over a period of time**.

#### 1. Cross-sectional Data

- Data collected at a single point in time from multiple sources or individuals.
- Shows a "snapshot" of a situation or phenomenon.
- Used for comparison across different groups or categories.

##### Examples:

- Income levels of 100 households in Delhi in the year 2025.
- Literacy rate across different Indian states in 2021.
- Number of students enrolled in different colleges in a given year.

#### 2. Time-series Data

- Data collected over a period of time, usually at regular intervals (daily, monthly, yearly).
- Helps identify trends, patterns, and changes over time.

##### Examples:

- Monthly unemployment rate in India from 2010 to 2020.
- Annual rainfall in Mumbai from 2000 to 2024.
- GDP growth of a country over the last 20 years.

### Comparison Table

Feature	Cross-sectional Data	Time-series Data
Time Frame	One point in time	Multiple time periods
Purpose	Compare between groups	Observe trends and patterns

Example	Test scores of students in 2025	Test scores from 2015–2025
---------	---------------------------------	----------------------------

### Illustrative Example

#### Cross-sectional Data Example (Students' Marks in 2025)

Table 1.4: Marks of Students in Different Subjects (2025)

Student	Math	Science	English
A	85	78	92
B	70	82	80
C	90	88	85
D	65	74	70

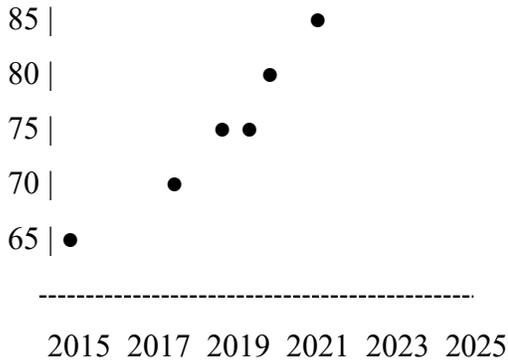
- This table shows the marks of four students **at one point in time (2025)**.
- It is a **cross-sectional dataset** because the data is not spread over years.

#### Time-series Data Example (Average Math Marks Over Years)

Table 1.5: Average Math Marks of Students (2015–2025)

Year	Average Marks
2015	65
2016	67
2017	70
2018	72
2019	75
2020	74
2021	78
2022	80
2023	82
2024	84
2025	85

**Graph: Time-series Data (2015–2025)**

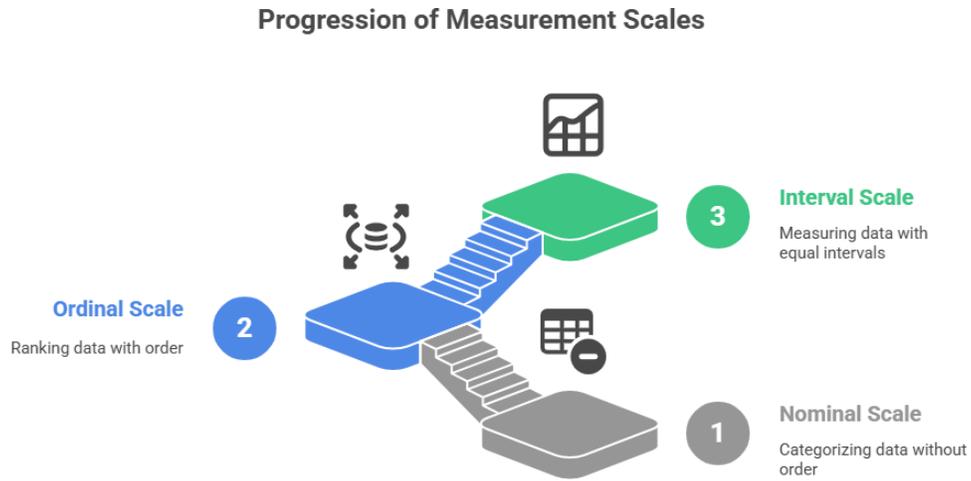


- This **line graph** shows how the average math marks changed over **11 years**.
- It is a **time-series dataset** since it tracks values over a continuous period.

### Did You Know?

“Did you know that **cross-sectional data** and **time-series data** serve very different purposes in statistics, even though they’re often mixed in business reports? Cross-sectional data gives you a **snapshot at a specific point in time**, comparing multiple subjects (like different cities, companies, or people). In contrast, time-series data captures how a **single subject changes over time**, creating a moving picture or trend line. Confusing the two can lead to **misinterpretation of patterns**, especially when making business or policy decisions. Being able to distinguish between the two helps analysts frame questions and solutions more accurately.”

## 1.3 Scales of Measurement



**Fig.1.2. Scales of Measurement**

In statistics, when we collect and organise data, it's important to understand **how that data is measured**. The **scale of measurement** tells us what kind of data we're working with and what mathematical or statistical operations can be performed on it.

There are **four main types of measurement scales**:

1. Nominal Scale
2. Ordinal Scale
3. Interval Scale
4. Ratio Scale

Each scale has different properties and levels of precision.

### 1.3.1 Nominal Scale

**Definition:**

The **nominal scale** is the simplest level of measurement. It is used for **categorising data into distinct groups that do not have any order or ranking**.

**Characteristics:**

- Data is **qualitative**.

- Categories are **mutually exclusive** (a value fits in one group only).
- **No mathematical operations** can be performed, except counting (frequency).

**Examples:**

- Gender: Male, Female, Other
- Nationality: Indian, American, Japanese
- Type of Vehicle: Car, Bike, Bus, Truck

**Uses:**

- Classifying responses in surveys
- Categorising populations in demographic studies

### 1.3.2 Ordinal Scale

**Definition:**

The **ordinal scale** is used when the data can be **classified and ranked** in a certain order, but the **differences between the ranks are not known or not equal**.

**Characteristics:**

- Data shows **order or ranking**, but **not the exact difference** between the ranks.
- Still **qualitative** or **semi-quantitative**.
- Mathematical operations are **limited to ranking or ordering**.

**Examples:**

- Customer satisfaction: Very satisfied, Satisfied, Neutral, Dissatisfied
- Education levels: High School, Bachelor's, Master's, Ph.D.
- Military ranks: Lieutenant, Captain, Major

**Uses:**

- Rating scales in questionnaires
- Socio-economic classification
- Measuring preference levels

### 1.3.3 Interval Scale

**Definition:**

The **interval scale** is a **quantitative** scale where the data has **ordered categories**, and the **difference between values is meaningful and equal**, but there is **no true zero point**.

**Characteristics:**

- Numerical data with **equal intervals** between points
- **No absolute zero**, so ratios are meaningless
- Addition and subtraction are possible; multiplication and division are not

**Examples:**

- Temperature in Celsius or Fahrenheit ( $0^{\circ}\text{C}$  does not mean 'no temperature')
- Dates in a calendar (difference between years is meaningful)
- IQ scores

**Uses:**

- Climate studies
- Psychological testing
- Education and academic research

### Did You Know?

“Did you know that the **interval scale**, unlike the **ratio scale**, does not have a true zero point? This means that while you can measure the **difference between values**, you **cannot make ratio-based comparisons**. For example, a temperature of  $20^{\circ}\text{C}$  is **not** twice as hot as  $10^{\circ}\text{C}$  because  **$0^{\circ}\text{C}$  is not an absence of temperature**, but rather an arbitrary point on the Celsius scale. This unique property makes interval scales perfect for analyzing data such as **IQ scores, temperature, and calendar years**, where meaningful intervals exist, but the concept of "twice as much" does not apply. Understanding this helps prevent common analytical errors when interpreting non-ratio data.”

### 1.3.4 Ratio Scale

**Definition:**

The **ratio scale** is the **highest and most precise** level of measurement. It has all the properties of an interval scale, **plus a meaningful and absolute zero point**. This allows for **all mathematical operations**, including ratios.

**Characteristics:**

- Numerical and **quantitative** data

- Equal intervals and **true zero** (zero means absence of the quantity)
- Allows comparison using all arithmetic operations, including percentages and ratios

**Examples:**

- Height, Weight, Age
- Income, Profit, Sales
- Distance, Speed, Time

**Uses:**

- Scientific and engineering measurements
- Financial analysis
- Health and medical studies

**Comparison Table of Measurement Scales:**

Feature	Nominal	Ordinal	Interval	Ratio
Type of Data	Qualitative	Qualitative	Quantitative	Quantitative
Order Present	No	Yes	Yes	Yes
Equal Intervals	No	No	Yes	Yes
True Zero	No	No	No	Yes
Example	Blood type	Class rank	Temperature (°C)	Income in ₹
Math Operations	Counting only	Ranking only	+ and -	+, -, ×, ÷

**1.3.5 Applications of Measurement Scales in Statistics**

Understanding measurement scales is essential because they determine:

- What **statistical tools and techniques** can be used
- What kind of **graphs or diagrams** are appropriate
- What **conclusions** can be drawn from data

**Applications:**

**1. Choosing statistical methods:**

- Nominal data → Mode, frequency, chi-square test
- Ordinal data → Median, rank correlation
- Interval/Ratio data → Mean, standard deviation, correlation, regression

## 2. Designing questionnaires and surveys:

- Ensures appropriate scale is used for measuring opinions, preferences, or behaviour.

## 3. Data analysis and interpretation:

- Ratio data allows for complex analysis like growth rates and elasticity.
- Ordinal data is used in ranking customer feedback.

## 4. Visual representation:

- Nominal/ordinal: Bar charts, pie charts
- Interval/ratio: Histograms, line graphs, scatter plots

## 5. Policy and decision-making:

- Governments and businesses use ratio and interval data for budgeting, forecasting, and evaluation.

## 1.4 Statistical Series

A **statistical series** is an **arrangement of data** in a **systematic order**, usually to show how values of a variable occur. It helps to summarise and present data in a compact and meaningful way for analysis and comparison.

Statistical series are formed **after classification** of data and can be grouped in different formats based on how data is recorded and how frequently each value appears.

There are **three main types** of statistical series:

1. Individual Series
2. Discrete Series
3. Continuous Series

Each type serves a different purpose based on the nature of data and the level of detail needed.

### 1.4.1 Individual Series

#### Definition:

An **individual series** is a series where each item or observation is recorded separately, without grouping or frequency.

#### Characteristics:

- Each value is listed individually.
- No frequencies are mentioned.

- Common when data is limited or needs to be seen in its original form.

**Example:**

Marks obtained by five students:

45, 52, 60, 47, 50

Or presented in tabular form:

*Table 1.6: Marks Obtained by Students*

Student Name	Marks
Rohan	45
Priya	52
Aarav	60
Meena	47
Karan	50

**Use:**

An individual series is used when the dataset is small and individual details are important, such as recording marks of a few students, incomes of a few households, or daily temperatures over a short period.

**1.4.2 Discrete Series**

**Definition:**

A **discrete series** is a series in which **data is presented along with its frequency**. The variable takes specific, separate values (not in ranges), and each value shows how **often** it occurs.

**Characteristics:**

- Data is presented in the form of value–frequency pairs.
- Used for **discrete variables** (whole numbers, not fractions).
- Frequencies show **how many times** each value occurs.

**Example:**

Number of students scoring a particular mark:

Marks (x)	No. of Students (f)
40	2

45	5
50	7
55	4
60	2

Here, “50 marks” was scored by 7 students.

**Use:**

Used when the data involves counting things that occur in whole numbers like number of children in families, goals in a football match, etc.

**1.4.3 Continuous Series**

**Definition:**

A **continuous series** is a series where data is grouped into class intervals (ranges), and the frequency shows how many observations fall within each interval.

**Characteristics:**

- Data is grouped in intervals such as 0–10, 10–20, etc.
- Used for **continuous variables** (which can take any value within a range).
- Class intervals are usually of equal width but may also vary.

**Example:**

*Distribution of Heights of Students*

Height (cm)	Number of Students (f)
140–150	3
150–160	6
160–170	10
170–180	5

Here, **10 students have heights between 160 cm and 170 cm.**

**Use:**

Continuous series is used when data is large and continuous in nature, such as height, weight, age, income, or marks. It helps organize wide ranges of data into manageable class intervals for easier analysis.

### “Activity”

List **five real-life examples of continuous variables** (such as height, weight, temperature, rainfall, or income). Then, identify what class intervals could be used if the data were to be grouped in a continuous series.

*Example Response:*

- Variable: Temperature → Possible Intervals: 0–10°C, 10–20°C, 20–30°C...
- Variable: Weight → Possible Intervals: 40–50 kg, 50–60 kg, 60–70 kg...

#### 1.4.4 Construction of Statistical Series

Constructing a statistical series involves several key steps that depend on the type of data and the purpose of the analysis.

##### Steps to Construct a Statistical Series:

###### 1. Collection of Data:

- Gather raw data from primary or secondary sources.
- Example: Test scores of students, daily temperatures, etc.

###### 2. Classification of Data:

- Organise the data into suitable groups or categories.
- Decide whether the data is individual, discrete, or continuous.

###### 3. Decide the Type of Series:

- **Individual Series:** When observations are few and no grouping is needed.
- **Discrete Series:** When data values are distinct and can be counted.
- **Continuous Series:** When data ranges are better suited for analysis.

###### 4. Determine Class Intervals (for continuous series):

- Choose the number and width of intervals.
- Ensure there are **no gaps** between intervals.
- Use the **inclusive** or **exclusive method** depending on the context.

##### Inclusive Method:

- Both lower and upper limits are included (e.g., 10–19, 20–29).

**Exclusive Method:**

- The lower limit is included, but the upper limit is excluded (e.g., 10–20 means  $10 \leq x < 20$ ).

**5. Tally and Frequency:**

- Count how many values fall into each class or category.
- Use **tally marks** for easy counting, then convert to frequency numbers.

**6. Tabulate the Data:**

- Prepare a table showing class intervals or values and corresponding frequencies.

**7. Check Totals:**

- Ensure the total frequency matches the total number of observations.

**1.4.5 Construction of Statistical Series****Definition:**

The construction of a statistical series involves arranging raw data into a meaningful form so that it can be easily understood and analyzed. Depending on the nature of the data, the series may be **individual, discrete, or continuous**.

**Steps in Construction of a Statistical Series:**

1. **Determine the Purpose** – Decide why the series is being constructed and what kind of data presentation is required.
2. **Decide the Type of Series** – Choose between an individual, discrete, or continuous series.
3. **Collect and Arrange Data** – Collect raw data and arrange it systematically.
4. **Decide on Class Intervals (if needed)** – For continuous series, determine suitable class intervals.
5. **Assign Frequencies** – Count and record how many observations fall in each class interval.
6. **Prepare the Table** – Present the data in tabular form with proper headings, titles, and sources.

**Working Example:**

Suppose we collected the **marks of 20 students** in a test (out of 50):  
15, 28, 32, 40, 22, 35, 45, 18, 25, 38, 30, 42, 20, 33, 27, 29, 31, 36, 24, 41

### Step 1: Decide the Type of Series

Since marks are continuous data, we will construct a **continuous series**.

### Step 2: Decide Class Intervals

Let's take intervals of width 10: 10–20, 20–30, 30–40, 40–50.

### Step 3: Tally the Data into Intervals

Marks (Class Interval)	Tally	Frequency (f)
10–20		
20–30		
30–40		
40–50		
<b>Total</b>		<b>20</b>

### Step 4: Interpret

- Most students (8) scored between **20 and 30 marks**.
- Fewest students (2) scored between **10 and 20 marks**.
- This arrangement helps us understand the overall performance distribution.

## Knowledge Check 1

### Choose the correct option:

1. Which of the following is an example of **primary data**?
  - A) Census data published by the government
  - B) A newspaper report on inflation
  - C) Marks recorded by a teacher during a classroom test
  - D) Information downloaded from an education website
2. Which measurement scale allows the calculation of **ratios** and has a **true zero point**?
  - A) Nominal
  - B) Ordinal
  - C) Interval
  - D) Ratio

3. Data collected **at a single point in time** from different units is called:
  - A) Time-series data
  - B) Cross-sectional data
  - C) Interval data
  - D) Discrete data
4. Which of the following is **qualitative data**?
  - A) Number of books in a library
  - B) Student's height in centimetres
  - C) Type of school (government, private)
  - D) Marks in a mathematics exam
5. In which type of statistical series is data organised into **class intervals**?
  - A) Individual series
  - B) Discrete series
  - C) Continuous series
  - D) Nominal series

## 1.5 Summary

- ❖ This unit introduced the foundational concepts of data handling in statistics. It began with the **organisation of data**, explaining how raw data must be structured through classification, tabulation, and visual presentation. Different **types of data** were discussed, including primary and secondary data, as well as qualitative vs. quantitative data, and time-based classifications like cross-sectional and time-series data.
- ❖ We then explored the **scales of measurement**, which are essential for choosing correct statistical tools: nominal, ordinal, interval, and ratio scales each define the properties of the data and the operations that can be applied. Lastly, the unit covered **statistical series**, which are systematic arrangements of data in individual, discrete, or continuous forms, and the steps involved in constructing them.
- ❖ This foundational understanding prepares learners to collect, organise, classify, and present data in a meaningful way for further statistical analysis.

## 1.6 Key Terms

1. **Data** - Raw facts and figures collected for analysis.

2. **Primary Data** - Data collected directly by the researcher for a specific purpose.
3. **Secondary Data** - Data previously collected and published by someone else.
4. **Qualitative Data** - Non-numeric data based on characteristics or attributes.
5. **Quantitative Data** - Numeric data that can be measured and analyzed statistically.
6. **Cross-sectional Data** - Data collected at a single point in time.
7. **Time-series Data** - Data collected over a period of time.
8. **Nominal Scale** - A scale used for labeling categories without any order.
9. **Ordinal Scale** - A scale that shows relative ranking or order.
10. **Interval Scale** - A numeric scale with equal intervals but no true zero.
11. **Ratio Scale** - A numeric scale with equal intervals and a true zero, allowing all mathematical operations.
12. **Individual Series** - A series in which each observation is recorded individually.
13. **Discrete Series** - A series where data values are distinct and associated with frequencies.
14. **Continuous Series** - A series where data is grouped into intervals and frequencies represent ranges.

## 1.7 Descriptive Questions

1. Define data organisation. Why is it important in statistics?
2. Differentiate between **primary** and **secondary** data with examples.
3. Explain the differences between **qualitative** and **quantitative** data.
4. What is the difference between **cross-sectional** and **time-series** data?
5. Define the four **scales of measurement** and provide one example for each.
6. What are the key characteristics of an **ordinal scale**?
7. Differentiate between **discrete** and **continuous** series with suitable examples.
8. List the steps involved in the **construction of a statistical series**.
9. Present the marks of 10 students using (a) individual series, (b) discrete series.
10. What is the role of classification in the organisation of data?

## 1.8 References

1. Gupta, S.P. (2014). *Statistical Methods*. Sultan Chand & Sons.
2. Sharma, J.K. (2018). *Business Statistics*. Pearson Education.
3. Goon, A.M., Gupta, M.K., & Dasgupta, B. (2013). *Fundamentals of Statistics*. World Press.

4. Indian Statistical Institute. *Introductory Statistics Course Notes*.
5. Government of India. *National Sample Survey Office (NSSO) Reports*.
6. Websites and data portals: [www.mospi.gov.in](http://www.mospi.gov.in), [www.data.gov.in](http://www.data.gov.in)

### Answers to Knowledge Check

#### ***Knowledge Check 1***

1. C) Marks recorded by a teacher during a classroom test
2. D) Ratio
3. B) Cross-sectional data
4. C) Type of school (government, private)
5. C) Continuous series

## 1.9 Case Study: Organising School Exam Results

### Introduction

In educational institutions, data is gathered from various sources such as student assessments, attendance records, and extracurricular participation logs. However, raw data in itself lacks meaning unless it is systematically organised and interpreted. Statistics plays a crucial role in transforming this raw information into actionable insights that guide academic planning, performance evaluation, and policy decisions.

This case study focuses on how a school used statistical techniques such as **data organisation, classification, tabulation, and statistical series** to improve student performance tracking and decision-making. The challenges included managing large volumes of unorganised student data, identifying performance trends, and presenting the data in a way that allowed easy interpretation by teachers and administrators.

### Background

Green Valley Public School, a mid-sized urban school, conducts regular assessments across subjects for all students. Until recently, the academic department maintained performance records in manual registers and spreadsheets, resulting in inconsistent data entry, duplication, and missing entries. Teachers struggled to extract meaningful insights, such as identifying students needing remedial support or recognising top performers.

The school management decided to apply basic statistical techniques to handle and analyse the data in a structured way. This included the use of **data organisation methods, identification of types of data, application of scales of measurement, and construction of statistical series.**

**Problem 1: Lack of Proper Data Organisation**

**Challenge:**

The school had scattered performance records that were not systematically organised. Without proper classification and tabulation, it was difficult to compare marks across subjects or identify academic trends.

**Dataset (Sample Marks of 8 Students – out of 100):**

Student	Gender	Math	Science	English	Total
Rohan	M	65	72	60	197
Priya	F	78	81	74	233
Aarav	M	55	60	58	173
Meena	F	88	92	85	265
Karan	M	72	70	68	210
Neha	F	60	65	70	195
Aditya	M	80	85	78	243
Simran	F	74	78	82	234

**Solution:**

The administration introduced classification by **gender, subject, and performance levels.** Data was tabulated and visually represented.

- **Bar Chart Example:** Comparing average marks of boys vs. girls showed that **girls outperformed boys across all subjects.**
- **Pie Chart Example:** Subject-wise totals were shown in a pie chart to identify which subjects students excelled in.

This organisation allowed quick identification of **weak performers (Aarav, Neha)** and **top performers (Meena, Aditya).**

**Problem 2: Unclear Understanding of Data Types and Measurement Scales**

**Challenge:**

Teachers were unsure how to treat different types of student data—some numerical (marks, attendance) and others categorical (gender, remarks). Misuse of inappropriate statistical methods led to flawed interpretations.

**Dataset Extract:**

Student	Gender	Attendance (%)	Grade (A–C)	Rank
Rohan	M	82	B	5
Priya	F	95	A	2
Aarav	M	78	C	7
Meena	F	96	A	1
Karan	M	85	B	4

**Solution:**

Workshops clarified:

- **Nominal scale** → Gender (M/F)
- **Ordinal scale** → Grade (A > B > C), Rank
- **Interval scale** → Lab temperature readings
- **Ratio scale** → Marks, Attendance (has absolute zero)

This distinction prevented misuse. For example, **averages were calculated only on ratio data (marks, attendance), not on nominal data (gender).**

**Problem 3: Difficulty in Constructing and Interpreting Statistical Series**

**Challenge:**

Raw scores of hundreds of students made it difficult to analyse trends without summarisation.

**Dataset (Science Marks of 20 Students):**

56, 72, 68, 60, 75, 80, 62, 70, 78, 65, 85, 90, 73, 66, 82, 74, 69, 71, 77, 88

**Solution:**

The IT team summarised this into **statistical series**:

1. **Individual Series** – Each student’s mark was listed separately.
2. **Discrete Series** – Frequency of exact scores. Example: 2 students scored 70.
3. **Continuous Series (Grouped):**

Marks (Interval)	Frequency (f)
------------------	---------------

55–64	4
65–74	8
75–84	5
85–94	3
<b>Total</b>	20

- **Histogram:** Showed that most students scored between **65 and 74 marks**.
- **Frequency Polygon:** Helped visualise performance trends more smoothly.

This made interpretation simple for both teachers and parents.

### Conclusion

By applying the foundational concepts of data organisation and measurement, Green Valley Public School transformed scattered information into actionable insights.

- **Problem 1:** Tabulation and charts highlighted group-wise and subject-wise performance.
- **Problem 2:** Correct identification of data types and scales prevented flawed analysis.
- **Problem 3:** Statistical series and graphs simplified large datasets into meaningful patterns.

Overall, statistics helped the school **identify weak students, reward top achievers, and plan remedial actions**, thereby improving academic outcomes.

## Unit 2: Frequency Distributions

### Learning Objectives

1. Understand the concept, purpose, and construction of frequency distributions, including how raw data is grouped into frequency tables for both discrete and continuous variables.
2. Develop the ability to calculate and interpret relative and percentage frequencies, and recognise their practical applications in comparing datasets of different sizes.
3. Construct and analyse cumulative frequency distributions, including both less-than and greater-than types, and understand how to interpret ogives in real-life data scenarios.
4. Explain the concept of frequency density and apply it effectively in situations involving unequal class intervals, especially in the creation of accurate histograms.
5. Explore the construction and interpretation of bivariate frequency distributions, and understand their importance in analysing the relationship between two variables simultaneously.
6. Learn various methods of graphical representation for both qualitative and quantitative data, including bar charts, pie charts, histograms, frequency polygons, and ogives, and identify the most suitable method for a given type of data.
7. Apply statistical tools to summarise, visualise, and interpret data through structured frequency tables and graphical techniques, building a foundation for further statistical analysis and decision-making.

### Content

- 2.0 Introductory Caselet
- 2.1 Construction of Frequency Distributions
- 2.2 Relative and Percentage Frequency Distribution
- 2.3 Cumulative Frequency Distribution
- 2.4 Frequency Density
- 2.5 Bivariate Frequency Distribution
- 2.6 Graphical Representations of Frequency Distributions
- 2.7 Summary
- 2.8 Key Terms
- 2.9 Descriptive Questions

2.10 References

2.11 Case Study

## 2.0 Introductory Caselet

### “Meena’s Data Discovery: Making Patterns Visible with Frequency Distributions”

Meena, a young mathematics teacher at a secondary school in Ahmedabad, had recently taken charge of compiling class-wise academic reports. The school had just completed mid-term exams for Classes 9 and 10, and Meena was expected to present performance summaries to the academic review committee.

She received a flood of raw marks from various teachers—some listed every student’s score individually, while others had grouped them into arbitrary ranges. There were no common class intervals, and many spreadsheets lacked any labels or headings. Meena felt overwhelmed by the inconsistencies and the sheer volume of data.

Determined to bring order to the chaos, Meena referred to a statistics module she had studied during her teacher training. She realised the first step was to **construct proper frequency distributions**—grouping the data systematically to make it readable and insightful. She sorted the scores into **discrete and continuous frequency tables**, depending on whether they represented whole numbers or ranged values.

When she noticed that some teachers had used unequal class intervals, Meena applied the concept of **frequency density** to correct the imbalance before plotting histograms. She also calculated **relative frequencies** and **percentage frequencies**, allowing for fair comparison across sections with different class sizes. To help identify score thresholds and medians, she created **cumulative frequency distributions** and drew **ogives**.

Later, she noticed a separate file that included student-reported **study hours**. This gave her an idea—she built a **bivariate frequency distribution table** comparing study hours to exam scores. The results were eye-opening: students who studied more than three hours daily consistently scored above 70 marks.

At the review meeting, Meena presented her findings using pie charts, histograms, and frequency polygons. The graphs showed not only how the marks were distributed but also pointed out sections where performance was lagging. Her bivariate analysis helped initiate a school-wide discussion on study habits and time management.

What started as an unmanageable pile of raw scores turned into a clear, visual representation of student performance trends. Meena's analytical approach helped teachers see where intervention was needed and gave the school leadership a roadmap for improving academic outcomes.

**Critical Thinking Question:**

If you were in Meena's position, how would you decide whether to use absolute frequency, relative frequency, or frequency density when presenting data to the school board? Explain your choice with reference to class size, class intervals, and the need for comparison.

## 2.1 Construction of Frequency Distributions

In statistics, raw data collected through surveys, tests, or observations is usually unstructured. A **frequency distribution** helps convert this raw data into an organised form by showing how often each value or group of values occurs. This makes it easier to analyse patterns, trends, and relationships in the data.

### 2.1.1 Meaning and Purpose of Frequency Distribution

A **frequency distribution** is a table that summarises data by showing the number of times (called the frequency) each distinct value or range of values appears in a dataset.

#### Purpose:

- To simplify large volumes of data
- To make data easier to visualise through charts or graphs
- To help identify patterns such as clustering, gaps, and outliers
- To support further statistical analysis such as mean, median, mode, etc.

#### Business Example:

Suppose a shop records the **daily sales (in ₹ thousands)** over 10 days as follows:

12, 15, 18, 12, 20, 15, 18, 22, 15, 20

We can count and organise them into a **frequency table**:

Daily Sales (₹ '000)	Frequency
12	2
15	3
18	2
20	2
22	1

This table shows:

- Sales of **₹15,000** occurred on **3 days** (the most frequent).
- Sales of **₹22,000** occurred only once (least frequent).

### 2.1.2 Frequency Distribution for Discrete Data

**Discrete data** consists of countable, often whole-number values (e.g., number of children, test scores, goals in a game).

In a **discrete frequency distribution**, each unique value is listed with the number of times it occurs in the dataset.

**Example:**

Books Borrowed	Number of Students
0	2
1	5
2	8
3	4
4	1

This format shows how many students borrowed how many books in a week.

### 2.1.3 Frequency Distribution for Continuous Data

**Continuous data** can take any value within a range, including decimals or fractions (e.g., height, temperature, income).

Since exact values rarely repeat in continuous data, we use **class intervals** (ranges) to group the data.

**Example:**

Marks (Interval)	Frequency
40 – 50	3
50 – 60	7
60 – 70	10
70 – 80	5
80 – 90	2

This distribution gives a clearer view of how the values are spread across intervals.

## 2.1.4 Steps in Constructing a Frequency Table



**Fig.2.1 Steps in Constructing a Frequency Table**

Follow these steps to create a well-structured frequency distribution:

### **Step 1: Collect Raw Data**

Gather the numerical or categorical data from your source (survey, test, etc.).

### **Step 2: Identify Data Type**

- If values are whole numbers: use **discrete frequency distribution**
- If values include ranges or decimals: use **continuous frequency distribution**

**Step 3: Decide Number of Classes (for continuous data)**

Choose a suitable number of class intervals. Typically, 5 to 10 classes are used for clarity.

**Step 4: Calculate Class Width**

Use the formula:

$$\text{Class Width} = (\text{Maximum Value} - \text{Minimum Value}) \div \text{Number of Classes}$$

Example: If scores range from 40 to 90 and we want 5 classes:

$$\text{Class Width} = (90 - 40) \div 5 = 10$$

**Step 5: Create Class Intervals**

Start from the minimum value and add the class width to form intervals. Ensure:

- Intervals do not overlap
- All values are covered
- The width is consistent (unless specified otherwise)

Example: 40–50, 50–60, 60–70, and so on.

**Step 6: Tally the Data**

Go through the raw data and mark (tally) each observation in the appropriate interval.

**Step 7: Count the Frequencies**

Convert tally marks into actual frequency numbers for each class.

**Step 8: Present in Table Format**

Class Interval	Tally	Frequency
40 – 50		
50 – 60		
60 – 70		

This completed table now becomes the basis for drawing histograms, frequency polygons, and further analysis.

**2.2 Relative and Percentage Frequency Distribution**

Once raw data has been organised into a frequency table, we can further analyse it using **relative** and **percentage frequencies**. These help in comparing data sets of different sizes and understanding the **proportion** of observations in each category.

**2.2.1 Definition of Relative Frequency**

**Relative frequency** shows how often a particular value or class appears in relation to the **total number of observations**. It is expressed as a **fraction or decimal** and provides insight into the **proportion** of each class.

**Formula:**

$$\text{Relative Frequency} = \text{Frequency of a class} \div \text{Total frequency}$$

**Example:**

If 8 students scored between 60–70 out of a total of 40 students:

$$\text{Relative Frequency} = 8 \div 40 = 0.20$$

This means that 20% of the students scored in that range.

**Purpose:**

- Makes it easier to compare datasets of different sizes
- Highlights the **importance or weight** of each class
- Useful in **probability** and **data analysis**

### 2.2.2 Calculation of Relative Frequency

To calculate relative frequencies:

1. **Create a standard frequency table**
2. **Find the total of all frequencies**
3. **Divide each individual frequency by the total frequency**
4. **Round the result to 2 or 3 decimal places, if required**

**Example Table:**

Class Interval	Frequency	Relative Frequency
10 – 20	5	$5 \div 25 = 0.20$
20 – 30	10	$10 \div 25 = 0.40$
30 – 40	6	$6 \div 25 = 0.24$
40 – 50	4	$4 \div 25 = 0.16$
<b>Total</b>	<b>25</b>	

This table shows the **proportion** of observations in each interval.

### 2.2.3 Percentage Frequency Distribution

A **percentage frequency distribution** shows the same information as relative frequency, but the results are expressed as **percentages** instead of decimals.

**Formula:**

$$\text{Percentage Frequency} = (\text{Frequency of a class} \div \text{Total frequency}) \times 100$$

**Example Table:**

Class Interval	Frequency	Percentage Frequency
10 – 20	5	$(5 \div 25) \times 100 = 20\%$
20 – 30	10	$(10 \div 25) \times 100 = 40\%$
30 – 40	6	$(6 \div 25) \times 100 = 24\%$
40 – 50	4	$(4 \div 25) \times 100 = 16\%$
<b>Total</b>	<b>25</b>	<b>100%</b>

**Advantages:**

- Easier to understand at a glance
- Good for creating **pie charts** and **bar graphs**
- Helps in **comparing data groups**, especially when dataset sizes differ

**2.2.4 Applications of Relative and Percentage Frequencies**

Relative and percentage frequencies are widely used in **statistics, business, education, and research** to understand **data proportions** and to **make comparisons** across categories or groups.

**Common Applications:**

1. **Market Research:**  
To find what proportion of consumers prefer a product category (e.g., 35% prefer tea, 45% prefer coffee).
2. **Elections and Polls:**  
To analyse voter preferences as percentages.
3. **Education Reports:**  
To determine the percentage of students scoring in specific ranges (e.g., 20% scored above 90%).
4. **Sales and Inventory Analysis:**  
To identify what share of total sales comes from each product.
5. **Healthcare:**  
To report disease incidence (e.g., 10% of patients show symptoms A).

## 6. Data Visualization:

Percentage frequencies are often used in **pie charts** and **stacked bar charts** for clearer interpretation.

### 2.3 Cumulative Frequency Distribution

A **cumulative frequency distribution** shows the **progressive total of frequencies** up to a certain point or above a certain point in the dataset. It helps in understanding **how data accumulates** across the class intervals, which is particularly useful when analysing **medians, percentiles, and graphical trends**.

Cumulative frequency is of two types:

- **Less-than cumulative frequency**
- **Greater-than cumulative frequency**

#### 2.3.1 Less-than Cumulative Frequency Distribution

A **less-than cumulative frequency distribution** is obtained by successively adding the frequencies from the first class interval up to a given class. It shows how many observations are **less than the upper boundary** of each class.

#### Steps to Construct:

1. Start from the lowest class.
2. Add the frequencies step by step as you move down the table.
3. The last cumulative frequency should equal the total frequency.

#### Business Example:

A company records the **monthly sales (in ₹ '000)** of its 30 sales representatives. The data is grouped into the following class intervals:

Monthly Sales (₹ '000)	Frequency	Less-than Cumulative Frequency
0 – 10	3	3
10 – 20	5	$3 + 5 = 8$
20 – 30	7	$8 + 7 = 15$
30 – 40	9	$15 + 9 = 24$
40 – 50	6	$24 + 6 = 30$

#### Interpretation:

- **8 salespersons** recorded sales of **less than ₹20,000**.

- 15 salespersons achieved less than ₹30,000 in sales.
- All 30 salespersons had sales less than ₹50,000 (which equals the total frequency).

### “Activity: Construct and Compare Cumulative Frequencies”

#### Instruction to Student:

You are given the following class intervals and frequencies showing the marks of students in a mathematics test:

Class Interval	Frequency
0 – 10	4
10 – 20	6
20 – 30	10
30 – 40	15
40 – 50	5

1. Calculate the **less-than cumulative frequency** for each class interval.
2. Present your results in a new column titled “Less-than Cumulative Frequency.”
3. Plot a **less-than ogive** using the upper class boundaries and cumulative frequencies.
4. Comment on the shape of the curve and estimate the **median mark** by locating the midpoint on the y-axis and projecting it onto the curve.

#### 2.3.2 Greater-than Cumulative Frequency Distribution

A **greater-than cumulative frequency distribution** is obtained by starting from the highest class interval and successively adding frequencies in reverse order. It tells us how many values are **greater than or equal to the lower boundary** of each class.

#### Steps to Construct:

1. Start from the highest class interval.
2. Subtract frequencies cumulatively from the total as you move upward.
3. The first cumulative frequency equals the total frequency.

### Business Example:

A company studied the **monthly sales performance (in ₹ ‘000)** of 30 employees. The distribution is as follows:

Monthly Sales (₹ ‘000)	Frequency	Greater-than Cumulative Frequency
0 – 10	3	30
10 – 20	5	$30 - 3 = 27$
20 – 30	7	$27 - 5 = 22$
30 – 40	9	$22 - 7 = 15$
40 – 50	6	$15 - 9 = 6$

### Interpretation:

- **27 employees** recorded sales **greater than or equal to ₹10,000.**
- **22 employees** recorded sales **greater than or equal to ₹20,000.**
- Only **6 employees** achieved sales **greater than or equal to ₹40,000.**

### 2.3.3 Ogives and Their Interpretation

An **ogive** is a graphical representation of **cumulative frequency data**. It helps to visualise how values are distributed across a dataset and is widely used in business and economics to study income levels, sales figures, or production outputs.

There are two types of ogives:

#### 1. Less-than Ogive

- Plots the **upper class boundaries** on the x-axis and **less-than cumulative frequencies** on the y-axis.
- The curve rises as cumulative frequency increases.

#### 2. Greater-than Ogive

- Plots the **lower class boundaries** on the x-axis and **greater-than cumulative frequencies** on the y-axis.
- The curve descends as the number of greater-than values reduces.

