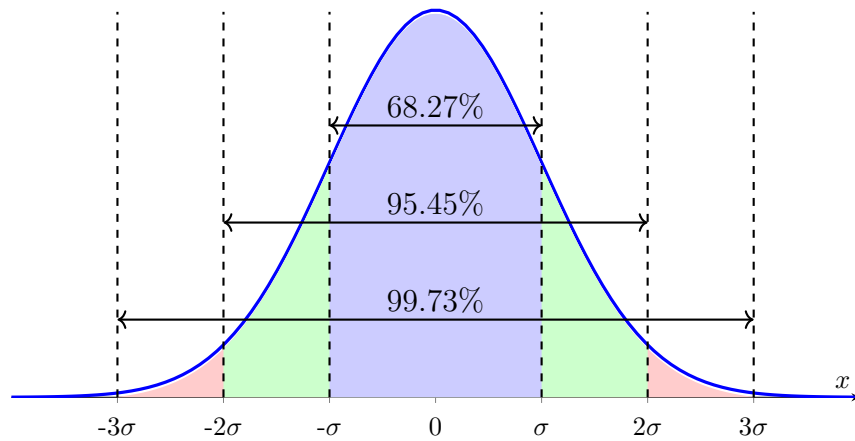


Basics of Probability and Statistics

An In-Depth Exploration of Core Concepts and Methods

Sandip Karar

July 1, 2025



Contents

1	Descriptive Statistics	9
1.1	Introduction	9
1.2	Presentation of Data	11
1.3	Measures of Central Tendency	18
1.4	Partition Values: Quartiles, Deciles, and Percentiles	31
1.5	Measures of Dispersion	32
1.6	Moments, Skewness and Kurtosis	37
2	Theory of Probability	43
2.1	Some Notation and Terminology	43
2.2	Definition of Probability	45
2.3	Axioms of Probability	45
2.4	Conditional Probability	50
2.5	Rule of Total Probability	51
2.6	Bayes' Theorem	52
2.7	Statistical Independence of Events	55
3	Random Variables and Probability Distributions	61
3.1	What is a Random Variable?	61
3.2	Probability Distribution	62
3.3	Mean and Variance of a Random Variable	65
3.4	Joint Distribution of Two random Variables	67
3.5	Conditional Probability Distribution	73
3.6	Functions of a Random Variable	75
3.7	Standardized Random Variable	78
3.8	Chebyshev's Inequality	78
3.9	Moments and Moment Generating Function	79
4	Common Distributions	85
4.1	Bernoulli Distribution	85
4.2	Binomial Distribution	86
4.3	Poisson Distribution	89
4.4	Uniform Distribution	93
4.5	Normal Distribution	95

5	Sampling Theory	107
5.1	Introduction	107
5.2	Sampling Methods	107
5.3	Sample Mean, Sample Variance and Sample Proportion	109
5.4	Sampling Distributions	110
5.5	The Sampling Distribution of the Sample Mean	111
5.6	The Sampling Distribution of the Sample Variance	116
5.7	Distribution of the Ratio of Two Sample Variances	118
5.8	The Sampling Distribution of the Sample Proportion	119
6	Theory of Estimation	121
6.1	Introduction	121
6.2	Point Estimation	121
6.3	Maximum Likelihood Estimation (MLE)	127
6.4	Bayesian Estimation	129
6.5	Estimation using Method of Moments	134
6.6	Interval Estimation	135
6.7	Confidence Interval for the Mean in a Normal Population	136
6.8	Confidence Interval for the Variance in a Normal Population	138
6.9	Confidence Interval for the Difference of Two Normally Distributed Population Means	139
6.10	Confidence Interval for the Ratio of Two Normally Distributed Population Variances	140
6.11	Confidence Interval for a Population Proportion	141
6.12	Determination of Sample Size	142
6.13	Frequentist vs Bayesian Approaches of Statistical Inference	145
7	Test of Hypothesis	149
7.1	What is Hypothesis Testing?	149
7.2	Test Statistic	151
7.3	Parametric vs. Non-Parametric Tests	153
7.4	Type I and Type II Error	154
7.5	General Procedure for Hypothesis Testing	158
7.6	Statistical Test for a Normally Distributed Population Mean μ	158
7.7	Statistical Test for a Normally Distributed Population Variance σ^2	161
7.8	Statistical Test for the Difference of Two Normally Distributed Population Means $\mu_1 - \mu_2$	162
7.9	Statistical Test for the Ratio of Two Normally Distributed Population Variances .	164
7.10	Statistical Test for a Population Proportion p	165
7.11	Relationship Between Confidence Interval and Hypothesis Tests	166
7.12	Power of a Test	167
8	Non-Parametric Tests and Chi-Square Tests	171

8.1	Introduction	171
8.2	Chi-Square Tests	172
8.3	Chi-Square Goodness-of-Fit Test	173
8.4	Chi-Square Test of Independence	173
8.5	Chi-Square Test of Homogeneity	175
8.6	Other Non-Parametric Tests	176
8.7	Assumptions and Limitations	176
8.8	Effect Size Measures	177
8.9	Power and Sample Size Calculations	177
8.10	Computational Examples	178
8.11	Advanced Topics	178
8.12	Practical Considerations	179
8.13	Conclusion	179
9	Analysis of Variance (ANOVA)	181
9.1	One-Way ANOVA	181
9.2	Mathematical Model for One-Way ANOVA	182
9.3	Analysis of One-Way ANOVA	183
9.4	Example of One-Way ANOVA	191
9.5	Effect Size and Practical Significance	193
9.6	Post-Hoc Analysis	194
9.7	Two-Way ANOVA	194
9.8	Mathematical Model for Two-Way ANOVA	195
9.9	Analysis of Two-Way Classified Data	197
9.10	Example of Two-Way ANOVA	201
9.11	Key Formulas Summary	204
10	Introduction to Design of Experiments	205
10.1	What is Design of Experiments?	205
10.2	Planning an Experiment	206
10.3	Types of Experimental Designs	207
10.4	Completely Randomized Design (CRD)	207
10.5	Randomized Block Design (RBD)	208
10.6	Latin Square Design (LSD)	209
10.7	Factorial Designs	211
10.8	Two-Factor Factorial Design (2^2 Design)	211
10.9	Three-Factor Factorial Design (2^3 Design)	212
10.10	Analysis of Variance for Factorial Designs	213
10.11	Assumptions and Model Checking	214
10.12	Summary	214

11 Correlation and Regression	217
11.1 Correlation	217
11.2 Pearson Correlation Coefficient	218
11.3 Effect of Linear Transformations on Correlation Coefficient	220
11.4 Correlation Is Not Causation	221
11.5 Regression Analysis	222
11.6 The Simple Linear Regression Model	222
11.7 Estimating Parameters Using Least Squares	223
11.8 Multiple Regression	227
12 Theory of Errors	235
12.1 Introduction	235
12.2 Random Errors	235
12.3 Systematic Errors	237
12.4 Gross Errors (Blunders)	238
12.5 Error Propagation	239
12.6 Practical Example: Pendulum Experiment	240
12.7 Conclusion	241
13 Statistical Quality Control	243
13.1 Introduction to Quality and Statistical Quality Control	243
13.2 Control Charts for Variables	245
13.3 Control Charts for Attributes	251
13.4 Pareto Chart in Statistical Quality Control	252
13.5 Process Capability Analysis	254
13.6 Acceptance Sampling	263
13.7 Reliability and Life Testing	264
13.8 Conclusion	267
14 Index Numbers	269
14.1 Introduction	269
14.2 Introduction to Index Numbers	269
14.3 Basic Concepts and Terminology	270
14.4 Types of Index Numbers	270
14.5 Methods of Constructing Price Index Numbers	271
14.6 Quantity Index Numbers	273
14.7 Properties and Tests of Index Numbers	273
14.8 Chain Index Numbers	273
14.9 Deflating and Real Values	274
14.10 Graphical Representation	275
14.11 Applications of Index Numbers	275

14.12	Limitations of Index Numbers	275
14.13	Worked Examples	276
14.14	Summary of Key Formulas	276

Chapter 1

Descriptive Statistics

1.1 Introduction

<https://bookdown.org/egarpor/inference/estmeth-mm.html>

Statistics is a branch of mathematics concerned with the systematic handling of data. It encompasses the processes of collecting, organizing, summarizing, and presenting data in meaningful ways. Whether working with numerical or categorical data, statistics helps process complex information into simpler, interpretable forms using tables, charts, and descriptive measures such as averages and variability.

Beyond description, statistics is fundamentally about **making inferences**. Using a sample from a larger population, statistical methods allow us to draw conclusions, test hypotheses, and estimate unknown parameters with a quantifiable level of confidence. This inferential power is crucial in fields where data are limited or difficult to obtain in full—such as medicine, economics, environmental science, and the social sciences.

At its core, statistics provides a rigorous framework for **decision-making under uncertainty**. By modeling randomness and variability, it enables researchers, policymakers, and professionals to make informed judgments backed by data. From predicting market trends to evaluating clinical trials, statistical reasoning underpins modern evidence-based practice across nearly every scientific and industrial domain.

1.1.1 Uses of Statistics

Statistics is widely used in various fields such as:

- **Business:** Companies use statistics to understand market trends, forecast sales, and study consumer behavior. It plays a key role in product development, marketing strategies, risk analysis, and making data-driven decisions that improve profitability and efficiency.
- **Healthcare:** In medicine and public health, statistics is essential for evaluating the effectiveness of new treatments, analyzing clinical trial data, and tracking the spread of diseases. It supports evidence-based decision-making and helps ensure patient safety and optimal care.
- **Government:** Governments rely on statistical data for policy-making, budget planning, and resource allocation. Population censuses, employment statistics, and health surveys are used to understand societal needs and implement programs for public welfare.
- **Education:** Statistics helps in analyzing student performance through test scores, attendance records, and learning outcomes. It supports curriculum planning, identifying gaps in learning, and measuring the effectiveness of teaching methods or interventions.

1.1.2 Types of Data

Data refers to raw information collected through observation, measurement, surveys, experiments, or other methods. It serves as the foundation for statistical analysis. In statistics, we use data to uncover patterns, test hypotheses, make predictions, and support decision-making in various fields such as education, healthcare, business, and science.

Data can be broadly classified into two main types based on their nature: **qualitative** and **quantitative**. Understanding these types is essential for choosing the appropriate statistical methods for analysis.

1. Qualitative (Categorical) Data

Qualitative data represent categories or labels that describe characteristics or qualities. These values are non-numerical and cannot be measured in terms of quantity, although they can be counted or classified.

- **Nominal data:** These are categories with no natural order or ranking. The categories are simply different from one another.
Examples: Types of fruits ('apple', 'banana', 'mango'), blood groups ('A', 'B', 'AB', 'O').
- **Ordinal data:** These categories have a meaningful order or ranking, but the intervals between the categories are not necessarily equal or known.
Examples: Customer satisfaction ratings ('poor', 'fair', 'good', 'excellent'), education level ('high school', 'undergraduate', 'postgraduate').

2. Quantitative (Numerical) data

Quantitative data represent numerical values that indicate amounts or measurements. These data can be subjected to arithmetic operations such as addition, subtraction, and averaging.

- **Discrete data:** These are countable numbers with distinct, separate values. Discrete data usually arise from counting processes.
Examples: Number of children in a family, number of cars in a parking lot, number of questions answered correctly on a test.
- **Continuous data:** These are measurable quantities that can take any value within a given range. Continuous data usually come from measurements and can include fractions and decimals.
Examples: Height of students, weight of a package, temperature of a city, time taken to complete a task.

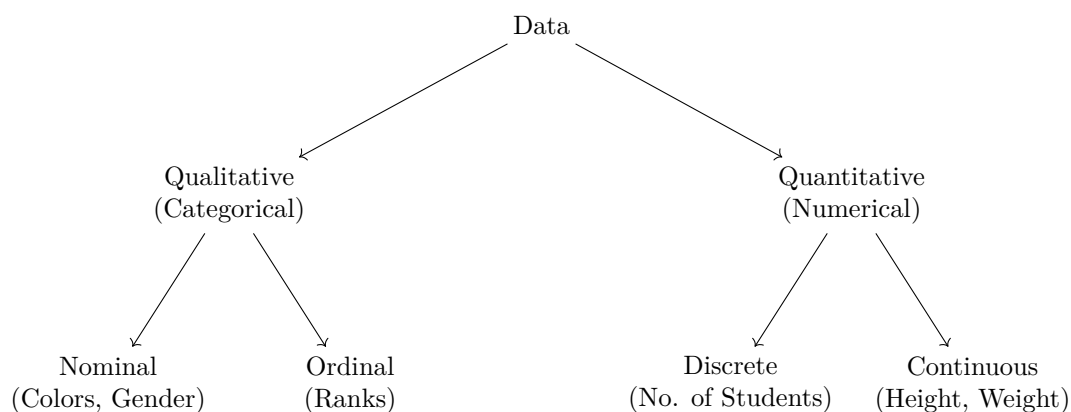


Figure 1.1: *Types of data.*

1.1.3 Data Collection Methods

Data collection is a fundamental step in any statistical investigation. The method used to collect data affects its reliability, accuracy, and suitability for analysis. Common methods include:

- **Surveys and questionnaires:** Used to gather information from individuals by asking a set of structured or semi-structured questions. Surveys can be conducted in person, over the phone, via email, or online platforms.
- **Interviews:** Data is collected through verbal interaction between the interviewer and the respondent. Interviews may be structured, semi-structured, or unstructured, and are useful for collecting in-depth qualitative information.
- **Observation:** Involves directly watching subjects in a natural or controlled environment. Useful for behavioral and environmental studies, especially when respondents may not reliably report information.
- **Experiments:** Data is collected under controlled conditions by manipulating variables and observing the outcomes. This method is common in natural and social sciences for studying causal relationships.
- **Document or record review:** Involves collecting data from existing documents, databases, reports, or archives. This method is typically used for gathering secondary data.

Broadly, data can be classified as either **primary** or **secondary** depending on how it is collected.

- **Primary data:** Data collected directly from first-hand sources by the investigator for a specific purpose. It is original, up-to-date, and tailored to the research objective, but can be time-consuming and costly to collect.
Examples: Conducting a survey, performing experiments, direct observations.
- **Secondary data:** Data that has already been collected, compiled, and published by others. It is economical and quick to access, but may not perfectly match the researcher's needs and might be outdated or biased.
Examples: Government census reports, published research papers, company financial statements.



Figure 1.2: *Types of data based on collection methods.*

1.2 Presentation of Data

Presentation of data refers to the various ways in which raw data can be organized and displayed to facilitate understanding, comparison, and interpretation. Broadly, these methods fall into three categories:

1. Textual Representation
2. Tabular Representation
3. Diagrammatic (Graphical) Representation

1.2.1 Textual Representation

Textual representation describes data in words or sentences. It is useful for small data sets or when contextual explanation is needed.

“In the month of June, the daily maximum temperatures recorded in City X ranged from 28°C to 35°C, with an average of 31°C. The lowest temperature (28°C) occurred on June 3 and June 17, while the highest (35°C) was on June 24.”

1.2.2 Tabular Representation

Tabular representation organizes data into rows and columns, providing a clear and concise way to display raw observations. Consider the daily maximum temperatures (in °C) recorded in City X for the first 12 days of June 2025:

30, 32, 28, 31, 33, 34, 29, 30, 32, 31, 30, 29.

We first present these observations in a simple day-value table.

Day	1	2	3	4	5	6	7	8	9	10	11	12
Temp (°C)	30	32	28	31	33	34	29	30	32	31	30	29

Table 1.1: *Daily maximum temperatures (°C) in city X, June 1–12, 2025.*

1. Simple Frequency Distribution

From the raw data table, it can be difficult to see how often each temperature occurs. To remedy this, we construct a *simple frequency distribution*, which counts the number of times each unique value appears.

In statistics, the **frequency** of a value is the number of times that value appears in a dataset, and a **frequency distribution** is a tabular summary that pairs each distinct observation with its frequency. To make it easy to see how often each temperature occurs, we convert our raw list of daily readings into a simple frequency distribution by counting the occurrences of each unique temperature:

Temp (°C)	28	29	30	31	32	33	34
Frequency	1	2	3	2	2	1	1

Table 1.2: *Simple frequency distribution of maximum temperatures (°C).*

This table immediately shows, for example, for example, that 30 °C was recorded on three days, while 28 °C and 33 °C each occurred only once. Sometimes, to express them proportionally, we also compute **relative frequencies** and accumulate them in a **cumulative frequency distribution**. The *relative frequency* is calculated as

$$\text{Relative frequency} = \frac{\text{Frequency}}{\text{Total no. of days}} \times 100\%,$$

which is generally expressed in percentages and the *cumulative frequency* at each temperature calculates the total number of observations less than that value. These additional columns help us understand both the proportion and the running total of the data.

Temp	28	29	30	31	32	33	34
Frequency	1	2	3	2	2	1	1
Relative frequency (%)	8.3	16.7	25.0	16.7	16.7	8.3	8.3
Cumulative frequency	1	3	6	8	10	11	12

Table 1.3: *Frequency, relative frequency, and cumulative frequency distribution.*

2. Grouped Frequency Distribution

When dealing with larger datasets or continuous measurements, listing every unique value becomes impractical. Instead, we group data into multiple groups or classes with definite boundaries, and then tabulate frequencies for each class. Suppose we choose class boundaries with an intervals of 2 °C each to cover the range of our data:

$$[28, 30), [30, 32), [32, 34), [34, 36)$$

Class Interval	[28,30)	[30,32)	[32,34)	[34,36)
Frequency	3	5	3	1
Relative frequency (%)	25.0	41.7	25.0	8.3
Cumulative frequency	3	8	11	12

Table 1.4: *Grouped frequency distribution with relative and cumulative frequencies.*

We notice the interval $[30, 32)$ contains $\{30, 30, 31, 31, 32\}$, giving a frequency of 5 which is the maximum among all classes. Its relative frequency is $5/12 \approx 41.7\%$, and its cumulative frequency of 8 indicates that 8 observations are below 32 °C.

1.2.3 Diagrammatic (Graphical) Representation

Graphical methods transform data into visual form. They allow quick identification of patterns, trends, and outliers.

1. Bar Graph

Bar graphs are used for categorical data.

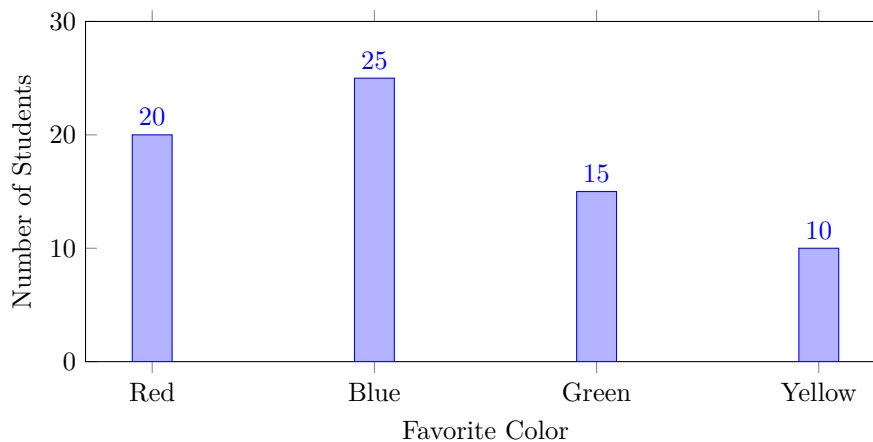


Figure 1.3: Favorite Colors of Students

2. Pie Chart

Pie charts show the proportion of categories within a whole.

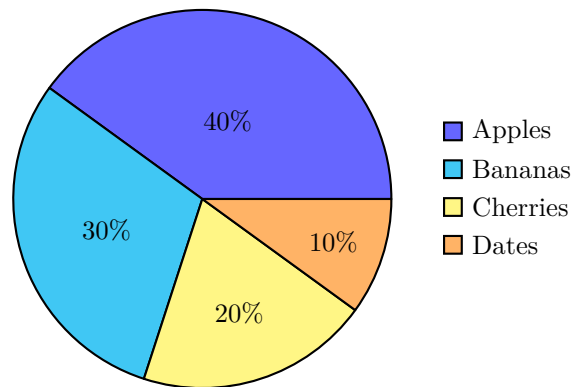


Figure 1.4: Fruit Preferences

3. Histogram

Histograms are used for continuous numerical data.

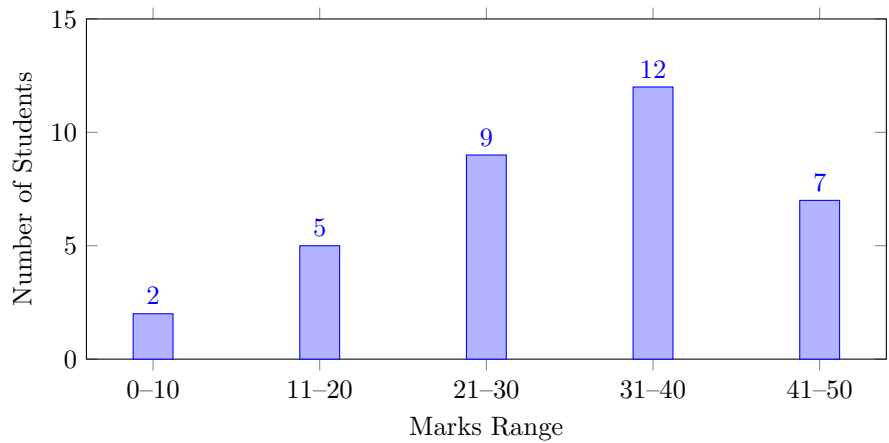


Figure 1.5: Distribution of Marks

4. Line Graph

Line graphs depict data trends over time.

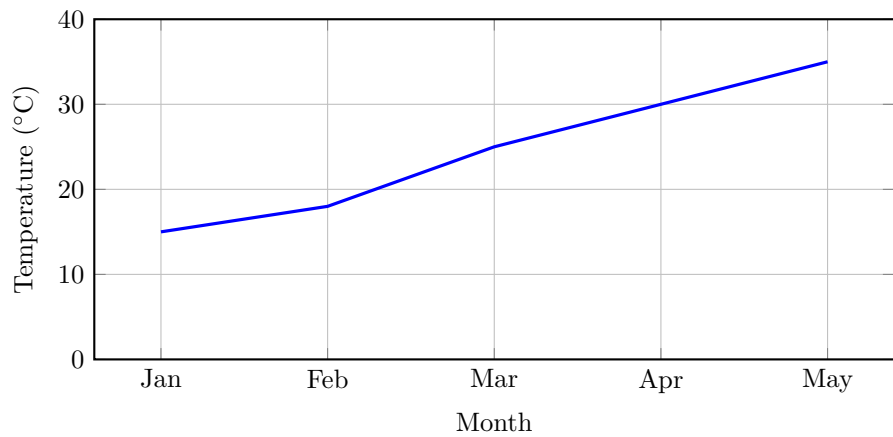


Figure 1.6: *Monthly average temperature.*

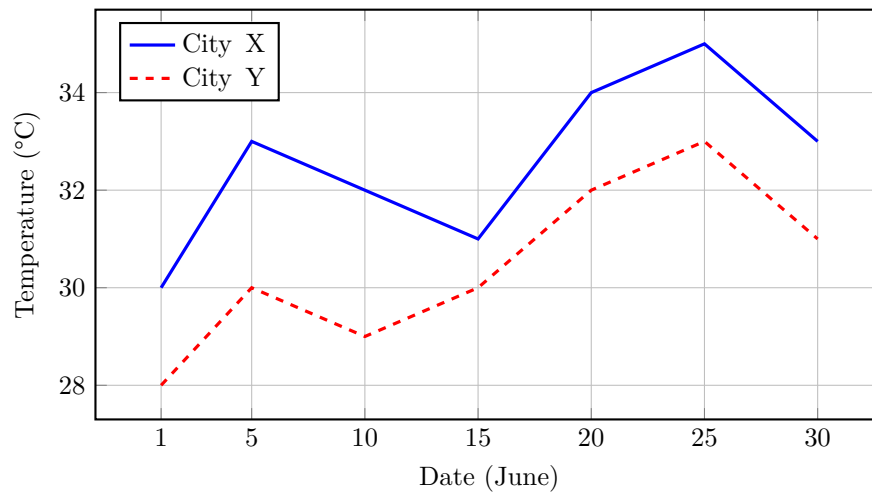


Figure 1.7: Multiple-line diagram of daily max temperatures.

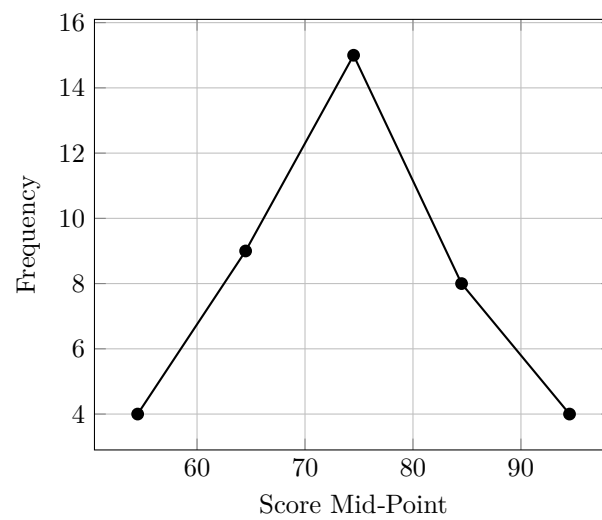


Figure 1.8: Frequency Polygon

Step (Frequency Polygon) Diagram

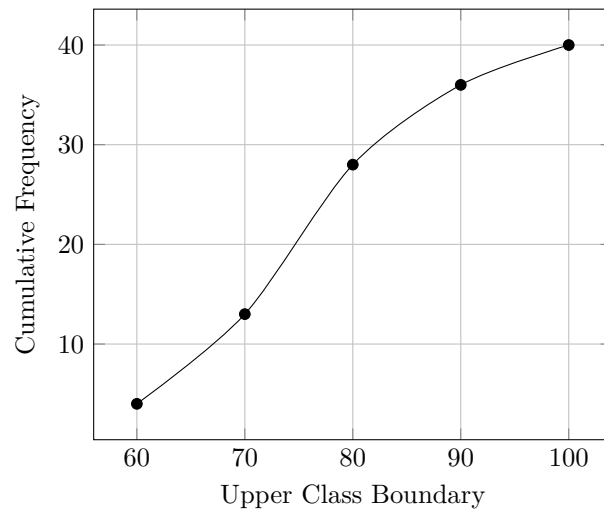


Figure 1.9: Ogive

Ogive (Cumulative Frequency Curve)

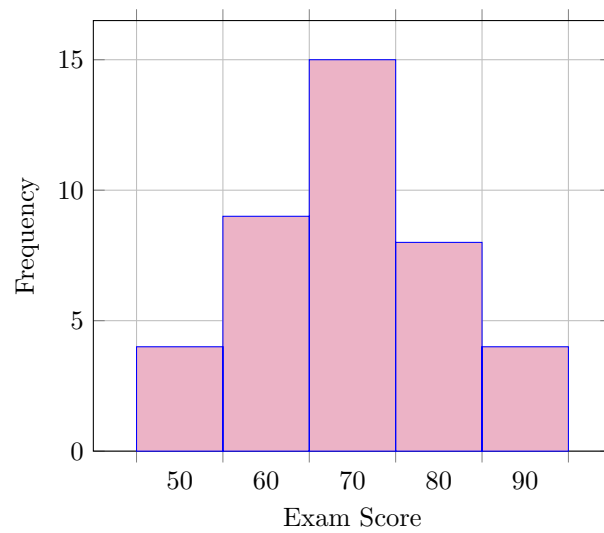


Figure 1.10: Histogram of Exam Scores

Histogram



Figure 1.11: Smooth Frequency Curve

Frequency Curve of Different Types

Tally Marks

Tally marks are a simple and intuitive method for counting and recording frequencies, especially useful during manual data collection or for small datasets. Each tally mark (|) represents one occurrence of a value. To improve readability, tallies are grouped in bundles of five, with the fifth mark crossing the previous four represents a count of 5.

Tally marks provide a quick visual count before converting to numerical frequencies in a formal table. They are especially helpful in classroom or field settings where data is collected in real time.

Suppose we record the maximum daily temperatures in City Y over 12 days:

30, 32, 32, 30, 32, 30, 32, 31, 32, 32, 32, 31

Using tally marks, we count how many times each temperature appears:

```

30: |||
31: ||
32: |||/|||

```

Stem-and-Leaf Diagram

Stem	Leaf
5	0 2 4 7
6	0 1 3 5 9
7	2 4 4 8
8	0 1 3
9	0 5

This completes the detailed presentation of data methods. Each technique has its strengths and is chosen based on the nature of the data and the purpose of analysis.

1.3 Measures of Central Tendency

Quite often, data exhibit a tendency to cluster around a central value. Measures of central tendency are numerical indicators that describe this central value of a data set. The most common measures include the mean, median, and mode. Each measure tells us something different about our data, and knowing when to use each one can really help us make sense of the numbers.

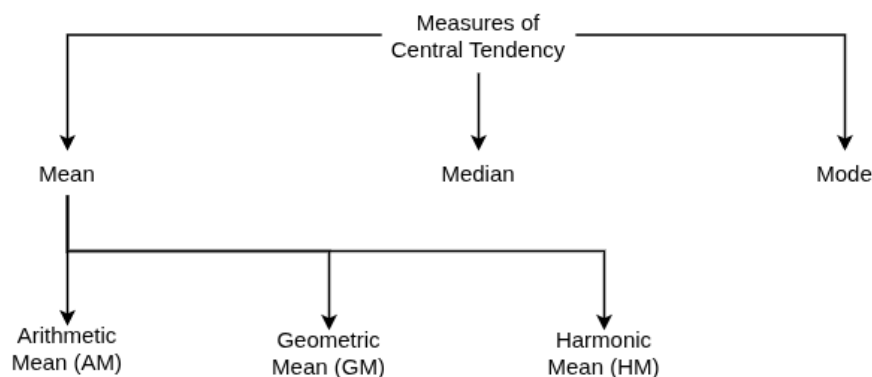


Figure 1.12: *Measures of central tendency.*

1.3.1 Mean

The mean (or average) is the most commonly used measure of central tendency and is defined in several forms:

- **Arithmetic Mean (AM):**

- Simple AM: For a dataset x_1, x_2, \dots, x_n , the arithmetic mean is given by:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

- Weighted AM: When each value x_i has an associated frequency f_i , the weighted arithmetic mean is:

$$\bar{x} = \frac{\sum_{i=1}^n f_i x_i}{\sum_{i=1}^n f_i}$$

Theorem: If $x_i = a$ (constant) for all i , then the arithmetic mean is also a , that is,

$$\bar{x} = a.$$

Proof:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \sum_{i=1}^n a = \frac{1}{n} \cdot na = a$$

■

Theorem: If $y_i = a + x_i$, then the mean of y is given by:

$$\bar{y} = a + \bar{x}.$$

Proof:

$$\begin{aligned}\bar{y} &= \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n (a + x_i) = \frac{1}{n} \left(\sum_{i=1}^n a + \sum_{i=1}^n x_i \right) \\ &= \frac{1}{n} \left(na + \sum_{i=1}^n x_i \right) = a + \bar{x}.\end{aligned}$$

■

Theorem: Let a dataset be composed of two distinct groups of observations:

- Group 1 consists of n_1 observations with arithmetic mean \bar{x}_1 ,
- Group 2 consists of n_2 observations with arithmetic mean \bar{x}_2 .

Then, the arithmetic mean \bar{x} of the combined dataset (of size $n_1 + n_2$) is given by:

$$\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}.$$

Proof: Total sum of group 1 is $n_1 \bar{x}_1$ and total sum of group 2 is $n_2 \bar{x}_2$.

Then total sum = $n_1 \bar{x}_1 + n_2 \bar{x}_2$

Total number of observations = $n_1 + n_2$

Therefore, combined AM = $\bar{x} = \frac{\text{Total sum}}{\text{Total number of observations}} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}$

■

- **Geometric Mean (GM):** Geometric mean of a set of n observation is the n th root of their product. It is only defined for positive values.

- Simple GM:

$$x_G = \left(\prod_{i=1}^n x_i \right)^{1/n} = \sqrt[n]{x_1 x_2 \dots x_n}$$

- Weighted GM:

$$x_G = \left(\prod_{i=1}^n x_i^{f_i} \right)^{1/N} = \sqrt[N]{x_1^{f_1} x_2^{f_2} \dots x_n^{f_n}}$$

Where

$$N = \sum_{i=1}^n f_i$$

Theorem: The GM of a set of positive values x_1, x_2, \dots, x_n is equal to the antilogarithm (exponential) of the AM of their logarithms:

$$\text{GM} = \exp \left(\frac{1}{n} \sum_{i=1}^n \log x_i \right)$$

Proof:

$$\begin{aligned} \text{GM} &= \left(\prod_{i=1}^n x_i \right)^{1/n} = \exp \left(\log \left(\prod_{i=1}^n x_i \right)^{1/n} \right) \\ &= \exp \left(\frac{1}{n} \log \left(\prod_{i=1}^n x_i \right) \right) = \exp \left(\frac{1}{n} \sum_{i=1}^n \log x_i \right) \end{aligned}$$

■

Theorem: Suppose we have two groups:

- Group 1 has N_1 positive values with geometric mean x_{G_1} ,
- Group 2 has N_2 positive values with geometric mean x_{G_2} .

Then the combined geometric mean GM of all $N_1 + N_2$ values is:

$$\text{GM} = \left(x_{G_1}^{N_1} \cdot x_{G_2}^{N_2} \right)^{1/(N_1+N_2)}$$

Proof: Let the product of values in group 1 be $P_1 = \prod_{i=1}^{N_1} x_i$, so that:

$$x_{G_1} = (P_1)^{1/N_1} \Rightarrow P_1 = x_{G_1}^{N_1}$$

Similarly, for group 2:

$$P_2 = \prod_{j=1}^{N_2} y_j = x_{G_2}^{N_2}$$

Then the overall product:

$$P = P_1 \cdot P_2 = x_{G_1}^{N_1} \cdot x_{G_2}^{N_2}$$

The combined GM is:

$$\text{GM} = (P)^{1/(N_1+N_2)} = \left(x_{G_1}^{N_1} \cdot x_{G_2}^{N_2} \right)^{1/(N_1+N_2)}$$

■

- **Harmonic Mean (HM):** Harmonic Mean of a set of observations is the reciprocal of the arithmetic mean of the reciprocals.

- Simple HM:

$$x_H = \frac{1}{\frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}} = \frac{n}{\sum_{i=1}^n \frac{1}{x_i}}$$

- Weighted HM:

$$x_H = \frac{1}{\frac{1}{N} \sum_{i=1}^n \frac{f_i}{x_i}} = \frac{N}{\sum_{i=1}^n \frac{f_i}{x_i}}$$

It is important to note where to use GM and where to use HM. GM is useful for averaging **ratios**, **rates** and **percentages**. As an illustration, we consider the following example.

Example: The ratio of the prices in 1994 and to those in 1982 for four commodities are 0.92, 1.25, 1.75 and 0.85. To get the average price ratio use geometric mean

$$\begin{aligned}\log x_G &= \frac{1}{n} (\log 0.92 + \log 1.25 + \log 1.75 + \log 0.85) \\ &= 0.5829 = \log 1.1436 \\ \Rightarrow x_G &= 1.1436\end{aligned}$$

GM is also useful if one wants to determine the values of a variable at the midpoint of a time interval when the variable changes over time exponentially. Thus if the value at two points 0 and t be a and ar^t , then its value at the midpoint $\frac{t}{2}$ is $(a \times ar^t)^{1/2} = ar^{t/2}$.

Now consider the following example:

Example: A person goes from X to Y on cycle at 20 km/h and returns at 24 km/h. What is the average speed for the entire trip?

If we use AM, then the average speed is

$$\frac{1}{2}(20 + 24) = 22 \text{ km/h}$$

But is this correct?

Consider the total distance between X and Y is 1 km for the sake of simplicity. So the total distance covered = 2 km. The time taken for the person to go from X to Y is $\frac{1}{20} = 0.05$ hr and the time taken to return is $\frac{1}{24} = 0.04166$ hr.

$$\text{Therefore average speed} = \frac{\text{Total distance}}{\text{Total time}} = \frac{2}{0.05 + 0.04166} = 21.8 \text{ km/h.}$$

Clearly, the AM value of 22 km/h overestimates the actual average speed. Now consider the harmonic mean (HM) of the two speeds:

$$\frac{2}{\frac{1}{20} + \frac{1}{24}} = 21.8 \text{ km/h}$$

This matches the correct value.

Now where to use the HM? The harmonic mean is particularly useful when dealing with quantities expressed in the form “ x per unit y ”, such as “km per hour”, “rupees per kg”, and similar rates. Here x and y are unit of measures, not numeric variables.

Rule of Thumb:

- Use the **harmonic mean (HM)** when equal quantities of x are involved.
- Use the **arithmetic mean (AM)** when equal quantities of y are involved.

This principle can be illustrated with the following example.

Example: Suppose a train covers n **equal distances**, each of s kilometers, with speeds v_1, v_2, \dots, v_n km/h. The average speed is the total distance divided by the total time taken. Thus,

$$\text{Average speed} = \frac{ns}{\frac{s}{v_1} + \frac{s}{v_2} + \dots + \frac{s}{v_n}} = \frac{n}{\frac{1}{v_1} + \frac{1}{v_2} + \dots + \frac{1}{v_n}} = \frac{n}{\sum_{i=1}^n \frac{1}{v_i}}$$

This is the **harmonic mean (HM)** of the given speeds.

On the other hand, if the train travels for n **equal time intervals**, each of duration t hours, at speeds v_1, v_2, \dots, v_n km/h, then the total distance covered is:

$$\text{Total distance} = v_1 t + v_2 t + \cdots + v_n t = t(v_1 + v_2 + \cdots + v_n)$$

and the total time is nt . So, the average speed is:

$$\text{Average speed} = \frac{t(v_1 + v_2 + \cdots + v_n)}{nt} = \frac{v_1 + v_2 + \cdots + v_n}{n} = \frac{1}{n} \sum_{i=1}^n v_i$$

which is the **arithmetic mean (AM)** of the given speeds.

Theorem: The sum of squared deviations from a constant A is minimized when A equals the arithmetic mean \bar{x} , i.e.,

$$\sum_{i=1}^n (x_i - A)^2 \text{ is minimized when } A = \bar{x}$$

Proof: Let $S(A) = \sum_{i=1}^n (x_i - A)^2$. Expand this:

$$S(A) = \sum_{i=1}^n (x_i^2 - 2Ax_i + A^2) = \sum x_i^2 - 2A \sum x_i + nA^2$$

To minimize $S(A)$, take derivative with respect to A and set it to zero:

$$\frac{dS}{dA} = -2 \sum x_i + 2nA = 0 \quad \Rightarrow \quad A = \frac{1}{n} \sum x_i = \bar{x}$$

Now, take the second derivative:

$$\frac{d^2S}{dA^2} = 2n > 0$$

Since the second derivative is positive, the function $S(A)$ has a minimum at $A = \bar{x}$. ■

Theorem: For two observations,

$$\text{GM}^2 = \text{AM} \times \text{HM}$$

Proof: Let a and b be two observations (positive numbers).

Compute left-hand side:

$$\text{GM}^2 = (\sqrt{ab})^2 = ab$$

Compute right-hand side:

$$\text{AM} \times \text{HM} = \left(\frac{a+b}{2} \right) \left(\frac{2ab}{a+b} \right) = ab$$

Hence,

$$\text{GM}^2 = \text{AM} \times \text{HM}$$
■

Theorem: For any set of n positive real numbers x_1, x_2, \dots, x_n , the following inequality holds:

$$\text{AM} \geq \text{GM} \geq \text{HM}$$

with equality if and only if $x_1 = x_2 = \dots = x_n$.

Proof: Let x_1, x_2, \dots, x_n be positive real numbers.

• **Step 1: Proving AM \geq GM**

We can prove this using the method of induction.

- **Base Case:** Let us first consider two observations $x_1 = a > 0$, $x_2 = b > 0$. We have to prove:

$$\frac{a+b}{2} \geq \sqrt{ab}$$

Consider the square of the difference:

$$\left(\frac{a-b}{2}\right)^2 \geq 0 \Rightarrow \frac{a^2 - 2ab + b^2}{4} \geq 0$$

$$\Rightarrow a^2 + b^2 \geq 2ab \Rightarrow (a+b)^2 \geq 4ab$$

$$\Rightarrow \left(\frac{a+b}{2}\right)^2 \geq ab \Rightarrow \frac{a+b}{2} \geq \sqrt{ab}$$

Equality holds if and only if $a = b$.

- **Inductive Step:** Assume the inequality holds for $n = k$, i.e., for all positive x_1, \dots, x_k :

$$\frac{x_1 + x_2 + \dots + x_k}{k} \geq (x_1 x_2 \dots x_k)^{1/k}$$

We must show it holds for $n = k + 1$ too.

Let $x_1, x_2, \dots, x_k, x_{k+1}$ be positive numbers. Define:

$$A = \frac{x_1 + x_2 + \dots + x_k}{k}, \quad G = (x_1 x_2 \dots x_k)^{1/k}$$

By the inductive hypothesis, $A \geq G$.

Now apply the $n = 2$ case to the numbers A and x_{k+1} :

$$\frac{A + x_{k+1}}{2} \geq \sqrt{Ax_{k+1}} \geq \sqrt{Gx_{k+1}}$$

Now note:

$$\frac{x_1 + \dots + x_k + x_{k+1}}{k+1} = \frac{kA + x_{k+1}}{k+1}$$

We now want to show:

$$\frac{kA + x_{k+1}}{k+1} \geq (x_1 x_2 \dots x_k x_{k+1})^{1/(k+1)}$$

Let us define:

$$B = (x_1 x_2 \dots x_k x_{k+1})^{1/(k+1)} = (G^k \cdot x_{k+1})^{1/(k+1)}$$

Use the inequality between arithmetic and geometric mean on A and x_{k+1} :

$$\frac{kA + x_{k+1}}{k+1} \geq (A^k \cdot x_{k+1})^{1/(k+1)}$$

Since $A \geq G$, and exponentiation preserves the inequality for positive values:

$$A^k \geq G^k \Rightarrow (A^k \cdot x_{k+1})^{1/(k+1)} \geq (G^k \cdot x_{k+1})^{1/(k+1)} = B$$

Therefore,

$$\frac{x_1 + \dots + x_{k+1}}{k+1} \geq B = (x_1 x_2 \dots x_k x_{k+1})^{1/(k+1)}$$

This therefore proves that:

$$\text{AM} \geq \text{GM}$$

• Step 2: Proving GM \geq HM

Recall the definition of the harmonic mean:

$$\text{HM} = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}$$

Now consider the reciprocals $\frac{1}{x_1}, \frac{1}{x_2}, \dots, \frac{1}{x_n}$, which are also positive. Thus we can apply the AM–GM inequality to the reciprocals:

$$\frac{1}{n} \left(\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n} \right) \geq \left(\frac{1}{x_1 x_2 \dots x_n} \right)^{1/n}$$

Taking reciprocals of both sides:

$$\frac{n}{\frac{1}{x_1} + \dots + \frac{1}{x_n}} \leq (x_1 x_2 \dots x_n)^{1/n}$$

$$\Rightarrow \text{HM} \leq \text{GM}$$

with equality if and only if $x_1 = x_2 = \dots = x_n$.

Combining both steps:

$$\text{AM} \geq \text{GM} \geq \text{HM}$$

with equality throughout if and only if all the x_i are equal. ■

1.3.2 Median

The **median** of a set of observation is the middlemost value when the observations are arranged in increasing or decreasing order of magnitude.

It is denoted by M_i or \tilde{x} . It divides the dataset into two equal halves: 50% of the values lie below the median and 50% lie above.

1. Median in a Simple Series (Ungrouped Data)

For a dataset with n observations arranged in ascending order:

- If n is odd, the median is the value at the $\left(\frac{n+1}{2}\right)^{\text{th}}$ position.
- If n is even, the median is the arithmetic mean of the values at the $\left(\frac{n}{2}\right)^{\text{th}}$ and $\left(\frac{n}{2} + 1\right)^{\text{th}}$ positions.

Example: Find the median of the dataset:

7, 2, 5, 9, 4

Arranging in ascending order: 2, 4, 5, 7, 9. Since there are 5 values (odd), the median is the 3rd value:

Median = 5

2. Median in a Simple Frequency Distribution

In a simple frequency distribution, each data value is associated with a frequency. The procedure is identical to that of a simple frequency distribution:

- Arrange the data in ascending order.
- Compute cumulative frequencies based on weights.
- Find total frequency N , then find the smallest value for which the cumulative frequency is greater than or equal to $\frac{N}{2}$.

Example: Consider the following table containing the values, frequencies and cumulative frequencies.

Value	Frequency	Cumulative Frequency
2	3	3
4	5	8
6	7	15
8	5	20

Table 1.5: Simple frequency distribution table to calculate median.

$$N = 3 + 5 + 7 + 5 = 20 \Rightarrow \frac{N}{2} = 10$$

Since 10 is between 8 and 15, the Median is 6. This works regardless of whether N is odd or even.

3. Median in a Grouped Frequency Distribution

For a grouped frequency distribution, the cumulative frequencies are used to locate the median.

- Compute cumulative frequencies.
- Find $N = \sum f_i$, the total number of observations.
- Find $\frac{N}{2}$.
- Locate the median class (the class whose cumulative frequency is greater than or equal to $\frac{N}{2}$).
- Use the formula:

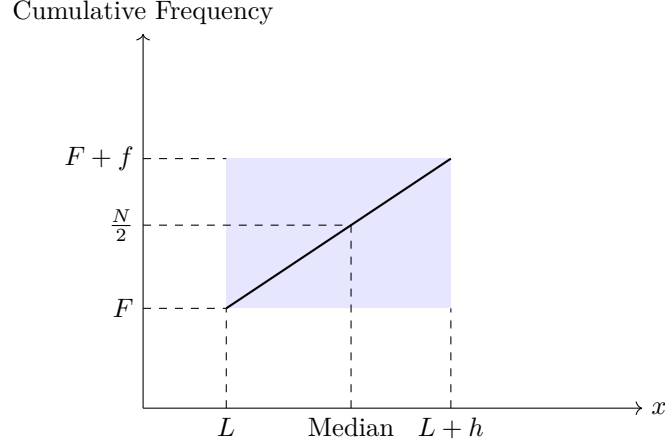
$$\text{Median} = L + \left(\frac{\frac{N}{2} - F}{f} \right) \cdot h$$

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 : width of the class interval

To arrive at this formula we assume that the cumulative frequency is a linear function of x within the class L and $L + h$. Then

$$\frac{\text{Median} - L}{h} = \frac{\frac{N}{2} - F}{f} \Rightarrow \text{Median} = L + \left(\frac{\frac{N}{2} - F}{f} \right) \cdot h$$



Example: Consider the following table containing the class values, frequencies and cumulative frequencies.

Class Interval	Frequency	Cumulative Frequency
0–10	5	5
10–20	8	13
20–30	12	25
30–40	6	31

Table 1.6: *Grouped frequency distribution table for calculating the median.*

$$N = 5 + 8 + 12 + 6 = 31, \quad \frac{N}{2} = 15.5$$

Since 15.5 is in between 13 and 25, the Median class is 20–30. Thus,

$$L = 20, F = 13, f = 12, h = 10$$

$$\text{Median} = 20 + \left(\frac{15.5 - 13}{12} \right) \cdot 10 = 20 + \left(\frac{2.5}{12} \right) \cdot 10 = 20 + 2.08 = 22.08$$

Theorem: Let x_1, x_2, \dots, x_n be a set of observations. Define the function:

$$S(A) = \sum_{i=1}^n |x_i - A|$$

Then $S(A)$ is minimized when $A = \text{Median}$.

Proof: Let us arrange the observations x_1, x_2, \dots, x_n in increasing order and denote the ordered sequence by y_1, y_2, \dots, y_n . Since this is just a rearrangement of the original data, we have:

$$\sum_{i=1}^n |x_i - A| = \sum_{i=1}^n |y_i - A|.$$

We now analyze the behavior of this sum in two cases:

- **Case 1: n is odd (say $n = 2m + 1$)**

$$\begin{aligned} \sum_{i=1}^n |x_i - A| &= \sum_{i=1}^{2m+1} |y_i - A| \\ &= |y_1 - A| + |y_2 - A| + \dots + |y_m - A| + |y_{m+1} - A| \\ &\quad + |y_{m+2} - A| + \dots + |y_{2m+1} - A|. \end{aligned}$$

There are $2m + 1$ terms in the sum. We consider them in symmetric pairs from both ends:

- The sum $|y_1 - A| + |y_{2m+1} - A|$ is minimized when $A \in [y_1, y_{2m+1}]$.
- The sum $|y_2 - A| + |y_{2m} - A|$ is minimized when $A \in [y_2, y_{2m}]$.
- Continuing this way, the sum $|y_m - A| + |y_{m+2} - A|$ is minimized when $A \in [y_m, y_{m+2}]$.

The unpaired central term is $|y_{m+1} - A|$, which attains its minimum value (zero) when $A = y_{m+1}$.

Since all these intervals of minimization overlap at y_{m+1} , we conclude that:

$$\sum_{i=1}^n |x_i - A| \text{ is minimized when } A = y_{m+1} = \text{Median.}$$

- **Case 2: n is even (say $n = 2m$)**

$$\begin{aligned} \sum_{i=1}^n |x_i - A| &= \sum_{i=1}^{2m} |y_i - A| \\ &= |y_1 - A| + |y_2 - A| + \dots + |y_m - A| + |y_{m+1} - A| + \dots + |y_{2m} - A|. \end{aligned}$$

Now there are $2m$ terms, which can be grouped into m symmetric pairs:

- $|y_1 - A| + |y_{2m} - A|$ is minimized when $A \in [y_1, y_{2m}]$,
- $|y_2 - A| + |y_{2m-1} - A|$ is minimized when $A \in [y_2, y_{2m-1}]$, and so on.
- The final pair $|y_m - A| + |y_{m+1} - A|$ is minimized when $A \in [y_m, y_{m+1}]$.

Thus, the entire sum is minimized when $A \in [y_m, y_{m+1}]$. A natural choice is:

$$A = \text{Median} = \frac{y_m + y_{m+1}}{2},$$

which lies within the minimizing interval and hence ensures that the sum is minimized.

In both cases—odd and even number of observations—the value of A that minimizes $\sum_{i=1}^n |x_i - A|$ is the **median** of the dataset.

■

The median is a better measure of central tendency than the mean (AM) in the presence of outliers in the observations.

The **mean** (AM) is sensitive to extreme values or outliers, while the **median** (the middle value) is more robust and resistant to such anomalies. This makes the median a better measure of central tendency in the presence of outliers or skewed data.

Consider the marks obtained by 5 students:

$$\text{Scores} = \{10, 70, 75, 80, 90\}$$

$$\text{Mean} = \frac{70 + 75 + 80 + 85 + 90}{5} = \frac{400}{5} = 80; \quad \text{Median} = \text{Middle value} = 75$$

Here, both the mean and the median are equal and representative of the data, as there are no extreme values.

Now suppose one student scored unusually low:

$$\text{Scores} = \{10, 70, 75, 80, 90\}$$

$$\text{Mean} = \frac{10 + 70 + 75 + 80 + 90}{5} = \frac{325}{5} = 65; \quad \text{Median} = \text{Middle value} = 75$$

The mean drops to 65 due to the outlier (10), even though most students scored 70 or above. The median stays at 75 and gives a better picture of the typical student performance.

1.3.3 Mode

The **mode** of a given set of observations is the value which occurs with maximum frequency. It represents the highest peak in the frequency distribution.

The mode is generally denoted by M_o .

1. Mode in a Simple Series (Ungrouped Data)

To determine the mode:

- Count the frequency of each data value.
- The mode is the value with the highest frequency.

Example: Find the mode of the dataset:

$$3, 7, 2, 3, 9, 3, 5$$

The number 3 occurs most frequently (3 times), so:

$$\text{Mode} = 3$$

2. Mode in a Simple Frequency Distribution

In a simple frequency distribution, the mode is the data value corresponding to the maximum frequency.

Example: Consider the table:

Value	Frequency
4	3
5	7
6	5
7	2

Table 1.7: *Simple frequency distribution table for calculating mode.*

The highest frequency is 7, corresponding to the value 5. Hence:

$$\text{Mode} = 5$$

3. Mode in a Grouped Frequency Distribution

When the data is grouped into intervals, it is very difficult to find the mode accurately. However if all the classes are of equal width, then it is possible to approximately calculate the mode using the formula:

$$\text{Mode} = L + \left(\frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right) \cdot h$$

where:

- L : lower boundary of the modal class
- f_1 : frequency of the modal class
- f_0 : frequency of the class preceding the modal class
- f_2 : frequency of the class succeeding the modal class
- h : class width

How do we arrive at the formula for mode? In addition to the modal class frequency f_1 , mode also depends on f_0 (the frequency of the class preceding the modal class) and f_2 (the frequency of the class following the modal class). If they are equal, then one would take the midpoint of the modal class $L + \frac{h}{2}$ as the mode. However, if $f_0 - f_1$ is greater (smaller) than $f_1 - f_2$, then one would suppose that the mode is nearer to (further from) the lower boundary (L) of the modal class than the upper boundary ($L + h$). Mathematically, if we assume the proportion is same, then

$$\frac{d}{f_1 - f_0} = \frac{h - d}{f_1 - f_2}$$

Cross-multiplying and simplifying:

$$d \cdot (2f_1 - f_0 - f_2) = (f_1 - f_0) \cdot h$$

Solving for d :

$$d = \frac{(f_1 - f_0) \cdot h}{2f_1 - f_0 - f_2}$$

Hence, the mode is:

$$\text{Mode} = L + d = L + \left(\frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right) \cdot h$$



Example: Consider the following grouped data:

Class Interval	Frequency
0–10	4
10–20	6
20–30	10
30–40	8

Table 1.8: *Grouped frequency distribution table for calculating the mode.*

Here, the modal class is 20–30 because it has the highest frequency ($f_1 = 10$). The required values are:

$$L = 20, \quad f_0 = 6, \quad f_1 = 10, \quad f_2 = 8, \quad h = 10$$

$$\text{Mode} = 20 + \left(\frac{10 - 6}{2(10) - 6 - 8} \right) \cdot 10 = 20 + \left(\frac{4}{6} \right) \cdot 10 = 20 + 6.67 = 26.67$$

1.3.4 Comparison and When to Use Each

- **Mean** is sensitive to outliers and skewed data. It is best used for symmetric, continuous data without extreme values.
- **Median** is more robust to outliers and skewed distributions. It is ideal when the data contain extreme values or are not symmetrically distributed.
- **Mode** is useful for categorical or discrete data, especially when identifying the most frequent category is of interest.
- **Geometric Mean** is appropriate when dealing with ratios, growth rates, or multiplicative processes (e.g., population growth, interest rates).
- **Harmonic Mean** is best for averaging rates, such as speed or price per unit when quantities vary.