### Steps to Construct an Ogive:

1. Prepare a cumulative frequency table (less-than or greater-than).

2. Use **class boundaries** on the x-axis and **cumulative frequencies** on the y-axis.
3. Plot the points and join them with a smooth curve or straight lines.

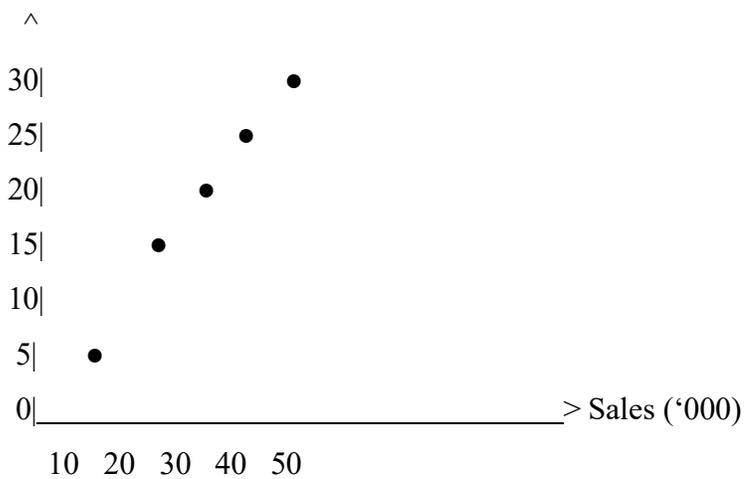
### Business Example

A company analyses the **monthly sales (in ₹ '000)** of 30 employees:

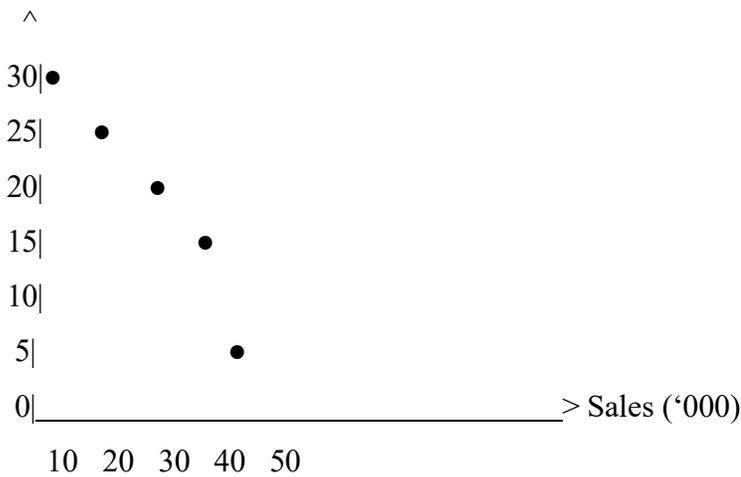
Sales (₹ '000)	Frequency	Less-than Cumulative Freq.	Greater-than Cumulative Freq.
0 – 10	3	3	30
10 – 20	5	8	27
20 – 30	7	15	22
30 – 40	9	24	15
40 – 50	6	30	6

### Graphical Representation (Ogives)

Less-than Ogive (rising curve)



Greater-than Ogive (falling curve)



(Both curves can also be drawn on the same graph for comparison. Their intersection point gives the **median**.)

### Interpreting the Ogive:

- From the **Less-than Ogive**, we see that **15 employees earned less than ₹30,000**, while all 30 earned less than ₹50,000.
- From the **Greater-than Ogive**, **22 employees earned more than or equal to ₹20,000**, while only **6 earned more than or equal to ₹40,000**.
- The **intersection point** of both ogives on the same graph indicates the **median sales level**.
- The **steeper section** between 20–40 indicates that most employees' sales are concentrated in this range.

## 2.4 Frequency Density

In statistics, especially when working with **histograms**, class intervals sometimes have **unequal widths**. In such cases, using only frequency to represent data visually may be misleading. To adjust for unequal widths and maintain visual accuracy, we use a concept called **frequency density**.

### 2.4.1 Concept of Frequency Density

**Frequency density** is a measure that allows for accurate graphical representation of data, especially when **class intervals are not of equal size**. It adjusts the frequency according to the width of the class interval so that the **area of each bar in a histogram** reflects the actual data distribution.

**Formula:**

$$\text{Frequency Density} = \text{Frequency} \div \text{Class Width}$$

Where:

- **Frequency** is the number of observations in the class.
- **Class Width** = Upper class boundary – Lower class boundary.

**Example:**

If a class has 20 students (frequency) and the interval width is 10 marks, then Frequency Density =  $20 \div 10 = 2$ .

This tells us that for each unit width, there are 2 observations on average.

**Did You Know?**

“Did you know that **frequency density** is the secret behind drawing **accurate histograms with unequal class intervals**? While many assume the height of a bar in a histogram always represents frequency, this only works when class intervals are equal. When class widths vary, using raw frequency distorts the data visually. That’s why we divide frequency by class width to get **frequency density**, ensuring the **area of the bar truly represents the data**. This keeps the histogram proportional and fair—even when class sizes aren’t.”

**2.4.2 Use of Frequency Density in Unequal Class Intervals**

When class intervals are **unequal**, directly comparing frequencies may give a misleading impression. A class with a larger interval width might naturally contain more observations, even though the data is less concentrated.

To avoid distortion, we use **Frequency Density**, which is calculated as:

**Formula:**

$$\text{Frequency Density} = \text{Frequency} \div \text{Class Width}$$

This ensures that the **height of bars** in a histogram is proportional to the density of data, not just the raw frequency.

### Business Example: Annual Income of Employees (₹ Lakhs)

A company analysed the annual income of 60 employees. The income groups were **unequal**, so frequency density was used.

Income Range (₹ Lakhs)	Frequency (No. of Employees)	Class Width	Frequency Density
0 – 5	6	5	1.2
5 – 10	12	5	2.4
10 – 20	20	10	2.0
20 – 40	15	20	0.75
40 – 60	7	20	0.35

#### Interpretation:

- Even though **20 employees** earn between ₹10–20 lakhs (highest frequency), the **densest group** is ₹5–10 lakhs (frequency density = 2.4).
- The income group **₹40–60 lakhs** has a relatively **low density (0.35)** despite covering a wide interval.
- Using frequency density prevents overestimation of higher-income ranges that have wider intervals.

### 2.4.3 Histograms Using Frequency Density

A **histogram** is a graphical representation where the **area of each bar** (not just the height) represents the frequency.

- When **class intervals are equal**, frequencies alone can be used as bar heights.
- When **class intervals are unequal**, we must use **frequency density** for the bar heights to ensure the **area of each bar** is proportional to the actual frequency.

#### Steps to Draw a Histogram with Frequency Density:

1. Calculate the class width for each interval.
2. Compute frequency density using the formula:

$$\text{Frequency Density} = \text{Frequency} \div \text{Class Width}$$

3. On the x-axis, mark the class boundaries.
4. On the y-axis, mark the frequency densities.
5. Draw bars with widths equal to the class intervals and heights equal to frequency density.

### Business Example: Weekly Sales (₹ Lakhs)

A retail company recorded weekly sales of 50 outlets. Since class intervals are unequal, frequency density is used.

Sales Range (₹ Lakhs)	Frequency (Outlets)	Class Width	Frequency Density
0 – 10	5	10	0.5
10 – 20	15	10	1.5
20 – 40	20	20	1.0
40 – 60	10	20	0.5

### Graphical Representation (Histogram Sketch)

Frequency Density (y-axis)



(Each bar's **width** = class interval and **height** = frequency density. Thus, the **area of each bar** reflects the actual frequency.)

### Why Use Frequency Density in Histograms?

- Ensures **proportional representation** when class intervals vary.
- Keeps the **area of each bar proportional** to actual frequency.
- Prevents **visual distortion** in data interpretation.

### Visual Interpretation:

- The sales range **10–20 lakhs** has the highest density (1.5), meaning sales are most concentrated in this range.
- The **20–40 lakhs** group has the highest number of outlets (20), but because the class is wide, its density (1.0) is lower.

- Without frequency density, the histogram would wrongly suggest 20–40 as the densest group.

### “Activity: Visualise Data Using Frequency Density”

#### Instruction to Student:

Below is the test score data for a group of students. The class intervals are **unequal**, so you need to use **frequency density** to draw an accurate histogram.

Class Interval	Frequency
0 – 10	5
10 – 30	10
30 – 40	6
40 – 60	9

1. Calculate the **class width** for each interval.
2. Use the formula: **Frequency Density = Frequency ÷ Class Width**
3. Prepare a table with a new column for frequency density.
4. Using graph paper or a digital tool, draw a histogram using **class intervals as bar widths** and **frequency density as bar heights**.
5. Reflect: How does using frequency density change the visual interpretation compared to using raw frequencies?

## 2.5 Bivariate Frequency Distribution

In many real-life situations, we are interested in analysing **the relationship between two variables** rather than just looking at one. For example, a school may want to know whether there is a connection between **hours of study** and **exam scores**. This is where **bivariate frequency distributions** become essential.

### 2.5.1 Concept and Importance of Bivariate Distributions

A **bivariate frequency distribution** is a table that presents the joint frequency distribution of two variables. It shows how many times combinations of values from two different variables occur together in a dataset.

Unlike **univariate distributions** (which involve only one variable), bivariate distributions allow us to:

- Explore relationships or associations between two factors
- Understand how one variable may influence another

- Prepare for further statistical tools like **correlation and regression analysis**

### Business Example: Advertising Spend vs. Sales Revenue

A company studies the relationship between **monthly advertising expenditure (₹ lakhs)** and **sales revenue (₹ lakhs)** across 12 months.

Advertising Spend (₹ Lakhs)	Sales Revenue (₹ Lakhs)
5	40
6	42
8	50
10	55
12	65
14	70
15	75
16	78
18	85
20	90
22	95
25	105

### Interpretation:

- The data shows a **positive relationship**: as advertising spend increases, sales revenue also increases.
- For example, spending **₹10 lakhs on advertising** corresponds to around **₹55 lakhs in sales**, while **₹25 lakhs in advertising** corresponds to **₹105 lakhs in sales**.
- This bivariate dataset prepares the ground for applying **correlation analysis** (to measure the strength of the relationship) and **regression analysis** (to predict sales based on advertising spend).

### 2.5.2 Construction of Bivariate Frequency Tables

To create a **bivariate frequency table**, follow these steps:

1. **Identify two variables** (e.g., Variable X = study hours, Variable Y = test scores).
2. **Group the data** for each variable into appropriate class intervals.

3. Create a **two-dimensional table** with:
  - One variable represented across the **rows**
  - The other variable along the **columns**
4. **Count the frequency** for each pair of class intervals and fill the cells accordingly.

**Example Table:**

Study Hours ↓ / Scores →	40–50	50–60	60–70	Total
0–2 hours	5	3	2	10
2–4 hours	2	6	5	13
4–6 hours	1	4	7	12
Total	8	13	14	35

This table shows how many students fall into each **combination** of study time and score ranges.

### 2.5.3 Applications of Bivariate Frequency Distribution

Bivariate frequency distributions are widely used in both academic and professional settings to understand **relationships between two variables**. Some key applications include:

1. **Education:**  
Comparing students' performance with attendance, study hours, or participation.
2. **Business and Marketing:**  
Studying the relationship between **advertising spend** and **sales volume**, or between **price** and **customer demand**.
3. **Healthcare:**  
Analysing the correlation between **age group** and **blood pressure**, or between **diet** and **BMI**.
4. **Social Research:**  
Exploring links between variables like **income level** and **education**, or **location** and **internet access**.
5. **Further Statistical Analysis:**

Bivariate tables serve as a foundation for calculating:

- **Correlation coefficients** (e.g., Pearson's  $r$ )
- **Regression equations**
- **Contingency tables** for categorical data

These applications help in **predictive modeling, decision-making, and hypothesis testing** across multiple fields.

### Did You Know?

“Did you know that **bivariate frequency tables** are the foundation of advanced statistical tools like **correlation and regression**? These simple two-way tables are often used in early research stages to explore relationships between two variables—like time spent studying and test scores, or advertising spend and product sales. They help identify patterns and can even guide future predictions, making them one of the most powerful tools in practical statistics.”

## 2.6 Graphical Representations of Frequency Distributions

Graphs and charts make statistical data easier to **understand, compare, and communicate**. While tables are essential for accuracy and detail, **graphical representation** helps highlight trends, proportions, and relationships at a glance. Depending on whether the data is **qualitative** or **quantitative**, different types of graphs are used.

### 2.6.1 Graphical Representation for Qualitative Data (Bar Charts, Pie Charts)

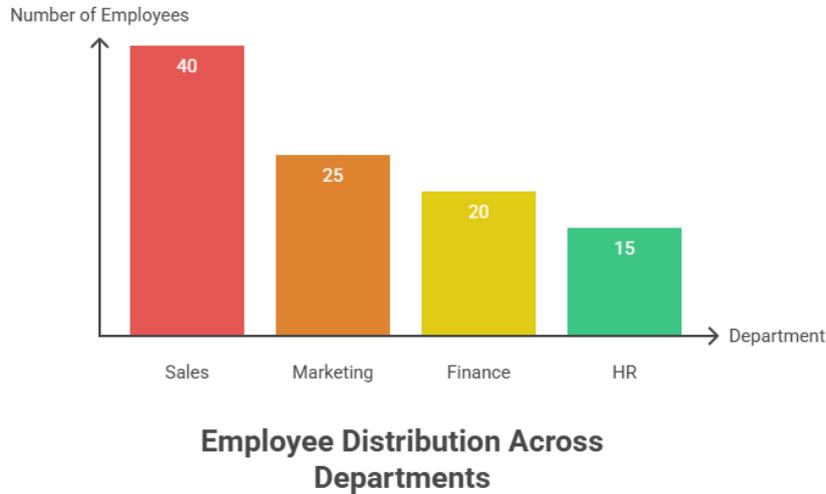
**Qualitative data** refers to non-numeric information based on categories, such as department, type of product, or customer preference. Since we can't calculate numerical averages for these categories, they are best represented through **visual methods** like bar charts and pie charts.

#### Bar Charts

- Display data using rectangular bars.
- Each bar represents a category, and the height or length of the bar represents the **frequency or percentage**.
- Bars are separated by spaces to show distinct categories.
- Useful for comparing frequencies across categories.

**Business Example:**

A company categorises its employees based on departments.



**Fig.2.2. Bar Chart**

Department	Number of Employees
Sales	40
Marketing	25
HR	15
Finance	20

- A **bar chart** of this data clearly shows that the **Sales department** has the highest number of employees, while **HR** has the fewest.

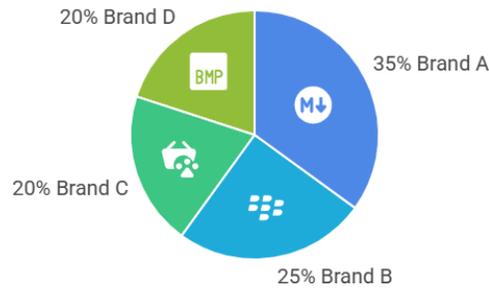
**Pie Charts**

- A circular chart divided into slices, where each slice represents a category’s proportion out of the total.
- The **angle or area** of each slice corresponds to the category’s percentage share.
- Useful for showing the **composition** of data.

**Business Example:**

A smartphone company analyses the **market share of different brands** in a city:

**Market Share Distribution Among Brands**



**Fig.2.3. Pie Chart**

Brand	Market Share (%)
Brand A	35%
Brand B	25%
Brand C	20%
Brand D	20%

- A **pie chart** of this data shows that **Brand A dominates the market (35%)**, followed by **Brand B (25%)**.

### 2.6.2 Graphical Representation for Quantitative Data (Histogram, Polygon, Ogive)

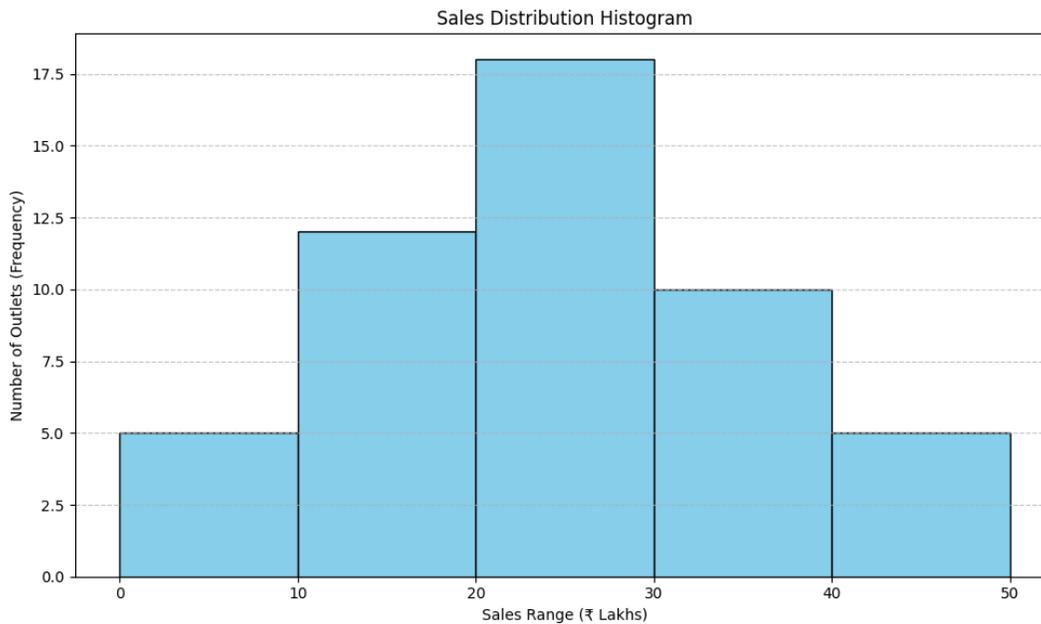
**Quantitative data** is numerical in nature and often involves **continuous or discrete values**. To represent such data visually, we use the following tools:

#### Histogram

- A type of bar graph for **continuous data**.
- Bars represent **class intervals**, and height represents **frequency or frequency density**.
- Unlike bar charts, the bars in a histogram are adjacent, showing **continuity**.
- Used to visualise the **distribution of data** (e.g., normal or skewed patterns).

#### Business Example:

A company records the **monthly sales (₹ lakhs)** of 50 outlets.



**Fig.2.4. Histogram**

Sales Range (₹ Lakhs)	Frequency (No. of Outlets)
0 – 10	5
10 – 20	12
20 – 30	18
30 – 40	10
40 – 50	5

A histogram of this data shows that most outlets achieved sales between **20–30 lakhs**, forming the peak of the distribution.

### Frequency Polygon

- A **line graph** constructed by plotting the **midpoints of class intervals** on the x-axis and frequencies on the y-axis.
- Points are connected using straight lines.
- Can be used **with or without a histogram** for comparison.
- Useful for comparing **two or more distributions** on the same graph.

### Business Example:

The company compares the **monthly sales of two product categories (A & B)**.

Sales Range (₹ Lakhs)	Frequency (Product A)	Frequency (Product B)
0 – 10	4	6
10 – 20	10	8
20 – 30	15	12
30 – 40	8	10
40 – 50	3	4

- A frequency polygon shows that **Product A peaked at 20–30 lakhs**, while **Product B had a more balanced spread between 10–40 lakhs**.

### Ogive (Cumulative Frequency Curve)

- A curve that represents **cumulative frequency distribution**.
- Two types: **Less-than ogive** and **Greater-than ogive**.
- Plotted using **class boundaries** on the x-axis and **cumulative frequencies** on the y-axis.
- Useful in estimating **median, quartiles, and percentiles**.

### Business Example:

Using the sales data of 50 outlets:

Sales Range (₹ Lakhs)	Frequency	Less-than Cumulative Frequency
0 – 10	5	5
10 – 20	12	17
20 – 30	18	35
30 – 40	10	45
40 – 50	5	50

- A **less-than ogive** shows that **35 outlets recorded sales below ₹30 lakhs**.
- The **median sales value** can be estimated from the point where the ogive crosses the 25th outlet ( $N/2$ ).

### 2.6.3 Comparison of Graphical Methods

Different graphical tools are chosen based on:

- **Type of data** (qualitative vs. quantitative)

- **Purpose of analysis** (comparison, distribution, proportion)
- **Ease of interpretation** for the audience

Graph Type	Suitable For	Features
Bar Chart	Qualitative data	Separate bars for categories
Pie Chart	Qualitative data	Shows percentage share of each category
Histogram	Quantitative data	Adjacent bars; used for continuous intervals
Frequency Polygon	Quantitative data	Line graph using midpoints
Ogive	Quantitative data	Cumulative frequency curve

**Choosing the right graph** depends on the nature of your data and what you want to communicate:

- Use **bar and pie charts** when dealing with **categories**.
- Use **histograms and polygons** to understand **distribution and frequency**.
- Use **ogives** to estimate **central tendencies and cut-off points** in cumulative data.

### Knowledge Check 1

**Choose the correct option:**

1. What is the purpose of a frequency distribution table?
  - A) To show the mean of data
  - B) To show the number of values below average
  - C) To organise data into classes and count frequencies
  - D) To eliminate all outliers
2. If a dataset has unequal class intervals, what should be used to construct a fair histogram?
  - A) Relative frequency
  - B) Cumulative frequency
  - C) Frequency density
  - D) Midpoints
3. What does a less-than cumulative frequency distribution show?
  - A) Frequencies from highest to lowest
  - B) The number of observations greater than a class boundary

- C) The total number of observations
- D) The number of observations less than the upper boundary of a class
4. Which of the following is **true** for a pie chart?
- A) It is best for showing changes over time
- B) It is suitable only for quantitative data
- C) It uses angles to represent percentage frequencies
- D) It shows cumulative frequencies
5. In a bivariate frequency distribution, the two variables are usually shown:
- A) Both in rows
- B) Both in columns
- C) One in rows and one in columns
- D) Only as percentages

## 2.7 Summary

- ❖ This unit introduced various methods of organising and representing frequency data to make analysis meaningful and efficient. Starting with the **construction of frequency distributions**, the unit explained how raw data can be summarised into discrete and continuous frequency tables. It then moved to **relative and percentage frequency distributions**, which help in comparing datasets and understanding proportions more clearly.
- ❖ The concept of **cumulative frequency** was discussed with both less-than and greater-than types, along with **ogives** to represent them graphically. The unit then covered **frequency density**, a tool used when class intervals are unequal, ensuring accurate histogram construction. In more complex scenarios, **bivariate frequency distributions** were introduced to examine the joint behaviour of two variables.
- ❖ Finally, the unit explained various **graphical methods** such as bar charts, pie charts, histograms, frequency polygons, and ogives, comparing their purposes and uses for both qualitative and quantitative data. Together, these tools provide a foundation for data summarisation, interpretation, and decision-making.

## 2.8 Key Terms

1. **Frequency** - Number of times a value or class occurs in a dataset.
2. **Class Interval** - A range of values used to group data in a frequency table.

3. **Relative Frequency** - Proportion of total observations that fall within a class ( $\text{Frequency} \div \text{Total}$ ).
4. **Percentage Frequency** - Relative frequency expressed as a percentage.
5. **Cumulative Frequency** - Running total of frequencies either less-than or greater-than a class boundary.
6. **Frequency Density** - Frequency divided by class width, used in histograms with unequal intervals.
7. **Histogram** - A bar graph representing the distribution of quantitative data.
8. **Ogive** - A line graph showing cumulative frequency distribution.
9. **Bivariate Distribution** - A frequency table representing the relationship between two variables.
10. **Bar Chart** - A graph with spaced bars used for categorical data.
11. **Pie Chart** - A circular chart divided into slices to illustrate proportions.

## 2.9 Descriptive Questions

1. Define a frequency distribution. How does it help in organising raw data?
2. What is the difference between discrete and continuous frequency distribution? Provide examples.
3. How is relative frequency calculated and why is it useful?
4. Explain the steps involved in drawing a histogram when class intervals are unequal.
5. What is cumulative frequency? How are less-than and greater-than cumulative frequencies different?
6. Describe how a bivariate frequency table is constructed. Mention any one application.
7. Compare histograms and bar charts in terms of structure and use.
8. Why is frequency density important in statistics? When must it be used?
9. What information can be interpreted from an ogive?
10. Which graph would you use to represent the popularity of different fruit juices in a survey? Why?

## 2.10 References

1. Gupta, S.P. (2014). *Statistical Methods*. Sultan Chand & Sons.
2. Sharma, J.K. (2018). *Business Statistics*. Pearson Education.
3. Goon, A.M., Gupta, M.K., & Dasgupta, B. (2013). *Fundamentals of Statistics*. World Press.
4. Indian Statistical Institute. *Introductory Statistics Course Notes*.
5. Government of India. *Ministry of Statistics and Programme Implementation* – [www.mospi.gov.in](http://www.mospi.gov.in)
6. NCERT. (2021). *Statistics for Economics* (Class XI Textbook).

7. Data sources and templates from [www.data.gov.in](http://www.data.gov.in) and official survey reports.

### Answers to Knowledge Check

#### *Knowledge Check 1*

1. C) To organise data into classes and count frequencies
2. C) Frequency density
3. D) The number of observations less than the upper boundary of a class
4. C) It uses angles to represent percentage frequencies
5. C) One in rows and one in columns

## 2.11 Case Study

### Analysing Class Performance through Frequency Distributions

#### Introduction

In business organisations, data is continuously collected from areas such as sales, marketing, finance, and operations. However, raw data remains unhelpful unless it is systematically organised and analysed. Frequency distribution techniques allow managers to transform large, scattered datasets into meaningful insights that guide decision-making.

In this case, the management of **BrightMart Retail Ltd.** sought to evaluate the **monthly sales performance of its sales representatives across three regions—North, South, and West.** The objective was to identify performance clusters, compare regional outcomes, and assess whether employee training hours influenced sales. The data analytics team was given the task of compiling and interpreting the sales data using **frequency distribution tools.**

#### Background

Each region had approximately **40 sales representatives.** The raw sales figures (in ₹ lakhs) were received in different formats:

- **North region:** Sales submitted as individual transaction records.
- **South region:** Data grouped but with unequal class intervals.
- **West region:** Summary with frequencies but no cumulative values.

The company wanted to:

- Identify the **sales ranges** in which most employees fell.
- Compare relative performance between the three regions.
- Explore whether **training hours** had a measurable impact on sales.

To achieve this, the analytics team applied statistical tools such as **relative frequency, cumulative frequency, frequency density, and bivariate frequency distribution.**

### Problem Statement 1: Inconsistency in Data Format and Lack of Summarisation

**Dataset (Sales in ₹ Lakhs – Sample Extract):**

Region	Raw Input Type	Example Data Provided
North	Individual sales data	12, 15, 18, 20, 25, 28, 30 ...
South	Grouped, unequal widths	0–10 (4 reps), 10–20 (12 reps), 20–40 (16 reps), 40–60 (8 reps)
West	Frequencies only	10–20 (6 reps), 20–30 (14 reps), 30–40 (12 reps), 40–50 (8 reps)

**Solution:**

- All datasets were **standardised** into continuous frequency tables.
- For South (unequal intervals), **frequency density** was computed:

Example (South Region):

Sales Range (₹ Lakhs)	Frequency	Class Width	Frequency Density
0–10	4	10	0.4
10–20	12	10	1.2
20–40	16	20	0.8
40–60	8	20	0.4

- Histograms based on **frequency density** corrected visual distortion.
- **Less-than ogives** were drawn to show cumulative sales performance.

**Interpretation:**

- South had the **largest spread** but most employees in the **10–20 lakhs band**.
- North’s ogive showed a **steeper rise near 25–30 lakhs**, indicating more concentration of employees in higher sales.

### Problem Statement 2: Need for Region-Wise Performance Comparison

### Dataset (Relative Frequencies – Simplified):

Region	Sales 20–30 Lakhs (%)	Sales 30–40 Lakhs (%)	Sales Above 40 Lakhs (%)
North	35%	25%	20%
South	20%	15%	10%
West	25%	30%	15%

#### Solution:

- **Relative frequency (%)** was calculated for each region.
- This allowed comparisons even though each region had slightly different staff totals.
- Visual aids (pie charts & histograms) made differences clear.

#### Interpretation:

- North had **more high achievers (above 40 lakhs)**.
- South underperformed, with most staff below 30 lakhs.
- West showed a balanced spread but fewer top performers than North.

### Problem Statement 3: Understanding the Role of Training Hours in Performance

#### Dataset (Bivariate Frequency Table – Training Hours vs. Sales):

Training Hours (per month)	Sales 10–20 Lakhs	Sales 20–30 Lakhs	Sales 30–40 Lakhs	Sales 40+ Lakhs	Total
0–5 hrs	8	6	2	1	17
5–10 hrs	4	10	6	3	23
10+ hrs	2	4	7	7	20

#### Solution:

- Constructed a **bivariate frequency distribution** linking **training hours** with **sales performance**.
- Results showed a **positive association** between more training and higher sales.

#### Interpretation:

- Among employees with **10+ hours training**, 70% achieved sales above 30 lakhs.
- Those with **0–5 hours** mostly remained below 20 lakhs.
- This insight supported management’s plan to **invest in advanced training programs**.

### MCQ

**Which statistical method helped the company compare performance between regions of different team sizes?**

- A) Cumulative frequency distribution
- B) Histogram with class width adjustment
- C) Relative and percentage frequency distribution
- D) Bivariate frequency table

**Answer:** C) Relative and percentage frequency distribution

**Explanation:** Relative frequencies express each class as a proportion of the total, making it easier to compare datasets of different sizes.

### **Conclusion**

By applying frequency distribution techniques, BrightMart Retail Ltd. converted unstructured sales data into actionable insights.

- **Problem 1:** Frequency density and ogives standardised inconsistent data formats.
- **Problem 2:** Relative frequencies allowed **region-wise comparisons** despite varying staff sizes.
- **Problem 3:** Bivariate analysis revealed the **impact of training on sales performance**.

This case shows how statistical tools such as **frequency density, cumulative frequency, relative frequency, and bivariate analysis** are indispensable in **business decision-making**—from identifying performance gaps to designing training and incentive strategies.

## Unit 3: Probability

### Learning Objectives

1. Understand the concept and real-world significance of probability, including its role in decision-making, risk assessment, and predictive modeling in various fields such as business, science, and economics.
2. Familiarise with key probability-related terminology such as experiment, outcome, sample space, event, mutually exclusive events, exhaustive events, and complementary events, for precise communication of statistical scenarios.
3. Differentiate between classical, empirical, and axiomatic definitions of probability, and apply the appropriate approach depending on the nature of the problem or dataset.
4. Apply basic theorems on probability, including the addition and multiplication rules, to calculate the likelihood of compound events, with and without replacement.
5. Understand and calculate conditional probability, recognising its importance in contexts where outcomes are dependent on prior occurrences, and apply it using structured methods like tree diagrams and formulas.
6. Apply the multiplicative rule for independent events, identifying when events are truly independent and using the rule to solve multi-stage probability problems effectively.

### Content

- 3.0 Introductory Caselet
- 3.1 Introduction to Probability
- 3.2 Important Terms and Concepts
- 3.3 Definitions of Probability
- 3.4 Theorems on Probability
- 3.5 Conditional Probability
- 3.6 Multiplicative Theorem for Independent Events
- 3.7 Bayes' Theorem
- 3.8 Summary
- 3.9 Key Terms
- 3.10 Descriptive Questions
- 3.11 References
- 3.12 Case Study

### 3.0 Introductory Caselet

#### “Raj’s Risk Radar: Making Smarter Predictions with Probability”

Raj, a 27-year-old data analyst at a logistics startup in Hyderabad, was facing a growing challenge. His job involved forecasting shipment delays across different regions. The company prided itself on its “on-time delivery promise,” but over the last quarter, customer complaints had started to rise. Shipments were delayed due to various unpredictable factors—weather, traffic, vehicle breakdowns, or staff shortages. Raj knew he needed a better way to understand and quantify these uncertainties.

During a weekend data science workshop, Raj was introduced to the world of **probability theory**. He learned that not all uncertainties are random guesses—some can be **measured, modeled, and predicted**. More importantly, he discovered how **conditional probability** and **Bayes’ Theorem** could help him revise predictions when new information was available.

Back at work, Raj started by identifying key events that impacted delays, such as **rainy weather, peak traffic hours, and driver availability**. He built **sample spaces** and defined **events** like “delay due to rain” and “delay during rush hour.” By applying **classical probability**, he estimated chances based on equal likelihood. But he quickly realised that not all events were equally likely.

Switching to the **relative frequency approach**, he pulled past delivery records and calculated that **35% of delays** occurred on rainy days. Then, using **conditional probability**, he asked more precise questions like: *“Given that it’s raining and traffic is heavy, what’s the probability of a delay?”*

Finally, Raj applied **Bayes’ Theorem** to update the chance of a vehicle being late once he received a weather alert. He discovered that **probability could not only predict risk—but adapt to changing information**. By using **addition and multiplication theorems**, he calculated combined probabilities and created a daily delay prediction dashboard for the operations team.

The results were remarkable. Within a month, on-time delivery improved by 15%, and the team could proactively reroute shipments on high-risk days. What started as vague unpredictability became a structured, data-informed decision process.

**Critical Thinking Question:**

If you were Raj, how would you use conditional probability or Bayes' Theorem to improve decision-making in a different domain—such as health, education, or sports? Share one example where an outcome depends on prior information.

### 3.1 Introduction to Probability

Probability is a fundamental concept in statistics and mathematics that deals with **uncertainty and chance**. Whether predicting the outcome of a coin toss, assessing the likelihood of rain, or calculating business risks, probability helps quantify how likely an event is to happen.

#### 3.1.1 Meaning and Importance of Probability in Statistics

##### Meaning:

Probability is a **numerical measure of the likelihood** that a particular event will occur. It ranges between **0 and 1**, where:

- 0 means the event is **impossible**.
- 1 means the event is **certain**.
- A value like 0.5 means the event is **equally likely** to happen or not.

In simple terms:

- If an event has a **high probability**, it is **very likely** to occur.
- If it has a **low probability**, it is **unlikely** to occur.

##### Importance in Statistics:

- Probability allows statisticians to **draw conclusions** about populations based on samples.
- It is the **foundation for inferential statistics**, helping estimate the reliability of predictions.
- Many statistical methods such as **hypothesis testing**, **confidence intervals**, and **risk analysis** are based on probability theory.
- It helps in dealing with **uncertainty** in real-world data.

#### 3.1.2 Applications of Probability in Real Life

Probability is not just a mathematical concept—it is widely applied in **real-life decision-making** and various fields.

Here are some examples:

##### 1. **Weather Forecasting:**

Meteorologists use probability to predict rain, storms, or sunshine based on historical and current data.

##### 2. **Insurance:**

Insurance companies assess **risk probabilities** to set premiums—for example, the probability of accidents, illness, or natural disasters.

### 3. **Business and Finance:**

Businesses use probability to forecast **demand, sales, or stock prices**. Risk analysis models use probability to calculate the likelihood of project failure or loss.

### 4. **Games and Sports:**

In card games or sports betting, probabilities are used to calculate odds and expected outcomes.

### 5. **Medicine and Healthcare:**

Doctors and researchers use probability to understand the **likelihood of disease occurrence, treatment success, or side effects**.

### 6. **Manufacturing:**

In quality control, probability helps estimate the chances of product defects in a production batch.

## 3.1.3 Random Experiments, Sample Space, and Events

To study probability, we need to understand three core terms: **random experiments, sample space, and events**.

### 1. **Random Experiment:**

A process or activity that leads to an outcome that **cannot be predicted with certainty**.

#### **Examples:**

- Tossing a coin
- Rolling a die
- Drawing a card from a deck

Each time the experiment is repeated, the outcome may vary.

### 2. **Sample Space (S):**

The **set of all possible outcomes** of a random experiment.

#### **Examples:**

- For tossing a coin:  
Sample Space,  $S = \{\text{Head, Tail}\}$
- For rolling a die:  
 $S = \{1, 2, 3, 4, 5, 6\}$

### 3. **Event (E):**

An **event** is a **subset of the sample space**. It consists of one or more outcomes.

#### **Examples:**

- Getting an even number when a die is rolled:

$$E = \{2, 4, 6\}$$

- Drawing a red card from a deck of cards:

$$E = \{\text{all 26 red cards from the deck}\}$$

### Types of Events:

- **Simple Event:** Consists of only one outcome (e.g., rolling a 3).
- **Compound Event:** Consists of more than one outcome (e.g., rolling an even number).
- **Certain Event:** Always occurs (e.g., getting a number between 1 and 6 when rolling a standard die).
- **Impossible Event:** Never occurs (e.g., rolling a 7 on a standard die).

## 3.2 Important Terms and Concepts

Before diving into probability calculations, it is essential to understand the nature and relationships of **events** within a sample space. Each term in this section describes how events interact or relate to one another, which influences how probabilities are calculated.

### 3.2.1 Mutually Exclusive Events

#### Definition:

Two or more events are said to be **mutually exclusive** if they **cannot occur at the same time**. In other words, the occurrence of one event **prevents** the occurrence of the other.

#### Example:

- When you toss a coin, the result can be either a **head** or a **tail**, but not both. So, the events "Head" and "Tail" are **mutually exclusive**.
- When rolling a die, the events “getting a 2” and “getting a 5” are also mutually exclusive.

#### Mathematically:

If A and B are mutually exclusive events:

$$P(A \cap B) = 0$$

### 3.2.2 Exhaustive Events

**Definition:**

Events are **exhaustive** if they **cover the entire sample space**, meaning **at least one of them must occur** when the experiment is conducted.

**Example:**

- When rolling a die, the events {1}, {2}, {3}, {4}, {5}, and {6} are exhaustive because one of these outcomes **must** occur.
- In tossing a coin, the events “Head” and “Tail” are exhaustive.

**Note:**

A set of events can be **mutually exclusive and exhaustive** at the same time.

### 3.2.3 Equally Likely Events

**Definition:**

Events are **equally likely** if each has the **same probability** of occurring.

**Example:**

- When tossing a fair coin, the events "Head" and "Tail" are equally likely, each with a probability of 0.5.
- When rolling a fair die, the outcomes 1 through 6 are all equally likely, each with a probability of  $1/6$ .

**Important****Point:**

Equally likely events often occur in situations involving **fair or unbiased instruments** (like coins, dice, cards).

### 3.2.4 Independent and Dependent Events

#### 1. Independent Events

Two events are **independent** if the occurrence of one **does not affect** the probability of the other.

**Example:**

- Tossing a coin and rolling a die are independent events.  
The result of the coin toss has no effect on the die roll.

**Mathematically:**

If A and B are independent,

$$P(A \cap B) = P(A) \times P(B)$$

---

## 2. Dependent Events

Two events are **dependent** if the outcome of one **affects** the probability of the other.

### Example:

- Drawing two cards from a deck **without replacement**:  
After drawing the first card, the total number of cards reduces, which affects the probability of the second draw.

## 3.2.5 Complementary Events

### Definition:

The **complement of an event A** (written as  $A'$  or  $A^c$ ) is the event that **A does not occur**.  
A and  $A'$  are **mutually exclusive and exhaustive**.

### Example:

- If event  $A$  = "getting an even number on a die",  
then  $A'$  = "getting an odd number".

### Mathematically:

$$P(A) + P(A') = 1$$

$$\text{So, } P(A') = 1 - P(A)$$

This relationship is often used to solve problems more efficiently by calculating the **complement** of the desired event.

## 3.3 Definitions of Probability

Probability can be defined in multiple ways depending on the context in which it is used. The three major approaches are:

- **Classical (Theoretical) Probability**
- **Relative Frequency (Empirical) Probability**
- **Axiomatic Probability**

Each approach has its own applications, limitations, and assumptions.

### 3.3.1 Classical Definition

**Definition:**

The **classical definition** (also known as theoretical probability) is based on the assumption that all **outcomes are equally likely**. It is used when each outcome in the sample space has the same chance of occurring.

**Formula:**

$$P(E) = (\text{Number of favourable outcomes}) \div (\text{Total number of outcomes})$$

**Example:**

- Tossing a fair coin:

$$P(\text{Head}) = 1 \div 2 = 0.5$$

- Rolling a fair six-sided die:

$$P(3) = 1 \div 6$$

$$P(\text{Even number}) = 3 \div 6 = 0.5$$

**Conditions:**

- The outcomes must be **mutually exclusive** and **equally likely**.
- The sample space must be **finite and countable**.

**Limitations:**

- Cannot be used when outcomes are **not equally likely**.
- Not suitable for **complex or real-world events** where symmetry doesn't exist.

### 3.3.2 Relative Frequency Definition

**Definition:**

The **relative frequency definition** (also known as empirical probability) defines probability based on **observed data or experiments**. It is calculated by **repeating an experiment** many times and observing the proportion of times the event occurs.

**Formula:**

$$P(E) = (\text{Number of times the event occurred}) \div (\text{Total number of trials})$$

**Example:**

If it rained on 30 out of the past 100 days:

$$P(\text{Rain}) = 30 \div 100 = 0.3$$

**Uses:**

- Used in **real-life experiments** and **observational studies**.

- Helps estimate probabilities where theoretical probabilities are difficult to assign.

**Advantages:**

- Reflects actual conditions.
- Applicable even when outcomes are **not equally likely**.

**Limitations:**

- Accuracy improves only with **large number of trials**.
- Not useful for **unique or one-time events**.

### 3.3.3 Axiomatic Definition

**Definition:**

The **axiomatic definition** of probability is a **modern and most general approach** introduced by **Andrey Kolmogorov**. It defines probability through a set of **logical rules (axioms)** rather than relying on experiments or equal likelihood.

This approach treats probability as a **function (P)** that assigns a number to each event in a sample space **S**, following three basic axioms:

**Kolmogorov's Axioms:****1. Non-negativity:**

$$P(E) \geq 0 \text{ for any event } E$$

**2. Certainty:**

$$P(S) = 1 \text{ (the probability of the entire sample space is 1)}$$

**3. Additivity (for mutually exclusive events):**

If A and B are mutually exclusive, then

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

**Advantages:**

- Can handle **infinite sample spaces**.
- Applies to **all types of events**, including **continuous probability**.
- Lays the foundation for **advanced statistical theory**.

**Example Use:**

- Axiomatic probability is used in **advanced mathematics, finance, machine learning, and statistical modeling** where classical assumptions may not hold.

### Did You Know?

“Did you know that the **axiomatic definition of probability** is the most **universal and mathematically rigorous** approach to probability? It doesn’t depend on physical experiments or equally likely outcomes. Instead, it uses a set of logical rules called **Kolmogorov’s axioms**, which apply to **finite, infinite**, and even **continuous sample spaces**. This is the foundation for most **modern probability theories** used in fields like **machine learning, quantum mechanics, and financial risk modeling.**”

### Summary Table:

Definition Type	Based On	Formula	Suitable For
Classical	Equally likely cases	$P(E) = \text{Favourable outcomes} \div \text{Total outcomes}$	Theoretical or ideal cases
Relative Frequency	Experimental results	$P(E) = \text{Trials with E} \div \text{Total trials}$	Real-life, data-based scenarios
Axiomatic	Logical postulates	Uses Kolmogorov's axioms	General and complex probability models

## 3.4 Theorems on Probability

In probability, events may happen **individually, together, or with overlap**. The **addition theorem** helps us calculate the probability of **either one event or another** occurring. It is especially useful when events are **not mutually exclusive**, meaning they can happen at the same time.

### 3.4.1 Addition Theorem

The **Addition Theorem** provides a way to calculate the probability of the **union of two events**—that is, the probability that **at least one** of them occurs.

#### General Formula:

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

Where:

- $P(A \cup B)$  is the probability that either event A or event B (or both) occur.
- $P(A \cap B)$  is the probability that both A and B occur simultaneously.

This formula corrects for **double counting** the overlap between A and B.

### 3.4.2 Special Case of Addition Theorem for Mutually Exclusive Events

If **A and B are mutually exclusive** (i.e., they **cannot occur together**), then:

$$P(A \cap B) = 0$$

So the addition theorem simplifies to:

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

**Example:**

- Let A = drawing a red card
- Let B = drawing a black card

Since a single card cannot be both red and black, A and B are mutually exclusive.

If  $P(A) = 26/52$  and  $P(B) = 26/52$ , then:

$$P(A \cup B) = 26/52 + 26/52 = 1$$

This makes sense, as the card must be either red or black.

### 3.4.3 General Case of Addition Theorem

If events A and B **can both occur together**, then they are **not mutually exclusive**, and the **general formula** must be used:

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

**Example:**

- Let A = student takes mathematics
  - Let B = student takes science
- Some students take both subjects.

If:

- $P(A) = 0.6$
- $P(B) = 0.5$
- $P(A \cap B) = 0.3$

Then:

$$P(A \cup B) = 0.6 + 0.5 - 0.3 = 0.8$$

So, there is an 80% chance that a student takes **either math or science** (or both).

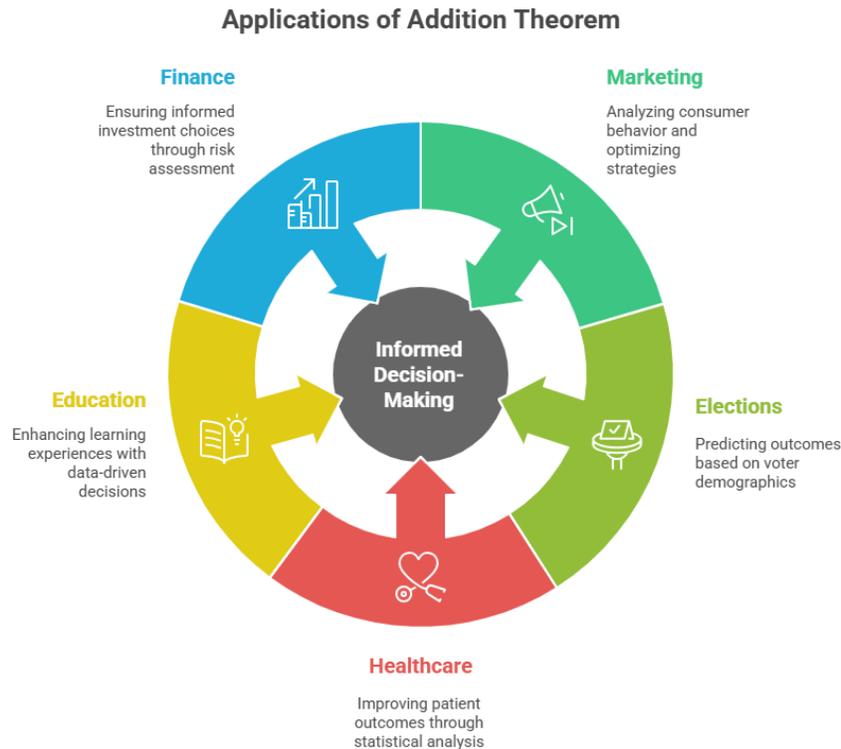
### “Activity Applying the Addition Theorem to Event Probabilities”

#### **Instruction to Student:**

Consider a school where 70% of students like playing football, 50% like playing cricket, and 30% like both sports.

1. Use the **general addition theorem** to calculate the probability that a randomly selected student likes **either football or cricket**.
2. Apply the formula:  
$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$
3. After calculating, represent the situation using a **Venn diagram**, clearly showing the overlapping region.
4. Submit your calculations and the diagram, and briefly explain why subtracting the intersection is necessary in this case.

#### **3.4.4 Applications of Addition Theorem**



**Fig.3.1. Applications of Addition Theorem**

The **addition rule** is widely used in both **academic problems** and **real-world scenarios** involving probability.

**Common applications include:**

1. **Marketing:**

Probability that a customer buys product A or product B.

2. **Elections:**

Probability that a voter supports either candidate A or candidate B.

3. **Healthcare:**

Probability that a patient shows symptom X or symptom Y.

4. **Education:**

Determining how many students take at least one of two electives.

5. **Finance:**

Calculating risk exposure when two market events could happen simultaneously.

### Why is it important?

- Helps avoid **overestimating probability** when events can overlap.
- Provides a **foundation for solving compound probability problems**.
- Often combined with **Venn diagrams** to visualise overlaps.

### Visual Tip:

In a Venn diagram, the addition theorem corresponds to the **combined area of both circles**, minus the **overlap**.

This helps learners visualise  $P(A \cup B)$  intuitively.

## 3.5 Conditional Probability

In many real-world scenarios, the probability of an event depends on whether another event has already occurred.

This is where **conditional probability** comes into play.

### 3.5.1 Concept of Conditional Probability

#### Definition:

Conditional probability refers to the **probability of an event occurring given that another event has already occurred**. It is denoted by  $P(A | B)$ , which means “**the probability of event A given that event B has occurred.**”

This concept is useful when:

- Events are **not independent**
- You have **partial information** about the outcome
- You want to **update** the probability based on new evidence

#### Example:

Suppose we know that a student has passed mathematics. What is the probability that the same student has also passed science?

This is a **conditional** question, where one event (passing math) affects our estimation of the second (passing science).

**Did You Know?**

“Did you know that in some cases, the **probability of an event can increase or decrease drastically once new information is known**? For example, the probability of someone having a disease might be low in the general population, but if they test positive on a medical test, the probability rises significantly. This is the basis for **diagnostic decision-making**, where **conditional probability** helps doctors interpret lab results based on patient risk factors and symptoms.”

### 3.5.2 Formula and Examples

#### Formula for Conditional Probability:

If A and B are two events, and  $P(B) \neq 0$ , then:

$$P(A | B) = P(A \cap B) \div P(B)$$

This means:

The probability that A occurs **given** that B has occurred =

The probability that **both A and B occur**, divided by the probability that **B occurs**.

#### Example 1:

A card is drawn from a deck.

Let:

- A = event that the card is a king
- B = event that the card is red

There are 2 red kings out of 26 red cards, so:

$$P(A \cap B) = 2/52$$

$$P(B) = 26/52$$

So:

$$P(A | B) = (2/52) \div (26/52) = 2/26 = 1/13$$

#### Interpretation:

Given that the card is red, the chance that it is a king is 1/13.

#### Example 2:

A box contains 3 red and 2 blue balls. A ball is drawn **without replacement**, and then a second ball is drawn.

Let:

- A = first ball is red
- B = second ball is red

We want to find  $P(B | A)$ : the probability that the second ball is red, **given** that the first was red.

If the first red ball is taken out, 2 red and 2 blue balls remain.

So,

$$P(B | A) = 2/4 = 0.5$$

### “Activity Exploring Conditional Probability through Classroom Data”

#### Instruction to Student:

Your class recently conducted a quiz in which 30 students participated. Out of them:

- 18 passed in Science
- 20 passed in Mathematics
- 12 passed in **both** Science and Mathematics

1. Use this data to calculate:

a)  $P(\text{Passed Science} | \text{Passed Math})$

b)  $P(\text{Passed Math} | \text{Passed Science})$

2. Apply the **conditional probability formula**:

$$P(A | B) = P(A \cap B) \div P(B)$$

3. Write a short reflection (4–5 lines) on how the probability changes when we already know a student passed one subject.

### 3.5.3 Properties of Conditional Probability

Here are some important properties:

1. **Multiplication Rule (Rewriting Joint Probability):**

$$P(A \cap B) = P(B) \times P(A | B)$$

or

$$P(A \cap B) = P(A) \times P(B | A)$$

(Choose based on known condition)

2. **If A and B are independent events:**

Then:

$$P(A | B) = P(A)$$

and

$$P(B | A) = P(B)$$

(Because one does not affect the other)

### 3. Probability of an Event Given Itself:

$$P(A | A) = 1$$

### 4. Zero Conditional Probability:

If events A and B **cannot happen together** (mutually exclusive):

$$P(A | B) = 0$$

### 5. Total Probability from Conditional Events:

The overall probability of an event can be calculated by **weighing its conditional probabilities** against different conditions (leads into Bayes' Theorem in next sections).

## 3.6 Multiplicative Theorem for Independent Events

The **multiplicative theorem** allows us to compute the probability of **multiple independent events occurring together**. It's a powerful tool in situations where one event does not influence the other—such as tossing a coin and rolling a die at the same time.

### 3.6.1 Statement of the Theorem

#### Definition of Independent Events:

Two events A and B are said to be **independent** if the occurrence of one **does not affect** the occurrence of the other.

#### Statement of the Multiplicative Theorem:

If **A and B are independent**, then:

$$P(A \cap B) = P(A) \times P(B)$$

This means the probability that **both A and B occur** is simply the product of their individual probabilities.

#### Example 1: Tossing a coin and rolling a die

- A = Getting a head ( $P(A) = 0.5$ )

- B = Rolling a 4 ( $P(B) = 1/6$ )

Since these are independent events:

$$P(A \cap B) = 0.5 \times 1/6 = 1/12$$

### 3.6.2 Applications

The multiplicative theorem is useful in a wide range of practical scenarios, especially when dealing with **repeated or simultaneous experiments**.

#### Common Applications:

##### 1. Business and Finance:

Probability of two independent market factors occurring, e.g., a rise in oil prices and a fall in the dollar.

##### 2. Quality Control in Manufacturing:

- If the probability that a single product is defective is 0.02, the probability that **two independently selected products** are defective is:

$$0.02 \times 0.02 = 0.0004$$

##### 3. Gaming and Probability Experiments:

- Rolling two dice and calculating the probability of getting a 6 on both.

##### 4. Weather Prediction:

- Probability that it rains in Delhi and simultaneously in Mumbai (assuming weather in both cities is independent).

##### 5. Medical Studies:

- If the probability of a side effect occurring from Drug A is 0.1 and from Drug B is 0.05, and the drugs act independently, then the probability of both side effects occurring is:

$$0.1 \times 0.05 = 0.005$$

### 3.6.3 Extension to More than Two Events

The multiplicative theorem can be extended to **more than two independent events**.

Let  $A_1, A_2, A_3, \dots, A_n$  be **n independent events**, then:

$$P(A_1 \cap A_2 \cap A_3 \cap \dots \cap A_n) = P(A_1) \times P(A_2) \times P(A_3) \times \dots \times P(A_n)$$

#### Example 2:

If you flip **three coins**, the probability of getting heads on all three is:

$$P(\text{Head on coin 1}) = 0.5$$

$$P(\text{Head on coin 2}) = 0.5$$

$$P(\text{Head on coin 3}) = 0.5$$

So:

$$P(\text{All heads}) = 0.5 \times 0.5 \times 0.5 = 0.125$$

### Important Note:

This theorem applies **only** when events are **independent**. If events are dependent (i.e., one affects the other), we must use **conditional probability** instead.

### Did You Know?

“Did you know that when you're tossing multiple coins or rolling multiple dice, the probability of a specific combined outcome is calculated using the **extended multiplication rule** for **independent events**? For example, the probability of getting **three heads in a row** is not 0.5—it's  $0.5 \times 0.5 \times 0.5 = 0.125$ . This principle is used in **genetics, cryptography**, and even **digital communication systems** to calculate the chances of exact sequences occurring.”

## 3.7 Bayes' Theorem

Bayes' Theorem is a powerful tool in probability that allows us to **revise existing probabilities** when new information becomes available. It is especially useful in **real-world scenarios where outcomes depend on prior events or conditions**.

### 3.7.1 Statement and Explanation of Bayes' Theorem

#### Statement:

If  $B_1, B_2, \dots, B_n$  are **mutually exclusive and exhaustive events**, and  $A$  is any event that has occurred, then the **conditional probability** of event  $B_i$  given that  $A$  has occurred is given by:

$$P(B_i | A) = [P(B_i) \times P(A | B_i)] \div \Sigma[P(B_j) \times P(A | B_j)]$$

Where:

- $P(B_i | A)$  is the **posterior probability** (updated probability of  $B_i$  after  $A$  occurs)

- $P(B_i)$  is the **prior probability** of  $B_i$  before A occurs
- $P(A | B_i)$  is the **likelihood** (probability of A occurring given  $B_i$  is true)
- The denominator is the **total probability** of A across all possible B events

### Why Use Bayes' Theorem?

- It allows us to **update beliefs** based on new evidence.
- It is used when **reverse conditional probabilities** are needed.
- It answers: "**Given that A has happened, how likely is B?**"

### Simple Example:

A medical test is used to detect a rare disease that affects 1% of the population.

- $P(\text{Disease}) = 0.01$
- $P(\text{No disease}) = 0.99$
- $P(\text{Test positive} | \text{Disease}) = 0.95$  (true positive rate)
- $P(\text{Test positive} | \text{No disease}) = 0.05$  (false positive rate)

### What is the probability that a person who tests positive actually has the disease?

Let:

- D = person has disease
- D' = person does not have disease
- T = person tests positive

We want  $P(D | T) = ?$

Apply Bayes' Theorem:

$$P(D | T) = \frac{P(D) \times P(T | D)}{P(D) \times P(T | D) + P(D') \times P(T | D')}$$

Substitute values:

$$= \frac{0.01 \times 0.95}{0.01 \times 0.95 + 0.99 \times 0.05}$$

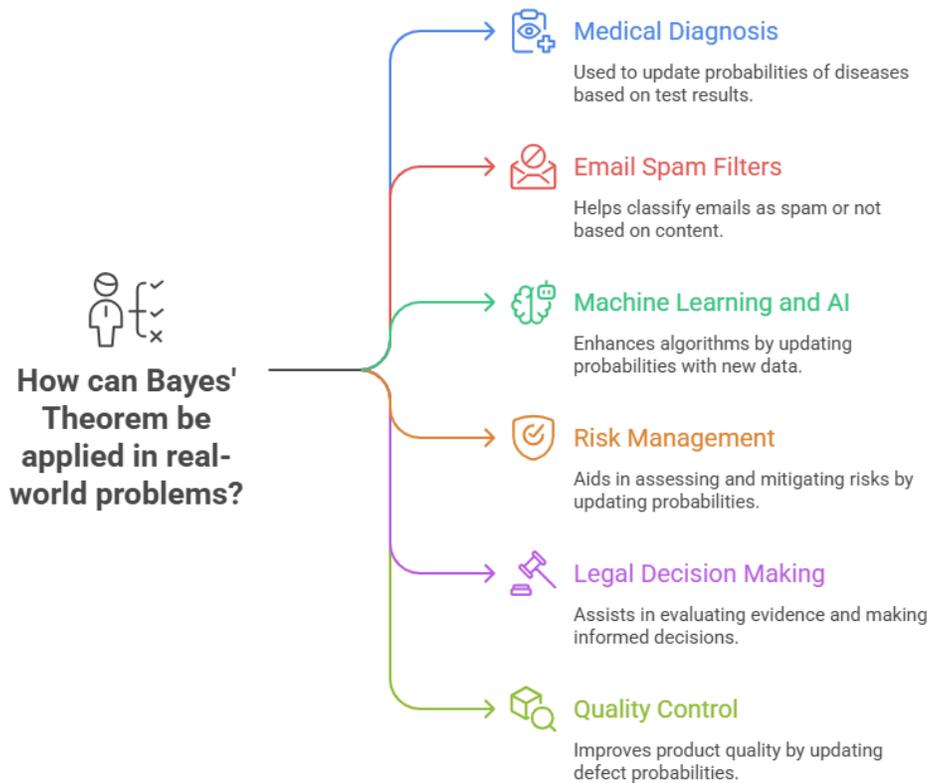
$$= 0.0095 \div (0.0095 + 0.0495)$$

$$= 0.0095 \div 0.059$$

$$\approx \mathbf{0.161} \text{ or } \mathbf{16.1\%}$$

So, even with a positive result, the person has only a **16.1% chance** of actually having the disease, due to the rarity of the disease and the false positive rate.

### 3.7.2 Applications of Bayes' Theorem in Real-world Problems



**Fig.3.2. Applications of Bayes' Theorem in Real-world Problems**

Bayes' Theorem is widely used in various fields to make **evidence-based predictions** and update prior beliefs. Here are some real-world applications:

#### 1. Medical Diagnosis

Used to calculate the **probability of a disease** given a positive or negative test result. Especially important when diseases are rare and test accuracy is not perfect.

#### 2. Email Spam Filters

Bayesian spam filters compute the **probability of an email being spam** based on the presence of certain words or phrases.

### 3. Machine Learning and AI

Bayes' Theorem is a core part of **Naive Bayes classifiers**, which are used in:

- Sentiment analysis
- Document classification
- Image recognition

### 4. Risk Management

Used in **finance and insurance** to assess how new market data affects the likelihood of risk events (e.g., market crashes, defaults).

### 5. Legal Decision Making

Used in evaluating **evidence in court cases**, where the probability of a suspect being guilty can be updated as new evidence (like DNA) is introduced.

### 6. Quality Control

In manufacturing, Bayes' Theorem helps determine the **likelihood of product defects** based on testing outcomes and batch history.

## Knowledge Check 1

**Choose the correct option:**

1. Which definition of probability is most appropriate when outcomes are not equally likely, and data is collected from actual observations?
  - A) Classical probability
  - B) Relative frequency probability
  - C) Axiomatic probability
  - D) Experimental error model

2. If two events A and B are mutually exclusive, what is the value of  $P(A \cap B)$ ?
  - A)  $P(A) + P(B)$
  - B)  $P(A) \times P(B)$
  - C) 1
  - D) 0
3. If the probability that a student passes in Maths is 0.7 and in English is 0.6, and the probability that the student passes in both subjects is 0.4, then what is the probability that the student passes in at least one subject?
  - A) 1.3
  - B) 0.7
  - C) 0.9
  - D) 1.0
4. Which of the following is the correct expression for conditional probability?
  - A)  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
  - B)  $P(A | B) = P(A \cap B) \div P(B)$
  - C)  $P(A \cap B) = P(A) + P(B)$
  - D)  $P(B | A) = P(B) \div P(A \cup B)$
5. Bayes' Theorem is mainly used to:
  - A) Measure the number of outcomes in a sample space
  - B) Predict future outcomes in completely random experiments
  - C) Calculate probability without any known prior information
  - D) Update a prior probability based on new evidence

### 3.8 Summary

- ❖ This unit introduced learners to the foundational concepts of **probability**, which is the study of measuring uncertainty and predicting the likelihood of events. It began by explaining the **meaning, importance, and applications** of probability in various fields, from statistics to healthcare and business.
- ❖ Key terms such as **mutually exclusive, exhaustive, complementary, independent, and dependent events** were defined to describe relationships between events. The three major definitions of probability—**classical,**

**relative frequency, and axiomatic**—provided different approaches to calculating and interpreting probability depending on the situation.

- ❖ The unit explored the **addition theorem** and its applications in calculating the probability of the union of two or more events. Next, it introduced **conditional probability**, which refines probability estimates when partial or prior information is available. The **multiplicative rule for independent events** and **Bayes' Theorem** demonstrated how to deal with joint and revised probabilities, making probability a powerful tool in prediction and decision-making.
- ❖ Through examples and formulas, learners are equipped to solve real-world problems involving uncertainty using systematic probability-based approaches.

### 3.9 Key Terms

1. **Experiment** - An action or process that leads to an outcome.
2. **Sample Space (S)** - The set of all possible outcomes of an experiment.
3. **Event (E)** - A subset of the sample space; the outcome(s) of interest.
4. **Mutually Exclusive Events** - Events that cannot happen at the same time.
5. **Exhaustive Events** - A set of events that cover all possible outcomes.
6. **Complementary Events** - Events where one is the "not" of the other.
7. **Independent Events** - Events where the outcome of one does not affect the other.
8. **Conditional Probability** - Probability of an event given that another has occurred.
9. **Classical Probability** - Probability based on equally likely outcomes.
10. **Relative Frequency** - Probability based on experimental data or observations.
11. **Axiomatic Probability** - Probability defined by formal rules (axioms).
12. **Addition Theorem** - Rule to calculate the probability of union of events.

### 3.10 Descriptive Questions

1. Define probability and explain its importance in real-life decision-making.
2. What are mutually exclusive and exhaustive events? Give examples.
3. Explain the difference between classical, relative frequency, and axiomatic definitions of probability.
4. State and explain the addition theorem of probability.
5. Define conditional probability and give an example from real life.

6. What is the multiplicative theorem for independent events? How is it applied?
7. State Bayes' Theorem and explain how it helps in updating probabilities.
8. A card is drawn from a deck. What is the probability that it is a red card or a king?
9. A box contains 3 red balls and 2 blue balls. If one ball is drawn and not replaced, what is the probability that the second ball is also red?
10. Differentiate between independent and dependent events with examples.

### 3.11 References

1. Gupta, S.P. (2014). *Statistical Methods*. Sultan Chand & Sons.
2. Goon, A.M., Gupta, M.K., & Dasgupta, B. (2013). *Fundamentals of Statistics*. World Press.
3. Hogg, R.V. & Tanis, E.A. (2010). *Probability and Statistical Inference*. Pearson Education.
4. NCERT (2021). *Mathematics: Probability (Class XI)*.
5. Khan Academy. (2023). *Probability and Statistics*. [<https://www.khanacademy.org>]
6. Statistics How To. (2023). *Bayes' Theorem Explained*. [<https://www.statisticshowto.com>]

### Answers to Knowledge Check

#### *Knowledge Check 1*

1. B) Relative frequency probability
2. D) 0
3. C)  $0.9 \rightarrow P(A \cup B) = 0.7 + 0.6 - 0.4$
4. B)  $P(A | B) = P(A \cap B) \div P(B)$
5. D) Update a prior probability based on new evidence

## 3.12 Case Study

### Bayes at Work: Diagnosing with Data

#### Background:

A private hospital in Mumbai introduced a **screening test** for a rare disease affecting 2 in every 1000 people. The test is **95% accurate**, meaning it correctly identifies 95% of those who have the disease (true positives) and incorrectly flags 5% of those who don't (false positives). A patient named Arjun tested positive. The doctor wanted to know: "What is the actual probability that Arjun has the disease?"

#### Problem:

The test result alone isn't enough because the disease is rare. There's a **high chance of false positives** due to the base rate of the disease being so low. The doctor must use **Bayes' Theorem** to determine the probability that Arjun truly has the disease, given the positive result.

#### Solution Using Bayes' Theorem:

Let:

- $D$  = has disease  $\rightarrow P(D) = 0.002$
- $\neg D$  = does not have disease  $\rightarrow P(\neg D) = 0.998$
- $T^+ | D$  = test positive given disease  $\rightarrow P(T^+ | D) = 0.95$
- $T^+ | \neg D$  = test positive given no disease  $\rightarrow P(T^+ | \neg D) = 0.05$

#### Apply Bayes' Theorem:

$$\begin{aligned}P(D | T^+) &= [P(D) \times P(T^+ | D)] \div [P(D) \times P(T^+ | D) + P(\neg D) \times P(T^+ | \neg D)] \\&= [0.002 \times 0.95] \div [0.002 \times 0.95 + 0.998 \times 0.05] \\&= 0.0019 \div (0.0019 + 0.0499) \\&\approx 0.0019 \div 0.0518 \\&\approx \mathbf{0.0367} \text{ or } \mathbf{3.67\%}\end{aligned}$$

#### Interpretation:

Even though Arjun tested positive, the chance he actually has the disease is only **3.67%**, due to the rarity of the disease and the rate of false positives.

#### Outcome:

The doctor avoids causing panic and advises further confirmatory testing. The hospital also reviews the screening program to improve decision-making using data-informed approaches.

## Unit 4: Random Variables

### Learning Objectives

1. Understand the concept of a random variable, distinguish between discrete and continuous random variables, and explain their roles in modeling real-world outcomes.
2. Define and interpret the Probability Mass Function (PMF) of a discrete random variable, and use it to calculate probabilities associated with specific values or ranges.
3. Understand and construct a Cumulative Distribution Function (CDF) for a discrete random variable, and explain how it reflects the accumulation of probabilities.
4. Apply rules of probability to solve problems involving discrete random variables, including expected value (mean), variance, and standard deviation.
5. Explore the concept of two-dimensional (joint) discrete random variables, and calculate joint, marginal, and conditional probabilities from a joint distribution table.
6. Compare and differentiate between PMF and CDF, and analyse how each is used in understanding distribution and cumulative behavior of discrete data.
7. Use discrete probability models to solve real-life problems in fields such as business, engineering, health sciences, and social sciences, making informed decisions under uncertainty.

### Content

- 4.0 Introductory Caselet
- 4.1 Introduction
- 4.2 Random Variables
- 4.3 Probability Mass Function (PMF)
- 4.4 Cumulative Distribution Function (CDF)
- 4.5 Two-Dimensional Discrete Random Variables
- 4.6 Summary
- 4.7 Key Terms
- 4.8 Descriptive Questions
- 4.9 References
- 4.10 Case Study

## 4.0 Introductory Caselet

### “Simran’s Sales Forecast: Making Sense of Daily Orders”

#### Background:

Simran, a young entrepreneur running an online handmade crafts store, noticed that her daily orders varied widely. Some days she would receive no orders; on other days, she’d get more than she could handle alone. She wanted to manage her inventory and time more efficiently, so she turned to probability for help.

A mentor introduced her to the concept of **discrete random variables**. Together, they defined **X** as the number of orders Simran receives in a day. Over a 30-day period, Simran recorded her order counts and calculated the **probability mass function (PMF)**. She discovered that the most frequent order volume was 3 per day.

Using this PMF, she built a **cumulative distribution function (CDF)** to estimate the probability of getting 2 or fewer orders, helping her plan low-volume days efficiently. Later, she extended her analysis by recording the number of **repeat customers (Y)** on the same days. She constructed a **joint probability table** for (X, Y) and used **marginal** and **conditional distributions** to understand how repeat customer patterns impacted daily sales.

Thanks to probability tools like PMF, CDF, and joint distributions, Simran could now:

- Predict inventory needs,
- Identify peak order days,
- And prioritize repeat customers for promotions.

What was once guesswork became **data-driven planning**, all because she learned to model her sales as a discrete random variable.

#### Critical Thinking Question:

If you were in Simran’s position, how would you use a joint probability distribution to improve marketing decisions? Give an example of two variables you would track and how you would analyse them.

## 4.1 Introduction

In probability and statistics, we often deal with uncertainty and unpredictable outcomes. To analyse such outcomes systematically, we introduce a concept called the **random variable**. It acts as a bridge between **outcomes of an experiment** and **numerical values** that can be analysed mathematically.

### 4.1.1 Concept of Random Variables in Probability Theory

A **random variable** is a **numerical outcome** of a random experiment. It assigns a real number to each outcome in the sample space.

There are two main types of random variables:

- **Discrete Random Variable:** Takes on **countable** values (e.g., 0, 1, 2, ...).
- **Continuous Random Variable:** Takes on values in a **continuous range** (e.g., any real number between 0 and 1).

#### **Definition:**

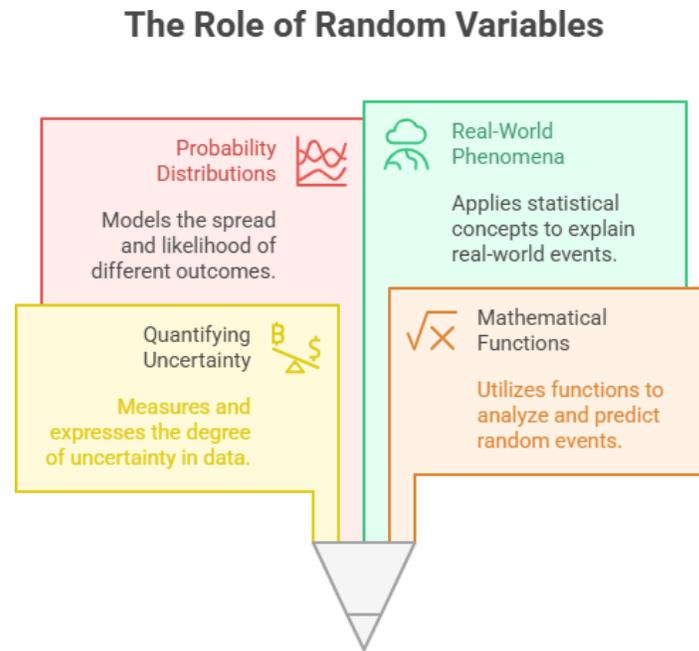
A random variable is a function that assigns a numerical value to each outcome of a random experiment.

#### **Example:**

- Tossing a coin → Let  $X = 1$  if Head,  $X = 0$  if Tail
- Rolling a die →  $X =$  outcome on the die (1 to 6)

The randomness comes from the experiment, but once an outcome occurs, the value of the random variable is fixed.

## 4.1.2 Importance of Random Variables in Statistics



**Fig.4.1. Importance of Random Variables in Statistics**

Random variables play a **central role in statistical analysis** because they:

- Allow us to **quantify uncertainty**.
- Make it possible to use **mathematical functions** (like mean, variance, probability functions) to study randomness.
- Are the basis for **probability distributions**, which model how outcomes are spread out.
- Help describe real-world phenomena in terms of **probabilities and numerical values**.

### **Applications in Statistics:**

- Estimating population averages (using expected value).
- Analyzing risk in finance (using variance and standard deviation).
- Modeling processes in engineering and medicine (e.g., failure rates, waiting times).

## 4.1.3 Examples of Random Variables in Real Life

Here are a few practical situations where random variables are used:

Scenario	Random Variable (X)	Type
Number of goals scored in a football match	$X = \text{Number of goals}$	Discrete
Time taken for a customer to be served	$X = \text{Time in minutes}$	Continuous
Number of defective items in a batch	$X = \text{Count of defective items}$	Discrete
Daily rainfall in a city	$X = \text{Rainfall in mm}$	Continuous
Result of a multiple-choice quiz	$X = \text{Number of correct answers}$	Discrete
Number of heads in 3 coin tosses	$X = 0, 1, 2, \text{ or } 3$	Discrete

**Key Idea:**

Each random variable is linked to an experiment or situation where the outcome is uncertain, but we can assign **numerical values** to those outcomes and **study their behavior** statistically.

**4.2 Random Variables**

A **random variable** is a core concept in probability and statistics that allows us to represent outcomes of random processes as **numerical values**. These variables help quantify uncertainty and form the basis for probability distributions, expectations, and statistical inference.

**4.2.1 Definition and Classification (Discrete and Continuous)**

**Definition:**

A **random variable** is a **numerical function** that assigns a real number to each outcome of a random experiment. It's called "random" because the value it takes on is determined by chance.

Let:

- $X$  be a random variable
- $S$  be the sample space

Then  $X: S \rightarrow \mathbb{R}$ , meaning  **$X$  maps outcomes to real numbers**.