Each measure gives a different perspective of the ‘center’ of the data. The choice of measure should be guided by the nature and scale of the data, and the specific analysis objective.

1.4 Partition Values: Quartiles, Deciles, and Percentiles

Just as the median divides a data set into two equal parts, there are other statistical measures that partition the data into a fixed number of equal segments — such as 4, 10, or 100 parts when the data is arranged in increasing order of magnitude. These measures are collectively referred to as **partition values** or **quantiles**. The most commonly used partition values are the **quartiles**, **deciles**, and **percentiles**, which divide the data into four, ten, and one hundred equal parts, respectively.

Partition values are useful in identifying the spread and concentration of data. For instance, if a student scores at the 90th percentile, they performed better than 90% of the population.

1.4.1 Quartiles

Quartiles divide a ordered data set into four equal parts. There are three quartiles:

- Q_1 (First Quartile): 25% of the data falls below Q_1 .
- Q_2 (Second Quartile or Median): 50% of the data falls below Q_2 .
- Q_3 (Third Quartile): 75% of the data falls below Q_3 .

$$Q_k = \left(\frac{k(n+1)}{4} \right) \text{th value, } k = 1, 2, 3$$

Example: Consider the ordered data: {5, 7, 8, 12, 13, 15, 16, 20, 21}.

- Number of observations $n = 9$
- $Q_1 = \left(\frac{1(9+1)}{4} \right) = 3\text{rd value} = 8$
- $Q_2 = \left(\frac{2(9+1)}{4} \right) = 5\text{th value} = 13$
- $Q_3 = \left(\frac{3(9+1)}{4} \right) = 7\text{th value} = 16$

Interquartile Range (IQR): The Interquartile Range (IQR) is a measure of statistical dispersion, which describes the spread of the middle 50% of a data set. It is the difference between the third quartile (Q_3) and the first quartile (Q_1).

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

1.4.2 Deciles

Deciles divide the ordered data into ten equal parts. There are nine deciles (D_1 to D_9), such that:

$$D_k = \left(\frac{k(n+1)}{10} \right) \text{th value, } k = 1, 2, \dots, 9$$

1.4.3 Percentiles

Percentiles divide the ordered data into one hundred equal parts. There are 99 percentiles (P_1 to P_{99}), commonly used to interpret standardized test scores and similar metrics.

$$P_k = \left(\frac{k(n+1)}{100} \right) \text{th value, } k = 1, 2, \dots, 99$$

1.5 Measures of Dispersion

Measures of dispersion describe the spread or variability within a data set. While measures of central tendency (such as the mean or median) indicate the typical value, measures of dispersion indicate how much the values in the dataset differ from the central value. A small dispersion means the data points are clustered close to the center, while a large dispersion indicates data points are spread out over a wide range.

Consider the following two data sets, each containing five values: $A = \{4, 5, 5, 5, 6\}$ and $B = \{1, 3, 5, 7, 9\}$. Both sets have the same mean, which is 5. For set A , the mean is

$$\frac{4 + 5 + 5 + 5 + 6}{5} = \frac{25}{5} = 5,$$

and for B , the mean is

$$\frac{1 + 3 + 5 + 7 + 9}{5} = \frac{25}{5} = 5.$$

However, their dispersions are quite different. The maximum value of A is 6 and minimum value is 4, whereas the maximum value of B is 9 and minimum value is 1. This means that although both sets center around the same average value, the values in set B are spread out much more widely around the mean compared to set A . Therefore, we say that the dispersion of B is greater than that of A .

Measures of dispersion are broadly classified into two types:

- **Absolute Measures of Dispersion:** These express dispersion in the same units as the original data.
- **Relative Measures of Dispersion:** These express dispersion as a ratio or percentage and are unit-free. They are useful for comparing variability between datasets with different units or magnitudes.

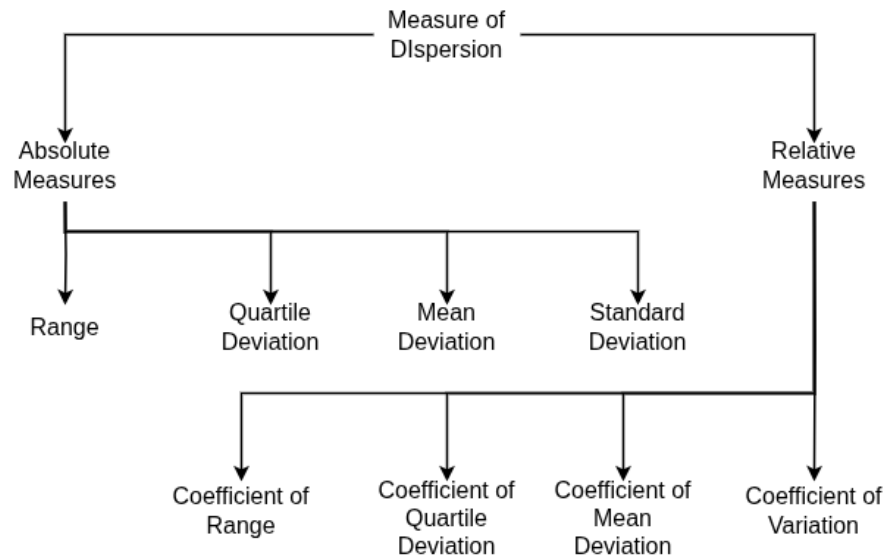


Figure 1.13: *Measures of dispersion.*

1.5.1 Absolute Measures of Dispersion

1. **Range:** Difference between the largest and smallest observations.

$$\text{Range} = L - S$$

where L is the largest value and S is the smallest value.

Consider the data set: $\{10, 15, 18, 22, 25\}$. Here, the largest value $L = 25$ and the smallest value $S = 10$.

$$\text{Range} = 25 - 10 = 15$$

2. **Quartile Deviation (Semi-Interquartile Range):** Quartile deviation is defined as half the difference between the lower and upper quartiles.

$$\text{Quartile Deviation} = \frac{Q_3 - Q_1}{2}$$

where Q_1 and Q_3 are the first and third quartiles.

Consider the ordered data set: $\{10, 15, 20, 25, 30, 35, 40\}$.

Here,

$$Q_1 = 15, \quad Q_3 = 35$$

$$\text{Quartile Deviation} = \frac{35 - 15}{2} = \frac{20}{2} = 10$$

3. **Mean Deviation:** Mean deviation is the arithmetic mean of absolute deviations from mean or any other specified value.

$$\text{Mean Deviation about } A = \frac{1}{n} \sum_{i=1}^n |x_i - A|$$

Generally mean deviation is taken from the arithmetic mean \bar{x} .

$$\text{Mean Deviation about mean} = \frac{1}{n} \sum_{i=1}^n |x_i - A|$$

Consider the data set: $\{2, 4, 6, 8, 10\}$. First we calculate the mean from the data:

$$\bar{x} = \frac{2 + 4 + 6 + 8 + 10}{5} = \frac{30}{5} = 6$$

Then we compute the absolute deviations from the mean:

$$|2 - 6| = 4, \quad |4 - 6| = 2, \quad |6 - 6| = 0, \quad |8 - 6| = 2, \quad |10 - 6| = 4$$

Finally we calculate the Mean Deviation about the mean:

$$\text{Mean Deviation about mean} = \frac{4 + 2 + 0 + 2 + 4}{5} = \frac{12}{5} = 2.4$$

For weighted data, the Mean Deviation about the mean is given by:

$$\text{Mean Deviation about mean} = \frac{\sum_{i=1}^n f_i |x_i - \bar{x}|}{\sum_{i=1}^n f_i}$$

where:

- x_i are the data values,
- f_i are their corresponding frequencies (weights),
- $\bar{x} = \frac{\sum_i f_i x_i}{\sum_i f_i}$ is the weighted mean.

4. **Standard Deviation:** In considering the deviations $x_i - A$ for obtaining a measure of dispersion, we may also get rid of their signs by taking their squares $(x_i - A)^2$, instead of taking their absolute values $|x_i - A|$. The square root of the arithmetic mean of these squares i.e. $\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - A)^2}$ which is called the **root mean square deviation** about A , may be accepted as a measure of dispersion. When $A = \bar{x}$, the measure of dispersion is called the standard deviation.

$$\text{Standard deviation} = \sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

The square of the standard deviation i.e. σ^2 is known as **variance**.

As an example, consider the data set: $\{2, 4, 4, 4, 5, 5, 7, 9\}$. First we calculate the mean:

$$\bar{x} = \frac{2 + 4 + 4 + 4 + 5 + 5 + 7 + 9}{8} = \frac{40}{8} = 5$$

Then we find the Squared Deviations from the mean:

$$(2 - 5)^2 = 9, \quad (4 - 5)^2 = 1, \quad (4 - 5)^2 = 1, \quad (4 - 5)^2 = 1, \\ (5 - 5)^2 = 0, \quad (5 - 5)^2 = 0, \quad (7 - 5)^2 = 4, \quad (9 - 5)^2 = 16$$

Finally compute the Standard Deviation:

$$\sigma = \sqrt{\frac{9 + 1 + 1 + 1 + 0 + 0 + 4 + 16}{8}} = \sqrt{\frac{32}{8}} = \sqrt{4} = 2$$

For weighted data, the standard deviation is calculated as:

$$\text{Standard deviation} = \sigma = \sqrt{\frac{\sum_{i=1}^n f_i (x_i - \bar{x})^2}{\sum_{i=1}^n f_i}}$$

where,

- x_i are the data values,
- f_i are the corresponding frequencies (weights),
- $\bar{x} = \frac{\sum_{i=1}^n f_i x_i}{\sum_{i=1}^n f_i}$ is the weighted mean,

Theorem: If $x = a$ (a constant), then $\sigma_x = 0$.

Proof: Since all observations are equal to a , the mean is

$$\bar{x} = a.$$

The standard deviation is

$$\sigma_x = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} = \sqrt{\frac{1}{n} \sum_{i=1}^n (a - a)^2} = \sqrt{0} = 0.$$

■

Theorem: If $y = a + bx$, where a, b are constants, then

$$\sigma_y = |b| \sigma_x$$

Proof: The mean of y is

$$\bar{y} = a + b\bar{x}$$

The standard deviation of y is

$$\begin{aligned} \sigma_y &= \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2} = \sqrt{\frac{1}{n} \sum_{i=1}^n (a + bx_i - (a + b\bar{x}))^2} \\ &= \sqrt{\frac{1}{n} \sum_{i=1}^n (b(x_i - \bar{x}))^2} = \sqrt{\frac{1}{n} \sum_{i=1}^n b^2 (x_i - \bar{x})^2} \\ &= |b| \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} = |b| \sigma_x \end{aligned}$$

■

Theorem (Pooled standard deviation): Let a dataset be composed of two groups:

- Group 1: n_1 observations, mean \bar{x}_1 , standard deviation σ_1 ,
- Group 2: n_2 observations, mean \bar{x}_2 , standard deviation σ_2 .

Then the combined (pooled) standard deviation σ of the dataset (size $n = n_1 + n_2$) is given by:

$$\sigma = \sqrt{\frac{n_1\sigma_1^2 + n_2\sigma_2^2}{n_1 + n_2} + \frac{n_1(\bar{x}_1 - \bar{x})^2 + n_2(\bar{x}_2 - \bar{x})^2}{n_1 + n_2}}$$

where

$$\bar{x} = \frac{n_1\bar{x}_1 + n_2\bar{x}_2}{n_1 + n_2}$$

Proof: Let us denote $x_{1j}, (j = 1, 2, \dots, n_1)$ and $x_{2j}, (j = 1, 2, \dots, n_2)$ the values of the two sets. Then

$$\sigma_1^2 = \sum_{j=1}^{n_1} (x_{1j} - \bar{x}_1)^2, \quad \sigma_2^2 = \sum_{j=1}^{n_2} (x_{2j} - \bar{x}_2)^2$$

The variance of the combined data set is

$$\sigma^2 = \frac{1}{n_1 + n_2} \left(\sum_{j=1}^{n_1} (x_{1j} - \bar{x})^2 + \sum_{j=1}^{n_2} (x_{2j} - \bar{x})^2 \right)$$

Expanding the first term:

$$\begin{aligned} \sum_{j=1}^{n_1} (x_{1j} - \bar{x})^2 &= \sum_{j=1}^{n_1} [(x_{1j} - \bar{x}_1) + (\bar{x}_1 - \bar{x})]^2 \\ &= \underbrace{\sum_{j=1}^{n_1} (x_{1j} - \bar{x}_1)^2}_{=n_1\sigma_1^2} + 2(\bar{x}_1 - \bar{x}) \underbrace{\sum_{j=1}^{n_1} (x_{1j} - \bar{x}_1)}_{=0} + \sum_{j=1}^{n_1} (\bar{x}_1 - \bar{x})^2 \\ &= n_1\sigma_1^2 + n_1(\bar{x}_1 - \bar{x})^2, \end{aligned}$$

Similarly,

$$\sum_{j=1}^{n_2} (x_{2j} - \bar{x})^2 = n_2\sigma_2^2 + n_2(\bar{x}_2 - \bar{x})^2.$$

Thus,

$$\sigma^2 = \frac{n_1\sigma_1^2 + n_2\sigma_2^2}{n_1 + n_2} + \frac{n_1(\bar{x}_1 - \bar{x})^2 + n_2(\bar{x}_2 - \bar{x})^2}{n_1 + n_2}$$

■

1.5.2 Relative Measures of Dispersion

Relative measures of dispersion are defined as ratio of absolute measures of dispersion to the corresponding measure of central tendency. The ratio is expressed in terms of a percentage.

1. Coefficient of Range:

$$\text{Coefficient of Range} = \frac{L - S}{L + S} \times 100\%$$

where L is the largest value and S is the smallest value.

2. Coefficient of Quartile Deviation:

$$\text{Coefficient of Q.D.} = \frac{Q_3 - Q_1}{Q_3 + Q_1} \times 100\%$$

3. Coefficient of Mean Deviation:

$$\text{Coefficient of M.D.} = \frac{\text{Mean Deviation}}{\bar{x}} \times 100\%$$

4. Coefficient of Variation (CV):

$$\text{CV} = \frac{\sigma}{\bar{x}} \times 100\%$$

1.6 Moments, Skewness and Kurtosis

1.6.1 Raw Moments and Central Moments

In descriptive statistics, **moments** are used to describe various characteristics of a dataset's distribution. Two important types of moments are:

- **Raw moments** (moments about the origin): The r -th **raw moment** of a dataset x_1, x_2, \dots, x_n is given by:

$$\mu'_r = \frac{1}{n} \sum_{i=1}^n x_i^r$$

- First raw moment: $\mu'_1 = \bar{x}$ (sample mean)
- Second raw moment: $\mu'_2 = \frac{1}{n} \sum x_i^2$

- **Central moments** (moments about the mean): The r -th **central moment** is the average of the r -th powers of deviations from the mean:

$$\mu_r = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^r$$

- First central moment: $\mu_1 = 0$ (since the mean deviation is zero)
- Second central moment: $\mu_2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ is the variance.

Example: Consider the dataset:

$$\{2, 4, 6, 8\}$$

- Raw moments:

$$\begin{aligned}\mu'_1 &= \frac{1}{4}(2 + 4 + 6 + 8) = 5 \\ \mu'_2 &= \frac{1}{4}(2^2 + 4^2 + 6^2 + 8^2) = \frac{120}{4} = 30\end{aligned}$$

- Central moments:

$$\begin{aligned}\mu_1 &= 0 \\ \mu_2 &= \frac{1}{4}[(2-5)^2 + (4-5)^2 + (6-5)^2 + (8-5)^2] \\ &= \frac{1}{4}(9 + 1 + 1 + 9) = \frac{20}{4} = 5\end{aligned}$$

1.6.2 Relationship Between Raw and Central Moments

The r -th central moment μ_r can be expressed in terms of raw moments μ'_k and powers of the mean \bar{x} :

$$\begin{aligned}\mu_r &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^r = \frac{1}{n} \sum_{i=1}^n \left[\sum_{k=0}^r \binom{r}{k} \bar{x}_i^{r-k} \bar{x}^k \right] \\ &= \frac{1}{n} \sum_{i=1}^n \left[x_i^r - \binom{r}{1} x_i^{r-1} \bar{x} + \binom{r}{2} x_i^{r-2} \bar{x}^2 + \cdots - \bar{x}^r \right] \\ &= \mu'_r - \binom{r}{1} \mu'_{r-1} \bar{x} + \binom{r}{2} \mu'_{r-2} \bar{x}^2 + \cdots - n \bar{x}^r\end{aligned}$$

Since $\mu'_1 = \bar{x}$, we can rewrite the above expression as:

$$\mu_r = \mu'_r - \binom{r}{1} \mu'_{r-1} \mu'_1 + \binom{r}{2} \mu'_{r-2} \mu_1'^2 + \cdots - n \mu_1'^r$$

- $r = 1$:

$$\begin{aligned}\mu_1 &= \mu'_1 - \binom{1}{1} \mu'_0 \mu'_1 \\ &= \mu'_1 - \mu'_1 = 0\end{aligned}$$

- $r = 2$:

$$\begin{aligned}\mu_2 &= \mu'_2 - \binom{2}{1} \mu'_1 \mu'_1 + \binom{2}{2} \mu'_0 (\mu'_1)^2 \\ &= \mu'_2 - 2(\mu'_1)^2 + (\mu'_1)^2 = \mu'_2 - (\mu'_1)^2\end{aligned}$$

- $r = 3$:

$$\begin{aligned}\mu_3 &= \mu'_3 - \binom{3}{1} \mu'_2 \mu'_1 + \binom{3}{2} \mu'_1 (\mu'_1)^2 - \binom{3}{3} \mu'_0 (\mu'_1)^3 \\ &= \mu'_3 - 3 \mu'_2 \mu'_1 + 3(\mu'_1)^3 - (\mu'_1)^3 = \mu'_3 - 3 \mu'_2 \mu'_1 + 2(\mu'_1)^3\end{aligned}$$

- $r = 4$:

$$\begin{aligned}\mu_4 &= \mu'_4 - \binom{4}{1} \mu'_3 \mu'_1 + \binom{4}{2} \mu'_2 (\mu'_1)^2 - \binom{4}{3} \mu'_1 (\mu'_1)^3 + \binom{4}{4} \mu'_0 (\mu'_1)^4 \\ &= \mu'_4 - 4 \mu'_3 \mu'_1 + 6 \mu'_2 (\mu'_1)^2 - 4(\mu'_1)^4 + (\mu'_1)^4 \\ &= \mu'_4 - 4 \mu'_3 \mu'_1 + 6 \mu'_2 (\mu'_1)^2 - 3(\mu'_1)^4\end{aligned}$$

Now let's derive the inverse relationship. Using the binomial expansion on $x_i = (x_i - \bar{x}) + \bar{x}$, the r -th raw moment can be written as

$$\begin{aligned}\mu'_r &= \frac{1}{n} \sum_{i=1}^n [\bar{x} + (x_i - \bar{x})]^r = \frac{1}{n} \sum_{i=1}^n \left[\sum_{k=0}^r \binom{r}{k} \bar{x}^{r-k} (x_i - \bar{x})^k \right] \\ &= \frac{1}{n} \sum_{i=1}^n \left[\bar{x}^r + \binom{r}{1} \bar{x}^{r-1} (x_i - \bar{x}) + \binom{r}{2} \bar{x}^{r-2} (x_i - \bar{x})^2 + \cdots + (x_i - \bar{x})^r \right] \\ &= \bar{x}^r + \binom{r}{1} \bar{x}^{r-1} \mu_1 + \binom{r}{2} \bar{x}^{r-2} \mu_2 + \cdots + \mu_r\end{aligned}$$

where by convention $\mu_1 = 0$.

$$\mu'_r = \bar{x}^r + \binom{r}{1} \bar{x}^{r-1} \mu_1 + \binom{r}{2} \bar{x}^{r-2} \mu_2 + \cdots + \mu_r$$

- $r = 1$:

$$\begin{aligned}\mu'_1 &= \bar{x}^1 + \binom{1}{1} \bar{x}^0 \mu_1 \\ &= \bar{x} + 0 = \bar{x}\end{aligned}$$

- $r = 2$:

$$\begin{aligned}\mu'_2 &= \bar{x}^2 + \binom{2}{1} \bar{x}^1 \mu_1 + \binom{2}{2} \bar{x}^0 \mu_2 \\ &= \bar{x}^2 + 0 + \mu_2 = \mu_2 + \bar{x}^2\end{aligned}$$

- $r = 3$:

$$\begin{aligned}\mu'_3 &= \bar{x}^3 + \binom{3}{1} \bar{x}^2 \mu_1 + \binom{3}{2} \bar{x}^1 \mu_2 + \binom{3}{3} \bar{x}^0 \mu_3 \\ &= \bar{x}^3 + 0 + 3 \bar{x} \mu_2 + \mu_3 = \mu_3 + 3 \bar{x} \mu_2 + \bar{x}^3\end{aligned}$$

- $r = 4$:

$$\begin{aligned}\mu'_4 &= \bar{x}^4 + \binom{4}{1} \bar{x}^3 \mu_1 + \binom{4}{2} \bar{x}^2 \mu_2 + \binom{4}{3} \bar{x}^1 \mu_3 + \binom{4}{4} \bar{x}^0 \mu_4 \\ &= \bar{x}^4 + 0 + 6 \bar{x}^2 \mu_2 + 4 \bar{x} \mu_3 + \mu_4 \\ &= \mu_4 + 4 \bar{x} \mu_3 + 6 \bar{x}^2 \mu_2 + \bar{x}^4\end{aligned}$$

1.6.3 Skewness

Skewness is a measure of the asymmetry of a frequency distribution about its mean. The frequency distribution of a dataset is called symmetrical about the value x_0 if the frequency of $x_0 - h$ is same as the frequency of $x_0 + h$, whatever h may be.

The sample skewness is defined as:

$$\text{Skewness}(\gamma_1) = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^3 = \frac{\mu_3}{\sigma^3}$$

where $\mu_3 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3$ is the third central moment.

The value of skewness determines the shape of the frequency curve:

- If skewness = 0, the distribution is **symmetric**.
- If skewness > 0, the distribution is **positively skewed** (long right tail).
- If skewness < 0, the distribution is **negatively skewed** (long left tail).

How to interpret the formula?

Skewness uses cubed deviations $(x_i - \bar{x})^3$. Cubing serves two purposes: it preserves the sign of the deviation — meaning values greater than the mean contribute positively and those less than the mean contribute negatively — and it exaggerates the impact of larger deviations, making the measure sensitive to extreme values in the tails. This helps identify whether the data are stretched more to the right or left.

Dividing by the cube of the standard deviation σ^3 standardizes the measure, removing units and allowing for meaningful comparisons across datasets with different scales. The result is a dimensionless quantity: positive skewness indicates a long right tail, negative skewness signals a long left tail, and zero skewness implies symmetry around the mean.



Example: Given data: $x = \{2, 3, 4, 5, 8\}$

- Mean: $\bar{x} = \frac{2 + 3 + 4 + 5 + 8}{5} = 4.4$
- Standard deviation: $s = \sqrt{\frac{1}{5} \sum_i (x_i - \bar{x})^2} \approx 2.058$
- Third central moment:

$$\mu_3 = \frac{1}{5} [(-2.4)^3 + (-1.4)^3 + (-0.4)^3 + 0.6^3 + 3.6^3] \approx 3.232$$

- Skewness:

$$\gamma_1 = \frac{3.232}{(2.058)^3} \approx 0.37$$

Since skewness > 0, the distribution is **positively skewed**.

In most unimodal distributions¹, the following “rule of thumb” holds regarding the ordering of Mean, Median, and Mode under skewness:

- **Positive skew (right-tailed):**

$$\text{Mode} < \text{Median} < \text{Mean}.$$

Extreme values on the right pull the mean farther out than the median, while the mode remains at the peak of the bulk of the data.

¹A dataset is said to have a **unimodal distribution** if its values tend to cluster around a single (not multiple) central peak when plotted as a histogram or a frequency curve.

- **Negative skew (left-tailed):**

$$\text{Mean} < \text{Median} < \text{Mode}.$$

Extreme values on the left drag the mean below the median, and the mode stays at the highest-density point on the right.

1.6.4 Kurtosis

Kurtosis measures the degree of ‘peakedness’ of a frequency distribution curve. Two frequency distributions may have the same mean, dispersion, and skewness, yet differ in how concentrated the values are around the mode. One distribution may have a sharper peak due to a higher concentration of values near the center, while the other appears flatter. This characteristic of a frequency distribution is known as kurtosis. It is calculated and reported either as an absolute or as a relative value. The absolute kurtosis is always a positive number.

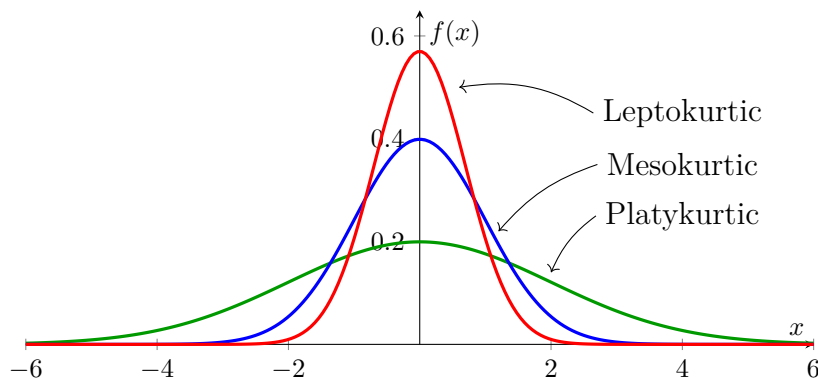
$$\text{Absolute Kurtosis} = \frac{\mu_4}{\sigma^4}$$

where $\mu_4 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4$ is the fourth central moment.

The absolute kurtosis of a normal distribution, which we will learn in later chapter, is 3. The value 3 is taken as a datum (reference point) to calculate the relative kurtosis.

$$\text{Absolute Kurtosis}(\gamma_2) = \text{Relative Kurtosis} - 3$$

- Relative kurtosis = 0: **Mesokurtic** (e.g. Normal).
- Relative kurtosis > 0: **Leptokurtic** (heavy tails, sharp peak).
- Relative kurtosis < 0: **Platykurtic** (light tails, flat top).



How to interpret the formula?

Kurtosis raises deviations from the mean to the fourth power. This has a distinct purpose: it emphasizes extreme values far from the mean far more than values closer to it. Unlike cubing (used in skewness), which preserves the sign of deviations to detect asymmetry, raising to the fourth power removes the sign, treating all deviations equally, but magnifying larger ones disproportionately.

Before applying the fourth power, each deviation is **first divided by the standard deviation σ** . This step is crucial: it standardizes the scale of deviations, ensuring that the measure reflects the *relative extremity* of values, not just their raw magnitude. Even if two distributions seem to have

similar tail weights, the distribution with a *smaller standard deviation* (i.e., a tighter central cluster) will yield *larger standardized deviations*, which get exaggerated further by the fourth power.

Also, dividing the fourth central moment by σ^4 makes kurtosis a **dimensionless** and **scale-invariant** quantity, allowing meaningful comparisons across datasets.

Example: Using the same data $x = \{2, 3, 4, 5, 8\}$:

- Mean $\bar{x} = 4.4$, Standard deviation $\sigma \approx 2.058$.
- Fourth central moment:

$$\mu_4 = \frac{1}{5} [(-2.4)^4 + (-1.4)^4 + (-0.4)^4 + 0.6^4 + 3.6^4] \approx 41.03.$$

- Kurtosis(γ_2) = $\frac{41.03}{(2.058)^4} \approx 2.29$, Relative Kurtosis ≈ -0.71 . This dataset is **platykurtic**.

Chapter 2

Theory of Probability

2.1 Some Notation and Terminology

2.1.1 Random Experiment

An **experiment** is generally defined as one or more actions that result in a specific outcome.

An experiment E is called a **random experiment** if it satisfies the following conditions:

- All possible outcomes of E are known in advance.
- It is not possible to predict with certainty which specific outcome will occur in any given trial^a of E .
- The experiment E can be repeated, at least conceptually, under identical conditions an infinite number of times.

^aA **trial** is a single performance or execution of an experiment. Tossing a coin once is a trial of the coin-tossing experiment.

A common example of a random experiment is the tossing of a coin. The possible outcomes—‘Head’ and ‘Tail’—are known in advance, but it is impossible to determine with certainty which of the two outcomes will occur on any single toss.

2.1.2 Event Space (a.k.a Sample Space)

The set of all possible outcomes of a random experiment E is called the **sample space** or **event space**, and it is denoted by S .

Each outcome, also known as an **elementary event point**, is an element of S .

For example, in the experiment of tossing a coin, the sample space is

$$S = \{H, T\}$$

where H represents ‘Head’ and T represents ‘Tail’.

If E is the experiment of rolling a pair of dice, the sample space consists of all ordered pairs of numbers from 1 to 6:

$$S = \{(1, 1), (1, 2), (1, 3), \dots, (6, 6)\}$$

This sample space contains 36 distinct outcomes, as each die can show any of 6 faces independently.

A sample space is **discrete** if it consists of a finite or countable infinite set of outcomes. A sample space is **continuous** if it contains an interval (either finite or infinite) of real numbers.

The sample space $S = \mathbb{R}^+$ is an example of a continuous sample space, whereas $S = \{H, T\}$ is a discrete sample space.

2.1.3 Events

We often focus on groups of related outcomes from a random experiment, which are represented as subsets of the sample space.

A subset of a sample space is called an **event**.

Consider the random experiment of rolling a die. The sample space is

$$S = \{1, 2, 3, 4, 5, 6\}$$

Let

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

be an event, which can be described as “an even number appears when the die is rolled.”

There are various types of events:

- **Impossible Event:** An event that contains no outcomes from the sample space is called an impossible event. For example, $A = \emptyset$ is an impossible event.
- **Certain Event:** An event that contains all outcomes of the sample space is called a certain or sure event. For example, $A = S$ is an impossible event.
- **Simple (Elementary) Event:** An event consisting of exactly one outcome of the sample space. For example, $A = \{4\}$ is a simple event when rolling a die.
- **Composite (Compound) Event:** An event that consists of more than one outcome of the sample space. For example, $A = \{2, 4, 6\}$ is a composite event when rolling a die.
- **Dependent and Independent Events:** Two events are considered dependent if the occurrence of one event influences the probability of the other event occurring. Conversely, they are independent if the occurrence of one event does not affect the probability of the other event.

2.1.4 Mutually Exclusive Events

Two events are said to be **mutually exclusive** if they cannot occur at the simultaneously. Mathematically, if events A_1 and A_2 are exhaustive, then:

$$A_1 \cap A_2 = \emptyset$$

When tossing a coin, the events ‘Head’ and ‘Tail’ are mutually exclusive because both cannot occur at the same time. If we get a head, we cannot get a tail in that toss.

2.1.5 Exhaustive Set of Events

A set of events is said to be **exhaustive** at least one of the events in the set must occur. The union of all the events in the set equals the entire sample space S . Mathematically, if events A_1, A_2, \dots, A_n are exhaustive, then:

$$A_1 \cup A_2 \cup \dots \cup A_n = S$$

When rolling a die, the events $\{1\}, \{2\}, \{3\}, \{4\}, \{5\}$, and $\{6\}$ form an exhaustive set because they cover all possible outcomes of the die roll. One of these events must occur when the die is thrown.

2.2 Definition of Probability

2.2.1 A Priori Probability

A priori probability, also known as **classical probability**, is the probability that is determined before an experiment is conducted. It is based on the knowledge of the system or experiment and is calculated using the total number of equally likely outcomes.

If there are n mutually exclusive, exhaustive and equally likely^a outcomes of a random experiment and out of them m outcomes are favorable to an event A , then the probability of the event A is defined as

$$P(A) = \frac{m}{n}$$

^aBy the phrase ‘**equally likely**’ it is meant that none of the outcomes is expected to occur in preference to other in any trial of the given random experiment.

For example, the probability of getting a head in a fair coin toss is

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

based on the assumption of equal likelihood of heads and tails.

2.2.2 A Posteriori Probability

A posteriori probability, also known as **empirical probability**, is the probability that is determined after an experiment is conducted. It is based on observed data or information obtained from the experiment.

Let A be an event of a given random experiment. Let event A occurs $N(A)$ number of times when the random experiment is repeated N times under identical conditions. The probability of the event A is defined as

$$P(A) = \lim_{N \rightarrow \infty} \frac{N(A)}{N}$$

A posteriori probability can be updated as new evidence becomes available. For example, after observing several rolls of a die, you may update the probability of rolling a particular number based on the outcomes observed.

2.3 Axioms of Probability

The subject of probability is based on three commonsense rules, known as axioms. They are:

- **First Axiom:** $P(S) = 1$ where S is the sample space.
- **Second Axiom:** $0 \leq P(E) \leq 1$ for any event E .
- **Third Axiom:** For two mutually exclusive events E_1 and E_2 ,

$$P(E_1 \cup E_2) = P(E_1) + P(E_2)$$

More generally, if E_1, E_2, \dots, E_n are mutually exclusive events,

$$P(E_1 \cup E_2 \cup \dots \cup E_n) = P(E_1) + P(E_2) + \dots + P(E_n)$$

The **first axiom** states that the probability of the entire sample space S , which represents all possible outcomes of an experiment, is 1. It reflects the certainty that something in the sample space will occur. For example, when flipping a fair coin, the sample space is $S = \{\text{Head}, \text{Tail}\}$, and the probability that the outcome is either ‘Head’ or ‘Tail’ is 1.

The **second axiom** ensures that probabilities are valid numerical values between 0 and 1. A probability of 0 means the event is impossible (e.g., rolling a 7 on a standard six-sided die), while a probability of 1 means the event is certain to happen. All other events fall somewhere in between these two extremes.

The **third axiom** applies when two events E_1 and E_2 are mutually exclusive—they cannot both happen at the same time. In such cases, the probability that either event occurs is the sum of their individual probabilities. For instance, when rolling a die, the probability of getting a 2 or a 5 is $P(2) + P(5)$, since a single die roll cannot result in both values.

These axioms imply the following theorems.

Theorem: $P(\bar{E}) = 1 - P(E)$ for any event E

Proof: Let S be a sample space and let E be an event. Then E and \bar{E} are mutually exclusive. So by axiom 3,

$$P(E \cup \bar{E}) = P(E) + P(\bar{E})$$

But $E \cup \bar{E} = S$, and by axiom 1, $P(S) = 1$. Therefore,

$$P(E) + P(\bar{E}) = 1$$

which implies

$$P(\bar{E}) = 1 - P(E)$$

■

Theorem: $P(\emptyset) = 0$ where \emptyset denotes the empty set.

Proof: Let S be a sample space. Then $\emptyset = \bar{S}$. Therefore

$$P(\emptyset) = 1 - P(S) = 1 - 1 = 0$$

■

Theorem: For any two events A and B (not necessarily mutually exclusive),

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

Proof: From the Venn diagram, we can see that the event $A \cup B$ consists of three mutually exclusive events (subsets) $A \cap \bar{B}$, $A \cap B$ and $\bar{A} \cap B$.



The event A consists of two mutually exclusive events $A \cap \bar{B}$ and $A \cap B$. Therefore

$$P(A) = P(A \cap \bar{B}) + P(A \cap B)$$

Similarly,

$$P(B) = P(\bar{A} \cap B) + P(A \cap B)$$

Now,

$$\begin{aligned} P(A \cup B) &= P(A \cap \bar{B}) + P(A \cap B) + P(\bar{A} \cap B) \\ &= [P(A) - P(A \cap B)] + P(A \cap B) + [P(B) - P(A \cap B)] \\ &= P(A) + P(B) - P(A \cap B) \end{aligned}$$

■

Example: Consider a fair six-sided die, and define two events:

- A : The event that the die shows an odd number (i.e., $A = \{1, 3, 5\}$)
- B : The event that the die shows a number greater than or equal to 4 (i.e., $A = \{4, 5, 6\}$)

The union of A and B is the event that either event A or event B occurs (or both). The union of A and B is denoted by:

$$A \cup B = \{1, 3, 5, 4, 6\} = \{1, 3, 4, 5, 6\}$$

To find the probability of the union $P(A \cup B)$, we use the formula:

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

- $P(A) = \frac{3}{6} = 0.5$ (since there are 3 odd numbers: 1, 3, 5)
- $P(B) = \frac{3}{6} = 0.5$ (since there are 3 numbers ≥ 4 : 4, 5, 6)
- $P(A \cap B) = \frac{1}{6} = \frac{1}{3}$ (since 5 is the only number that is both odd and ≥ 4)

So, the probability of $A \cup B$ is:

$$P(A \cup B) = 0.5 + 0.5 - \frac{1}{3} = 1 - \frac{1}{3} = \frac{2}{3}$$

Thus, the probability of rolling a die and getting either an odd number or a number greater than or equal to 4 is $\frac{2}{3}$.

Theorem: For any three events A , B and C (not necessarily mutually exclusive),

$$\begin{aligned} P(A \cup B \cup C) &= P(A) + P(B) + P(C) - P(A \cap B) - P(B \cap C) \\ &\quad - P(A \cap C) + P(A \cap B \cap C) \end{aligned}$$

Proof: We begin by applying the two-event formula to $A \cup (B \cup C)$:

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

Now, apply the two-event formula to $P(B \cup C)$:

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

Also, apply distributivity to expand $P(A \cap (B \cup C))$:

$$\begin{aligned} P(A \cap (B \cup C)) &= P((A \cap B) \cup (A \cap C)) \\ &= P(A \cap B) + P(A \cap C) - P(A \cap B \cap C) \end{aligned}$$

Substituting back into the original expression:

$$\begin{aligned} P(A \cup B \cup C) &= P(A) + [P(B) + P(C) - P(B \cap C)] \\ &\quad - [P(A \cap B) + P(A \cap C) - P(A \cap B \cap C)] \\ &= P(A) + P(B) + P(C) - P(B \cap C) \\ &\quad - P(A \cap B) - P(A \cap C) + P(A \cap B \cap C) \end{aligned}$$

■

Theorem: Let A_1, A_2, \dots, A_n be n number of events of a random experiment. Then the probability of their union is given by:

$$\begin{aligned} P\left(\bigcup_{i=1}^n A_i\right) &= \sum_{i=1}^n P(A_i) - \sum_{1 \leq i < j \leq n} P(A_i \cap A_j) \\ &+ \sum_{1 \leq i < j < k \leq n} P(A_i \cap A_j \cap A_k) - \dots \\ &+ (-1)^{n+1} P(A_1 \cap A_2 \cap \dots \cap A_n) \end{aligned}$$

The proof of this theorem is left as an exercise.

Theorem: If A and B are two events of a random experiment, then the probability that exactly one of them occurs is given by

$$P(\text{exactly one of } A \text{ or } B) = P(A) + P(B) - 2P(A \cap B)$$

Proof: The event “exactly one of A or B occurs” means either A happens and B doesn’t, or B happens and A doesn’t. That means the event:

$$(A \cap \bar{B}) \cup (\bar{A} \cap B)$$



We can see it from the Venn diagram:

$$\begin{aligned} P((A \cap \bar{B}) \cup (\bar{A} \cap B)) &= P(A \cup B) - P(A \cap B) \\ &= (P(A) + P(B) - P(A \cap B)) - P(A \cap B) \\ &= P(A) + P(B) - 2P(A \cap B) \end{aligned}$$

Therefore,

$$P(\text{exactly one of } A \text{ or } B) = P(A) + P(B) - 2P(A \cap B)$$

■

Boole’s Inequality: Let A_1, A_2, \dots, A_n be n events of a random experiment. Then:

$$P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i)$$

Proof: We prove the inequality by induction on n . Base Case: For $n = 1$,

$$P(A_1) = P(A_1)$$

so the inequality holds with equality. Inductive Step: Assume the inequality holds for $n = k$, i.e.,

$$P\left(\bigcup_{i=1}^k A_i\right) \leq \sum_{i=1}^k P(A_i)$$

Now consider $n = k + 1$. Let $B = \bigcup_{i=1}^k A_i$. Then:

$$P\left(\bigcup_{i=1}^{k+1} A_i\right) = P(B \cup A_{k+1})$$

Using the formula for the union of two events:

$$P(B \cup A_{k+1}) = P(B) + P(A_{k+1}) - P(B \cap A_{k+1}) \leq P(B) + P(A_{k+1})$$

Applying the induction hypothesis to $P(B)$:

$$P(B \cup A_{k+1}) \leq \sum_{i=1}^k P(A_i) + P(A_{k+1}) = \sum_{i=1}^{k+1} P(A_i)$$

Thus, the inequality holds for $n = k + 1$. By the principle of mathematical induction, the inequality holds for all $n \in \mathbb{N}$. ■

Boole's inequality provides a simple and conservative upper bound for the probability of the union of multiple events. This is important because, without detailed knowledge of the relationships between the events (e.g., how much they overlap), we can still estimate the probability that at least one of the events occurs by adding their individual probabilities.

Example: If you have three events with probabilities:

$$P(A_1) = 0.3, \quad P(A_2) = 0.5, \quad P(A_3) = 0.7,$$

but you don't know the intersections between them, Boole's inequality will tell you that the probability of at least one occurring is at most:

$$P(A_1 \cup A_2 \cup A_3) \leq 0.3 + 0.5 + 0.7 = 1.5$$

Since probabilities cannot exceed 1, this shows that the bound is very loose, but it is still useful for getting a rough estimate.

Bonferroni's inequality: Let A_1, A_2, \dots, A_n be events of a random experiment. Then:

$$P\left(\bigcap_{i=1}^n A_i\right) \geq \sum_{i=1}^n P(A_i) - (n - 1)$$

Proof: We proceed by induction on n .

Base case: For $n = 2$, we begin with the identity:

$$P(A_1 \cup A_2) = P(A_1) + P(A_2) - P(A_1 \cap A_2)$$

Using the axiom that $P(A_1 \cup A_2) \leq 1$, we substitute:

$$P(A_1) + P(A_2) - P(A_1 \cap A_2) \leq 1$$

Rearranging:

$$P(A_1 \cap A_2) \geq P(A_1) + P(A_2) - 1$$

So the inequality holds for $n = 2$. Inductive step: Assume the result holds for $n = k$, i.e.,

$$P\left(\bigcap_{i=1}^k A_i\right) \geq \sum_{i=1}^k P(A_i) - (k - 1)$$

Let $B = \bigcap_{i=1}^k A_i$. Then:

$$P\left(\bigcap_{i=1}^{k+1} A_i\right) = P(B \cap A_{k+1}) \geq P(B) + P(A_{k+1}) - 1$$

Using the induction hypothesis:

$$P(B \cap A_{k+1}) \geq \left(\sum_{i=1}^k P(A_i) - (k-1) \right) + P(A_{k+1}) - 1 = \sum_{i=1}^{k+1} P(A_i) - k$$

Therefore, the inequality holds for $n = k + 1$. By induction, it holds for all $n \in \mathbb{N}$.

The Bonferroni inequality for intersections provides a lower bound for the probability of the intersection of multiple events. ■

Example: Consider three events A_1 , A_2 , and A_3 with probabilities:

$$P(A_1) = 0.6, \quad P(A_2) = 0.5, \quad P(A_3) = 0.7$$

Using Bonferroni's inequality, the lower bound for the probability that all three events occur is:

$$P(A_1 \cap A_2 \cap A_3) \geq 0.6 + 0.5 + 0.7 - 2 = 0.8$$

Thus, the probability that all three events occur simultaneously is at least 0.8.

2.4 Conditional Probability

Conditional probability is the probability of an event occurring given that another event has already occurred.

Let A and B be two events of a random experiment. The **conditional probability** of event A given that event B has already occurred is defined as:

$$P(A | B) = \frac{P(A \cap B)}{P(B)}, \quad \text{provided } P(B) \neq 0$$

This formula tells us how likely event A is, given that we know event B has happened. The idea is that we restrict our sample space to the outcomes where B occurs, and then we compute the probability of A within this restricted sample space.

Example: Let's define two events when rolling a fair six-sided die:

- A : The event that the die shows an *even face*, i.e., $A = \{2, 4, 6\}$.
- B : The event that the die shows a *multiple of 3*, i.e., $B = \{3, 6\}$.

$$\begin{aligned} P(A) &= P(\{2, 4, 6\}) = \frac{3}{6} = \frac{1}{2} \\ P(B) &= P(\{3, 6\}) = \frac{2}{6} = \frac{1}{3} \\ P(A \cap B) &= P(\{6\}) = \frac{1}{6} \end{aligned}$$

The conditional probability of A given B is:

$$P(A | B) = \frac{P(A \cap B)}{P(B)} = \frac{\frac{1}{6}}{\frac{1}{3}} = \frac{1}{2}$$

The conditional probability of B given A is:

$$P(B | A) = \frac{P(A \cap B)}{P(A)} = \frac{\frac{1}{6}}{\frac{1}{2}} = \frac{1}{3}$$

Multiplication Rule of Probabilities: If A and B are any events in the sample space S , then:

$$\begin{aligned} P(A \cap B) &= P(A) \cdot P(B | A), \quad \text{if } P(A) \neq 0 \\ &= P(B) \cdot P(A | B), \quad \text{if } P(B) \neq 0 \end{aligned}$$

The second rule follows directly from the definition of conditional probability by multiplying both sides by $P(B)$. The first rule is obtained from the second by simply switching the roles of A and B .

2.5 Rule of Total Probability

The Rule of Total Probability allows us to compute the probability of an event based on a partition of the sample space.

2.5.1 For Two Events

If B and its complement \bar{B} form a partition of the sample space (i.e., mutually exclusive and collectively exhaustive events), then for any event A :

$$P(A) = P(A | B)P(B) + P(A | \bar{B})P(\bar{B})$$

Proof: Let B and \bar{B} form a partition (i.e., mutually exclusive and exhaustive events) of the sample space. Then any event A can be expressed as:

$$A = (A \cap B) \cup (A \cap \bar{B})$$

Since the sets $A \cap B$ and $A \cap \bar{B}$ are disjoint, the two parts are mutually exclusive. So:

$$P(A) = P(A \cap B) + P(A \cap \bar{B})$$

Using the definition of conditional probability:

$$P(A \cap B) = P(A | B)P(B), \quad P(A \cap \bar{B}) = P(A | \bar{B})P(\bar{B})$$

Substituting:

$$P(A) = P(A | B)P(B) + P(A | \bar{B})P(\bar{B})$$

■

2.5.2 For Multiple Events

Let B_1, B_2, \dots, B_n be a partition of the sample space (i.e., mutually exclusive and collectively exhaustive events). Then for any event A :

$$P(A) = \sum_{i=1}^n P(A | B_i)P(B_i)$$

Proof: Let B_1, B_2, \dots, B_n be a partition of the sample space, i.e.,

- The events B_i are mutually exclusive: $B_i \cap B_j = \emptyset$ for $i \neq j$
- The events B_i are collectively exhaustive: $\bigcup_{i=1}^n B_i = S$

Then for any event $A \subseteq S$:

$$A = A \cap S = A \cap \left(\bigcup_{i=1}^n B_i \right) = \bigcup_{i=1}^n (A \cap B_i)$$

Since the B_i are disjoint, so are the $A \cap B_i$, so:

$$P(A) = \sum_{i=1}^n P(A \cap B_i)$$

Using conditional probability:

$$P(A \cap B_i) = P(A | B_i)P(B_i)$$

Therefore:

$$P(A) = \sum_{i=1}^n P(A | B_i)P(B_i)$$

■

Example

Suppose a factory has three machines:

- Machine M_1 produces 30% of the items, with a defect rate of 1%.
- Machine M_2 produces 50% of the items, with a defect rate of 2%.
- Machine M_3 produces 20% of the items, with a defect rate of 3%.

Let D be the event that an item is defective. We are asked to find $P(D)$, the total probability that a randomly selected item is defective.

Using the Rule of Total Probability:

$$P(D) = P(D | M_1)P(M_1) + P(D | M_2)P(M_2) + P(D | M_3)P(M_3)$$

Substituting the known values:

$$\begin{aligned} P(D) &= (0.01 \times 0.30) + (0.02 \times 0.50) + (0.03 \times 0.20) \\ &= 0.003 + 0.010 + 0.006 = 0.019 \end{aligned}$$

The probability that a randomly chosen item is defective is 0.019 or 1.9%.

2.6 Bayes' Theorem

Bayes' Theorem is a fundamental result in probability theory that arises directly from the definition of conditional probability and the Rule of Total Probability.

Bayes' Theorem: Let B_1, B_2, \dots, B_n be a partition of the sample space (i.e., mutually exclusive and collectively exhaustive events) and none of which has zero probability i.e. $P(B_i) > 0$ for all i , then for any event A with $P(A) > 0$, the probability of B_r given A is:

$$P(B_r | A) = \frac{P(B_r) \cdot P(A | B_r)}{\sum_{i=1}^n P(B_i) \cdot P(A | B_i)}$$

for $r = 1, 2, \dots, n$.

Proof: Let B_1, B_2, \dots, B_n be a partition of the sample space S such that:

- The events B_i are mutually exclusive: $B_i \cap B_j = \emptyset$ for $i \neq j$,
- The union of the B_i 's covers the whole sample space: $\bigcup_{i=1}^n B_i = S$,
- $P(B_i) > 0$ for all i .

Let A be any event with $P(A) > 0$. By the definition of conditional probability:

$$P(B_r | A) = \frac{P(B_r \cap A)}{P(A)}$$

We apply the Multiplication Rule of Probabilities:

$$P(B_r \cap A) = P(B_r) \cdot P(A | B_r)$$

So:

$$P(B_r | A) = \frac{P(B_r) \cdot P(A | B_r)}{P(A)}$$

Now, using the Rule of Total Probability:

$$P(A) = \sum_{i=1}^n P(B_i) \cdot P(A | B_i)$$

Substitute this into the denominator:

$$P(B_r | A) = \frac{P(B_r) \cdot P(A | B_r)}{\sum_{i=1}^n P(B_i) \cdot P(A | B_i)}$$

■

Example

Suppose in a dataset of 1000 emails, 200 are identified as spam and 800 as non-spam. Among the spam emails, 80 contain the word 'discount', while 80 of the non-spam emails also contain this word. We need to calculate the probability that an email is spam given that it contains the word 'discount'.

- $P(S) = \frac{200}{1000} = 0.2$: Probability that an email is spam.
- $P(\bar{S}) = \frac{800}{1000} = 0.8$: Probability that an email is not spam.
- $P(D | S) = \frac{80}{200} = 0.4$: Probability that the word "discount" appears in a spam email.
- $P(D | \bar{S}) = \frac{80}{800} = 0.1$: Probability that "discount" appears in a non-spam email.

We want to find the probability that an email is spam given that it contains the word "discount", i.e., $P(S | D)$.

Using Bayes' Theorem:

$$P(S | D) = \frac{P(S) \cdot P(D | S)}{P(S) \cdot P(D | S) + P(\bar{S}) \cdot P(D | \bar{S})}$$

Substituting the values:

$$P(S | D) = \frac{0.2 \times 0.4}{0.2 \cdot 0.4 + 0.8 \cdot 0.1} = \frac{0.08}{0.08 + 0.08} = \frac{0.08}{0.16} = 0.5$$

So, the probability that the email is spam given it contains the word "discount" is 50%. Even though only 20% of all emails are spam, once we see the word "discount" the chance the email is spam rises to 50%, because that word is much more common in spam than in legitimate messages.

—

2.6.1 Importance of Bayes' Theorem and Updating Probability

Bayes' Theorem is important because it provides a mathematical framework for updating probabilities when new information becomes available. In real-world terms, it helps us refine our beliefs or predictions as we gather more data.

Suppose we want to determine the probability of an event A , such as an email being spam. Initially, we rely on prior knowledge, which is represented by the prior probability $P(A)$. Now, imagine we observe some new evidence B , such as the presence of a specific word like “discount” in the email.

Bayes' Theorem helps us compute the updated probability $P(A | B)$, known as the posterior probability, using the following formula:

$$P(A | B) = \frac{P(A) \cdot P(B | A)}{P(B)}$$

Here:

- $P(A)$ is the **prior probability** — our initial belief (prediction) about the event A .
- $P(B | A)$ is the **likelihood** — the probability of observing evidence B given that A is true.
- $P(B)$ is the **marginal probability** of observing evidence B under all possible conditions.
- $P(A | B)$ is the **posterior probability** — our updated belief (prediction) about A after observing B .

This formula enables us to revise our estimate of the probability of A whenever new information B becomes available. Depending on the relationship between A and B , the posterior probability may be higher or lower than the prior, reflecting the impact of the new evidence.

2.6.2 The Base Rate Fallacy and Bayes' Theorem

The **base rate fallacy** is a cognitive bias which occurs when people wrongly judge the probability of an event (like having a disease) by focusing on new information (such as a test result) while ignoring how rare or common the event is in the general population.

Imagine you're told that a COVID-19 test is 95-99% accurate. You take the test, and it comes back **positive**. Does that mean you're almost certainly infected with 95% possibility?

Not necessarily. This is where the **base rate fallacy** comes into play—a common error in reasoning where people ignore the underlying probability of an event (known as the *base rate*) and instead focus too heavily on new information, like a test result.

To illustrate this idea, suppose we are testing a population where only 1% of individuals are actually infected with COVID-19. Consider the characteristics of the test:

- **Prevalence (Base Rate):** $P(\text{COVID}) = 0.01$
- **Sensitivity (True Positive Rate):** $P(\text{Positive} | \text{COVID}) = 0.99$
- **Specificity (True Negative Rate):** $P(\text{Negative} | \text{No COVID}) = 0.95$

Now the question is, if a person tests positive, what is the probability that they actually have COVID-19?

To answer this, we apply **Bayes' Theorem**:

$$P(\text{COVID} | \text{Positive}) = \frac{P(\text{Positive} | \text{COVID}) \cdot P(\text{COVID})}{P(\text{Positive})}$$

We compute the total probability of a positive test as:

$$\begin{aligned}
 P(\text{Positive}) &= \underbrace{P(\text{Positive} \mid \text{COVID})}_{0.99} \cdot \underbrace{P(\text{COVID})}_{0.01} + \underbrace{P(\text{Positive} \mid \text{No COVID})}_{1-0.95=0.05} \cdot \underbrace{P(\text{No COVID})}_{0.99} \\
 &= 0.99 \times 0.01 + 0.05 \times 0.99 \\
 &= 0.0099 + 0.0495 = 0.0594
 \end{aligned}$$

Now, substituting into Bayes' formula:

$$P(\text{COVID} \mid \text{Positive}) = \frac{0.99 \times 0.01}{0.0594} \approx \frac{0.0099}{0.0594} \approx 0.167$$

That is, even with a **positive test result**, the probability of actually having COVID is only around **16.7%**. Most people wrongly assume that a positive result from a highly accurate test means they almost certainly have the disease. However, because COVID-19 is relatively rare in the population (a low base rate), even highly accurate tests produce many false positives. This is the essence of the **base rate fallacy**: ignoring the background prevalence (base rate) of the disease and overestimating the significance of the test result.