## Classification of Random Variables:

### Classification of Random Variables

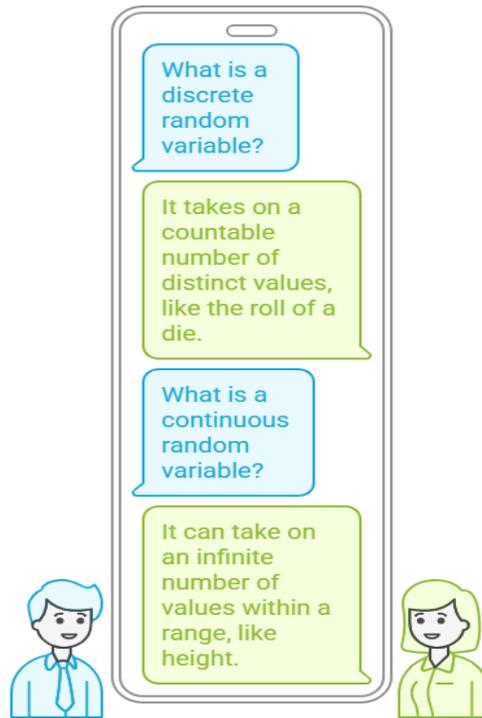


Fig.4.2. Classification of Random Variables

Random variables are classified into two main types:

#### A. Discrete Random Variable

- Takes **finite or countable** number of distinct values.
- Often associated with **counting**.
- Examples:
  - Number of heads in 3 coin tosses  $\rightarrow X \in \{0, 1, 2, 3\}$
  - Number of cars in a parking lot
  - Score on a 10-mark test

#### B. Continuous Random Variable

- Takes **infinite, uncountable** values within an interval.
- Usually associated with **measuring**.
- Examples:
  - Time taken to finish a race ( $X \in [0, \infty)$ )
  - Height of students in a class
  - Temperature of a city at noon

## 4.2.2 Properties of Random Variables

Random variables have certain mathematical properties that make them useful in statistics and probability modeling:

### 1. Probability Assignment

For a **discrete random variable X**, each value  $x_i$  has a probability  $P(X = x_i)$  such that:

- $0 \leq P(X = x_i) \leq 1$
- $\sum P(X = x_i) = 1$  (sum over all possible values)

### 2. Expected Value (Mean)

The **expected value** of a random variable gives the **long-run average** outcome.

For discrete X:

$$E(X) = \sum x_i \times P(X = x_i)$$

This is like a **weighted average** of possible values.

### 3. Variance and Standard Deviation

- **Variance (Var(X))** measures how spread out the values are from the mean.

$$\text{Var}(X) = E[(X - \mu)^2] = \sum (x_i - \mu)^2 \times P(X = x_i)$$

- **Standard Deviation** is the square root of variance:
  - $\sigma = \sqrt{\text{Var}(X)}$

These tell us how **reliable or variable** the outcomes of the random process are.

### 4. Linearity of Expectation

If  $X$  and  $Y$  are two random variables, and  $a, b$  are constants:

$$E(aX + bY) = aE(X) + bE(Y)$$

This property is **always true**, even if  $X$  and  $Y$  are dependent.

## 5. Indicator Random Variables

A special type of random variable that takes only values 0 or 1.

Used to represent whether a specific event occurred.

### Example:

Let  $A$  = event student passes.

Define:

**$X = 1$  if student passes;  $X = 0$  if not**

Then  $E(X) = P(\text{student passes})$

## 4.3 Probability Mass Function (PMF)

A **Probability Mass Function (PMF)** is a mathematical tool used to describe the **distribution of a discrete random variable**. It tells us the **probability that a random variable takes on a specific value**.

### 4.3.1 Definition of PMF

The **Probability Mass Function (PMF)** of a discrete random variable  $X$  is a function that gives the **probability of each possible value** that  $X$  can take.

Formally, if  $X$  is a discrete random variable, then its PMF is defined as:

$$P(X = x_i) = p(x_i)$$

Where:

- $x_i$  is a value that  $X$  can take
- $p(x_i)$  is the probability that  $X$  equals  $x_i$
- $p(x_i) \geq 0$ , for all  $i$
- The sum of all  $p(x_i)$  over the sample space is 1

### 4.3.2 Properties of PMF

A valid PMF must satisfy the following properties:

1. **Non-negativity:**

For every value  $x_i$ ,

$$p(x_i) \geq 0$$

2. **Normalization (Total Probability = 1):**

The sum of the probabilities of all possible values must be equal to 1:

$$\sum p(x_i) = 1$$

3. **Defined Only for Discrete Values:**

The PMF is only defined for the values that the discrete random variable can actually take. For other values, the probability is zero.

**Probability of an Exact Value:**

Since PMFs are for discrete variables, you can directly compute:

$$P(X = x) = p(x)$$

### Did You Know?

“Did you know that the **Probability Mass Function (PMF)** can never assign a probability greater than **1** or less than **0** to any event? This rule applies strictly because PMF represents **actual probabilities**, not just relative frequencies. Also, the **sum of all PMF values for a discrete random variable must always equal 1**, even if the variable has many possible values—this ensures that all possible outcomes are accounted for.”

### 4.3.3 Examples of PMFs

#### Example 1: Tossing a Fair Coin

Let  $X$  = number of heads when one coin is tossed.

Possible values:

$$X \in \{0, 1\}$$

PMF:

- $P(X = 0) = 0.5$
- $P(X = 1) = 0.5$

This satisfies:

- $p(x) \geq 0$
- $\sum p(x) = 1$

### Example 2: Rolling a Fair Die

Let  $X$  = outcome when one fair die is rolled.

Possible values:  $X \in \{1, 2, 3, 4, 5, 6\}$

PMF:

- $P(X = 1) = 1/6$
- $P(X = 2) = 1/6$
- $P(X = 3) = 1/6$
- $P(X = 4) = 1/6$
- $P(X = 5) = 1/6$
- $P(X = 6) = 1/6$

Check:

- All probabilities are  $\geq 0$
- Sum =  $6 \times (1/6) = 1$

### Example 3: Custom PMF

Let  $X$  = number of calls received at a customer service desk in one hour.

Assume probabilities are as follows:

<b>x (Calls)</b>	<b>0</b>	<b>1</b>	<b>2</b>	<b>3</b>
<b>p(x)</b>	0.1	0.3	0.4	0.2

Check:

- All values are  $\geq 0$
- $0.1 + 0.3 + 0.4 + 0.2 = 1$

Thus, this is a valid PMF.

### Visual Representation (Optional in class):

PMFs can be represented as **bar charts**, where:

- X-axis = values of the random variable
- Y-axis = corresponding probabilities  $p(x)$

### “Activity: Constructing a PMF and CDF Table”

**Title:** *Daily Orders – PMF and CDF Practice*

**Instruction to Student:**

Your team manages a small e-commerce store. Over 10 days, you recorded the following number of orders:

{2, 3, 3, 1, 4, 2, 2, 3, 1, 2}

1. List all unique values of the random variable  $X$  (number of orders).
2. Compute the **frequency** and **relative frequency** (PMF) for each value.
3. Using the PMF, calculate the **Cumulative Distribution Function (CDF)** for each value of  $X$ .
4. Present your final table with three columns:  $X$ ,  $P(X)$ , and  $F(X)$ .
5. Submit a short reflection: What is the probability of receiving **3 or fewer orders** on a given day?

## 4.4 Cumulative Distribution Function (CDF)

A **Cumulative Distribution Function (CDF)** gives the **cumulative probability** that a random variable takes a value **less than or equal to a given number**. It is a powerful tool to summarise the behavior of a random variable over a range of values, rather than at one point.

### 4.4.1 Concept and Definition

The **CDF of a random variable  $X$** , denoted by  **$F(x)$** , is the **probability that  $X$  takes a value less than or equal to  $x$** .

**Mathematically:**

$$F(x) = P(X \leq x)$$

This means  $F(x)$  is the sum of the probabilities for all values of  $X$  that are less than or equal to  $x$ .

### 4.4.2 CDF for Discrete Random Variables

For a **discrete random variable**, the CDF is calculated by **summing the PMF values** up to a given point.

If  $X$  is a discrete random variable with possible values  $x_1, x_2, x_3, \dots, x_n$  and PMF  $p(x)$ , then:

$$F(x) = \sum p(x_i), \text{ for all } x_i \leq x$$

The CDF is:

- **Stepwise increasing**

- Always between **0 and 1**
- Right-continuous (value jumps at each defined x)

**Example:**

Let X be the number of heads in two coin tosses. Possible values:  $X \in \{0, 1, 2\}$

**x P(X = x)**

0 0.25

1 0.50

2 0.25

Then the CDF F(x) is:

- $F(0) = P(X \leq 0) = 0.25$
- $F(1) = P(X \leq 1) = 0.25 + 0.50 = 0.75$
- $F(2) = P(X \leq 2) = 0.25 + 0.50 + 0.25 = 1.0$

For values between these points:

- $F(x) = 0$  when  $x < 0$
- $F(x) = 0.25$  when  $0 \leq x < 1$
- $F(x) = 0.75$  when  $1 \leq x < 2$
- $F(x) = 1$  when  $x \geq 2$

**4.4.3 Relationship between PMF and CDF**

The **PMF and CDF** are closely connected:

- The **PMF** shows the probability of a specific value:  
 $p(x) = P(X = x)$
- The **CDF** shows the cumulative probability up to a value:  
 $F(x) = P(X \leq x)$

To go from **PMF to CDF**, you **add up** the PMF values.

To go from **CDF to PMF**, you **subtract**:

$$p(x) = F(x) - F(x-1)$$

This allows you to **reconstruct one from the other** if needed.

#### 4.4.4 Graphical Representation of Discrete Distributions

For a **discrete CDF**, the graph is a **step function**.

**Key characteristics of the graph:**

- X-axis = values of the random variable
- Y-axis = cumulative probabilities  $F(x)$
- The graph **increases in steps** at each value of  $X$
- The height of the step = probability from the PMF
- It is **flat** between steps and **jumps** at the  $x$ -values

**Example:**

Using the coin toss example:

- Plot points at  $(0, 0.25)$ ,  $(1, 0.75)$ , and  $(2, 1.0)$
- Draw horizontal lines between the steps

This graph helps quickly visualise:

- How much probability is **accumulated** up to each value
- Where the **most likely values** lie
- Whether the distribution is **skewed or centered**

#### Did You Know?

“Did you know that a **Cumulative Distribution Function (CDF)** graph for a discrete variable looks like a **step function**? Unlike continuous functions which flow smoothly, CDF graphs for discrete variables “jump” at each possible value of the variable. The size of each step corresponds to the probability assigned by the PMF, and the final step always reaches 1 (or 100%)”

#### 4.5 Two-Dimensional Discrete Random Variables

In many real-world situations, we are interested in **two or more random variables at once**. For example, in a family, we might want to study the number of boys and girls; in business, sales and marketing expenses; in a survey, education level and employment status.

When we consider **two random variables simultaneously**, we deal with **two-dimensional (bivariate) discrete distributions**.

Let  $X$  and  $Y$  be two discrete random variables. Their combined behavior is studied using a **joint probability distribution**.

#### 4.5.1 Joint Probability Mass Function (Joint PMF)

The **Joint PMF** gives the probability that  $X = x_i$  and  $Y = y_j$  occur **together**.

It is denoted as:

$$P(X = x_i, Y = y_j) = p(x_i, y_j)$$

The joint PMF satisfies the following:

- $p(x_i, y_j) \geq 0$  for all  $i, j$
- The sum of all joint probabilities is **1**:

$$\sum \sum p(x_i, y_j) = 1$$

These values are typically arranged in a **joint probability table** with rows and columns representing values of  $X$  and  $Y$ .

#### 4.5.2 Marginal Distributions

The **marginal distribution** of a variable is its **individual distribution**, derived from the joint distribution by **summing over the other variable**.

- **Marginal PMF of X:**

$$P(X = x_i) = \sum p(x_i, y_j) \text{ (sum over all } y_j)$$

- **Marginal PMF of Y:**

$$P(Y = y_j) = \sum p(x_i, y_j) \text{ (sum over all } x_i)$$

These show how **each variable behaves on its own**, regardless of the value of the other Variable.

### “Activity: Joint and Marginal Distribution Analysis”

**Title:** *Tracking Sales: Bread vs Milk*

**Instruction to Student:**

You are given the following joint probability table showing the number of customers who bought **bread (X)** and **milk (Y)** on a given day:

<b>X\Y</b>	<b>0</b>	<b>1</b>	<b>2</b>
<b>0</b>	0.05	0.10	0.05
<b>1</b>	0.10	0.20	0.10
<b>2</b>	0.05	0.15	0.20

1. Compute the **marginal distribution** of X and Y.
2. What is the probability that a customer bought at least 1 unit of **either product**?
3. Calculate  $P(X = 2 | Y = 2)$  — the probability that a customer bought 2 breads **given** they bought 2 milks.
4. Submit your calculations along with a brief interpretation of what this data suggests about buying behavior.

### 4.5.3 Conditional Distributions

The **conditional distribution** of one variable **given** the other represents the updated probability **after partial information is known**.

- $P(X = x_i | Y = y_j) = p(x_i, y_j) \div P(Y = y_j)$
- $P(Y = y_j | X = x_i) = p(x_i, y_j) \div P(X = x_i)$

It allows us to study how **one variable behaves under a fixed value of the other**.

This is especially useful when events are **not independent**.

### 4.5.4 Examples on Marginal and Conditional Distributions

**Example: Joint Probability Table**

Let X = Number of books sold in a day (0, 1)

Let Y = Number of customer complaints (0, 1)

	<b>Y = 0</b>	<b>Y = 1</b>	<b>Total</b>
<b>X = 0</b>	0.1	0.2	0.3
<b>X = 1</b>	0.4	0.3	0.7

<b>Total</b>	0.5	0.5	1.0
--------------	-----	-----	-----

**Marginal PMFs:**

- $P(X = 0) = 0.1 + 0.2 = 0.3$
- $P(X = 1) = 0.4 + 0.3 = 0.7$
- $P(Y = 0) = 0.1 + 0.4 = 0.5$
- $P(Y = 1) = 0.2 + 0.3 = 0.5$

**Conditional PMFs:**

- $P(X = 0 | Y = 1) = 0.2 \div 0.5 = 0.4$
- $P(X = 1 | Y = 1) = 0.3 \div 0.5 = 0.6$
- $P(Y = 0 | X = 1) = 0.4 \div 0.7 \approx 0.571$
- $P(Y = 1 | X = 1) = 0.3 \div 0.7 \approx 0.429$

**Knowledge Check 1**

**Choose the correct option:**

1. What is the essential condition for a valid **Probability Mass Function (PMF)**?
  - A) It must have at least 10 values
  - B) All values must be negative
  - C) The sum of all probabilities must be 1
  - D) The cumulative probability must be 0
2. Which of the following statements is **true** about a **Cumulative Distribution Function (CDF)**?
  - A) CDF decreases as the value of the random variable increases
  - B) CDF always remains constant
  - C) CDF is always equal to PMF
  - D) CDF is a non-decreasing function that reaches 1
3. A discrete random variable X has the following PMF:  
 $P(0) = 0.2, P(1) = 0.5, P(2) = 0.3$

What is the  $P(X \leq 1)$ ?

- A) 0.5
- B) 0.7
- C) 0.9
- D) 1.0

4. In a **joint probability distribution**, the **marginal probability** of variable  $X$  is obtained by:
- A) Dividing joint probabilities by totals
  - B) Multiplying all joint probabilities
  - C) Subtracting  $Y$ 's probabilities from  $X$ 's
  - D) Summing over all values of  $Y$  for each value of  $X$
5. What does the value  $P(X = 2 \mid Y = 3)$  represent?
- A) Probability of both  $X$  and  $Y$  being 2
  - B) Probability that  $Y$  equals 3, regardless of  $X$
  - C) Probability that  $X$  is 2, given that  $Y$  is 3
  - D) Total probability of  $X$  and  $Y$  being less than 5

## 4.6 Summary

- ❖ This unit introduced the key concepts of **discrete probability distributions**—a critical part of understanding random behavior in statistics. It began with the **concept of random variables**, followed by their **classification** into discrete and continuous types.
- ❖ We then studied the **Probability Mass Function (PMF)**, which describes the probability of each possible value of a **discrete random variable**. Building on this, the **Cumulative Distribution Function (CDF)** was introduced to track the cumulative probability up to a certain value.
- ❖ The unit then extended into **two-dimensional (bivariate) discrete random variables**, including how to construct and interpret **joint, marginal, and conditional distributions**. Through examples and step-by-step tables, learners gained the ability to explore relationships between two variables at once.
- ❖ These foundational tools allow statisticians, data scientists, and researchers to model uncertainty, analyze patterns, and make data-driven decisions in real-world applications.

## 4.7 Key Terms

1. **Random Variable (RV)** - A variable that takes numerical values based on outcomes of a random experiment.
2. **Discrete Random Variable** - A random variable that takes a finite or countable number of values.
3. **PMF (Probability Mass Function)** - A function that gives the probability for each value of a discrete random variable.
4. **CDF (Cumulative Distribution Function)** - A function that gives the cumulative probability up to a value  $x$ , i.e.,  $P(X \leq x)$ .
5. **Joint PMF** - Probability distribution of two discrete random variables occurring together.
6. **Marginal Distribution** - Probability distribution of one variable obtained by summing the joint probabilities over the other variable.
7. **Conditional Distribution** - Probability of one variable given that another has a specific value.

#### 4.8 Descriptive Questions

1. What is a discrete random variable? Give an example.
2. Define the Probability Mass Function (PMF) and list its properties.
3. Explain how the CDF is calculated from the PMF.
4. Differentiate between discrete and continuous random variables.
5. How do you construct a joint probability distribution table?
6. What is the difference between marginal and conditional distributions?
7. Explain the relationship between PMF and CDF using an example.
8. A die is rolled once. Define a suitable random variable and write its PMF.
9. In a classroom survey, students were asked about time spent on homework and number of subjects taken. How would you model this using two-dimensional random variables?
10. Why must the sum of all PMF values equal 1?

#### 4.9 References

1. Hogg, R.V. & Tanis, E.A. (2015). *Probability and Statistical Inference*. Pearson.
2. Sheldon Ross (2014). *Introduction to Probability Models*. Academic Press.
3. Walpole, R.E., Myers, R.H., et al. (2012). *Probability and Statistics for Engineers and Scientists*. Pearson.
4. Goon, A.M., Gupta, M.K., & Dasgupta, B. (2013). *Fundamentals of Statistics, Vol I*.
5. NCERT (2021). *Statistics Textbook, Class XI and XII*

6. Khan Academy. (2023). *Random Variables and Probability Distributions*. [<https://www.khanacademy.org>]

### Answers to Knowledge Check

#### ***Knowledge Check 1***

1. C) The sum of all probabilities must be 1
2. D) CDF is a non-decreasing function that reaches 1
3. C) 0.9 (i.e.,  $P(0) + P(1) = 0.2 + 0.5$ )
4. D) Summing over all values of Y for each value of X
5. C) Probability that X is 2, given that Y is 3

## 4.10 Case Study

### “Using Discrete Probability Distributions in Inventory Management”

#### Introduction

In today’s competitive retail environment, businesses rely heavily on data to drive operational efficiency. One key challenge faced by retail managers is **stock management**—deciding how much inventory to maintain daily to avoid shortages or overstocking. This decision is often based on sales patterns that fluctuate due to demand uncertainty. Understanding and applying **discrete probability distributions** can help managers forecast better and reduce costs.

This caselet explores how a retail store manager uses **random variables, probability mass functions (PMF), and cumulative distribution functions (CDF)** to optimize her stock planning. It demonstrates the practical use of two-dimensional discrete variables and how **joint and conditional probabilities** can lead to smarter inventory strategies.

#### Background

Anjali, the operations manager of a mid-sized convenience store in Nagpur, noticed inconsistencies in daily sales of certain perishable items like bread and milk. On some days, the store would run out of stock by evening, while on other days, leftover items had to be discarded. This imbalance was hurting both sales and profitability.

To address this, Anjali began recording the **number of bread packets sold daily** (denoted as random variable  $X$ ) over a period of one month. She classified the outcomes and calculated the **PMF**, noting that certain quantities (e.g., 4 or 5 packets) occurred most frequently. She then constructed a **CDF** to determine the probability of selling a maximum of ‘ $k$ ’ packets on any given day.

As her analysis matured, Anjali introduced a second random variable  $Y$ , representing the number of milk bottles sold. She created a **joint probability table** of  $X$  and  $Y$  to study their relationship and calculated **marginal and conditional probabilities**. This allowed her to answer questions like:

*“If 4 bread packets were sold, what is the likelihood that milk sales were also high?”*

#### Problem Statement 1: Poor Forecasting of Demand for Bread

Due to unpredictable daily demand, Anjali either overstocked or understocked bread, leading to wastage or lost sales.

**Solution:**

Anjali defined a discrete random variable  $X = \text{number of bread packets sold per day}$ . Using recorded data, she built a **PMF** and then a **CDF** to find cumulative probabilities such as  $P(X \leq 3)$  or  $P(X \geq 5)$ . This helped her decide reorder points and safety stock levels more efficiently.

**MCQ 1:**

What is the role of a cumulative distribution function (CDF) in Anjali's stock planning?

- A) To calculate total profit from bread sales
- B) To estimate the probability of exact demand
- C) To find the cumulative probability of selling up to a certain number of packets
- D) To measure average shelf life

**Answer:** C) To find the cumulative probability of selling up to a certain number of packets

**Problem Statement 2: Inability to Link Sales of Bread and Milk**

Anjali suspected a link between bread and milk sales but had no structured method to analyse them together.

**Solution:**

She introduced a second variable,  $Y = \text{number of milk bottles sold}$ , and created a **joint probability distribution table** for  $(X, Y)$ . She then calculated:

- **Marginal probabilities** (e.g.,  $P(X = 3)$ )
- **Conditional probabilities** (e.g.,  $P(Y = 2 | X = 3)$ )

This analysis showed that high bread sales often coincided with moderate milk sales, which helped in **coordinated inventory planning**.

**MCQ 2:**

What is the purpose of using a joint probability distribution?

- A) To calculate only average sales
- B) To study sales of unrelated products
- C) To analyse relationships between two variables like bread and milk sales
- D) To randomly allocate inventory

**Answer:** C) To analyse relationships between two variables like bread and milk sales

### **Problem Statement 3: Lack of Predictive Insight from Daily Sales Data**

Before her analysis, Anjali treated daily sales data as isolated numbers rather than part of a statistical pattern.

#### **Solution:**

She treated sales as **discrete random variables**, enabling her to apply **statistical functions** like PMF, mean, and variance. She also visualized these distributions to understand which sales outcomes were most likely and how spread out they were. This allowed her to shift from reactive to **predictive inventory control**.

#### **MCQ 3:**

What is the benefit of treating sales counts as discrete random variables?

- A) It helps randomize supply
- B) It avoids the need for data
- C) It enables probability-based analysis of outcomes
- D) It simplifies purchase orders

**Answer:** C) It enables probability-based analysis of outcomes

#### **Conclusion**

By using **discrete probability distributions**, Anjali transformed her inventory planning process from guesswork to **data-driven forecasting**. Through **PMF, CDF, and joint distribution analysis**, she could optimise her stocking decisions, reduce waste, and respond better to customer demand. This case illustrates how **basic statistical tools**, when applied thoughtfully, can improve operational efficiency and support smarter business decisions in everyday retail management.

## Unit 5: Measures of Central Tendency

### Learning Objectives

1. **Define and distinguish** between different measures of average, including arithmetic mean, median, and mode.
2. **Calculate the arithmetic mean** for both ungrouped and grouped data sets.
3. **Determine the median and mode** for various data types and interpret their significance in real-world contexts.
4. **Compare and contrast** the characteristics, strengths, and limitations of arithmetic mean, median, and mode.
5. **Apply empirical formulas** to estimate the mode and understand the relationship between the three central tendency measures.
6. **Interpret and analyze** statistical data sets using appropriate measures of central tendency for decision-making.
7. **Evaluate case studies** to identify the most suitable average measure and justify its application with evidence and reasoning.

### Content

- 5.0 Introductory Caselet
- 5.1 Measures of Average
- 5.2 Arithmetic Mean
- 5.3 Positional Averages: Median and Mode
- 5.4 Empirical Analysis of Central Tendency
- 5.5 Summary
- 5.6 Key Terms
- 5.7 Descriptive Questions
- 5.8 References
- 5.9 Case Study

## 5.0 Introductory Caselet

### “Rahul’s Rental Records: Finding the Typical Customer”

#### Background:

Rahul manages a mid-sized car rental agency in a metropolitan city. Over the past year, he noticed significant fluctuations in the number of days cars were rented, depending on seasons, customer profiles, and ongoing promotions. To optimize pricing, fleet availability, and customer targeting, Rahul realized he needed to identify what constituted a "typical" rental.

He began recording the number of rental days for each customer. After collecting data from 150 rentals, Rahul observed a wide variation—some customers rented for just one day, while others extended up to ten days or more. A data consultant recommended analyzing **measures of central tendency** to better understand rental behavior. Together, they calculated:

- The **arithmetic mean** (average), which indicated the general rental duration across all customers.
- The **median**, representing the midpoint of the dataset.
- The **mode**, showing the most frequently occurring rental length.

Rahul discovered that while the **mean** rental duration was **4.8 days**, the **median** was **4 days**, and the **mode** was **3 days**. This revealed that although some customers rented for long durations, the majority rented for shorter periods.

By interpreting these findings, Rahul was able to:

- Design pricing packages based on common rental durations,
- Create predictive models for vehicle availability,
- And develop promotional offers targeting different customer segments.

Understanding and applying **central tendency measures** transformed Rahul's decision-making from reactive to data-driven, improving both operational efficiency and customer satisfaction.

#### Critical Thinking Question:

If you were in Rahul’s position and noticed a large gap between the mean and the median rental duration, what would that suggest about the data distribution? How would this influence your business strategy in terms of pricing or promotions?

## 5.1 Measures of Average

In statistics, a **measure of average** is a value that represents or summarizes a set of data. It gives us a single number that describes the **center point** or **typical value** in a data set. Averages help simplify large amounts of data by providing a number that best represents the whole group.

There are different types of averages used in statistics, depending on the type of data and the purpose of analysis.

The most common measures of average are:

- **Arithmetic Mean** (usually referred to simply as the "mean")
- **Median**
- **Mode**

These measures are also known as **measures of central tendency** because they describe where the center of the data lies.

Understanding averages is useful in everyday life and in various fields such as economics, business, education, and social sciences. For example, the average temperature helps predict weather, the average income shows economic conditions, and the average score helps assess student performance.

### 5.1.1 Concept and Importance of Central Tendency

#### Concept of Central Tendency

The **central tendency** of a data set refers to the idea of finding the central or middle value around which other values are distributed. It gives a snapshot of a typical value in a group of numbers.

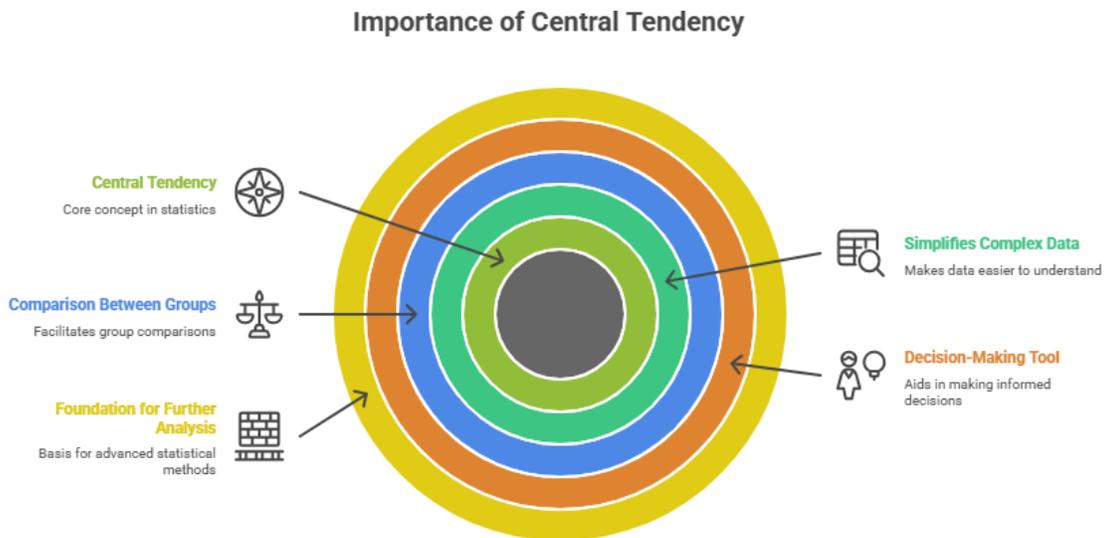
For example, if we collect the ages of 50 students in a class, it may not be useful to list all the ages every time. Instead, we might say the "average age" is 19 years. This single value gives us a good sense of what the general age of the students is.

There are three main ways to measure central tendency:

1. **Mean:** The sum of all values divided by the number of values.
2. **Median:** The middle value when the data is arranged in order.
3. **Mode:** The value that appears most frequently in the data.

Each of these measures gives a different perspective on what is considered "typical" in a data set.

## Importance of Central Tendency

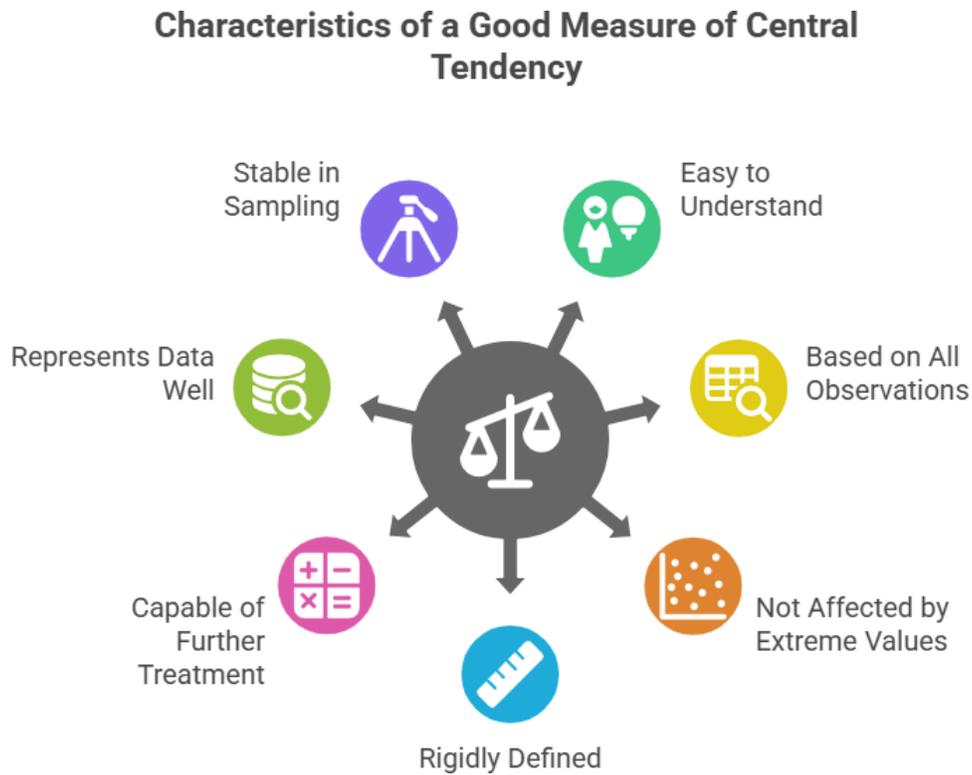


**Fig.5.1. Importance of Central Tendency**

1. **Simplifies Complex Data:** Instead of dealing with hundreds or thousands of individual numbers, a single average gives a summary of the whole set.
2. **Comparison Between Groups:** Central tendency helps compare two or more groups. For instance, comparing average incomes of two cities to determine which is more prosperous.
3. **Decision-Making Tool:** In business, averages help managers make decisions, such as determining average sales, profits, or customer preferences.
4. **Foundation for Further Statistical Analysis:** Measures like standard deviation and variance rely on the mean. Many statistical models use measures of central tendency as their base.
5. **Used in Every Field:** Whether it's a government analyzing population data, a teacher calculating average scores, or a doctor looking at average recovery times, central tendency is used universally.
6. **Identifying Patterns:** Central tendency helps recognize trends or patterns over time, such as average rainfall in a region, which can influence agriculture and planning.

Understanding central tendency is a fundamental step in working with data. It is the starting point for exploring and interpreting numerical information.

### 5.1.2 Characteristics of a Good Measure of Central Tendency



**Fig.5.2. Characteristics of a Good Measure of Central Tendency**

A good measure of central tendency should meet several important criteria to ensure it accurately and effectively represents a data set. These characteristics include:

**1. Easy to Understand and Calculate**

The method of calculation should be simple, and the result should be easy to interpret. This makes it more practical and useful in real-world applications.

**2. Based on All Observations**

A good measure should take into account every value in the data set. This ensures that no part of the data is ignored, making the measure more accurate.

**3. Not Affected Much by Extreme Values**

If a data set contains extremely high or low values (called outliers), a good measure should not be overly influenced by them. For example, the median is more resistant to extreme values than the mean.

#### **4. Rigidly Defined**

The method of calculating the measure should be clearly and precisely defined so that different people will arrive at the same result when applying it to the same data.

#### **5. Capable of Further Mathematical Treatment**

The measure should be useful for further statistical analysis. For instance, the mean can be used in formulas for standard deviation, correlation, and regression analysis.

#### **6. Should Represent the Data Set Well**

The chosen measure should lie within the range of the data and reflect the overall pattern or distribution of the data accurately.

#### **7. Should Be Stable in Sampling**

When multiple samples are taken from the same population, the value of a good central tendency measure should not vary significantly from sample to sample.

### **5.1.3 Types of Averages (Mean, Median, Mode, etc.)**

There are several types of averages, each serving a specific purpose depending on the nature of the data. The three most common types are the **mean**, **median**, and **mode**. Additional types include the **geometric mean** and **harmonic mean**.

#### **1. Arithmetic Mean (or simply "Mean")**

The mean is the most commonly used average. It is calculated by dividing the total of all values by the number of values.

##### **Formula (for ungrouped data):**

$$\text{Mean} = (\Sigma x) \div n$$

Where:

- $\Sigma x$  = sum of all values
- $n$  = number of values

**Example:**

If the values are: 10, 20, 30, 40, 50

Then, Mean =  $(10 + 20 + 30 + 40 + 50) \div 5 = 150 \div 5 = 30$

**Note:** The mean is affected by very high or low values (outliers).

**2. Median**

The median is the middle value of the data when arranged in ascending or descending order. It divides the data into two equal parts.

**Steps to calculate the median (for ungrouped data):**

- Arrange the data in order.
- If the number of values (n) is odd, the median is the middle value.
- If n is even, the median is the average of the two middle values.

**Example:**

Data: 15, 20, 25, 30, 35

Median = 25 (middle value)

Data: 12, 16, 20, 24

Median =  $(16 + 20) \div 2 = 36 \div 2 = 18$

**Note:** The median is not affected by extreme values.

**3. Mode**

The mode is the value that appears most frequently in the data set.

**Example:**

Data: 5, 8, 8, 10, 12

Mode = 8 (appears twice)

- A data set may have:
  - No mode (if all values are unique)
  - One mode (unimodal)
  - Two modes (bimodal)
  - More than two modes (multimodal)

**Note:** The mode is useful for categorical data and is not affected by extreme values.

#### 4. Geometric Mean

The geometric mean is used to calculate the average of values that are multiplied together or involve percentages and rates.

##### Formula (for n values):

$$\text{Geometric Mean (GM)} = \sqrt[n]{(x_1 \times x_2 \times x_3 \times \dots \times x_n)}$$

##### Example:

For values: 2, 4, 8

$$\text{GM} = \sqrt[3]{(2 \times 4 \times 8)} = \sqrt[3]{64} = 4$$

**Use Cases:** Compound interest, population growth, and growth rates.

#### 5. Harmonic Mean

The harmonic mean is used when data values are expressed as rates (e.g., speed, efficiency).

##### Formula (for n values):

$$\text{Harmonic Mean (HM)} = n \div (1/x_1 + 1/x_2 + 1/x_3 + \dots + 1/x_n)$$

##### Example:

If the values are: 2, 4, 6

$$\text{HM} = 3 \div (1/2 + 1/4 + 1/6)$$

$$= 3 \div (0.5 + 0.25 + 0.1667)$$

$$= 3 \div 0.9167 \approx 3.27$$

**Use Cases:** Average speed, pricing per unit, or efficiency rates.

### 5.2 Arithmetic Mean

The **arithmetic mean**, commonly referred to as the **mean**, is one of the most widely used measures of central tendency. It is calculated by dividing the **sum of all values** in a data set by the **number of values**. It gives a general idea of the "average" or "typical" value.

#### 5.2.1 Simple Arithmetic Mean (Ungrouped and Grouped Data)

##### A. Ungrouped Data

For ungrouped data, all individual values are listed without being classified into intervals.

**Formula:**

$$\text{Mean } (\bar{x}) = (\Sigma x) \div n$$

Where:

- $\Sigma x$  = sum of all observations
- $n$  = number of observations

**Example:**

If the marks obtained by five students are: 50, 60, 70, 80, and 90

$$\text{Mean} = (50 + 60 + 70 + 80 + 90) \div 5 = 350 \div 5 = 70$$

**B. Grouped Data**

Grouped data is presented in the form of a frequency distribution.

There are two main methods for calculating the mean of grouped data:

**1. Direct Method**

**Formula:**

$$\text{Mean } (\bar{x}) = (\Sigma f \times x) \div \Sigma f$$

Where:

- $f$  = frequency of each class
- $x$  = mid-point of each class
- $\Sigma f$  = total frequency

**Steps:**

1. Find the mid-point ( $x$ ) of each class:  $(\text{Lower limit} + \text{Upper limit}) \div 2$
2. Multiply each mid-point by the corresponding frequency
3. Add all the results to get  $\Sigma f \times x$
4. Divide by the total frequency ( $\Sigma f$ )

**Example:**

Class Interval	Frequency (f)	Mid-point (x)	f × x
10 – 20	3	15	45
20 – 30	5	25	125
30 – 40	7	35	245
40 – 50	5	45	225

<b>Total</b>	<b>20</b>		<b>640</b>
--------------	-----------	--	------------

$$\text{Mean} = 640 \div 20 = 32$$

### “Activity: Finding the Class Average”

You are provided with the following frequency distribution, representing the marks scored by students in a statistics examination. The class intervals are 0–10, 10–20, 20–30, 30–40, and 40–50, with corresponding frequencies of 2, 5, 8, 12, and 3. Your task is to calculate the mid-point for each class interval and then use the **direct method** to compute the **arithmetic mean** of the marks. After calculating the mean, write a short interpretation explaining what this average tells you about the performance of the class. Does the mean accurately reflect a “typical” score, or do you observe any signs of skewed distribution? Include both your calculation steps and a brief reflection in your submission.