Example

Imagine testing 10,000 people:

- 1% are infected \Rightarrow 100 people have COVID.
- 99% are not infected \Rightarrow 9,900 people do not have COVID.

Now apply the test:

- Among 100 infected people 99 test positive (99% true positives).
- Among 9,900 uninfected people 9,405 test negative (95% true positives) meaning 495 test positive (5% false positives).

So, total positive tests = 99 (true) + 495 (false) = 594

Probability of actually having COVID given a positive test = $\frac{99}{594} \approx 16.7\%$

—

2.7 Statistical Independence of Events

If A and B are any two events in a sample space S , we say that “ A is independent of B ” if:

$$P(A \mid B) = P(A)$$

This means that knowing whether or not B has occurred “does not change” the probability of A occurring.

From the definition of conditional probability:

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)},$$

if we substitute and rearrange the condition $P(A \mid B) = P(A)$, we get:

$$\frac{P(A \cap B)}{P(B)} = P(A) \quad \Rightarrow \quad P(A \cap B) = P(A) \cdot P(B)$$

Thus, another equivalent definition of independence is:

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

Now, to check whether B is independent of A , we look at:

$$P(B | A) = \frac{P(B \cap A)}{P(A)} = \frac{P(A \cap B)}{P(A)}$$

Since we just showed that $P(A \cap B) = P(A) \cdot P(B)$, substitute:

$$P(B | A) = \frac{P(A) \cdot P(B)}{P(A)} = P(B)$$

Thus, B is also independent of A .

If A is independent of B , then B is also independent of A . Therefore, we say that A and B are **mutually independent**.

Multiplication Rule for Independent Events: Two events A and B are (mutually) independent events if and only if

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

Example: Consider the experiment of rolling two fair six-sided dice, and consider the following events:

- A : The first die shows a 1.

$$A = \{(1, 1), (1, 2), (1, 3), (1, 4), (1, 5), (1, 6)\}$$

- B : The second die shows a 2.

$$B = \{(1, 2), (2, 2), (3, 2), (4, 2), (5, 2), (6, 2)\}$$

The total number of outcomes in the sample space is $6 \times 6 = 36$.

We compute:

$$P(A) = \frac{6}{36} = \frac{1}{6}, \quad P(B) = \frac{6}{36} = \frac{1}{6}$$

$$A \cap B = \{(1, 2)\} \Rightarrow P(A \cap B) = \frac{1}{36}$$

Since

$$P(A \cap B) = P(A) \cdot P(B) = \frac{1}{6} \cdot \frac{1}{6} = \frac{1}{36},$$

we conclude that the events A and B are **independent**.

Theorem: If the events A and B are independent, then the following pairs of events are also independent:

$$(1) A \text{ and } \bar{B}, \quad (2) \bar{A} \text{ and } B, \quad (3) \bar{A} \text{ and } \bar{B}$$

Proof:

1. Since A and B are independent, we have:

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

Then,

$$P(A \cap \bar{B}) = P(A) - P(A \cap B) = P(A) - P(A) \cdot P(B) = P(A)(1 - P(B)) = P(A) \cdot P(\bar{B})$$

So, A and \bar{B} are independent.

2. Similarly,

$$P(\overline{A} \cap B) = P(B) - P(A \cap B) = P(B) - P(A) \cdot P(B) = P(B)(1 - P(A)) = P(\overline{A}) \cdot P(B)$$

So, \overline{A} and B are independent.

3. Finally,

$$\begin{aligned} P(\overline{A} \cap \overline{B}) &= P(\overline{A \cup B}) = 1 - P(A \cup B) = 1 - [P(A) + P(B) - P(A \cap B)] \\ &= 1 - [P(A) + P(B) - P(A) \cdot P(B)] = (1 - P(A))(1 - P(B)) = P(\overline{A}) \cdot P(\overline{B}) \end{aligned}$$

So, \overline{A} and \overline{B} are also independent. ■

Theorem: If A and B are independent events, then:

$$P(A \cup B) = 1 - P(\overline{A}) \cdot P(\overline{B})$$

Proof: Using the formula for the union of two events:

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

Since A and B are independent, $P(A \cap B) = P(A) \cdot P(B)$, so:

$$\begin{aligned} P(A \cup B) &= P(A) + P(B) - P(A) \cdot P(B) \\ &= 1 - (1 - P(A) - P(B) + P(A) \cdot P(B)) \\ &= 1 - (1 - P(A))(1 - P(B)) \\ &= 1 - P(\overline{A})P(\overline{B}) \end{aligned}$$
■

2.7.1 Pairwise vs. Mutual Independence of Multiple Events

Let A , B , and C be three events in a sample space S . The events A , B , and C are said to be **pairwise independent** if:

$$\begin{aligned} P(A \cap B) &= P(A) \cdot P(B), \\ P(A \cap C) &= P(A) \cdot P(C), \\ P(B \cap C) &= P(B) \cdot P(C) \end{aligned}$$

The events A , B , and C are said to be **mutually independent** if:

- They are pairwise independent, and
- The joint probability satisfies:

$$P(A \cap B \cap C) = P(A) \cdot P(B) \cdot P(C)$$

Note: Mutual independence *implies* pairwise independence, but the converse is not necessarily true.

2.7.2 Mutually Exclusive and Independent Events

Let A and B be two events in a sample space.

- If A and B are **mutually exclusive**, then:

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

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

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

If both conditions are true, then:

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

This implies that either $P(A) = 0$, $P(B) = 0$, or both.

Two events A and B cannot be both mutually exclusive and independent unless at least one of them has probability zero.

In terms of a Venn diagram, the independence of events implies that the overlap between sets A and B (i.e., $A \cap B$) should be such that its area (probability) equals the product of the areas (probabilities) of each individual circle.



This is different from mutually exclusive events, where the sets do not overlap at all. Both the condition will be satisfied if the area of at least one of the circle is zero.

2.7.3 Reliability Analysis using Statistical Independence

Reliability analysis is the branch of engineering concerned with estimating the failure rates of systems. If a machine has a reliability of 0.95 over 1 year, it means there is a 95% chance that it will work without failure for the entire year. In many systems, components are arranged either in *series* or in *parallel*, and the system's reliability depends on the configuration.

1. System with Components in Series:

Consider a system with two components, A and B , connected in **series**. In this setup, the system will work only if **both** components function properly.



Suppose the probability that A functions is given by $P(A) = 0.96$, and the probability that B functions is given by $P(B) = 0.92$. Assuming that the components function independently, the probability that the system works is:

$$\begin{aligned} P(\text{System functions}) &= P(A \cap B) \\ &= P(A) \cdot P(B) \\ &= 0.96 \times 0.92 = 0.8832 \end{aligned}$$

Since both components must function, the system's reliability is lower than that of either component individually. The more components in series, the more chances for failure.

2. System with Components in Parallel:

Now consider a system with two components, C and D , connected in **parallel**. In this case, the system will function as long as **at least one** component functions.



Suppose the probabilities that the components function are:

$$P(C) = 0.88, \quad P(D) = 0.85$$

Assuming independence, the probability that the system functions is given by:

$$\begin{aligned} P(\text{System functions}) &= P(C \cup D) \\ &= P(C) + P(D) - P(C \cap D) \\ &= P(C) + P(D) - P(C) \cdot P(D) \\ &= 0.88 + 0.85 - (0.88 \times 0.85) \\ &= 1.73 - 0.748 = 0.982 \end{aligned}$$

In parallel systems, the system is more reliable than the individual components. This configuration adds redundancy, improving fault tolerance.

Chapter 3

Random Variables and Probability Distributions

3.1 What is a Random Variable?

In most cases, we can associate a real number with each elementary event in a sample space. For example, in a coin toss, we may assign the number 1 to the outcome ‘Head’ and 0 to the outcome ‘Tail’. Similarly, when a die is thrown, the outcomes correspond to the numbers $1, 2, 3, \dots, 6$, depending on which face appears on top.

This assignment of numerical values to outcomes allows us to define a function on the sample space. A real-valued function defined on the sample space is called a **random variable** (also referred to as a **stochastic variable**).

Let S be a sample space of a random experiment. A **random variable** is a function

$$X : S \rightarrow \mathbb{R}$$

where each outcome $s \in S$ is mapped to a real number $X(s)$.

Note: A random variable is denoted by an uppercase letter such as X and Y . After experiment is conducted, the measured value of the random variable is denoted by a lowercase letter such as x and y .

3.1.1 Random Variable Types

There are two main types:

1. **Discrete Random Variable:**

A discrete random variable takes on a countable¹ number of distinct values. Let X be the number of heads obtained when a fair coin is tossed three times. The possible values of X are 0, 1, 2, 3. Since these values are countable and finite, X is a discrete random variable.

2. **Continuous Random Variable:**

A continuous random variable takes any value within a certain range of real numbers. The possible values are uncountable and include fractions and decimals. Let Y be the amount of time (in minutes) a customer waits in a queue at a bank. The variable Y can take any real

¹This means the values can be finite (like a die roll) or countably infinite (like the number of coin tosses until the first head).

value within a range, such as $0 \leq Y \leq 30$, including fractions like 3.5 or 12.75. Hence, Y is a continuous random variable.

3.2 Probability Distribution

To each value of a random variable X , there corresponds a definite probability. Let x_1, x_2, \dots, x_n be the possible values of X , and let p_1, p_2, \dots, p_n be the corresponding probabilities such that:

$$P(X = x_i) = p_i \quad \text{for } i = 1, 2, \dots, n$$

A statement of these values along with their associated probabilities defines the probability distribution of the random variable X .

3.2.1 Probability Mass Function (PMF)

The Probability Mass Function (PMF) is used for **discrete random variables**. It gives the probability that a discrete random variable X takes a specific value x_k . The PMF is defined as:

$$p_X(x_k) = P(X = x_k)$$

where x_k is a specific value of the random variable X . The PMF satisfies the following properties:

1. $0 \leq p_X(x_k) \leq 1$ for all k
2. $\sum_k p_X(x_k) = 1$

Example: Consider a discrete random variable X with the following distribution:

$$p_X(x) = P(X = x) = \begin{cases} \frac{1}{4}, & \text{if } x = 1 \\ \frac{1}{2}, & \text{if } x = 2 \\ \frac{1}{4}, & \text{if } x = 3 \\ 0, & \text{otherwise} \end{cases}$$

The PMF is visualized as follows:



3.2.2 Probability Density Function (PDF)

The Probability Density Function (PDF) is used for **continuous random variables**. It is defined as a function f_X such that the probability that a continuous random variable X lies within an interval $[a, b]$ is given by the integral of f_X over that interval:

$$P(a \leq X \leq b) = \int_a^b f_X(x) dx$$

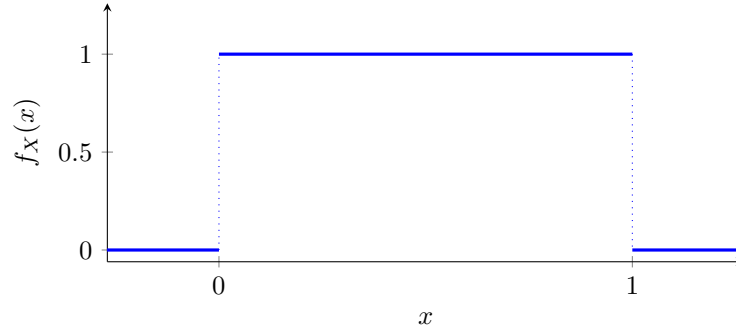
where $f_X(x)$ is the PDF of X . The PDF satisfies the following properties:

1. $f_X(x) \geq 0$ for all x
2. $\int_{-\infty}^{\infty} f_X(x) dx = 1$

Example: Consider a continuous random variable with the following distribution:

$$f_X(x) = \begin{cases} 1, & 0 \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

The PDF can be visualised as follows:



3.2.3 Cumulative Distribution

The cumulative distribution gives the probability that a random variable takes a value less than or equal to a specified value. Mathematically **Cumulative Distribution Function** (CDF) of a random variable X is defined as:

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

This function gives the probability that X takes a value less than or equal to x .

CDF for Discrete Random Variables

For a discrete random variable X , the CDF is given by the sum of the probabilities for all values less than or equal to x . If X takes the values x_1, x_2, \dots, x_n , the CDF is:

$$F_X(x) = \sum_{x_k \leq x} P(X = x_k)$$

This sum includes all the probabilities up to and including x .

Example: Consider the random variable X representing the outcome of a fair six-sided die roll. The CDF for X is:

$$F_X(x) = \begin{cases} 0, & \text{if } x < 1 \\ \frac{1}{6}, & \text{if } 1 \leq x < 2 \\ \frac{2}{6}, & \text{if } 2 \leq x < 3 \\ \frac{3}{6}, & \text{if } 3 \leq x < 4 \\ \frac{4}{6}, & \text{if } 4 \leq x < 5 \\ \frac{5}{6}, & \text{if } 5 \leq x < 6 \\ 1, & \text{if } x \geq 6 \end{cases}$$



CDF for Continuous Random Variables

For a continuous random variable X , the CDF is obtained by integrating the probability density function (PDF) up to x :

$$F_X(x) = \int_{-\infty}^x f_X(t) dt$$

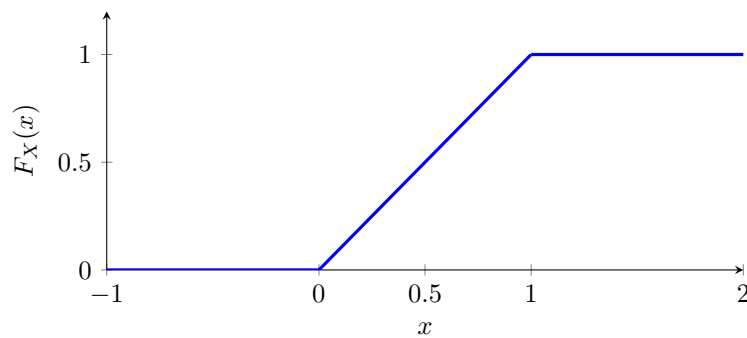
This represents the area under the PDF curve from $-\infty$ to x , giving the cumulative probability up to x .

Example: Consider a continuous random variable X with a probability density function (PDF) $f_X(x) = 1$ for $0 \leq x \leq 1$ (uniform distribution). The CDF is given by:

$$F_X(x) = \begin{cases} 0, & \text{if } x < 0 \\ x, & \text{if } 0 \leq x \leq 1 \\ 1, & \text{if } x > 1 \end{cases}$$

This is obtained by integrating the PDF:

$$F_X(x) = \int_0^x 1 \cdot dt = x \quad \text{for } 0 \leq x \leq 1$$



Properties of the CDF

1. The CDF $F_X(x)$ is a non-decreasing function, meaning:

$$F_X(x_1) \leq F_X(x_2) \quad \text{for } x_1 \leq x_2$$

2. The CDF is bounded between 0 and 1:

$$0 \leq F_X(x) \leq 1 \quad \text{for all } x$$

3. The CDF approaches 1 as $x \rightarrow \infty$ and 0 as $x \rightarrow -\infty$:

$$\lim_{x \rightarrow -\infty} F_X(x) = 0, \quad \lim_{x \rightarrow \infty} F_X(x) = 1$$

4. For a continuous random variable, the CDF is continuous, while for a discrete random variable, the CDF has jumps at the values taken by the random variable.
5. For continuous random variable, the derivative of the CDF gives the PDF (if the derivative exists):

$$f_X(x) = \frac{d}{dx} F_X(x)$$

3.3 Mean and Variance of a Random Variable

The mean and variance are important measures that describe the central tendency and spread of a random variable, respectively.

3.3.1 Mean of a Random Variable

The **mean** or **expected value** of a random variable X , denoted as $\mathbb{E}(X)$, provides the long-run average of the outcomes of the random variable. It is defined as the weighted sum of all possible values of X , weighted by their probabilities.

1. **Discrete Random Variable:** For a discrete random variable X with possible values x_1, x_2, \dots, x_n and corresponding probabilities $p_X(x_1), p_X(x_2), \dots, p_X(x_n)$, the expected value (mean) is denoted by is given by:

$$\mu = \mathbb{E}(X) = \sum_{i=1}^n x_i \cdot p_X(x_i)$$

Example: Suppose X is the outcome of a roll of a fair die. The possible values of X are 1, 2, 3, 4, 5, 6, and each value has a probability of $\frac{1}{6}$. The expected value is:

$$\mathbb{E}(X) = \sum_{i=1}^6 x_i \cdot \frac{1}{6} = \frac{1}{6} \cdot (1 + 2 + 3 + 4 + 5 + 6) = \frac{21}{6} = 3.5$$

2. **Continuous Random Variable:** For a continuous random variable X with probability density function $f_X(x)$, the expected value is given by the integral of x weighted by the probability density function:

$$\mathbb{E}(X) = \int_{-\infty}^{\infty} x \cdot f_X(x) dx$$

Example: Suppose X is a continuous random variable with a uniform distribution between 0 and 1. The probability density function is:

$$f_X(x) = 1 \quad \text{for } 0 \leq x \leq 1$$

The expected value is:

$$\mathbb{E}(X) = \int_0^1 x \cdot 1 dx = \left[\frac{x^2}{2} \right]_0^1 = \frac{1}{2}$$

Theorem If $X = a$, where $a \in \mathbb{R}$ is a constant, then

$$\mathbb{E}(X) = a$$

Proof. Since X is always equal to a ,

- In the discrete case:

$$\mathbb{E}(X) = \sum_x x \cdot P(X = x) = a \cdot P(X = a) = a$$

- In the continuous case:

$$\mathbb{E}(X) = \int_{-\infty}^{\infty} x f_X(x) dx = \int_{-\infty}^{\infty} a \cdot \delta(x - a) dx = a$$

■

Theorem: If $Y = bX$, where $b \in \mathbb{R}$, then

$$\mathbb{E}(Y) = b \cdot \mathbb{E}(X)$$

Proof.

- Discrete case:

$$\mathbb{E}(bX) = \sum_x bx \cdot P(X = x) = b \sum_x x \cdot P(X = x) = b \cdot \mathbb{E}(X)$$

- Continuous case:

$$\mathbb{E}(bX) = \int_{-\infty}^{\infty} bx \cdot f_X(x) dx = b \int_{-\infty}^{\infty} x f_X(x) dx = b \cdot \mathbb{E}(X)$$

■

3.3.2 Variance of a Random Variable

The **variance** of a random variable X , denoted as σ_X^2 or $\text{Var}(X)$, measures the spread or dispersion of the random variable around its mean. It is defined as the expected squared deviation from the mean:

$$\begin{aligned} \text{Var}(X) &= \mathbb{E}[(X - \mathbb{E}(X))^2] \\ &= \mathbb{E}[X^2 - 2X\mathbb{E}(X) + (\mathbb{E}(X))^2] \\ &= \mathbb{E}(X^2) - 2(\mathbb{E}(X))^2 + (\mathbb{E}(X))^2 \\ &= \mathbb{E}(X^2) - (\mathbb{E}(X))^2 \end{aligned}$$

1. **Discrete Random Variable:** For a discrete random variable X with possible values x_1, x_2, \dots, x_n and corresponding probabilities $p_X(x_1), p_X(x_2), \dots, p_X(x_n)$, the variance is given by:

$$\text{Var}(X) = \sum_{i=1}^n (x_i - \mathbb{E}(X))^2 \cdot p_X(x_i)$$

Alternatively, it can be computed as:

$$\text{Var}(X) = \sum_{i=1}^n x_i^2 \cdot p_X(x_i) - (\mathbb{E}(X))^2$$

Example: For the fair die roll example where $\mathbb{E}(X) = 3.5$, the variance is:

$$\begin{aligned} \text{Var}(X) &= \sum_{i=1}^6 (x_i - 3.5)^2 \cdot \frac{1}{6} = \frac{1}{6} ((1 - 3.5)^2 + (2 - 3.5)^2 + \dots + (6 - 3.5)^2) \\ &= \frac{1}{6} \cdot (6.25 + 2.25 + 0.25 + 0.25 + 2.25 + 6.25) = \frac{17.5}{6} \approx 2.92 \end{aligned}$$

2. **Continuous Random Variable:** For a continuous random variable X with probability density function $f_X(x)$, the variance is given by:

$$\text{Var}(X) = \int_{-\infty}^{\infty} (x - \mathbb{E}(X))^2 \cdot f_X(x) dx$$

Alternatively, it can be computed as:

$$\text{Var}(X) = \int_{-\infty}^{\infty} x^2 \cdot f_X(x) dx - (\mathbb{E}(X))^2$$

Example: For the continuous uniform random variable X between 0 and 1, $\mathbb{E}(X) = \frac{1}{2}$, the variance is:

$$\text{Var}(X) = \int_0^1 x^2 \cdot 1 dx - \left(\frac{1}{2}\right)^2 = \left[\frac{x^3}{3}\right]_0^1 - \frac{1}{4} = \frac{1}{3} - \frac{1}{4} = \frac{1}{12}$$

Theorem: If $X = a$, where $a \in \mathbb{R}$ is constant, then

$$\text{Var}(X) = 0$$

Proof.

$$\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}(X))^2] = \mathbb{E}[(a - a)^2] = \mathbb{E}[0] = 0$$

■

Theorem: If $Y = bX$, where $b \in \mathbb{R}$, then

$$\text{Var}(Y) = b^2 \cdot \text{Var}(X)$$

Proof.

$$\begin{aligned} \text{Var}(Y) &= \text{Var}(bX) = \mathbb{E}[(bX - \mathbb{E}(bX))^2] \\ &= \mathbb{E}[(bX - b\mathbb{E}(X))^2] = \mathbb{E}[b^2(X - \mathbb{E}(X))^2] \\ &= b^2 \cdot \text{Var}(X) \end{aligned}$$

■

3.4 Joint Distribution of Two random Variables

Let X and Y be two random variables defined on the same probability space. The joint distribution of X and Y describes the probability behavior of the pair (X, Y) . It can be either discrete or continuous depending on the nature of X and Y .

- **Discrete Case:**

If X and Y are discrete random variables, then their joint distribution is defined by the **joint probability mass function (PMF)**:

$$p_{X,Y}(x, y) = P(X = x, Y = y)$$

which gives the probability that $X = x$ and $Y = y$ simultaneously.

The PMF must satisfy:

- $p_{X,Y}(x, y) \geq 0$ for all x, y
- $\sum_x \sum_y p_{X,Y}(x, y) = 1$

The **marginal distributions** can be obtained by summing out the other variable:

$$p_X(x) = \sum_y p_{X,Y}(x, y), \quad p_Y(y) = \sum_x p_{X,Y}(x, y)$$

The value of $p_X(x)$ gives the probability $P(X = x)$ irrespective of Y . Similarly, $p_Y(y)$ gives the probability $P(Y = y)$ irrespective of X .

Suppose the possible values of X are x_1, x_2, \dots, x_m and possible values of Y are y_1, y_2, \dots, y_n . The joint distribution can be depicted in a table format as shown Table 3.1.

$\mathbf{Y} \setminus \mathbf{X}$	x_1	x_2	\cdots	x_m	$p_Y(y_j)$
y_1	$p_{X,Y}(x_1, y_1)$	$p_{X,Y}(x_2, y_1)$	\cdots	$p_{X,Y}(x_m, y_1)$	$p_Y(y_1)$
y_2	$p_{X,Y}(x_1, y_2)$	$p_{X,Y}(x_2, y_2)$	\cdots	$p_{X,Y}(x_m, y_2)$	$p_Y(y_2)$
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
y_n	$p_{X,Y}(x_1, y_n)$	$p_{X,Y}(x_2, y_n)$	\cdots	$p_{X,Y}(x_m, y_n)$	$p_Y(y_n)$
$p_X(x_i)$	$p_X(x_1)$	$p_X(x_2)$	\cdots	$p_X(x_m)$	1

Table 3.1: Joint probability distribution of discrete random variables X and Y .

The **joint cumulative probability distribution function (CDF)** can be obtained from joint probability mass function:

$$F_{X,Y}(x, y) = P(X \leq x, Y \leq y) = \sum_{x_i \leq x} \sum_{y_j \leq y} p_{X,Y}(x_i, y_j)$$

It sums over all values (x_i, y_j) such that $x_i \leq x$ and $y_j \leq y$. This gives the total probability that the random pair (X, Y) falls within the region $X \leq x, Y \leq y$.

- **Continuous Case:**

If X and Y are continuous random variables, then their joint distribution is described by the **joint probability density function**:

$$f_{X,Y}(x, y)$$

which satisfies:

$$P((X, Y) \in A) = \iint_A f_{X,Y}(x, y) dx dy,$$

for any region $A \subset \mathbb{R}^2$.

The joint PDF must satisfy:

- $f_{X,Y}(x, y) \geq 0$ for all x, y
- $\iint_{\mathbb{R}^2} f_{X,Y}(x, y) dx dy = 1$

The **marginal PDFs** are obtained by integrating out the other variable:

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy, \quad f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx$$

The term $f_X(x)$ gives the PDF of X irrespective of Y . Similarly, $f_Y(y)$ gives the PDF of Y irrespective of X .

The **joint cumulative distribution function (CDF)** is defined as:

$$F_{X,Y}(x, y) = P(X \leq x, Y \leq y) = \int_{-\infty}^x \int_{-\infty}^y f_{X,Y}(u, v) dv du$$

This function gives the probability that the random vector (X, Y) falls within the region defined by $X \leq x, Y \leq y$.

Theorem: If $Z = X + Y$, then

$$\mathbb{E}(Z) = \mathbb{E}(X) + \mathbb{E}(Y)$$

Proof.

- Discrete case:

$$\begin{aligned}\mathbb{E}(X + Y) &= \sum_{x,y} (x + y) \cdot p_{X,Y}(x, y) \\ &= \sum_{x,y} x \cdot p_{X,Y}(x, y) + \sum_{x,y} y \cdot p_{X,Y}(x, y) \\ &= \sum_x x \cdot p_X(x) + \sum_y y \cdot p_Y(y) \\ &= \mathbb{E}(X) + \mathbb{E}(Y)\end{aligned}$$

- Continuous case:

$$\begin{aligned}\mathbb{E}(X + Y) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x + y) f_{X,Y}(x, y) dx dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x f_{X,Y}(x, y) dx dy + \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f_{X,Y}(x, y) dx dy \\ &= \int_{-\infty}^{\infty} x f_X(x) dx + \int_{-\infty}^{\infty} y f_Y(y) dy = \mathbb{E}(X) + \mathbb{E}(Y)\end{aligned}$$

■

3.4.1 Independence of Two Random Variables

Two random variables X and Y are **independent** if and only if:

- In the discrete case:

$$p_{X,Y}(x, y) = p_X(x) \cdot p_Y(y) \text{ for all } (x, y)$$

- In the continuous case:

$$f_{X,Y}(x, y) = f_X(x) \cdot f_Y(y) \text{ for all } (x, y)$$

Theorem: If X and Y are independent random variables, then

$$\mathbb{E}(XY) = \mathbb{E}(X) \cdot \mathbb{E}(Y)$$

Proof:

- Discrete Case:

$$\begin{aligned}\mathbb{E}(XY) &= \sum_x \sum_y xy \cdot p_{X,Y}(x, y) \\ &= \sum_x \sum_y xy \cdot p_X(x) \cdot p_Y(y) \\ &= \sum_x x p_X(x) \cdot \sum_y y p_Y(y) \\ &= \mathbb{E}(X) \cdot \mathbb{E}(Y)\end{aligned}$$

- Continuous Case:

$$\begin{aligned}
\mathbb{E}(XY) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy \cdot f_{X,Y}(x,y) dy dx \\
&= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy \cdot f_X(x) \cdot f_Y(y) dy dx \\
&= \int_{-\infty}^{\infty} x f_X(x) dx \cdot \int_{-\infty}^{\infty} y f_Y(y) dy \\
&= \mathbb{E}(X) \cdot \mathbb{E}(Y)
\end{aligned}$$

3.4.2 Covariance of Two Random Variables

One important feature of the joint distribution of X and Y is their covariance, which is used to measure the degree of association between X and Y . The **covariance** of two random variables X and Y , denoted by $\text{Cov}(X, Y)$ is defined as:

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}(X)) \cdot (Y - \mathbb{E}(Y))]$$

Now,

$$\begin{aligned}
\text{Cov}(X, Y) &= \mathbb{E}[(X - \mathbb{E}(X)) \cdot (Y - \mathbb{E}(Y))] \\
&= \mathbb{E}[XY - X \cdot \mathbb{E}(Y) - \mathbb{E}(X) \cdot Y + \mathbb{E}(X) \cdot \mathbb{E}(Y)] \\
&= \mathbb{E}(XY) - \mathbb{E}(X) \cdot \mathbb{E}(Y) - \mathbb{E}(X) \cdot \mathbb{E}(Y) + \mathbb{E}(X) \cdot \mathbb{E}(Y) \\
&= \mathbb{E}(XY) - \mathbb{E}(X) \cdot \mathbb{E}(Y)
\end{aligned}$$

Thus we get an alternative shortcut formula for covariance:

$$\text{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}(X) \cdot \mathbb{E}(Y)$$

Interpretation:

- $\text{Cov}(X, Y) > 0$: X and Y tend to increase (or decrease) together.
- $\text{Cov}(X, Y) < 0$: X increases as Y decreases (or vice versa).
- $\text{Cov}(X, Y) = 0$: No linear relationship between X and Y .

Theorem: If X and Y are random variables, then

$$\text{Cov}(X, Y) = \text{Cov}(Y, X)$$

Proof: By definition of covariance,

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

Using the commutative property of multiplication,

$$(X - \mathbb{E}[X])(Y - \mathbb{E}[Y]) = (Y - \mathbb{E}[Y])(X - \mathbb{E}[X])$$

Therefore,

$$\text{Cov}(X, Y) = \mathbb{E}[(Y - \mathbb{E}[Y])(X - \mathbb{E}[X])] = \text{Cov}(Y, X)$$

■

Theorem: If X and Y are independent random variables, then

$$\text{Cov}(X, Y) = 0$$

Proof: If X and Y are independent, then

$$\mathbb{E}(XY) = \mathbb{E}(X) \cdot \mathbb{E}(Y)$$

Thus,

$$\begin{aligned} \text{Cov}(X, Y) &= \mathbb{E}(XY) - \mathbb{E}(X) \cdot \mathbb{E}(Y) \\ &= \mathbb{E}(X) \cdot \mathbb{E}(Y) - \mathbb{E}(X) \cdot \mathbb{E}(Y) \\ &= 0 \end{aligned}$$

■

Note: The converse is not necessarily true. Two random variables may have $\text{Cov}(X, Y) = 0$ and yet not be independent.

Theorem: For any random variable X ,

$$\text{Cov}(X, X) = \text{Var}(X)$$

Proof:

$$\begin{aligned} \text{Cov}(X, X) &= \mathbb{E}(X \cdot X) - \mathbb{E}(X) \cdot \mathbb{E}(X) \\ &= \mathbb{E}(X^2) - (\mathbb{E}(X))^2 \\ &= \text{Var}(X) \end{aligned}$$

■

Theorem: Let $a, b \in \mathbb{R}$ be constants. Then for any random variables X and Y ,

$$\text{Cov}(aX, bY) = ab \cdot \text{Cov}(X, Y)$$

Proof.

$$\begin{aligned} \text{Cov}(aX, bY) &= \mathbb{E}(aX \cdot bY) - \mathbb{E}(aX) \cdot \mathbb{E}(bY) \\ &= ab \cdot \mathbb{E}(XY) - ab \cdot \mathbb{E}(X) \cdot \mathbb{E}(Y) \\ &= ab \cdot (\mathbb{E}(XY) - \mathbb{E}(X) \cdot \mathbb{E}(Y)) \\ &= ab \cdot \text{Cov}(X, Y) \end{aligned}$$

■

Theorem: If X, Y are two random variables, then

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2 \cdot \text{Cov}(X, Y)$$

Proof.

$$\text{Var}(X + Y) = \mathbb{E}[(X + Y - \mathbb{E}(X + Y))^2]$$

By linearity of expectation, $\mathbb{E}(X + Y) = \mathbb{E}(X) + \mathbb{E}(Y)$, so:

$$= \mathbb{E}[(X - \mathbb{E}(X) + Y - \mathbb{E}(Y))^2]$$

Expanding the square:

$$= \mathbb{E}[(X - \mathbb{E}(X))^2 + (Y - \mathbb{E}(Y))^2 + 2(X - \mathbb{E}(X))(Y - \mathbb{E}(Y))]$$

Using linearity of expectation:

$$\begin{aligned} &= \mathbb{E}[(X - \mathbb{E}(X))^2] + \mathbb{E}[(Y - \mathbb{E}(Y))^2] + 2 \cdot \mathbb{E}[(X - \mathbb{E}(X))(Y - \mathbb{E}(Y))] \\ &= \text{Var}(X) + \text{Var}(Y) + 2 \cdot \text{Cov}(X, Y) \end{aligned}$$

■

Theorem: For random variables X_1, X_2, \dots, X_n ,

$$\text{Var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \text{Var}(X_i) + 2 \sum_{1 \leq i < j \leq n} \text{Cov}(X_i, X_j)$$

Proof: Start with the definition of variance:

$$\text{Var}(Y) = \mathbb{E}[(Y - \mathbb{E}[Y])^2]$$

Let

$$Y = \sum_{i=1}^n X_i$$

Then

$$\text{Var}\left(\sum_{i=1}^n X_i\right) = \mathbb{E}\left[\left(\sum_{i=1}^n X_i - \mathbb{E}\left[\sum_{i=1}^n X_i\right]\right)^2\right]$$

Since expectation is linear,

$$\mathbb{E}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \mathbb{E}[X_i]$$

So,

$$\text{Var}\left(\sum_{i=1}^n X_i\right) = \mathbb{E}\left[\left(\sum_{i=1}^n (X_i - \mathbb{E}[X_i])\right)^2\right]$$

Expanding the square,

$$= \mathbb{E}\left[\sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \sum_{j=1}^n (X_j - \mathbb{E}[X_j])\right] = \mathbb{E}\left[\sum_{i=1}^n \sum_{j=1}^n (X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])\right]$$

By linearity of expectation,

$$= \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])]$$

Recall the definition of covariance:

$$\text{Cov}(X_i, X_j) = \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])]$$

Thus,

$$\text{Var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \sum_{j=1}^n \text{Cov}(X_i, X_j)$$

Split the double sum into terms where $i = j$ and $i \neq j$:

$$= \sum_{i=1}^n \text{Var}(X_i) + \sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n \text{Cov}(X_i, X_j).$$

Since the covariance terms are symmetric, the sum over $i \neq j$ counts each pair twice, so

$$\sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n \text{Cov}(X_i, X_j) = 2 \sum_{1 \leq i < j \leq n} \text{Cov}(X_i, X_j)$$

Hence,

$$\text{Var} \left(\sum_{i=1}^n X_i \right) = \sum_{i=1}^n \text{Var}(X_i) + 2 \sum_{1 \leq i < j \leq n} \text{Cov}(X_i, X_j).$$

■

3.5 Conditional Probability Distribution

When working with two random variables X and Y , it is often important to understand the distribution and expected value of one variable given the other. This is captured by **conditional probability distributions**.

3.5.1 Conditional Probability Mass Function

If X and Y are discrete random variables with joint PMF $p_{X,Y}(x, y)$, the conditional PMF of Y given $X = x$ is:

$$p_{Y|X}(y | x) = P(Y = y | X = x) = \frac{p_{X,Y}(x, y)}{p_X(x)}, \quad \text{for } p_X(x) > 0$$

If X and Y are **independent**, then the joint PMF factorizes as:

$$p_{X,Y}(x, y) = p_X(x) \cdot p_Y(y)$$

and the conditional PMF reduces to the marginal PMF of Y :

$$p_{Y|X}(y | x) = p_Y(y)$$

In other words, knowing $X = x$ does not change the probability distribution of Y .

3.5.2 Conditional Probability Density Function

If X and Y are continuous random variables with joint PDF $f_{X,Y}(x, y)$, the conditional PDF of Y given $X = x$ is:

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x, y)}{f_X(x)}, \quad \text{for } f_X(x) > 0$$

If X and Y are **independent**, then the joint PDF factorizes as:

$$f_{X,Y}(x, y) = f_X(x) \cdot f_Y(y)$$

and the conditional PDF reduces to the marginal PDF of Y :

$$f_{Y|X}(y|x) = f_Y(y)$$

Thus, knowing $X = x$ does not affect the distribution of Y .

3.5.3 Conditional Expectation

The **conditional expectation** of Y given $X = x$, denoted $\mathbb{E}[Y | X = x]$, is the expected value of Y calculated using the conditional distribution of Y given $X = x$. It provides the average or mean value of Y when X is known.

- **Discrete case:**

$$\mathbb{E}[Y | X = x] = \sum_y y \cdot P(Y = y | X = x)$$

- **Continuous case:**

$$\mathbb{E}[Y | X = x] = \int_{-\infty}^{\infty} y \cdot f_{Y|X}(y|x) dy$$

The conditional expectation is a function of x and can be viewed as the “best guess” of Y given the value of X .

Law of Total Expectation: Let X and Y be random variables (discrete or continuous). Then

$$\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X | Y]]$$

This means that the overall expectation of X can be obtained by first computing the conditional expectation of X given Y , and then taking the expectation of that conditional expectation over the distribution of Y .

Proof:

- **Discrete case:** Assume X and Y are discrete random variables with joint PMF $p_{X,Y}(x, y)$:

$$\begin{aligned} \mathbb{E}[X] &= \sum_x x p_X(x) = \sum_x x \sum_y p_{X,Y}(x, y) \\ &= \sum_y \sum_x x p_{X,Y}(x, y) = \sum_y \left(\sum_x x p_{X|Y}(x|y) \right) p_Y(y) \\ &= \sum_y \mathbb{E}[X | Y = y] p_Y(y) = \mathbb{E}[\mathbb{E}[X | Y]] \end{aligned}$$

- **Continuous case:** If X, Y are continuous with joint PDF $f_{X,Y}(x, y)$:

$$\begin{aligned} \mathbb{E}[X] &= \int_{-\infty}^{\infty} x f_X(x) dx = \int_{-\infty}^{\infty} x \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy dx \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x f_{X|Y}(x|y) f_Y(y) dx dy \\ &= \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} x f_{X|Y}(x|y) dx \right) f_Y(y) dy \\ &= \int_{-\infty}^{\infty} \mathbb{E}[X | Y = y] f_Y(y) dy = \mathbb{E}[\mathbb{E}[X | Y]] \end{aligned}$$

■

Example: Suppose two factories supply light bulbs to the market.

- Factory X produces bulbs that last on average 5500 hours.
- Factory Y produces bulbs that last on average 4200 hours.

- Factory X supplies 70% of all bulbs, and factory Y supplies 30%.

What is the expected lifetime of a randomly chosen bulb?

Let the random variable L denote the lifetime of a randomly chosen bulb. Let $F \in \{X, Y\}$ be the factory that produced the bulb.

We are asked to find the expected lifetime $\mathbb{E}[L]$.

Given,

$$\begin{aligned}\mathbb{E}[L \mid F = X] &= 2000 \\ \mathbb{E}[L \mid F = Y] &= 3000 \\ P(F = X) &= \frac{70}{100} = 0.7 \\ P(F = Y) &= \frac{30}{100} = 0.3\end{aligned}$$

Using the **Law of Total Expectation**:

$$\mathbb{E}[L] = \mathbb{E}[\mathbb{E}[L \mid F]]$$

We compute:

$$\begin{aligned}\mathbb{E}[L] &= P(F = X) \cdot \mathbb{E}[L \mid F = X] + P(F = Y) \cdot \mathbb{E}[L \mid F = Y] \\ &= 0.7 \times 2000 + 0.3 \times 3000 \\ &= 1400 + 900 = 2300\end{aligned}$$

So, the expected lifetime of a randomly chosen bulb is 2300 hours.

3.6 Functions of a Random Variable

In many practical scenarios, we are interested not just in a random variable X , but in some transformation or function of X , such as $Y = g(X)$. Y is also a random variable. Understanding how the distribution of X affects the distribution of Y is a key part of probability theory.

3.6.1 Discrete Case

If X is a discrete random variable with known probability mass function (PMF) $p_X(x) = P(X = x)$, and $Y = g(X)$, then the PMF of Y is computed as:

$$p_Y(y) = P(Y = y) = \sum_{\{x \mid g(x)=y\}} p_X(x)$$

That is, for each possible value y of Y , sum the probabilities of all values x of X that are mapped to y by the function g . This summation holds for any function g : whether it is one-to-one or many-to-one.

Example: Let X be the outcome of a fair six-sided die, so $X \in \{1, 2, 3, 4, 5, 6\}$ with $P(X = x) = \frac{1}{6}$. Let $Y = X \bmod 2$ (i.e., Y is the parity of X).

Then Y takes values in $\{0, 1\}$, where:

$$\begin{aligned}f_Y(0) &= P(Y = 0) = P(X \in \{2, 4, 6\}) = \frac{3}{6} = 0.5 \\ f_Y(1) &= P(Y = 1) = P(X \in \{1, 3, 5\}) = \frac{3}{6} = 0.5\end{aligned}$$

3.6.2 Continuous Case

If X is a continuous random variable with probability density function (PDF) $f_X(x)$, and $Y = g(X)$ is a function of X , then the PDF of Y depends on whether g is monotonic² or not.

Monotonic Transformation

Theorem: If g is a strictly monotonic and differentiable function, and $Y = g(X)$, then the PDF of Y is:

$$f_Y(y) = f_X(g^{-1}(y)) \cdot \left| \frac{d}{dy} g^{-1}(y) \right|$$

Or equivalently:

Theorem. If g is a strictly monotonic and differentiable function, and $Y = g(X)$, then the CDF of Y given by:

- If g is strictly increasing, then

$$F_Y(y) = F_X(g^{-1}(y))$$

- If g is strictly decreasing, then

$$F_Y(y) = 1 - F_X(g^{-1}(y))$$

This formula ensures that the total probability is preserved under the transformation.

Proof: We treat separately the cases when g is strictly increasing and strictly decreasing.

- **Case 1: g is strictly increasing.**

Since g is strictly increasing, it has a well-defined inverse

$$x = g^{-1}(y)$$

Let $F_X(x)$ and $F_Y(y)$ be the CDFs of X and Y . Then

$$F_Y(y) = P(Y \leq y) = P(g(X) \leq y)$$

Because g is increasing,

$$g(X) \leq y \iff X \leq g^{-1}(y)$$

Hence

$$F_Y(y) = P(X \leq g^{-1}(y)) = F_X(g^{-1}(y))$$

Now differentiate both sides using the chain rule;

$$f_Y(y) = \frac{d}{dy} F_Y(y) = \frac{d}{dy} F_X(g^{-1}(y)) = f_X(g^{-1}(y)) \frac{d}{dy} g^{-1}(y)$$

Since g^{-1} is increasing, its derivative is positive i.e. $\frac{d}{dy} g^{-1}(y) > 0$, and so

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|$$

²A function $g(x)$ is called **monotonic** if it never “switches direction” as x moves along its domain. In plain english it is either strictly increasing or strictly decreasing function. If g is strictly increasing or decreasing, then it’s **one-to-one and onto**, so there is a well-defined inverse function g^{-1} that ‘undoes’ g i.e.

$$g^{-1}(g(x)) = x$$

- **Case 2: g is strictly decreasing.**

Again g^{-1} exists, but now g^{-1} is decreasing. For a decreasing g ,

$$g(X) \leq y \iff X \geq g^{-1}(y)$$

Thus

$$F_Y(y) = P(Y \leq y) = P(X \geq g^{-1}(y)) = 1 - P(X < g^{-1}(y)) = 1 - F_X(g^{-1}(y))$$

Differentiate to get the PDF:

$$f_Y(y) = \frac{d}{dy} F_Y(y) = \frac{d}{dy} [1 - F_X(g^{-1}(y))] = -f_X(g^{-1}(y)) \frac{d}{dy} g^{-1}(y)$$

Since g^{-1} is decreasing, $\frac{d}{dy} g^{-1}(y) < 0$, so the two negatives cancel:

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|$$

In both cases we arrive at the same formula:

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|$$

■

Example: Let X has the following PDF:

$$f_X(x) = \begin{cases} 1, & 0 \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

and define $Y = -\log(X)$.

To calculate the PDF of Y , note that $g(x) = -\log(x)$ is strictly decreasing on $(0, 1)$. The inverse function is

$$g^{-1}(y) = e^{-y}$$

The absolute value of the derivative is:

$$\left| \frac{d}{dy} e^{-y} \right| = e^{-y}$$

The PDF of X is $f_X(x) = 1$ for $x \in (0, 1)$, so:

$$f_Y(y) = f_X(e^{-y}) \cdot e^{-y} = 1 \cdot e^{-y} = e^{-y}, \quad y > 0$$

Non-Monotonic Transformation

If g is not monotonic, the formula for $f_Y(y)$ generalizes to:

$$f_Y(y) = \sum_{x \in g^{-1}(y)} \frac{f_X(x)}{|g'(x)|}$$

Here, the sum runs over all x values such that $g(x) = y$ (as g can be a many-to-one function such that multiple values of x can be mapped to same y). The proof is omitted as it is beyond the scope of this text.

Example: Let the random variable X follows a standard normal distribution³ i.e. $X \sim \mathcal{N}(0, 1)$ and define $Y = X^2$. The function $g(x) = x^2$ is not one-to-one, but has two inverse branches: $x = \sqrt{y}$ and $x = -\sqrt{y}$.

³The discussion on standard normal distribution will be coming in a later section.

Then:

$$f_Y(y) = \frac{1}{\sqrt{2\pi}} \left(\frac{1}{\sqrt{y}} \exp\left(-\frac{y}{2}\right) + \frac{1}{\sqrt{y}} \exp\left(-\frac{y}{2}\right) \right) = \frac{1}{\sqrt{2\pi y}} \exp\left(-\frac{y}{2}\right), \quad y > 0$$

This is the PDF of a chi-squared distribution with 1 degree of freedom.

3.7 Standardized Random Variable

Let X be a random variable with mean $\mu = \mathbb{E}(X)$ and variance $\sigma^2 = \text{Var}(X)$, where $\sigma = \sqrt{\text{Var}(X)}$ must satisfy $\sigma > 0$. The **standardized random variable** X^* is defined by

$$X^* = \frac{X - \mu}{\sigma}$$

By construction, X^* has

$$\text{Mean} = \mathbb{E}[X^*] = \mathbb{E}\left[\frac{X - \mu}{\sigma}\right] = \frac{\mathbb{E}(X) - \mu}{\sigma} = 0,$$

and

$$\text{Variance} = \text{Var}(X^*) = \text{Var}\left(\frac{X - \mu}{\sigma}\right) = \frac{1}{\sigma^2} \text{Var}(X) = 1$$

3.8 Chebyshev's Inequality

Chebyshev's Inequality is a fundamental result in probability theory that provides an upper bound on the probability that the value of a random variable deviates from its mean by more than a certain number of standard deviations. It is particularly useful when the distribution of the random variable is unknown.

Chebyshev's Inequality: Let X be a random variable with finite expected value $\mu = \mathbb{E}(X)$ and finite variance $\sigma^2 = \text{Var}(X)$. Then for any $k > 0$,

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$$

Proof⁴:

- **Discrete Random Variables**

Assume X takes values in a countable set (sample space) $S \subset \mathbb{R}$ with probability mass function $p_X(x) = P(X = x)$. The variance of X is given by:

$$\sigma^2 = \sum_{x \in S} (x - \mu)^2 p_X(x)$$

Let $A = \{x \in S : |x - \mu| \geq k\sigma\}$ and $\bar{A} = S \setminus A = \{x : |x - \mu| < k\sigma\}$. Then,

$$\sigma^2 = \sum_{x \in A} (x - \mu)^2 p_X(x) + \sum_{x \in \bar{A}} (x - \mu)^2 p_X(x)$$

⁴An equivalent version of Chebyshev's Inequality states that the probability that X lies within $k\sigma$ is bounded by the inequality:

$$P(|X - \mu| < k\sigma) \geq 1 - \frac{1}{k^2}$$

On the set A , we have $(x - \mu)^2 \geq k^2\sigma^2$, hence

$$\sum_{x \in A} (x - \mu)^2 p_X(x) \geq k^2\sigma^2 \sum_{x \in A} p(x) = k^2\sigma^2 P(|X - \mu| \geq k\sigma)$$

Therefore,

$$\sigma^2 \geq k^2\sigma^2 P(|X - \mu| \geq k\sigma)$$

Dividing both sides by $k^2\sigma^2$ gives:

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$$

• Continuous Random Variables

Suppose X is a continuous random variable with probability density function $f_X(x)$. The variance is:

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f_X(x) dx$$

Let $B = \{x \in \mathbb{R} : |x - \mu| \geq k\sigma\}$ and $\bar{B} = \{x : |x - \mu| < k\sigma\}$. Then,

$$\sigma^2 = \int_B (x - \mu)^2 f_X(x) dx + \int_{\bar{B}} (x - \mu)^2 f_X(x) dx$$

On B , we have $(x - \mu)^2 \geq k^2\sigma^2$, so

$$\int_B (x - \mu)^2 f_X(x) dx \geq k^2\sigma^2 \int_B f_X(x) dx = k^2\sigma^2 P(|X - \mu| \geq k\sigma)$$

Thus,

$$\sigma^2 \geq k^2\sigma^2 P(|X - \mu| \geq k\sigma),$$

and dividing both sides by $k^2\sigma^2$ yields:

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$$

■

Properties:

- Chebyshev's Inequality holds for any distribution with finite mean and variance, regardless of its shape.
- The bound provided by the inequality is not always tight; in many cases, the actual probability is much smaller than the upper bound.
- This inequality is particularly useful when the distribution is unknown.

Example: Suppose X is a random variable with mean $\mu = 50$ and standard deviation $\sigma = 5$. Using Chebyshev's Inequality, the probability that X lies outside the interval $[35, 65]$ (which is $\mu \pm 3\sigma$) is bounded by

$$P(|X - 50| \geq 15) \leq \frac{1}{3^2} = \frac{1}{9} \approx 0.1111.$$

Thus, we can say that at least $1 - \frac{1}{9} = \frac{8}{9} \approx 88.89\%$ of the probability mass lies within three standard deviations of the mean.

3.9 Moments and Moment Generating Function

Moments are quantitative measures that capture various aspects of the shape of a probability distribution—its central location, spread, asymmetry, and tail heaviness.

3.9.1 Raw Moments and Central Moments

Moments can be calculated about the origin (raw moments) or about the mean (central moments).

The **k -th raw moment** (also called the moment about the origin) of a random variable X is defined as:

$$\mu'_k = \mathbb{E}[X^k]$$

The **k -th central moment** of a random variable X is defined as:

$$\mu_k = \mathbb{E}[(X - \mu)^k]$$

where $\mu = \mathbb{E}(X)$ is the mean of the distribution.

Moments provide insight into the shape of a distribution. In particular:

- **Mean:** First raw moment:

$$\mu'_1 = \mathbb{E}(X) = \mu$$

- **Variance:** The second central moment:

$$\mu_2 = \mathbb{E}[(X - \mu)^2] = \sigma^2$$

- **Skewness:** Measures the asymmetry of the distribution. The coefficient of skewness as the third central moment of the standardized random variable $X^* = \frac{X - \mu}{\sigma}$:

$$\gamma_1 = \mathbb{E}[(X^*)^3] = \frac{\mathbb{E}[(X - \mu)^3]}{\sigma^3} = \frac{\mu_3}{\sigma^3}$$

- **Kurtosis:** Measures the ‘tailedness’ or ‘peakedness’ of the distribution. The coefficient of kurtosis is defined as the fourth central moment of the standardized random variable $X^* = \frac{X - \mu}{\sigma}$ minus 3:

$$\gamma_2 = \mathbb{E}[(X^*)^4] = \frac{\mathbb{E}[(X - \mu)^4]}{\sigma^4} - 3 = \frac{\mu_4}{\sigma^4} - 3$$

3.9.2 Relation Between Raw and Central Moments

Let X be a random variable with mean $\mu = \mathbb{E}(X)$. The k -th central moment of X is defined as:

$$\mu_k = \mathbb{E}[(X - \mu)^k]$$

To relate this to the raw moments $\mu'_r = \mathbb{E}[X^r]$, we expand $(X - \mu)^k$ using the binomial theorem:

$$(X - \mu)^k = \sum_{r=0}^k \binom{k}{r} (-\mu)^{k-r} X^r$$

Taking expectations on both sides:

$$\mu_k = \mathbb{E}[(X - \mu)^k] = \mathbb{E} \left[\sum_{r=0}^k \binom{k}{r} (-\mu)^{k-r} X^r \right] = \sum_{r=0}^k (-1)^{k-r} \binom{k}{r} \mu^{k-r} \mathbb{E}[X^r]$$

Hence, the central moment is:

$$\mu_k = \sum_{r=0}^k (-1)^{k-r} \binom{k}{r} \mu^{k-r} \mu'_r$$

This formula expresses the k -th central moment μ_k as a linear combination of raw moments $\mu'_r = \mathbb{E}[X^r]$ for $r = 0, 1, \dots, k$.

- **First central moment:**

$$\mu_1 = \mu'_1$$

- **Second central moment⁵:**

$$\mu_2 = \mu'_2 - \mu^2$$

- **Third central moment:**

$$\mu_3 = \mu'_3 - 3\mu\mu'_2 + 3\mu^2\mu'_1 - \mu^3$$

- **Fourth central moment:**

$$\mu_4 = \mu'_4 - 4\mu\mu'_3 + 6\mu^2\mu'_2 - 4\mu^3\mu'_1 + \mu^4$$

Now to get the expression of the raw moment μ'_k in terms of central moments μ_r , we expand $X^k = (\mu + (X - \mu))^k$ as:

$$X^k = \sum_{r=0}^k \binom{k}{r} \mu^{k-r} (X - \mu)^r$$

Taking expectation on both sides:

$$\mu'_k = \mathbb{E}[X^k] = \sum_{r=0}^k \binom{k}{r} \mu^{k-r} \mathbb{E}[(X - \mu)^r] = \sum_{r=0}^k \binom{k}{r} \mu^{k-r} \mu_r$$

$$\mu'_k = \sum_{r=0}^k \binom{k}{r} \mu^{k-r} \mu_r$$

This gives the expression of the raw moment μ'_k in terms of central moments μ_r for $r = 0, 1, \dots, k$, where:

$$\mu_r = \mathbb{E}[(X - \mu)^r], \quad \mu_0 = 1$$

- **First raw moment:**

$$\mu'_1 = \mu$$

- **Second raw moment:**

$$\mu'_2 = \mu^2 + \mu_2$$

- **Third raw moment:**

$$\mu'_3 = \mu^3 + 3\mu\mu_2 + \mu_3$$

- **Fourth raw moment:**

$$\mu'_4 = \mu^4 + 6\mu^2\mu_2 + 4\mu\mu_3 + \mu_4$$

⁵An easier way to calculate the second central moment:

$$\mu_2 = \mathbb{E}[(X - \mu)^2] = \mathbb{E}[X^2 - 2\mu X + \mu^2] = \mathbb{E}[X^2] - 2\mu\mathbb{E}(X) + \mu^2 = \mu'_2 - \mu^2$$

3.9.3 Moment Generating Function

There is a clever way of organizing all the moments into one mathematical object, and that object is called the moment generating function.