### 5.2.2 Weighted Arithmetic Mean

The **weighted arithmetic mean** is used when different items in a data set carry different levels of importance or "weights".

#### Formula:

$$\text{Weighted Mean } (\bar{x}) = (\sum w \times x) \div \sum w$$

Where:

- $x$  = value
- $w$  = weight of each value
- $\sum w$  = total of weights

#### Example:

A student scored 80 in theory (weight = 3) and 90 in practical (weight = 2).

$$\text{Weighted Mean} = (80 \times 3 + 90 \times 2) \div (3 + 2) = (240 + 180) \div 5 = 420 \div 5 = 84$$

This method is useful in academic grading, economic indices, and situations where certain elements have more significance.

### Did You Know?

“The **weighted mean** is used in **stock market indices** like the **S&P 500**, where each company’s weight in the index is based on its market capitalization. This means larger companies (like Apple or Microsoft) have a bigger impact on the index value than smaller ones. Unlike the simple mean, the weighted mean reflects the real-world significance of each component.”

### 5.2.3 Properties and Applications of Arithmetic Mean

#### Properties:

1. **Algebraic Sum of Deviations from the Mean is Zero**

$$\Sigma(x - \bar{x}) = 0$$

This means the positive and negative deviations from the mean balance each other out.

2. **Mean is Affected by Every Value**

Changing any value in the data set changes the mean.

3. **Mean Lies Between the Smallest and Largest Values**

The arithmetic mean is always somewhere within the data range.

4. **Mean is Unique**

For a given data set, there is only one mean.

5. **Mean Can Be Combined for Multiple Groups**

A combined mean can be calculated for multiple groups using:

$$\bar{x} \text{ (combined)} = (n_1\bar{x}_1 + n_2\bar{x}_2) \div (n_1 + n_2)$$

#### Applications:

- Business: Average sales, revenue, cost analysis
- Education: Average marks, grades, performance analysis
- Economics: Average income, prices, and growth
- Engineering: Measurement and quality control
- Daily Life: Average expenditure, average speed

### 5.2.4 Merits and Limitations of Arithmetic Mean

#### Merits:

1. **Simple to Calculate and Understand**

The formula is easy and requires basic arithmetic.

2. **Uses All Observations**

Every value in the data set contributes to the mean.

3. **Mathematically Useful**

The mean is used in many statistical formulas and further analysis.

4. **Consistent and Rigidly Defined**

The result is unique and not open to interpretation.

5. **Suitable for Further Statistical Operations**

It allows for algebraic treatment such as in standard deviation and correlation.

#### **Limitations:**

1. **Affected by Extreme Values**

Very high or low values can distort the mean.

2. **Not Always a Realistic Value**

The mean may be a value that doesn't actually exist in the data (e.g., average family size of 4.2).

3. **Not Suitable for Qualitative Data**

The mean cannot be used with categories like gender or eye color.

4. **May Mislead if Data Is Skewed**

In skewed distributions, the mean may not represent the typical value accurately.

5. **Requires Numerical Data**

Cannot be used for non-numeric or ordinal data without assigning numeric values.

## **5.3 Positional Averages: Median and Mode**

**Positional averages** are statistical measures that depend on the position of values in a data set, rather than their arithmetic operations. Unlike the mean, these measures are **not affected by extreme values**. The most commonly used positional averages are the **median** and **mode**, along with other related measures like **quartiles**, **deciles**, and **percentiles**.

### **5.3.1 Median for Ungrouped and Grouped Data**

### A. Median for Ungrouped Data

The **median** is the value that divides the data set into two equal parts after arranging the data in ascending or descending order.

#### Steps:

1. Arrange the data in order.
2. Count the number of observations ( $n$ ).
  - If  $n$  is odd:  
Median = value at position  $(n + 1) \div 2$
  - If  $n$  is even:  
Median = average of the values at positions  $n \div 2$  and  $(n \div 2 + 1)$

#### Example (odd number of values):

Data: 25, 30, 35, 40, 45

$n = 5 \rightarrow$  Median = value at position  $(5 + 1) \div 2 = 3$ rd position  $\rightarrow$  Median = 35

#### Example (even number of values):

Data: 20, 30, 40, 50

$n = 4 \rightarrow$  Median =  $(30 + 40) \div 2 = 35$

### B. Median for Grouped Data

For **continuous grouped data**, the median is calculated using the following formula:

#### Formula:

$$\text{Median} = L + [(N \div 2 - F) \div f] \times h$$

Where:

- $L$  = lower boundary of the median class
- $N$  = total frequency
- $F$  = cumulative frequency before the median class
- $f$  = frequency of the median class
- $h$  = class width

#### Steps:

1. Calculate  $N = \Sigma f$
2. Find  $N \div 2$

3. Identify the class where the cumulative frequency just exceeds  $N \div 2 \rightarrow$  this is the **median class**
4. Apply the formula

**Example:**

Class Interval	Frequency (f)	Cumulative Frequency
0 – 10	4	4
10 – 20	6	10
20 – 30	10	20
30 – 40	5	25
40 – 50	5	30

- $N = 30, N \div 2 = 15$
- Median class = 20 – 30
- $L = 20, F = 10, f = 10, h = 10$

$$\text{Median} = 20 + [(15 - 10) \div 10] \times 10 = 20 + (5 \div 10) \times 10 = 20 + 5 = 25$$

### 5.3.2 Quartiles, Deciles, and Percentiles

These are **positional measures** that divide the data into equal parts:

#### A. Quartiles

Quartiles divide the data into **four equal parts**.

- $Q_1$  (First Quartile) = value below which 25% of data lies
- $Q_2$  (Second Quartile) = median = 50%
- $Q_3$  (Third Quartile) = value below which 75% of data lies

#### Grouped data formula for $Q_1, Q_3$ :

$$Q_1 = L + [(N \div 4 - F) \div f] \times h$$

$$Q_3 = L + [(3N \div 4 - F) \div f] \times h$$

#### B. Deciles

Deciles divide the data into **ten equal parts**.

- $D_1, D_2, \dots, D_9$
- $D_5 = \text{Median}$

**Formula:**

$$D_k = L + [(kN \div 10 - F) \div f] \times h$$

Where  $k = 1$  to  $9$

**C. Percentiles**

Percentiles divide the data into **100 equal parts**.

- $P_1, P_2, \dots, P_{99}$
- $P_{50} = \text{Median}$

**Formula:**

$$P_k = L + [(kN \div 100 - F) \div f] \times h$$

Where  $k = 1$  to  $99$

These measures are commonly used in test scores, income distributions, and rankings.

**Did You Know?**

“In **competitive exams** like SAT, GRE, or national aptitude tests, your score is often reported in **percentiles**, not just marks. A percentile score of 90 means you scored **better than 90% of the test-takers**—a positional statistic used to rank performance without disclosing raw scores.”

**5.3.3 Mode for Ungrouped and Grouped Data****A. Mode for Ungrouped Data**

The **mode** is the value that appears **most frequently** in a data set.

**Example:**

Data: 12, 14, 14, 18, 19

Mode = 14 (as it appears twice)

There can be:

- No mode (all values appear once)
- One mode (unimodal)
- Two modes (bimodal)
- More than two modes (multimodal)

## B. Mode for Grouped Data

For **continuous grouped data**, the mode is calculated using the following formula:

### Formula:

$$\text{Mode} = L + [(f_1 - f_0) \div (2f_1 - f_0 - f_2)] \times h$$

Where:

- L = lower boundary of the modal class
- $f_1$  = frequency of the modal class
- $f_0$  = frequency of the class before modal class
- $f_2$  = frequency of the class after modal class
- h = class width

### Steps:

1. Identify the class with the highest frequency → this is the **modal class**
2. Plug values into the formula

### Example:

Class Interval	Frequency
10 – 20	5
20 – 30	12
30 – 40	18 ← Modal class
40 – 50	10
50 – 60	5

$$L = 30, f_1 = 18, f_0 = 12, f_2 = 10, h = 10$$

$$\text{Mode} = 30 + [(18 - 12) \div (2 \times 18 - 12 - 10)] \times 10$$

$$= 30 + (6 \div 14) \times 10$$

$$= 30 + 4.29 \approx 34.29$$

### 5.3.4 Comparison between Mean, Median, and Mode

Basis of Comparison	Mean	Median	Mode
Definition	Sum of values ÷ Number of values	Middle value of ordered data	Most frequently occurring value

Use of All Data	Uses all values	Uses position only	Uses only frequent values
Affected by Outliers	Yes	No	No
Type of Data	Quantitative	Quantitative	Both quantitative and categorical
Mathematical Use	Useful in further analysis	Limited mathematical use	Not used in advanced calculations
Stability	Stable across samples	Less stable	Can be unstable
Real-World Example	Average salary	Median household income	Most common shoe size

**Key Insight:**

- **Mean** is best for symmetric distributions without outliers.
- **Median** is better when data is skewed or contains outliers.
- **Mode** is ideal for identifying the most common value, especially in categorical data.

**“Activity: Choosing the Best Average”**

Select a dataset containing at least **ten numerical values** from any real-life scenario of your choice—such as daily expenses, step counts from a fitness tracker, monthly rainfall, or scores from a recent cricket or football series. Calculate the **mean**, **median**, and **mode** for this data. Carefully observe how similar or different these three values are. Based on your results, analyze whether the dataset is **symmetrical, positively skewed, or negatively skewed**. In a paragraph, explain which measure of central tendency best represents your data and justify your reasoning. Make sure your answer reflects a clear comparison among the three measures and demonstrates understanding of when each is most appropriate.

**5.4 Empirical Analysis of Central Tendency**

The empirical analysis of central tendency focuses on how **mean**, **median**, and **mode** relate to one another in different types of data distributions. It also examines **which measure is most appropriate** in a given situation and how these measures apply in real-world fields like business and economics.

**5.4.1 Relationship among Mean, Median, and Mode**

In a **perfectly symmetrical (normal) distribution**, all three measures—mean, median, and mode—are equal:

$$\text{Mean} = \text{Median} = \text{Mode}$$

However, in **asymmetrical (skewed) distributions**, the relationship among them changes:

- In a **positively skewed** distribution (tail on the right):

$$\text{Mean} > \text{Median} > \text{Mode}$$

- In a **negatively skewed** distribution (tail on the left):

$$\text{Mean} < \text{Median} < \text{Mode}$$

Understanding the pattern of skewness helps in choosing the appropriate measure of central tendency. In skewed distributions, the **median** is often a better representative of the central value because it is not affected by extreme values.

### 5.4.2 Karl Pearson's Empirical Formula

When either the **mode** or **median** is not available or difficult to calculate, Karl Pearson proposed an empirical formula to estimate the missing value based on known relationships.

#### A. Empirical Relationship Formula:

$$\text{Mode} \approx 3 \times \text{Median} - 2 \times \text{Mean}$$

This formula is useful in moderately skewed distributions where the exact mode is not clear but the mean and median are known.

#### B. Rearranged Forms:

- $\text{Median} \approx (\text{Mode} + 2 \times \text{Mean}) \div 3$
- $\text{Mean} \approx (3 \times \text{Median} - \text{Mode}) \div 2$

These formulas are **approximations**, not exact values, and should be used only when the distribution is not highly skewed.

#### Example:

If Mean = 60, Median = 55,

$$\text{Then Mode} \approx 3 \times 55 - 2 \times 60 = 165 - 120 = 45$$

This suggests the data is positively skewed, with the mean being greater than the mode.

**Did You Know?**

“Karl Pearson’s empirical formula—**Mode**  $\approx$  **3**  $\times$  **Median**  $-$  **2**  $\times$  **Mean**—was developed before the widespread use of computers. It provided a **quick mental estimate** of skewness in distribution and was often used in early social and economic research to interpret census or survey data with limited tools.”

### 5.4.3 Situational Use of Different Measures

Choosing the correct measure of central tendency depends on the nature of the data and the purpose of the analysis.

Below are some situations and the preferred measures:

Situation	Preferred Measure	Reason
Data with extreme values (outliers)	Median	Not affected by outliers
Categorical data (e.g., favorite brand)	Mode	Only measure that applies to qualitative data
Symmetrical distribution	Mean	Uses all data and supports further calculations
Skewed income or wealth data	Median	Provides a more realistic picture
Most common product size or choice	Mode	Shows the highest frequency
Financial or scientific calculations	Mean	Required for mathematical analysis

### 5.4.4 Practical Applications in Business and Economics

Central tendency measures are widely used in business and economics for data analysis, forecasting, planning, and decision-making.

#### A. Applications of Mean

- **Marketing:** Average customer spending helps set pricing strategies.
- **Finance:** Average returns on investment are used for risk analysis.
- **Production:** Average production per worker helps measure efficiency.
- **Sales:** Average monthly sales help in goal setting and performance tracking.

#### B. Applications of Median

- **Income Statistics:** Median income gives a better picture than mean when income is highly skewed.
- **Real Estate:** Median home price is used to represent the typical market price.
- **Public Policy:** Median household consumption helps in welfare planning.

#### C. Applications of Mode

- **Inventory Management:** Mode helps identify the most commonly sold product size or color.
- **Fashion Industry:** Used to determine the most popular clothing size.

- **Healthcare:** Mode can show the most frequent diagnosis in a hospital.

Each measure serves different purposes. Understanding when and how to use them is essential for drawing accurate conclusions from data.

### Knowledge Check 1

#### Choose the correct option:

1. Which of the following measures is most affected by extreme values (outliers)?
  - A) Median
  - B) Mode
  - C) Mean
  - D) None of the above
2. In a positively skewed distribution, the correct order of central tendency measures is:
  - A) Mean < Median < Mode
  - B) Mode < Median < Mean
  - C) Mean = Median = Mode
  - D) Median < Mode < Mean
3. The arithmetic mean of 10 numbers is 45. If one number is removed and the new mean becomes 40, what is the value of the number removed?
  - A) 45
  - B) 40
  - C) 95
  - D) 100
4. Which of the following is the correct formula to find the mode in grouped data?
  - A)  $\text{Mode} = L + [(N/2 - F) \div f] \times h$
  - B)  $\text{Mode} = L + [(f_1 - f_0) \div (2f_1 - f_0 - f_2)] \times h$
  - C)  $\text{Mode} = L + [(3 \times \text{Median} - 2 \times \text{Mean})]$
  - D)  $\text{Mode} = L + [(F - f) \div h] \times N$
5. Which of the following best describes the median?
  - A) The value that occurs most frequently
  - B) The average of all values

- C) The middle value when data is ordered
- D) The sum of frequencies divided by the number of classes

## 5.5 Summary

- ❖ This unit explored the concept of **central tendency**, a foundational idea in descriptive statistics used to represent large datasets with a single, representative value.
  - **Measures of Average** provide insight into the central value of a data set and include the **arithmetic mean**, **median**, and **mode**.
  - The **arithmetic mean** is the sum of all values divided by the number of values and is most useful for symmetrical distributions and further statistical analysis.
  - **Positional averages**, such as **median** and **mode**, are based on the position or frequency of data and are better suited for skewed distributions or non-numeric data.
  - The relationship between the three measures can indicate the **skewness** of the distribution. **Karl Pearson's empirical formula** helps estimate one measure using the others.
  - Practical applications span across business, economics, healthcare, and social science fields, where choosing the appropriate measure improves data-driven decisions.

## 5.6 Key Terms

1. **Central Tendency** - A statistical measure that identifies a single value representing a dataset
2. **Arithmetic Mean** - The sum of all observations divided by the number of observations
3. **Median** - The middle value when data is arranged in order
4. **Mode** - The value that occurs most frequently in a dataset
5. **Grouped Data** - Data organized into class intervals with corresponding frequencies
6. **Weighted Mean** - Mean where each value is assigned a weight based on importance or frequency
7. **Skewness** - A measure of asymmetry in a data distribution
8. **Modal Class** - The class interval with the highest frequency
9. **Empirical Formula** - An approximate relationship:  $\text{Mode} \approx 3 \times \text{Median} - 2 \times \text{Mean}$
10. **Quartiles** - Values that divide the data into four equal parts
11. **Percentiles** - Values that divide the data into 100 equal parts

## 5.7 Descriptive Questions

1. Define central tendency and explain its importance in statistical analysis.
2. Describe the steps to calculate the arithmetic mean for ungrouped data.
3. Differentiate between the arithmetic mean and the median with suitable examples.
4. Explain the formula and steps involved in calculating the mode for grouped data.
5. In what types of data distributions is the median preferred over the mean?
6. State Karl Pearson's empirical relationship among mean, median, and mode.
7. Write short notes on:
  - Quartiles
  - Weighted Mean
  - Modal Class
8. How does skewness affect the values of mean, median, and mode?
9. Explain one real-life situation where mode is the most appropriate measure.
10. Why is the arithmetic mean considered suitable for further mathematical treatment?

## 5.8 References

1. Gupta, S. C. (2014). *Fundamentals of Statistics*. Himalaya Publishing House.
2. Levin, R. I., & Rubin, D. S. (2012). *Statistics for Management*. Pearson Education.
3. Sharma, J. K. (2018). *Business Statistics*. Vikas Publishing House.
4. Spiegel, M. R., & Stephens, L. J. (2018). *Schaum's Outline of Statistics*. McGraw Hill.
5. UGC e-Pathshala. (n.d.). *Descriptive Statistics Modules*. Retrieved from <https://epgp.inflibnet.ac.in/>
6. Government of India. (n.d.). *National Statistical Handbook*. Ministry of Statistics and Programme Implementation.

### Answers to Knowledge Check

#### ***Knowledge Check 1***

1. C) Mean
2. B) Mode < Median < Mean
3. C) 95

4. B) Mode =  $L + [(f_1 - f_0) \div (2f_1 - f_0 - f_2)] \times h$
5. C) The middle value when data is ordered

## 5.9 Case Study

### “The Role of Central Tendency in Retail Strategy Optimization”

#### Introduction

In the fast-paced world of retail, data-driven decision-making is essential for staying competitive. One of the core tools used by retail analysts and managers is **central tendency**, which allows them to summarize large volumes of customer, sales, and pricing data into meaningful insights. Measures such as **mean**, **median**, and **mode** help businesses identify typical values and trends, enabling strategic planning and efficient resource allocation.

This caselet explores how a national retail chain applied concepts of central tendency to optimize store operations and marketing efforts. It highlights real-world challenges such as misleading averages due to outliers, the need for positional averages in skewed data, and the importance of choosing the right statistical measure depending on business context.

#### Background

RetailMart, a chain of 200 stores across India, observed significant differences in customer purchasing behavior across regions. Some urban stores showed very high average monthly sales, while many rural stores lagged behind. The company’s central management initially relied on **arithmetic mean** to assess store performance. However, this led to unrealistic targets for many branch managers whose stores were not comparable to metro locations.

A deeper analysis revealed that the **mean** was skewed by a few high-performing outlets. To address this, the analytics team began comparing **median** and **mode** values alongside the mean. The **median monthly sales** gave a better understanding of what a "typical" store achieved, while the **mode of footfall** helped identify common traffic levels and plan staffing accordingly.

RetailMart’s use of all three central tendency measures allowed them to:

- Adjust regional targets based on **median sales**, not just the mean
- Identify **popular product categories** using mode
- Reallocate promotional budgets to stores closer to the **typical performance** levels

#### Problem Statement 1: Misleading Store Performance Benchmarks

RetailMart's reliance on the arithmetic mean led to **unrealistic sales targets** for smaller outlets, especially in tier-2 and tier-3 cities.

**Solution:**

Use the **median** instead of the mean for performance benchmarking. The median excludes extreme values and reflects a more accurate central figure in a skewed data set.

**Problem Statement 2: Inaccurate Staffing Based on Average Footfall**

The mean footfall data suggested evenly distributed traffic, but in reality, most stores saw **peak volumes** on weekends and very low weekday traffic.

**Solution:**

Use the **mode** to identify the most frequent footfall pattern. Store managers then scheduled staff based on peak traffic days, improving efficiency.

**Problem Statement 3: One-size-fits-all Promotion Strategy**

Using only mean sales data, RetailMart launched promotional campaigns targeting stores with high average performance, overlooking those near the national average.

**Solution:**

Combine all three measures—**mean, median, and mode**—to segment stores and tailor promotions. High performers received different incentives, while median stores were offered promotional support to raise performance.

**Conclusion**

This case study illustrates the practical value of understanding and applying different measures of central tendency in a business context. By moving beyond a single metric and choosing the right average based on data distribution, RetailMart improved planning, fairness in performance reviews, and operational decision-making. Proper application of statistical tools like **mean, median, and mode** can transform raw data into actionable business intelligence.

## Unit 6: Measures of Dispersion

### Learning Objectives

1. Explain the purpose and importance of measuring dispersion in statistical data, and distinguish it from central tendency.
2. Calculate the range and interpret it as a basic measure of variability within a dataset.
3. Understand and compute quartile and percentile measures, and use them to analyze data distribution across different segments.
4. Calculate and interpret mean deviation for both ungrouped and grouped data, using both actual and assumed means.
5. Compute standard deviation and variance, and understand their role as key indicators of spread and consistency in a dataset.
6. Compare and contrast different measures of dispersion, and evaluate their appropriateness based on data characteristics and analytical needs.
7. Apply the concepts of dispersion to real-life business and economic scenarios, interpreting results for better decision-making and risk assessment.

### Content

- 6.0 Introductory Caselet
- 6.1 Objectives of Measuring Dispersion
- 6.2 Range
- 6.3 Quartile and Percentile Measures
- 6.4 Mean Deviation
- 6.5 Standard Deviation and Variance
- 6.6 Summary
- 6.7 Key Terms
- 6.8 Descriptive Questions
- 6.9 References
- 6.10 Case Study

## 6.0 Introductory Caselet

### “Ananya’s Apparel Analytics: Understanding Customer Spending Spread”

#### Background:

Ananya runs a mid-sized clothing brand with both online and offline stores across three cities. After a year of operations, she observed that while average customer spending per visit was ₹1,200, sales still fluctuated unpredictably. On some days, most customers made small purchases; on others, a few high-spending buyers pushed sales figures up dramatically.

Though she had already calculated the **average purchase value**, Ananya realized it didn’t tell the full story. She needed to understand **how widely customer spending varied** from that average. Her marketing consultant introduced her to the concept of **dispersion**, which measures the spread or variability in a data set.

Ananya began by calculating the **range** and noticed the difference between the lowest and highest purchase amounts was significant. She then moved on to more reliable indicators like **standard deviation** and **variance** to determine how consistent customer spending was.

She learned:

- A **low standard deviation** indicated that most customers spent close to the average,
- A **high standard deviation** suggested irregular spending patterns,
- **Quartile analysis** revealed that the top 25% of customers accounted for more than half of total revenue.

By analyzing dispersion:

- She **redesigned loyalty offers** for the top spending quartile,
- Created **price-sensitive promotions** for the bottom quartile,
- And set **sales targets** more realistically based on spending variability, not just averages.

For Ananya, understanding the **spread of data** became more valuable than relying on averages alone.

#### Critical Thinking Question:

If you were Ananya, and you found that your average sales were stable but your standard deviation was increasing every month, what might that indicate? What actions would you take in terms of marketing or customer segmentation?

## 6.1 Objectives of Measuring Dispersion

While measures of central tendency such as the mean, median, and mode provide a summary of where data is centered, they do not reveal **how much the data varies**. Two datasets can have the same average but show very different patterns in terms of consistency, variability, or spread. This is where **measures of dispersion** play a key role.

Dispersion tells us **how much the data values differ from each other and from the central value**. It reflects the degree of consistency or variability in a dataset and helps us interpret data more accurately.

### 6.1.1 Meaning of Dispersion

**Dispersion** refers to the **extent to which the values in a data set are spread out or scattered**. It indicates the **degree of variation or diversity** within a dataset.

If all values in a dataset are the same or close to each other, the dispersion is low. If values are spread far apart, the dispersion is high.

#### Example:

Consider these two sets:

- Set A: 45, 46, 47, 48, 49
- Set B: 20, 35, 50, 65, 80

Both may have similar means, but Set B clearly has more variability. That difference is captured by dispersion.

## 6.1.2 Importance of Dispersion in Statistics

### Unveiling the Multifaceted Role of Dispersion



**Fig.6.1. Importance of Dispersion in Statistics**

Understanding dispersion is essential for several reasons:

**1. Reveals Data Consistency:**

Dispersion shows how consistent or reliable the data is. Lower dispersion indicates more uniformity, which may be desirable in quality control, for example.

**2. Enhances Interpretation of Central Tendency:**

Averages alone can be misleading. Dispersion provides the context needed to interpret the mean or median effectively.

**3. Helps in Risk Analysis:**

In finance or business, higher dispersion in returns means higher risk. Investors rely on standard deviation and variance to assess risk.

**4. Supports Comparison Between Data Sets:**

Two datasets with the same average can behave very differently. Dispersion helps identify which one is more stable or predictable.

### 5. Forms Basis for Advanced Statistical Tools:

Many statistical models (e.g., regression, correlation, hypothesis testing) use dispersion as a key component.

#### 6.1.3 Comparison of Central Tendency and Dispersion

Aspect	Central Tendency	Dispersion
Purpose	Identifies the central or average value	Measures how spread out the values are
Common Measures	Mean, Median, Mode	Range, Variance, Standard Deviation
Sensitivity	Sensitive to values, especially mean	Sensitive to data variation
Completeness	Does not reveal data variability	Complements central tendency
Interpretation	Shows typical value	Shows reliability or consistency

**Conclusion:** Central tendency tells "where" the data is centered; dispersion tells "how much" the data varies around that center. Both are needed for a complete statistical understanding.

#### 6.1.4 Applications of Dispersion in Business and Economics



**Fig.6.2. Applications of Dispersion in Business and Economics**

Dispersion plays a vital role in real-world decision-making:

1. **Finance and Investment:**

Standard deviation is used to measure volatility of stock returns. A higher standard deviation means higher investment risk.

2. **Production and Quality Control:**

Consistent production output with low dispersion is a sign of process stability. Wide variation can lead to waste or customer dissatisfaction.

3. **Human Resources and Payroll:**

Understanding wage dispersion helps analyze income inequality, employee satisfaction, or the fairness of compensation policies.

4. **Marketing and Customer Analysis:**

Dispersion in customer spending reveals whether a business relies on a few high-paying customers or has a broad, consistent customer base.

5. **Economic Policy and Planning:**

Governments assess regional disparities using measures like income or employment dispersion to guide policy decisions and resource allocation.

## 6.2 Range

The **range** is the simplest measure of dispersion. It shows the **difference between the highest and lowest values** in a dataset. Although basic, it provides a quick snapshot of variability, especially useful in comparing the spread of two or more data sets.

### 6.2.1 Definition and Formula for Range

**Definition:**

The **range** is a measure of dispersion that indicates the **spread between the maximum and minimum values** in a dataset. It tells us how far apart the extreme values are.

**Formula (for ungrouped data):**

**Range = Largest value – Smallest value**

**Example:**

If the temperatures recorded over a week are:

28°C, 32°C, 35°C, 30°C, 27°C, 29°C, 34°C

Then,

$$\text{Range} = 35 - 27 = 8^\circ\text{C}$$

**Coefficient of Range (Relative Range):**

To compare datasets with different units or scales, we use the **coefficient of range**:

$$\text{Coefficient of Range} = (L - S) \div (L + S)$$

Where:

L = Largest value

S = Smallest value

**Example:**

If L = 80 and S = 20

$$\text{Coefficient of Range} = (80 - 20) \div (80 + 20) = 60 \div 100 = 0.6$$

## 6.2.2 Merits and Limitations of Range

**Merits:**

1. **Simple and Easy to Calculate**

Requires only two values—the maximum and minimum. Useful for quick comparisons.

2. **Gives a Quick Sense of Spread**

Helpful in early data exploration or when only the extremes are of interest.

3. **Applicable to Both Qualitative and Quantitative Data**

Can be used for numerical values (like sales or temperature) or categorical data that can be ordered (like ranks).

4. **Useful in Quality Control**

Especially in manufacturing, where detecting variation between extremes is critical.

**Limitations:**

1. **Ignores All Intermediate Values**

It does not consider how data behaves between the smallest and largest values.

2. **Highly Affected by Outliers**

One extreme value can distort the range significantly.

### 3. Not Reliable for Large or Grouped Datasets

Provides no insight into the distribution or consistency of the majority of the values.

### 4. Not Suitable for Statistical Analysis

Range alone cannot support further statistical methods like standard deviation or hypothesis testing.

## Did You Know?

“Although **range** is the simplest measure of dispersion, it was historically one of the first statistical tools used in **meteorology**. Scientists used it to identify climate variability in different regions long before more advanced measures like standard deviation were developed.”

### 6.2.3 Applications of Range

The range is widely used across various fields, especially where quick assessments of variability or risk are required.

#### 1. Weather Forecasting:

Meteorologists use the range to compare temperature variation between cities or across seasons.

#### 2. Stock Market Analysis:

Traders use daily price ranges (high – low) to assess volatility. A wide range indicates a more volatile asset.

#### 3. Sports Performance:

In games like cricket or football, range in scores or distances can reflect consistency or variability among players.

#### 4. Quality Control in Manufacturing:

Engineers monitor the range of product measurements (e.g., lengths, weights) to ensure standard compliance.

#### 5. Education and Test Analysis:

Teachers may look at the range of exam scores to determine how much student performance varies in a class.

#### 6. Business Decision-Making:

Comparing sales ranges across branches can help identify performance variability.

## 6.3 Quartile and Percentile Measures

While the range considers only the extreme values, **quartiles and percentiles** focus on **positions within the dataset**, providing a more refined understanding of how values are distributed. These measures divide the data into equal parts and are especially useful in identifying **central blocks of data** and spotting **outliers or skewness**.

### 6.3.1 Quartile Deviation (QD)

**Definition:**

**Quartile Deviation** (also called **Semi-Interquartile Range**) is a measure of dispersion based on the **middle 50% of the data**. It uses the **first quartile ( $Q_1$ )** and **third quartile ( $Q_3$ )** to calculate how spread out the central portion of data is.

**Formula:**

$$\text{Quartile Deviation (QD)} = (Q_3 - Q_1) \div 2$$

**Coefficient of QD:**

$$\text{Coefficient of QD} = (Q_3 - Q_1) \div (Q_3 + Q_1)$$

**Example:**

If  $Q_1 = 45$  and  $Q_3 = 65$ ,

$$\text{Then QD} = (65 - 45) \div 2 = 20 \div 2 = \mathbf{10}$$

$$\text{Coefficient} = (65 - 45) \div (65 + 45) = 20 \div 110 \approx \mathbf{0.18}$$

**Note:** QD is less affected by extreme values and is particularly useful in **skewed distributions**.

### 6.3.2 Interquartile Range (IQR)

**Definition:**

The **Interquartile Range (IQR)** is the **range of the middle 50%** of the dataset. It measures the spread between the first and third quartiles.

**Formula:**

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

**Use:**

IQR is a **robust measure** of variability because it **ignores outliers** at the extremes. It is commonly used in **box plots** and to detect **outliers** (values beyond  $1.5 \times \text{IQR}$  from  $Q_1$  or  $Q_3$ ).

**Example:**

If  $Q_1 = 40$  and  $Q_3 = 80$

$$\text{Then IQR} = 80 - 40 = \mathbf{40}$$

**“Activity: Analyzing Income Segments Using Interquartile Range (IQR)”**

**Instruction to Student:**

You are provided with the monthly income data (in ₹) of 15 families in a neighborhood:

₹18,000, ₹22,000, ₹24,000, ₹25,000, ₹27,000, ₹28,000, ₹30,000, ₹32,000, ₹33,000, ₹34,000, ₹35,000, ₹37,000, ₹40,000, ₹42,000, ₹48,000

1. Arrange the data in ascending order (if not already).
2. Identify  $Q_1$  (first quartile),  $Q_2$  (median), and  $Q_3$  (third quartile).
3. Calculate the **interquartile range (IQR)**.
4. Based on the IQR, analyze the **income inequality** in this neighborhood.
5. Briefly comment on whether the central 50% of families have closely grouped incomes or not.

Submit your quartile calculations and a 4–5 line interpretation.

### 6.3.3 Percentile Measures of Dispersion

**Definition:**

**Percentiles** divide a dataset into **100 equal parts**. They are used to understand **relative standing** and **distribution** within large datasets.

- The **10th percentile ( $P_{10}$ )** is the value below which 10% of the observations fall.
- The **90th percentile ( $P_{90}$ )** is the value below which 90% of the observations fall.
- The **percentile range** between any two percentiles can be used to measure spread.

**Formula (for percentile range):**

$$\text{Percentile Range} = P_{90} - P_{10}$$

This is sometimes used as an **alternative to the IQR**, especially when the focus is on a broader central range.

**Example:**

If  $P_{10} = 20$  and  $P_{90} = 90$ ,

Then Percentile Range =  $90 - 20 = 70$

Percentile-based ranges are frequently used in **educational testing**, **demographic studies**, and **health indicators**.

**Did You Know?**

“In **large-scale educational testing** (like SAT, GRE, and India’s NEET), students are not ranked by marks but by **percentiles**. A student with a percentile of 95 didn’t score 95 out of 100, but instead **scored better**”

than 95% of all candidates. This makes percentile-based dispersion measures essential for interpreting competitive exam results.”

### 6.3.4 Uses of Quartiles and Percentiles

Quartiles and percentiles are widely used in various fields for analysis and comparison:

#### 1. Education and Testing:

Percentile ranks show how a student performed relative to peers (e.g., "You scored in the 85th percentile").

#### 2. Human Resources and Salary Analysis:

Quartiles are used to evaluate wage distribution and to develop compensation benchmarks. For example, companies may target top performers in the 75th percentile.

#### 3. Medical and Health Research:

Percentiles are used to interpret growth charts (e.g., child height or weight in the 50th percentile).

#### 4. Market Research and Customer Segmentation:

Used to analyze spending patterns or customer value by dividing customers into quartiles or deciles.

#### 5. Outlier Detection:

Values falling outside the range  $Q_1 - 1.5 \times IQR$  to  $Q_3 + 1.5 \times IQR$  are often considered outliers in data visualization and cleaning.

#### 6. Data Summarization:

Quartiles help in generating **box plots**, a key tool in exploratory data analysis.

## 6.4 Mean Deviation

**Mean Deviation** is a measure of dispersion that indicates the **average of the absolute deviations** of each observation from a central value (mean, median, or mode). Unlike standard deviation, it does not square the deviations, making it easier to interpret and less sensitive to extreme values.

### 6.4.1 Definition and Calculation of Mean Deviation

#### Definition:

Mean deviation (also called **average absolute deviation**) is the **average of the absolute differences** between each data point and a central value (typically the mean or median).

### Formula (Ungrouped Data):

$$MD = (\Sigma|x - A|) \div n$$

Where:

- $x$  = individual observations
- $A$  = central value (mean, median, or mode)
- $n$  = number of observations
- $|x - A|$  = absolute deviation (ignoring negative signs)

### Example:

Data: 5, 7, 9

$$\text{Mean } (\bar{x}) = (5 + 7 + 9) \div 3 = 21 \div 3 = 7$$

$$MD = (|5-7| + |7-7| + |9-7|) \div 3 = (2 + 0 + 2) \div 3 = \mathbf{1.33}$$

### Grouped Data Formula:

$$MD = (\Sigma f \times |x - A|) \div \Sigma f$$

Where:

- $f$  = frequency of each class
- $x$  = mid-point of class
- $A$  = central value (mean, median, or mode)

## 6.4.2 Mean Deviation about Mean, Median, and Mode

Mean deviation can be calculated **with respect to** any of the three central values:

### A. Mean Deviation about Mean

Most commonly used in mathematical/statistical operations.

#### Formula:

$$MD_{\bar{x}} = (\Sigma|x - \bar{x}|) \div n$$

### B. Mean Deviation about Median

Often used when data is skewed because the **median is less affected by outliers**.

#### Formula:

$$MD_{me} = (\Sigma|x - M|) \div n$$

### C. Mean Deviation about Mode

Used occasionally, especially in categorical or modal distributions.

**Formula:**

$$MD_{mo} = (\sum|x - Mo|) \div n$$

**Note:**

Mean deviation about the **median** usually results in **the smallest value** compared to other central points, making it useful for minimizing overall deviation.

### 6.4.3 Merits and Limitations of Mean Deviation

**Merits:**

1. **Easy to Understand and Interpret**

Since deviations are taken in absolute terms, results are easier to grasp.

2. **Based on All Observations**

Every value in the dataset is considered in the calculation.

3. **Better than Range**

Unlike range, it gives a clearer picture of average dispersion rather than focusing only on extremes.

4. **Useful for Skewed Distributions**

Especially when calculated from the **median**, it provides a reliable measure for asymmetrical data.

5. **Can Be Used for Both Individual and Grouped Data**

Flexible across data types.

**Limitations:**

1. **Ignores Direction of Deviation**

By taking absolute values, it doesn't distinguish whether data points are above or below the central value.

2. **Not Suitable for Further Algebraic Treatment**

Unlike variance and standard deviation, it lacks properties useful for advanced statistical operations.