The **moment generating function** (MGF) of a random variable X is a function $M_X : \mathbb{R} \rightarrow [0, \infty)$ given by

$$M_X(t) = \mathbb{E}(e^{tX})$$

provided the expectation exists in an open neighborhood^a of $t = 0$.

^aAn **open neighborhood** of $t = 0$ is an open interval around 0, say $(-\varepsilon, \varepsilon)$ for some $\varepsilon > 0$. The condition says that the expectation $\mathbb{E}[e^{tX}]$ must *converge* (i.e., be finite) for all values of t in some interval around 0.

More explicitly, the moment generating function (MGF) of a random variable X can be written as:

- If X is a **discrete random variable** with probability mass function $p_X(x_i) = P(X = x_i)$, then

$$M_X(t) = \sum_{x_i} e^{tx_i} p_X(x_i)$$

- If X is a **continuous random variable** with probability density function $f_X(x)$, then

$$M_X(t) = \int_{-\infty}^{\infty} e^{tx} f_X(x) dx$$

The method to generate moments is given in the following theorem.

Theorem: Let $M_X(t)$ be the moment generating function (MGF) of a random variable X . Then the k th raw moment μ'_k of X is given by the k th derivative of $M_X(t)$ evaluated at $t = 0$, i.e.,

$$\mu'_k = M_X^{(k)}(0) = \left. \frac{d^k}{dt^k} M_X(t) \right|_{t=0}$$

Proof: Let us start by expanding e^{tX} using its Taylor series about $t = 0$:

$$e^{tX} = \sum_{n=0}^{\infty} \frac{(tX)^n}{n!} = \sum_{n=0}^{\infty} \frac{t^n X^n}{n!}$$

Taking expectation on both sides:

$$M_X(t) = \mathbb{E}(e^{tX}) = \mathbb{E}\left(\sum_{n=0}^{\infty} \frac{t^n X^n}{n!}\right)$$

We can interchange summation and expectation⁶:

$$\begin{aligned} M_X(t) &= \sum_{n=0}^{\infty} \frac{t^n}{n!} \mathbb{E}(X^n) = \sum_{n=0}^{\infty} \frac{t^n}{n!} \mu'_n \\ &= 1 + \frac{t}{1} \mu'_1 + \frac{t^2}{2!} \mu'_2 + \frac{t^3}{3!} \mu'_3 + \cdots + \frac{t^k}{k!} \mu'_k + \frac{t^{k+1}}{(k+1)!} \mu'_{k+1} + \cdots \end{aligned}$$

This is the Taylor expansion of $M_X(t)$, and by definition of the derivative:

$$\frac{d^k}{dt^k} M_X(t) = \mu'_k + t \mu'_{k+1} + \text{terms with higher orders of } t \dots$$

⁶Provided the series converges absolutely which it does in some neighborhood around $t = 0$ due to the MGF assumption.

Thus,

$$\mu'_k = \left. \frac{d^k}{dt^k} M_X(t) \right|_{t=0}$$

■

Chapter 4

Common Distributions

4.1 Bernoulli Distribution

A **Bernoulli trial** is a random experiment that has exactly two possible outcomes:

1. **Success**, with probability p ,
2. **Failure**, with probability $1 - p$.

For any Bernoulli trial, we define a random variable X such that if the experiment results in success, then $X = 1$. Otherwise, $X = 0$. It follows that X is a discrete random variable, with probability mass function $p_X(x)$ defined by

$$p_X(x) = P(X = x) = \begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0 \\ 0, & \text{otherwise} \end{cases}$$

This can be compactly written as:

$$p_X(x) = p^x(1 - p)^{1-x}, \quad \text{for } x \in \{0, 1\}$$

The random variable X is said to follow a **Bernoulli distribution** with parameter p , written as:

$$X \sim \text{Bernoulli}(p)$$

Example

In a fair coin toss, the outcomes can be either ‘Head’ or ‘Tail’. If we define a success as getting ‘Head’, then $p = 0.5$. The PMF of the distribution is then given by

$$f_X(x) = P(X = x) = \begin{cases} 0.5, & \text{if } x = 1 \\ 0.5, & \text{if } x = 0 \\ 0, & \text{otherwise} \end{cases}$$



Figure 4.1: *Bernoulli distribution: $p = 0.5$.*

4.1.1 Mean and Variance of the Bernoulli Distribution

Let $X \sim \text{Bernoulli}(p)$.

- **Mean:**

$$\mathbb{E}(X) = \sum_x x \cdot P(X = x) = 0 \cdot (1 - p) + 1 \cdot p = p$$

$$\mathbb{E}(X) = p$$

The mean of a Bernoulli distribution is simply the probability of success, p .

- **Variance:**

$$\text{Var}(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2$$

Note that for a Bernoulli variable, $X^2 = X$ (since X is either 0 or 1), so:

$$\mathbb{E}(X^2) = \mathbb{E}(X) = p$$

$$\text{Var}(X) = p - p^2 = p(1 - p)$$

$$\text{Var}(X) = p(1 - p)$$

The variance of a Bernoulli distribution depends on both the probability of success and failure. It is maximum when $p = 0.5$.

4.2 Binomial Distribution

The Binomial distribution arises from repeating a Bernoulli trial independently n number of times, where each trial has the same probability of success p .

Let X denote the number of successes in n independent Bernoulli trials, where each trial has two outcomes: success (with probability p) and failure (with probability $1 - p$). Then the discrete random variable X follows a **Binomial distribution** with parameters n and p , written as:

$$X \sim \text{Binomial}(n, p)$$

The probability mass function (PMF) of X is given by:

$$p_X(x) = P(X = x) = \begin{cases} \binom{n}{x} p^x (1-p)^{n-x}, & \text{for } x = 0, 1, 2, \dots, n \\ 0, & \text{otherwise} \end{cases}$$

To derive this PMF, consider the following:

- We want the probability of getting exactly x successes (and hence $n - x$ failures) in n trials.
- Each specific sequence of outcomes with x successes and $n - x$ failures has probability:

$$p^x (1-p)^{n-x}$$

because of the independence of trials.

- However, there are multiple ways (distinct sequences) to arrange x successes among n trials. The number of such arrangements is given by the binomial coefficient:

$$\binom{n}{x} = \frac{n!}{x!(n-x)!}$$

Thus, the total probability of getting exactly x successes is:

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

This expression defines the PMF of the binomial distribution.

A Binomial random variable is the sum of independent Bernoulli random variables i.e. if

$$X = \sum_{i=1}^n X_i, \text{ with } X_i \sim \text{Bernoulli}(p)$$

for all $i = 1, 2, \dots, n$, then,

$$X \sim \text{Binomial}(n, p)$$

Example

Suppose a fair coin (with $p = 0.5$) is tossed 4 times. Let X be the number of heads observed. Then $X \sim \text{Binomial}(4, 0.5)$. The PMF is:

$$p_X(x) = \binom{4}{x} \times (0.5)^x \times (0.5)^{4-x}, \quad x = 0, 1, 2, 3, 4$$

Evaluating:

$$p_X(0) = \binom{4}{0} \times (0.5)^0 \times (0.5)^4 = 1 \times 1 \times 0.0625 = 0.0625$$

$$p_X(1) = \binom{4}{1} \times (0.5)^1 \times (0.5)^3 = 4 \times 0.5 \times 0.125 = 0.25$$

$$p_X(2) = \binom{4}{2} \times (0.5)^2 \times (0.5)^2 = 6 \times 0.25 \times 0.25 = 0.375$$

$$p_X(3) = \binom{4}{3} \times (0.5)^3 \times (0.5)^1 = 4 \times 0.125 \times 0.5 = 0.25$$

$$p_X(4) = \binom{4}{4} \times (0.5)^4 \times (0.5)^0 = 1 \times 0.0625 \times 1 = 0.0625$$



Figure 4.2: *Binomial distribution: $n = 4$, $p = 0.5$.*

Example

A fair six-sided die is rolled 8 times. What is the probability that the number 3 or 4 appears exactly 3 times?

Let a ‘success’ be defined as getting either a 3 or a 4 in a single roll. The probability of success on one roll is:

$$p = P(3 \text{ or } 4) = \frac{1}{6} + \frac{1}{6} = \frac{2}{6} = \frac{1}{3}$$

Let X be the number of successes (i.e., times 3 or 4 occurs) in $n = 8$ independent die rolls. Then X follows a binomial distribution:

$$X \sim \text{Binomial}\left(n = 8, p = \frac{1}{3}\right)$$

We want to find the probability three successes i.e. $X = 3$:

$$P(X = 3) = \binom{8}{3} \times \left(\frac{1}{3}\right)^3 \times \left(\frac{2}{3}\right)^5$$

Now, compute the values:

$$\binom{8}{3} = 56, \quad \left(\frac{1}{3}\right)^3 = \frac{1}{27}, \quad \left(\frac{2}{3}\right)^5 = \frac{32}{243}$$

$$P(X = 3) = 56 \times \frac{1}{27} \times \frac{32}{243} = \frac{1792}{6561} \approx 0.273$$

4.2.1 Mean and Variance of the Binomial Distribution

Let $X \sim \text{Binomial}(n, p)$.

- **Mean:**

Consider X as the sum of n independent Bernoulli random variables:

$$X = X_1 + X_2 + \cdots + X_n, \quad \text{where } X_i \sim \text{Bernoulli}(p)$$

Since expectation is linear:

$$\mathbb{E}(X) = \mathbb{E}[X_1 + X_2 + \cdots + X_n] = \mathbb{E}[X_1] + \mathbb{E}[X_2] + \cdots + \mathbb{E}[X_n]$$

Each X_i has expected value p , so:

$$\mathbb{E}(X) = n \cdot p$$

Thus on average, we can expect $n \cdot p$ successes in n trials.

- **Variance:**

Since the X_i 's are independent:

$$\text{Var}(X) = \text{Var}(X_1 + X_2 + \cdots + X_n) = \text{Var}(X_1) + \text{Var}(X_2) + \cdots + \text{Var}(X_n)$$

Each X_i is Bernoulli with variance $p(1 - p)$, so:

$$\text{Var}(X) = n \cdot p(1 - p)$$

The variance increases with the number of trials and depends on both the probability of success and failure.

4.3 Poisson Distribution

The Poisson distribution is commonly used to model the number of occurrences of an event in a fixed interval of time or space, under the following assumptions:

- Events occur independently.
- The average rate (λ) of occurrence is constant over the interval.
- Two events cannot occur at exactly the same instant.

Let X denote the number of such events occurring in a fixed interval with an average value λ , then we say that the discrete random variable X follows a **Poisson distribution** with parameter $\lambda > 0$, and write:

$$X \sim \text{Poisson}(\lambda)$$

The probability mass function (PMF) of X is given by:

$$p_X(x) = P(X = x) = \begin{cases} \frac{e^{-\lambda} \lambda^x}{x!}, & \text{for } x = 0, 1, 2, \dots \\ 0, & \text{otherwise} \end{cases}$$

Example

Suppose a call center receives an average of 4 calls per minute. What is the probability that exactly 2 calls are received in a particular minute?

Let X be the number of calls per minute. Then:

$$X \sim \text{Poisson}(\lambda = 4)$$

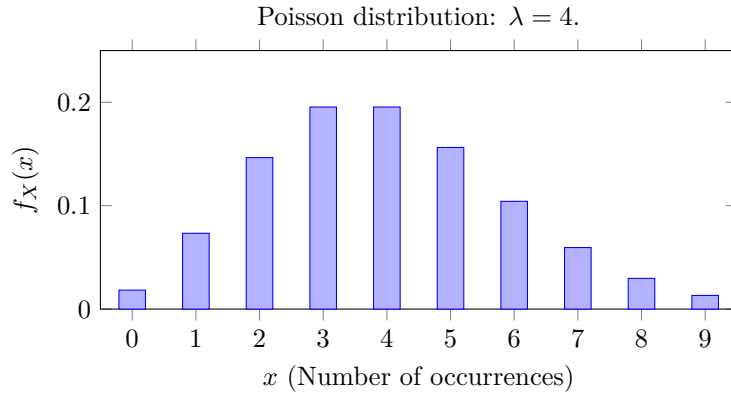


Figure 4.3: *Poisson distribution: $\lambda = 4$.*

We want to find $P(X = 2)$:

$$P(X = 2) = \frac{e^{-4} \cdot 4^2}{2!} = \frac{e^{-4} \cdot 16}{2} = 8e^{-4} \approx 0.1465$$

Example

In a football league, the number of goals scored by a team in a match is modeled using a Poisson distribution. Based on historical performance, Team A scores an average of 2.1 goals per match.

1. What is the probability that Team A scores exactly 3 goals in an upcoming match?
2. What is the probability that Team A scores fewer than 2 goals?
3. What is the probability that Team A scores at least 2 goals?
4. What is the expected number of goals Team A will score over their next 5 matches?

The Poisson probability mass function (PMF) is given by:

$$f_X(x) = P(X = x) = \frac{e^{-\lambda} \lambda^x}{x!} \quad \text{where } \lambda = 2.1, \quad x = 0, 1, 2, \dots$$

1. Probability that Team A scores exactly 3 goals:

$$f_X(3) = \frac{e^{-2.1} \cdot 2.1^3}{3!} = \frac{e^{-2.1} \cdot 9.261}{6} \approx \frac{0.1225 \cdot 9.261}{6} \approx 0.189$$

2. Probability that Team A scores fewer than 2 goals:

$$P(X < 2) = P(X = 0) + P(X = 1) = f_X(0) + f_X(1)$$

$$f_X(0) = e^{-2.1} \approx 0.1225, \quad f_X(1) = \frac{e^{-2.1} \cdot 2.1}{1!} \approx 0.2573$$

$$P(X < 2) \approx 0.1225 + 0.2573 = 0.3798$$

3. Probability that Team A scores at least 2 goals:

$$P(X \geq 2) = 1 - P(X < 2) = 1 - 0.3798 = 0.6202$$

4. Expected number of goals over 5 matches:

$$5 \times \mathbb{E}(X) = 5 \times \lambda = 5 \times 2.1 = 10.5$$

Theorem: The Poisson distribution can be obtained as the limiting distribution of the Binomial distribution when the number of trials $n \rightarrow \infty$, the success probability $p \rightarrow 0$, while the expected number of successes $\lambda = np$ remains constant. Formally,

$$\text{Binomial}(n, p) \longrightarrow \text{Poisson}(\lambda) \quad \text{as } n \rightarrow \infty, p \rightarrow 0 \text{ such that } np = \lambda(\text{constant})$$

Proof: Let $X \sim \text{Binomial}(n, p)$. The probability mass function (PMF) is:

$$p_X(k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k = 0, 1, 2, \dots, n.$$

We assume $n \rightarrow \infty$ and $p \rightarrow 0$ such that the product $\lambda = np$ remains fixed and finite. We can write the binomial coefficient as:

$$\begin{aligned} \binom{n}{x} &= \frac{n(n-1) \cdots (n-x+1)}{x!} \\ &= \frac{n^x}{x!} \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{x-1}{n}\right) \end{aligned}$$

As $n \rightarrow \infty$, each of the product terms approaches 1: Thus,

$$\binom{n}{x} \rightarrow \frac{n^x}{x!}$$

Now using this limiting expression of $\binom{n}{x}$ and replacing p with $\frac{\lambda}{n}$, we get:

$$p_X(x) \approx \frac{n^x}{x!} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x} = \frac{\lambda^x}{x!} \left(1 - \frac{\lambda}{n}\right)^{n-x}$$

Rewrite the last term as:

$$\left(1 - \frac{\lambda}{n}\right)^{n-x} = \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x}$$

As $n \rightarrow \infty$,

$$\left(1 - \frac{\lambda}{n}\right)^n \rightarrow e^{-\lambda} \quad \text{and} \quad \left(1 - \frac{\lambda}{n}\right)^{-x} \rightarrow 1,$$

since x is fixed.

Therefore in the limit $n \rightarrow \infty$:

$$\lim_{n \rightarrow \infty} p_X(x) = \frac{\lambda^x}{x!} e^{-\lambda} \cdot 1 = \frac{e^{-\lambda} \lambda^x}{x!}$$

This matches the PMF of the Poisson distribution with parameter λ . ■

4.3.1 Mean and Variance of the Poisson Distribution

Let $X \sim \text{Poisson}(\lambda)$.

- **Mean:**

$$\begin{aligned}
\mathbb{E}(X) &= \sum_{x=0}^{\infty} x \cdot P(X=x) = \sum_{x=0}^{\infty} x \cdot \frac{e^{-\lambda} \lambda^x}{x!} \\
&= e^{-\lambda} \sum_{x=1}^{\infty} x \cdot \frac{\lambda^x}{x!} = e^{-\lambda} \sum_{x=1}^{\infty} \frac{\lambda^x}{(x-1)!} \\
&= e^{-\lambda} \cdot \lambda \sum_{x=1}^{\infty} \frac{\lambda^{x-1}}{(x-1)!} \\
&= \lambda \cdot e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} \quad (\text{let } k = x-1) \\
&= \lambda \cdot e^{-\lambda} \cdot e^{\lambda} = \lambda
\end{aligned}$$

$$\mathbb{E}(X) = \lambda$$

- **Variance:**

$$\text{Var}(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2$$

First we evaluate

$$\mathbb{E}(X^2) = \sum_{x=0}^{\infty} x^2 \cdot \frac{e^{-\lambda} \lambda^x}{x!}$$

We use the identity $x^2 = x(x-1) + x$, giving:

$$\mathbb{E}(X^2) = \sum_{x=0}^{\infty} [x(x-1) + x] \cdot \frac{e^{-\lambda} \lambda^x}{x!} = e^{-\lambda} \left(\sum_{x=0}^{\infty} \frac{x(x-1)\lambda^x}{x!} + \sum_{x=0}^{\infty} \frac{x\lambda^x}{x!} \right)$$

We compute each sum:

$$\sum_{x=0}^{\infty} \frac{x(x-1)\lambda^x}{x!} = \sum_{x=2}^{\infty} \frac{\lambda^x}{(x-2)!} = \lambda^2 \sum_{j=0}^{\infty} \frac{\lambda^j}{j!} = \lambda^2 e^{\lambda}$$

And,

$$\sum_{x=0}^{\infty} \frac{x\lambda^x}{x!} = \sum_{x=1}^{\infty} \frac{\lambda^x}{(x-1)!} = \lambda \sum_{j=0}^{\infty} \frac{\lambda^j}{j!} = \lambda e^{\lambda}$$

Therefore,

$$\mathbb{E}(X^2) = e^{-\lambda} (\lambda^2 e^{\lambda} + \lambda e^{\lambda}) = \lambda^2 + \lambda$$

Therefore:

$$\text{Var}(X) = (\lambda + \lambda^2) - \lambda^2 = \lambda$$

$$\text{Var}(X) = \lambda$$

There is an alternative way of calculating the variance using the moment generating function (MGF).

The moment generating function (MGF) of X is defined as:

$$M_X(t) = \mathbb{E}(e^{tX}) = \sum_{k=0}^{\infty} e^{tx} \cdot \frac{e^{-\lambda} \lambda^x}{x!}$$

Factor out the constant $e^{-\lambda}$:

$$M_X(t) = e^{-\lambda} \sum_{x=0}^{\infty} \frac{(\lambda e^t)^x}{x!}$$

This is the exponential series:

$$\sum_{x=0}^{\infty} \frac{(\lambda e^t)^x}{x!} = e^{\lambda e^t}$$

Therefore,

$$M_X(t) = e^{-\lambda} \cdot e^{\lambda e^t} = e^{\lambda(e^t-1)}$$

To compute the variance $\text{Var}(X)$, we use the identity:

$$\text{Var}(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2 = \mu'_2 - \mu^2 = \mu'_2 - (\mu'_1)^2$$

We compute the first and second raw moments using derivatives of the MGF:

First raw moment (mean):

$$M'_X(t) = \frac{d}{dt} \left[e^{\lambda(e^t-1)} \right] = \lambda e^t \cdot e^{\lambda(e^t-1)}$$

Evaluating at $t = 0$:

$$\mu'_1 = M'_X(0) = \lambda \cdot 1 \cdot e^{\lambda(1-1)} = \lambda$$

Second raw moment:

$$M''_X(t) = \frac{d}{dt} \left[\lambda e^t \cdot e^{\lambda(e^t-1)} \right] = \lambda e^t \left[\lambda e^t \cdot e^{\lambda(e^t-1)} + e^{\lambda(e^t-1)} \right] = \lambda e^t e^{\lambda(e^t-1)} (\lambda e^t + 1)$$

Evaluating at $t = 0$:

$$\mu'_2 = M''_X(0) = \lambda \cdot 1 \cdot 1 \cdot (\lambda \cdot 1 + 1) = \lambda(\lambda + 1) = \lambda^2 + \lambda$$

Now compute the variance:

$$\text{Var}(X) = \mu'_2 - (\mu'_1)^2 = (\lambda^2 + \lambda) - \lambda^2 = \lambda$$

4.4 Uniform Distribution

The Uniform distribution is the simplest continuous probability distribution, where all outcomes in a given interval are equally likely.

Let X be a continuous random variable that is uniformly distributed on the interval $[a, b]$, where $a < b$. This means that X has constant probability density over this interval. We write:

$$X \sim \text{Uniform}(a, b)$$

The probability density function (PDF) of X is given by:

$$f_X(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{otherwise} \end{cases}$$

That is, the probability density is constant between a and b , and zero elsewhere. The total area under the curve is 1, ensuring it satisfies the definition of a probability density function.

A continuous uniform distribution models situations where every outcome in an interval is equally likely—such as the exact time (within an hour) a bus arrives, or the position of a point randomly placed on a stick.

Example

Suppose that a variable X is uniformly distributed over the interval $[2, 5]$. Then:

$$f_X(x) = \begin{cases} \frac{1}{5-2} = \frac{1}{3}, & 2 \leq x \leq 5 \\ 0, & \text{otherwise} \end{cases}$$



Figure 4.4: Uniform distribution: $a = 2$, $b = 5$.

We can compute probabilities over intervals by integrating the density. For example:

$$P(3 \leq X \leq 4) = \int_3^4 \frac{1}{3} dx = \frac{1}{3}(4 - 3) = \frac{1}{3}$$

—

4.4.1 Mean and Variance of the Uniform Distribution

Let $X \sim \text{Uniform}(a, b)$.

- **Mean:**

$$\begin{aligned} \mathbb{E}(X) &= \int_a^b x \cdot \frac{1}{b-a} dx = \frac{1}{b-a} \int_a^b x dx \\ &= \frac{1}{b-a} \cdot \left[\frac{x^2}{2} \right]_a^b = \frac{1}{b-a} \cdot \left(\frac{b^2 - a^2}{2} \right) \\ &= \frac{1}{b-a} \cdot \frac{(b-a)(b+a)}{2} = \frac{a+b}{2} \end{aligned}$$

$$\mathbb{E}(X) = \frac{a+b}{2}$$

- **Variance:**

$$\text{Var}(X) = \mathbb{E}(X^2) - [\mathbb{E}(X)]^2$$

with:

$$\mathbb{E}(X^2) = \int_a^b x^2 \cdot \frac{1}{b-a} dx = \frac{1}{b-a} \left[\frac{x^3}{3} \right]_a^b = \frac{b^3 - a^3}{3(b-a)}$$

$$\text{Var}(X) = \frac{(b-a)^2}{12}$$

Example

For $X \sim \text{Uniform}(2, 5)$:

$$\mathbb{E}(X) = \frac{2+5}{2} = 3.5, \quad \text{Var}(X) = \frac{(5-2)^2}{12} = \frac{9}{12} = 0.75$$

—

4.5 Normal Distribution

The **Normal distribution**, also known as the **Gaussian distribution**, is one of the most important continuous probability distributions in statistics. It models many naturally occurring phenomena such as heights, test scores, measurement errors, etc. When a continuous random variable X is said to follow a Normal distribution with parameter μ and σ^2 , we denote it as:

$$X \sim \mathcal{N}(\mu, \sigma^2)$$

The probability density function (PDF) of the Normal distribution is given by:

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad \text{for } -\infty \leq x \leq \infty$$

The distribution is completely determined by the parameters μ and σ^2 . In future, we will show that the mean and variance of the normal distribution are those parameters μ and σ^2 respectively.

Example

Let $X \sim \mathcal{N}(2, 1^2)$, i.e., a normal distribution with mean $\mu = 2$ and standard deviation $\sigma = 1$. We want to compute the value of the probability density function (PDF) at $x = 1.5$.

The PDF of a normal distribution is:

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

The plot of the PDF of $X \sim \mathcal{N}(2, 1)$ is shown in the figure below:

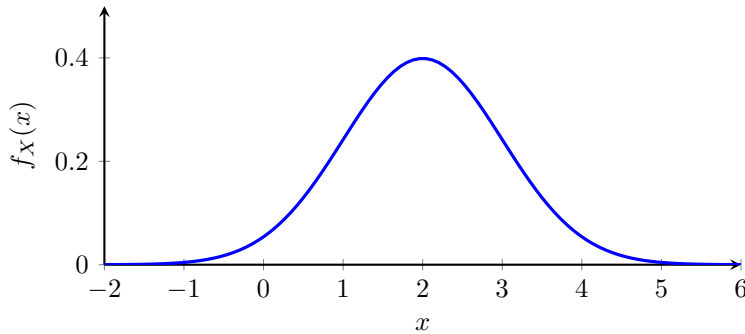


Figure 4.5: Normal distribution $\mathcal{N}(2, 1)$.

Substitute $\mu = 2$, $\sigma = 1$, and $x = 1.5$:

$$f_X(1.5) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(1.5-2)^2}{2}\right) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{0.25}{2}\right) = \frac{1}{\sqrt{2\pi}} \exp(-0.125)$$

Now compute numerically:

$$\begin{aligned}\frac{1}{\sqrt{2\pi}} &\approx 0.3989, & \exp(-0.125) &\approx 0.8825 \\ f_X(1.5) &\approx 0.3989 \times 0.8825 \approx 0.3521\end{aligned}$$

—

4.5.1 Mean and Variance of the Normal Distribution

Let $X \sim \mathcal{N}(\mu, \sigma^2)$.

- **Mean:**

$$\mathbb{E}(X) = \int_{-\infty}^{\infty} x f_X(x) dx = \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx$$

$$\text{Let } u = \frac{x-\mu}{\sigma} \implies x = \sigma u + \mu, \quad dx = \sigma du,$$

$$\begin{aligned}\mathbb{E}(X) &= \int_{-\infty}^{\infty} (\sigma u + \mu) \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{u^2}{2}\right) \sigma du \\ &= \int_{-\infty}^{\infty} (\sigma u + \mu) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) du \\ &= \underbrace{\int_{-\infty}^{\infty} \sigma u \frac{1}{\sqrt{2\pi}} e^{-u^2/2} du}_{=0 \text{ (odd integrand)}} + \mu \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-u^2/2} du \\ &= 0 + \mu \cdot 1 = \mu.\end{aligned}$$

$$\mathbb{E}(X) = \mu$$

- **Variance:**

$$\begin{aligned}\text{Var}(X) &= \mathbb{E}[(X - \mu)^2] = \int_{-\infty}^{\infty} (x - \mu)^2 f_X(x) dx \\ &= \int_{-\infty}^{\infty} (x - \mu)^2 \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx\end{aligned}$$

$$\text{Let } u = \frac{x-\mu}{\sigma} \implies x - \mu = \sigma u, \quad dx = \sigma du,$$

$$\begin{aligned}\text{Var}(X) &= \int_{-\infty}^{\infty} (\sigma u)^2 \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{u^2}{2}\right) \sigma du \\ &= \int_{-\infty}^{\infty} \sigma^2 u^2 \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) du \\ &= \sigma^2 \underbrace{\int_{-\infty}^{\infty} u^2 \frac{1}{\sqrt{2\pi}} e^{-u^2/2} du}_{\text{We need to prove this equals 1}} \\ &= \frac{\sigma^2}{\sqrt{2\pi}} \int_{-\infty}^{\infty} u^2 e^{-u^2/2} du \\ &= \frac{2\sigma^2}{\sqrt{2\pi}} \int_0^{\infty} u^2 e^{-u^2/2} du\end{aligned}$$

Now, using the Gamma integral formula¹,

$$\begin{aligned}\int_0^\infty u^2 e^{-u^2/2} du &= \frac{1}{2} \left(\frac{1}{2}\right)^{-\frac{3}{2}} \Gamma\left(\frac{3}{2}\right) \\ &= \frac{1}{2} \cdot 2^{3/2} \cdot \frac{1}{2} \sqrt{\pi} \quad \left(\left(\frac{1}{2}\right)^{-3/2} = 2^{3/2}, \Gamma\left(\frac{3}{2}\right) = \frac{1}{2} \Gamma\left(\frac{1}{2}\right) = \frac{1}{2} \sqrt{\pi}\right) \\ &= \frac{\sqrt{2\pi}}{2}\end{aligned}$$

Therefore,

$$\text{Var}(X) = \sigma^2 \frac{2}{\sqrt{2\pi}} \times \frac{\sqrt{2\pi}}{2} = \sigma^2$$

$$\text{Var}(X) = \sigma^2$$

4.5.2 Moment Generating Function of the Normal Distribution

Let $X \sim \mathcal{N}(\mu, \sigma^2)$ be a normal random variable with mean μ and variance σ^2 . The moment generating function (MGF) of X is defined as

$$\begin{aligned}M_X(t) &= \mathbb{E}(e^{tX}) = \int_{-\infty}^\infty e^{tx} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^\infty \exp\left(tx - \frac{(x-\mu)^2}{2\sigma^2}\right) dx\end{aligned}$$

Rewrite the exponent as:

$$\begin{aligned}tx - \frac{(x-\mu)^2}{2\sigma^2} &= -\frac{1}{2\sigma^2}(x-\mu)^2 + tx \\ &= -\frac{1}{2\sigma^2}(x^2 - 2\mu x + \mu^2) + tx \\ &= -\frac{x^2}{2\sigma^2} + \frac{\mu x}{\sigma^2} - \frac{\mu^2}{2\sigma^2} + tx \\ &= -\frac{x^2}{2\sigma^2} + \left(\frac{\mu}{\sigma^2} + t\right)x - \frac{\mu^2}{2\sigma^2} \quad (\text{Group terms involving } x) \\ &= -\frac{1}{2\sigma^2} [x^2 - 2(\mu + \sigma^2 t)x] - \frac{\mu^2}{2\sigma^2} \\ &= -\frac{1}{2\sigma^2} [(x - (\mu + \sigma^2 t))^2 - (\mu + \sigma^2 t)^2] - \frac{\mu^2}{2\sigma^2} \\ &= -\frac{(x - (\mu + \sigma^2 t))^2}{2\sigma^2} + \frac{(\mu + \sigma^2 t)^2}{2\sigma^2} - \frac{\mu^2}{2\sigma^2} \\ &= -\frac{(x - (\mu + \sigma^2 t))^2}{2\sigma^2} + \frac{\mu^2 + 2\mu\sigma^2 t + \sigma^4 t^2 - \mu^2}{2\sigma^2} \\ &= -\frac{(x - (\mu + \sigma^2 t))^2}{2\sigma^2} + \mu t + \frac{1}{2}\sigma^2 t^2\end{aligned}$$

Substitute this back into the integral for the MGF:

$$M_X(t) = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^\infty \exp\left(-\frac{(x - (\mu + \sigma^2 t))^2}{2\sigma^2} + \mu t + \frac{1}{2}\sigma^2 t^2\right) dx$$

¹Gamma integral formula:

$$\begin{aligned}\int_0^\infty x^n e^{-ax^2} dx &= \frac{1}{2} a^{-\frac{n+1}{2}} \Gamma\left(\frac{n+1}{2}\right) \\ \Gamma(n+1) &= n\Gamma(n), \quad \Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}\end{aligned}$$

Factor out terms that do not depend on x :

$$M_X(t) = \exp\left(\mu t + \frac{1}{2}\sigma^2 t^2\right) \cdot \underbrace{\frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} \exp\left(-\frac{(x - (\mu + \sigma^2 t))^2}{2\sigma^2}\right) dx}_{=1}$$

The integral is the integral of a normal pdf with mean $\mu + \sigma^2 t$ and variance σ^2 , and is equal to 1. Hence,

$$M_X(t) = \exp\left(\mu t + \frac{1}{2}\sigma^2 t^2\right)$$

This moment generating function exists for all real values of t and uniquely characterizes the normal distribution.

4.5.3 Properties of the Normal Distribution

1. The **support** of a normal random variable $X \sim \mathcal{N}(\mu, \sigma^2)$ is the entire real line:

$$(-\infty \leq X \leq +\infty)$$

This reflects that, however unlikely, arbitrarily large positive or negative values can occur.

2. The probability density function is perfectly **symmetric** about its mean μ since,

$$f_X(\mu + x_0) = f_X(\mu - x_0) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{x_0^2}{2\sigma^2}\right) \quad \forall x_0 \in \mathbb{R}.$$

As a result, the left and right tails of the distribution mirror each other.

3. Since the distribution is symmetrical about μ , its mean and median coincide. To get the mode, we need to calculate the peak point of the distribution.

$$\begin{aligned} f_X(x) &= \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \\ \frac{d}{dx} f_X(x) &= \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \cdot \frac{d}{dx} \left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \\ &= f_X(x) \left(-\frac{2(x - \mu)}{2\sigma^2}\right) = -\frac{x - \mu}{\sigma^2} f_X(x) \end{aligned}$$

Setting the derivative to zero for a stationary point:

$$\begin{aligned} -\frac{x - \mu}{\sigma^2} f_X(x) = 0 &\implies x - \mu = 0 \\ &\implies x = \mu \end{aligned}$$

Now,

$$\begin{aligned} f_X''(x) &= -\frac{1}{\sigma^2} f_X(x) + \left(-\frac{x - \mu}{\sigma^2}\right) f_X'(x) \\ &= -\frac{1}{\sigma^2} f_X(x) + \left(-\frac{x - \mu}{\sigma^2}\right) \left(-\frac{x - \mu}{\sigma^2} f_X(x)\right) \\ &= -\frac{1}{\sigma^2} f_X(x) + \frac{(x - \mu)^2}{\sigma^4} f_X(x) \\ &= \frac{(x - \mu)^2 - \sigma^2}{\sigma^4} f_X(x) \end{aligned}$$

Evaluating at the stationary point $x = \mu$:

$$f_X''(\mu) = \frac{(\mu - \mu)^2 - \sigma^2}{\sigma^4} f_X(\mu) = -\frac{1}{\sigma^2} f_X(\mu) < 0,$$

Hence the peak (mode) of the normal density occurs at $x = \mu$.

For normal distribution, all measures of central tendency coincide:

$$\text{Mean} = \text{Median} = \text{Mode} = \mu$$

4. A normal distribution $X \sim \mathcal{N}(\mu, \sigma^2)$ satisfies the following empirical rules:

Empirical Rule (68-95-99.7 Rule):

- About 68.27% of the values lie within σ of the mean ($\mu \pm \sigma$).
- About 95.45% of the values lie within 2σ of the mean ($\mu \pm 2\sigma$).
- About 99.73% of the values lie within 3σ of mean ($\mu \pm 3\sigma$).



Figure 4.6: 68-95-99.7 rule.

5. The normal distribution curve has two **points of inflection**² at a distance σ on either side of μ i.e. at

$$x = \mu \pm \sigma.$$

At these points the second derivative of $f_X(x)$ vanishes, marking the transition between “concave down” near the center and “concave up” in the tails. To show that let’s take the second derivative of $f_X(x)$ and equate it to zero:

$$\begin{aligned} f_X''(x) &= \frac{(x - \mu)^2 - \sigma^2}{\sigma^4} f_X(x) = 0 \\ \Rightarrow (x - \mu)^2 &= \sigma^2 \\ \Rightarrow x - \mu &= \pm \sigma \\ \Rightarrow x &= \mu \pm \sigma \end{aligned}$$

²A **point of inflection** of a function $f(x)$ is a point $x = a$ such that

$$f''(a) = 0,$$

At the point of inflexion, the second derivative $f''(x)$ changes sign as x passes through a , meaning the curve switches between concave-up and concave-down at $x = a$.

6. All odd central moments are zero (due to symmetry), and the even central moments have closed-form expressions:

$$E[(X - \mu)^{2n+1}] = 0, \quad E[(X - \mu)^{2n}] = \sigma^{2n} (2n - 1)!!, \quad n = 1, 2, \dots$$

In particular, the variance is $E[(X - \mu)^2] = \sigma^2$, and the fourth central moment is $3\sigma^4$, etc.

7. The kurtosis of normal distribution is 3. (prove it)
 8. The linear combination of two independent normal variables is also normal.

Theorem: Let $X \sim N(\mu_1, \sigma_1^2)$ and $Y \sim N(\mu_2, \sigma_2^2)$ be two independent normal random variables. Then for any real constants a and b , the linear combination

$$Z = aX + bY$$

is also normally distributed with

$$Z \sim N(a\mu_1 + b\mu_2, a^2\sigma_1^2 + b^2\sigma_2^2)$$

Proof: Since $X \sim N(\mu_1, \sigma_1^2)$, its moment generating function (MGF) is

$$M_X(t) = \exp\left(\mu_1 t + \frac{1}{2}\sigma_1^2 t^2\right)$$

Similarly, the MGF of $Y \sim N(\mu_2, \sigma_2^2)$ is

$$M_Y(t) = \exp\left(\mu_2 t + \frac{1}{2}\sigma_2^2 t^2\right)$$

Consider $Z = aX + bY$. Since X and Y are independent, the MGF of Z is:

$$M_Z(t) = \mathbb{E}\left(e^{t(aX+bY)}\right) = \mathbb{E}\left(e^{taX}\right) \cdot \mathbb{E}\left(e^{tbY}\right) = M_X(at) \cdot M_Y(bt)$$

Substituting the MGFs:

$$M_Z(t) = \exp\left(a\mu_1 t + \frac{1}{2}a^2\sigma_1^2 t^2\right) \cdot \exp\left(b\mu_2 t + \frac{1}{2}b^2\sigma_2^2 t^2\right)$$

Combining the exponents:

$$M_Z(t) = \exp\left((a\mu_1 + b\mu_2)t + \frac{1}{2}(a^2\sigma_1^2 + b^2\sigma_2^2)t^2\right)$$

This is the MGF of a normal distribution with mean $a\mu_1 + b\mu_2$ and variance $a^2\sigma_1^2 + b^2\sigma_2^2$. Therefore,

$$Z \sim N(a\mu_1 + b\mu_2, a^2\sigma_1^2 + b^2\sigma_2^2)$$

■

4.5.4 Standard Normal Distribution:

When $\mu = 0$ and $\sigma^2 = 1$, the normal distribution is called the **standard normal distribution**, denoted as:

$$Z \sim \mathcal{N}(0, 1)$$

Its PDF becomes:



Figure 4.7: *Standard normal distribution function.*

$$f_Z(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right)$$

Cumulative distribution:

$$F_Z(z) = P(Z \leq z) = \int_{-\infty}^z f_Z(t) dt$$

In many standard textbook, the CDF of the standard normal distribution is denoted by a special symbol $\Phi(\cdot)$ such that

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z \exp\left(-\frac{t^2}{2}\right) dt$$

Theorem. Let $X \sim \mathcal{N}(\mu, \sigma^2)$ be a normally distributed random variable and let $Z \sim \mathcal{N}(0, 1)$ be a standard normal random variable. Then, for any real number k ,

$$F_X(k) = \Phi\left(\frac{k - \mu}{\sigma}\right)$$

where F_X and Φ denote the cumulative distribution functions of X and Z , respectively.

Proof: Define the function

$$z = g(x) = \frac{x - \mu}{\sigma},$$

which is strictly increasing (since $\sigma > 0$) and differentiable, with inverse

$$g^{-1}(y) = x = \sigma y + \mu$$

Set

$$Y = g(X) = \frac{X - \mu}{\sigma}$$

Then Y has the same distribution as Z , i.e. $Y \sim \mathcal{N}(0, 1)$. By the change-of-variable theorem for CDFs (strictly increasing case),

$$\Phi(z) = P(Z \leq z) = F_X(g^{-1}(z))$$

Hence,

$$\Phi(z) = F_X(\sigma z + \mu)$$

Now replace z by $\frac{k - \mu}{\sigma}$. Since $\sigma \cdot \frac{k - \mu}{\sigma} + \mu = k$, we obtain

$$\Phi\left(\frac{k - \mu}{\sigma}\right) = F_X(k)$$

as required.

Any normal random variable $X \sim \mathcal{N}(\mu, \sigma^2)$ can be converted to a **standard normal variable** using the transformation:

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

Theorem: Let $Z \sim \mathcal{N}(0, 1)$ be a standard normal random variable with CDF $\Phi(z)$. Then, for any real number k ,

$$\Phi(-k) = 1 - \Phi(k)$$

Since $f_Z(z)$ is symmetrical about zero, so for any k

$$f_Z(-k) = f_Z(k)$$

Now,

$$\begin{aligned}\Phi(-k) &= P(Z \leq -k) \\ &= \int_{-\infty}^{-k} f_Z(z) dz\end{aligned}$$

Change variable $t = -z$, so when $z = -\infty \rightarrow t = +\infty$, and $z = -k \rightarrow t = k$, with $dz = -dt$:

$$\begin{aligned}\Phi(-k) &= \int_{\infty}^k f_Z(-t) (-dt) \\ &= - \int_{\infty}^k f_Z(t) dt \\ &= \int_k^{\infty} f_Z(t) dt \\ &= P(Z \geq k) \\ &= 1 - P(Z < k) \\ &= 1 - \Phi(k)\end{aligned}$$

4.5.5 Standard Normal Table

The **standard normal table** (or Z -table) shown in Table 4.1 is used to quickly find cumulative probabilities for the standard normal distribution $Z \sim \mathcal{N}(0, 1)$ without evaluating the integral

$$\int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$$

by hand. By converting any normal random variable $X \sim \mathcal{N}(\mu, \sigma^2)$ into the standard form $Z = (X - \mu)/\sigma$, one can look up probabilities such as $P(X \leq x)$ in a single, universal table, greatly simplifying calculations in statistical inference and hypothesis testing.

To look up $\Phi(k)$, follow the following steps:

1. Write k to **two decimal places**, e.g. $k = 1.23$.
2. Split into

$$\text{row part} = 1.2, \quad \text{column part} = 0.03$$

3. In Table 4.1, go to the row labeled “1.2” and the column labeled “0.03”. The entry at their intersection is

$$\Phi(1.23) = P(Z \leq 1.23) = 0.89065$$

z	0	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0	0.5	0.50399	0.50798	0.51197	0.51595	0.51994	0.52392	0.5279	0.53188	0.53586
0.1	0.53983	0.5438	0.54776	0.55172	0.55567	0.55962	0.56356	0.56749	0.57142	0.57535
0.2	0.57926	0.58317	0.58706	0.59095	0.59483	0.59871	0.60257	0.60642	0.61026	0.61409
0.3	0.61791	0.62172	0.62552	0.6293	0.63307	0.63683	0.64058	0.64431	0.64803	0.65173
0.4	0.65542	0.6591	0.66276	0.6664	0.67003	0.67364	0.67724	0.68082	0.68439	0.68793
0.5	0.69146	0.69497	0.69847	0.70194	0.7054	0.70884	0.71226	0.71566	0.71904	0.7224
0.6	0.72575	0.72907	0.73237	0.73565	0.73891	0.74215	0.74537	0.74857	0.75175	0.7549
0.7	0.75804	0.76115	0.76424	0.7673	0.77035	0.77337	0.77637	0.77935	0.7823	0.78524
0.8	0.78814	0.79103	0.79389	0.79673	0.79955	0.80234	0.80511	0.80785	0.81057	0.81327
0.9	0.81594	0.81859	0.82121	0.82381	0.82639	0.82894	0.83147	0.83398	0.83646	0.83891
1	0.84134	0.84375	0.84614	0.84849	0.85083	0.85314	0.85543	0.85769	0.85993	0.86214
1.1	0.86433	0.8665	0.86864	0.87076	0.87286	0.87493	0.87698	0.879	0.881	0.88298
1.2	0.88493	0.88686	0.88877	0.89065	0.89251	0.89435	0.89617	0.89796	0.89973	0.90147
1.3	0.9032	0.9049	0.90658	0.90824	0.90988	0.91149	0.91309	0.91466	0.91621	0.91774
1.4	0.91924	0.92073	0.9222	0.92364	0.92507	0.92647	0.92785	0.92922	0.93056	0.93189
1.5	0.93319	0.93448	0.93574	0.93699	0.93822	0.93943	0.94062	0.94179	0.94295	0.94408
1.6	0.9452	0.9463	0.94738	0.94845	0.9495	0.95053	0.95154	0.95254	0.95352	0.95449
1.7	0.95543	0.95637	0.95728	0.95818	0.95907	0.95994	0.9608	0.96164	0.96246	0.96327
1.8	0.96407	0.96485	0.96562	0.96638	0.96712	0.96784	0.96856	0.96926	0.96995	0.97062
1.9	0.97128	0.97193	0.97257	0.9732	0.97381	0.97441	0.975	0.97558	0.97615	0.9767
2	0.97725	0.97778	0.97831	0.97882	0.97932	0.97982	0.9803	0.98077	0.98124	0.98169
2.1	0.98214	0.98257	0.983	0.98341	0.98382	0.98422	0.98461	0.985	0.98537	0.98574
2.2	0.9861	0.98645	0.98679	0.98713	0.98745	0.98778	0.98809	0.9884	0.9887	0.98899
2.3	0.98928	0.98956	0.98983	0.9901	0.99036	0.99061	0.99086	0.99111	0.99134	0.99158
2.4	0.9918	0.99202	0.99224	0.99245	0.99266	0.99286	0.99305	0.99324	0.99343	0.99361
2.5	0.99379	0.99396	0.99413	0.9943	0.99446	0.99461	0.99477	0.99492	0.99506	0.9952
2.6	0.99534	0.99547	0.9956	0.99573	0.99585	0.99598	0.99609	0.99621	0.99632	0.99643
2.7	0.99653	0.99664	0.99674	0.99683	0.99693	0.99702	0.99711	0.9972	0.99728	0.99736
2.8	0.99744	0.99752	0.9976	0.99767	0.99774	0.99781	0.99788	0.99795	0.99801	0.99807
2.9	0.99813	0.99819	0.99825	0.99831	0.99836	0.99841	0.99846	0.99851	0.99856	0.99861
3	0.99865	0.99869	0.99874	0.99878	0.99882	0.99886	0.99889	0.99893	0.99896	0.999
3.1	0.99903	0.99906	0.9991	0.99913	0.99916	0.99918	0.99921	0.99924	0.99926	0.99929
3.2	0.99931	0.99934	0.99936	0.99938	0.9994	0.99942	0.99944	0.99946	0.99948	0.9995
3.3	0.99952	0.99953	0.99955	0.99957	0.99958	0.9996	0.99961	0.99962	0.99964	0.99965
3.4	0.99966	0.99968	0.99969	0.9997	0.99971	0.99972	0.99973	0.99974	0.99975	0.99976
3.5	0.99977	0.99978	0.99978	0.99979	0.9998	0.99981	0.99981	0.99982	0.99983	0.99983
3.6	0.99984	0.99985	0.99985	0.99986	0.99986	0.99987	0.99987	0.99988	0.99988	0.99989
3.7	0.99989	0.9999	0.9999	0.9999	0.99991	0.99991	0.99992	0.99992	0.99992	0.99992
3.8	0.99993	0.99993	0.99993	0.99994	0.99994	0.99994	0.99994	0.99995	0.99995	0.99995
3.9	0.99995	0.99995	0.99996	0.99996	0.99996	0.99996	0.99996	0.99996	0.99997	0.99997
4	0.99997	0.99997	0.99997	0.99997	0.99997	0.99997	0.99998	0.99998	0.99998	0.99998

Table 4.1: *Standard normal CDF values* $\Phi(z) = P(Z \leq z)$.

4. For negative z , use symmetry:

$$\Phi(-k) = P(Z \leq -k) = 1 - P(Z \leq k) = 1 - \Phi(k)$$

5. For right-tail probabilities,

$$P(Z > k) = 1 - \Phi(k)$$

6. For probabilities within a specified interval of Z values,

$$P(k_1 \leq Z \leq k_2) = P(Z \leq k_2) - P(Z \leq k_1) = \Phi(k_2) - \Phi(k_1)$$

Example

Suppose the heights of adult males are normally distributed with mean $\mu = 175$ cm and standard deviation $\sigma = 10$ cm. Let X denote the height of a randomly chosen male. Then:

$$X \sim \mathcal{N}(175, 100)$$

What is the probability that a randomly chosen male is taller than 190 cm?

We standardize:

$$Z = \frac{190 - 175}{10} = 1.5$$

Using the standard normal table:

$$P(X > 190) = P(Z > 1.5) = 1 - \Phi(1.5) \approx 1 - 0.9332 = 0.0668$$

Thus, approximately 6.68% of adult males are taller than 190 cm.



Figure 4.8: *Distribution of the heights of adult males.*

4.5.6 Critical Points of the Standard Normal Distribution

In the standard normal distribution, certain values on the horizontal axis divide the distribution into regions with specified probabilities. These special values are called **critical points**.

One-sided critical point

For a given probability α , the one-sided critical point z_α is such that

$$P(Z > z_\alpha) = \alpha$$

or equivalently,

$$P(Z \leq z_\alpha) = 1 - \alpha$$

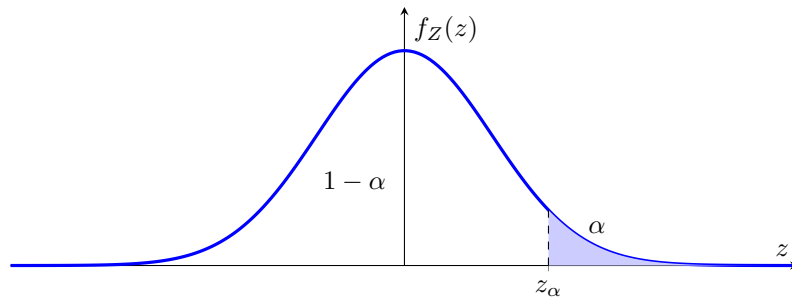


Figure 4.9: *One-sided critical point.*

- For $\alpha = 0.05$, the one-sided critical point is approximately $z_{0.05} \approx 1.645$.
- For $\alpha = 0.01$, the one-sided critical point is approximately $z_{0.01} \approx 2.326$.

Two-sided critical points

For a given probability α , the two-sided critical points $\pm z_{\alpha/2}$ are such that

$$P(Z > z_{\alpha/2}) = \alpha/2, \quad P(Z < -z_{\alpha/2}) = \alpha/2$$

or equivalently,

$$P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) = 1 - \alpha$$

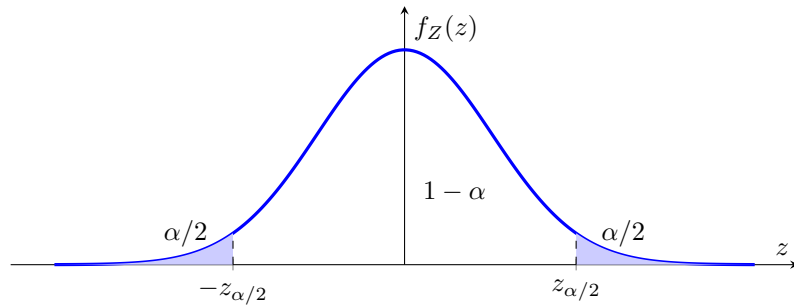


Figure 4.10: *Two-sided critical points.*

- For $\alpha = 0.05$, the two-sided critical points are approximately ± 1.96 .
- For $\alpha = 0.01$, the two-sided critical points are approximately ± 2.576 .

These critical points help us identify how extreme a value is relative to the overall distribution, and they are fundamental when constructing intervals to estimate population parameters. These critical points are used in hypothesis tests and to construct confidence intervals which we will see in future chapters.

Chapter 5

Sampling Theory

5.1 Introduction

Sampling theory is a part of statistics that helps us understand how to learn about a large group by looking at just a small part of it. For example, imagine a company makes a new type of battery and wants to know how long the batteries last. Testing every single battery would take too much time and money, so the company picks a few batteries to test. These few batteries are called a **sample**, and all the batteries made by the company are called the **population**.

The company really wants to know the average lifetime of all batteries—this is called a **parameter**. But since they can't test them all, they use the average lifetime from the sample—this is called a **statistic**.

Sampling theory helps us understand how close this statistic is likely to be to the real average. It also helps us decide how many batteries to test and how to choose them so we get useful, reliable results.

Population: The entire group of individuals or items that we want to learn about. Example: All the batteries produced by a company.

Sample: A smaller group taken from the population, which is actually tested or studied. Example: 100 batteries chosen from the whole production batch.

Parameter: A numerical value that describes a characteristic of the population (usually unknown). Example: The true average lifetime of all the batteries.

Statistic: A numerical value that describes a characteristic of the sample (used to estimate the parameter). Example: The average lifetime of the 100 batteries tested.

5.2 Sampling Methods

5.2.1 Simple Random Sampling

A **simple random sampling** is one in which every member of the population has an equal chance of being selected in the sample.

Mathematically, this means that each possible sample of size n from a population has the same probability of being chosen. For example, suppose a factory produces 10,000 batteries in a day. To estimate the average lifespan, 100 batteries are selected randomly so that each battery has the same

chance of inclusion. An important advantage of simple random sampling is that it is straightforward to analyze using statistical theory, which makes inference about the population simpler

In this text, we will limit our discussion to **simple random sampling**. Before a random sample of size n is selected, the observations are modeled as the random variables X_1, X_2, \dots, X_n . For example, if we randomly select 5 light bulbs from a production batch, their lifespans can be represented by the random variables X_1, X_2, X_3, X_4, X_5 , each denoting the lifespan (in hours) of a selected bulb.

$$\begin{aligned} X_1 &= \text{Lifespan (in hours) of 1st selected bulb} \\ X_2 &= \text{Lifespan (in hours) of 2nd selected bulb} \\ &\dots \\ X_5 &= \text{Lifespan (in hours) of 5th selected bulb} \end{aligned}$$

Assume a first draw yields the following lifespans (in hours) for $n = 5$ randomly selected light bulbs:

$$\text{Draw 1: } \{X_1 = 1200, X_2 = 1140, X_3 = 1180, X_4 = 1300, X_5 = 1250\}$$

Because each sample is chosen at random, a fresh draw of five bulbs would almost surely yield different numerical values for X_1, X_2, \dots, X_5 . Assume a second draw produces:

$$\text{Draw 2: } \{X_1 = 1400, X_2 = 1550, X_3 = 1200, X_4 = 1420, X_5 = 1380\}$$

In this way, each X_i behaves as a genuine random variable, capturing the uncertainty inherent in the sampling process.

There are two main types of simple random sampling:

1. **Simple Random Sampling With Replacement (SRSWR)**: This is a method of selecting a sample of size n from a population of size N one by one such that after each stage of selection, the element is returned to the population before the next draw. Because each selection is made from the full population, the sample observations X_1, X_2, \dots, X_n are *independent and identically distributed (i.i.d.)*¹ random variables following the population distribution.
2. **Simple Random Sampling Without Replacement (SRSWOR)**: This is a method of selecting a sample of size n from a population of size N one by one such that after each stage of selection, the element is not returned to the population. So there is no chance of a particular item being selected twice in the sample. Although the sample observations X_1, X_2, \dots, X_n are identically distributed (each has the same marginal distribution), they are *not independent*, due to the changing composition of the population after each draw.

5.2.2 Other Sampling Methods

Stratified Sampling

In stratified sampling, the population is divided into distinct subgroups or strata based on a specific characteristic (e.g., age, income, region), and a random sample is drawn from each stratum. This method ensures representation from all key subgroups.

- *Example*: A company wants to sample employee opinions. Employees are divided into departments (e.g., HR, Sales, R&D), and a random sample is taken from each department.
- *Advantages*: Increases accuracy by reducing variability; ensures important groups are represented.
- *Disadvantages*: Requires knowledge of strata and population characteristics in advance.

¹**Independent and Identically Distributed (i.i.d.)** is a fundamental assumption in statistics. *Identically distributed* means that each random variable X_i follows the same probability distribution (e.g., normal, binomial). *Independent* means the outcome of one observation does not influence or provide information about the others; knowing X_1 gives no information about X_2, X_3 , etc.

Systematic Sampling

Systematic sampling selects every k -th individual from a population list after a random starting point. The interval k is calculated by dividing the population size by the desired sample size.

- *Example:* If a company has a list of 1,000 employees and wants to survey 100, it selects a random starting point between 1 and 10, then picks every 10th employee on the list.
- *Advantages:* Simple and quick to implement; useful when population is ordered.
- *Disadvantages:* Can introduce bias if there is a hidden pattern in the population that coincides with the sampling interval.

Cluster Sampling

In cluster sampling, the population is divided into clusters (often based on geography or natural groupings). A few clusters are randomly selected, and then all individuals within those clusters are included in the sample.

- *Example:* A research team wants to survey households in a city. The city is divided into neighborhoods (clusters), a few neighborhoods are selected at random, and all households in those neighborhoods are surveyed.
- *Advantages:* Cost-effective and practical for large, dispersed populations.
- *Disadvantages:* Can lead to higher sampling error if clusters are not homogeneous.

Multistage Sampling

Multistage sampling combines several sampling techniques. Typically, it begins with cluster sampling to select large groups, and then simple random or stratified sampling is used within those groups.

- *Example:* In a national education survey, schools are randomly selected (cluster sampling), then students within each selected school are randomly chosen (simple random sampling).
- *Advantages:* Flexible and practical for large-scale surveys; reduces cost and time.
- *Disadvantages:* More complex design and analysis; potential for increased sampling error if stages are not carefully planned.

5.3 Sample Mean, Sample Variance and Sample Proportion

Let X_1, X_2, \dots, X_n be a random sample of size n drawn from a population which are modeled as random variables. The **sample mean** is defined as:

$$\text{Sample Mean} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

It represents the average of the observed sample values.

The **sample variance** is defined as:

$$\text{Sample Variance} = S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

This measures the spread or variability of the sample values around the sample mean. The denominator $n - 1$ (instead of n) ensures that S^2 is an *unbiased estimator* of the population variance σ^2 . We will discuss the concept of unbiased estimator in later chapter.

The **sample standard deviation** is the positive square root of the sample variance:

$$\text{Sample Standard Deviation} = S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2}$$

Now let's suppose X_i is modeled as a binary indicator variable where

$$X_i = \begin{cases} 1, & \text{if the } i\text{-th observation has the characteristic of interest} \\ 0, & \text{otherwise} \end{cases}$$

The **sample proportion** for the characteristic of interest is defined as:

$$\text{Sample Proportion} = \hat{p} = \frac{1}{n} \sum_{i=1}^n X_i = \frac{X}{n}$$

It represents the fraction of the sample exhibiting the characteristic of interest and serves as an estimator of the population proportion p .

5.4 Sampling Distributions

The value of any statistic (e.g. sample mean) will vary from sample to sample.

The **sampling distribution** of a statistic is the probability distribution of the statistic's values computed from all possible random samples of the same size taken from a given population.

Suppose a factory produces thousands of batteries, and the lifetimes of these batteries follow a distribution with a population mean $\mu = 100$ hours and a population standard deviation $\sigma = 20$ hours.

Now, imagine taking a random sample of 5 batteries and computing the average lifetime (sample mean). You repeat this process many times—each time taking a new random sample of 5 batteries and calculating its mean. Each of these sample means will be a bit different due to natural variation in the samples. If you plot all these sample means on a graph, the result is the **sampling distribution of the sample mean**.

The standard deviation of the sampling distribution of a statistic is given a specific name — it is called the **standard error** of that sample statistic.

The **standard error** of a sample statistic is the standard deviation of its sampling distribution. It measures how much the statistic is expected to vary from sample to sample due to random chance.

5.5 The Sampling Distribution of the Sample Mean

Theorem: Let X_1, X_2, \dots, X_n be random samples of size n *chosen with replacements* from a population with mean μ and variance σ^2 , then

$$\mathbb{E}(\bar{X}) = \mu, \text{ and } \text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

Proof: We assume the population consists of N elements $\{y_1, y_2, \dots, y_N\}$. The population mean and population variance are defined as:

$$\mu = \frac{1}{N} \sum_{j=1}^N y_j, \quad \sigma^2 = \frac{1}{N} \sum_{j=1}^N (y_j - \mu)^2$$

The sample mean is defined as:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

By the linearity of expectation:

$$\mathbb{E}(\bar{X}) = \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i)$$

Since each X_i is drawn from the population $\{y_1, y_2, \dots, y_N\}$ each with probability $\frac{1}{N}$. Hence, we have for all i :

$$\mathbb{E}(X_i) = \sum_{j=1}^N y_j \cdot \underbrace{(X_i = y_j)}_{1/N} = \frac{1}{N} \sum_{j=1}^N y_j = \mu$$

So:

$$\mathbb{E}(\bar{X}) = \frac{1}{n} \cdot n \cdot \mu = \mu$$

Using the formula for the variance of a sum of independent random variables:

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i)$$

Now,

$$\text{Var}(X_i) = \sum_{j=1}^N (y_j - \mu)^2 \cdot \underbrace{P(X_i = y_j)}_{1/N} = \frac{1}{N} \sum_{j=1}^N (y_j - \mu)^2 = \sigma^2$$

Therefore,

$$\text{Var}(\bar{X}) = \frac{1}{n^2} \cdot n \cdot \sigma^2 = \frac{\sigma^2}{n}$$

■

Theorem: Let X_1, X_2, \dots, X_n be random samples of size n *chosen without replacement* from a population of size N with mean μ and variance σ^2 , then

$$\mathbb{E}(\bar{X}) = \mu, \text{ and } \text{Var}(\bar{X}) = \frac{\sigma^2}{n} \left(\frac{N-n}{N-1} \right)$$

We skip the proof as it is beyond the scope of this text. The term

$$\frac{N-n}{N-1}$$

is often called **finite population correction factor**, is close to 1 (and can be omitted for most practical purposes) unless the sample constitutes a substantial portion of the population.

5.5.1 Sampling from a Normal Distribution

In the preceding discussion, no specific assumptions were made about the actual distribution of the population from which the observations X_1, X_2, \dots, X_n were sampled. Nevertheless, we know two key characteristics of the sampling distribution of the sample mean \bar{X} :

- Its expected value: $\mathbb{E}(\bar{X})$
- Its variance: $\text{Var}(\bar{X})$

But what about the shape of the sampling distribution? If the population itself is normally distributed, then the sampling distribution of \bar{X} is also normal, regardless of the sample size.

Theorem: When sampling is done from a normal distribution with mean μ and standard deviation σ , the sample mean \bar{X} follows a normal distribution:

$$\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

Proof: Let $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$ be a random sample of size n from a normal population with mean μ and variance σ^2 . Define the sample mean as:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

Since each X_i is normally distributed and is independent, the sample mean \bar{X} is a linear combination of independent normal random variables:

$$\bar{X} = \frac{1}{n}X_1 + \frac{1}{n}X_2 + \dots + \frac{1}{n}X_n$$

A linear combination of independent normal variables is also normally distributed. Therefore, $\bar{X} \sim \mathcal{N}(\mathbb{E}(\bar{X}), \text{Var}(\bar{X}))$.

From the previous theorem, we already know,

$$\mathbb{E}[\bar{X}] = \mu \quad \text{and} \quad \text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

Thus we conclude:

$$\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

■

5.5.2 Central Limit Theorem (CLT)

In the previous section, we saw that when we are sampling from a normal distribution, \bar{X} is also normally distributed. However, there are many situations where we cannot determine the exact

form of the distribution of X . In such circumstances, we may appeal to the central limit theorem and obtain an approximate distribution.

Central Limit Theorem: If \bar{X} is the mean of a random sample of size n taken from a population having the mean μ and the finite variance σ^2 , then \bar{X} approximately follows $\mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$ as $n \rightarrow \infty$.

In other words, the statistic

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

is a random variable whose distribution function approaches to that of the standard normal distributions as $n \rightarrow \infty$.

Proof: Let X_1, X_2, \dots, X_n be independent and identically distributed (i.i.d.) random variables with mean $\mu = \mathbb{E}(X_i)$ and variance $\sigma^2 = \text{Var}(X_i) < \infty$. Define the standardized sum:

$$Z_n = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)$$

Then, as $n \rightarrow \infty$, we have to prove²

$$Z_n \xrightarrow{d} \mathcal{N}(0, 1)$$

Now define

$$Y_i = \frac{X_i - \mu}{\sigma} \quad \text{so that} \quad \mathbb{E}[Y_i] = 0, \quad \text{Var}(Y_i) = 1$$

Then we can write:

$$Z_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i$$

Let $M_Y(t)$ be the moment generating function (MGF) of Y_i . Then, since Y_1, \dots, Y_n are i.i.d., the MGF of Z_n is:

$$M_{Z_n}(t) = \mathbb{E}(e^{tZ_n}) = \left(M_Y\left(\frac{t}{\sqrt{n}}\right) \right)^n$$

Using a Taylor expansion of $M_Y(t)$ around $t = 0$, we have:

$$M_Y(t) = 1 + \frac{t^2}{2} + \frac{\kappa_3 t^3}{6} + \dots$$

where κ_3 is the third central moment of Y_i .

Substituting t/\sqrt{n} into this expansion:

$$M_Y\left(\frac{t}{\sqrt{n}}\right) = 1 + \frac{t^2}{2n} + \frac{\kappa_3 t^3}{6n^{3/2}} + o\left(\frac{1}{n}\right)$$

Therefore,

$$M_{Z_n}(t) = \left(1 + \frac{t^2}{2n} + o\left(\frac{1}{n}\right) \right)^n \xrightarrow{n \rightarrow \infty} e^{t^2/2}$$

But $e^{t^2/2}$ is the MGF of the standard normal distribution $\mathcal{N}(0, 1)$. By the uniqueness theorem for MGFs, this implies:

$$Z_n \xrightarrow{d} \mathcal{N}(0, 1)$$

■

²The notation $Z_n \xrightarrow{d} Z$ (read as “converges in distribution”) means: the distribution of a sequence of random variables Z_n converges to the distribution of another random variable Z .

The Central Limit Theorem says that even if the population distribution is not normal, the sampling distribution of the sample mean will be approximately normal when the sample size is sufficiently large.

A common question is “how large does n have to be before the normality of \bar{X} is reasonable?” The answer depends on the degree of non-normality of the underlying distribution from which the sample has been drawn. The more non-normal the population distribution is, the larger n needs to be.

A useful **rule-of-thumb** is that n should be at least 30 for the central limit theorem to take effect.

If the number of observations n increases, the expected value of the sample mean \bar{X} remains fixed at μ , but its variance decreases, approaching zero. In other words,

$$\text{Var}(\bar{X}) = \mathbb{E}[(\bar{X} - \mu)^2] \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

This means the sample mean \bar{X} converges to μ in mean square, since the average squared deviation from the true mean diminishes with larger samples. This result represents one form of the **Law of Large Numbers (LLN)**, which formalizes the idea that

$$\bar{X} \rightarrow \mu \quad \text{as } n \rightarrow \infty.$$

Together with the Central Limit Theorem, the Law of Large Numbers assures us that the sample mean not only becomes approximately normally distributed but also increasingly concentrates around the true mean μ .

5.5.3 The Sampling Distribution of the Sample Mean When σ is Unknown

In the preceding subsections, we assume that the population variance σ^2 is known. If n is large, this does not pose any problems even when σ is unknown, as it is reasonable in that case to substitute for it the sample standard deviation S . However, for small sample sizes, the distribution of the sample mean \bar{X} is not known unless we assume that the sample comes from a normal population. Under this assumption, one can prove the following:

Theorem: Let \bar{X} be the sample mean of a random sample of size n drawn from a normal population with mean μ . Define the sample variance as

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Then the statistic (standardized sample mean)

$$t = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

follows a **t -distribution** with $\nu = n - 1$ degrees of freedom.



As illustrated in the figure above, the overall shape of the t -distribution closely resembles that of the standard normal distribution: both are bell-shaped and symmetric about the mean. Like the standard normal distribution, the t -distribution has a mean of 0. However, its variance depends on the parameter ν , known as the **degrees of freedom**. The variance of the t -distribution is greater than 1 but decreases as ν increases, approaching 1 in the limit.

The t -distribution with ν degrees of freedom converges to the standard normal distribution as $\nu \rightarrow \infty$. As a general rule of thumb, the standard normal distribution provides a good approximation to the t -distribution when the sample size is 30 or larger.



The critical point $t_{\alpha, \nu}$ is defined as the point so that the area to its right under the t -distribution with ν degrees of freedom equals α i.e.

$$P(t > t_{\alpha, \nu}) = \alpha$$

By symmetry,

$$t_{1-\alpha, \nu} = -t_{\alpha, \nu}$$

So the critical point for a left-tail area of α is $-t_{\alpha, \nu}$.

Example: Suppose we take a random sample of $n = 10$ measurements of battery lifespans (in hours) from a normally distributed population. The data are:

$$\{42, 38, 41, 39, 40, 37, 44, 36, 38, 40\}$$

We want to estimate the population mean μ and test whether the mean battery life is significantly different from 40 hours.

This is a case where the population standard deviation σ is unknown, so we use the sample standard deviation S , and apply the following test statistic:

$$t = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$$

where:

- \bar{X} is the sample mean,
- S is the sample standard deviation,
- $\mu_0 = 40$ is the hypothesized population mean,
- $n = 10$ is the sample size.

The statistic t follows a t -distribution with $\nu = n - 1 = 9$ degrees of freedom under the assumption that the population is normal.

Sample mean:

$$\bar{X} = \frac{1}{10}(42 + 38 + 41 + 39 + 40 + 37 + 44 + 36 + 38 + 40) = \frac{395}{10} = 39.5$$

Sample variance:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^{10} (X_i - \bar{X})^2 = \frac{1}{9} \sum_{i=1}^{10} (X_i - 39.5)^2 \approx 6.17$$

$$S = \sqrt{6.17} \approx 2.48$$

Test statistic:

$$t = \frac{\bar{X} - \mu}{S/\sqrt{n}} = \frac{39.5 - 40}{2.48/\sqrt{10}} \approx \frac{-0.5}{0.784} \approx -0.637$$

The computed t -value is approximately -0.637 , and it follows a t -distribution with $\nu = 9$ degrees of freedom. We can compare this value to critical values from the t -table or compute a p -value to make inference about μ .

5.6 The Sampling Distribution of the Sample Variance

When we take a random sample from a population, not only the sample mean but also the **sample variance**

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

behaves as a random variable. That is, the value of the sample variance S^2 will vary from one sample to another.

Theorem: If S^2 is the variance of a random sample of size n taken (with replacements) from a population of variance σ^2 , then

$$\mathbb{E}(S^2) = \sigma^2$$

Proof: Let X_1, X_2, \dots, X_n be a random sample from a population with mean μ and variance σ^2 . The sample variance is defined as:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

Now,

$$\begin{aligned}\sum_{i=1}^n (X_i - \bar{X})^2 &= \sum_{i=1}^n ((X_i - \mu) - (\bar{X} - \mu))^2 \\ &= \sum_{i=1}^n (X_i - \mu)^2 - n(\bar{X} - \mu)^2\end{aligned}$$

Taking expectations on both sides:

$$\mathbb{E} \left[\sum_{i=1}^n (X_i - \bar{X})^2 \right] = \mathbb{E} \left[\sum_{i=1}^n (X_i - \mu)^2 \right] - n \cdot \mathbb{E} [(\bar{X} - \mu)^2]$$

Now, observe:

$$\begin{aligned}\mathbb{E} \left[\sum_{i=1}^n (X_i - \mu)^2 \right] &= \sum_{i=1}^n \mathbb{E} (X_i - \mu)^2 = n\sigma^2 \\ \mathbb{E} [(\bar{X} - \mu)^2] &= \text{Var}(\bar{X}) = \frac{\sigma^2}{n}\end{aligned}$$

Therefore:

$$\mathbb{E} \left[\sum_{i=1}^n (X_i - \bar{X})^2 \right] = n\sigma^2 - n \cdot \frac{\sigma^2}{n} = n\sigma^2 - \sigma^2 = (n-1)\sigma^2$$

Dividing both sides by $n-1$, we get:

$$\mathbb{E}(S^2) = \frac{1}{n-1} \cdot (n-1)\sigma^2 = \sigma^2$$

■

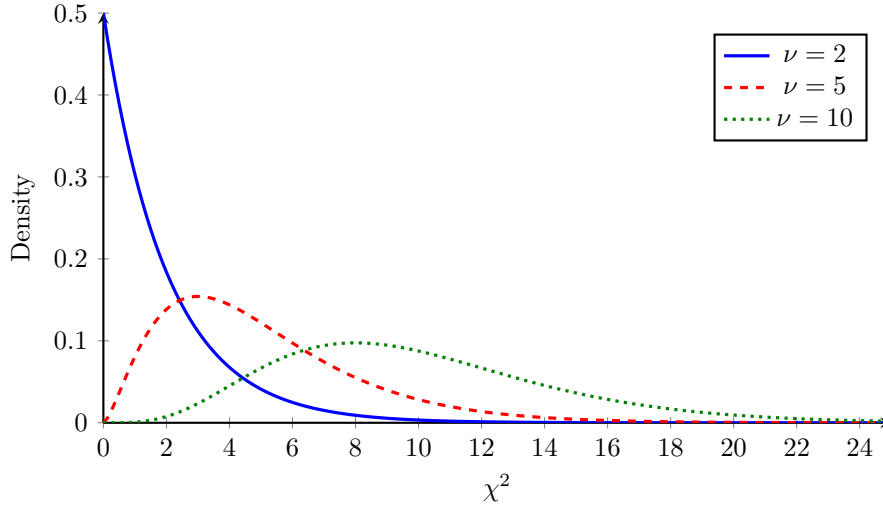
Thus the expectation value of the sample variance is the population variance. This result illustrates why the term $n-1$ is used in the denominator of the definition of sample variance, rather than n . But we still don't know the exact shape of the sampling distribution. To describe the exact sampling distribution of S^2 , we require the additional assumption that the population is normally distributed. Under this assumption, the following result holds:

Theorem: If S^2 is the variance of a random sample of size n taken from a normal population of variance σ^2 , then the statistic

$$\chi^2 = \frac{(n-1)S^2}{\sigma^2}$$

follows a **chi-squared distribution** with $\nu = n-1$ degrees of freedom.

This theorem tells us how the sample variance S^2 is distributed around the true population variance σ^2 .



A Chi-squared distribution has the following properties:

- The chi-squared distribution is not symmetric; it is skewed to the right, especially for small degrees of freedom.
- As the sample size increases ($n \rightarrow \infty$), the distribution becomes more symmetric and approaches normality.
- The expected value of the distribution is $\mathbb{E}[\chi^2] = n - 1$.
- The critical point $\chi_{\alpha, \nu}^2$ is defined as a point such that

$$P(\chi_\nu^2 > \chi_{\alpha, \nu}^2) = \alpha$$

where χ_ν^2 denotes a chi-square random variable with ν degrees of freedom.

5.7 Distribution of the Ratio of Two Sample Variances

A problem closely related to that of finding the distribution of the sample variance is that of determining the distribution of the ratio of the variances of two independent random samples. This problem is of considerable importance because it arises in hypothesis testing situations where we want to assess whether two samples come from populations with equal variances.

Theorem: Let S_1^2 and S_2^2 be the sample variances of two independent random samples of sizes n_1 and n_2 , respectively, drawn from two normal populations with equal variances σ_1^2 and σ_2^2 respectively. Then the statistic

$$F = \frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} = \frac{S_1^2/S_2^2}{\sigma_1^2/\sigma_2^2}$$

follows an F -distribution ($F \sim F_{\nu_1, \nu_2}$) with $\nu_1 = n_1 - 1$ and $\nu_2 = n_2 - 1$ degrees of freedom.

This theorem tells us how the ration of sample variances S_1^2/S_2^2 is distributed around the true population variance σ_1^2/σ_2^2 .

Properties of F -distribution:

- The F -distribution is characterized by two parameters:
 - ν_1 : **numerator degrees of freedom**,

– ν_2 : **denominator degrees of freedom.**

- Reciprocal property:

$$\text{If } X \sim F_{\nu_1, \nu_2}, \text{ then the reciprocal of } X \text{ i.e. } Y = \frac{1}{X} \sim F_{\nu_2, \nu_1}.$$

- The critical point $F_{\alpha; \nu_1, \nu_2}$ is defined as a point such that

$$P(X > F_{\alpha; \nu_1, \nu_2}) = \alpha$$

Equivalently³,

$$P\left(\frac{1}{X} < \frac{1}{F_{\alpha; \nu_1, \nu_2}}\right) = \alpha$$

But $\frac{1}{X} = Y \sim F_{\nu_2, \nu_1}$. Under the upper-tail critical point definition convention, if $F_{1-\alpha; \nu_2, \nu_1}$ is the critical value satisfying

$$P(Y > F_{1-\alpha; \nu_2, \nu_1}) = 1 - \alpha$$

then equivalently

$$P(Y \leq F_{1-\alpha; \nu_2, \nu_1}) = \alpha$$

Comparison with $P(Y < 1/F_{\alpha; \nu_1, \nu_2}) = \alpha$ shows

$$\frac{1}{F_{\alpha; \nu_1, \nu_2}} = F_{1-\alpha; \nu_2, \nu_1}$$

and hence

$$F_{\alpha; \nu_1, \nu_2} = \frac{1}{F_{1-\alpha; \nu_2, \nu_1}}$$

5.8 The Sampling Distribution of the Sample Proportion

Let X_1, X_2, \dots, X_n be a random sample from a Bernoulli population with proportion parameter p . Then the **sample proportion** is defined as:

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i = \frac{X}{n}$$

where X is the number of successes in the sample of size n . Since X is a binomial random variable with parameters n and p , i.e., $X \sim \text{Bin}(n, p)$, it follows that:

- The mean of \hat{p} is:

$$\mathbb{E}(\hat{p}) = \frac{1}{n} \mathbb{E}(X) = \frac{np}{n} = p$$

- The variance of \hat{p} is:

$$\text{Var}(\hat{p}) = \frac{1}{n^2} \text{Var}(X) = \frac{npq}{n^2} = \frac{pq}{n}, \quad \text{where } q = 1 - p$$

³This inverse relationship is maintained because $f(x) = \frac{1}{x}$ is a **strictly decreasing bijection** on $(0, \infty)$. Hence the event $\{A < B\}$ is exactly the same as the event $\{f(A) > f(B)\}$ or the event $\left\{\frac{1}{A} > \frac{1}{B}\right\}$.

Theorem: For large sample size n , the distribution of sample proportion \hat{p} can be approximated by a normal distribution according to the Central Limit Theorem:

$$\hat{p} \sim \mathcal{N}\left(p, \frac{pq}{n}\right)$$

This approximation is generally considered valid when both $np \geq 5$ and $nq \geq 5$.

The standardized form of \hat{p} is:

$$\frac{\hat{p} - p}{\sqrt{pq/n}}$$

which follows the standard normal distribution $\mathcal{N}(0, 1)$ for large sample.

Example: Suppose the true population proportion is $p = 0.6$ and a sample of size $n = 100$ is taken. Then:

$$\begin{aligned}\mathbb{E}(\hat{p}) &= 0.6 \\ \text{Var}(\hat{p}) &= \frac{0.6 \cdot 0.4}{100} = 0.0024\end{aligned}$$

Thus, the distribution of \hat{p} can be approximated as $N(0.6, 0.0024)$.

Chapter 6

Theory of Estimation

6.1 Introduction

Estimation is a fundamental component of **statistical inference**, which deals with drawing conclusions about population parameters from the analysis of sample data. There are two primary types of statistical inference:

1. **Estimation of parameters:** The true value of a population parameter is an unknown constant. The goal of estimation is to make informed guesses about this parameter using sample data, along with an assessment of the accuracy of these guesses.
2. **Hypothesis testing:** Sometimes, preliminary or tentative information about a population parameter is available. The objective of hypothesis testing is to use sample data to either support or reject such information about the parameter.

In this chapter, we focus on the first type—estimation of parameters. Hypothesis testing will be addressed in the next chapter.

Statistical estimation techniques are broadly classified into two categories:

1. **Point estimation:** A point estimation provides a single best guess of the unknown population parameter.
2. **Interval estimation:** An interval estimation gives a range of plausible values for the parameter, along with a specified level of confidence that the interval contains the true value.

Both point and interval estimation play crucial roles in quantifying uncertainty and guiding decision-making in the presence of incomplete information.

6.2 Point Estimation

Let θ be an unknown parameter (e.g. the population mean) associated with a particular variable. For estimating θ on the basis of random samples X_1, X_2, \dots, X_n , we may use a particular statistic T . This statistic T is called the **point estimator** of θ and the value of T obtained from a given sample is referred to as an **estimate** of θ .

Example: When we estimate the population mean $\theta = \mu$, the most intuitive estimator is the sample mean $\bar{X} = \frac{1}{N} \sum_{i=1}^n X_i$. Similarly sample variance (S^2) estimates population variance (σ^2) and sample proportion (\hat{p}) estimates population proportion (p)

6.2.1 Desirable Properties of a Good Estimator

There are often multiple point estimates available for any given parameter. So it is important to develop some evaluating criteria to judge the performance of each estimator and compare their performance. A good estimator should possess following desirable properties that make it reliable in estimating the true parameter value.

1. Unbiasedness:

An estimator T is said to be an **unbiased** estimator of θ if

$$\mathbb{E}(T) = \theta$$

Otherwise, T is said to be biased. The **bias** (\mathcal{B}) is given by

$$\mathcal{B} = \mathbb{E}(T) - \theta$$

Example: The sample mean \bar{X} and sample variance S^2 are unbiased estimator of the population mean μ and population variance σ^2 respectively, because

$$\mathbb{E}(\bar{X}) = \mu, \quad \mathbb{E}(S^2) = \sigma^2$$

The **mean-square-error** of the estimator T , denoted by $\text{MSE}(T)$ is defined as

$$\text{MSE}(T) = \mathbb{E}[(T - \theta)^2]$$

MSE measures, on average, how close an estimator comes to the true value of the parameter.

Theorem: Let T be an estimator of a population parameter θ . Then, the Mean Squared Error (MSE) of T is given by:

$$\text{MSE}(T) = \text{Var}(T) + \mathcal{B}^2(T)$$

Proof:

$$\begin{aligned} \text{MSE}(T) &= \mathbb{E}[(T - \theta)^2] \\ &= \mathbb{E}[(T - \mathbb{E}(T)) + (\mathbb{E}(T) - \theta)]^2 \\ &= \mathbb{E}[(T - \mathbb{E}(T))^2 + 2(T - \mathbb{E}(T))(\mathbb{E}(T) - \theta) + (\mathbb{E}(T) - \theta)^2] \\ &= \mathbb{E}[(T - \mathbb{E}(T))^2] + \underbrace{2(\mathbb{E}(T) - \mathbb{E}(T))\mathbb{E}(T - \mathbb{E}(T))}_{=0} + (\mathbb{E}(T) - \theta)^2 \\ &= \text{Var}(T) + \mathcal{B}^2(T) \end{aligned}$$

■

For an unbiased estimator $\mathcal{B} = 0$, and therefore $\text{MSE}(T) = \text{Var}(T)$.

In this context, the **standard error (SE)** of T is defined as the standard deviation of T i.e. $\sqrt{\text{Var}(T)}$ which is different from $\text{MSE}(T)$.

2. Consistency:

It is desirable that the estimator should behave more and more satisfactorily as the sample size n becomes larger. Consistency provides the criteria.

An estimator T_n (from a sample of size n) of a parameter θ is said to be **consistent** if, as the sample size n grows, T_n converges in probability to the true parameter value. Which means that for every $\varepsilon > 0$,

$$P(|T_n - \theta| > \varepsilon) \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

Consistency as defined above is sometimes called **weak consistency**. If we replace convergence in probability with almost sure convergence, i.e.

$$P\left(\lim_{n \rightarrow \infty} T_n = \theta\right) = 1 \quad \text{as } n \rightarrow \infty,$$

then the estimator is said to be **strongly consistent**¹.

Sufficient Conditions for Consistency: An estimator T_n of a parameter θ is said to be consistent if it satisfies the following two conditions:

(a) If T_n is an *asymptotically unbiased* estimator of θ i.e.

$$\mathbb{E}(T_n) \rightarrow \theta \quad \text{as } n \rightarrow \infty$$

(b) The variance of estimator T_n decreases with increasing sample size i.e.

$$\text{Var}(T_n) \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

Proof: By Chebyshev's inequality, if both conditions hold, then

$$\begin{aligned} P(|T_n - \theta| \geq \varepsilon) &\leq \frac{\mathbb{E}(T_n - \theta)^2}{\varepsilon^2}, \quad \text{for every } \varepsilon > 0 \\ &= \frac{1}{\varepsilon^2} \left(\mathbb{E}[(T_n - \mathbb{E}(T_n)) + (\mathbb{E}(T_n) - \theta)]^2 \right) \\ &= \frac{1}{\varepsilon^2} \left(\underbrace{\mathbb{E}(T_n - \mathbb{E}(T_n))^2}_{\text{Var}(T_n)} + (\mathbb{E}(T_n) - \theta)^2 \right) \\ &= \frac{1}{\varepsilon^2} \left(\text{Var}(T_n) + (\mathbb{E}(T_n) - \theta)^2 \right) \\ &\rightarrow 0 \quad \text{as } n \rightarrow \infty \end{aligned}$$

■

An estimator can be consistent even if it is biased for each finite n , provided the bias vanishes as $n \rightarrow \infty$.

Example: Let X_1, X_2, \dots, X_n be i.i.d. with mean μ and finite variance σ^2 . The sample mean

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

satisfies

$$\mathbb{E}[\bar{X}_n] = \mu, \quad \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n} \rightarrow 0.$$

Hence, by the two conditions above, \bar{X}_n is a consistent estimator of μ .

3. Efficiency:

¹**Weak consistency** says “in the long run, most of your estimates will be good,” but allows for occasional wildly bad estimates—even when n is very large.

Strong consistency rules out even those rare catastrophes: it guarantees that once you’ve accumulated enough data, your estimator will stay arbitrarily close to θ for every subsequent sample.

Among all unbiased estimators, the one with the smallest variance is said to be most **efficient**.

Unbiasedness is certainly a desirable property for point estimators but the criterion of unbiasedness does not generally provide a unique statistic for a given problem of estimation. For example, for symmetric population distribution, the sample median is also unbiased estimator for all sample sizes. Clearly, we need a further criterion to decide among different candidates.

One natural refinement is to compare their variances which measures the spread of the sampling distribution. Although both the sample mean and the sample median of a normal population are unbiased and have bell-shaped sampling distributions centered at μ , the variance of the sample mean is

$$\text{Var}(\bar{X}) = \frac{\sigma^2}{n},$$

whereas the variance of the sample median is approximately

$$\text{Var}(X_{\text{median}}) \approx 1.5708 \frac{\sigma^2}{n}.$$

Because the mean's distribution is more concentrated around μ , it will, on average, provide estimates closer to the truth. In other words, among unbiased estimators we favor the one with the smaller variance—and we call it, the most **efficient** estimator.

This leads to the concept of the *Minimum Variance Unbiased Estimator (MVUE)*:

The unbiased estimator T^* is a **Minimum Variance Unbiased Estimator (MVUE)** of a parameter if it has the smallest variance among all unbiased estimators of the parameter.

Formally, if \mathcal{U} is the class of all unbiased estimators of θ , then the MVUE T^* satisfies

$$\text{Var}(T^*) = \inf_{T \in \mathcal{U}} \text{Var}(T)$$

If two unbiased estimators T_1 and T_2 estimate the same parameter θ , the **relative efficiency** of T_1 with respect to T_2 is defined as:

$$\text{Relative Efficiency} = \frac{\text{Var}(T_2)}{\text{Var}(T_1)}.$$

An estimator is more efficient if it has a smaller variance. If the relative efficiency is close to 1, both estimators are equally good in terms of variance.

4. Sufficiency:

A statistic is said to be **sufficient** for a parameter if it captures all the information in the sample about that parameter.

Sufficiency is a key concept because it allows us to summarize the data without losing any relevant information about the parameter of interest.

Let $\mathbf{X} = (X_1, X_2, \dots, X_n)$ be a random sample from a distribution with conditional joint probability density function (or joint probability mass function)

$$f_{\mathbf{X}}(x_1, x_2, \dots, x_n; \theta),$$

where θ is an unknown parameter².

²Why we write the joint PDF in the form $f_{\mathbf{X}}(x_1, x_2, \dots, x_n; \theta)$ and not $f_{\mathbf{X}}(x_1, x_2, \dots, x_n | \theta)$?

Because in frequentist statistics, the parameter θ is treated as a fixed (but unknown) constant, while the data $\mathbf{X} = (X_1, X_2, \dots, X_n)$ are considered random variables. Therefore, we write the joint probability density (or mass)

A statistic $T(\mathbf{X})$ is said to be **sufficient** for θ if the conditional distribution of X_1, X_2, \dots, X_n given $T(\mathbf{X}) = t$ does not depend on θ . That is,

$$f_{\mathbf{X}}(x_1, x_2, \dots, x_n \mid T = t; \theta) = f_{\mathbf{X}}(x_1, x_2, \dots, x_n \mid T = t)$$

for all θ .

In other words, once you know $T = t$, the probability (or density) of seeing any particular arrangement of the raw observations does not depend on θ . After you condition on $T = t$, you look at the probability of different possible datasets that all share that same T -value. If that conditional probability still changes with θ , then those leftover items are carrying extra clues about θ .

A useful tool to verify sufficiency is the Neyman–Fisher Factorization Theorem, which states:

Neyman–Fisher Factorization Theorem: A necessary and sufficient condition for the statistic $T(\mathbf{X})$ to be a sufficient statistic for θ is that the joint PDF (or joint PMF) function of the sample can be factorized as:

$$f_{\mathbf{X}}(x_1, x_2, \dots, x_n; \theta) = g(T; \theta) \cdot h(x_1, x_2, \dots, x_n)$$

where

- $g(T; \theta)$ is a function that depends on the data (x_1, x_2, \dots, x_n) only through the function $T(x_1, x_2, \dots, x_n)$.
- $h(x_1, x_2, \dots, x_n)$ is a function of the data that does not depend on the parameter θ .

Example: Let X_1, X_2, \dots, X_n be independent Bernoulli random variables with common success probability p . Therefore

$$X_i = \begin{cases} 1, & \text{with probability } p \\ 0, & \text{with probability } 1 - p \end{cases}$$

We wish to show that the total proportion of successes in the sample,

$$T(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n X_i$$

is a sufficient statistic for p .

The joint probability mass function of the sample $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is

$$f_{\mathbf{X}}(x_1, x_2, \dots, x_n; p) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i}$$

Observe that this can be written in the form

$$f(x_1, x_2, \dots, x_n; p) = \underbrace{p^{nT(\mathbf{x})} (1-p)^{n-nT(\mathbf{x})}}_{g(T(\mathbf{x}), p)} \times \underbrace{1}_{h(\mathbf{x})}$$

where

$$T(\mathbf{x}) = \sum_{i=1}^n x_i$$

Here:

function as $f_{\mathbf{X}}(x_1, x_2, \dots, x_n; \theta)$ which emphasizes that this is a function of the data and a given fixed parameter.

On the other hand, the notation $f_{\mathbf{X}}(x_1, x_2, \dots, x_n \mid \theta)$ is typically reserved for conditional distributions, where θ is treated as a random variable—as in Bayesian statistics. In the frequentist context, θ is not random, so we avoid the conditional notation.

- $g(T(\mathbf{x}); p) = p^{T(\mathbf{x})}(1-p)^{n-T(\mathbf{x})}$ depends only on the statistic $T(\mathbf{x})$ and the parameter p .
- $h(x) = 1$ depends on the full data $\mathbf{x} = (x_1, x_2, \dots, x_n)$ but not on p .

Therefore, the statistic $T(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n X_i$ is sufficient for the parameter p .

Example: Let X_1, X_2, \dots, X_n be independent random variables drawn from a normal distribution with unknown mean μ and known variance σ^2 . That is,

$$X_i \sim \mathcal{N}(\mu, \sigma^2),$$

so each density is

$$f_{X_i}(x_i; \mu) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right)$$

We wish to show that the sample mean

$$T(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n X_i$$

is a sufficient statistic for μ .

The joint density of the sample $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is

$$\begin{aligned} f_{\mathbf{X}}(x_1, x_2, \dots, x_n; \mu) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right) \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right) \end{aligned}$$

Rewrite the exponential term:

$$\sum_{i=1}^n (x_i - \mu)^2 = \sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + n\mu^2 = \sum_{i=1}^n x_i^2 - n\mu^2$$

Thus

$$\exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right) = \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n x_i^2\right) \times \exp\left(\frac{\mu}{\sigma^2} \sum_{i=1}^n x_i - \frac{n\mu^2}{2\sigma^2}\right)$$

Hence the joint density factors as

$$f_{\mathbf{X}}(x_1, x_2, \dots, x_n; \mu) = \underbrace{\exp\left(\frac{n\mu}{\sigma^2} T(\mathbf{x}) - \frac{n\mu^2}{2\sigma^2}\right)}_{g(T(\mathbf{x}), \mu)} \times \underbrace{(2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n x_i^2\right)}_{h(\mathbf{x})}$$

where $T(\mathbf{x}) = \sum_{i=1}^n x_i$. Therefore, by the Neyman–Fisher factorization theorem, because the statistic $T(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n X_i$ is sufficient for μ .

Example: Let X_1 and X_2 be two independent and identically distributed random variables from the Poisson distribution with parameter $\lambda > 0$, i.e.,

$$X_1, X_2 \stackrel{iid}{\sim} \text{Poisson}(\lambda)$$

The probability mass function of a Poisson random variable is given by

$$P(X_i = x_i; \lambda) = \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}, \quad x_i = 0, 1, 2, \dots$$

Since X_1 and X_2 are independent, the joint PMF of the sample is

$$p_{X_1, X_2}(x_1, x_2; \lambda) = \frac{e^{-\lambda} \lambda^{x_1}}{x_1!} \cdot \frac{e^{-\lambda} \lambda^{x_2}}{x_2!} = \frac{e^{-2\lambda} \lambda^{x_1+x_2}}{x_1! x_2!}$$

Now consider the statistic

$$T = X_1 + 2X_2$$

To apply the Neyman–Fisher Factorization Theorem, we attempt to factor the joint PMF in the form

$$p_{X_1, X_2}(x_1, x_2; \lambda) = g(T(x_1, x_2), \lambda) \cdot h(x_1, x_2)$$

i.e., express the λ -dependence entirely through the statistic T .

However, in our case the joint PMF is

$$p_{X_1, X_2}(x_1, x_2; \lambda) = \frac{e^{-2\lambda} \lambda^{x_1+x_2}}{x_1! x_2!} = \frac{e^{-2\lambda} \lambda^{T(x_1, x_2) - x_2}}{x_1! x_2!}$$

where the λ -dependent part is $\lambda^{T(x_1, x_2) - x_2}$, not only a function of $T(x_1, x_2)$ but also a function of x_2 . Therefore, the statistic $T = X_1 + 2X_2$ is not sufficient.

6.3 Maximum Likelihood Estimation (MLE)

Maximum Likelihood Estimation (MLE) is one of the most widely used methods for estimating the parameters of a statistical model. The basic idea is to choose the parameter values that make the observed data most probable.

Let $\mathbf{X} = (X_1, X_2, \dots, X_n)$ be a random sample from a population with joint probability density function (or probability mass function) $f_{\mathbf{X}}(x_1, x_2, \dots, x_n; \theta)$, where θ is the unknown parameter to be estimated. Given $\mathbf{x} = (x_1, x_2, \dots, x_n)$, it may be looked upon as a function of θ , called the **likelihood function** of θ and is denoted by $L(\theta)$.

$$L(\theta) = f_{\mathbf{X}}(x_1, x_2, \dots, x_n; \theta)$$

The value of θ that maximizes this function is called the **maximum likelihood estimator (MLE)** of θ :

$$\hat{\theta} = \arg \max_{\theta} L(\theta)$$

In practice, it is often more convenient to work with the **log-likelihood function**:

$$\log L(\theta) = \log f_{\mathbf{X}}(x_1, x_2, \dots, x_n; \theta)$$

Maximizing the log-likelihood yields the same estimator as maximizing the likelihood.

Example: Poisson distribution

Suppose X_1, X_2, \dots, X_n are i.i.d. from a Poisson distribution with parameter $\lambda > 0$. The PMF for the Poisson distribution is:

$$f(x; \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}, \quad x = 0, 1, 2, \dots$$

The likelihood function is:

$$\begin{aligned} L(\lambda) &= f_{\mathbf{X}}(x_1, x_2, \dots, x_n; \lambda) \\ &= \prod_{i=1}^n f(x_i; \lambda) = \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{x_i!} \\ &= e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i} \prod_{i=1}^n \frac{1}{x_i!} \end{aligned}$$

The log-likelihood is:

$$\log L(\lambda) = -n\lambda + \left(\sum_i x_i \right) \log \lambda - \sum_i \log(x_i!)$$

Differentiating and setting the derivative to zero:

$$\frac{d}{d\lambda} \log L(\lambda) = -n + \frac{1}{\lambda} \sum_i x_i = 0 \Rightarrow \hat{\lambda} = \frac{1}{n} \sum_i x_i = \bar{x}$$

Hence, the MLE of λ is $\hat{\lambda} = \bar{X}$.

Example: Normal distribution

Assume $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$. The PDF is:

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

Likelihood function:

$$L(\mu, \sigma^2) = \prod_{i=1}^n f(x_i; \mu, \sigma) = (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right)$$

Log-likelihood function:

$$\log L(\mu, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

(i) **Case 1: μ unknown, σ known ($= \sigma_0$)**

Log-likelihood function:

$$\log L(\mu) = -\frac{n}{2} \log(2\pi\sigma_0^2) - \frac{1}{2\sigma_0^2} \sum_{i=1}^n (x_i - \mu)^2$$

Taking derivative and setting to zero:

$$\frac{d}{d\mu} \log L(\mu) = \frac{1}{\sigma_0^2} \sum_{i=1}^n (x_i - \mu) = 0 \Rightarrow \hat{\mu} = \bar{x}$$

Hence, the MLE of μ is $\hat{\mu} = \bar{X}$.

(ii) **Case 2: μ known ($= \mu_0$), σ unknown**

Log-likelihood:

$$\log L(\sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_i (x_i - \mu_0)^2$$

Taking derivative w.r.t. σ^2 and setting to zero:

$$\begin{aligned} \frac{d}{d\sigma^2} \log L(\sigma^2) &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_i (x_i - \mu_0)^2 = 0 \\ \Rightarrow \hat{\sigma}^2 &= \frac{1}{n} \sum_i (x_i - \mu_0)^2 \end{aligned}$$

Hence, the MLE of σ^2 is $\hat{\sigma}^2 = \frac{1}{n} \sum_i (x_i - \mu_0)^2$.

(iii) **Case 3: μ and σ both unknown**

Log-likelihood:

$$\log L(\mu, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2$$

Taking partial derivatives and solving the system:

$$\begin{aligned} \frac{\partial}{\partial \mu} \log L(\mu, \sigma^2) &= 0 \Rightarrow \hat{\mu} = \bar{x} \\ \frac{\partial}{\partial \sigma^2} \log L(\mu, \sigma^2) &= 0 \Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_i (x_i - \bar{x})^2 \end{aligned}$$

Thus, the MLEs for μ and σ^2 are the sample mean \bar{X} and sample variance $\frac{1}{n} \sum_i (x_i - \bar{x})^2$ (without Bessel's correction), respectively.

It is important to note that the maximum likelihood estimator (MLE) of the population variance σ^2 is **not an unbiased estimator**. The MLE is given by

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2,$$

where \bar{X} is the sample mean. But we have already seen that the unbiased estimator of σ^2 is:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{n-1}{n} \hat{\sigma}^2$$

The MLE estimator tends to underestimate the true variance.

6.4 Bayesian Estimation

Bayesian estimation is a method of statistical inference in which the unknown parameter θ is modeled as a random variable Θ with a probability distribution $\pi_{\Theta}(\theta)$, known as the **prior distribution**. It is intended to reflect our knowledge of the parameter θ , before we gather data.

Let $\mathbf{X} = (X_1, X_2, \dots, X_n)$ be random variables representing the observations in the sample data with joint PDF (PMF) given $\Theta = \theta$ is given by $f_{\mathbf{X}}(\mathbf{x} | \theta)$ which is also known as the **likelihood**. When data $\mathbf{X} = \mathbf{x}$ are observed, the extra information about θ is combined with the prior distribution to obtain the **posterior distribution** $\pi_{\Theta}(\theta | \mathbf{x})$ for given $\mathbf{X} = \mathbf{x}$ using Bayes theorem as follows:

$$\pi_{\Theta}(\theta | \mathbf{X}) = \frac{f_{\mathbf{X}}(\mathbf{x} | \theta) \pi_{\Theta}(\theta)}{f_{\mathbf{X}}(\mathbf{x})}$$

where,

$$f_{\mathbf{X}}(\mathbf{x}) = \begin{cases} \sum f_{\mathbf{X}}(\mathbf{x} | \theta) \pi_{\Theta}(\theta), & \text{in the discrete case,} \\ \int_{-\infty}^{\infty} f_{\mathbf{X}}(\mathbf{x} | \theta) \pi_{\Theta}(\theta) d\theta, & \text{in the continuous case.} \end{cases}$$

Thus,

$$\underbrace{\pi_{\Theta}(\theta | \mathbf{x})}_{\text{posterior}} \propto \underbrace{f_{\mathbf{X}}(\mathbf{x} | \theta)}_{\text{likelihood}} \times \underbrace{\pi_{\Theta}(\theta)}_{\text{prior}}$$

In practice the constant of proportionality chosen in such a way that it makes the total mass of the posterior distribution equal to one. The posterior distribution reflects our updated belief about the parameter after seeing the data.

Example: Suppose you have three coins in your pocket:

- Coin 1: Biased in favour of tails with head probability $\theta = 0.25$
- Coin 2: A fair coin with $\theta = 0.5$
- Coin 3: Biased in favour of heads with $\theta = 0.75$

You randomly select one coin and flip it once. You observe a head. What is the posterior probability that you chose Coin 3?

This is a classic example of Bayesian inference with a discrete parameter space. The **population** consists of three types of coins, each with a different probability of producing a head: $\theta \in \{0.25, 0.5, 0.75\}$.

We assume one of these coins is selected at random. The **sample** is a single coin toss from the selected coin, which results in observing a head. Using this one data point, we update our belief (prior distribution) over the possible values of θ to obtain a posterior distribution.

Let $X = 1$ denote the event that you observe a head, and $X = 0$ for a tail.

Let θ denote the probability of heads. Then $\theta \in \{0.25, 0.5, 0.75\}$.

The **prior probabilities** are:

$$P(\theta = 0.25) = P(\theta = 0.5) = P(\theta = 0.75) = \frac{1}{3}$$

Because the probability of selecting any coin at random is same before we have the knowledge of the sample observation.

The **likelihood** is given by the Bernoulli probability mass function:

$$P(X = x | \theta) = \theta^x (1 - \theta)^{1-x}$$

Since we observed $X = 1$, the likelihood becomes $P(X = 1 | \theta) = \theta$.

We now calculate the unnormalized and normalized **posterior probabilities** using Bayes' Theorem:

$$\underbrace{P(\theta | X = 1)}_{\text{posterior}} \propto \underbrace{P(X = 1 | \theta)}_{\text{likelihood}} \times \underbrace{P(\theta)}_{\text{prior}}$$

Coin	θ	Prior	Likelihood	Unnorm. Posterior	Norm. Posterior
		$P(\theta)$	$P(X = 1 \theta)$	$P(\theta X = 1)$	$\frac{P(\theta X = 1)}{\sum_{\theta} P(\theta X = 1)}$
1	0.25	0.33	0.25	0.0825	0.167
2	0.50	0.33	0.50	0.1650	0.333
3	0.75	0.33	0.75	0.2475	0.500
Sum		1.00		0.495	1.000

Table 6.1: Table for calculating the posterior probabilities.

The posterior probability that the coin chosen was Coin 3, given that a head was observed, is:

$$P(\theta = 0.75 | X = 1) = 0.5$$

This illustrates how Bayesian estimation updates our belief about which coin was selected based on the observed outcome.

6.4.1 Bayesian Approach to Point Estimation

When you have your posterior density $\pi_{\Theta}(\theta \mid \mathbf{X})$, you still need a single “best-guess” $\hat{\theta}$. Bayesian decision theory tells us that the choice of $\hat{\theta}$ depends on a loss function $L(\theta, \hat{\theta})$, which quantifies how “bad” it is to decide $\hat{\theta}$ when the true parameter is θ . When our estimate is $\hat{\theta}$, the **expected posterior loss** is

$$h(\hat{\theta}) = \int_{\theta} L(\theta, \hat{\theta}) \pi_{\Theta}(\theta \mid \mathbf{X}) d\theta$$

The Bayes estimator $\hat{\theta}$ minimises the expected posterior loss i.e.

$$\begin{aligned} \hat{\theta} &= \arg \min_{\hat{\theta}} h(\hat{\theta}) \\ &= \arg \min_{\hat{\theta}} \int_{\theta} L(\theta, \hat{\theta}) \pi_{\Theta}(\theta \mid \mathbf{X}) d\theta \end{aligned}$$

The form of the minimiser depends on the choice of L . Some common cases are:

1. **Squared-error loss:**

$$L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$$

The posterior expected loss

$$\begin{aligned} h(\hat{\theta}) &= \int (\theta - \hat{\theta})^2 \pi_{\Theta}(\theta \mid \mathbf{X}) d\theta \\ &= \int (\theta^2 - 2\hat{\theta} \cdot \theta + \hat{\theta}^2) \pi_{\Theta}(\theta \mid \mathbf{X}) d\theta \\ &= \int \theta^2 \pi_{\Theta}(\theta \mid \mathbf{X}) d\theta - 2\hat{\theta} \int \theta \pi_{\Theta}(\theta \mid \mathbf{X}) d\theta + \hat{\theta}^2 \int \pi_{\Theta}(\theta \mid \mathbf{X}) d\theta \\ &= \mathbb{E}[\Theta^2 \mid \mathbf{X}] - 2\hat{\theta} \mathbb{E}[\Theta \mid \mathbf{X}] + \hat{\theta}^2, \end{aligned}$$

using $\int \pi(\theta \mid \mathbf{X}) d\theta = 1$

To find the minimiser, differentiate $h(\hat{\theta})$ with respect to $\hat{\theta}$:

$$\frac{dh(\hat{\theta})}{d\hat{\theta}} = -2\mathbb{E}[\Theta \mid \mathbf{X}] + 2\hat{\theta}.$$

Setting this derivative to zero gives

$$-2\mathbb{E}[\Theta \mid \mathbf{X}] + 2\hat{\theta} = 0 \implies \hat{\theta} = \mathbb{E}[\Theta \mid \mathbf{X}]$$

Finally, check the second derivative:

$$\frac{d^2 h(\hat{\theta})}{d\hat{\theta}^2} = 2 > 0,$$

so this critical point is indeed a minimum.

Hence, the Bayes estimator $\hat{\theta}$ under absolute-error loss is the **mean of the posterior distribution**:

$$\hat{\theta} = \mathbb{E}[\Theta \mid \mathbf{X}]$$

2. **Absolute-error loss:**

$$L(\theta, \hat{\theta}) = |\theta - \hat{\theta}|$$

The posterior expected loss

$$h(\hat{\theta}) = \int_{\Theta} |\theta - \hat{\theta}| \pi_{\Theta}(\theta \mid \mathbf{X}) d\theta$$

Split the integral at $\hat{\theta}$:

$$h(\hat{\theta}) = \int_{-\infty}^{\hat{\theta}} (\hat{\theta} - \theta) \pi_{\Theta}(\theta | \mathbf{X}) d\theta + \int_{\hat{\theta}}^{\infty} (\theta - \hat{\theta}) \pi_{\Theta}(\theta | \mathbf{X}) d\theta$$

Differentiate with respect to $\hat{\theta}$. By Leibniz's rule³:

$$\frac{dh(\hat{\theta})}{d\hat{\theta}} = \int_{-\infty}^{\hat{\theta}} \pi_{\Theta}(\theta | \mathbf{X}) d\theta - \int_{\hat{\theta}}^{\infty} \pi_{\Theta}(\theta | \mathbf{X}) d\theta = \Pi_{\Theta}(\hat{\theta}) - (1 - \Pi_{\Theta}(\hat{\theta})),$$

where $\Pi_{\Theta}(t) = \int_{-\infty}^t \pi_{\Theta}(\theta | \mathbf{X}) d\theta$ is the posterior CDF. Setting this derivative to zero:

$$\Pi_{\Theta}(\hat{\theta}) - (1 - \Pi_{\Theta}(\hat{\theta})) = 0 \implies \Pi_{\Theta}(\hat{\theta}) = \frac{1}{2}$$

Hence the Bayes estimator $\hat{\theta}$ under absolute-error loss is the **median (midpoint of the CDF) of the posterior distribution**, satisfying

$$\int_{-\infty}^{\hat{\theta}} \pi_{\Theta}(\theta | \mathbf{X}) d\theta = \frac{1}{2}$$

Finally, the second derivative is

$$\frac{d^2 h(\hat{\theta})}{d\hat{\theta}^2} = 2 \pi_{\Theta}(\hat{\theta} | \mathbf{X}) \geq 0,$$

so $h(\hat{\theta})$ is convex at the solution, confirming that $\hat{\theta}$ indeed minimizes the expected loss.