3. **Less Common in Inferential Statistics**

Rarely used in models, hypothesis testing, or predictive analysis.

4. **Manual Calculation Can Be Time-Consuming**

Especially for grouped data where mid-points and frequencies must be considered.

## 6.5 Standard Deviation and Variance

**Standard deviation** and **variance** are the most widely used and reliable measures of dispersion in statistics. They not only reflect the average spread of values in a dataset but also form the basis for many advanced statistical tools and analyses. These measures are especially valuable when assessing **consistency, risk, and variability** in data.

### 6.5.1 Concept of Standard Deviation

#### Definition:

Standard deviation (SD) is the **square root of the average of the squared deviations** from the mean. It measures how much each data point differs from the mean and gives a precise idea of **data spread**.

- A **low standard deviation** means values are close to the mean (less spread).
- A **high standard deviation** means values are spread out over a wide range.

#### Symbol:

- Standard Deviation:  $\sigma$  (population),  $s$  (sample)

#### Conceptual Formula:

$$\text{Standard Deviation } (\sigma) = \sqrt{[\Sigma(x - \bar{x})^2 \div n]}$$

Where:

- $x$  = each data point
- $\bar{x}$  = mean of the data
- $n$  = total number of observations

### 6.5.2 Calculation of Standard Deviation (Ungrouped and Grouped Data)

#### A. For Ungrouped Data (Raw Data)

##### Step-by-step:

1. Find the mean ( $\bar{x}$ ).
2. Subtract the mean from each value ( $x - \bar{x}$ ).
3. Square each result.
4. Find the average of these squared differences.
5. Take the square root.

#### Formula:

$$\sigma = \sqrt{[\Sigma(x - \bar{x})^2 \div n]}$$

**Example:**

Data: 4, 6, 8

$$\text{Mean } \bar{x} = (4 + 6 + 8) \div 3 = 6$$

$$\text{Squared deviations} = (4-6)^2 + (6-6)^2 + (8-6)^2 = 4 + 0 + 4 = 8$$

$$\text{Variance} = 8 \div 3 = 2.67$$

$$\text{Standard deviation} = \sqrt{2.67} \approx 1.63$$

**B. For Grouped Data**

**Formula (Direct Method):**

$$\sigma = \sqrt{[\Sigma f(x - \bar{x})^2 \div \Sigma f]}$$

Where:

- x = mid-point of class
- f = frequency
- $\bar{x}$  = mean
- $\Sigma f$  = total frequency

**Shortcut (Assumed Mean) Method:**

If numbers are large, use:

$$\sigma = \sqrt{\{[\Sigma f(d^2) \div \Sigma f] - (\Sigma f \times d \div \Sigma f)^2\} \times h}$$

Where:

- d = (x - A)  $\div$  h (deviation from assumed mean A)
- h = class width

**“Activity: Exploring Output Consistency through Standard Deviation”**

**Instruction to Student:**

You are given the following daily output data (in units) for two machines over 7 days:

**Machine A:** 120, 118, 122, 121, 119, 120, 121

**Machine B:** 105, 115, 140, 90, 130, 95, 135

1. Calculate the **mean output** for each machine.
2. Compute the **standard deviation** for both machines using the appropriate formula for ungrouped data.
3. Interpret the results: Which machine is more consistent?

4. Write a short paragraph explaining why standard deviation is more informative than just comparing the mean output.

Submit your full calculations and interpretation.

### 6.5.3 Concept and Calculation of Variance

**Variance** is the **average of the squared deviations** from the mean. It is the **square of standard deviation** and is used to quantify the spread of the data points.

**Formula for Variance:**

$$\text{Variance } (\sigma^2) = \Sigma(x - \bar{x})^2 \div n$$

**Relationship:**

$$\text{Standard Deviation} = \sqrt{\text{Variance}}$$

**Example (continued):**

From earlier: Variance = 2.67

Then, SD =  $\sqrt{2.67} \approx 1.63$

**Note:** Variance is useful in statistical modeling, while SD is easier to interpret.

#### Did You Know?

“The term **"variance"** was first introduced by **Ronald A. Fisher** in 1918, not as a theoretical idea, but as a practical tool for analyzing agricultural crop experiments. Variance is now one of the most important tools in risk analysis, genetics, psychology, and finance, where it's used to calculate **volatility in stock returns** and **variability in genetic traits.**”

### 6.5.4 Properties and Applications of Standard Deviation

**Properties:**

1. **Non-negative**

Standard deviation is always zero or positive.

2. **Minimum Value**

If all values are the same, SD = 0.

3. **Uses All Observations**

It is based on every data point, making it comprehensive.

4. **Mathematically Treatable**

Suitable for algebraic manipulations in further statistical analysis.

5. **Affected by Extreme Values**

Since deviations are squared, large differences have a big impact.

**Applications:**

1. **Finance:**

Used to measure the **volatility** of investment returns (risk analysis).

2. **Business:**

Helps monitor **consistency in production**, quality control, and cost variability.

3. **Education:**

Helps compare **student performance spread** across different subjects or exams.

4. **Healthcare:**

Assists in understanding **variability in treatment effects**, lab test results, or response rates.

5. **Operations Management:**

Identifies fluctuations in **demand forecasting, inventory usage, and process control.**

**6.5.5 Merits and Limitations of Standard Deviation**

**Merits:**

1. **Most Accurate Measure of Dispersion**

Unlike range or mean deviation, SD reflects true variability with precision.

2. **Based on All Observations**

Makes full use of available data.

3. **Useful in Inferential Statistics**

Forms the foundation for z-scores, confidence intervals, regression, etc.

4. **Mathematically Efficient**

Can be easily manipulated in equations.

**Limitations:**

1. **Affected by Extreme Values**

Since it squares deviations, outliers have a larger effect.

2. **Complex for Beginners**

Requires more steps and understanding of squared values.

3. **Less Intuitive Interpretation**

Compared to range or mean deviation, SD is harder to understand at a glance.

4. **Sensitive to Scaling**

Changing units or measurement scale affects the SD directly.

### Knowledge Check 1

**Choose the correct option:**

1. Which of the following measures is **most affected by extreme values**?

- A) Interquartile Range
- B) Standard Deviation
- C) Median
- D) Mode

2. What does a **low standard deviation** indicate about a dataset?

- A) The data is skewed to the left
- B) The data values are widely spread
- C) The data values are tightly clustered around the mean
- D) The dataset has outliers

3. Which formula is used to calculate **quartile deviation (QD)**?

- A)  $(Q_1 + Q_3) \div 2$
- B)  $Q_3 - Q_1$
- C)  $(Q_3 - Q_1) \div 2$
- D)  $Q_3 \div Q_1$

4. If the mean of a data set is 60, and the variance is 25, what is the standard deviation?

- A) 5
- B) 12

- C) 35
- D) 8
5. In which of the following scenarios would **interquartile range (IQR)** be the most appropriate measure of dispersion?
- A) When you want to include all values
- B) When the data has extreme outliers
- C) When data is normally distributed
- D) When the mean is equal to the median

## 6.6 Summary

- ❖ This chapter introduced **dispersion**, a key concept in descriptive statistics used to measure how data values deviate from a central point. While measures of central tendency describe the center of a dataset, **dispersion helps us understand the consistency, spread, or variability** of data.
  - **Range** offers a simple measure by showing the gap between the highest and lowest values.
  - **Quartile and Percentile Measures** divide data into equal segments and are particularly helpful in comparing positions and identifying outliers.
  - **Mean Deviation** calculates the average of absolute differences from a central value and provides a more balanced view of variability.
  - **Standard Deviation and Variance** are the most mathematically robust measures, widely used in business, finance, science, and research to assess risk, consistency, and variability.
- ❖ Together, these tools equip analysts and decision-makers with a deeper understanding of data reliability, spread, and underlying patterns.

## 6.7 Key Terms

1. **Dispersion** - The degree to which data values vary from the average or central value
2. **Range** - Difference between the largest and smallest value in a dataset
3. **Quartile** - Values that divide a data set into four equal parts
4. **Interquartile Range** - The range between the third and first quartile ( $Q_3 - Q_1$ )
5. **Percentile** - A value below which a given percentage of observations fall
6. **Mean Deviation** - Average of the absolute differences from a central value

7. **Standard Deviation** - Square root of the average of squared deviations from the mean
8. **Variance** - The average of squared deviations from the mean
9. **Coefficient of Dispersion** - A relative measure of variability, independent of units
10. **Outlier** - An observation significantly different from other data points

## 6.8 Descriptive Questions

1. Define dispersion. Why is it important to study dispersion in statistics?
2. What is range? How is it calculated, and what are its limitations?
3. Explain the concept of quartile deviation and interquartile range with examples.
4. Differentiate between quartiles and percentiles.
5. How is mean deviation calculated for ungrouped and grouped data?
6. Compare mean deviation about the mean, median, and mode.
7. Define standard deviation. How is it calculated using the shortcut method for grouped data?
8. Distinguish between standard deviation and variance.
9. Explain two real-life business applications of standard deviation.
10. Compare the merits and limitations of range, mean deviation, and standard deviation.

## 6.9 References

1. Gupta, S. C. (2014). *Fundamentals of Statistics*. Himalaya Publishing House.
2. Levin, R. I., & Rubin, D. S. (2012). *Statistics for Management*. Pearson Education.
3. Sharma, J. K. (2018). *Business Statistics*. Vikas Publishing House.
4. Spiegel, M. R., & Stephens, L. J. (2018). *Schaum's Outline of Statistics*. McGraw Hill Education.
5. Ministry of Statistics and Programme Implementation (MoSPI), Government of India
6. UGC e-Pathshala: Online modules on Statistics and Data Analysis

### Answers to Knowledge Check

#### *Knowledge Check 1*

1. B) Standard Deviation
2. C) The data values are tightly clustered around the mean

3. C)  $(Q_3 - Q_1) \div 2$
4. A) 5
5. B) When the data has extreme outliers

## 6.10 Case Study

### “Identifying Variability in Production: Rohan’s Factory Dilemma”

#### Introduction

In any production-based industry, tracking the consistency of output is just as important as tracking the average. For Rohan, an operations manager at a mid-sized electronics manufacturing firm, relying only on average daily production data led to misleading conclusions. His factory operated multiple machines that produced LED circuit drivers. Initially, the focus was on daily output averages—usually around 120 units per machine. However, product quality issues and unpredictable shipment volumes indicated that more was happening beneath the surface.

When Rohan began looking beyond averages, he turned to **statistical measures of dispersion**—range, mean deviation, and standard deviation. These helped him evaluate variability in output and assess whether each machine’s performance was consistent or erratic. He realized that high variability, even with a stable average, could impact quality and delivery timelines. This case explores how Rohan used dispersion measures to uncover hidden inconsistencies in production and make data-driven operational improvements.

#### Background

Rohan collected data from two machines (Machine A and Machine B) for a 10-day production cycle. Although both machines reported an average output of 120 units per day, the daily figures for Machine B fluctuated significantly—ranging from 90 to 150 units. Meanwhile, Machine A’s outputs hovered steadily between 115 and 125.

He calculated:

- **Range:** Machine A had a range of 10 units; Machine B had a range of 60.
- **Standard Deviation:** Machine A = 4.2 units; Machine B = 18.5 units.
- **Mean Deviation:** Machine A = 3.9 units; Machine B = 17.2 units.

Rohan discovered that Machine B’s inconsistency was affecting overall production flow, leading to unanticipated stockouts and delays in packaging. He shared the findings with the maintenance team, which diagnosed intermittent component overheating in Machine B.

#### Problem Statement 1: Misleading Reliance on Averages

Focusing solely on mean output made both machines appear equally productive. However, it concealed serious differences in consistency.

**Solution:**

Rohan introduced regular calculation of **standard deviation** alongside the average in daily reports. This helped management identify inconsistencies early and plan maintenance accordingly.

**Problem Statement 2: Inventory Planning Challenges Due to Output Variability**

The wide fluctuation in Machine B's output made it difficult to plan raw material purchases and labor shifts efficiently.

**Solution:**

Rohan worked with the supply chain team to use the **interquartile range (IQR)** and standard deviation to set safety stock levels and buffer time, reducing last-minute inventory shortages.

**Problem Statement 3: Poor Performance Benchmarking Across Machines**

Machines were being evaluated on average output alone, overlooking reliability and variance.

**Solution:**

A new performance metric was introduced that considered both **mean output** and **standard deviation**. Machines with high average and low variability were rated higher.

**MCQ:**

Which of the following best helps in identifying output inconsistency when the mean is the same?

- A) Mode
- B) Median
- C) Range and Standard Deviation
- D) Total Output

**Answer:** C) Range and Standard Deviation

**Explanation:**

These measures capture the spread of data and highlight how consistent the machine's performance is, even if the average remains unchanged.

## **Conclusion**

This case highlights the importance of not just knowing **how much** is produced, but also **how consistently** it is produced. Rohan's decision to incorporate **dispersion measures** helped uncover hidden inefficiencies that the average alone could not reveal. As a result, he was able to reduce production issues, improve quality control, and better allocate resources. Understanding dispersion empowered the team to make smarter, data-backed decisions that improved operational efficiency.

## Unit 7: Probability Distributions

### Learning Objectives

1. **Explain the concept of probability distributions**, including the distinction between discrete and continuous distributions.
2. **Understand the characteristics, assumptions, and probability structure of the Binomial Distribution**, and apply it to real-life problems involving success/failure outcomes.
3. **Define and apply the Poisson Distribution**, identifying suitable conditions (such as rare events over fixed intervals of time or space) and computing relevant probabilities.
4. **Recognize the key features of the Normal Distribution**, including symmetry, bell-shaped curve, mean-variance relationship, and the empirical rule.
5. **Use the Standard Normal Distribution (Z-distribution)** to compute probabilities and understand standardization of raw scores.
6. **Differentiate between Binomial, Poisson, and Normal Distributions**, and select the appropriate distribution model based on data conditions.
7. **Apply distribution models in business, economics, operations, and quality control**, interpreting outcomes to support decision-making under uncertainty.

### Content

- 7.0 Introductory Caselet
- 7.1 Introduction
- 7.2 Binomial Distribution
- 7.3 Poisson Distribution
- 7.4 Normal Distribution
- 7.5 Summary
- 7.6 Key Terms
- 7.7 Descriptive Questions
- 7.8 References
- 7.9 Case Study

## 7.0 Introductory Caselet

### “Maya’s Delivery Dilemma: Predicting Outcomes with Probability Distributions”

#### Background:

Maya manages operations for a local express courier service that specializes in same-day deliveries. Her company serves both individual and business clients across the city. Over time, Maya noticed a puzzling trend—some days, her drivers were overwhelmed with delivery requests, while on other days, they had barely anything to do. Despite having a consistent monthly average of 200 deliveries, the daily numbers fluctuated unpredictably.

Initially, Maya tried forecasting deliveries based on averages, but the results were inconsistent and unhelpful. She needed a better way to **model the uncertainty and variability in daily delivery volumes**, so she consulted a data analyst. The analyst introduced Maya to the world of **probability distributions**.

They began by using the **Binomial Distribution** to model success/failure scenarios—like whether a package would be delivered on time or not—based on historical delivery performance. Next, for the number of delivery requests per day, especially those that occurred **randomly and independently**, they applied the **Poisson Distribution**.

For example:

- On average, 12 orders came in per hour, but the **exact number varied**. Using Poisson helped her estimate the **probability of getting 15 or more orders in a given hour**.
- To predict broader customer behavior trends over a month, the analyst used the **Normal Distribution**, which fit well for modeling the **distribution of total monthly delivery distances**.

By applying probability distributions:

- Maya could now **anticipate driver workload more accurately**,
- Allocate backup vehicles during **peak Poisson-projected hours**,
- And **forecast late delivery risks** using binomial probability models.

What seemed like randomness in daily operations began to follow statistical patterns once she applied the correct distribution models.

#### Critical Thinking Question:

If you were Maya, and you found that customer order arrivals varied randomly each hour, how would using a Poisson Distribution help you plan staffing and fleet availability? Can you think of another business process where the Poisson model might be applicable?

## 7.1 Introduction

In the study of probability and statistics, we often encounter situations where outcomes are uncertain but follow recognizable patterns. These patterns are captured through **probability distributions**—mathematical functions that describe the likelihood of different outcomes in a random experiment.

There are two main types of probability distributions:

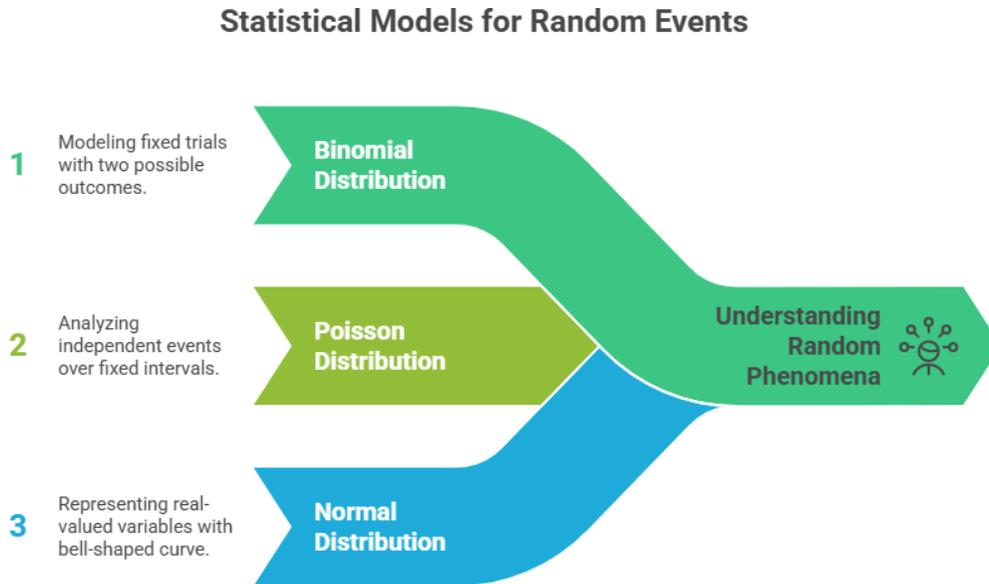
- **Discrete distributions**, where outcomes take on countable values (e.g., number of defective items in a batch).
- **Continuous distributions**, where outcomes can take any value within a given range (e.g., height of individuals, time taken to complete a task).

Understanding these distributions is essential for predicting outcomes, assessing risks, and making informed decisions in uncertain conditions. Many business, industrial, and scientific problems involve random variables that follow well-known distributions. Among these, three particularly important ones are:

- **Binomial Distribution** (discrete)
- **Poisson Distribution** (discrete)
- **Normal Distribution** (continuous)

Each of these distributions applies to specific types of data and conditions. Learning how and when to use them is foundational in probability theory and statistical modeling.

### 7.1.1 Introduction and Applications of Some Special Probability Distributions



**Fig.7.1. Introduction and Applications of Some Special Probability Distributions**

Some probability distributions occur so frequently in real-world problems that they are referred to as **special probability distributions**. These include:

#### 1. Binomial Distribution

- **Nature:** Discrete
- **Applies when:** The outcome of an experiment is binary (success or failure), repeated under the same conditions.
- **Example Applications:**
  - Predicting the number of defective products in a batch.
  - Estimating the number of successful sales calls in a day.
  - Modeling pass/fail outcomes in quality control.

#### 2. Poisson Distribution

- **Nature:** Discrete

- **Applies when:** The event is rare, occurs randomly, and independently over a fixed interval of time or space.
- **Example Applications:**
  - Modeling the number of customer calls received in an hour.
  - Estimating the number of accidents at a junction per week.
  - Predicting arrivals at a service point (e.g., patients in a clinic).

### 3. Normal Distribution

- **Nature:** Continuous
- **Applies when:** The data is symmetrically distributed around a mean and most values are close to the average.
- **Example Applications:**
  - Analyzing exam scores or employee performance ratings.
  - Forecasting demand and supply.
  - Quality control in manufacturing (e.g., measuring product dimensions).

These special distributions are powerful tools in statistical analysis. They help businesses:

- **Estimate probabilities**
- **Set control limits**
- **Model uncertainties**
- **Optimize resource allocation**

As we explore each distribution in detail, it's important to understand the **assumptions, parameters, and formulas** involved, as well as the types of problems where each is most suitable.

## 7.2 Binomial Distribution

The **Binomial Distribution** is one of the most widely used **discrete probability distributions**. It models the number of successes in a fixed number of independent trials, where each trial has only two possible outcomes: **success** or **failure**.

### 7.2.1 Concept and Characteristics of Binomial Distribution

#### Concept:

A **binomial distribution** arises when a random experiment is:

- Repeated **n** times (fixed number of trials),

- Each trial results in a **success (S)** or **failure (F)**,
- The **probability of success (p)** remains the same in each trial,
- All trials are **independent** of each other.

If  $X$  represents the number of successes in  $n$  trials, then  $X$  follows a **Binomial Distribution** with parameters  $n$  and  $p$ , written as:

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

#### Key Characteristics:

- **Discrete distribution:**  $X$  takes integer values from 0 to  $n$ .
- Two parameters:  $n$  (**number of trials**) and  $p$  (**probability of success**).
- Probability of failure  $q = 1 - p$ .
- Symmetrical if  $p = 0.5$ , otherwise skewed.
- As  $n$  increases and  $p$  is not too close to 0 or 1, the binomial distribution tends to resemble the **normal distribution**.

### 7.2.2 Probability Mass Function (PMF) of Binomial Distribution

The **PMF** gives the probability of getting exactly  $k$  successes in  $n$  trials:

$$P(X = k) = C(n, k) \times p^k \times q^{n-k}$$

Where:

- $C(n, k) = n! \div [k!(n - k)!]$  (binomial coefficient)
- $p$  = probability of success
- $q = 1 - p$  = probability of failure
- $k$  = number of successes ( $0 \leq k \leq n$ )

#### Example:

If a coin is tossed 4 times ( $n = 4$ ), and the probability of heads (success) is 0.5, the probability of getting exactly 2 heads is:

$$P(X = 2) = C(4, 2) \times (0.5)^2 \times (0.5)^2 = 6 \times 0.25 \times 0.25 = 0.375$$

#### Did You Know?

“The **Binomial PMF** is not just used in statistics—it's also used in **genetics**. In Mendelian inheritance, the probability of inheriting dominant or recessive traits follows a binomial pattern. For example, if two

heterozygous parents cross ( $Aa \times Aa$ ), the probability of their child inheriting a recessive gene ( $aa$ ) follows the same logic as calculating binomial probabilities.”

### 7.2.3 Mean and Variance of Binomial Distribution

The binomial distribution has the following statistical properties:

- **Mean ( $\mu$ )** =  $n \times p$
- **Variance ( $\sigma^2$ )** =  $n \times p \times q$
- **Standard Deviation ( $\sigma$ )** =  $\sqrt{(n \times p \times q)}$

#### Interpretation:

- The **mean** represents the expected number of successes.
- The **variance and standard deviation** indicate how much variability there is around the mean.

#### Example:

For a binomial distribution with  $n = 10$  and  $p = 0.4$ :

- Mean =  $10 \times 0.4 = 4$
- Variance =  $10 \times 0.4 \times 0.6 = 2.4$
- Standard deviation =  $\sqrt{2.4} \approx 1.55$

### 7.2.4 Applications and Examples of Binomial Distribution

The binomial distribution is widely used in situations involving repeated, independent binary trials. Common applications include:

#### Business and Quality Control:

- Estimating the probability that **a batch has a certain number of defective items**.
- Measuring the **number of successful sales calls** out of total attempts.

#### Healthcare:

- Modeling the probability that **a certain number of patients recover** after treatment.

#### Education:

- Estimating the probability of **a specific number of students passing** an exam.

#### Finance:

- Analyzing the number of **successful investments** or **profitable trading days**.

#### Example:

A factory has a 5% defect rate. In a batch of 10 items ( $n = 10$ ), what is the probability that exactly one item is defective?

Here,

- $p = 0.05$  (defective),
- $q = 0.95$  (non-defective),
- $n = 10$ ,
- $k = 1$

$$P(X = 1) = C(10, 1) \times (0.05)^1 \times (0.95)^9 \approx 10 \times 0.05 \times 0.630 = \mathbf{0.315}$$

## 7.3 Poisson Distribution

The **Poisson distribution** is a discrete probability distribution used to model the number of times an event occurs in a **fixed interval of time or space, when the events occur randomly and independently**.

### 7.3.1 Concept and Characteristics of Poisson Distribution

#### Concept:

A **Poisson distribution** describes the probability of a given number of events occurring in a **fixed interval of time, distance, area, or volume**, assuming:

- The events are **rare**, and occur **independently** of one another,
- The **average rate of occurrence ( $\lambda$ )** is constant,
- Two events cannot occur at the same instant.

#### Notation:

If  $X$  is the number of events, then:

$$X \sim \text{Poisson}(\lambda)$$

Where:

- $\lambda$  (lambda) = average number of events in the given interval

#### Characteristics:

- It is a **discrete** distribution ( $X = 0, 1, 2, 3, \dots$ ).
- It models **count data** (e.g., number of emails per hour).
- Theoretically, there is **no upper limit** to the number of occurrences.
- The events are **independent and memoryless**.

- The mean and variance are **equal** ( $\lambda$ ).

### 7.3.2 Probability Mass Function (PMF) of Poisson Distribution

The **PMF** of the Poisson distribution is:

$$P(X = k) = (e^{-\lambda} \times \lambda^k) \div k!$$

Where:

- $\lambda$  = mean number of occurrences in the interval
- $k$  = actual number of occurrences ( $k = 0, 1, 2, \dots$ )
- $e \approx 2.71828$
- $k!$  = factorial of  $k$

#### Example:

Suppose  $\lambda = 4$  (average 4 calls per hour at a call center).

What is the probability that **exactly 2 calls** are received in an hour?

$$P(X = 2) = (e^{-4} \times 4^2) \div 2!$$

$$= (0.0183 \times 16) \div 2 = 0.1464$$

#### Did You Know?

“The **Poisson distribution** was actually developed as a way to approximate the **binomial distribution** for rare events. In fact, it was first used in 1837 by Simeon Denis Poisson to analyze **the number of soldiers accidentally killed by horse kicks** in the Prussian army—making it one of the earliest data-based applications of probability theory.”

#### “Activity: Modeling Customer Support Calls with Poisson Distribution”

#### Instruction to Student:

A customer support center receives an average of 5 calls per hour. Using the **Poisson distribution**, calculate the probability that:

- a) Exactly 3 calls are received in an hour
- b) 5 or more calls are received in an hour

Use the formula:

$$P(X = k) = (e^{-\lambda} \times \lambda^k) \div k!, \text{ where } \lambda = 5$$

Show all steps and round your final probabilities to four decimal places. Comment on how useful this method would be for planning the number of support staff needed per shift.

### 7.3.3 Mean and Variance of Poisson Distribution

Both the **mean** and **variance** of the Poisson distribution are equal to  $\lambda$ .

- **Mean ( $\mu$ )** =  $\lambda$
- **Variance ( $\sigma^2$ )** =  $\lambda$
- **Standard deviation ( $\sigma$ )** =  $\sqrt{\lambda}$

#### Example:

If  $\lambda = 5$ , then:

- Mean = 5
- Variance = 5
- Standard deviation =  $\sqrt{5} \approx 2.24$

This property makes the Poisson distribution easy to interpret, especially in contexts like queueing systems and inventory analysis.

### 7.3.4 Relationship between Poisson and Binomial Distribution

The **Poisson distribution is a limiting case of the binomial distribution** under specific conditions:

- The number of trials  $n \rightarrow \infty$  (very large),
- The probability of success  $p \rightarrow 0$  (very small),
- The product  $n \times p = \lambda$  remains finite.

This relationship allows the Poisson distribution to be used as an **approximation to the binomial distribution** when:

- $n$  is large (typically  $> 20$ ), and
- $p$  is small (typically  $< 0.05$ )

#### Example Use Case:

If we want to model the number of typing errors per page, and each page has many words ( $n$  is large) but a very low probability of error per word ( $p$  is small), the Poisson model is more efficient than the binomial.

### 7.3.5 Applications and Examples of Poisson Distribution

The Poisson distribution is used widely across sectors to model **random, independent events occurring over time or space**.

#### **Business & Operations:**

- Number of customer service calls per hour
- Number of defects per 100 meters of cable

#### **Healthcare:**

- Number of patients arriving in an emergency room per night
- Occurrence of rare diseases in a population

#### **Manufacturing & Quality Control:**

- Number of flaws per square meter in fabric
- Number of defective parts per shipment

#### **Traffic and Transport:**

- Number of cars passing through a toll plaza per minute
- Number of train delays per week

#### **Example:**

A web server receives an average of 3 hits per minute. What is the probability that it receives exactly 5 hits in a given minute?

$$\lambda = 3, k = 5$$

$$P(X = 5) = (e^{-3} \times 3^5) \div 5!$$

$$= (0.0498 \times 243) \div 120 \approx 0.1008$$

So, there is a 10.08% chance the server receives exactly 5 hits in one minute.

## 7.4 Normal Distribution

The **Normal Distribution** is a continuous probability distribution that describes how the values of a random variable are distributed. It is also called the **Gaussian distribution** and is one of the most important concepts in probability and statistics due to its wide range of applications.

### 7.4.1 Concept and Characteristics of Normal Distribution

#### Concept:

The Normal Distribution models **continuous data** where values are symmetrically distributed around a central mean. The classic "bell-shaped curve" is used to represent data where:

- Most values cluster around the mean,
- Fewer values appear as we move further from the mean.

It is defined by **two parameters**:

- $\mu$  (**mu**) = Mean (center of the distribution),
- $\sigma$  (**sigma**) = Standard deviation (controls the spread).

#### Key Characteristics:

- Symmetrical about the mean ( $\mu$ ),
- Mean = Median = Mode,
- The total area under the curve = 1,
- The curve approaches, but never touches, the x-axis,
- Defined for all real values of  $x$  ( $-\infty$  to  $\infty$ ),
- The shape is completely determined by  $\mu$  and  $\sigma$ .

### 7.4.2 Probability Density Function (PDF) of Normal Distribution

The **Probability Density Function (PDF)** for a normal distribution is:

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

Where:

- $x$  = any real number,
- $\mu$  = mean,
- $\sigma$  = standard deviation,
- $e \approx 2.71828$  (base of natural logarithm),
- $\pi \approx 3.14159$

This function defines the bell-shaped curve. The highest point of the curve occurs at  $x = \mu$ , and the spread of the curve is determined by  $\sigma$ .

### Did You Know?

“The **Poisson distribution** was actually developed as a way to approximate the **binomial distribution** for rare events. In fact, it was first used in 1837 by Simeon Denis Poisson to analyze **the number of soldiers accidentally killed by horse kicks** in the Prussian army—making it one of the earliest data-based applications of probability theory.”

### 7.4.3 Properties of Normal Curve

1. **Symmetry:** The curve is symmetric about the mean. The left and right halves are mirror images.
2. **Mean = Median = Mode:** All central tendency measures lie at the center.
3. **Bell Shape:** The distribution is unimodal and smoothly curved.
4. **Tails Extend to Infinity:** The curve never touches the x-axis; it gets closer and closer but never equals zero.
5. **Area under the Curve:**
  - About **68.27%** of data lies within  $\pm 1\sigma$  from the mean,
  - About **95.45%** lies within  $\pm 2\sigma$ ,
  - About **99.73%** lies within  $\pm 3\sigma$ .
6. **Empirical Rule:** These percentages (68–95–99.7 rule) are useful for estimating probabilities.

### Did You Know?

“Over **99% of human heights**, test scores, and even errors in manufacturing **follow a normal distribution**. That’s why it’s often called the “**natural law of error**.” The fact that the mean, median, and mode are all the same in a normal distribution makes it uniquely useful for modeling **balanced systems and populations**.”

### 7.4.4 Standard Normal Distribution (Z-Distribution)

A **Standard Normal Distribution** is a special case of the normal distribution where the **mean ( $\mu$ ) is 0** and the **standard deviation ( $\sigma$ ) is 1**. It is commonly referred to as the **Z-distribution**.

#### Z-Score Formula

The **Z-score** tells you how many standard deviations a value ( $X$ ) is from the mean ( $\mu$ ):

$$Z = (X - \mu) \div \sigma$$

Where:

- $X$  = Individual value
- $\mu$  = Mean of the distribution
- $\sigma$  = Standard deviation of the distribution

### Interpretation of Z-scores

- A **positive Z-score** indicates a value **above the mean**
- A **negative Z-score** indicates a value **below the mean**
- A **Z-score of 0** means the value is **equal to the mean**

### Example

A test has the following statistics:

- Mean ( $\mu$ ) = 70
- Standard deviation ( $\sigma$ ) = 10
- A student scores  $X = 85$

Using the Z-score formula:

$$Z = (85 - 70) \div 10 = 1.5$$

This means the student scored **1.5 standard deviations above the mean**.

### Using the Z-Table

A **Z-table** (Standard Normal Table) shows the **cumulative probability** from the far left of the distribution ( $Z = -\infty$ ) **up to the given Z-score**.

For  $Z = 1.5$ , the cumulative probability is approximately **0.9332**.

This means:

- 93.32% of scores fall **below 85**
- The student is at the **93rd percentile**

### Standard Normal Table (Z-Table)

*Cumulative probabilities from the left of the distribution*

Z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
1.2	0.8849	0.8869	0.8888	0.8907	0.8925	0.8944	0.8962	0.8980	0.8997	0.9015
1.3	0.9032	0.9049	0.9066	0.9082	0.9099	0.9115	0.9131	0.9147	0.9162	0.9177
1.4	0.9192	0.9207	0.9222	0.9236	0.9251	0.9265	0.9279	0.9292	0.9306	0.9319
1.5	0.9332	0.9345	0.9357	0.9370	0.9382	0.9394	0.9406	0.9418	0.9429	0.9441

**Example:**

To find the cumulative probability for  $Z = 1.53$ , look at the row **1.5** and the column **0.03** → Value = **0.9370**

**Applications of Standard Normal Distribution**

- **Comparing scores** from different normal distributions
- **Finding percentiles** of individual scores
- **Calculating probabilities** for hypothesis testing
- **Detecting outliers**, typically using Z-scores beyond  $\pm 2$  or  $\pm 3$

**“Activity: Finding Probabilities Using Z-Scores”**

**Instruction to Student:**

The scores on a standardized test are normally distributed with a **mean of 500** and a **standard deviation of 100**.

1. Calculate the **Z-scores** for students who scored:
  - a) 620
  - b) 450
  - c) 700
2. Use a Z-table (or standard normal distribution table) to find the **probability** of scoring:
  - a) Less than 620
  - b) More than 700
  - c) Between 450 and 620

Write a short interpretation of each result. Discuss how Z-scores can help identify performance levels (e.g., above average, average, below average).

### 7.4.5 Applications and Examples of Normal Distribution

The normal distribution is extensively used in both theoretical statistics and real-life situations.

#### Applications:

1. **Business and Finance:**

- Stock price fluctuations,
- Demand forecasting,
- Return on investment (ROI) modeling.

2. **Manufacturing and Quality Control:**

- Measurement of parts (e.g., diameter, weight),
- Tolerances and control charts.

3. **Education:**

- Standardized test scores (e.g., SAT, IQ),
- Grading on a curve.

4. **Healthcare:**

- Blood pressure and cholesterol levels,
- Patient recovery times.

5. **Social Sciences:**

- Human behavior measurements,
- Population studies.

#### Example:

The heights of adult men in a city are normally distributed with a mean of 170 cm and a standard deviation of 6 cm.

- What percentage of men are taller than 182 cm?

First, calculate the Z-score:

$$Z = (182 - 170) \div 6 = 2.0$$

From the Z-table,  $P(Z < 2.0) \approx 0.9772$

So,  $P(Z > 2.0) = 1 - 0.9772 = \mathbf{0.0228}$  or **2.28%** of men are taller than 182 cm.

## 7.5 Summary

- ❖ This chapter introduced three of the most important probability distributions used in statistics:
  - The **Binomial Distribution** is used to model the number of successes in a fixed number of independent trials, where each trial has only two possible outcomes (success or failure). It is suitable for discrete data with a fixed number of repeated experiments.
  - The **Poisson Distribution** is ideal for modeling the number of occurrences of a rare event within a fixed time or space interval, assuming the events are random, rare, and independent. It is often applied in service systems, traffic studies, and quality control.
  - The **Normal Distribution** is a continuous, symmetric, bell-shaped curve widely used for naturally occurring data. It is defined by its mean and standard deviation and forms the foundation for many statistical methods, including hypothesis testing and confidence intervals.
- ❖ Each of these distributions has its own set of assumptions, formulas, and use cases, and understanding when and how to apply them is essential for effective data analysis and decision-making.

## 7.6 Key Terms

1. **Probability Distribution** - A function that describes the likelihood of different outcomes
2. **Binomial Distribution** - Discrete distribution modeling the number of successes in fixed trials
3. **Poisson Distribution** - Discrete distribution modeling rare events over fixed intervals
4. **Normal Distribution** - Continuous, symmetric distribution shaped like a bell curve
5. **Probability Mass Function (PMF)** - Function for discrete distributions defining probabilities for specific values
6. **Probability Density Function (PDF)** - Function for continuous distributions where area under curve represents probability
7. **Z-Score** - Number of standard deviations a value is from the mean
8. **Mean ( $\mu$ )** - Average value of the distribution
9. **Variance ( $\sigma^2$ )** - Measure of spread in a dataset
10. **Standard Deviation ( $\sigma$ )** - Square root of variance; indicates consistency or variability

## 7.7 Descriptive Questions

1. Define and explain the key assumptions of a binomial distribution.

2. How is the Poisson distribution related to the binomial distribution? Give an example.
3. State the formula for the probability mass function of the binomial distribution.
4. Explain how the standard normal distribution is derived and used.
5. Describe the 68–95–99.7 rule and its relevance to the normal curve.
6. Compare the binomial and Poisson distributions. When would you use each?
7. Give two practical applications of the normal distribution in business or science.
8. A machine has a 10% defect rate. What is the probability that 2 out of 5 items are defective?
9. What are the limitations of using the normal distribution in real-life data analysis?
10. A Poisson process averages 4 customer arrivals per hour. What is the probability of exactly 2 arrivals in a given hour?

## 7.8 References

1. Gupta, S. C. (2014). *Fundamentals of Statistics*. Himalaya Publishing House.
2. Levin, R. I., & Rubin, D. S. (2012). *Statistics for Management*. Pearson Education.
3. Sharma, J. K. (2018). *Business Statistics*. Vikas Publishing House.
4. Anderson, D. R., Sweeney, D. J., & Williams, T. A. (2016). *Statistics for Business and Economics*. Cengage Learning.
5. UGC e-Pathshala. *Modules on Probability Distributions and Statistical Theory*.
6. MoSPI – Ministry of Statistics and Programme Implementation, Government of India.

## 7.9 Case Study

### “Demand Forecasting with Distributions: Priya’s Inventory Planning”

#### **Background:**

Priya is a supply chain manager for a grocery distribution company that supplies perishable products to 60 retail outlets. One of her key responsibilities is to forecast the demand for certain items, like milk packets, fruits, and bread, that must be replenished daily. Despite having access to historical sales data, her team struggled with frequent understocking and overstocking.

She began exploring **probability distribution models** to make the forecasts more accurate.

- For items like **milk packets**, which either sell or do not (success/failure), she used the **binomial distribution** to estimate how many outlets would sell out on a given day.
- To model the **number of unexpected order requests during the day**, she used the **Poisson distribution**, as these arrivals were random and typically low in volume.
- For **daily total sales volume**, which followed a predictable trend with minor fluctuations, she used the **normal distribution** to determine average demand and apply safety stock buffers based on standard deviations.

#### **Problem Statement 1: Inconsistent Demand Patterns**

Despite stable average sales, stockouts occurred frequently at specific outlets.

#### **Solution:**

Priya used binomial probability to determine the likelihood of an outlet selling out and reallocated stock dynamically based on real-time success probability.

#### **Problem Statement 2: Unpredictable Same-Day Orders**

Midday order requests from retail outlets were random and difficult to plan for.

#### **Solution:**

Using the Poisson model, Priya estimated the number of extra orders expected per hour and adjusted delivery schedules accordingly.

#### **Problem Statement 3: Overestimated Safety Stock**

Holding too much inventory led to waste and increased cost.

**Solution:**

Priya applied the **normal distribution** to model demand and used **Z-scores** to set appropriate safety stock levels, covering 95% of the expected demand range.

**MCQ:**

Which distribution should Priya use to model the probability that 3 outlets will sell out of stock on a day when each has a 20% chance of doing so?

- A) Poisson Distribution
- B) Binomial Distribution
- C) Normal Distribution
- D) Uniform Distribution

**Answer:** B) Binomial Distribution

**Conclusion:**

By applying the right probability distribution to the right business scenario, Priya was able to reduce wastage, improve stock availability, and build a responsive supply chain. What initially seemed like randomness became manageable through structured statistical modeling.

## Unit 8: Correlation

### Learning Objectives

1. **Define correlation** and explain its significance in measuring the strength and direction of a relationship between two variables.
2. **Identify and differentiate between various types of correlation** (positive, negative, zero; linear and non-linear), using both graphical and numerical methods.
3. **Interpret scatter diagrams and simple correlation graphs** to visualize data patterns and preliminary relationships between variables.
4. **Calculate Karl Pearson's Coefficient of Correlation ( $r$ )** for both ungrouped and grouped data, and interpret the result in terms of strength and direction.
5. **Understand the properties and limitations of the correlation coefficient**, including its range, sensitivity to outliers, and lack of causality.
6. **Apply Spearman's Rank Correlation method** for ranked or ordinal data and evaluate correlation where data is not suitable for Pearson's method.
7. **Use correlation analysis in real-world contexts** such as economics, business, psychology, and social sciences to study relationships (e.g., income vs. expenditure, advertising vs. sales).

### Content

- 8.0 Introductory Caselet
- 8.1 Introduction
- 8.2 Types of Correlation
- 8.3 Scatter Diagram and Simple Graph
- 8.4 Karl Pearson's Coefficient of Correlation
- 8.5 Properties of Coefficient of Correlation
- 8.6 Spearman's Rank Correlation
- 8.7 Summary
- 8.8 Key Terms
- 8.9 Descriptive Questions
- 8.10 References
- 8.11 Case Study

## 8.0 Introductory Caselet

### “Ravi’s Research on Return Rates: Finding Hidden Connections”

#### Background:

Ravi, a data analyst working for an e-commerce platform, was asked to investigate an increase in product return rates across categories. His initial reports showed no clear pattern—some customers returned expensive electronics, while others returned inexpensive apparel. The marketing team suspected that advertising campaigns and aggressive discounts might be responsible, but there was no concrete evidence.

Ravi decided to explore **correlation analysis** to see if the return rates were linked to other measurable variables.

He started by collecting data on:

- Product return rate (%),
- Price of the product,
- Discount offered (%),
- Customer satisfaction score (from feedback surveys).

He began with **scatter diagrams** and noticed that for some product types, higher discounts correlated with higher return rates. He then calculated **Karl Pearson’s correlation coefficient (r)** to quantify the strength of these relationships. For certain categories, such as clothing and accessories, **r exceeded +0.75**, showing a strong **positive correlation** between high discounts and return rates.

For customer satisfaction and returns, he used **Spearman’s Rank Correlation** and found a **strong negative correlation** ( $r_s \approx -0.85$ ), indicating that products with lower satisfaction scores were more likely to be returned.

Thanks to correlation tools, Ravi could now:

- Identify which marketing strategies were triggering excessive returns,
- Recommend quality improvement for low-rated products,
- Help the company reduce return costs while protecting customer trust.

What once seemed like random consumer behavior became a pattern of **data-driven insight** powered by correlation analysis.

**Critical Thinking Question:**

If you were Ravi, and you discovered a high positive correlation between discounts and return rates in only one product category, how would you advise the marketing team? What additional variables might you analyze to ensure your recommendation is reliable?

## 8.1 Introduction

**Correlation** is a statistical tool used to measure and describe the strength and direction of a **relationship between two variables**. When two variables tend to move together in a predictable pattern, they are said to be **correlated**.

Correlation analysis helps us answer questions like:

- Do advertising expenses influence sales?
- Does employee satisfaction impact productivity?
- Is there a link between inflation and interest rates?

It is important to note that **correlation measures association**, not causation. It quantifies how closely two variables move in relation to each other but does not imply that one causes the other.

### 8.1.1 Meaning and Importance of Correlation

#### Meaning:

Correlation refers to the **degree of linear relationship** between two quantitative variables. It is measured using a **correlation coefficient**, most commonly **Karl Pearson's  $r$** , which ranges from **-1 to +1**.

- **+1** indicates a perfect positive correlation
- **-1** indicates a perfect negative correlation
- **0** indicates no linear relationship

#### Importance:

##### 1. Identifies Relationships:

Correlation helps identify whether two variables are related and the direction of that relationship.

##### 2. Supports Decision-Making:

In business, knowing the relationship between factors like sales and marketing spend helps in strategic planning.

##### 3. Reduces Uncertainty:

Understanding relationships helps reduce guesswork and allows for data-driven insights.

##### 4. Forms Basis for Further Analysis:

Correlation is the foundation for advanced statistical tools like regression analysis.

### 8.1.2 Distinction between Correlation and Causation

A **common misconception** is that correlation implies causation. This is not true.

Feature	Correlation	Causation
Definition	Measures degree of association	Indicates that one variable affects another
Direction	Can be positive, negative, or zero	Has a specific direction of influence
Proof Needed	No proof of cause-effect is needed	Requires experimental or theoretical proof
Example	Ice cream sales ↑ and drowning cases ↑	Ice cream doesn't cause drowning

**Key Point:**

Just because two variables are correlated does not mean one causes the other. They may both be influenced by a **third (lurking) variable** or be coincidentally related.

**Did You Know?**

“The famous example of **ice cream sales and drowning incidents** often cited in statistics classes is used to explain why correlation  $\neq$  causation. Although these two variables may show a strong positive correlation during summer months, one **does not cause the other**. Instead, a **lurking variable**—hot weather— influences both.”

**8.1.3 Applications of Correlation in Business and Economics**

Correlation is widely used across business functions and economic analysis. Some practical applications include:

**In Business:**

- **Marketing Analysis:** Understanding if more advertising leads to higher sales.
- **Human Resources:** Studying the link between employee engagement and performance.
- **Operations Management:** Exploring the relationship between production volume and defect rate.

**In Economics:**

- **Macroeconomic Analysis:** Examining how inflation is correlated with interest rates or unemployment.
- **Consumer Behavior:** Understanding the relationship between income levels and spending patterns.
- **Stock Market Analysis:** Studying the correlation between two stock returns to build investment portfolios.

By using correlation, professionals can make informed assumptions, plan resource allocations, and optimize strategic decisions.

**8.2 Types of Correlation**

Correlation can be classified into different types based on the **direction, nature, and number of variables involved**. Understanding these distinctions helps in selecting the right method for analysis and interpreting the results correctly.

### 8.2.1 Positive and Negative Correlation

This classification is based on the **direction of the relationship** between two variables.

#### A. Positive Correlation

Two variables are said to have a **positive correlation** when an increase in one variable leads to an increase in the other, and vice versa.

##### Examples:

- Height and weight
- Advertising spend and sales revenue
- Education level and income

**Graphically:** The points on a scatter diagram slope **upward from left to right**.

#### B. Negative Correlation

Two variables show **negative correlation** when an increase in one variable leads to a decrease in the other.

##### Examples:

- Price and quantity demanded (law of demand)
- Fuel efficiency and car weight
- Interest rate and investment levels

**Graphically:** The points slope **downward from left to right**.

### 8.2.2 Linear and Non-linear Correlation

This classification is based on the **pattern or form** of the relationship between variables.

#### A. Linear Correlation

A correlation is **linear** when the change in one variable is **proportional** to the change in another. The relationship can be represented by a **straight line**.

##### Example:

- Sales increasing proportionally with the number of sales representatives

**Equation Form:**  $y = a + bx$

### B. Non-linear (Curvilinear) Correlation

If the change in one variable is **not proportional** to the change in the other, the correlation is **non-linear**. The data points form a **curved line**.

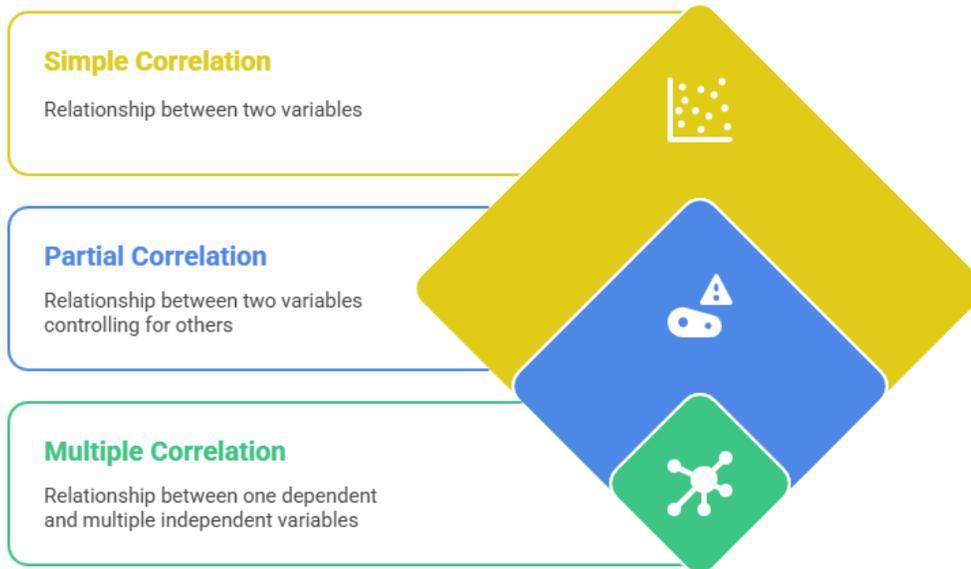
**Examples:**

- Learning curve: Speed of learning increases initially but plateaus later
- Relationship between age and income: Income may rise, then fall after retirement

**Graphically:** The scatter plot forms a **curve** rather than a straight line.

### 8.2.3 Simple, Partial, and Multiple Correlation

#### Correlation Types in Statistics



**Fig.8.1 Simple, Partial, and Multiple Correlation**

This classification is based on the **number of variables** under consideration.

### A. Simple Correlation

Involves **two variables only**, studying the relationship between one independent and one dependent variable.

#### Example:

- Relationship between temperature and electricity consumption

### B. Partial Correlation

Analyzes the relationship between **two variables** while **controlling for the effect of one or more additional variables**.

#### Example:

- Studying the effect of study time on grades, while keeping sleep hours constant

**Used when:** We want to isolate the influence of specific variables by eliminating the effect of others.

### C. Multiple Correlation

Examines the relationship between **one dependent variable and two or more independent variables** simultaneously.

#### Example:

- Predicting sales based on advertising expenditure and market size
- Analyzing job performance using experience, education level, and age

**Mathematical form:**  $y = a + b_1X_1 + b_2X_2 + \dots + b_nX_n$

### Summary Table: Types of Correlation

Basis	Types	Example
Direction	Positive / Negative	Sales & Ads / Price & Demand
Form	Linear / Non-linear	Salary & Experience / Age & Income
No. of Variables	Simple / Partial / Multiple	Rainfall & Yield / Yield & (Rain, Fertilizer)

## 8.3 Scatter Diagram and Simple Graph

Graphical methods are often the **first step** in understanding the **nature of the relationship** between two variables. Scatter diagrams and simple graphs help visualize whether a relationship exists and what kind of correlation (if any) is present.

### 8.3.1 Concept of Scatter Diagram

A **scatter diagram** (also called a **scatter plot**) is a **graphical representation of paired data**, used to determine the **existence, direction, and strength of correlation** between two variables.

- Each **point** on the graph represents a pair of values (x, y).
- The horizontal axis (X-axis) represents the **independent variable**, while the vertical axis (Y-axis) represents the **dependent variable**.

**Purpose:**

- To identify patterns or relationships,
- To judge whether a linear or non-linear trend exists,
- To observe outliers or clusters.

### 8.3.2 Method of Plotting Scatter Diagram

To create a scatter diagram, follow these steps:

1. **Collect paired data:** e.g., marks obtained in two subjects by a group of students.
2. **Label the axes:** X-axis for the first variable, Y-axis for the second.
3. **Choose appropriate scales:** Based on the range of values.
4. **Plot the points:** For each pair (x, y), mark a dot at the intersection.
5. **Observe the pattern:** This helps in understanding the type of correlation.

### 8.3.3 Interpretation of Scatter Diagrams

Based on how the points are distributed, correlation can be interpreted as follows:

Pattern of Dots	Interpretation
Dots slope <b>upward</b> to the right	<b>Positive correlation</b>
Dots slope <b>downward</b> to the right	<b>Negative correlation</b>
Dots show <b>no visible pattern</b>	<b>No correlation</b>
Dots form a <b>tight straight line</b>	<b>Perfect correlation (<math>r = \pm 1</math>)</b>
Dots are <b>loosely spread</b>	<b>Weak or moderate correlation</b>
Dots follow a <b>curved path</b>	<b>Non-linear (curvilinear) correlation</b>

**Example:**

If the number of hours studied and exam scores are plotted, and points trend upwards, it shows **positive correlation**—more study hours generally result in higher scores.

**“Activity: Exploring Correlation through Real Business Data”**

**Instruction to the Student:**

You are given a table showing monthly data of **advertising expenses (₹)** and **monthly sales (₹)** for a company over one year.

- Plot the data on a **scatter diagram** using graph paper or Excel.
- Observe the pattern and describe whether the correlation is **positive, negative, or none**.
- Write a short paragraph interpreting the relationship.
- Based on the scatter pattern, suggest whether Karl Pearson’s method is appropriate or if a non-linear analysis might be more suitable.

**8.3.4 Simple Graph Method of Studying Correlation**

While scatter diagrams plot **paired values**, **simple graphs** plot **two separate curves** on the same axes to show movement over time or another common variable.

**Steps:**

1. Collect time-series data for two related variables (e.g., monthly sales and monthly advertising expense).
2. Plot both variables on the **same graph**, usually against a common X-axis (time).
3. Use **two lines or curves**, one for each variable.
4. Compare the trends visually. If the curves move **in the same direction**, it suggests **positive correlation**. If they move **in opposite directions**, it indicates **negative correlation**.

**Advantages:**

- Useful for comparing **trends over time**,
- Can help identify **lagging relationships** (e.g., advertising effect appearing with a delay on sales).

**Comparison: Scatter Diagram vs. Simple Graph**

Feature	Scatter Diagram	Simple Graph

Data Structure	Paired observations	Time-series or trend-based data
Visualization	Points on a graph	Curves or lines
Correlation Type	Linear / Non-linear	Directional movement
Use Case	Measuring degree of association	Comparing trends over time

## 8.4 Karl Pearson’s Coefficient of Correlation

**Karl Pearson’s Coefficient of Correlation**, often denoted as **r**, is the most widely used statistical method to measure the **strength and direction of a linear relationship** between two variables. It is a **quantitative tool** that complements graphical methods like the scatter diagram.

### 8.4.1 Definition and Formula

#### Definition:

Karl Pearson’s coefficient of correlation measures the **degree of linear correlation** between two variables X and Y. The value of **r** lies between **-1 and +1**.

- **r = +1**: Perfect positive linear correlation
- **r = -1**: Perfect negative linear correlation
- **r = 0**: No linear correlation

#### Formula (for ungrouped data):

$$r = \frac{\Sigma[(x - \bar{x})(y - \bar{y})]}{\sqrt{[\Sigma(x - \bar{x})^2 \times \Sigma(y - \bar{y})^2]}}$$

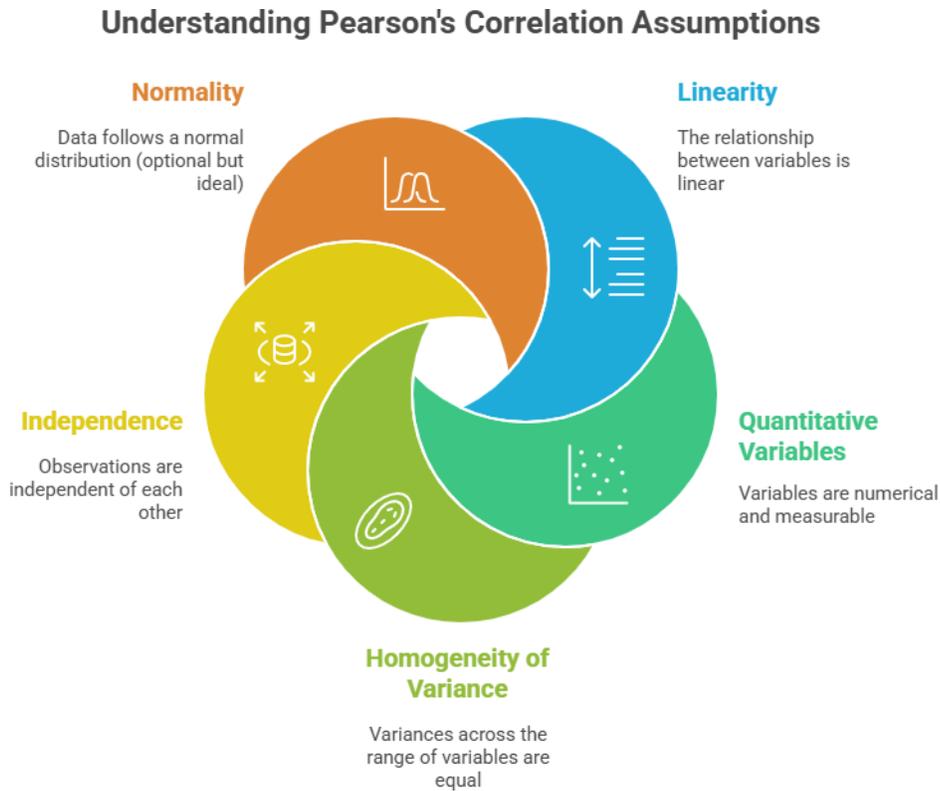
Or using raw scores:

$$r = \frac{[n\Sigma xy - (\Sigma x)(\Sigma y)]}{\sqrt{\{[n\Sigma x^2 - (\Sigma x)^2] \times [n\Sigma y^2 - (\Sigma y)^2]\}}}$$

Where:

- x and y are individual data points,
- $\bar{x}$  and  $\bar{y}$  are the means of X and Y respectively,
- n is the number of observations.

## 8.4.2 Assumptions of Pearson's Correlation



**Fig.8.2. Assumptions of Pearson's Correlation**

To ensure that Pearson's correlation gives valid and reliable results, the following **assumptions** must hold:

1. **Linearity:** The relationship between the two variables should be linear.
2. **Quantitative Variables:** Both variables must be numerical (interval or ratio scale).
3. **Homogeneity of Variance:** The spread of values should be roughly constant across the range.
4. **Independence:** Each pair of observations should be independent.
5. **Normality (optional but ideal):** For inference, both variables should be approximately normally distributed.

## 8.4.3 Computation of Pearson's Correlation

### A. For Ungrouped Data

Use the raw score formula:

**Step-by-Step:**

1. Calculate the sum of x, y, x<sup>2</sup>, y<sup>2</sup>, and xy.
2. Substitute into the formula:

$$r = \frac{[n\sum xy - (\sum x)(\sum y)]}{\sqrt{\{[n\sum x^2 - (\sum x)^2] \times [n\sum y^2 - (\sum y)^2]\}}}$$

**Example:**

Suppose we have scores of 5 students in Math (X): 40, 50, 60, 70, 80

and in Science (Y): 42, 49, 65, 68, 75

Compute r using the formula above (values plugged in during actual computation).

**B. For Grouped Data**

When data is given in class intervals, use the formula:

$$r = \frac{\sum f(dx \times dy)}{\sqrt{[\sum f(dx)^2 \times \sum f(dy)^2]}}$$

Where:

- dx = deviation of X from assumed mean (A) divided by class width,
- dy = deviation of Y from assumed mean (B),
- f = frequency of class.

**Steps:**

1. Calculate midpoints, dx and dy,
2. Multiply dx and dy by f,
3. Compute all required sums,
4. Substitute into the formula.

**Did You Know?**

“If all data points lie exactly on a straight line with a **positive slope**, Karl Pearson’s coefficient **r becomes +1**, indicating a **perfect positive correlation**. However, such perfect correlation is **extremely rare in real-world business data**, due to noise and variation in measurements.”

**8.4.4 Merits and Limitations of Pearson’s Correlation**

**Merits:**

1. **Quantifies the Relationship:** Provides an exact numerical value showing strength and direction.
2. **Mathematically Rigid:** Based on an exact formula—results are consistent and replicable.
3. **Widely Applicable:** Suitable for a variety of fields—economics, psychology, business, etc.
4. **Foundation for Further Analysis:** Serves as a basis for regression, prediction, and hypothesis testing.

**Limitations:**

1. **Assumes Linearity:** Cannot detect or accurately describe non-linear relationships.
2. **Affected by Outliers:** Extreme values can significantly distort the result.
3. **Only Measures Association, Not Causation:** A high correlation does not imply one variable causes the other.
4. **Requires Interval or Ratio Data:** Not suitable for ordinal or nominal variables.
5. **Sensitive to Scale:** Changing units of measurement affects the covariance but not  $r$  (though it can affect interpretation).

## 8.5 Properties of Coefficient of Correlation

The **correlation coefficient ( $r$ )**, most commonly measured using **Karl Pearson’s method**, is a **statistical index** used to describe the **strength and direction** of the linear relationship between two variables. While easy to compute, interpreting and applying it accurately requires understanding its key properties and limitations.

### 8.5.1 Range of Correlation Coefficient

The value of the Pearson correlation coefficient  $r$  always lies in the range:

$$-1 \leq r \leq +1$$

- $r = +1$ : Perfect **positive** linear correlation
- $r = -1$ : Perfect **negative** linear correlation
- $r = 0$ : **No** linear correlation

These boundaries are important for identifying the extent to which two variables are linearly associated.

### 8.5.2 Interpretation of Values of ‘ $r$ ’

The **magnitude and sign** of  $r$  indicate the **strength and direction** of the relationship.

Value of $r$	Interpretation
--------------	----------------

+0.90 to +1.00	Very strong positive correlation
+0.70 to +0.89	Strong positive correlation
+0.40 to +0.69	Moderate positive correlation
+0.10 to +0.39	Weak positive correlation
0	No linear correlation
-0.10 to -0.39	Weak negative correlation
-0.40 to -0.69	Moderate negative correlation
-0.70 to -0.89	Strong negative correlation
-0.90 to -1.00	Very strong negative correlation

**Note:** These categories are guidelines. The interpretation also depends on the **context** (e.g., in social sciences, even  $r = 0.3$  might be meaningful).

### 8.5.3 Mathematical Properties of ‘r’

#### 1. Symmetry Property:

Correlation is **symmetric** in nature:

$$r(x, y) = r(y, x)$$

This means the correlation between X and Y is the same as between Y and X.

#### 2. Unit-Free Measure:

The coefficient **r is independent of the unit of measurement**, because it is based on standardized values.

Changing from kilograms to grams or rupees to dollars does not affect the value of **r**.

#### 3. Covariance-Based:

**r** is derived from covariance and standard deviations:

$$r = \text{Cov}(X, Y) \div (\sigma_x \times \sigma_y)$$

Where  $\text{Cov}(X, Y)$  is the covariance of variables X and Y.

#### 4. Linear Transformation Invariance:

If you apply a linear transformation to variables (e.g.,  $X' = aX + b$ ), the value of **r** remains unchanged, as long as **a > 0**.

#### 5. No Directional Causality:

The coefficient **r** only reflects association, not direction or causality.

### 8.5.4 Common Misinterpretations

Despite its simplicity, the correlation coefficient is often misunderstood or misused. Here are common pitfalls:

1. **Correlation Implies Causation (False Assumption):**

Just because two variables are correlated does **not** mean that one causes the other.

*Example:* Ice cream sales and drowning rates may be positively correlated due to hot weather—not because one causes the other.

2. **Zero Correlation Means No Relationship (Incorrect in Non-linear Cases):**

$r = 0$  only implies **no linear** relationship. A strong **non-linear relationship** may still exist.

3. **Ignoring Outliers:**

A few extreme values can significantly **distort** the correlation coefficient, giving a misleading picture of the relationship.

4. **Assuming Linearity Always Holds:**

Pearson's  $r$  assumes a linear relationship. Using it with **non-linear** data gives invalid results.

5. **Correlation vs. Agreement:**

High correlation does **not** necessarily mean **good agreement**. For example, two measuring devices can have high  $r$  but different average readings.

## 8.6 Spearman's Rank Correlation

**Spearman's Rank Correlation Coefficient** is a **non-parametric measure** of the strength and direction of association between **two ranked variables**. It is especially useful when the data are **ordinal** or when Pearson's correlation assumptions (e.g., linearity, normality) are not met.

### 8.6.1 Concept and Formula

**Concept:**

Spearman's Rank Correlation (denoted by  $\rho$  or  $r_s$ ) measures how well the relationship between two variables can be described using a **monotonic function** (a consistently increasing or decreasing relationship, not necessarily linear).

It is based on the **rankings** (positions) of the data values, not the actual values.

**Formula (when there are no tied ranks):**

$$r_s = 1 - [6 \times \Sigma d^2 \div n(n^2 - 1)]$$

Where:

- $r_s$  = Spearman's rank correlation coefficient
- $d$  = difference between the ranks of each pair
- $n$  = number of observations

**Step-by-step:**

1. Rank the values of each variable separately (smallest = rank 1).
2. Compute the difference  $d$  between ranks.
3. Square each  $d$  to get  $d^2$ .
4. Sum up all  $d^2$  values.
5. Plug into the formula.

**Result range:**

- **+1**: Perfect positive rank correlation
- **-1**: Perfect negative rank correlation
- **0**: No correlation in ranks

**“Activity: Spearman’s Rank Correlation in Customer Preferences”**

**Instruction to the Student:**

Imagine you surveyed 8 customers to rank five product features (Price, Quality, Durability, Packaging, Brand) in order of preference.

- Collect rankings from two customers.
- Assign ranks for each feature and compute  $d$  and  $d^2$ .
- Apply the Spearman’s Rank Correlation formula:

$$r_s = 1 - [6\sum d^2 \div n(n^2 - 1)]$$

- Interpret the result: Does one customer’s preferences align with the other’s?

Prepare a brief note summarizing your correlation value and explaining what it reveals about customer behavior similarity.

### 8.6.2 Tied Ranks and Adjustments

**Tied ranks** occur when two or more items have the **same value** in a dataset. Instead of assigning separate ranks, they are given the **average of the ranks they would occupy**.

**Example:**

If three students tie for 2nd place, we take average of ranks 2, 3, and 4:

$$\text{Rank assigned} = (2 + 3 + 4) \div 3 = 3$$

**Adjustment in formula:**

The standard formula still applies, but you must ensure ranks are **adjusted for ties** before calculating **d** and **d<sup>2</sup>**.

**Did You Know?**

“In Spearman’s rank correlation, **tied ranks** are not an error—they are a natural part of ordinal data. The formula remains valid if you assign the **average of the ranks** for tied observations. Ignoring ties or assigning them arbitrarily can significantly distort the correlation result.”

**8.6.3 Advantages of Rank Correlation****1. No Assumption of Normality:**

Works well even when the data is not normally distributed.

**2. Handles Non-linear Relationships:**

Suitable when the relationship is monotonic but not linear.

**3. Ordinal Data Friendly:**

Can be used when data is in **ranks or preferences** rather than numerical form.

**4. Resistant to Outliers:**

Since it uses ranks, extreme values have less influence compared to Pearson’s *r*.

**5. Easy to Compute for Small Data Sets:**

Especially suitable for competitions, surveys, or preference rankings.

**8.6.4 Limitations of Rank Correlation****1. Less Precise with Numerical Data:**

When actual numerical values are available and assumptions are met, Pearson’s correlation provides more precise measurement.

**2. Inefficient for Large Tied Data Sets:**

Too many tied ranks can distort the correlation strength or require complex adjustments.

**3. Cannot Detect Linear Strength:**

It detects monotonic trends but doesn't measure **how strong the linear relationship** is.

**4. Not Suitable for Interval/Ratio Data Analysis:**

For statistical modeling or predictive analysis, rank correlation is less informative than numerical correlation.

**Summary Table: Spearman's vs. Pearson's**

Feature	Pearson's r	Spearman's r <sub>s</sub>
Data Type	Interval/Ratio	Ordinal or Ranked
Assumes Normality?	Yes	No
Handles Non-linear Trends?	No	Yes (monotonic)
Impact of Outliers	High	Low
Measures	Linear correlation	Rank correlation

**Knowledge Check 1**

**Choose the correct option:**

- If the value of Karl Pearson's correlation coefficient (r) is  $-0.95$ , it indicates:
  - No correlation
  - Weak positive correlation
  - Strong negative correlation
  - Perfect correlation
- A scatter diagram shows points sloping downward from left to right. What type of correlation is indicated?
  - No correlation
  - Positive correlation
  - Perfect correlation
  - Negative correlation
- Which of the following is **true** about Spearman's Rank Correlation?
  - It assumes a linear relationship
  - It cannot handle tied ranks

- C) It is based on data values, not ranks
  - D) It is useful for ordinal or preference data
4. The correlation coefficient is always:
- A) Greater than or equal to 0
  - B) Between  $-1$  and  $+1$
  - C) Equal to or more than 1
  - D) Less than  $-1$
5. Which of the following statements is **incorrect**?
- A) Correlation measures association, not causation
  - B) A high correlation always implies a strong cause-effect relationship
  - C) A zero correlation means no linear relationship
  - D) Karl Pearson's correlation is sensitive to outliers

## 8.7 Summary

- ❖ In this chapter, we explored the concept of **correlation**, which measures the degree and direction of relationship between two variables.
  - We began with the **basic meaning of correlation** and its distinction from causation.
  - Different **types of correlation** were discussed: positive vs. negative, linear vs. non-linear, and simple, partial, and multiple correlation.
  - Through **scatter diagrams** and **simple graphs**, we learned how to visually interpret the presence and strength of relationships.
  - The **Karl Pearson's Coefficient of Correlation ( $r$ )** was introduced as a numerical measure of linear correlation, along with its formula, assumptions, and computation methods.
  - We studied the **properties and limitations of Pearson's  $r$** , including common misinterpretations.
  - Finally, we explored **Spearman's Rank Correlation**, a non-parametric measure useful for ordinal or ranked data, and discussed how to handle tied ranks.
- ❖ Together, these tools form a strong foundation for understanding relationships in statistical and business data, which is crucial for forecasting, decision-making, and pattern recognition.

## 8.8 Key Terms

1. **Correlation** - A statistical measure showing how two variables move in relation to each other
2. **Positive Correlation** - Both variables increase or decrease together
3. **Negative Correlation** - One variable increases as the other decreases
4. **Linear Correlation** - A relationship where change in one variable is proportional to the other
5. **Karl Pearson's Coefficient (r)** - A numerical measure of linear correlation, ranging from  $-1$  to  $+1$
6. **Spearman's Rank Correlation** - A measure of monotonic relationships based on the ranking of data
7. **Scatter Diagram** - A graphical method to show the relationship between two variables
8. **Tied Ranks** - Equal values assigned an average rank in ranking-based methods
9. **Partial Correlation** - Relationship between two variables after removing the effect of a third
10. **Multiple Correlation** - Relationship between one dependent and two or more independent variables

## 8.9 Descriptive Questions

1. Define correlation and explain its significance in business and economics.
2. Differentiate between positive and negative correlation with suitable examples.
3. What are the key differences between linear and non-linear correlation?
4. Explain the steps involved in constructing a scatter diagram. What can it tell us?
5. Write down the formula for Karl Pearson's coefficient of correlation and explain each term.
6. What assumptions must be satisfied to apply Pearson's correlation method?
7. Discuss the merits and limitations of Karl Pearson's method of correlation.
8. Explain the concept of rank correlation. How is it useful in the absence of quantitative data?
9. How do you deal with tied ranks in Spearman's rank correlation?
10. Compare Karl Pearson's correlation and Spearman's rank correlation. When is each method appropriate?

## 8.10 References

1. Gupta, S. P. (2014). *Statistical Methods*. Sultan Chand & Sons.
2. Sharma, J. K. (2018). *Business Statistics*. Vikas Publishing House.
3. Levin, R. I., & Rubin, D. S. (2013). *Statistics for Management*. Pearson Education.
4. Anderson, D. R., Sweeney, D. J., & Williams, T. A. (2016). *Statistics for Business and Economics*. Cengage Learning.
5. UGC e-Pathshala Modules – *Statistical Analysis*

6. MoSPI (Government of India) – *National Statistics Handbook*

**Answers to Knowledge Check**

***Knowledge Check 1***

1. C) Strong negative correlation
2. D) Negative correlation
3. D) It is useful for ordinal or preference data
4. B) Between  $-1$  and  $+1$
5. B) A high correlation always implies a strong cause-effect relationship

## 8.11 Case Study

### “Analyzing Sales and Advertising Spend: A Correlation Study”

#### Introduction

In today's competitive business environment, companies are constantly trying to understand the factors that influence their sales performance. One of the most commonly assumed relationships is between **advertising expenditure** and **sales revenue**. Marketing departments often argue for higher budgets, citing a direct link between the two. However, to make informed decisions, business managers must rely on data rather than assumptions.

This caselet explores how correlation analysis—specifically **Karl Pearson's and Spearman's methods**—can help evaluate the **strength and direction** of relationships between advertising spend and sales. It also touches upon visual tools like **scatter diagrams** and explains how misinterpretation of correlation can lead to misleading business strategies.

#### Background

Imagine a mid-sized retail company that operates across five metro cities. The marketing team increased its ad budget over the past six months with the expectation that sales would rise accordingly. However, while sales did improve in some cities, others showed no significant change despite higher advertising costs.

The regional manager, Priya, decided to dig deeper. She gathered monthly data from all five cities for:

- Amount spent on digital ads (₹)
- Monthly sales revenue (₹)

She plotted the data using a **scatter diagram**, which showed an upward trend for three cities but a scattered or downward pattern for the others. This prompted her to compute the **Karl Pearson's coefficient of correlation (r)** to measure the linear relationship between the two variables.

The results:

- City A:  $r = +0.87$  (strong positive)
- City B:  $r = +0.65$  (moderate positive)
- City C:  $r = -0.10$  (very weak negative)
- City D:  $r = +0.02$  (no linear correlation)
- City E:  $r = +0.90$  (strong positive)

She then used **Spearman's Rank Correlation** for customer satisfaction rankings and repeat purchases to see if ranking variables followed similar patterns.

### **Problem Statement 1: Misinterpreting Zero Correlation**

Many departments assumed that a correlation close to zero meant **no relationship at all**. In City D, the ad spend had no linear pattern with sales, but customer footfall increased significantly.

#### **Solution:**

The team was educated on the **difference between zero linear correlation and no relationship**. A scatter plot helped show that a **non-linear relationship** might still exist. They considered alternative models like polynomial regression and focused more on engagement metrics.

#### **MCQ:**

What does a correlation coefficient ( $r$ ) of 0 imply?