3. Zero-one loss:

$$L(\theta, \hat{\theta}) = \begin{cases} 0, & \theta = \hat{\theta}, \\ 1, & \theta \neq \hat{\theta}. \end{cases}$$

The expected posterior loss:

$$h(\hat{\theta}) = \int_{\Theta} L(\theta, \hat{\theta}) \pi_{\Theta}(\theta | \mathbf{X}) d\theta.$$

Since $L(\theta, \hat{\theta}) = 0$ only at $\theta = \hat{\theta}$ and equals 1 elsewhere, we can simplify the expected posterior loss:

$$\begin{aligned} h(\hat{\theta}) &= \int_{\Theta} L(\theta, \hat{\theta}) \pi_{\Theta}(\theta | \mathbf{X}) d\theta \\ &= \int_{\theta \neq \hat{\theta}} \pi_{\Theta}(\theta | \mathbf{X}) d\theta \\ &= 1 - \pi_{\Theta}(\hat{\theta} | \mathbf{X}) \end{aligned}$$

Therefore,

$$\hat{\theta} = \arg \min_{\hat{\theta}} (1 - \pi_{\Theta}(\hat{\theta} | \mathbf{X})) = \arg \max_{\hat{\theta}} \pi_{\Theta}(\hat{\theta} | \mathbf{X})$$

Hence, the Bayes estimator under zero-one loss is the value of θ that maximizes the posterior density, i.e., the **maximum a posteriori (MAP) estimator**. In other words, it is the **mode of the posterior distribution**.

³The **Leibniz integral rule** (differentiation under the integral sign) states that if

$$H(t) = \int_{a(t)}^{b(t)} g(x, t) dx,$$

then

$$\frac{d}{dt} \int_{a(t)}^{b(t)} g(x, t) dx = g(b(t), t) b'(t) - g(a(t), t) a'(t) + \int_{a(t)}^{b(t)} \frac{\partial}{\partial t} g(x, t) dx.$$

Example: Suppose we have:

- Observations: $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, 1)$,
- Prior distribution: $\mu \sim \mathcal{N}(0, \tau^{-2})$

Likelihood:

$$f_{\mathbf{X}}(\mathbf{x} \mid \mu) = \frac{1}{\sqrt{2\pi}} \exp \left(- \sum_i \frac{(x_i - \mu)^2}{2} \right)$$

Prior distribution:

$$\pi_M(\mu) = \frac{1}{\tau\sqrt{2\pi}} \exp \left(- \frac{\mu^2 \tau^2}{2} \right)$$

Then the posterior distribution is given by

$$\begin{aligned} \pi_M(\mu \mid \mathbf{x}) &\propto f_{\mathbf{X}}(\mathbf{x} \mid \mu) \cdot \pi_M(\mu) \\ &\propto \exp \left(- \frac{1}{2} \sum_i (x_i - \mu)^2 \right) \cdot \exp \left(- \frac{\mu^2 \tau^2}{2} \right) \end{aligned}$$

Now expand the squared term in the sum:

$$\sum_i (x_i - \mu)^2 = \sum_i x_i^2 - 2\mu \sum_i x_i + n\mu^2$$

Ignoring terms not involving μ , we write:

$$\begin{aligned} \pi_M(\mu \mid \mathbf{x}) &\propto \exp \left(\mu \sum_i x_i - \frac{n}{2} \mu^2 \right) \cdot \exp \left(- \frac{1}{2} \mu^2 \tau^2 \right) \\ &\propto \exp \left[\mu \sum_i x_i - \frac{1}{2} (n + \tau^2) \mu^2 \right] \end{aligned}$$

Now,

$$\begin{aligned} \mu \sum_i x_i - \frac{1}{2} (n + \tau^2) \mu^2 &= - \frac{1}{2} (n + \tau^2) \left[\mu^2 - \frac{2 \sum_i x_i}{n + \tau^2} \mu \right] \\ &= - \frac{1}{2} (n + \tau^2) \left[\left(\mu - \frac{\sum_i x_i}{n + \tau^2} \right)^2 - \left(\frac{\sum_i x_i}{n + \tau^2} \right)^2 \right] \end{aligned}$$

Drop the constant term (independent of μ):

$$\pi_M(\mu \mid \mathbf{x}) \propto \exp \left[- \frac{1}{2} (n + \tau^2) \left(\mu - \frac{\sum_i x_i}{n + \tau^2} \right)^2 \right]$$

So the posterior distribution of μ given data \mathbf{x} is a Normal distribution with mean and variance given by:

$$\text{Mean} = \frac{\sum_i x_i}{n + \tau^2} = \frac{n\bar{x}}{n + \tau^2}, \quad \text{Variance} = \frac{1}{n + \tau^2}$$

The normal density is symmetric, and so the posterior mean, median and mode have the same value. Thus the optimal Bayes estimate of μ under squared, absolute and zero-one error loss is given by

$$\hat{\theta} = \frac{n\bar{X}}{n + \tau^2}$$

6.5 Estimation using Method of Moments

The **method of moments** is a classical technique used to estimate unknown parameters of a probability distribution using sample data. The core idea is simple: it equates the theoretical moments of a distribution (which depend on the parameters) to the corresponding sample moments computed from the data.

Let X_1, X_2, \dots, X_n be a random sample from a population with a distribution that depends on one or more parameters $\theta_1, \theta_2, \dots, \theta_k$.

- The **r-th population moment** about the origin is defined as:

$$\mu'_r = \mathbb{E}(X^r)$$

which is a function of the unknown parameters $\theta_1, \dots, \theta_k$.

- The **r-th sample moment** about the origin is defined as:

$$m'_r = \frac{1}{n} \sum_{i=1}^n X_i^r$$

which is computable from observed data.

The **method of moments estimator** $\hat{\theta}$ is obtained by equating the sample moment to the corresponding population moment. Let's equate the first sample moment to the first population moment:

$$m'_1 = \mu(\theta)$$

Solving this equation for θ yields the estimator:

$$\hat{\theta} = \mu^{-1}(m'_1)$$

assuming $\mu(\theta)$ is invertible.

This approach can be extended to multiple parameters. If the distribution depends on multiple parameters $\theta_1, \dots, \theta_k$, then the first k theoretical moments are equated to the first k sample moments:

$$\begin{aligned} m'_1 &= \mu'_1(\theta_1, \theta_2, \dots, \theta_k), \\ m'_2 &= \mu'_2(\theta_1, \theta_2, \dots, \theta_k), \\ &\vdots \\ m'_k &= \mu'_k(\theta_1, \theta_2, \dots, \theta_k) \end{aligned}$$

Solving these k equations yields the method of moments estimators $\hat{\theta}_1, \theta_2, \dots, \hat{\theta}_k$.

Example: Exponential Distribution

Suppose X_1, X_2, \dots, X_n is a sample from an exponential distribution with parameter $\lambda > 0$, having density

$$f_X(x; \lambda) = \lambda e^{-\lambda x}, \quad x \geq 0$$

The first population moment (mean) is:

$$\mu'_1 = \mathbb{E}(X) = \frac{1}{\lambda}$$

The first sample moment is:

$$m'_1 = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}$$

Equating the moments:

$$\bar{X} = \frac{1}{\lambda} \Rightarrow \hat{\lambda} = \frac{1}{\bar{X}}$$

Hence, the method of moments estimator for λ is:

$$\hat{\lambda} = \frac{1}{\bar{X}}$$

6.6 Interval Estimation

A **point estimate**, being a single number, gives no indication of how accurate or reliable it is. For example, suppose a sample of batteries yields an average lifetime of $\bar{X} = 218.2$ hours. While this provides a best guess for the true average lifetime μ , it says nothing about how close it is to μ . Due to sampling variability, the point estimate is almost never exactly equal to the true value. A more informative approach is to report an **interval estimate**, which provides a range of plausible values for μ and reflects the uncertainty in the estimation process.

Basically, the purpose of interval estimation for a population parameter θ is to find two values L and R from the random sample such that

$$L \leq \theta \leq R$$

with some specific probability. Information about the precision of an interval estimate is conveyed by the width of the interval $R - L$. Because L and R depend on sample values, they will be random. The interval $[L, R]$ should have the following properties:

1. The probability that θ lies within $[L, R]$, i.e., $P(L \leq \theta \leq R)$, should be high.
2. The length of the interval, $R - L$, should be as short as possible to ensure precision.

In addition to providing the interval $[L, R]$, we also specify a measure of confidence in the accuracy of the estimate. This leads to the concept of a **confidence interval (CI)**.

- The **confidence interval** is the interval estimate $[L, R]$ of the parameter θ .
- The **confidence level** is the probability that the confidence interval contains the true value of θ .
- The endpoints L and R are called the **lower** and **upper confidence limits**, respectively.

Definition:

- A **confidence interval** for a parameter is a range of values, computed from sample data, within which the true parameter is believed to lie.
- The **confidence level** is the probability that the confidence interval contains the true parameter value. It is typically chosen close to 1, such as 0.95 or 0.99.

We can write for the interval estimate of θ

$$P(L \leq \theta \leq R) = 1 - \alpha$$

We read this as we are $100(1 - \alpha)\%$ confident that θ lies within the interval $[L, R]$. The interval is called the $100(1 - \alpha)\%$ **confidence interval** or simply the $100(1 - \alpha)\%$ **CI**.

- For a 95% CI, $\alpha = 0.05$,
- For a 99% CI, $\alpha = 0.01$.

A 95% confidence interval means that if we repeated the sampling procedure many times, approximately 95% of the calculated intervals would contain the true parameter value.

6.6.1 Pivotal Method

One of the most widely used methods for constructing confidence intervals is the **pivotal quantity method**, also known as the **pivotal method**. This approach is particularly useful when we can identify a function of the sample data and the parameter, called a *pivotal quantity*, whose distribution does not depend on the unknown parameter.

A **pivotal quantity** is a function $T(X_1, X_2, \dots, X_n; \theta)$ of the sample data and the parameter θ , such that the probability distribution of T is independent of θ .

The steps to construct a confidence interval using the pivotal method are as follows:

1. Identify a suitable pivotal quantity $T(X_1, \dots, X_n; \theta)$ whose distribution is known and does not depend on θ .
2. Find constants a and b such that

$$P(a \leq T(X_1, \dots, X_n; \theta) \leq b) = 1 - \alpha$$

where $1 - \alpha$ is the desired confidence level. The constants a and b are called **critical values**.

3. Solve the inequality to find the interval

$$a \leq T(X_1, \dots, X_n; \theta) \leq b$$

for θ in terms of the sample data.

4. The resulting interval gives the $100(1 - \alpha)\%$ confidence interval for θ .

6.7 Confidence Interval for the Mean in a Normal Population

6.7.1 Population Variance Known

Suppose X_1, X_2, \dots, X_n is a random sample from a normal distribution $N(\mu, \sigma^2)$, where the variance σ^2 is known. The sample mean \bar{X} follows

$$\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

The pivotal quantity

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

has a standard normal distribution, $Z \sim \mathcal{N}(0, 1)$, independent of the unknown parameter μ .

To find the confidence interval, we have to find critical values a and b such that

$$P(a \leq Z \leq b) = 1 - \alpha$$

Since the distribution is symmetric about zero, we can choose the critical values $a = -q, b = q$. For standard normal distribution, $q = z_{\alpha/2}$ represents the value of z with tail area $\alpha/2$.



Thus,

$$P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) = 1 - \alpha$$

which translates to

$$P\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

Therefore, the $100(1 - \alpha)\%$ confidence interval for μ is

$$\left[\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \quad \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$$

- For a 95% confidence level, $1 - \alpha = 0.95$. From table $z_{\alpha/2} = z_{0.025} = 1.96$.
- For a 99% confidence level, $1 - \alpha = 0.99$. From table $z_{\alpha/2} = z_{0.005} = 2.576$.

Example: Suppose the lifetimes of a certain type of battery are normally distributed with unknown mean μ and known standard deviation $\sigma = 10$ hours. A random sample of $n = 25$ batteries yields a sample mean of $\bar{X} = 100$ hours. Construct a 95% confidence interval for the population mean lifetime.

- Given: $\sigma = 10$, $n = 25$, $\bar{X} = 100$, and confidence level = 95%
- For 95% confidence, $z_{\alpha/2} = z_{0.025} = 1.96$

Compute the standard error:

$$\frac{\sigma}{\sqrt{n}} = \frac{10}{\sqrt{25}} = \frac{10}{5} = 2$$

Compute the margin of error:

$$z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} = 1.96 \times 2 = 3.92$$

Construct the confidence interval:

$$[100 - 3.92, \quad 100 + 3.92] = [96.08, 103.92]$$

We are 95% confident that the true mean lifetime μ of the batteries lies between 96.08 and 103.92 hours.

6.7.2 Population Variance Unknown

Let X_1, X_2, \dots, X_n be a random sample from a normal distribution $N(\mu, \sigma^2)$, where both μ and σ^2 are unknown. Let \bar{X} be the sample mean and S^2 be sample variance.

The pivotal quantity

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

follows a Student's t -distribution with $\nu = n - 1$ degrees of freedom:

$$T \sim t_{n-1}$$

Let $t_{\alpha/2, n-1}$ be the critical value of the t -distribution such that $P(t > t_{\alpha/2, n-1}) = \alpha/2$. Hence, we can write

$$P(-t_{\alpha/2, n-1} \leq T \leq t_{\alpha/2, n-1}) = 1 - \alpha$$

which implies

$$P\left(\bar{X} - t_{\alpha/2, n-1} \frac{S}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{\alpha/2, n-1} \frac{S}{\sqrt{n}}\right) = 1 - \alpha$$

Therefore, the $100(1 - \alpha)\%$ confidence interval for μ is

$$\left[\bar{X} - t_{\alpha/2, n-1} \frac{S}{\sqrt{n}}, \quad \bar{X} + t_{\alpha/2, n-1} \frac{S}{\sqrt{n}} \right]$$

6.8 Confidence Interval for the Variance in a Normal Population

Consider a random sample X_1, X_2, \dots, X_n from a normal distribution $N(\mu, \sigma^2)$ with unknown variance σ^2 . Let S^2 be sample variance.

The pivotal quantity

$$\chi^2 = \frac{(n-1)S^2}{\sigma^2}$$

follows a chi-square distribution with $\nu = n - 1$ degrees of freedom:

$$\chi^2 \sim \chi_{n-1}^2$$

We now need to find the critical values L and U such that

$$P\left(L \leq \frac{(n-1)S^2}{\sigma^2} \leq U\right) = 1 - \alpha$$



Figure 6.1: Confidence interval for χ^2 distribution.

We take areas to the left of L and to the right of R to be equal to $\alpha/2$ i.e.

$$L = \chi^2_{1-\alpha/2, n-1}, \quad R = \chi^2_{\alpha/2, n-1}$$

Thus,

$$P\left(\chi^2_{1-\alpha/2, n-1} \leq \frac{(n-1)S^2}{\sigma^2} \leq \chi^2_{\alpha/2, n-1}\right) = 1 - \alpha$$

Rearranging to solve for σ^2 , we get

$$P\left(\frac{(n-1)S^2}{\chi^2_{\alpha/2, n-1}} \leq \sigma^2 \leq \frac{(n-1)S^2}{\chi^2_{1-\alpha/2, n-1}}\right) = 1 - \alpha$$

Therefore, the $100(1 - \alpha)\%$ confidence interval for σ^2 is

$$\left[\frac{(n-1)S^2}{\chi^2_{1-\alpha/2, n-1}}, \frac{(n-1)S^2}{\chi^2_{\alpha/2, n-1}} \right]$$

6.9 Confidence Interval for the Difference of Two Normally Distributed Population Means

When comparing the means of two normally distributed populations, it is often of interest to estimate the difference between the population means, say $\mu_1 - \mu_2$, using data from independent random samples from each population.

Let X_1, X_2, \dots, X_{n_1} be a random sample from population 1 with mean μ_1 and variance σ_1^2 , and Y_1, Y_2, \dots, Y_{n_2} be a random sample from population 2 with mean μ_2 and variance σ_2^2 . Let \bar{X} and \bar{Y} denote the respective sample means, and S_1^2, S_2^2 the respective sample variances.

We consider two cases:

6.9.1 Population Variances Known

Assume that σ_1^2 and σ_2^2 are known. Then the sampling distribution of $\bar{X} - \bar{Y}$ is normal:

$$\bar{X} - \bar{Y} \sim N\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)$$

Define the pivotal quantity:

$$Z = \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1)$$

Then, a $100(1 - \alpha)\%$ confidence interval for $\mu_1 - \mu_2$ is given by:

$$\left[(\bar{X} - \bar{Y}) - z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, (\bar{X} - \bar{Y}) + z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right]$$

6.9.2 Population Variances Unknown

Assume the populations are normally distributed but σ_1^2 and σ_2^2 are unknown. The confidence interval depends on whether variances can be assumed equal.

(i) **Equal but Unknown Variances:**

Assume $\sigma_1^2 = \sigma_2^2 = \sigma^2$, and estimate with pooled variance:

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$

Then, the pivotal quantity

$$T = \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1 + n_2 - 2}$$

Thus, the $100(1 - \alpha)\%$ confidence limits is:

$$(\bar{X} - \bar{Y}) \pm t_{\alpha/2, n_1 + n_2 - 2} \cdot S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

(ii) **Unequal and Unknown Variances (Welch's Approximation):**

If equal variance cannot be assumed, use:

$$T = \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}$$

This approximately follows a t -distribution with ν degrees of freedom, where:

$$\nu \approx \frac{\left(\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}\right)^2}{\frac{1}{n_1 - 1} \left(\frac{S_1^2}{n_1}\right)^2 + \frac{1}{n_2 - 1} \left(\frac{S_2^2}{n_2}\right)^2}$$

Then, the approximate $100(1 - \alpha)\%$ confidence interval for $\mu_1 - \mu_2$ is:

$$\left[(\bar{X} - \bar{Y}) - t_{\alpha/2, \nu} \cdot \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}, (\bar{X} - \bar{Y}) + t_{\alpha/2, \nu} \cdot \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}} \right]$$

6.10 Confidence Interval for the Ratio of Two Normally Distributed Population Variances

Let X_1, X_2, \dots, X_{n_1} be a random sample from population 1 with mean μ_1 and variance σ_1^2 , and Y_1, Y_2, \dots, Y_{n_2} be a random sample from population 2 with mean μ_2 and variance σ_2^2 . Let S_1^2, S_2^2 denote the respective sample variances.

The pivotal quantity

$$F = \frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} = \frac{S_1^2}{S_2^2} \cdot \frac{\sigma_2^2}{\sigma_1^2}$$

follows an F -distribution with $(n_1 - 1, n_2 - 1)$ degrees of freedom.

$$F \sim F_{n_1-1, n_2-1}$$

Let $F_{1-\alpha/2; n_1-1, n_2-1}$ and $F_{\alpha/2; n_1-1, n_2-1}$ be the lower and upper $\alpha/2$ -quantiles of this F -distribution:

$$P(F_{1-\alpha/2; n_1-1, n_2-1} \leq F \leq F_{\alpha/2; n_1-1, n_2-1}) = 1 - \alpha$$

or,

$$P\left(F_{1-\alpha/2; n_1-1, n_2-1} \leq \frac{S_1^2}{S_2^2} \cdot \frac{\sigma_2^2}{\sigma_1^2} \leq F_{\alpha/2; n_1-1, n_2-1}\right) = 1 - \alpha$$

Rewriting in terms of σ_1^2/σ_2^2 gives

$$P\left(\frac{1}{F_{\alpha/2; n_1-1, n_2-1}} \frac{S_1^2}{S_2^2} \leq \frac{\sigma_1^2}{\sigma_2^2} \leq \frac{1}{F_{1-\alpha/2; n_1-1, n_2-1}} \frac{S_1^2}{S_2^2}\right) = 1 - \alpha$$

Since for F -distribution,

$$F_{1-\alpha/2; n_1-1, n_2-1} = \frac{1}{F_{\alpha/2; n_2-1, n_1-1}}$$

Therefore, a $100(1 - \alpha)\%$ confidence interval for the ratio σ_1^2/σ_2^2 is

$$\left[\frac{1}{F_{\alpha/2; n_1-1, n_2-1}} \frac{S_1^2}{S_2^2}, F_{\alpha/2; n_2-1, n_1-1} \frac{S_1^2}{S_2^2} \right]$$

6.11 Confidence Interval for a Population Proportion

Let X_1, X_2, \dots, X_n be a random sample from a Bernoulli population with unknown proportion parameter p . Define the sample proportion as

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i = \frac{X}{n}$$

where X is the number of successes in the sample of size n .

For large n , the sampling distribution of \hat{p} is approximately normal by the Central Limit Theorem:

$$\hat{p} \sim N\left(p, \frac{p(1-p)}{n}\right)$$

Since p is unknown, we approximate the standard deviation using \hat{p} , leading to the following pivotal quantity:

$$Z = \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}}$$

For sufficiently large n , Z approximately follows the standard normal distribution:

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

This approximation is generally considered valid when both $n\hat{p} \geq 5$ and $n(1 - \hat{p}) \geq 5$.

Let $z_{\alpha/2}$ be the critical value from the standard normal distribution such that

$$P(Z > z_{\alpha/2}) = \frac{\alpha}{2}$$

Then,

$$P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) = 1 - \alpha$$

This implies

$$P\left(\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right) = 1 - \alpha$$

Therefore, the $100(1 - \alpha)\%$ confidence interval for the population proportion p is

$$\left[\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \quad \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right]$$

Example: A political poll surveys $n = 400$ registered voters to estimate the proportion who support Proposition X. Out of the 400 respondents, $X = 180$ say they support the measure. Construct a 95% confidence interval for the true support proportion p .

- Sample size: $n = 400$
- Number of successes: $X = 180$
- Sample proportion:

$$\hat{p} = \frac{X}{n} = \frac{180}{400} = 0.45$$

- Check normal approximation conditions:

$$n\hat{p} = 400 \times 0.45 = 180 \geq 5, \quad n(1 - \hat{p}) = 400 \times 0.55 = 220 \geq 5.$$

- For 95% confidence, $z_{\alpha/2} = z_{0.025} = 1.96$.

Compute the margin of error:

$$z_{\alpha/2} \times \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = 1.96 \times \sqrt{\frac{0.45 \times 0.55}{400}} \approx 1.96 \times 0.0249 \approx 0.0488$$

Thus the 95% confidence interval for p is

$$[0.45 - 0.0488, 0.45 + 0.0488] = [0.4012, 0.4988]$$

We are 95% confident that the true proportion of all registered voters who support Proposition X lies between 40.12% and 49.88%.

6.12 Determination of Sample Size

Sample size determination is a fundamental aspect of planning any statistical study. If the sample is too small, the results may not be trustworthy. But if the sample is too large, it can waste time, money, and effort. The goal is to find a sample size that is just right: big enough to give accurate and useful results, but not bigger than necessary. There is no universal ideal sample size for all problems. The required sample size depends on:

- The **parameter** being estimated,
- A specified **margin of error (E)**,
- A desired **confidence level** (commonly 90%, 95%, or 99%).

6.12.1 Determining Sample Size for Estimating a Population Mean

Let μ denote the population mean and σ the population standard deviation (assumed known). To estimate μ within a margin of error E at confidence level $1 - \alpha$, the sample size n must satisfy:

$$z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq E$$

Solving for n :

$$n \geq \left(\frac{z_{\alpha/2} \cdot \sigma}{E} \right)^2$$

Example: A researcher wants to estimate the average height of adult males in a city with a 95% confidence level and a margin of error of 2 cm. Suppose previous studies suggest $\sigma = 7$ cm.

$$\begin{aligned} z_{\alpha/2} &= z_{0.025} = 1.96 \\ n &\geq \left(\frac{1.96 \cdot 7}{2} \right)^2 = \left(\frac{13.72}{2} \right)^2 = (6.86)^2 = 47.06 \end{aligned}$$

Thus, a sample size of at least 48 individuals is needed.

6.12.2 Determining Sample Size for Estimating a Population Proportion

Suppose the goal is to estimate a population proportion p . The margin of error in this case must satisfy:

$$z_{\alpha/2} \cdot \sqrt{\frac{p(1-p)}{n}} \leq E$$

Solving for n :

$$n \geq \left(\frac{z_{\alpha/2}}{E} \right)^2 \cdot p(1-p)$$

If the value of p is known roughly in the neighbourhood of a value p^* , then n can be determined as:

$$n \geq \left(\frac{z_{\alpha/2}}{E} \right)^2 \cdot p^*(1-p^*)$$

Without prior knowledge of p , $p(1-p)$ can be replaced by its maximum possible value⁴ $\frac{1}{4}$ corresponding to $p = \frac{1}{2}$ and n is determined from the relation:

$$n \geq \frac{1}{4} \left(\frac{z_{\alpha/2}}{E} \right)^2$$

⁴To calculate the maximum value of $f = p(1-p) = p - p^2$, we need to differentiate it w.r.t. p and set it to zero.

$$\frac{df}{dp} = 1 - 2p = 0 \Rightarrow p = \frac{1}{2}$$

$$\frac{d^2f}{dp^2} = -2$$

The double derivative is negative indicating f is maximum at $p = \frac{1}{2}$.

The maximum value of f is

$$f_{max} = \frac{1}{2} \left(1 - \frac{1}{2} \right) = \frac{1}{4}$$

Example: A political analyst wants to estimate the proportion of voters who support a candidate with 95% confidence and a margin of error of 5%. If no prior estimate for p is available, the conservative choice is $p = 0.5$:

$$z_{\alpha/2} = z_{0.025} = 1.96, \quad E = 0.05, \quad p = 0.5$$

$$n \geq \left(\frac{1.96}{0.05} \right)^2 \times 0.5 \times 0.5 = 1536.64 \times 0.25 = 384.16$$

Thus, a sample size of at least 385 respondents is needed.

6.12.3 Finite Population Correction (FPC)

When sampling without replacement from a finite population of size N , and the sample size n exceeds 5% of the population ($n > 0.05N$), the effective sample size should be adjusted using the finite population correction (FPC):

$$n_{\text{adj}} = \frac{n_0}{1 + \frac{n_0 - 1}{N}}$$

where n_0 is the sample size calculated assuming an infinite population.

To prove this formula let us rewrite the formula for the variance of the sample mean:

- Finite population (without replacement):

$$\text{Var}(\bar{X}) = \frac{\sigma^2}{n} \times \frac{N - n}{N - 1}$$

- Infinite population (or with replacement):

$$\text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

To achieve the same precision (equal variances), set

$$\frac{\sigma^2}{n_0} = \frac{\sigma^2}{n_{\text{adj}}} \times \frac{N - n_{\text{adj}}}{N - 1}$$

where n_0 is the sample size required to achieve the variance for the infinite population and n_{adj} is the sample size required to achieve the variance for the finite population.

Cancel σ^2 and rearrange:

$$\begin{aligned} \frac{1}{n_0} &= \frac{1}{n_{\text{adj}}} \times \frac{N - n_{\text{adj}}}{N - 1} \\ \Rightarrow n_{\text{adj}} \cdot \frac{N - 1}{n_0} &= N - n_{\text{adj}} \\ \Rightarrow n_{\text{adj}} \left(1 + \frac{N - 1}{n_0} \right) &= N \\ \Rightarrow n_{\text{adj}}(n_0 + N - 1) &= Nn_0 \end{aligned}$$

Hence, the adjusted sample size is

$$n_{\text{adj}} = \frac{Nn_0}{n_0 + N - 1} = \frac{n_0}{1 + \frac{n_0 - 1}{N}}$$

Example: From the earlier example, if the population consists of only 2,000 individuals:

$$n_{\text{adj}} = \frac{384.16}{1 + \frac{384.16-1}{2000}} = \frac{384.16}{1 + 0.1916} = \frac{384.16}{1.1916} \approx 322.4$$

Thus, only about 323 individuals are needed when the population is finite.

6.13 Frequentist vs Bayesian Approaches of Statistical Inference

The foundations of statistical inference are built on two major schools of thought: the frequentist and the Bayesian.

1. The **frequentist approach**, formalized in the early 20th century by statisticians like Ronald Fisher, Jerzy Neyman, and Egon Pearson, views parameters such as population means as fixed but unknown constants. In this framework, the probability is interpreted as the long-run frequency of events, and statistical inference is based on the behavior of data across repeated random samples.
2. In contrast, the **Bayesian approach**—rooted in the 18th-century work of Thomas Bayes but revived and expanded in modern times—treats unknown parameters as random variables with probability distributions that reflect prior beliefs. These beliefs are updated using Bayes' theorem after observing data.

The key distinction lies in interpretation: frequentists quantify uncertainty through the sampling distribution of estimators, while Bayesians express uncertainty directly about parameters using probability based on both data and prior information.

To illustrate the fundamental difference between the frequentist and Bayesian approaches we consider an example where we collect the heights (in cm) of 5 randomly selected adult men from a town:

170, 172, 168, 174, 176

The sample mean is:

$$\bar{x} = \frac{170 + 172 + 168 + 174 + 176}{5} = 172$$

We assume that the population standard deviation is known and equal to $\sigma = 4$ cm.

6.13.1 Frequentist Approach

The frequentist treats the true mean μ as a fixed but unknown constant. The sample mean \bar{x} is used as a point estimator for μ . The goal is to make inference based on the behavior of \bar{x} under repeated sampling.

Confidence Interval

A 95% confidence interval for the population mean μ is given by:

$$\bar{x} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

Substituting values:

$$172 \pm 1.96 \cdot \frac{4}{\sqrt{5}} \approx 172 \pm 3.51$$

So the 95% confidence interval is:

$$[168.49, 175.51]$$

Interpretation (Frequentist):

If this sampling process were repeated many times, approximately 95% of the confidence intervals constructed this way would contain the true population mean μ . Note that in this view, μ is fixed and not random — the randomness lies in the interval.

6.13.2 Bayesian Approach

The Bayesian treats the true mean μ as a random variable and assigns a prior distribution to reflect belief about μ before seeing data. This prior is then updated using Bayes' theorem to obtain a posterior distribution after observing the data.

Prior Distribution

Suppose we believe, before observing the data, that the average male height is centered around 170 cm, but we are not very confident. We use a normal prior:

$$\mu \sim \mathcal{N}(170, 5^2)$$

Posterior Distribution

Given the normal prior and normal likelihood (due to known σ), the posterior distribution for μ is also normal:

$$\mu \mid \text{data} \sim \mathcal{N}(\mu_{\text{post}}, \sigma_{\text{post}}^2)$$

where

$$\mu_{\text{post}} = \frac{\frac{170}{25} + \frac{172}{(4^2/5)}}{\frac{1}{25} + \frac{1}{(4^2/5)}} \approx 171.4, \quad \sigma_{\text{post}} \approx 1.6$$

Credible Interval

A 95% credible interval is:

$$\mu_{\text{post}} \pm 1.96 \cdot \sigma_{\text{post}} \approx 171.4 \pm 3.14 = [168.26, 174.54]$$

Interpretation (Bayesian):

Given the prior belief and the observed data, there is a 95% probability that the true mean μ lies within the interval.

6.13.3 Summary of Differences

Aspect	Frequentist Approach	Bayesian Approach
Nature of μ	Fixed but unknown	Random variable with prior
Use of prior knowledge	Not used	Incorporated through prior
Interval interpretation	Probability attached to the confidence interval, not parameter	Probability applies to the parameter

Table 6.2: Comparison between frequentist and Bayesian interpretations.

Parameter	Assumptions	Confidence Interval Formula
Population mean μ (known σ)	Normal population or large n (CLT)	$\bar{X} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$
Population mean μ (unknown σ)	Normal population	$\bar{X} \pm t_{\alpha/2, n-1} \cdot \frac{s}{\sqrt{n}}$
Difference of means $\mu_1 - \mu_2$ (equal variances)	Normal populations, independent samples	$(\bar{X}_1 - \bar{X}_2) \pm t_{\alpha/2, n_1+n_2-2} \cdot s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$ where $s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}$
Proportion p	Large n , with $np \geq 5$, $n(1-p) \geq 5$	$\hat{p} \pm z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$
Difference of proportions $p_1 - p_2$	Large n_1, n_2	$(\hat{p}_1 - \hat{p}_2) \pm z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$
Population variance σ^2	Normal population	$\left(\frac{(n-1)s^2}{\chi_{1-\alpha/2}^2}, \frac{(n-1)s^2}{\chi_{\alpha/2}^2} \right)$
Ratio of variances $\frac{\sigma_1^2}{\sigma_2^2}$	Normal, independent samples	$\left(\frac{s_1^2/s_2^2}{F_{1-\alpha/2}}, \frac{s_1^2/s_2^2}{F_{\alpha/2}} \right)$

Table 6.3: Important confidence interval formulas.

Chapter 7

Test of Hypothesis

7.1 What is Hypothesis Testing?

There are many problems in which, rather than estimating the value of a parameter, we must decide whether a statement concerning a parameter is true or false. Statistically speaking, we test a hypothesis about a parameter.

A **statistical hypothesis** is a statement about the parameter of a population.

For example, a factory produces screws that are supposed to be 5 cm long. A quality inspector takes a sample of 50 screws and finds that the average length is 4.8 cm. Here, the hypothesis is the statement:

“The average length of the screws is 5 cm.”

Based on the sample’s average length, the inspector tests whether this difference is simply due to random variation or if the machine requires adjustment.

Hypothesis testing is a formal procedure for testing a claim or hypothesis about a population parameter using sample data.

It helps us to determine whether the evidence in a sample supports a certain belief or hypothesis about the population.

7.1.1 Null and Alternative Hypothesis

The hypothesis that will actually be tested is called the **null hypothesis**, denoted by H_0 . This is a particular claim about a population parameter. The null hypothesis is assumed to be true unless there is any strong evidence to the contrary.

The **alternative hypothesis**, denoted by H_1 or H_a , is a hypothesis that contradicts the null hypothesis. These two hypotheses are mutually exclusive and exhaustive so that one is true to the exclusion of the other. For the above example, we define the hypotheses as:

- **Null Hypothesis (H_0):** The average length of screws is 5 cm.
- **Alternative Hypothesis (H_1 or H_a):** The average length is not 5 cm.

In simple expression we write:

$$H_0 : \mu = 5$$

$$H_1 : \mu \neq 5$$

The standard procedure is to assume that H_0 is true. The burden of proof is placed on those who believe in the alternative claim¹ (alternative hypothesis). This initially favored claim (H_0) will not be rejected in favor of the alternative claim (H_1 or H_a) unless the sample evidence provides significant support for the alternative claim. If the sample does not strongly contradict H_0 , we will continue to believe in the plausibility of the null hypothesis.



Statisticians avoid saying “accept H_0 ” because it wrongly implies that the null hypothesis has been proven true. When there is no strong evidence against H_0 , we retain it—but this does not confirm its truth; instead, we say we “fail to reject H_0 ”. On the other hand, when the evidence is strong, we confidently “reject H_0 ”. Thus, rejecting H_0 is a **strong decision**, supported by sufficient evidence, while retaining H_0 is a **weak decision**, reflecting inconclusive results.

To test the hypothesis in the above example, we assume that the standard deviation $\sigma = 0.3$ is known. Since the sample size is $n = 50$, the sampling distribution of the sample mean is approximately normal:

$$\bar{X} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

or,

$$\bar{X} \sim N(50, 0.042)$$



Figure 7.1: *Sampling distribution of sample mean for the screw length.*

The distribution of \bar{X} (sampling distribution) is shown in Figure 7.1. From the figure, we can see that there is a 95% chance that the value of \bar{X} measured from a random sample will lie in the region $4.917 \leq \bar{X} \leq 5.083$. However, the observed value is 4.8 cm, which falls well outside of this interval. The chance of obtaining such an extreme value when the true mean is actually 5 cm is less than 5%.

¹A close analogy can be made to a criminal court trial, where the jury holds to the null hypothesis of “Not guilty” unless there is convincing evidence of guilt. The purpose of the hearing is to establish the assertion that the accused is guilty rather than to prove that he or she is innocent.

This provides strong statistical evidence against the null hypothesis. Therefore, we reject the null hypothesis $H_0 : \mu = 5$ in favor of the alternative hypothesis $H_1 : \mu \neq 5$. It suggests that the mean length of the screws is significantly different from 5 cm, and that the manufacturing process may need to be adjusted.

The probability of the tails of the distribution, which determines the threshold for making a decision, is called the **level of significance**, denoted by α .

In our example, we chose the level of significance $\alpha = 0.05$, which is equally split between the two tails of the distribution, allocating 0.025 to each side.

The range of values for which the null hypothesis is rejected is called the **critical region**.

For the above example, the critical region is characterized by $\bar{X} < 4.917$ and $\bar{X} > 5.083$.

7.1.2 One-Sided and Two-Sided Hypothesis Testing

In hypothesis testing, the form of the alternative hypothesis determines whether the test is **one-sided** (one-tailed) or **two-sided** (two-tailed).

- **Two-Sided Test:** Used when we are interested in detecting any difference from the null hypothesis value, whether it is an increase or a decrease. For example:

$$H_0 : \mu = 5 \quad \text{vs.} \quad H_1 : \mu \neq 5$$

This test considers deviations on both sides of the hypothesized mean as we saw in the previous example. The significance level α is split between the two tails of the sampling distribution.

- **One-Sided Test:** Used when we are only interested in deviations in one direction.
 - To test if the mean is *less than* 5:

$$H_0 : \mu = 5 \quad \text{vs.} \quad H_1 : \mu < 5$$

- To test if the mean is *greater than* 5:

$$H_0 : \mu = 5 \quad \text{vs.} \quad H_1 : \mu > 5$$

In this case, the entire significance level α is placed in one tail of the distribution.

The choice between one-sided and two-sided testing depends on the research question. If deviations in both directions are meaningful, a two-sided test is appropriate. If only an increase or a decrease is relevant, a one-sided test is more powerful.

7.2 Test Statistic

The **test statistic** is a function of the sample data that forms the basis for making the statistical decision to either reject or not reject the null hypothesis.

The main purpose of the test statistic is to provide a measure of how far the sample statistic (such as the sample mean) deviates from the hypothesized value under the null hypothesis. The further this value is from the hypothesized value, the stronger the evidence against the null hypothesis.

Depending on the type of hypothesis test being conducted, the test statistic can take various forms. For instance, a **z-test statistic** is used to test hypotheses about a population mean when the

population standard deviation (σ) is known and the sample size is large ($n > 30$). The z -test compares the sample mean to the population mean, and the test statistic is given by:

$$Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$$

where \bar{X} is the sample mean, μ_0 is the population mean under the null hypothesis, σ is the population standard deviation, and n is the sample size. If the null hypothesis is true, $\mathbb{E}(\bar{X}) = \mu_0$, and it follows that the distribution of Z is the standard normal distribution i.e. $Z \sim \mathcal{N}(0, 1)$.

For the example with a sample of 50 screws, where the sample mean is $\bar{X} = 4.8$ cm, the population mean under the null hypothesis is $\mu_0 = 5$ cm, and the population standard deviation is $\sigma = 0.3$ cm. The z -test statistic is then calculated as follows:

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} = \frac{4.8 - 5}{0.3/\sqrt{50}} \approx -4.72$$

Once the test statistic is calculated, it is compared with a **critical value** calculated from the relevant probability distribution of the test statistic under null hypothesis (e.g., the standard normal distribution). The critical value c is chosen so that

$$P(\text{Test Statistic is more extreme than } c \mid H_0) = \alpha$$

where α is the level of significance. Equivalently, the set of all values of the test statistic that lead to rejection of H_0 is called the **critical region** or **rejection region**. Therefore

$$P(\text{Test statistic} \in \text{rejection region} \mid H_0) = \alpha$$

Test Type	Critical Value(s)	Rejection Region
Upper-tailed z -test	z_α	$\{Z > z_\alpha\}$
Lower-tailed z -test	$-z_\alpha$	$\{Z < -z_\alpha\}$
Two-tailed z -test	$\pm z_{\alpha/2}$	$\{ Z > z_{\alpha/2}\}$

Table 7.1: Critical values and rejection regions for z -tests at significance level α .

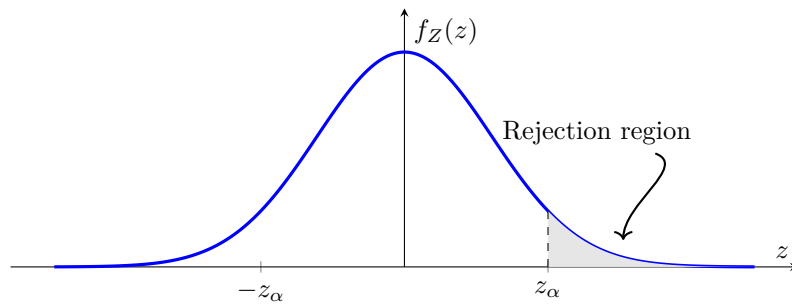


Figure 7.2: Rejection region for upper-tailed z -test.



Figure 7.3: *Rejection region for lower-tailed z-test.*



Figure 7.4: *Rejection region for two-tailed z-test.*

Now we set the decision rule:

- If the observed test statistic z_{obs} lies in the rejection region, then we reject H_0 .
- Otherwise, we fail to reject H_0 .

In our example, the computed z -value is -4.72 , which lies in the left-hand critical region $\{Z < z_{\alpha/2} = -1.96\}$ for the two-tailed test with $\alpha = 0.05$. Therefore, we reject the null hypothesis H_0 . In general, the larger the magnitude of the test statistic, the stronger the evidence against the null hypothesis.

7.3 Parametric vs. Non-Parametric Tests

In statistical inference, tests are often classified according to the assumptions they make about the population distribution. Two broad families are **parametric** and **non-parametric** tests.

- **Parametric Tests:** Parametric tests rely on specific assumptions about the form of the underlying population distribution, typically that it belongs to a known family (e.g. the normal distribution). They are characterized by:
- **Non-Parametric Tests:** Non-parametric tests make *fewer* assumptions about the population distribution. They often use the ranking or signs of the data rather than their numerical values. Key features include:

Characteristic	Parametric Tests	Non-Parametric Tests
Distribution assumption	Assume a known distribution (e.g. normal)	Do not assume a specific distribution
Data usage	Use parameters like mean (μ) and variance (σ^2)	Use ranks, signs, or counts instead of raw values
Data scale	Require interval or ratio data	Work with ordinal or non-normal interval/ratio data
When to use	For well-behaved numeric data meeting assumptions	When assumptions fail or data are ordinal
Examples	Z-test, t-test, ANOVA, Pearson's correlation	Mann-Whitney U, Wilcoxon signed-rank, chi-square

Table 7.2: Comparison of Parametric and Non-Parametric Tests

In this and the following chapters, we will focus on parametric tests. A dedicated chapter on non-parametric tests will be presented later.

7.4 Type I and Type II Error

This decision procedure can lead to either of two incorrect conclusions, known as Type I and Type II errors.

For example, suppose the true average length of the screws is indeed 5 cm. However, due to random variation in the sample, we might observe a test statistic that falls into the critical region. In this case, we would reject the null hypothesis H_0 in favor of the alternative H_1 , even though H_0 is actually true. This mistake is called a **Type I error**. The probability of a Type I error is also called the **level of significance**, denoted by α . It is usually set at $\alpha = 0.05$ or $\alpha = 0.01$.

$$P(\text{Type I error}) = P(\text{Rejecting } H_0 \mid H_0 \text{ is true}) = \alpha$$

On the other hand, suppose the true average length of the screws has actually changed (for example, to 4.8 cm), but the observed sample does not provide enough evidence (such as 4.92 cm) to reject H_0 . In this case, we fail to reject the null hypothesis, even though it is false. This mistake is called a **Type II error**. The probability of Type II error is denoted by β .

$$P(\text{Type II error}) = P(\text{Not rejecting } H_0 \mid H_0 \text{ is false}) = \beta$$

	H_0 is True	H_1 is True
Reject H_0	Type I Error (α)	Correct Decision
Fail to Reject H_0	Correct Decision	Type II Error (β)

7.4.1 Tradeoff between Type I and Type II Errors

These error probabilities are inherently related through the choice of the rejection region. Generally, reducing α increases β , and vice versa. This trade-off can be visualized by the overlapping distributions of the test statistic (sample mean in the picture) under H_0 and H_1 as shown in Figure 7.5.



Figure 7.5: *Type I and Type II error.*

Here we assume the population parameter (i.e. population mean μ) under H_0 is μ_0 and under H_1 is a specific value² μ_1 .

Mathematically, for a given critical value c :

$$\alpha = P(\text{Reject } H_0 \mid H_0 \text{ true}) = P(T > c \mid H_0)$$

$$\beta = P(\text{Fail to reject } H_0 \mid H_1 \text{ true}) = P(T \leq c \mid H_1)$$

Changing the value of c affects α and β in opposite ways, showing the trade-off between making Type I and Type II errors.

- If we make the critical value c higher, the rejection region becomes smaller (more stringent), so α will **decrease**, but β will **increase**.
- If we lower the value c , the rejection region becomes larger (more liberal), so α will **increase**, but β will **decrease**.

The Figure 7.6 shows one typical example of the relationship between α and β .



Figure 7.6: *Illustration of the inverse relationship between α and β .*

²In practice, the specific value μ_1 is unknown when the alternative hypothesis is specified as $H_1 : \mu \neq \mu_0$. However, to compute the probability of a Type II error (β), we must assume a particular value of μ under the alternative hypothesis. This means that β is not a single number, but rather a function of the true value of μ under H_1 . For every possible value of μ_1 under H_1 , there is a corresponding value of β .

Hence, a balance must be struck depending on the consequences of making Type I vs. Type II errors. To arrive at a fair compromise we should know the cost of each type of error. In practice, these costs depend on the true (but unknown) parameter and are hard to quantify exactly, so we adopt one of three conventional values—1%, 5%, or 10%—based on the relative severity of each error.

For example, consider the design of an email spam filter. A **Type I error** here means flagging a legitimate email as spam, potentially causing the user to miss important messages. A **Type II error** would allow a spam message into the inbox, usually a minor nuisance. Thus, in this case, the cost of a Type I error is considered more serious than that of a Type II error, and the filter should be designed to minimize α , the probability of Type I error, even at the expense of a slightly higher β . A common choice might be $\alpha = 1\%$, to ensure that almost no valid emails are incorrectly filtered out.

On the other hand, imagine a public health agency conducting virus screening at airports during an outbreak of a contagious disease. A **Type I error** in this context would mean falsely identifying a healthy traveler as infected, leading to temporary quarantine and inconvenience. A **Type II error** would allow an infected traveler to enter the general population, potentially triggering a wider outbreak. In this scenario, the cost of a Type II error is far more serious. To minimize this risk, we may choose a higher significance level, such as $\alpha = 10\%$, accepting more false positives in order to reduce the probability of missing actual infections.

Finally, in cases where the consequences of Type I and Type II errors are either unknown or roughly same, such as in exploratory scientific studies, a balanced approach is often taken. The conventional significance level $\alpha = 5\%$ serves as a compromise that moderately controls both types of errors in the absence of more specific cost information.

7.4.2 How β Depends on the True Parameter Value?

The probability of a Type II error (β) is not fixed; it varies depending on how far the actual parameter value is from the hypothesized value under H_0 . Specifically:

- If the true value is close to H_0 , then the evidence against H_0 is weak, and β is high (it's harder to detect the difference).
- If the true value is far from H_0 , then the evidence becomes stronger, and β is lower (easier to detect the difference).

We generally write β as a function of the parameter value θ_1 under alternative hypothesis denoted by $\beta(\theta)$.

7.4.3 Dependence on Sample Size

Earlier, we discussed that one can choose a low α or a low β depending on which type of error is more serious. But what if both Type I and Type II errors are costly, and we want to minimize both? The only effective way to achieve this is by improving the reliability of the evidence—primarily by increasing the sample size. As the sample size n increases, the sampling distributions under both H_0 and H_1 become narrower (i.e., have smaller variance), which makes it easier to distinguish between them. This shrinkage leads to a reduction in both α and β . Therefore, when minimizing both types of error is important, increasing the sample size is the most practical solution.



Figure 7.7: Effect of sample size on the inverse relationship between α and β . Larger sample size reduces β for the same α .

7.4.4 p -Value in Hypothesis Testing

The **p -value** is the probability, under null hypothesis, that the test statistic takes a value equal to or more extreme than the value actually observed.

Let T be the test statistic used to assess H_0 versus H_1 , and let t_{obs} be its observed value from the sample. Then for **one-sided test**,

$$p\text{-value} = P(T \geq t_{\text{obs}} \mid H_0)$$

For a **two-sided test** (where both large and small values of T count against H_0),

$$p\text{-value} = P(T \geq t_{\text{obs}} \mid H_0) + P(T \leq -t_{\text{obs}} \mid H_0)$$

Assuming the null distribution is symmetric:

$$p\text{-value} = 2P(T \geq t_{\text{obs}} \mid H_0)$$

- A **small** p -value (below the chosen significance level α) indicates that observing t_{obs} would be very unlikely if H_0 were true, so we **reject** H_0 .
- A **large** p -value means the data are consistent with H_0 , so we **fail to reject** H_0 .

Condition	Decision
$p\text{-value} \leq \alpha$	Reject H_0
$p\text{-value} > \alpha$	Fail to reject H_0

Table 7.3: Decision rule based on p -value.

Example. Suppose $X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$ with known $\sigma = 0.3$, and we test

$$H_0 : \mu = 5 \quad \text{vs.} \quad H_1 : \mu > 5$$

using

$$T = \frac{\bar{X} - 5}{\sigma/\sqrt{n}} \sim N(0, 1) \quad \text{under } H_0$$

If $\bar{X} = 5.1$, $n = 36$, and $\sigma = 0.3$, then

$$t_{\text{obs}} = \frac{5.1 - 5}{0.3/\sqrt{36}} = 2.00 \Rightarrow p\text{-value} = P(Z \geq 2.00) \approx 0.0228$$

Since $p\text{-value} = 0.0228 < 0.05$, we reject H_0 at 5% level of significance.

7.5 General Procedure for Hypothesis Testing

Step 1. Formulate the hypotheses

Formulate the null hypothesis (H_0), which represents the default claim about the population parameter, and the alternative hypothesis (H_1), which represents the effect or difference you aim to detect.

Step 2. Choose a significance level

Select a level of significance α (commonly 0.01, 0.05, or 0.10) that defines the maximum probability of committing a Type I error (rejecting H_0 when it is true).

Step 3. Define and calculate the test statistic

Identify an appropriate test statistic based on the hypotheses and assumption of population distribution. Then collect a sample of size n and compute the observed value of the test statistic from the data.

Step 4. Calculate the p -value or determine the critical value(s) from data

- *p-value approach:* Compute the probability, under H_0 , of observing a test statistic at least as extreme as the one obtained.
- *Critical value approach:* Identify the cutoff point(s) from the null distribution that correspond to the chosen α , and define the rejection region(s).

Step 5. Make a decision

- If the p -value is less than or equal to α , **reject** H_0 ; otherwise, **do not reject** H_0 .
- State the result in context, indicating whether there is sufficient evidence to support the alternative hypothesis at the chosen significance level.

7.6 Statistical Test for a Normally Distributed Population Mean μ

7.6.1 Variance σ^2 Known

1. Hypothesis:

$$H_0 : \mu = \mu_0$$

$$H_1 : \begin{cases} \mu > \mu_0, & \text{Right tail test} \\ \mu < \mu_0, & \text{Left tail test} \\ \mu \neq \mu_0, & \text{Two-sided test} \end{cases}$$

2. Select α (commonly 0.01, 0.05, or 0.10) as the maximum probability of a Type I error.

3. Test statistic

$$Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$$

where \bar{X} is the sample mean, σ is known, and n is the sample size. Under H_0 , $Z \sim N(0, 1)$. We then calculate the test statistic from the sample data. Let z_{obs} is the observed value of Z from sample data.

4. We now calculate the p -value from the distribution of test statistic:

$$p\text{-value} = \begin{cases} P(Z \geq z_{\text{obs}}), & \text{for } H_1 : \mu > \mu_0 \\ P(Z \leq z_{\text{obs}}), & \text{for } H_1 : \mu < \mu_0 \\ 2P(|Z| \geq |z_{\text{obs}}|), & \text{for } H_1 : \mu \neq \mu_0 \end{cases}$$

Alternatively, we can identify the critical values and the rejection region at level α :

H_1	Critical Value(s)	Rejection Region
$\mu > \mu_0$	z_α	$\{Z > z_\alpha\}$
$\mu < \mu_0$	$-z_\alpha$	$\{Z < -z_\alpha\}$
$\mu \neq \mu_0$	$\pm z_{\alpha/2}$	$\{ Z > z_{\alpha/2}\}$

5. If $p\text{-value} \leq \alpha$, **reject** H_0 ; otherwise, **do not reject** H_0 .

Alternatively, if z_{obs} lies in the rejection region, **reject** the H_0 ; otherwise, **do not reject** H_0 .

Example: A factory claims that their lightbulbs last an average of 1,000 hours. To verify this claim, a sample of lightbulbs is tested. The population standard deviation is assumed known as 50 hours.

- Let μ be the true mean lifetime of the lightbulbs (in hours).
- The null and alternative hypotheses are:

$$H_0 : \mu = 1000 \quad (\text{the claim is true})$$

$$H_1 : \mu \neq 1000 \quad (\text{the claim is false})$$

- A random sample of $n = 30$ lightbulbs is selected. The sample mean is $\bar{X} = 980$ and the population standard deviation is assumed known as $\sigma = 50$.
- Under the null hypothesis, the test statistic is:

$$z_{\text{obs}} = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} = \frac{980 - 1000}{50/\sqrt{30}} = \frac{-20}{9.13} \approx -2.19$$

- From standard normal distribution tables, the two-tailed p -value is:

$$p = 2 \cdot P(Z < -2.19) \approx 2 \cdot 0.0143 = 0.0286$$

- Since $p < 0.05$, we reject the null hypothesis at the 5% significance level. Thus there is statistically significant evidence at the 5% level to suggest that the true average lifetime of the lightbulbs differs from 1,000 hours.

7.6.2 Variance σ^2 Unknown

1. Hypothesis:

$$H_0 : \mu = \mu_0$$

$$H_1 : \begin{cases} \mu > \mu_0, & \text{Right tail test} \\ \mu < \mu_0, & \text{Left tail test} \\ \mu \neq \mu_0, & \text{Two-sided test} \end{cases}$$

2. Select α (commonly 0.01, 0.05, or 0.10) as the maximum probability of a Type I error.
3. Test statistic:

- **Large sample** ($n \geq 30$):

$$Z = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$$

where S is the sample standard deviation. By the Central Limit Theorem, $Z \sim N(0, 1)$ approximately under H_0 . Let z_{obs} is the observed value of Z from sample data.

- **Small sample** ($n < 30$):

$$T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$$

Here, T follows the t -distribution with $n - 1$ degrees of freedom under H_0 . Let t_{obs} is the observed value of T from sample data.

4. We now calculate the p -value from the distribution of the test statistic:

- **Large sample** ($n \geq 30$) — Use standard normal distribution:

$$p\text{-value} = \begin{cases} P(Z \geq z_{\text{obs}}), & \text{for } H_1 : \mu > \mu_0, \\ P(Z \leq z_{\text{obs}}), & \text{for } H_1 : \mu < \mu_0, \\ 2 P(|Z| \geq |z_{\text{obs}}|), & \text{for } H_1 : \mu \neq \mu_0 \end{cases}$$

- **Small sample** ($n < 30$) — Use t -distribution with $n - 1$ degrees of freedom:

$$p\text{-value} = \begin{cases} P(T \geq t_{\text{obs}}), & \text{for } H_1 : \mu > \mu_0, \\ P(T \leq t_{\text{obs}}), & \text{for } H_1 : \mu < \mu_0, \\ 2 P(|T| \geq |t_{\text{obs}}|), & \text{for } H_1 : \mu \neq \mu_0 \end{cases}$$

Alternatively, identify critical values and rejection regions at level α :

- **Large sample** ($n \geq 30$)

H_1	Critical Value(s)	Rejection Region
$\mu > \mu_0$	z_α	$\{Z > z_\alpha\}$
$\mu < \mu_0$	$-z_\alpha$	$\{Z < -z_\alpha\}$
$\mu \neq \mu_0$	$\pm z_{\alpha/2}$	$\{ Z > z_{\alpha/2}\}$

- **Small sample** ($n < 30$)

H_1	Critical Value(s)	Rejection Region
$\mu > \mu_0$	$t_{\alpha, n-1}$	$\{T > t_{\alpha, n-1}\}$
$\mu < \mu_0$	$-t_{\alpha, n-1}$	$\{T < -t_{\alpha, n-1}\}$
$\mu \neq \mu_0$	$\pm t_{\alpha/2, n-1}$	$\{ T > t_{\alpha/2, n-1}\}$

5. If $p\text{-value} \leq \alpha$, **reject** H_0 ; otherwise, **do not reject** H_0 .

Alternatively, if t_{obs} (for large sample) or t_{obs} (for small sample) lies in the rejection region, **reject** the H_0 ; otherwise, **do not reject** H_0 .

7.7 Statistical Test for a Normally Distributed Population Variance σ^2

1. Hypothesis:

$$H_0 : \sigma^2 = \sigma_0^2$$

$$H_1 : \begin{cases} \sigma^2 > \sigma_0^2, & \text{Right tail test} \\ \sigma^2 < \sigma_0^2, & \text{Left tail test} \\ \sigma^2 \neq \sigma_0^2, & \text{Two-sided test} \end{cases}$$

2. Select α (commonly 0.01, 0.05, or 0.10) as the maximum probability of a Type I error.

3. Test statistic

$$\chi^2 = \frac{(n-1)S^2}{\sigma_0^2}$$

where S^2 is the sample variance and n is the sample size. Under H_0 , $\chi^2 \sim \chi_{n-1}^2$ (chi-squared distribution with $n-1$ degrees of freedom). Let χ_{obs}^2 be the observed value from the sample.

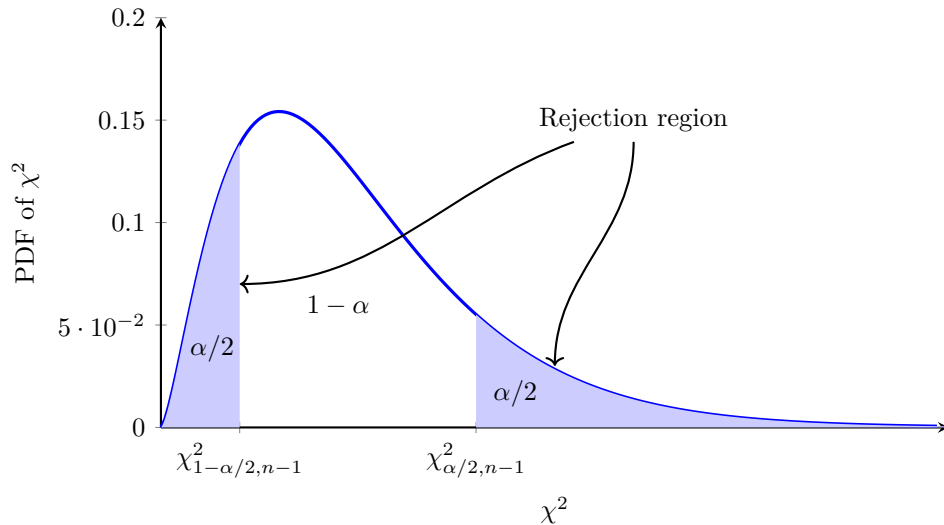


Figure 7.8: Critical values and rejection region for χ^2 test statistic for two-sided test.

4. We now calculate the p -value from the distribution of the test statistic:

$$p\text{-value} = \begin{cases} P(\chi^2 \geq \chi_{\text{obs}}^2), & \text{for } H_1 : \sigma^2 > \sigma_0^2 \\ P(\chi^2 \leq \chi_{\text{obs}}^2), & \text{for } H_1 : \sigma^2 < \sigma_0^2 \\ 2 \min\{P(\chi^2 \leq \chi_{\text{obs}}^2), P(\chi^2 \geq \chi_{\text{obs}}^2)\}, & \text{for } H_1 : \sigma^2 \neq \sigma_0^2 \end{cases}$$

In the two-sided test $H_1 : \sigma^2 \neq \sigma_0^2$, we are interested in deviations of the sample variance from σ_0^2 in both directions. Since the chi-squared distribution is not symmetric, the two-tailed p -value is computed by taking the smaller of the two tail probabilities and double it, forming a conservative and valid two-sided p -value. This formula captures the probability of observing a value as extreme as χ_{obs}^2 in either tail of the distribution.

Alternatively, we can use critical values and define the rejection region based on α :

H_1	Critical Value(s)	Rejection Region
$\sigma^2 > \sigma_0^2$	$\chi_{\alpha, n-1}^2$	$\{\chi^2 > \chi_{\alpha, n-1}^2\}$
$\sigma^2 < \sigma_0^2$	$\chi_{\alpha, n-1}^2$	$\{\chi^2 < \chi_{\alpha, n-1}^2\}$
$\sigma^2 \neq \sigma_0^2$	$\chi_{1-\alpha/2, n-1}^2, \chi_{\alpha/2, n-1}^2$	$\{\chi^2 < \chi_{1-\alpha/2, n-1}^2\} \cup \{\chi^2 > \chi_{\alpha/2, n-1}^2\}$

5. If $p\text{-value} \leq \alpha$, **reject** H_0 ; otherwise, **do not reject** H_0 .

Alternatively, if χ_{obs}^2 lies in the rejection region, **reject** H_0 ; otherwise, **do not reject** H_0 .

7.8 Statistical Test for the Difference of Two Normally Distributed Population Means $\mu_1 - \mu_2$

7.8.1 Variances σ_1^2 and σ_2^2 Known

1. Hypothesis:

$$H_0 : \mu_1 - \mu_2 = \Delta_0$$

$$H_1 : \begin{cases} \mu_1 - \mu_2 > \Delta_0, & \text{Right tail test} \\ \mu_1 - \mu_2 < \Delta_0, & \text{Left tail test} \\ \mu_1 - \mu_2 \neq \Delta_0, & \text{Two-sided test} \end{cases}$$

where Δ_0 is the hypothesized difference (often 0).

2. Select α (commonly 0.01, 0.05, or 0.10) as the maximum probability of a Type I error.
3. Test statistic

$$Z = \frac{\bar{X}_1 - \bar{X}_2 - \Delta_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

where \bar{X}_1 and \bar{X}_2 are the sample means, σ_1^2 and σ_2^2 are the known population variances, and n_1, n_2 are the respective sample sizes. Under H_0 , $Z \sim \mathcal{N}(0, 1)$. Let z_{obs} be the observed value of Z from the sample data.

4. We now calculate the p -value from the distribution of the test statistic:

$$p\text{-value} = \begin{cases} P(Z \geq z_{\text{obs}}), & \text{for } H_1 : \mu_1 - \mu_2 > \Delta_0 \\ P(Z \leq z_{\text{obs}}), & \text{for } H_1 : \mu_1 - \mu_2 < \Delta_0 \\ 2 P(|Z| \geq |z_{\text{obs}}|), & \text{for } H_1 : \mu_1 - \mu_2 \neq \Delta_0 \end{cases}$$

Alternatively, we can identify the critical values and rejection regions at level α :

H_1	Critical Value(s)	Rejection Region
$\mu_1 - \mu_2 > \Delta_0$	z_α	$\{Z > z_\alpha\}$
$\mu_1 - \mu_2 < \Delta_0$	$-z_\alpha$	$\{Z < -z_\alpha\}$
$\mu_1 - \mu_2 \neq \Delta_0$	$\pm z_{\alpha/2}$	$\{ Z > z_{\alpha/2}\}$

5. If $p\text{-value} \leq \alpha$, **reject** H_0 ; otherwise, **do not reject** H_0 .

Alternatively, if z_{obs} lies in the rejection region, **reject** H_0 ; otherwise, **do not reject** H_0 .

7.8.2 Variances σ_1^2 and σ_2^2 Unknown

1. Hypothesis:

$$H_0 : \mu_1 - \mu_2 = \Delta_0$$

$$H_1 : \begin{cases} \mu_1 - \mu_2 > \Delta_0, & \text{Right-tail test} \\ \mu_1 - \mu_2 < \Delta_0, & \text{Left-tail test} \\ \mu_1 - \mu_2 \neq \Delta_0, & \text{Two-sided test} \end{cases}$$

where Δ_0 is the hypothesized difference (often 0).

2. Select α (e.g. 0.01, 0.05, 0.10).

3. Test statistic and degrees of freedom:

(i) **Equal variances assumed** $\sigma_1^2 = \sigma_2^2$ (**pooled t -test**)

$$T = \frac{\bar{X}_1 - \bar{X}_2 - \Delta_0}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

where

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$

is the pooled variance. Under H_0 , T follows a t -distribution with ν degrees freedom where $\nu = n_1 + n_2 - 2$.

(ii) **Unequal variances assumed** $\sigma_1^2 \neq \sigma_2^2$ (**Smith-Satterthwaite test**)

$$T = \frac{\bar{X}_1 - \bar{X}_2 - \Delta_0}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}$$

Under H_0 , T approximately follows a t -distribution with ν degrees of freedom (**Welch approximation**) where

$$\nu \approx \frac{\left(\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}\right)^2}{\frac{(S_1^2/n_1)^2}{n_1 - 1} + \frac{(S_2^2/n_2)^2}{n_2 - 1}}$$

4. We now calculate the p -value using the appropriate t -distribution for test statistic:

$$p\text{-value} = \begin{cases} P(T \geq t_{\text{obs}}), & H_1 : \mu_1 - \mu_2 > \Delta_0 \\ P(T \leq t_{\text{obs}}), & H_1 : \mu_1 - \mu_2 < \Delta_0 \\ 2P(|T| \geq |t_{\text{obs}}|), & H_1 : \mu_1 - \mu_2 \neq \Delta_0 \end{cases}$$

Alternatively, we can identify the critical values and rejection regions at level α :

H_1	Critical Value(s)	Rejection Region
$\mu_1 - \mu_2 > \Delta_0$	$t_{1-\alpha, \nu}$	$\{T > t_{\alpha, \nu}\}$
$\mu_1 - \mu_2 < \Delta_0$	$-t_{\alpha, \nu}$	$\{T < -t_{\alpha, \nu}\}$
$\mu_1 - \mu_2 \neq \Delta_0$	$\pm t_{\alpha/2, \nu}$	$\{ T > t_{\alpha/2, \nu}\}$

Here $\nu = n_1 + n_2 - 2$ for the pooled case, and ν is given by the Welch–Satterthwaite formula for unequal variances.

5. If $p\text{-value} \leq \alpha$, **reject** H_0 ; otherwise, **do not reject** H_0 .

7.9 Statistical Test for the Ratio of Two Normally Distributed Population Variances

1. Let σ_1^2 and σ_2^2 be the variances of two independent normal populations. We want to test whether the ratio of these variances is equal to a specified positive value δ . The null and alternative hypotheses are:

$$H_0 : \frac{\sigma_1^2}{\sigma_2^2} = \delta$$

$$H_1 : \begin{cases} \sigma_1^2/\sigma_2^2 > \delta, & \text{Right tail test} \\ \sigma_1^2/\sigma_2^2 < \delta, & \text{Left tail test} \\ \sigma_1^2/\sigma_2^2 \neq \delta, & \text{Two-sided test} \end{cases}$$

2. Select α , the level of significance (e.g., 0.01, 0.05, or 0.10).
3. Let S_1^2 and S_2^2 be the sample variances from independent random samples of sizes n_1 and n_2 , respectively. The test statistic is:

$$F = \frac{S_1^2}{\delta S_2^2}$$

Under H_0 , $F \sim F_{\nu_1, \nu_2}$, i.e., an F-distribution with $\nu_1 = n_1 - 1$ and $\nu_2 = n_2 - 1$ degrees of freedom. Let F_{obs} be the observed value of the test statistic from the sample data.

4. We now calculate the p -value from the F-distribution:

$$p\text{-value} = \begin{cases} P(F \geq F_{\text{obs}}), & \text{for } H_1 : \sigma_1^2/\sigma_2^2 > \delta \\ P(F \leq F_{\text{obs}}), & \text{for } H_1 : \sigma_1^2/\sigma_2^2 < \delta \\ 2 \cdot \min\{P(F \geq F_{\text{obs}}), P(F \leq F_{\text{obs}})\}, & \text{for } H_1 : \sigma_1^2/\sigma_2^2 \neq \delta \end{cases}$$

Alternatively, use the critical values from the F-distribution at level α :

H_1	Critical Value(s)	Rejection Region
$\sigma_1^2/\sigma_2^2 > \delta$	$F_{\alpha; \nu_1, \nu_2}$	$\{F > F_{\alpha; \nu_1, \nu_2}\}$
$\sigma_1^2/\sigma_2^2 < \delta$	$F_{1-\alpha; \nu_1, \nu_2}$	$\{F < F_{1-\alpha; \nu_1, \nu_2}\}$
$\sigma_1^2/\sigma_2^2 \neq \delta$	$F_{1-\alpha/2; \nu_1, \nu_2}, F_{\alpha/2; \nu_1, \nu_2}$	$\{F < F_{1-\alpha/2; \nu_1, \nu_2}\} \cup \{F > F_{\alpha/2; \nu_1, \nu_2}\}$

5. If $p\text{-value} \leq \alpha$, **reject** H_0 ; otherwise, **do not reject** H_0 .

Alternatively, if F_{obs} lies in the rejection region, **reject** H_0 ; otherwise, **do not reject** H_0 .

7.10 Statistical Test for a Population Proportion p

1. Let p be the true proportion of success in a Bernoulli population. We need to test whether this population proportion is equal to a specified value p_0 . The null and alternative hypotheses can be specified as:

$$H_0 : p = p_0$$

$$H_1 : \begin{cases} p > p_0, & \text{Right tail test} \\ p < p_0, & \text{Left tail test} \\ p \neq p_0, & \text{Two-sided test} \end{cases}$$

2. Select α (e.g., 0.01, 0.05, or 0.10), the maximum allowable probability of a Type I error.
3. Let \hat{p} be the sample proportion, based on a sample of size n . Under H_0 , the test statistic is:

$$Z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

which **approximately** follows the standard normal distribution $\mathcal{N}(0, 1)$ for sufficiently large n (typically $np_0 \geq 5$ and $n(1-p_0) \geq 5$). Let z_{obs} be the observed value of the test statistic.

4. Compute the p -value using the standard normal distribution:

$$p\text{-value} = \begin{cases} P(Z \geq z_{\text{obs}}), & \text{for } H_1 : p > p_0 \\ P(Z \leq z_{\text{obs}}), & \text{for } H_1 : p < p_0 \\ 2P(|Z| \geq |z_{\text{obs}}|), & \text{for } H_1 : p \neq p_0 \end{cases}$$

Alternatively, determine the rejection region using critical values from the standard normal distribution:

H_1	Critical Value(s)	Rejection Region
$p > p_0$	z_{α}	$\{Z > z_{\alpha}\}$
$p < p_0$	$-z_{\alpha}$	$\{Z < -z_{\alpha}\}$
$p \neq p_0$	$\pm z_{\alpha/2}$	$\{ Z > z_{\alpha/2}\}$

5. If $p\text{-value} \leq \alpha$, **reject** H_0 ; otherwise, **do not reject** H_0 .

Alternatively, if z_{obs} lies in the rejection region, **reject** H_0 ; otherwise, **do not reject** H_0 .

Example: You suspect that a coin is biased. To test this, you flip the coin 100 times and observe 60 heads.

- Let p be the true probability of getting heads.
- The null and alternative hypotheses are:

$$\begin{aligned} H_0 : p &= 0.5 & (\text{the coin is fair}) \\ H_1 : p &\neq 0.5 & (\text{the coin is biased}) \end{aligned}$$

- The sample proportion is $\hat{p} = \frac{60}{100} = 0.6$.
- Under the null hypothesis, the test statistic is:

$$z_{\text{obs}} = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{0.6 - 0.5}{\sqrt{\frac{0.5 \times 0.5}{100}}} = \frac{0.1}{0.05} = 2.0$$

- From standard normal distribution tables, the two-tailed p -value is:

$$p = 2 \cdot P(Z > 2.0) \approx 2 \cdot 0.0228 = 0.0456$$

- Since $p < 0.05$, we reject the null hypothesis at the 5% significance level. Therefore there is statistically significant evidence at the 5% level to suggest that the coin may be biased.

7.11 Relationship Between Confidence Interval and Hypothesis Tests

Let X_1, X_2, \dots, X_n be a random sample from a normal population with unknown mean μ and known variance σ^2 . A $(1 - \alpha)100\%$ confidence interval for μ is

$$\left[\bar{X} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \right]$$

because,

$$P\left(\bar{X} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

If are interested in testing the null hypothesis:

$$H_0 : \mu = \mu_0 \quad \text{vs} \quad H_1 : \mu \neq \mu_0$$

The test statistic is:

$$Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$$

We reject H_0 at level α if:

$$|Z| > z_{\alpha/2} \quad \Longleftrightarrow \quad |\bar{X} - \mu_0| > z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

In contrast we cannot reject H_0 at level α if:

$$|\bar{X} - \mu_0| \leq z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

which translates to

$$\bar{X} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu_0 \leq \bar{X} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

This prompts us to state the following theorem:

Theorem: The null hypothesis $H_0 : \mu = \mu_0$ cannot be rejected in favor of the alternative hypothesis $H_1 : \mu \neq \mu_0$ at significance level α if and only if μ_0 lies within $(1 - \alpha)100\%$ CI for μ .

Generalized version: The null hypothesis $H_0 : \theta = \theta_0$ about an unknown parameter θ (such as mean, proportion, variance, etc.) cannot be rejected in favor of the alternative hypothesis $H_1 : \theta \neq \theta_0$ at significance level α if and only if the hypothesized value θ_0 lies within the $(1 - \alpha)100\%$ CI for the parameter θ .

7.12 Power of a Test

One very important concept related to error probabilities is the notion of the *power* of a test. Intuitively, power measures a test's ability to detect when the null hypothesis is false. It is the probability of rejecting H_0 given that a specific alternative is true.

The **power** of a test at a specific alternative parameter value $\theta = \theta_1$ is defined as

$$\begin{aligned}\text{Power}(\theta_1) &= \gamma(\theta_1) = P(\text{Rejecting } H_0 \mid \theta = \theta_1) \\ &= 1 - P(\text{Not rejecting } H_0 \mid \theta = \theta_1) \\ &= 1 - \beta(\theta_1)\end{aligned}$$

Here $\beta(\theta_1)$ is the probability of a Type II error when the true parameter equals θ_1 .

In other words, **it quantifies your test's sensitivity to detect real departures from the null hypothesis.**

7.12.1 Calculating Power in a One-Sample Z-Test

One-sided Z-test: Consider testing

$$H_0 : \mu = \mu_0 \quad \text{vs.} \quad H_1 : \mu > \mu_0,$$

with known population standard deviation σ , and a significance level α . The rejection region is:

$$\bar{X} > \mu_0 + z_\alpha \cdot \frac{\sigma}{\sqrt{n}}.$$

When the true mean is $\mu = \mu_1 > \mu_0$, the sampling distribution of $\bar{X} \sim N(\mu_1, \sigma^2/n)$, and the power is:

$$\gamma(\mu_1) = P\left(\bar{X} > \mu_0 + z_\alpha \cdot \frac{\sigma}{\sqrt{n}}\right)$$

Standardizing under H_1 , we get:

$$\begin{aligned}\gamma(\mu_1) &= P\left(\frac{\bar{X} - \mu_1}{\sigma/\sqrt{n}} > \frac{1}{\sigma/\sqrt{n}} \left(\mu_0 + z_\alpha \cdot \frac{\sigma}{\sqrt{n}} - \mu_1\right)\right) \\ &= P\left(Z > z_\alpha - \frac{(\mu_1 - \mu_0)\sqrt{n}}{\sigma}\right)\end{aligned}$$

Two-Sided Z-Test: Now consider the two-sided alternative:

$$H_0 : \mu = \mu_0 \quad \text{vs.} \quad H_1 : \mu \neq \mu_0.$$

with known population standard deviation σ , and a significance level α . The rejection region is:

$$\bar{X} < \mu_0 - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \quad \text{or} \quad \bar{X} > \mu_0 + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

If the true mean is $\mu = \mu_1 \neq \mu_0$, then:

$$\gamma(\mu_1) = P\left(\bar{X} < \mu_0 - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right) + P\left(\bar{X} > \mu_0 + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right)$$

Standardizing:

$$\gamma(\mu_1) = P\left(Z < -z_{\alpha/2} - \frac{(\mu_1 - \mu_0)\sqrt{n}}{\sigma}\right) + P\left(Z > z_{\alpha/2} - \frac{(\mu_1 - \mu_0)\sqrt{n}}{\sigma}\right),$$

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

In both cases, we observed the following factors affecting power in the Z -test:

- **Significance level α :** Increasing α increases power.
- **Effect size $\mu_1 - \mu_0$:** Larger differences are easier to detect.
- **Population standard deviation σ :** Smaller σ increases power.
- **Sample size n :** Increasing n increases power by reducing standard error.

Example: The power of a two-sided Z -test ($H_0 : \mu = 5$ vs. $H_1 : \mu \neq 5$) with known $\sigma = 0.3$, $n = 50$, and $\alpha = 0.05$ is the probability it will reject H_0 when μ truly equals 4.8. Here the critical region is

$$\bar{X} < 4.9168 \quad \text{or} \quad \bar{X} > 5.0832$$

If in reality $\mu = 4.8$, then

$$\bar{X} \sim N(4.8, \text{SE}^2),$$

and

$$P(\bar{X} < 4.9168) = F_Z\left(\frac{4.9168 - 4.8}{0.04243}\right) \approx 0.9971$$

while

$$P(\bar{X} > 5.0832) = 1 - F_Z\left(\frac{5.0832 - 4.8}{0.04243}\right) \approx 0$$

Thus the test's power at $\mu = 4.8$ is

$$\gamma(4.8) = P(\bar{X} < 4.9168) + P(\bar{X} > 5.0832) = 0.9971 \quad (\approx 99.7\%)$$

meaning there's a 99.7% chance (very sensitive) of detecting this 0.2-unit shift.



Figure 7.9: Density curves for $H_0 : \mu = 5$ and $H_1 : \mu = 4.8$ with $\sigma = 0.3$, $n = 50$. The dashed lines at 4.9168 and 5.0832 are the critical cutoffs; the shaded regions under the H_1 curve indicate the power ($\approx 99.7\%$).

Higher power can be achieved by increasing the sample size, accepting a larger α , or targeting a larger effect size. Before running an experiment, one often conducts a power analysis to choose n (or α) so that the power exceeds a desired threshold (commonly 0.8 or 0.9) for a specified effect size.

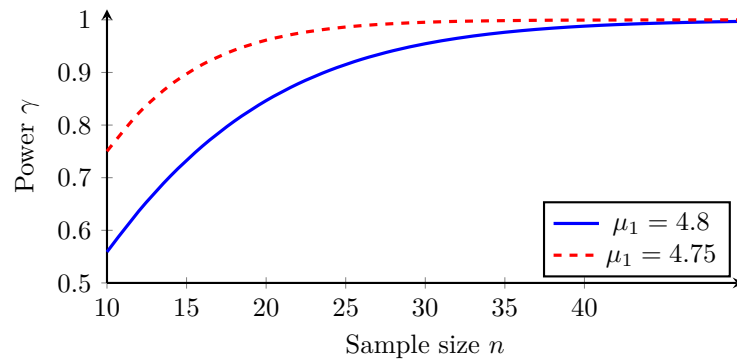


Figure 7.10: Power curve for the two-sided Z-test in the above example with $\mu_0 = 5$, $\sigma = 0.3$, $\alpha = 0.05$, showing power curves for $\mu_1 = 4.8$ and $\mu_1 = 4.75$ as sample size n varies from 10 to 40.

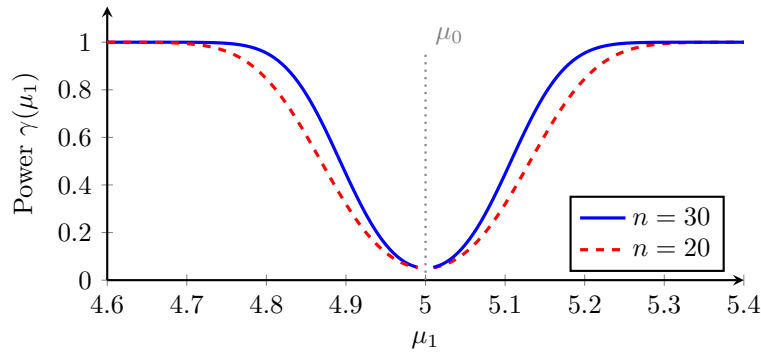


Figure 7.11: Power curve of the two-sided Z-test in the above example as the true mean μ_1 varies with $\mu_0 = 5$, $\sigma = 0.3$, and $\alpha = 0.05$. The vertical dashed line marks the null value $\mu_0 = 5$. The value of $\gamma(\mu_1)$ as $\mu_1 \rightarrow \mu_0$ is $\alpha = 0.05$. Because this is the probability of rejecting H_0 when it is actually true which happens when $\mu_1 \rightarrow \mu_0$.

7.12.2 Required Sample Size for a One-Sample Z-Test

One-sided Z-test: When planning a one-sided Z-test

$$H_0 : \mu = \mu_0 \quad \text{versus} \quad H_1 : \mu > \mu_0,$$

we wish to choose n so that the power at $\mu > \mu_1 \neq \mu_0$ is $\gamma = 1 - \beta$. Recall that

$$\gamma(\mu_1) = P\left(Z > z_\alpha - \frac{(\mu_1 - \mu_0)\sqrt{n}}{\sigma}\right)$$

By the definition of z_β ,

$$P(Z > z_\beta) = \beta \implies 1 - \beta = 1 - P(Z > z_\beta) = P(Z < z_\beta) = P(Z > -z_\beta)$$

Thus we can write

$$z_\alpha - \frac{(\mu_1 - \mu_0)\sqrt{n}}{\sigma} = -z_\beta$$

Solving for n gives the minimum sample size needed:

$$n = \left(\frac{z_\alpha + z_\beta}{(\mu_1 - \mu_0)/\sigma} \right)^2 = \left(\frac{z_\alpha + z_\beta}{\delta} \right)^2$$

where $\delta = \frac{\mu_1 - \mu_0}{\sigma}$ is the standardize effect size. The same formula applies when the alternative hypothesis is $H_1 : \mu < \mu_0$.

Thus, to achieve significance level α and power $1 - \beta$ against a shift of $\mu_1 - \mu_0$, one needs

$$n \geq \left(\frac{z_\alpha + z_\beta}{\delta} \right)^2$$

Two-sided Z-test: In two-sided tests, a closed-form expression for the required sample size is generally not available and is therefore best performed using statistical software.

7.12.3 Receiver Operating Characteristic (ROC) Curve

The ROC curve is the plot of the true positive rate (TPR) against the false positive rate (FPR) for different thresholds of a test statistic.

- **True Positive Rate (TPR) or Sensitivity:** $1 - \beta = P(\text{Reject } H_0 \mid H_1 \text{ true})$.
- **False Positive Rate (FPR):** $\alpha = P(\text{Reject } H_0 \mid H_0 \text{ true})$.

By varying the critical value c , we obtain different pairs $(\alpha, 1 - \beta)$, which can be plotted with α on the x-axis and power on the y-axis to produce the ROC curve.

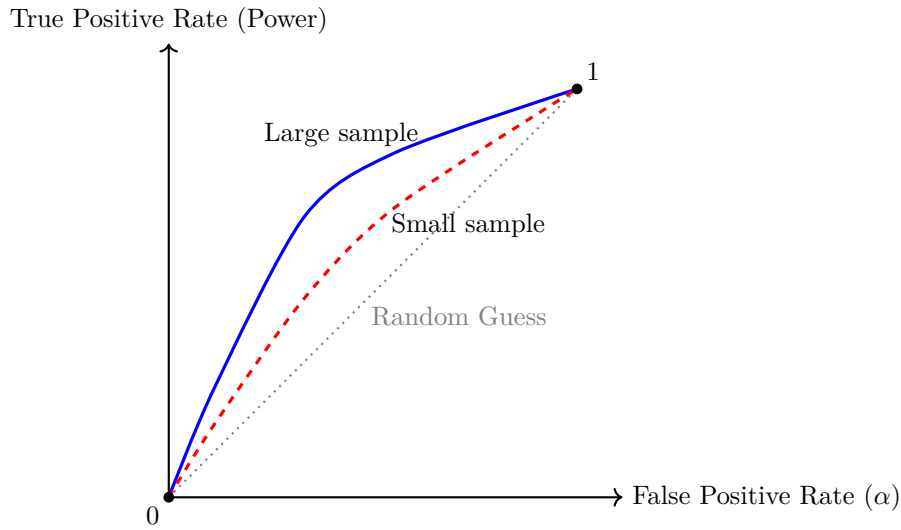


Figure 7.12: ROC curves for two sample sizes.

The ROC curve helps us visualize and compare the effectiveness of tests across different thresholds. The closer the ROC curve approaches the top-left corner, the better the test's ability to discriminate between H_0 and H_1 . From the Figure 7.12, we can see as the sample size increases, the ROC curve moves closer to the top-left corner, indicating improved test power and discrimination ability.

The **Area Under the Curve (AUC)** quantifies overall test performance: an AUC of 1 indicates perfect discrimination, whereas an AUC of 0.5 corresponds to random guessing.

Chapter 8

Non-Parametric Tests and Chi-Square Tests

8.1 Introduction

Statistical tests are broadly classified into two categories: parametric and non-parametric tests. While parametric tests make specific assumptions about the population distribution (typically assuming normality), non-parametric tests make minimal assumptions about the underlying population distribution. This chapter explores the fundamental concepts, applications, and implementations of non-parametric tests, with special emphasis on chi-square tests.

8.1.1 Parametric vs. Non-Parametric Tests

Parametric Tests:

- Assume specific probability distributions (usually normal)
- Use population parameters (μ, σ)
- Generally more powerful when assumptions are met
- Examples: t-test, ANOVA, Pearson correlation

Non-Parametric Tests:

- Make minimal distributional assumptions
- Distribution-free methods
- Robust to outliers
- Less powerful when parametric assumptions are met
- Examples: Mann-Whitney U, Wilcoxon, Kruskal-Wallis, Chi-square

8.1.2 When to Use Non-Parametric Tests

Non-parametric tests are preferred when:

1. Data violates normality assumptions

2. Sample sizes are small
3. Data contains outliers
4. Variables are ordinal or categorical
5. Distribution shape is unknown

8.2 Chi-Square Tests

The chi-square (χ^2) test is one of the most commonly used non-parametric tests. It is used to test hypotheses about categorical data and examine relationships between categorical variables.

8.2.1 Chi-Square Distribution

The chi-square distribution is a continuous probability distribution defined by:

$$f(x) = \frac{1}{2^{k/2}\Gamma(k/2)} x^{k/2-1} e^{-x/2} \quad (8.1)$$

where k is the degrees of freedom and $\Gamma(k/2)$ is the gamma function.



8.2.2 Chi-Square Test Statistic

The general form of the chi-square test statistic is:

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \quad (8.2)$$

where:

- O_i = observed frequency in category i
- E_i = expected frequency in category i
- k = number of categories

8.3 Chi-Square Goodness-of-Fit Test

The goodness-of-fit test determines whether a sample comes from a population with a specific distribution.

8.3.1 Hypotheses

$$H_0 : \text{The data follows the specified distribution} \quad (8.3)$$

$$H_1 : \text{The data does not follow the specified distribution} \quad (8.4)$$

8.3.2 Test Procedure

1. State the null and alternative hypotheses 2. Choose significance level α 3. Calculate expected frequencies: $E_i = np_i$ 4. Compute test statistic: $\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$ 5. Determine degrees of freedom: $df = k - 1 - \text{number of estimated parameters}$ 6. Find critical value or p-value 7. Make decision

8.3.3 Example 1: Testing Fair Die

A die is rolled 60 times with the following results:

Face	1	2	3	4	5	6
Observed	8	12	10	15	7	8

Test whether the die is fair at $\alpha = 0.05$.

Solution:

H_0 : The die is fair (each face has probability $\frac{1}{6}$) H_1 : The die is not fair

Expected frequency for each face: $E_i = 60 \times \frac{1}{6} = 10$

Face	Observed (O_i)	Expected (E_i)	$(O_i - E_i)^2$	$\frac{(O_i - E_i)^2}{E_i}$
1	8	10	4	0.4
2	12	10	4	0.4
3	10	10	0	0.0
4	15	10	25	2.5
5	7	10	9	0.9
6	8	10	4	0.4
Total	60	60	46	4.6

$\chi^2 = 4.6$ with $df = 6 - 1 = 5$

Critical value at $\alpha = 0.05$: $\chi_{0.05,5}^2 = 11.07$

Since $4.6 < 11.07$, we fail to reject H_0 . There is insufficient evidence to conclude the die is unfair.

8.4 Chi-Square Test of Independence

This test determines whether two categorical variables are independent.

8.4.1 Hypotheses

$$H_0 : \text{The two variables are independent} \quad (8.5)$$

$$H_1 : \text{The two variables are dependent} \quad (8.6)$$

8.4.2 Expected Frequencies

For a contingency table with r rows and c columns:

$$E_{ij} = \frac{R_i \times C_j}{n} \quad (8.7)$$

where:

- R_i = total for row i
- C_j = total for column j
- n = grand total

Degrees of freedom: $df = (r - 1)(c - 1)$

8.4.3 Example 2: Gender and Preference

A survey asks 200 people about their preference for coffee or tea, categorized by gender:

	Coffee	Tea	Total
Male	60	40	100
Female	30	70	100
Total	90	110	200

Test for independence at $\alpha = 0.01$.

Solution:

H_0 : Gender and beverage preference are independent H_1 : Gender and beverage preference are dependent

Expected frequencies:

$$E_{11} = \frac{100 \times 90}{200} = 45 \quad (8.8)$$

$$E_{12} = \frac{100 \times 110}{200} = 55 \quad (8.9)$$

$$E_{21} = \frac{100 \times 90}{200} = 45 \quad (8.10)$$

$$E_{22} = \frac{100 \times 110}{200} = 55 \quad (8.11)$$

Cell	Observed	Expected	$(O - E)^2$	$\frac{(O - E)^2}{E}$
Male-Coffee	60	45	225	5.00
Male-Tea	40	55	225	4.09
Female-Coffee	30	45	225	5.00
Female-Tea	70	55	225	4.09
Total	200	200	900	18.18

$$\chi^2 = 18.18 \text{ with } df = (2 - 1)(2 - 1) = 1$$

Critical value at $\alpha = 0.01$: $\chi_{0.01,1}^2 = 6.635$

Since $18.18 > 6.635$, we reject H_0 . There is strong evidence that gender and beverage preference are dependent.

8.5 Chi-Square Test of Homogeneity

This test compares the distribution of a categorical variable across different populations.

8.5.1 Example 3: Comparing Three Populations

Three different teaching methods are used, and student performance is categorized as Pass or Fail:

Method	Pass	Fail	Total
A	25	15	40
B	30	10	40
C	20	20	40
Total	75	45	120

Test whether the three methods have the same success rate at $\alpha = 0.05$.

Solution:

H_0 : The three methods have the same success rate H_1 : At least one method has a different success rate

Expected frequencies (assuming equal success rates):

$$E_{A,Pass} = \frac{40 \times 75}{120} = 25 \quad (8.12)$$

$$E_{A,Fail} = \frac{40 \times 45}{120} = 15 \quad (8.13)$$

$$E_{B,Pass} = \frac{40 \times 75}{120} = 25 \quad (8.14)$$

$$E_{B,Fail} = \frac{40 \times 45}{120} = 15 \quad (8.15)$$

$$E_{C,Pass} = \frac{40 \times 75}{120} = 25 \quad (8.16)$$

$$E_{C,Fail} = \frac{40 \times 45}{120} = 15 \quad (8.17)$$

Cell	Observed	Expected	$(O - E)^2$	$\frac{(O-E)^2}{E}$
A-Pass	25	25	0	0.00
A-Fail	15	15	0	0.00
B-Pass	30	25	25	1.00
B-Fail	10	15	25	1.67
C-Pass	20	25	25	1.00
C-Fail	20	15	25	1.67
Total	120	120	100	5.34

$$\chi^2 = 5.34 \text{ with } df = (3 - 1)(2 - 1) = 2$$

Critical value at $\alpha = 0.05$: $\chi_{0.05,2}^2 = 5.991$

Since $5.34 < 5.991$, we fail to reject H_0 . There is insufficient evidence to conclude that the methods have different success rates.

8.6 Other Non-Parametric Tests

8.6.1 Mann-Whitney U Test

The Mann-Whitney U test (also known as the Wilcoxon rank-sum test) is used to compare two independent samples.

Test Statistic:

$$U = n_1 n_2 + \frac{n_1(n_1 + 1)}{2} - R_1 \quad (8.18)$$

where R_1 is the sum of ranks for sample 1.

8.6.2 Wilcoxon Signed-Rank Test

Used for comparing two related samples or repeated measurements.

Test Statistic:

$$W = \min(W^+, W^-) \quad (8.19)$$

where W^+ and W^- are the sum of positive and negative ranks, respectively.

8.6.3 Kruskal-Wallis Test

Non-parametric alternative to one-way ANOVA for comparing multiple independent groups.

Test Statistic:

$$H = \frac{12}{N(N+1)} \sum_{i=1}^k \frac{R_i^2}{n_i} - 3(N+1) \quad (8.20)$$

where R_i is the sum of ranks for group i .

8.7 Assumptions and Limitations

8.7.1 Chi-Square Test Assumptions

1. **Independence:** Observations must be independent 2. **Expected Frequencies:** Each expected frequency should be at least 5 3. **Sample Size:** Total sample size should be reasonably large 4. **Categorical Data:** Variables must be categorical

8.7.2 Dealing with Small Expected Frequencies

When expected frequencies are less than 5:

- Combine categories when meaningful
- Use Fisher's exact test for 2×2 tables
- Increase sample size
- Use alternative tests

8.8 Effect Size Measures

8.8.1 Cramér's V

For measuring the strength of association in contingency tables:

$$V = \sqrt{\frac{\chi^2}{n \times \min(r-1, c-1)}} \quad (8.21)$$

where r and c are the number of rows and columns, respectively.

Interpretation:

- $V = 0$: No association
- $V = 1$: Perfect association
- $0.1 \leq V < 0.3$: Small effect
- $0.3 \leq V < 0.5$: Medium effect
- $V \geq 0.5$: Large effect

8.8.2 Phi Coefficient

For 2×2 contingency tables:

$$\phi = \sqrt{\frac{\chi^2}{n}} \quad (8.22)$$

8.9 Power and Sample Size Calculations

8.9.1 Power of Chi-Square Tests

The power of a chi-square test depends on:

- Effect size
- Sample size
- Significance level
- Degrees of freedom

8.9.2 Sample Size Calculation

For a chi-square goodness-of-fit test, the required sample size is:

$$n = \frac{(z_{\alpha/2} + z_{\beta})^2 \sum_{i=1}^k \frac{(p_i - p_{0i})^2}{p_{0i}}}{\left(\sum_{i=1}^k \frac{(p_i - p_{0i})^2}{p_{0i}} \right)^2} \quad (8.23)$$

where p_i and p_{0i} are the true and hypothesized proportions, respectively.

8.10 Computational Examples

8.10.1 Example 4: Large Contingency Table

A study examines the relationship between education level and job satisfaction:

Education	Very Satisfied	Satisfied	Dissatisfied	Total
High School	20	30	10	60
Bachelor's	35	40	15	90
Master's	25	20	5	50
Total	80	90	30	200

Expected Frequencies:

Education	Very Satisfied	Satisfied	Dissatisfied
High School	$\frac{60 \times 80}{200} = 24$	$\frac{60 \times 90}{200} = 27$	$\frac{60 \times 30}{200} = 9$
Bachelor's	$\frac{90 \times 80}{200} = 36$	$\frac{90 \times 90}{200} = 40.5$	$\frac{90 \times 30}{200} = 13.5$
Master's	$\frac{50 \times 80}{200} = 20$	$\frac{50 \times 90}{200} = 22.5$	$\frac{50 \times 30}{200} = 7.5$

Chi-Square Calculation:

$$\chi^2 = \frac{(20 - 24)^2}{24} + \frac{(30 - 27)^2}{27} + \frac{(10 - 9)^2}{9} \quad (8.24)$$

$$+ \frac{(35 - 36)^2}{36} + \frac{(40 - 40.5)^2}{40.5} + \frac{(15 - 13.5)^2}{13.5} \quad (8.25)$$

$$+ \frac{(25 - 20)^2}{20} + \frac{(20 - 22.5)^2}{22.5} + \frac{(5 - 7.5)^2}{7.5} \quad (8.26)$$

$$= 0.667 + 0.333 + 0.111 + 0.028 + 0.006 + 0.167 \quad (8.27)$$

$$+ 1.250 + 0.278 + 0.833 \quad (8.28)$$

$$= 3.673 \quad (8.29)$$

With $df = (3 - 1)(3 - 1) = 4$ and $\alpha = 0.05$, critical value is $\chi_{0.05,4}^2 = 9.488$.

Since $3.673 < 9.488$, we fail to reject H_0 . There is insufficient evidence of an association between education level and job satisfaction.

8.11 Advanced Topics

8.11.1 Yates' Continuity Correction

For 2×2 contingency tables with small expected frequencies, Yates' correction is applied:

$$\chi_{Yates}^2 = \sum_{i=1}^4 \frac{(|O_i - E_i| - 0.5)^2}{E_i} \quad (8.30)$$

8.11.2 McNemar's Test

For paired categorical data (e.g., before/after comparisons):

$$\chi^2 = \frac{(|b - c| - 1)^2}{b + c} \quad (8.31)$$

where b and c are the discordant pairs in a 2×2 table.

8.11.3 Cochran's Q Test

Extension of McNemar's test for more than two related samples:

$$Q = \frac{k(k-1) \sum_{j=1}^k (C_j - \bar{C})^2}{k \sum_{i=1}^n R_i - \sum_{i=1}^n R_i^2} \quad (8.32)$$

8.12 Practical Considerations

8.12.1 Choosing the Right Test

Data Type	Samples	Appropriate Test
One categorical variable	One sample	Chi-square goodness-of-fit
Two categorical variables	Independent samples	Chi-square test of independence
One categorical variable	Multiple independent samples	Chi-square test of homogeneity
Two categorical variables	Paired samples	McNemar's test
One categorical variable	Multiple paired samples	Cochran's Q test

8.12.2 Reporting Results

When reporting chi-square test results, include:

- Test statistic and degrees of freedom
- P-value or critical value comparison
- Effect size measure
- Interpretation in context

Example: "A chi-square test of independence revealed a significant association between gender and beverage preference, $\chi^2(1, N = 200) = 18.18, p < 0.001$. The effect size was large (Cramér's $V = 0.30$)."

8.13 Conclusion

Non-parametric tests, particularly chi-square tests, are essential tools in statistical analysis when dealing with categorical data or when parametric assumptions are violated. They provide robust methods for testing hypotheses about proportions, independence, and homogeneity. While generally less powerful than their parametric counterparts when assumptions are met, they offer valuable alternatives that are widely applicable across various fields of research.

The key to successful application of these tests lies in understanding their assumptions, choosing the appropriate test for the research question, and correctly interpreting the results in the context of the study. As with all statistical procedures, these tests should be used as part of a comprehensive analytical approach that includes careful data exploration, assumption checking, and thoughtful interpretation of results.

Chapter 9

Analysis of Variance (ANOVA)

9.1 One-Way ANOVA

Analysis of Variance (ANOVA) is a statistical method used to determine whether there are significant differences (that is unlikely to be due to chance alone) between the means of three or more independent groups. Specifically, **one-way ANOVA**—also called **one-factor ANOVA**—examines the impact of a single categorical independent variable (called a *factor*) on a continuous dependent variable.

The key question one-way ANOVA addresses is:

Are the population means of all groups equal?

This can be stated formally using hypotheses:

$$\begin{aligned} H_0 : \mu_1 = \mu_2 = \cdots = \mu_k & \quad (\text{All group means are equal}) \\ H_1 : \text{At least one } \mu_i \text{ is different} & \quad (\text{At least one group differs}) \end{aligned}$$

Here, μ_i represents the population mean of the i -th group, for $i = 1, 2, \dots, k$.

When comparing just two groups, a standard **pooled t -test** is sufficient. However, when there are more than two groups ($k > 2$), performing all possible pairwise t -tests becomes inefficient and statistically problematic. This is because the number of comparisons increases rapidly:

$$\binom{k}{2} = \frac{k!}{2!(k-2)!}$$

For example, with $k = 5$ groups, there are $\binom{5}{2} = 10$ comparisons; with $k = 10$, there are $\binom{10}{2} = 45$ comparisons.

Moreover, each pairwise t -test carries a risk of a false positive (Type I error), typically 5%. When many such tests are conducted, the overall probability of making *at least one* false claim rises dramatically. For $k = 10$, the false alarm probability is

$$1 - 0.95^{\binom{k}{2}} = 1 - 0.95^{45} \approx 1 - 0.10 = 0.90$$

This means there's a 90% chance of a false positive if 45 tests are run independently at the 5% level.

One-way ANOVA solves this issue by combining all comparisons into a single test, controlling the overall error rate and providing a more reliable answer to whether the group means differ.

Suppose we want to compare three teaching methods (A, B, and C) using student test scores:

Method	Student 1	Student 2	Student 3	Student 4
A	70	72	68	75
B	80	82	78	85
C	65	60	62	63

Table 9.1: *Test scores under three teaching methods.*

Instead of conducting three separate t -tests (A vs. B, B vs. C, A vs. C), we use a single one-way ANOVA to answer:

Is there evidence that at least one teaching method leads to a different average score?

This approach simplifies analysis, avoids excessive error inflation, and gives a more holistic view of group differences.

9.1.1 How One-way ANOVA Works?

Essentially, one-way ANOVA works by comparing two types of variation:

- **Between-group variation**, which measures differences among the group means. This variation is attributed to the effect of different **treatments** or conditions, often referred to as *controlled causes*.
- **Within-group variation**, which captures the natural variability among individuals within the same group. This variation is considered to be due to random chance.

In the context of our teaching methods example, the between-group variation arises from using different teaching approaches (Methods A, B, and C), while the within-group variation reflects individual differences among students who received the same method.

To illustrate:

- If the average scores of the three groups differ substantially compared to the variation within each group, this suggests that the teaching method has a significant effect.
- On the other hand, if the group means are similar and the within-group variability is large, we are less likely to conclude that the method has any real impact.

This comparison of between-group and within-group variances forms the basis of the ANOVA test.

9.2 Mathematical Model for One-Way ANOVA

The one-way ANOVA model begins by expressing each observation as the sum of a group-specific mean and a random error component. Specifically, for the j -th observation in the i -th group, we write:

$$Y_{ij} = \mu_i + \epsilon_{ij}$$

where:

- μ_i is the true mean of the i -th group,
- ϵ_{ij} represents random error,

- $i = 1, 2, \dots, k$ (number of groups),
- $j = 1, 2, \dots, n_i$ (number of observations in group i).

An equivalent and commonly used formulation decomposes the group-specific mean μ_i into an overall (grand) mean and a treatment effect. This gives the model:

$$Y_{ij} = \mu + \alpha_i + \epsilon_{ij}$$

where:

- μ is the grand mean (the average across all groups and observations),
- $\alpha_i = \mu_i - \mu$ is the treatment effect, i.e., the deviation of the i -th group mean from the grand mean,
- $\sum_{i=1}^k n_i \alpha_i = 0$ is a constraint imposed on the treatment effect.

This additive decomposition clearly separates the contribution of the overall mean, the group-specific deviation, and random noise, which is central to the analysis of variance.

9.2.1 Key Assumptions

For valid inference using one-way ANOVA, the following assumptions must be satisfied:

1. **Independence:** Observations must be independent both within and across groups. This means that the value of one observation should not influence or predict another.
2. **Homogeneity of Variances (Homoscedasticity):** The variance of the response should be the same across all groups:

$$\text{Var}(Y_{1j}) = \text{Var}(Y_{2j}) = \dots = \text{Var}(Y_{kj}) = \sigma^2$$

3. **Normality of Errors:** The error terms are independently and identically distributed as:

$$\epsilon_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

Consequently, $Y_{ij} \sim \mathcal{N}(\mu + \alpha_i, \sigma^2)$ for each group i . This assumption is especially important for small samples; with larger samples, ANOVA is robust due to the Central Limit Theorem.

9.3 Analysis of One-Way ANOVA

Suppose we have k independent treatment groups, where i th group has n_i observations. The total number of observations is $N = \sum_{i=1}^k n_i$. The following table 9.2 shows the observations for the one-way ANOVA model.

						Total	Mean
Treatment 1	Y_{11}	Y_{12}	Y_{13}	\dots	Y_{1n_1}	$Y_{1\bullet}$	$\bar{Y}_{1\bullet}$
Treatment 2	Y_{21}	Y_{22}	Y_{23}	\dots	Y_{2n_2}	$Y_{2\bullet}$	$\bar{Y}_{2\bullet}$
Treatment 3	Y_{31}	Y_{32}	Y_{33}	\dots	Y_{3n_2}	$Y_{3\bullet}$	$\bar{Y}_{3\bullet}$
\vdots	\vdots	\vdots			\vdots	\vdots	\vdots
Treatment k	Y_{k1}	Y_{k2}	Y_{k3}	\dots	Y_{kn_k}	$Y_{k\bullet}$	$\bar{Y}_{k\bullet}$
Overall						$Y_{\bullet\bullet}$	$\bar{Y}_{\bullet\bullet}$

Table 9.2: Random observations for one-way ANOVA.

In this table:

- Y_{ij} is the response variable corresponding to the j th observation in the i th group.
- $Y_{i\bullet}$ is the total of all observations in the i th group:

$$Y_{i\bullet} = \sum_{j=1}^{n_i} Y_{ij}$$

- $\bar{Y}_{i\bullet}$ is the sample mean of the i th group:

$$\bar{Y}_{i\bullet} = \frac{Y_{i\bullet}}{n_i} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}$$

- $Y_{\bullet\bullet}$ is the **grand total** of all observations:

$$Y_{\bullet\bullet} = \sum_{i=1}^k Y_{i\bullet} = \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}$$

- $\bar{Y}_{\bullet\bullet}$ is the **overall (grand) mean**:

$$\bar{Y}_{\bullet\bullet} = \frac{Y_{\bullet\bullet}}{N} = \frac{1}{N} \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij} = \frac{1}{N} \sum_{i=1}^k n_i \bar{Y}_{i\bullet}$$

9.3.1 Sum of Squares

As mentioned before, the analysis of variance partitions the total variability in the sample data into two component parts: between-group variation and within-group variation. It is customary to denote

- total variability as **Total Sum of Squares (TSS)**, where

$$\text{TSS} = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{\bullet\bullet})^2$$

- between-group variation as **Sum of Squares Between treatment groups (SST)**, where

$$\text{SST} = \sum_{i=1}^k (n_i \bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2$$

- within-group variation as **Sum of Squares of Errors (SSE)**, where

$$\text{SSE} = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2$$

Mathematically, it can be shown that

$$\text{TSS} = \text{SST} + \text{SSE}$$

Proof: Each observed data Y_{ij} can be decomposed as

$$\underbrace{Y_{ij}}_{\text{observation}} = \underbrace{\bar{Y}_{\bullet\bullet}}_{\text{grand mean}} + \underbrace{(\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})}_{\text{deviation due to treatment}} + \underbrace{(Y_{ij} - \bar{Y}_{i\bullet})}_{\text{error}}$$

Overall deviation of each observed data can be written as

$$Y_{ij} - \bar{Y}_{..} = (Y_{ij} - \bar{Y}_{i.}) + (\bar{Y}_{i.} - \bar{Y}_{..})$$

$$\begin{aligned} \text{TSS} &= \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{..})^2 \\ &= \sum_{i=1}^k \sum_{j=1}^{n_i} [(Y_{ij} - \bar{Y}_{i.}) + (\bar{Y}_{i.} - \bar{Y}_{..})]^2 \\ &= \sum_{i=1}^k \sum_{j=1}^{n_i} \left[(Y_{ij} - \bar{Y}_{i.})^2 + 2(Y_{ij} - \bar{Y}_{i.})(\bar{Y}_{i.} - \bar{Y}_{..}) + (\bar{Y}_{i.} - \bar{Y}_{..})^2 \right] \\ &= \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.})^2 + 2 \sum_{i=1}^k (\bar{Y}_{i.} - \bar{Y}_{..}) \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.}) + \sum_{i=1}^k \sum_{j=1}^{n_i} (\bar{Y}_{i.} - \bar{Y}_{..})^2 \end{aligned}$$

Since for each group i , we have:

$$\sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.}) = 0,$$

the middle term vanishes.