- A) There is no relationship of any kind
- B) There is no **linear** relationship
- C) The variables are independent
- D) A causal relationship exists

**Answer:** B) There is no **linear** relationship

### **Problem Statement 2: Handling Tied Ranks in Spearman's Method**

When customer feedback was collected using star ratings, many customers gave identical ratings, leading to **tied ranks** in the dataset.

#### **Solution:**

The data analyst applied **average ranks** for tied observations before computing Spearman's coefficient. This allowed the analysis to reflect monotonic trends in satisfaction and repeat purchase behavior.

#### **MCQ:**

How are tied ranks handled in Spearman's rank correlation?

- A) Ties are ignored
- B) Assign the lowest rank
- C) Assign the average of tied ranks
- D) Assign the highest rank

**Answer:** C) Assign the average of tied ranks

### **Problem Statement 3: Assuming Correlation Means Causation**

Some teams used a high correlation to justify marketing campaigns, assuming **ads directly caused sales to increase**.

#### **Solution:**

Priya trained teams on the **distinction between correlation and causation**, showing examples where a third factor (like seasonal demand or competitor discounting) influenced both ad spend and sales.

#### **MCQ:**

Which of the following is **true** about correlation?

- A) Correlation always implies causation
- B) Correlation shows the strength and direction of relationship
- C) A high  $r$  value proves one variable causes the other
- D) Zero correlation proves variables are unrelated in all forms

**Answer:** B) Correlation shows the strength and direction of relationship

### **Conclusion**

Through statistical tools such as Pearson's and Spearman's correlation coefficients, the company gained clarity on where ad spending was effective and where other factors were at play. The case also highlighted the importance of **visual tools** like scatter diagrams and the **need for statistical literacy** to interpret data correctly. By understanding the **true nature of correlation**, businesses can make better strategic decisions and avoid misleading conclusions.

## Unit 9: Regression

### Learning Objectives

1. **Define regression** and understand its importance in predicting the value of one variable based on another.
2. **Differentiate between correlation and regression**, especially in terms of direction and application.
3. **Identify various types of regression** such as simple, multiple, linear, and non-linear regression, and recognize their appropriate use cases.
4. **Understand and apply the algebraic methods** used to study regression, including the least squares method.
5. **Construct regression lines** (Y on X and X on Y) and interpret their significance in a real-world dataset.
6. **Calculate regression coefficients** and understand their meaning in terms of the rate of change and the direction of relationship.
7. **Explain the mathematical properties of regression lines**, such as how they intersect and how they relate to correlation coefficients.
8. **Apply regression analysis to real-life business and economics problems**, such as sales forecasting, cost prediction, and performance analysis.

### Content

- 9.0 Introductory Caselet
- 9.1 Introduction
- 9.2 Types of Regression
- 9.3 Methods of Studying Regression
- 9.4 Lines of Regression
- 9.5 Regression Coefficients
- 9.6 Properties of Lines of Regression (Linear Regression)
- 9.7 Summary
- 9.8 Key Terms
- 9.9 Descriptive Questions
- 9.10 References
- 9.11 Case Study

## 9.0 Introductory Caselet

### “Mira’s Marketing Spend: Forecasting ROI with Regression”

#### Background:

Mira is the marketing head of a fast-growing online grocery startup operating in three metro cities. Her team regularly invests in different advertising channels—Google Ads, Instagram promotions, and local influencer campaigns. Despite consistent marketing budgets, Mira observed that the **return on investment (ROI)** varied significantly each month. Some months saw a steep rise in sales, while others barely broke even.

To address this inconsistency, Mira consulted a data analyst who introduced her to **regression analysis**. Unlike correlation, which only measures the strength of a relationship, regression allowed her to **quantify the influence** of each advertising platform on total monthly sales and make predictions for future months.

Over the past 10 months, Mira's team compiled the following data:

- Amount spent on Google Ads (₹),
- Number of influencer posts,
- Engagement rate on Instagram,
- Total monthly sales (₹).

The analyst used this dataset to build a **multiple linear regression model**, where **sales** was the dependent variable and the other factors were independent predictors. The resulting equation enabled Mira to estimate how changes in ad strategy could impact sales outcomes.

Using the regression model, she found:

- **Google Ads** had the strongest positive impact on monthly sales,
- **Influencer posts** were effective, but only above a certain threshold,
- **Instagram engagement** had a weaker, but still positive, correlation with sales.

Now, Mira can:

- Forecast sales based on planned ad spend,

- Allocate budgets across channels more efficiently,
- Justify marketing decisions with data.

What was once a **trial-and-error marketing strategy** became a **predictive, data-driven planning process**, all because Mira applied regression analysis to her business data.

**Critical Thinking Question:**

If you were in Mira's position and had to choose **one variable to cut from the model** to reduce campaign costs, how would you use regression coefficients and  $R^2$  values to guide your decision? What would you watch out for when removing a variable from a multiple regression model?

## 9.1 Introduction

**Regression analysis** is a powerful statistical technique used to examine the relationship between two or more variables. While correlation measures the **degree of association**, regression goes a step further to **predict or estimate the value** of one variable based on another. This makes regression particularly valuable in forecasting, decision-making, and business modeling.

### 9.1.1 Meaning and Importance of Regression

#### Meaning:

Regression refers to a statistical process for estimating the relationships among variables. In its simplest form, **regression analysis studies the effect of one independent variable (X) on a dependent variable (Y).**

It helps answer questions like:

- If advertising expense increases, by how much can we expect sales to increase?
- How does education level affect income?
- What is the expected cost of production based on output levels?

#### Importance:

1. **Predictive Power:** Regression helps predict future outcomes based on past data.
2. **Quantitative Relationship:** It provides an equation to estimate one variable using another.
3. **Business Decision-Making:** Supports pricing, budgeting, forecasting, and risk assessment.
4. **Understanding Dependencies:** Helps in identifying cause-effect-like relationships in controlled scenarios.

### 9.1.2 Distinction between Correlation and Regression

Though both correlation and regression deal with relationships between variables, they are fundamentally different in purpose and interpretation.

Feature	Correlation	Regression
Purpose	Measures <b>degree of association</b>	<b>Estimates or predicts</b> one variable based on another
Symmetry	$r(x, y) = r(y, x)$	Regression of Y on X $\neq$ Regression of X on Y
Direction	No direction implied	One variable is <b>dependent</b> , the other is <b>independent</b>
Output	One single value (correlation coefficient)	Regression <b>equation or line</b>
Use Case	Understanding strength of association	Forecasting and cause-effect modeling

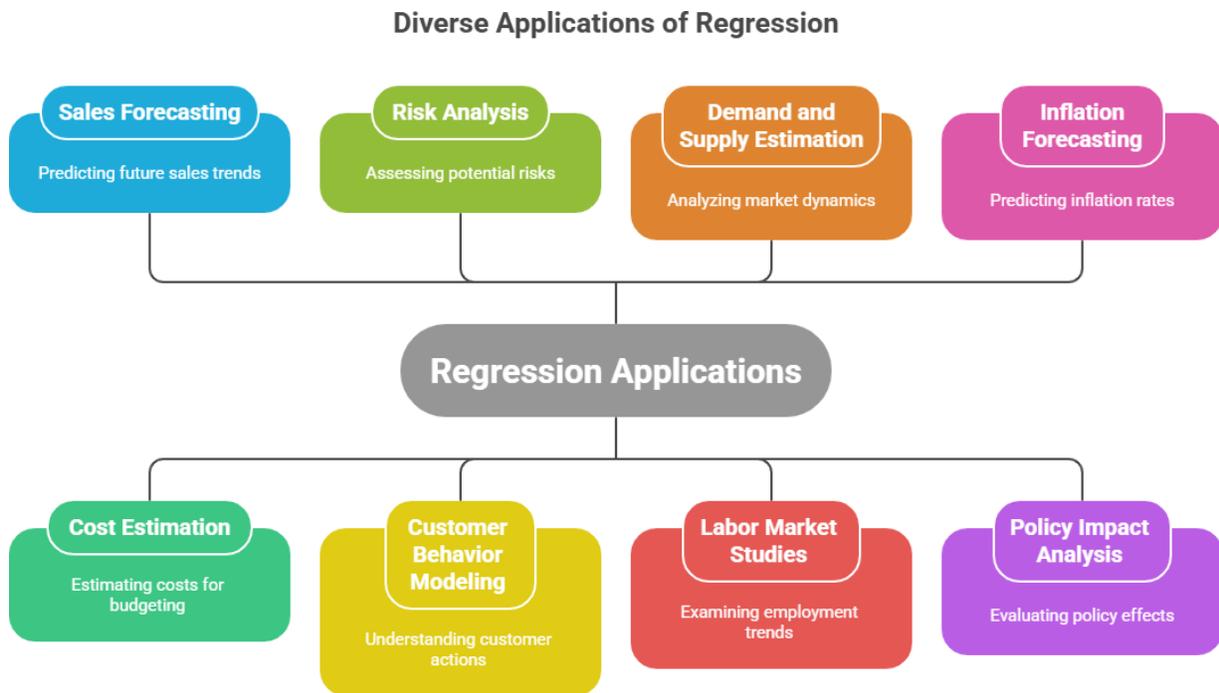
#### Example:

- **Correlation** tells us that advertising and sales are related.
- **Regression** gives us the equation:  $Sales = a + b \times Advertising\ Spend$ .

**Did You Know?**

“Correlation tells you **how strong** the relationship is between two variables, but **regression tells you how much** one variable changes with another. In other words, correlation is about **association**, while regression is about **prediction**.”

### 9.1.3 Applications of Regression in Business and Economics



**Fig.9.1. Applications of Regression in Business and Economics**

Regression is widely used in business analytics, economics, and management sciences to support **data-driven decisions**.

#### **Applications in Business:**

- **Sales Forecasting:** Predicting future sales based on past sales and marketing spend.
- **Cost Estimation:** Estimating production costs based on output levels.
- **Risk Analysis:** Understanding how risk factors influence financial returns.
- **Customer Behavior Modeling:** Predicting purchasing decisions using demographics and income.

#### **Applications in Economics:**

- **Demand and Supply Estimation:** Analyzing how price changes affect demand.
- **Labor Market Studies:** Examining how education and skills affect wages.
- **Inflation Forecasting:** Predicting future inflation based on interest rates and currency values.
- **Policy Impact Analysis:** Assessing the economic impact of taxation or subsidies

## **9.2 Types of Regression**

Regression analysis can take different forms depending on the **number of variables involved** and the **nature of the relationship** between them. The major types include **simple, multiple, linear, and non-linear regression**. Understanding these types helps in selecting the most suitable model for a given data set or problem.

### **9.2.1 Simple Regression**

#### **Definition:**

Simple regression (or **bivariate regression**) involves **two variables**: one **independent variable (X)** and one **dependent variable (Y)**. The objective is to predict the value of Y using X.

#### **Equation:**

$$Y = a + bX$$

Where:

- **Y** = Dependent variable
- **X** = Independent variable
- **a** = Intercept
- **b** = Regression coefficient (slope)

**Example:**

Predicting a student's exam score (Y) based on the number of study hours (X).

**Use Case:**

When there is only **one influencing factor** for prediction, such as predicting profit based on sales.

**“Activity: Forecasting Sales Using Simple Regression”**

**Title:** *"Predicting Next Month's Sales from Advertising Spend"*

**Instruction to the student:**

You are provided with monthly data for the past 8 months, including total advertising spend (₹ in lakhs) and corresponding sales revenue (₹ in lakhs).

- Use the **simple linear regression formula** to compute the regression line:  
 $Y = a + bX$ , where Y = Sales and X = Ad Spend.
- Calculate the values of **a** and **b** using the least squares method.
- Based on your model, predict the sales revenue if next month's ad spend is ₹5.5 lakhs.
- Submit your regression equation and prediction in a short report.

**9.2.2 Multiple Regression****Definition:**

Multiple regression involves **more than one independent variable** used to predict the **value of a single dependent variable**. It is suitable when multiple factors influence an outcome.

**Equation:**

$$Y = a + b_1X_1 + b_2X_2 + \dots + b_nX_n$$

Where:

- $X_1, X_2, \dots, X_n$  = Multiple independent variables
- $b_1, b_2, \dots, b_n$  = Respective regression coefficients

**Example:**

Predicting house prices (Y) based on variables like size ( $X_1$ ), location rating ( $X_2$ ), and number of bedrooms ( $X_3$ ).

**Use**

Used in marketing (e.g., predicting sales using price, promotion, and advertising spend), HR (predicting employee performance), or finance (predicting risk).

**Case:**

### 9.2.3 Linear Regression

**Definition:**

Linear regression assumes that the relationship between the variables can be **represented by a straight line**. It applies to both **simple** and **multiple** regression as long as the pattern of relationship is linear.

**Equation (Simple Linear):**

$$Y = a + bX$$

**Equation (Multiple Linear):**

$$Y = a + b_1X_1 + b_2X_2 + \dots + b_nX_n$$

**Graph:** A straight line in two dimensions (for simple linear) or a plane/hyperplane in higher dimensions.

**Use**

Useful when data shows a **consistent rate of change**, e.g., predicting monthly utility bills based on usage.

**Case:**

### 9.2.4 Non-linear Regression

**Definition:**

Non-linear regression models situations where the relationship between dependent and independent variables **cannot be adequately captured by a straight line**. The regression curve could be exponential, logarithmic, polynomial, etc.

**Example Equations:**

- $Y = a \times e^{(bX)}$  (Exponential)
- $Y = a + bX + cX^2$  (Quadratic)

**Example:**

Modeling the learning curve of a new employee, where performance improves rapidly at first and then slows down.

**Use**

Applicable in fields like biology (growth models), economics (diminishing returns), and marketing (response to increasing advertising).

**Case:**

**Summary Table: Comparison of Regression Types**

Type of Regression	No. of Independent Variables	Relationship Form	Common Use Case
Simple Regression	One	Linear	Sales prediction from marketing spend
Multiple Regression	Two or more	Linear	House price prediction using multiple factors
Linear Regression	One or more	Linear	Predictive analytics in most domains
Non-linear Regression	One or more	Curved/Complex	Customer behavior modeling, biological growth

### 9.3 Methods of Studying Regression

Regression analysis involves determining the **best-fit relationship** between two or more variables. The methods used to study regression are based on either **visual interpretation** or **mathematical computation**. This section introduces the two most commonly used methods: the **scatter diagram method** and the **method of least squares**, along with a discussion of the key **assumptions** behind regression analysis.

#### 9.3.1 Scatter Diagram Method

**Definition:**

A **scatter diagram** is a **graphical method** used to visually examine the relationship between two variables. It is the simplest way to understand whether a **linear or non-linear trend** exists.

**Steps to Draw:**

1. Plot the **independent variable (X)** on the horizontal axis and the **dependent variable (Y)** on the vertical axis.
2. For each pair of values, plot a point (x, y).
3. Observe the overall pattern formed by the points.

**Interpretation:**

- **Upward trend** → Positive correlation
- **Downward trend** → Negative correlation
- **No clear pattern** → No correlation

- **Points form a straight line** → Strong linear correlation
- **Points curve or cluster** → Possible non-linear correlation

#### Use Case:

Used as a **preliminary tool** to decide whether regression analysis is appropriate and what type of model might fit the data.

### 9.3.2 Method of Least Squares

#### Definition:

The **method of least squares** is a **mathematical technique** used to find the **best-fitting regression line** through a set of data points by **minimizing the sum of the squares of the vertical deviations** (errors) from the actual data points to the line.

#### Objective:

To find a line of the form:

$$Y = a + bX$$

Where:

- **a** = Y-intercept
- **b** = Slope of the regression line
- The line minimizes the value of:  $\Sigma(Y - \hat{Y})^2$ , where  $\hat{Y}$  = predicted value

#### Formulas:

- $b = \frac{[n\Sigma XY - (\Sigma X)(\Sigma Y)]}{[n\Sigma X^2 - (\Sigma X)^2]}$
- $a = \bar{Y} - b\bar{X}$

#### Steps:

1. Compute  $\Sigma X$ ,  $\Sigma Y$ ,  $\Sigma XY$ , and  $\Sigma X^2$ .
2. Calculate b (slope).
3. Calculate a (intercept).
4. Construct the regression equation.

#### Use Case:

This method is widely used in economics, business forecasting, production planning, and any context that requires **predictive modeling** based on past data.

### 9.3.3 Assumptions Underlying Regression Analysis

Regression analysis is built on several assumptions that must be met for the results to be valid and reliable.

Assumption	Explanation
<b>Linearity</b>	The relationship between the dependent and independent variable is linear.
<b>Independence</b>	Observations are independent of one another.
<b>Homoscedasticity</b>	The variance of errors (residuals) is constant across all levels of X.
<b>Normality of Errors</b>	The residuals (errors) are normally distributed.
<b>No Multicollinearity</b>	In multiple regression, independent variables are not highly correlated.

#### Importance of Assumptions:

- Violations can lead to **biased or misleading estimates**, such as **underestimating error** or **overstating significance**.
- Assumptions should be checked using **diagnostic tools** like residual plots and statistical test

## 9.4 Lines of Regression

In regression analysis, the **line of regression** is the line that best fits a set of data points on a scatter plot. There are two lines:

- One that estimates **Y for a given X**, called the **regression line of Y on X**
- One that estimates **X for a given Y**, called the **regression line of X on Y**

Each line minimizes the sum of squared deviations **in the dependent variable** it is trying to predict.

### 9.4.1 Regression Line of X on Y

#### Purpose:

This line is used to **estimate X** (independent variable) based on known values of Y (dependent variable).

#### Equation format:

$$X = a + bY$$

Where:

- **X** = Estimated value of the independent variable
- **Y** = Known value of the dependent variable
- **a** = Intercept
- **b** = Regression coefficient of X on Y

**Formula for slope (b):**

$$b_{xy} = r \times (\sigma_x \div \sigma_y)$$

Here:

- $r$  = Pearson correlation coefficient
- $\sigma_x$  = Standard deviation of X
- $\sigma_y$  = Standard deviation of Y

**Interpretation:**

This line is rarely used unless there is a special interest in predicting X from Y (for example, estimating advertising budget based on sales target).

**9.4.2 Regression Line of Y on X****Purpose:**

This is the more commonly used line, used to **estimate Y** (dependent variable) for given values of X.

**Equation format:**

$$Y = a + bX$$

Where:

- $Y$  = Estimated value of the dependent variable
- $X$  = Known value of the independent variable
- $a$  = Intercept
- $b$  = Regression coefficient of Y on X

**Formula for slope (b):**

$$b_{yx} = r \times (\sigma_y \div \sigma_x)$$

**Interpretation:**

This line is used when the objective is **prediction**, for example, predicting sales based on advertising expenditure, or forecasting expenses based on output.

**9.4.3 Properties of Regression Lines****1. Two Lines Exist:**

There are two separate regression lines unless the correlation is perfect ( $r = \pm 1$ ), in which case both lines **coincide**.

**2. Intersection at Mean Point:**

The two regression lines **always intersect at the point**  $(\bar{X}, \bar{Y})$ , which is the **mean** of X and Y.

**3. Relation with Correlation Coefficient:**

The product of the two regression coefficients equals the square of the correlation coefficient:

$$b_{yx} \times b_{xy} = r^2$$

**4. Directionality Matters:**

The regression line of Y on X is not the same as the regression line of X on Y. They serve **different purposes**.

**5. Least Squares Criterion:**

Each line minimizes the **sum of squared deviations** in its own predicted variable (X or Y).

**Did You Know?**

“The two regression lines—Y on X and X on Y—**always intersect at the mean point**  $(\bar{X}, \bar{Y})$  of the data. This is true regardless of the direction or strength of the correlation.”

**9.4.4 Geometric Representation of Regression Lines**

On a **scatter diagram**, the two regression lines provide **best-fit representations** of the trend of the data.

**Visual Representation:**

- The **line of Y on X** minimizes vertical deviations  $(Y - \hat{Y})^2$ .
- The **line of X on Y** minimizes horizontal deviations  $(X - \hat{X})^2$ .
- Both lines **intersect at the point**  $(\bar{X}, \bar{Y})$ , representing the **average of both variables**.
- The **angle between the two lines** depends on the strength of correlation:
  - If  $r = \pm 1$ , both lines merge into a **single straight line**.
  - If  $r = 0$ , they intersect **at right angles (90°)**.

**Summary Table: Regression Line Comparison**

Feature	Y on X	X on Y
Dependent Variable	Y	X
Independent Variable	X	Y

Equation Format	$Y = a + bX$	$X = a + bY$
Minimizes Deviation In	Y (vertical distances)	X (horizontal distances)
Slope Formula	$b = r \times (\sigma_y \div \sigma_x)$	$b = r \times (\sigma_x \div \sigma_y)$
Common Use	Forecasting, prediction models	Estimating X from known Y

## 9.5 Regression Coefficients

Regression coefficients are central to understanding the nature of the relationship between variables in regression analysis. These coefficients quantify **how much the dependent variable changes** when the independent variable increases by one unit.

### 9.5.1 Concept and Calculation of Regression Coefficients

#### Concept:

A **regression coefficient** represents the **rate of change** in the dependent variable as a result of a one-unit change in the independent variable, assuming all other factors are constant.

In simple linear regression, we usually calculate:

- $b_{yx}$ : Regression coefficient of **Y on X**
- $b_{xy}$ : Regression coefficient of **X on Y**

These coefficients appear in the equations:

- $Y = a + b_{yx}X$
- $X = a + b_{xy}Y$

#### Formulas:

- $b_{yx} = r \times (\sigma_y \div \sigma_x)$
- $b_{xy} = r \times (\sigma_x \div \sigma_y)$

Where:

- $r$  = Pearson's correlation coefficient
- $\sigma_x, \sigma_y$  = Standard deviations of X and Y respectively

Alternatively, using the raw data:

- $b_{yx} = \Sigma(xy) \div \Sigma(x^2)$
- $b_{xy} = \Sigma(xy) \div \Sigma(y^2)$

### 9.5.2 Relationship between Correlation Coefficient and Regression Coefficients

There is a **mathematical link** between the correlation coefficient ( $r$ ) and the two regression coefficients ( $b_{yx}$  and  $b_{xy}$ ):

$$r = \sqrt{b_{yx} \times b_{xy}}$$

Or:

$$b_{yx} \times b_{xy} = r^2$$

**Interpretation:**

- If both regression coefficients are **positive**, **r is positive**.
- If both regression coefficients are **negative**, **r is negative**.
- If one is positive and one is negative, the relationship is **inconsistent**, which typically shouldn't occur in well-structured data.

### 9.5.3 Properties of Regression Coefficients

#### 1. Two Coefficients:

For any two variables, there are two regression coefficients:  $b_{yx}$  and  $b_{xy}$ .

#### 2. Signs Match Correlation Coefficient:

Both regression coefficients carry the **same sign** as the correlation coefficient  $r$ .

#### 3. Geometric Mean Relationship:

The product of the two coefficients equals the square of the correlation coefficient:

$$b_{yx} \times b_{xy} = r^2$$

#### 4. Unit Sensitivity:

Regression coefficients are **not unit-free**. Their values change when the **units of measurement** of the variables change (e.g., from kg to g).

#### 5. No Fixed Range:

Unlike correlation (which ranges between  $-1$  and  $+1$ ), regression coefficients can have **any real value** (positive or negative).

#### 6. When $r = \pm 1$ :

The two regression lines **coincide**, and both regression coefficients equal the ratio of the standard deviations:

- $b_{yx} = \sigma_y \div \sigma_x$  (for  $r = +1$ )
- $b_{yx} = -(\sigma_y \div \sigma_x)$  (for  $r = -1$ )

### Did You Know?

“The **product of the two regression coefficients** ( $b_{yx} \times b_{xy}$ ) is always equal to the **square of the correlation coefficient** ( $r^2$ ). This provides a mathematical bridge between correlation and regression.”

### “Activity: Identifying Key Predictors Using Multiple Regression”

**Title:** *"Which Factors Drive Revenue? A Data-Driven Approach"*

**Instruction to the student:**

A startup has collected data over 10 weeks, including the following for each week:

- Number of social media posts,
- Average daily website traffic,
- Weekly ad spend,
- Weekly revenue.

You are required to:

1. Run a **multiple linear regression** to predict revenue using the three predictors.
2. Identify the **regression coefficients** for each variable and interpret their meanings.
3. Determine which variable has the **most significant impact** on revenue based on the magnitude of the coefficient.
4. Write a short explanation (100–150 words) advising the company where to focus its marketing efforts for maximum return.

#### 9.5.4 Interpretation of Regression Coefficients

The **slope (b)** in a regression line equation has a practical interpretation:

- $b_{yx} = 2.5$  means:  
"For every 1 unit increase in X, Y increases by 2.5 units on average."
- $b_{yx} = -0.75$  means:  
"For every 1 unit increase in X, Y decreases by 0.75 units on average."

This helps **predict outcomes**, **assess impact**, and **guide decision-making**.

**Real-life**

**Example:**

In marketing, if  $b_{yx} = 5.2$  in the equation  $\text{Sales} = a + 5.2 \times \text{AdSpend}$ , this suggests that for each ₹1,000 increase in advertising spend, sales increase by ₹5,200 on average.

**Summary Table: Regression Coefficient Essentials**

Feature	Value or Explanation
Number per pair of variables	2 (Y on X and X on Y)
Relation to correlation coefficient	$r^2 = b_{yx} \times b_{xy}$
Sign	Same as that of <b>r</b>
Unit-dependence	Yes
Range	No fixed limits
Interpretation	Rate of change in dependent variable

**9.6 Properties of Lines of Regression (Linear Regression)**

Linear regression is based on fitting a **straight line** through a set of data points that best represents the relationship between the independent and dependent variables. This section outlines the key properties of the regression line, the role of least squares, the standard error involved, and its limitations.

**9.6.1 Best Fit Line and Principle of Least Squares**

**Best Fit Line:**

The **best fit line** in linear regression is the line that **minimizes the vertical distances** (errors) between the actual data points and the predicted values on the line. These vertical distances are also known as **residuals**.

**Principle of Least Squares:**

The least squares method determines the best-fitting line by **minimizing the sum of the squares of the residuals** (i.e., differences between observed and estimated values of the dependent variable).

**Objective:**

Minimize:  $\Sigma(Y - \hat{Y})^2$

Where:

- $Y$  = Actual value
- $\hat{Y}$  = Predicted value from regression equation

**Equation form (for Y on X):**

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

- $a$  = intercept
- $b$  = slope (regression coefficient)

This method ensures that **total error is as small as possible**, which makes the resulting regression equation more reliable for prediction.

### 9.6.2 Properties of Linear Regression Line

#### 1. Passes through Mean Point:

The regression line always **passes through the point**  $(\bar{X}, \bar{Y})$ , the means of the X and Y variables.

#### 2. Two Separate Lines:

Unless the correlation is perfect ( $r = \pm 1$ ), the regression line of Y on X **is not the same** as the regression line of X on Y.

#### 3. Minimizes Sum of Squares:

The least squares regression line **minimizes**  $\Sigma(Y - \hat{Y})^2$ , not  $\Sigma(Y - \hat{Y})$  or any other measure of error.

#### 4. Linear in Parameters:

Even if the variables are transformed, the regression equation remains **linear in the parameters** ( $a, b$ ).

#### 5. Slope Determined by Correlation and Standard Deviations:

The slope of the regression line of Y on X is:

$$b = r \times (\sigma_y \div \sigma_x)$$

#### 6. Direction Determined by Sign of r:

If  $r > 0$ , the regression line slopes upward;

If  $r < 0$ , it slopes downward.

### 9.6.3 Errors in Estimation and Standard Error of Regression

#### Errors in Estimation (Residuals):

Each prediction made by the regression line may differ from the actual observed value. This **difference is called the residual**:

$$\text{Residual (e)} = Y - \hat{Y}$$

Residuals are important in evaluating how well the regression model fits the data.

### Standard Error of Estimate ( $S_e$ ):

This measures the **average size of the residuals**. A smaller standard error means the predictions are more accurate.

### Formula:

$$S_e = \sqrt{[\Sigma(Y - \hat{Y})^2 \div n]}$$

Where:

- $Y$  = Actual values
- $\hat{Y}$  = Predicted values from regression
- $n$  = Number of observations

### Interpretation:

- A **low**  $S_e$  indicates that the regression line is a good fit.
- A **high**  $S_e$  suggests the model may not be reliable for prediction.

## 9.6.4 Limitations of Linear Regression

### 1. Assumes Linearity:

The model assumes a straight-line relationship. If the actual relationship is non-linear, the results will be misleading.

### 2. Sensitive to Outliers:

Extreme values can distort the slope and intercept, reducing accuracy.

### 3. Dependent on Assumptions:

Results are valid only if assumptions like normality, homoscedasticity, and independence are satisfied.

### 4. Does Not Prove Causation:

Even with strong correlation and a predictive model, **causality cannot be inferred**.

### 5. Only One Dependent Variable:

Basic linear regression only explains one outcome variable. In real-world cases, multiple variables often influence outcomes simultaneously.

## 6. Limited in Multicollinearity Scenarios:

In multiple regression (beyond simple linear), if independent variables are highly correlated, it can **distort** the regression coefficients.

### Knowledge Check 1

#### Choose the correct option:

- In **simple linear regression**, the dependent variable (Y) is predicted using:
  - Only one independent variable (X)
  - Multiple independent variables
  - The mean of the data
  - The correlation coefficient
- What is the key difference between **correlation** and **regression**?
  - Correlation measures how one variable causes another
  - Regression predicts the value of one variable based on another
  - Regression only applies to nominal data
  - Correlation calculates the slope of a line
- In the regression equation  $Y = a + bX$ , what does the coefficient **b** represent?
  - The predicted value of Y
  - The slope of the regression line (change in Y per unit change in X)
  - The value of X when Y = 0
  - The intercept of the regression line
- The formula for calculating the regression coefficient of Y on X is:
  - $b_{yx} = r \times (\sigma_x \div \sigma_y)$
  - $b_{yx} = (\Sigma X \times \Sigma Y) \div (\Sigma X^2 \times \Sigma Y^2)$
  - $b_{yx} = \Sigma(Y - \hat{Y}) \div \Sigma X$
  - $b_{yx} = (\sigma_x \div \sigma_y) \times r$
- What is the **interpretation** of a regression coefficient of **0.5** for X in a model predicting Y?
  - For every 1 unit increase in Y, X increases by 0.5 units
  - For every 1 unit increase in X, Y increases by 0.5 units

- C) Y has no effect on X
- D) The regression model is not valid

## 9.7 Summary

- ❖ In this chapter, we explored the fundamentals and practical applications of **regression analysis**, a core statistical tool used to understand and model relationships between variables.
  - We began with the **meaning and importance of regression**, followed by its distinction from correlation.
  - We examined **types of regression**, including simple, multiple, linear, and non-linear forms, identifying when each type is appropriate.
  - Through the **scatter diagram method** and the **method of least squares**, we learned how to graphically and algebraically study regression.
  - The focus then shifted to the **lines of regression**—Y on X and X on Y—and how they relate to one another and the mean point  $(\bar{X}, \bar{Y})$ .
  - We calculated **regression coefficients**, understood their interpretation, and explored their relationship with the correlation coefficient.
  - Finally, we studied the **properties of regression lines**, including assumptions, error estimation, and limitations of linear regression models.
- ❖ Regression analysis plays a key role in business, economics, and data science by enabling predictions, trend analysis, and decision-making based on historical data.

## 9.8 Key Terms

1. **Regression** - A statistical method to model the relationship between dependent and independent variables
2. **Simple Regression** - Regression involving one independent and one dependent variable
3. **Multiple Regression** - Regression involving two or more independent variables
4. **Linear Regression** - Regression where the relationship is modeled by a straight line
5. **Non-linear Regression** - Regression where the relationship is curved or non-linear
6. **Least Squares Method** - Technique that minimizes the sum of squared residuals
7. **Regression Coefficient** - Rate of change of the dependent variable with respect to the independent variable
8. **Residual** - The difference between actual and predicted values
9. **Standard Error of Estimate** - A measure of the accuracy of predictions made by the regression model

10. **Best Fit Line** - The regression line that minimizes the error of prediction

## 9.9 Descriptive Questions

1. Define regression and explain its importance in business forecasting.
2. Distinguish between simple and multiple regression with examples.
3. What is the difference between linear and non-linear regression? Give one use case for each.
4. Describe the method of least squares. How is it used to derive the regression line?
5. Explain the regression lines of Y on X and X on Y. When do they coincide?
6. What is the relationship between correlation coefficient and regression coefficients?
7. Discuss any three properties of a linear regression line.
8. How is the standard error of estimate interpreted in regression analysis?
9. Explain any four limitations of linear regression in practical application.
10. What assumptions must hold true for linear regression to be valid?

## 9.10 References

1. Gupta, S. C. (2014). *Fundamentals of Statistics*. Himalaya Publishing House.
2. Sharma, J. K. (2018). *Business Statistics*. Vikas Publishing House.
3. Levin, R. I., & Rubin, D. S. (2017). *Statistics for Management*. Pearson Education.
4. Anderson, D. R., Sweeney, D. J., & Williams, T. A. (2016). *Statistics for Business and Economics*. Cengage Learning.
5. UGC e-Pathshala – *Statistical Methods Modules*
6. National Statistical Office (NSO) Reports, Govt. of India

### Answers to Knowledge Check

#### **Knowledge Check 1**

1. A) Only one independent variable (X)
2. B) Regression predicts the value of one variable based on another
3. B) The slope of the regression line (change in Y per unit change in X)
4. A)  $b_{yx} = r \times (\sigma_x \div \sigma_y)$

5. B) For every 1 unit increase in X, Y increases by 0.5 units

## 9.11 Case Study

### “Predicting Sales Performance: Applying Regression in Retail Forecasting”

#### Introduction

In today’s dynamic retail environment, data-driven decision-making is critical for staying competitive. Businesses often rely on historical data to predict future performance and align their operations accordingly. One such method is **regression analysis**, a statistical technique that helps understand and quantify the relationship between dependent and independent variables.

This case study explores how a mid-sized retail company, **TrendMart**, applied regression analysis to predict monthly sales based on advertising expenditure, foot traffic, and online engagement. Through a systematic approach to analyzing past data, the company was able to refine its marketing strategies, optimize resource allocation, and improve profitability forecasts.

#### Background

TrendMart operates across multiple urban locations and maintains a significant presence online. Over the last year, the management noticed inconsistencies between their advertising investments and actual sales performance. Although certain campaigns led to spikes in foot traffic or social media engagement, the effect on monthly revenue was unpredictable.

To address this, the analytics team proposed the use of **multiple linear regression** to model the relationship between **monthly sales (Y)** and the following predictors:

- **Ad Spend (X<sub>1</sub>)** in ₹,
- **Foot Traffic (X<sub>2</sub>)** measured as average daily store visitors,
- **Social Media Engagement (X<sub>3</sub>)** in terms of likes, shares, and comments.

Using 12 months of historical data, the team conducted regression analysis and derived the following equation:

$$\text{Sales} = 2.5 + 3.1(\text{Ad Spend}) + 2.8(\text{Foot Traffic}) + 0.6(\text{Engagement})$$

The R<sup>2</sup> value of the model was **0.89**, indicating that 89% of the variation in sales could be explained by the three predictors.

#### Problem Statement 1: Misalignment Between Ad Spend and Sales

Despite consistent ad spending, the management found that not all campaigns yielded a proportional increase in sales.

**Solution:**

The regression model revealed that **foot traffic** had a slightly higher influence on sales than ad spend alone. Therefore, instead of increasing ad budgets blindly, TrendMart redirected spending to **local targeting** and **in-store promotions**, which had a higher correlation with store visits, and in turn, with sales.

**Problem Statement 2: Underestimating the Impact of Online Engagement**

Marketing teams previously undervalued social media performance, assuming it had minimal influence on revenue.

**Solution:**

Regression analysis showed that although **online engagement** had the smallest coefficient, it was still a **statistically significant** predictor. Based on this, the digital team began **timing promotions to coincide with social media peaks**, resulting in improved campaign alignment and better conversion rates.

**Problem Statement 3: Forecasting Next Quarter Sales with Limited Data**

TrendMart needed to project sales for the next quarter but had incomplete data for the ongoing month.

**Solution:**

The team used the regression model to plug in the available values for Ad Spend and Foot Traffic, and estimated Engagement based on prior trends. This enabled the creation of a **preliminary forecast**, which could be updated as new data came in.

**MCQ (Knowledge Check)**

1. In the regression equation  $\text{Sales} = 2.5 + 3.1X_1 + 2.8X_2 + 0.6X_3$ , what does 3.1 represent?

- A) Intercept
- B) Error term
- C) Regression coefficient of Ad Spend
- D) Total variance

**Answer:** C) Regression coefficient of Ad Spend

2. What does an  $R^2$  value of 0.89 indicate in a regression model?

- A) The model has low accuracy
- B) The variables are not related
- C) 89% of variation in sales is explained by predictors
- D) The prediction error is 89%

**Answer:** C) 89% of variation in sales is explained by predictors

3. If the regression coefficient of Social Media Engagement is positive but small, what does it imply?

- A) It should be removed
- B) It has no influence
- C) It contributes positively, but less than other variables
- D) It is negatively correlated

**Answer:** C) It contributes positively, but less than other variables

### Conclusion

This case illustrates the practical value of regression analysis in business settings. By quantifying the influence of multiple variables on sales, TrendMart was able to shift from intuition-based decisions to **evidence-based forecasting**. As a result, they improved campaign efficiency, budget allocation, and revenue predictability.

Regression empowered the company to:

- Identify the most impactful predictors of sales,
- Forecast future revenue more accurately,
- Adjust marketing strategies based on data-driven insights.