Therefore,

$$\begin{aligned} \text{TSS} &= \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.})^2 + n \sum_{i=1}^k (\bar{Y}_{i.} - \bar{Y}_{..})^2 \\ &= \text{SSE} + \text{SST} \end{aligned}$$

9.3.2 Computational Formulas

For computational efficiency, the following formulas are often preferred:

$$\text{TSS} = \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}^2 - \frac{Y_{..}^2}{N}$$

Proof:

$$\begin{aligned} \text{TSS} &= \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{..})^2 \\ &= \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij}^2 - 2Y_{ij}\bar{Y}_{..} + \bar{Y}_{..}^2) \\ &= \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}^2 - 2\bar{Y}_{..} \underbrace{\sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}}_{=N\bar{Y}_{..}} + \sum_{i=1}^k \sum_{j=1}^{n_i} \bar{Y}_{..}^2 \\ &= \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}^2 - 2N\bar{Y}_{..}^2 + N\bar{Y}_{..}^2 \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}^2 - N \bar{Y}_{..}^2 \\
&= \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}^2 - N \left(\frac{Y_{..}^2}{N} \right) \\
&= \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}^2 - \frac{Y_{..}^2}{N}
\end{aligned}$$

■

$$\text{SST} = \sum_{i=1}^k \frac{Y_{i.}^2}{n_i} - \frac{Y_{..}^2}{N}$$

Proof:

$$\begin{aligned}
\text{SST} &= \sum_{i=1}^k n_i (\bar{Y}_{i.} - \bar{Y}_{..})^2 \\
&= \sum_{i=1}^k n_i (\bar{Y}_{i.}^2 - 2\bar{Y}_{i.}\bar{Y}_{..} + \bar{Y}_{..}^2) \\
&= \sum_{i=1}^k n_i \bar{Y}_{i.}^2 - 2\bar{Y}_{..} \underbrace{\sum_{i=1}^k n_i \bar{Y}_{i.}}_{=N\bar{Y}_{..}} + \bar{Y}_{..}^2 \underbrace{\sum_{i=1}^k n_i}_{=N} \\
&= \sum_{i=1}^k n_i \bar{Y}_{i.}^2 + N\bar{Y}_{..}^2 - 2N\bar{Y}_{..}^2 \\
&= \sum_{i=1}^k n_i \bar{Y}_{i.}^2 - N\bar{Y}_{..}^2 \\
&= \sum_{i=1}^k n_i \left(\frac{Y_{i.}}{n_i} \right)^2 - N \left(\frac{Y_{..}}{N} \right)^2 \\
&= \sum_{i=1}^k \frac{Y_{i.}^2}{n_i} - \frac{Y_{..}^2}{N}
\end{aligned}$$

■

Once we have TSS and SST, we can obtain SSE by subtracting SST from TSS.

9.3.3 Degrees of Freedom in ANOVA

Each component has an associated **degree of freedom** (df), which reflects the number of independent quantities involved in estimating a particular sum of squares.

- **Total degrees of freedom:** The total number of observations is $N = kn$, where k is the number of groups and n is the number of observations per group. Since we compute the grand mean while calculating TSS and it imposes only one constraint. Thus one degree of freedom is lost. The total degrees of freedom is:

$$df_{\text{TSS}} = N - 1 = kn - 1$$

- **Degrees of freedom between groups (Treatment):** We estimate k group means. These k group means are considered k independent values. However, computing the grand mean from these k means imposes one constraint, so one degree of freedom is lost. Thus:

$$df_{\text{SST}} = k - 1$$

- **Degrees of freedom within groups (Error):** Within each group, we compute deviations from the group mean. Since each group contributes $n - 1$ independent deviations and there are k groups:

$$df_{\text{SSE}} = k(n - 1) = kn - k$$

These degrees of freedom satisfy the identity:

$$df_{\text{TSS}} = df_{\text{SST}} + df_{\text{SSE}}$$

because

$$(kn - 1) = (k - 1) + (kn - k)$$

9.3.4 Mean Squares

Mean squares are obtained by dividing the sums of squares by their respective degrees of freedom:

- **Mean Square for Treatments (MST):**

$$\text{MST} = \frac{\text{SST}}{df_{\text{SST}}} = \frac{\text{SST}}{k - 1}$$

- **Mean Square for error (MSE):**

$$\text{MSE} = \frac{\text{SSE}}{df_{\text{SSE}}} = \frac{\text{SSE}}{N - k}$$

Theorem: Under the null hypothesis $H_0 : \mu_1 = \mu_2 = \dots = \mu_k$:

$$\mathbb{E}[\text{MST}] = \sigma^2, \quad \mathbb{E}[\text{MSE}] = \sigma^2$$

Under the alternative hypothesis:

$$\mathbb{E}[\text{MST}] = \sigma^2 + \frac{\sum_{i=1}^k n_i (\mu_i - \mu)^2}{k - 1}, \quad \mathbb{E}[\text{MSE}] = \sigma^2$$

where $\mu = \frac{1}{N} \sum_{i=1}^k n_i \mu_i$ is the weighted grand mean.

Proof: Consider the one-way ANOVA model:

$$Y_{ij} = \mu_i + \epsilon_{ij} = \mu + \alpha_i + \epsilon_{ij}$$

where Y_{ij} is the j -th observation in the i -th group, $i = 1, 2, \dots, k$, $j = 1, 2, \dots, n_i$, μ is the overall mean, α_i is the effect of the i -th group, $\epsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$ are independent error terms, and $N = \sum_{i=1}^k n_i$ is the total sample size. Define the sample means:

$$\bar{Y}_{i\bullet} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}, \quad \bar{Y}_{\bullet\bullet} = \frac{1}{N} \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}$$

The mean squares are:

$$\text{MST} = \frac{\sum_{i=1}^k n_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2}{k - 1}, \quad \text{MSE} = \frac{\sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2}{N - k}$$

- **Expected value of MSE:**

From the model, we have:

$$\bar{Y}_{i\bullet} = \mu + \alpha_i + \bar{\epsilon}_{i\bullet}, \quad \text{where } \bar{\epsilon}_{i\bullet} = \frac{1}{n_i} \sum_{j=1}^{n_i} \epsilon_{ij}$$

Therefore:

$$Y_{ij} - \bar{Y}_{i\bullet} = (\mu + \alpha_i + \epsilon_{ij}) - (\mu + \alpha_i + \bar{\epsilon}_{i\bullet}) = \epsilon_{ij} - \bar{\epsilon}_{i\bullet}$$

Thus:

$$\begin{aligned} \text{SSE} &= \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (\epsilon_{ij} - \bar{\epsilon}_{i\bullet})^2 \\ &= \sum_{i=1}^k \sum_{j=1}^{n_i} \epsilon_{ij}^2 - 2 \sum_{i=1}^k \sum_{j=1}^{n_i} \epsilon_{ij} \bar{\epsilon}_{i\bullet} + \sum_{i=1}^k \sum_{j=1}^{n_i} \bar{\epsilon}_{i\bullet}^2 \\ &= \sum_{i=1}^k \sum_{j=1}^{n_i} \epsilon_{ij}^2 - \sum_{i=1}^k n_i \bar{\epsilon}_{i\bullet}^2 \end{aligned}$$

Since $\epsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$, we can write¹

$$\mathbb{E}(\epsilon_{ij}^2) = \sigma^2$$

Also, since

$$\bar{\epsilon}_{i\bullet} = \frac{1}{n_i} \sum_{j=1}^{n_i} \epsilon_{ij}$$

we can write

$$\bar{\epsilon}_{i\bullet} \sim \mathcal{N}\left(0, \frac{\sigma^2}{n_i}\right)$$

Hence,

$$\mathbb{E}(\bar{\epsilon}_{i\bullet}^2) = \frac{\sigma^2}{n_i}$$

Now,

$$\begin{aligned} \mathbb{E}(\text{SSE}) &= \sum_{i=1}^k \sum_{j=1}^{n_i} \mathbb{E}(\epsilon_{ij}^2) - \sum_{i=1}^k n_i \mathbb{E}(\bar{\epsilon}_{i\bullet}^2) \\ &= N\sigma^2 - kn_i \frac{\sigma^2}{n_i} = (N - k)\sigma^2 \end{aligned}$$

Thus,

$$\mathbb{E}[\text{MSE}] = \mathbb{E}\left[\frac{\text{SSE}}{N - k}\right] = \sigma^2$$

- **Expected Value of MST:**

¹This comes from the theorem that if $X \sim N(0, \sigma^2)$, then $\mathbb{E}(X^2) = \sigma^2$. To prove this we write the formula of variance:

$$\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}(X))^2] = \mathbb{E}[X^2 - 2X\mathbb{E}(X) + (\mathbb{E}(X))^2]$$

Since $\mathbb{E}(X) = 0$, this simplifies to:

$$\text{Var}(X) = \mathbb{E}(X^2)$$

Hence,

$$\mathbb{E}(X^2) = \sigma^2$$

– **Case 1: Under $H_0 : \mu_1 = \mu_2 = \dots = \mu_k = 0$**

The equivalent hypothesis in terms of α 's can be written as:

$$H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_k = 0$$

Under H_0 , thus we have $Y_{ij} = \mu + \epsilon_{ij}$, so:

$$\bar{Y}_{i\bullet} = \mu + \bar{\epsilon}_{i\bullet}, \quad \bar{Y}_{\bullet\bullet} = \mu + \bar{\epsilon}_{\bullet\bullet}$$

where $\bar{\epsilon}_{\bullet\bullet} = \frac{1}{N} \sum_{i=1}^k \sum_{j=1}^{n_i} \epsilon_{ij} = \frac{1}{N} \sum_{i=1}^k n_i \bar{\epsilon}_{i\bullet}$. Thus:

$$\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet} = \bar{\epsilon}_{i\bullet} - \bar{\epsilon}_{\bullet\bullet}$$

Therefore:

$$\begin{aligned} SST &= \sum_{i=1}^k n_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2 = \sum_{i=1}^k n_i (\bar{\epsilon}_{i\bullet} - \bar{\epsilon}_{\bullet\bullet})^2 \\ &= \sum_{i=1}^k n_i \bar{\epsilon}_{i\bullet}^2 - 2 \underbrace{\sum_{i=1}^k n_i \bar{\epsilon}_{\bullet\bullet}}_{=N\bar{\epsilon}_{\bullet\bullet}} + \underbrace{\sum_{i=1}^k n_i}_{=N} \bar{\epsilon}_{\bullet\bullet}^2 \\ &= \sum_{i=1}^k n_i \bar{\epsilon}_{i\bullet}^2 - N \bar{\epsilon}_{\bullet\bullet}^2 \end{aligned}$$

Since $\bar{\epsilon}_{\bullet\bullet} = \frac{1}{N} \sum_{i=1}^k \sum_{j=1}^{n_i} \epsilon_{ij}$ is the average value of ϵ_{ij} over all the observations N , we can write:

$$\mathbb{E}(\bar{\epsilon}_{\bullet\bullet}) = \frac{\sigma^2}{N}$$

Now,

$$\begin{aligned} \mathbb{E}(SST) &= \sum_{i=1}^k n_i \mathbb{E}(\bar{\epsilon}_{i\bullet}^2) - N \mathbb{E}(\bar{\epsilon}_{\bullet\bullet}^2) \\ &= \sum_{i=1}^k n_i \frac{\sigma^2}{n-i} - N \frac{\sigma^2}{N} = (k-1)\sigma^2 \end{aligned}$$

Hence:

$$\mathbb{E}[MST] = \mathbb{E}\left[\frac{SST}{k-1}\right] = \sigma^2$$

– **Case 2: Under H_1 (Not All μ_i 's are equal)**

The equivalent hypothesis in terms of α 's can be written as:

$$H_1 : \text{Not All } \alpha_i = 0$$

Under the alternative hypothesis:

$$\begin{aligned} \bar{Y}_{i\bullet} &= \mu + \alpha_i + \bar{\epsilon}_{i\bullet} \\ \bar{Y}_{\bullet\bullet} &= \mu + \bar{\alpha} + \bar{\epsilon}_{\bullet\bullet} = \mu + \bar{\epsilon}_{\bullet\bullet} \end{aligned}$$

where $\bar{\alpha} = \frac{1}{N} \sum_{i=1}^k n_i \alpha_i = 0$. Therefore:

$$\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet} = \alpha_i + (\bar{\epsilon}_{i\bullet} - \bar{\epsilon}_{\bullet\bullet})$$

Thus:

$$\begin{aligned} \text{SST} &= \sum_{i=1}^k n_i [\alpha_i + (\bar{\epsilon}_{i\bullet} - \bar{\epsilon}_{\bullet\bullet})]^2 \\ &= \sum_{i=1}^k n_i \alpha_i^2 + 2 \sum_{i=1}^k n_i \alpha_i (\bar{\epsilon}_{i\bullet} - \bar{\epsilon}_{\bullet\bullet}) + \sum_{i=1}^k n_i (\bar{\epsilon}_{i\bullet} - \bar{\epsilon}_{\bullet\bullet})^2 \end{aligned}$$

Taking expectations:

$$\begin{aligned} \mathbb{E}[\text{SST}] &= \sum_{i=1}^k n_i \alpha_i^2 + 2 \sum_{i=1}^k n_i \alpha_i \left(\underbrace{\mathbb{E}(\bar{\epsilon}_{i\bullet})}_{=0} - \underbrace{\mathbb{E}(\bar{\epsilon}_{\bullet\bullet})}_{=0} \right) + (k-1)\sigma^2 \\ &= \sum_{i=1}^k n_i \alpha_i^2 + 0 + (k-1)\sigma^2 = \sum_{i=1}^k n_i \alpha_i^2 + (k-1)\sigma^2 \\ &= \sum_{i=1}^k n_i (\mu_i - \mu)^2 + (k-1)\sigma^2 \end{aligned}$$

Therefore:

$$\mathbb{E}[\text{MST}] = \frac{\mathbb{E}[\text{SST}]}{k-1} = \sigma^2 + \frac{\sum_{i=1}^k n_i (\mu_i - \mu)^2}{k-1}$$

■

9.3.5 The F-Statistic

The above theorem tells us that the Mean Square Treatment (MST) and Mean Square Error (MSE) have the same expectation under the null hypothesis. However, under the alternative, $\mathbb{E}[\text{MST}] > \mathbb{E}[\text{MSE}]$ due to the added term involving group mean deviations.

Hence, a natural test statistic is the ratio:

$$F = \frac{\text{MST}}{\text{MSE}}$$

Intuitively:

- If all group means are equal (i.e., under H_0), both MST and MSE estimate the same error variance σ^2 , so F is expected to be close to 1.
- If the group means differ significantly, MST becomes larger, and hence F tends to exceed 1.

Under the null hypothesis:

$$F \sim F_{k-1, N-k}$$

That is, the F follows an F -distribution with:

- $k-1$ degrees of freedom in the numerator (corresponding to the number of groups minus one), and
- $N-k$ degrees of freedom in the denominator (corresponding to the total number of observations minus the number of groups).

This distribution provides the critical region for hypothesis testing. Large values of F suggest significant differences between the group means, leading us to reject the null hypothesis.

9.3.6 Hypothesis Testing Procedure

1. **State the Hypotheses**

$$H_0 : \mu_1 = \mu_2 = \cdots = \mu_k$$

H_1 : At least one μ_i differs from the others

2. **Choose the Significance Level**

Select α (commonly 0.05, 0.01, or 0.10).

3. **Compute the Test Statistic**

Calculate

$$F = \frac{\text{MST}}{\text{MSE}}$$

4. **Determine the Critical Value**

Find the critical value $F_{\alpha; k-1, N-k}$ for the given α from the F -distribution table.

5. **Make the Decision**

- If $F > F_{\alpha; k-1, N-k}$, reject H_0 .
- Otherwise, fail to reject H_0 .

Alternatively, compare the p -value to α :

- If $p\text{-value} < \alpha$, reject H_0 .
- Otherwise, fail to reject H_0 .



Figure 9.1: *Rejection region for right-tailed F -test in one-way ANOVA.*

9.3.7 ANOVA Table

The results of a one-way ANOVA are typically summarized in an ANOVA table:

9.4 Example of One-Way ANOVA

Source	SS	df	MS	F	p-value
Between Groups	SST	$k - 1$	$\frac{SST}{k - 1}$	$\frac{MST}{MSE}$	$P(F_{k-1, N-k} > F)$
Within Groups	SSE	$N - k$	$\frac{SSE}{N - k}$		
Total	TSS	$N - 1$			

Table 9.3: *One-Way ANOVA Table.*

9.4.1 Problem

A researcher wants to compare the effectiveness of three different teaching methods on student test scores. Random samples of students were assigned to each teaching method, and their test scores were recorded.

Method A	Method B	Method C
85	79	92
87	82	94
83	78	88
91	85	96
89	81	90

Table 9.4: *Test score of file students for different teaching methods.*

Test at $\alpha = 0.05$ whether there is a significant difference in mean test scores among the three teaching methods.

9.4.2 Solution

- **Calculate Sample Statistics:**

- Number of groups:

$$k = 3$$

- Number of observations in each group and total:

$$n_1 = n_2 = n_3 = 5, \quad N = 15$$

- Group totals:

$$Y_{1\bullet} = 85 + 87 + 83 + 91 + 89 = 435$$

$$Y_{2\bullet} = 79 + 82 + 78 + 85 + 81 = 405$$

$$Y_{3\bullet} = 92 + 94 + 88 + 96 + 90 = 460$$

- Grand total:

$$Y_{\bullet\bullet} = 435 + 405 + 460 = 1300$$

- **Calculate Sums of Squares:**

- Total sum of squares:

$$\begin{aligned}
 \text{TSS} &= \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}^2 - \frac{Y_{\bullet\bullet}^2}{N} \\
 &= (85^2 + 87^2 + \dots + 90^2) - \frac{(1300)^2}{15} \\
 &= 113080 - 112666.67 = 413.33
 \end{aligned}$$

- Sum of squares between treatment groups:

$$\begin{aligned} \text{SST} &= \sum_{i=1}^k \frac{Y_{i\bullet}^2}{n_i} - \frac{Y_{\bullet\bullet}^2}{N} \\ &= \left(\frac{435^2}{5} + \frac{405^2}{5} + \frac{460^2}{5} \right) - \frac{(1300)^2}{15} \\ &= 303.33 \end{aligned}$$

- Sum of squares for error:

$$\text{SSE} = \text{TSS} - \text{SST} = 413.33 - 303.33 = 110$$

- **Calculate Mean Squares:**

$$\begin{aligned} \text{MSB} &= \frac{\text{SSB}}{k-1} = \frac{303.33}{3-1} = \frac{303.33}{2} = 151.67 \\ \text{MSW} &= \frac{\text{SSW}}{N-k} = \frac{110}{15-3} = \frac{110}{12} = 9.17 \end{aligned}$$

- **Calculate F-Statistic:**

$$F = \frac{\text{MSB}}{\text{MSW}} = \frac{151.67}{9.17} = 16.54$$

- **Determine critical value and make decision:**

With $\alpha = 0.05$,

$$df_1 = k - 1 = 2, \text{ and } df_2 = N - k = 12$$

From F-table,

$$F_{0.05;2,12} = 3.89$$

Since $F = 16.54 > 3.89$, we reject H_0 .

- **ANOVA Table:**

Source	SS	df	MS	F	p-value
Between Groups	303.33	2	151.67	16.54	< 0.001
Within Groups	110.00	12	9.17		
Total	413.33	14			

Table 9.5: ANOVA Table for Teaching Methods Example.

- **Conclusion:**

At the 0.05 significance level, there is sufficient evidence to conclude that there is a significant difference in mean test scores among the three teaching methods.

9.5 Effect Size and Practical Significance

9.5.1 Eta-squared (η^2)

Eta-squared measures the proportion of total variance explained by the treatment:

$$\eta^2 = \frac{\text{SSB}}{\text{TSS}} \tag{9.1}$$

For our example:

$$\eta^2 = \frac{303.33}{413.33} = 0.734 \quad (9.2)$$

This indicates that approximately 73.4% of the variance in test scores is explained by the teaching method.

9.5.2 Omega-squared (ω^2)

Omega-squared provides a less biased estimate of effect size:

$$\omega^2 = \frac{SSB - (k - 1)MSW}{TSS + MSW} \quad (9.3)$$

For our example:

$$\omega^2 = \frac{303.33 - (2)(9.17)}{413.33 + 9.17} = \frac{284.99}{422.50} = 0.675 \quad (9.4)$$

9.6 Post-Hoc Analysis

When the F-test indicates significant differences, post-hoc tests are used to determine which specific groups differ from each other.

9.6.1 Tukey's Honestly Significant Difference (HSD)

For equal sample sizes, the critical difference is:

$$HSD = q_{\alpha, k, N-k} \sqrt{\frac{MSW}{n}} \quad (9.5)$$

where $q_{\alpha, k, N-k}$ is the critical value from the studentized range distribution.

For our example with $\alpha = 0.05$, $k = 3$, $N - k = 12$, and $n = 5$: $q_{0.05, 3, 12} = 3.77$

$$HSD = 3.77 \sqrt{\frac{9.17}{5}} = 3.77 \sqrt{1.834} = 3.77 \times 1.354 = 5.10 \quad (9.6)$$

Pairwise comparisons:

- $|\bar{X}_1 - \bar{X}_2| = |87 - 81| = 6 > 5.10$ (Significant)
- $|\bar{X}_1 - \bar{X}_3| = |87 - 92| = 5 < 5.10$ (Not significant)
- $|\bar{X}_2 - \bar{X}_3| = |81 - 92| = 11 > 5.10$ (Significant)

9.7 Two-Way ANOVA

Two-way ANOVA is an extension of one-way ANOVA that allows us to examine the effects of two categorical independent variables (called *factors*) on a continuous dependent variable simultaneously. This statistical method is particularly powerful because it can detect not only the individual effects of each factor but also their *interaction effect*.

The key questions two-way ANOVA addresses are:

1. Is there a significant main effect of Factor A?

2. Is there a significant main effect of Factor B?
3. Is there a significant interaction effect between Factor A and Factor B?

This can be stated formally using hypotheses:

$$\begin{aligned} \text{Main effect of Factor A: } H_{0A} : \alpha_1 = \alpha_2 = \dots = \alpha_a = 0 \\ \text{Main effect of Factor B: } H_{0B} : \beta_1 = \beta_2 = \dots = \beta_b = 0 \\ \text{Interaction effect: } H_{0AB} : (\alpha\beta)_{ij} = 0 \text{ for all } i, j \end{aligned}$$

where α_i represents the effect of the i -th level of Factor A, β_j represents the effect of the j -th level of Factor B, and $(\alpha\beta)_{ij}$ represents the interaction effect between the i -th level of Factor A and the j -th level of Factor B.

Consider an experiment to study the effect of both teaching method (Factor A: Methods 1, 2, 3) and class size (Factor B: Small, Medium, Large) on student test scores:

Teaching Method	Small Class	Medium Class	Large Class
Method 1	85, 87, 83	80, 82, 78	75, 77, 73
Method 2	90, 92, 88	85, 87, 83	80, 82, 78
Method 3	88, 90, 86	88, 90, 86	87, 89, 85

Table 9.6: *Test scores for different teaching methods and class sizes.*

Two-way ANOVA allows us to determine:

- Whether teaching method affects test scores (main effect of Factor A)
- Whether class size affects test scores (main effect of Factor B)
- Whether the effect of teaching method depends on class size (interaction effect)

9.7.1 Advantages of Two-Way ANOVA

- **Efficiency:** Examines multiple factors simultaneously, reducing the number of separate experiments needed.
- **Interaction Detection:** Can identify synergistic or antagonistic effects between factors.
- **Control for Confounding:** Accounts for the effect of a second factor, providing a clearer picture of each factor's individual impact.
- **Reduced Error:** By explaining more variation in the data, it often provides more precise estimates of effects.

9.8 Mathematical Model for Two-Way ANOVA

The two-way ANOVA model with interaction can be expressed as:

$$Y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk}$$

where:

- Y_{ijk} is the k -th observation in the i -th level of Factor A and j -th level of Factor B
- μ is the overall (grand) mean
- α_i is the main effect of the i -th level of Factor A ($i = 1, 2, \dots, a$)
- β_j is the main effect of the j -th level of Factor B ($j = 1, 2, \dots, b$)
- $(\alpha\beta)_{ij}$ is the interaction effect between the i -th level of Factor A and j -th level of Factor B
- ϵ_{ijk} is the random error term ($k = 1, 2, \dots, n$, where n is the number of replications per cell)

9.8.1 Constraints on the Model

To ensure unique parameter estimates, the following constraints are imposed:

- The main effects of Factor A are centered around zero.

$$\sum_{i=1}^a \alpha_i = 0$$

- The main effects of Factor B are centered around zero.

$$\sum_{j=1}^b \beta_j = 0$$

- For each level j of Factor B, the interaction effects over all levels of Factor A average to zero.

$$\sum_{i=1}^a (\alpha\beta)_{ij} = 0 \quad \text{for all } j$$

- For each level i of Factor A, the interaction effects over all levels of Factor B average to zero.

$$\sum_{j=1}^b (\alpha\beta)_{ij} = 0 \quad \text{for all } i$$

9.8.2 Key Assumptions

For valid inference using two-way ANOVA, the following assumptions must be satisfied:

1. **Independence:** All observations must be independent of each other.
2. **Normality:** The error terms are normally distributed:

$$\epsilon_{ijk} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

3. **Homogeneity of Variances (Homoscedasticity):** The variance is constant across all treatment combinations:

$$\text{Var}(Y_{ijk}) = \sigma^2 \quad \text{for all } i, j, k$$

4. **Additivity:** The effects of the factors are additive (this assumption is relaxed when interaction terms are included in the model).

9.9 Analysis of Two-Way Classified Data

9.9.1 Data Structure

Consider a two-way factorial experiment with:

- Factor A having a levels
- Factor B having b levels
- n replications in each cell
- Total number of observations: $N = abn$

The data can be organized as shown in Table 9.7:

Factor A	Factor B				Row Total
	B_1	B_2	\cdots	B_b	
A_1	Y_{11k}	Y_{12k}	\cdots	Y_{1bk}	$Y_{1\bullet\bullet}$
A_2	Y_{21k}	Y_{22k}	\cdots	Y_{2bk}	$Y_{2\bullet\bullet}$
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
A_a	Y_{a1k}	Y_{a2k}	\cdots	Y_{abk}	$Y_{a\bullet\bullet}$
Column Total	$Y_{\bullet 1\bullet}$	$Y_{\bullet 2\bullet}$	\cdots	$Y_{\bullet b\bullet}$	$Y_{\bullet\bullet\bullet}$

Table 9.7: Data structure for two-way ANOVA.

9.9.2 Notation and Sample Statistics

- Cell totals and means:

$$Y_{ij\bullet} = \sum_{k=1}^n Y_{ijk} \quad (\text{total for cell } (i, j))$$

$$\bar{Y}_{ij\bullet} = \frac{Y_{ij\bullet}}{n} = \frac{1}{n} \sum_{k=1}^n Y_{ijk} \quad (\text{mean for cell } (i, j))$$

- Row totals and means (Factor A):

$$Y_{i\bullet\bullet} = \sum_{j=1}^b Y_{ij\bullet} = \sum_{j=1}^b \sum_{k=1}^n Y_{ijk} \quad (\text{total for level } i \text{ of Factor A})$$

$$\bar{Y}_{i\bullet\bullet} = \frac{Y_{i\bullet\bullet}}{bn} = \frac{1}{bn} \sum_{j=1}^b \sum_{k=1}^n Y_{ijk} \quad (\text{mean for level } i \text{ of Factor A})$$

- Column totals and means (Factor B):

$$Y_{\bullet j\bullet} = \sum_{i=1}^a Y_{ij\bullet} = \sum_{i=1}^a \sum_{k=1}^n Y_{ijk} \quad (\text{total for level } j \text{ of Factor B})$$

$$\bar{Y}_{\bullet j\bullet} = \frac{Y_{\bullet j\bullet}}{an} = \frac{1}{an} \sum_{i=1}^a \sum_{k=1}^n Y_{ijk} \quad (\text{mean for level } j \text{ of Factor B})$$

- Grand total and mean:

$$Y_{\bullet\bullet\bullet} = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^n Y_{ijk} \quad (\text{grand total})$$

$$\bar{Y}_{\bullet\bullet\bullet} = \frac{Y_{\bullet\bullet\bullet}}{abn} = \frac{1}{abn} \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^n Y_{ijk} \quad (\text{grand mean})$$

9.9.3 Sum of Squares Decomposition

The total variation in two-way ANOVA is partitioned into four components:

$$\text{TSS} = \text{SSA} + \text{SSB} + \text{SSAB} + \text{SSE}$$

where:

- Total Sum of Squares (TSS):

$$\text{TSS} = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^n (Y_{ijk} - \bar{Y}_{\bullet\bullet\bullet})^2$$

- Sum of Squares for Factor A (SSA):

$$\text{SSA} = bn \sum_{i=1}^a (\bar{Y}_{i\bullet\bullet} - \bar{Y}_{\bullet\bullet\bullet})^2$$

- Sum of Squares for Factor B (SSB):

$$\text{SSB} = an \sum_{j=1}^b (\bar{Y}_{\bullet j\bullet} - \bar{Y}_{\bullet\bullet\bullet})^2$$

- Sum of Squares for Interaction (SSAB):

$$\text{SSAB} = n \sum_{i=1}^a \sum_{j=1}^b (\bar{Y}_{ij\bullet} - \bar{Y}_{i\bullet\bullet} - \bar{Y}_{\bullet j\bullet} + \bar{Y}_{\bullet\bullet\bullet})^2$$

- Sum of Squares for Error (SSE):

$$\text{SSE} = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^n (Y_{ijk} - \bar{Y}_{ij\bullet})^2$$

9.9.4 Computational Formulas

For computational efficiency, the following formulas are preferred:

$$\begin{aligned}
\text{TSS} &= \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^n Y_{ijk}^2 - \frac{Y_{\dots}^2}{abn} \\
\text{SSA} &= \frac{1}{bn} \sum_{i=1}^a Y_{i\bullet\bullet}^2 - \frac{Y_{\dots}^2}{abn} \\
\text{SSB} &= \frac{1}{an} \sum_{j=1}^b Y_{\bullet j\bullet}^2 - \frac{Y_{\dots}^2}{abn} \\
\text{SSAB} &= \frac{1}{n} \sum_{i=1}^a \sum_{j=1}^b Y_{ij\bullet}^2 - \frac{Y_{\dots}^2}{abn} - \text{SSA} - \text{SSB} \\
\text{SSE} &= \text{TSS} - \text{SSA} - \text{SSB} - \text{SSAB}
\end{aligned}$$

9.9.5 Degrees of Freedom

Each sum of squares has associated degrees of freedom:

- $df_{\text{TSS}} = abn - 1$ (total observations minus 1)
- $df_{\text{SSA}} = a - 1$ (levels of Factor A minus 1)
- $df_{\text{SSB}} = b - 1$ (levels of Factor B minus 1)
- $df_{\text{SSAB}} = (a - 1)(b - 1)$ (interaction degrees of freedom)
- $df_{\text{SSE}} = ab(n - 1)$ (error degrees of freedom)

These satisfy the identity:

$$df_{\text{TSS}} = df_{\text{SSA}} + df_{\text{SSB}} + df_{\text{SSAB}} + df_{\text{SSE}}$$

9.9.6 Mean Squares

Mean squares are calculated by dividing sums of squares by their corresponding degrees of freedom:

$$\begin{aligned}
\text{MSA} &= \frac{\text{SSA}}{a - 1} \\
\text{MSB} &= \frac{\text{SSB}}{b - 1} \\
\text{MSAB} &= \frac{\text{SSAB}}{(a - 1)(b - 1)} \\
\text{MSE} &= \frac{\text{SSE}}{ab(n - 1)}
\end{aligned}$$

9.9.7 F-Statistics

Three F-statistics are computed to test the respective hypotheses:

$$F_A = \frac{MSA}{MSE} \quad (\text{test for main effect of Factor A})$$

$$F_B = \frac{MSB}{MSE} \quad (\text{test for main effect of Factor B})$$

$$F_{AB} = \frac{MSAB}{MSE} \quad (\text{test for interaction effect})$$

Under the null hypotheses, these F-statistics follow F-distributions:

$$F_A \sim F_{a-1, ab(n-1)}$$

$$F_B \sim F_{b-1, ab(n-1)}$$

$$F_{AB} \sim F_{(a-1)(b-1), ab(n-1)}$$

9.9.8 Hypothesis Testing

Test for Main Effect of Factor A

$$H_{0A} : \alpha_1 = \alpha_2 = \cdots = \alpha_a = 0$$

$$H_{1A} : \text{At least one } \alpha_i \neq 0$$

Test Statistic:

$$F_A = \frac{MSA}{MSE}$$

Decision Rule: Reject H_{0A} if $F_A > F_{\alpha; a-1, ab(n-1)}$

Test for Main Effect of Factor B

$$H_{0B} : \beta_1 = \beta_2 = \cdots = \beta_b = 0$$

$$H_{1B} : \text{At least one } \beta_j \neq 0$$

Test Statistic:

$$F_B = \frac{MSB}{MSE}$$

Decision Rule: Reject H_{0B} if $F_B > F_{\alpha; b-1, ab(n-1)}$

Test for Interaction Effect

$$H_{0AB} : (\alpha\beta)_{ij} = 0 \text{ for all } i, j$$

$$H_{1AB} : \text{At least one } (\alpha\beta)_{ij} \neq 0$$

Test Statistic:

$$F_{AB} = \frac{MSAB}{MSE}$$

Decision Rule: Reject H_{0AB} if $F_{AB} > F_{\alpha; (a-1)(b-1), ab(n-1)}$

9.9.9 Two-Way ANOVA Table

The results are typically summarized in the following ANOVA table:

Source	SS	df	MS	F
Factor A	SSA	$a - 1$	$\frac{SSA}{a - 1}$	$\frac{MSA}{MSE}$
Factor B	SSB	$b - 1$	$\frac{SSB}{b - 1}$	$\frac{MSB}{MSE}$
A \times B	SSAB	$(a - 1)(b - 1)$	$\frac{SSAB}{(a - 1)(b - 1)}$	$\frac{MSAB}{MSE}$
Error	SSE	$ab(n - 1)$	$\frac{SSE}{ab(n - 1)}$	
Total	TSS	$abn - 1$		

Table 9.8: *Two-Way ANOVA Table.*

9.9.10 Interpretation of Results

Main Effects

- **Significant Main Effect of Factor A:** The means of Factor A differ significantly across its levels, averaging over all levels of Factor B.
- **Significant Main Effect of Factor B:** The means of Factor B differ significantly across its levels, averaging over all levels of Factor A.

Interaction Effect

- **Significant Interaction:** The effect of one factor depends on the level of the other factor. This means:
 - The effect of Factor A is not the same at all levels of Factor B
 - The effect of Factor B is not the same at all levels of Factor A
 - Simple main effects should be examined instead of overall main effects
- **Non-significant Interaction:** The factors act independently. Main effects can be interpreted directly.

9.10 Example of Two-Way ANOVA

9.10.1 Problem

A researcher wants to study the effect of teaching method (Factor A: Methods 1, 2, 3) and class size (Factor B: Small, Large) on student test scores. Two students were randomly assigned to each combination of teaching method and class size.

Teaching Method	Small Class	Large Class
Method 1	85, 87	78, 80
Method 2	90, 92	83, 85
Method 3	88, 90	86, 88

Table 9.9: *Test scores for teaching methods and class sizes.*

Test at $\alpha = 0.05$ whether there are significant main effects and interaction effect.

9.10.2 Solution

Setup:

- $a = 3$ (teaching methods), $b = 2$ (class sizes), $n = 2$ (replications)
- $N = abn = 3 \times 2 \times 2 = 12$ total observations

Calculate Cell, Row, Column, and Grand Totals:

Method	Small	Large	Row Total
1	172	158	330
2	182	168	350
3	178	174	352
Column Total	532	500	1032

Table 9.10: *Cell, row, column, and grand totals.*

Calculate Cell Means:

$$\begin{aligned}
 \bar{Y}_{11\bullet} &= \frac{172}{2} = 86.0 & \bar{Y}_{12\bullet} &= \frac{158}{2} = 79.0 \\
 \bar{Y}_{21\bullet} &= \frac{182}{2} = 91.0 & \bar{Y}_{22\bullet} &= \frac{168}{2} = 84.0 \\
 \bar{Y}_{31\bullet} &= \frac{178}{2} = 89.0 & \bar{Y}_{32\bullet} &= \frac{174}{2} = 87.0
 \end{aligned}$$

Calculate Row and Column Means:

$$\begin{aligned}
 \bar{Y}_{1\bullet\bullet} &= \frac{330}{4} = 82.5 & \bar{Y}_{\bullet 1\bullet} &= \frac{532}{6} = 88.67 \\
 \bar{Y}_{2\bullet\bullet} &= \frac{350}{4} = 87.5 & \bar{Y}_{\bullet 2\bullet} &= \frac{500}{6} = 83.33 \\
 \bar{Y}_{3\bullet\bullet} &= \frac{352}{4} = 88.0 & \bar{Y}_{\bullet \bullet\bullet} &= \frac{1032}{12} = 86.0
 \end{aligned}$$

Calculate Sums of Squares:

$$\begin{aligned}
 \text{TSS} &= (85^2 + 87^2 + \cdots + 88^2) - \frac{1032^2}{12} \\
 &= 89048 - 88,776 = 272
 \end{aligned}$$

$$\begin{aligned}
 \text{SSA} &= \frac{1}{4}(330^2 + 350^2 + 352^2) - \frac{1032^2}{12} \\
 &= \frac{1}{4}(378,404) - 88,776 = 94,601 - 88,776 = 94
 \end{aligned}$$

$$\begin{aligned}
 \text{SSB} &= \frac{1}{6}(532^2 + 500^2) - \frac{1032^2}{12} \\
 &= \frac{1}{6}(533,024) - 88,776 = 88,837.33 - 88,776 = 85.33
 \end{aligned}$$

$$\begin{aligned} \text{SSAB} &= \frac{1}{2}(172^2 + 158^2 + 182^2 + 168^2 + 178^2 + 174^2) - \frac{1032^2}{12} - 94 - 85.33 \\ &= \frac{1}{2}(178,396) - 88,776 - 179.33 = 89,198 - 88,955.33 = 6.67 \end{aligned}$$

$$\text{SSE} = 272 - 94 - 85.33 - 6.67 = 86$$

Calculate Mean Squares:

$$\begin{aligned} \text{MSA} &= \frac{94}{2} = 47 \\ \text{MSB} &= \frac{85.33}{1} = 85.33 \\ \text{MSAB} &= \frac{6.67}{2} = 3.34 \\ \text{MSE} &= \frac{86}{6} = 14.33 \end{aligned}$$

Calculate F-Statistics:

$$\begin{aligned} F_A &= \frac{47}{14.33} = 3.28 \\ F_B &= \frac{85.33}{14.33} = 5.95 \\ F_{AB} &= \frac{3.34}{14.33} = 0.23 \end{aligned}$$

Critical Values at $\alpha = 0.05$:

$$\begin{aligned} F_{0.05;2,6} &= 5.14 \quad (\text{for Factor A}) \\ F_{0.05;1,6} &= 5.99 \quad (\text{for Factor B}) \\ F_{0.05;2,6} &= 5.14 \quad (\text{for Interaction}) \end{aligned}$$

ANOVA Table:

Source	SS	df	MS	F	p-value
Teaching Method (A)	94.00	2	47.00	3.28	> 0.05
Class Size (B)	85.33	1	85.33	5.95	< 0.05
A \times B	6.67	2	3.34	0.23	> 0.05
Error	86.00	6	14.33		
Total	272.00	11			

Table 9.11: ANOVA table for teaching methods and class sizes example.

Conclusions:

- **Teaching Method (Factor A):** $F_A = 3.28 < 5.14$, so we fail to reject H_{0A} . There is no significant difference between teaching methods.
- **Class Size (Factor B):** $F_B = 5.95 < 5.99$ (marginally), so we fail to reject H_{0B} , but it's very close to significance.
- **Interaction:** $F_{AB} = 0.23 < 5.14$, so we fail to reject H_{0AB} . There is no significant interaction between teaching method and class size.

At the 0.05 significance level, only class size shows a marginally significant effect on test scores, with small classes generally producing higher scores than large classes.

9.11 Key Formulas Summary

Table 9.12: Summary of Key ANOVA Formulas

Statistic	Formula
Group Mean	$\bar{X}_{i.} = \frac{1}{n_i} \sum_{j=1}^{n_i} X_{ij}$
Grand Mean	$\bar{X}_{..} = \frac{1}{N} \sum_{i=1}^k \sum_{j=1}^{n_i} X_{ij}$
TSS	$\sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_{..})^2$
SSB	$\sum_{i=1}^k n_i (\bar{X}_{i.} - \bar{X}_{..})^2$
SSW	$\sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_{i.})^2$
MSB	$\frac{SSB}{k-1}$
MSW	$\frac{SSW}{N-k}$
F-statistic	$\frac{MSB}{MSW}$
Eta-squared	$\frac{SSB}{TSS}$

Chapter 10

Introduction to Design of Experiments

10.1 What is Design of Experiments?

Design of Experiments (DOE) is a systematic approach to planning, conducting, and analyzing experiments to obtain maximum information with minimum resources. Consider a simple medical example. Suppose a pharmaceutical company wants to test if a new drug is effective in reducing blood pressure. They could simply give the drug to patients and measure results, but this approach has problems:

- No comparison group (control)
- Patient differences may affect results
- Other factors (diet, exercise, age) may influence outcomes
- Results may not be reliable or generalizable

A properly designed experiment would:

- Include a control group receiving placebo
- Randomly assign patients to drug or placebo groups
- Control for patient characteristics (age, gender, severity)
- Use sufficient sample size for reliable conclusions

10.1.1 Why Design Experiments?

Proper experimental design ensures:

1. **Valid Results:** Conclusions are statistically sound
2. **Efficiency:** Maximum information with minimum resources
3. **Reliability:** Results can be reproduced
4. **Generalizability:** Findings apply to broader populations

10.1.2 Fundamental Principles

The three cornerstone principles of experimental design are:

1. Randomization

Random assignment of experimental units to treatments eliminates systematic bias.

Example: In our drug study, patients are randomly assigned to receive either the new drug or placebo. This ensures that factors like age, gender, and disease severity are equally distributed between groups.

2. Replication

Multiple observations under the same conditions provide:

- Estimate of experimental error
- Increased precision of results
- Ability to detect smaller effects

Example: Instead of testing the drug on just one patient per group, we test it on multiple patients (say 20 per group) to get reliable results.

3. Blocking

Grouping similar experimental units controls for known sources of variation.

Example: If we suspect that age affects drug response, we can group patients into age blocks (20-40, 40-60, 60+ years) and ensure both drug and placebo are tested within each age group.

10.2 Planning an Experiment

10.2.1 Step 1: Define the Problem

Clearly specify:

- **Response Variable:** What you want to measure (e.g., blood pressure reduction)
- **Factors:** Variables that might affect the response (e.g., drug type, dosage)
- **Levels:** Values each factor can take (e.g., 10mg, 20mg, 30mg dosage)
- **Objectives:** What questions you want to answer

10.2.2 Step 2: Choose Factors and Levels

Types of Factors:

- **Quantitative:** Numerical values (dosage, temperature, time)
- **Qualitative:** Categories (drug A vs drug B, male vs female)

Factor Classification:

- **Fixed Effects:** Specific levels chosen deliberately (specific drugs)
- **Random Effects:** Levels randomly selected from population (random patients)

10.2.3 Step 3: Select Experimental Design

Choose the most appropriate design based on:

- Number of factors
- Available resources
- Sources of variation to control
- Research objectives

10.3 Types of Experimental Designs

10.3.1 Classification by Structure

Experimental designs can be classified into four main types:

1. **Completely Randomized Design (CRD)**
2. **Randomized Block Design (RBD)**
3. **Latin Square Design (LSD)**
4. **Factorial Design**

10.4 Completely Randomized Design (CRD)

10.4.1 Description

The simplest experimental design where treatments are assigned completely at random to experimental units.

10.4.2 When to Use CRD

- Experimental units are homogeneous
- No obvious sources of systematic variation
- Complete randomization is feasible
- Resources are limited

10.4.3 Example: Drug Effectiveness Study

Problem: Compare effectiveness of three treatments: Drug A, Drug B, and Placebo on reducing blood pressure.

Design:

- 30 patients with similar blood pressure levels
- Randomly assign 10 patients to each treatment
- Measure blood pressure reduction after 4 weeks

Table 10.1: CRD Layout for Drug Study

Drug A	Drug B	Placebo
Patient 3	Patient 1	Patient 2
Patient 7	Patient 4	Patient 5
Patient 9	Patient 6	Patient 8
\vdots	\vdots	\vdots

10.4.4 Statistical Model

$$y_{ij} = \mu + \tau_i + \epsilon_{ij} \quad (10.1)$$

where:

- y_{ij} = response of j th unit in i th treatment
- μ = overall mean
- τ_i = effect of i th treatment
- ϵ_{ij} = random error

10.4.5 ANOVA for CRD

$$SS_{Total} = SS_{Treatment} + SS_{Error} \quad (10.2)$$

$$SS_{Total} = \sum_{i=1}^a \sum_{j=1}^{n_i} y_{ij}^2 - \frac{T^2}{N} \quad (10.3)$$

$$SS_{Treatment} = \sum_{i=1}^a \frac{T_i^2}{n_i} - \frac{T^2}{N} \quad (10.4)$$

$$SS_{Error} = SS_{Total} - SS_{Treatment} \quad (10.5)$$

Table 10.2: ANOVA Table for CRD

Source	df	SS	MS	F
Treatment	$a - 1$	$SS_{Treatment}$	$MS_{Treatment}$	$\frac{MS_{Treatment}}{MS_{Error}}$
Error	$N - a$	SS_{Error}	MS_{Error}	
Total	$N - 1$	SS_{Total}		

10.5 Randomized Block Design (RBD)

10.5.1 Description

Used when experimental units are not homogeneous and there's a known source of variation that can be controlled through blocking.

10.5.2 When to Use RBD

- Experimental units are heterogeneous
- There's an identifiable source of variation to control
- Want to increase precision of treatment comparisons

10.5.3 Example: Drug Study with Age Groups

Problem: Same drug effectiveness study, but patients vary significantly in age, which affects drug response.

Design:

- Create 4 age blocks: 20-30, 30-40, 40-50, 50-60 years
- Within each block, randomly assign patients to Drug A, Drug B, and Placebo
- Each treatment appears once in each block

Table 10.3: RBD Layout for Drug Study

Age Block	Drug A	Drug B	Placebo
20-30 years	Patient 3	Patient 7	Patient 12
30-40 years	Patient 15	Patient 2	Patient 8
40-50 years	Patient 9	Patient 18	Patient 5
50-60 years	Patient 21	Patient 11	Patient 14

10.5.4 Statistical Model

$$y_{ij} = \mu + \tau_i + \beta_j + \epsilon_{ij} \quad (10.6)$$

where:

- β_j = effect of j th block
- Other terms as in CRD

10.5.5 ANOVA for RBD

$$SS_{Total} = SS_{Treatment} + SS_{Block} + SS_{Error} \quad (10.7)$$

Table 10.4: ANOVA Table for RBD

Source	df	SS	MS	F
Treatment	$t - 1$	$SS_{Treatment}$	$MS_{Treatment}$	$\frac{MS_{Treatment}}{MS_{Error}}$
Block	$b - 1$	SS_{Block}	MS_{Block}	$\frac{MS_{Block}}{MS_{Error}}$
Error	$(t - 1)(b - 1)$	SS_{Error}	MS_{Error}	
Total	$tb - 1$	SS_{Total}		

10.6 Latin Square Design (LSD)

10.6.1 Description

Used when there are two sources of variation to be controlled simultaneously, and the number of levels for each source equals the number of treatments.

10.6.2 When to Use LSD

- Two identifiable sources of variation
- Number of treatments = number of levels of each blocking factor
- Want to control both sources of variation

10.6.3 Example: Drug Study with Age and Gender

Problem: Test effectiveness of 4 drug formulations, controlling for both age groups and gender.

Design:

- 4 treatments: Drug A, Drug B, Drug C, Placebo
- 4 age groups (rows): 20-30, 30-40, 40-50, 50-60
- 4 gender/severity combinations (columns)
- Each treatment appears once in each row and column

Table 10.5: Latin Square Design Layout

Age Group	Gender/Severity Combination			
	1	2	3	4
20-30	A	B	C	D
30-40	B	C	D	A
40-50	C	D	A	B
50-60	D	A	B	C

10.6.4 Statistical Model

$$y_{ijk} = \mu + \alpha_i + \beta_j + \tau_k + \epsilon_{ijk} \quad (10.8)$$

where:

- α_i = effect of i th row
- β_j = effect of j th column
- τ_k = effect of k th treatment

10.6.5 ANOVA for LSD

$$SS_{Total} = SS_{Row} + SS_{Column} + SS_{Treatment} + SS_{Error} \quad (10.9)$$

Table 10.6: ANOVA Table for Latin Square Design

Source	df	SS	MS	F
Rows	$p - 1$	SS_{Row}	MS_{Row}	$\frac{MS_{Row}}{MS_{Error}}$
Columns	$p - 1$	SS_{Column}	MS_{Column}	$\frac{MS_{Column}}{MS_{Error}}$
Treatments	$p - 1$	$SS_{Treatment}$	$MS_{Treatment}$	$\frac{MS_{Treatment}}{MS_{Error}}$
Error	$(p - 2)(p - 1)$	SS_{Error}	MS_{Error}	
Total	$p^2 - 1$	SS_{Total}		

10.7 Factorial Designs

10.7.1 Introduction

Factorial designs study multiple factors simultaneously, allowing examination of:

- **Main Effects:** Individual effect of each factor
- **Interaction Effects:** Combined effect of factors
- **Efficiency:** More information per experimental run

10.7.2 Advantages of Factorial Design

1. Study multiple factors in single experiment
2. Detect interactions between factors
3. More efficient than one-factor-at-a-time experiments
4. Results apply over wider range of conditions

10.8 Two-Factor Factorial Design (2^2 Design)

10.8.1 Description

Studies two factors, each at two levels, requiring $2 \times 2 = 4$ treatment combinations.

10.8.2 Example: Drug Dosage and Timing Study

Problem: Study effect of drug dosage (Low/High) and timing (Morning/Evening) on blood pressure reduction.

Factors:

- Factor A (Dosage): Low (10mg), High (20mg)
- Factor B (Timing): Morning, Evening

Treatment Combinations:

1. a_1b_1 : Low dosage, Morning
2. a_1b_2 : Low dosage, Evening
3. a_2b_1 : High dosage, Morning
4. a_2b_2 : High dosage, Evening

Table 10.7: 2^2 Factorial Design Layout

Factor A (Dosage)	Factor B (Timing)	
	Morning (b_1)	Evening (b_2)
Low (a_1)	a_1b_1	a_1b_2
High (a_2)	a_2b_1	a_2b_2

10.8.3 Effects in 2^2 Design

Main Effect of Factor A (Dosage):

$$\text{Main Effect of A} = \frac{(a_2b_1 + a_2b_2) - (a_1b_1 + a_1b_2)}{2} \quad (10.10)$$

Main Effect of Factor B (Timing):

$$\text{Main Effect of B} = \frac{(a_1b_2 + a_2b_2) - (a_1b_1 + a_2b_1)}{2} \quad (10.11)$$

Interaction Effect AB:

$$\text{AB Interaction} = \frac{(a_1b_1 + a_2b_2) - (a_1b_2 + a_2b_1)}{2} \quad (10.12)$$

10.8.4 Interpretation of Interaction

- **No Interaction:** Effect of one factor is same regardless of level of other factor
- **Interaction Present:** Effect of one factor depends on level of other factor

10.8.5 Statistical Model for 2^2 Design

$$y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk} \quad (10.13)$$

where:

- $(\alpha\beta)_{ij}$ = interaction effect between factors A and B

10.9 Three-Factor Factorial Design (2^3 Design)

10.9.1 Description

Studies three factors, each at two levels, requiring $2 \times 2 \times 2 = 8$ treatment combinations.

10.9.2 Example: Drug, Dosage, and Exercise Study

Problem: Study effects of drug type (A/B), dosage (Low/High), and exercise (No/Yes) on blood pressure.

Factors:

- Factor A (Drug): Type A, Type B
- Factor B (Dosage): Low, High
- Factor C (Exercise): No, Yes

Treatment Combinations:

1. $a_1b_1c_1$: Drug A, Low dosage, No exercise
2. $a_1b_1c_2$: Drug A, Low dosage, Exercise

3. $a_1b_2c_1$: Drug A, High dosage, No exercise
4. $a_1b_2c_2$: Drug A, High dosage, Exercise
5. $a_2b_1c_1$: Drug B, Low dosage, No exercise
6. $a_2b_1c_2$: Drug B, Low dosage, Exercise
7. $a_2b_2c_1$: Drug B, High dosage, No exercise
8. $a_2b_2c_2$: Drug B, High dosage, Exercise

10.9.3 Effects in 2^3 Design

The 2^3 design allows estimation of:

- 3 main effects: A, B, C
- 3 two-factor interactions: AB, AC, BC
- 1 three-factor interaction: ABC

Main Effects:

$$\text{Main Effect of A} = \frac{1}{4}[\text{sum of responses with } a_2 - \text{sum of responses with } a_1] \quad (10.14)$$

$$\text{Main Effect of B} = \frac{1}{4}[\text{sum of responses with } b_2 - \text{sum of responses with } b_1] \quad (10.15)$$

$$\text{Main Effect of C} = \frac{1}{4}[\text{sum of responses with } c_2 - \text{sum of responses with } c_1] \quad (10.16)$$

10.9.4 Statistical Model for 2^3 Design

$$\begin{aligned} y_{ijkl} = & \mu + \alpha_i + \beta_j + \gamma_k \\ & + (\alpha\beta)_{ij} + (\alpha\gamma)_{ik} + (\beta\gamma)_{jk} \\ & + (\alpha\beta\gamma)_{ijk} + \epsilon_{ijkl} \end{aligned} \quad (10.17)$$

10.10 Analysis of Variance for Factorial Designs

10.10.1 ANOVA Table for 2^2 Design

Table 10.8: ANOVA Table for 2^2 Factorial Design

Source	df	SS	MS	F
Factor A	1	SS_A	MS_A	$\frac{MS_A}{MS_E}$
Factor B	1	SS_B	MS_B	$\frac{MS_B}{MS_E}$
Interaction AB	1	SS_{AB}	MS_{AB}	$\frac{MS_{AB}}{MS_E}$
Error	$ab(n-1)$	SS_E	MS_E	
Total	$abn-1$	SS_T		

Table 10.9: ANOVA Table for 2^3 Factorial Design

Source	df	SS	MS	F
Factor A	1	SS_A	MS_A	$\frac{MS_A}{MS_E}$
Factor B	1	SS_B	MS_B	$\frac{MS_B}{MS_E}$
Factor C	1	SS_C	MS_C	$\frac{MS_C}{MS_E}$
AB Interaction	1	SS_{AB}	MS_{AB}	$\frac{MS_{AB}}{MS_E}$
AC Interaction	1	SS_{AC}	MS_{AC}	$\frac{MS_{AC}}{MS_E}$
BC Interaction	1	SS_{BC}	MS_{BC}	$\frac{MS_{BC}}{MS_E}$
ABC Interaction	1	SS_{ABC}	MS_{ABC}	$\frac{MS_{ABC}}{MS_E}$
Error	$abc(n-1)$	SS_E	MS_E	
Total	$abcn-1$	SS_T		

10.10.2 ANOVA Table for 2^3 Design

10.11 Assumptions and Model Checking

10.11.1 ANOVA Assumptions

For valid ANOVA results, the following assumptions must hold:

1. **Normality:** Residuals are normally distributed
2. **Independence:** Observations are independent
3. **Homoscedasticity:** Equal variances across treatments
4. **Additivity:** Treatment and block effects are additive

10.11.2 Checking Assumptions

- **Normality:** Normal probability plots, Shapiro-Wilk test
- **Independence:** Randomization, knowledge of experimental procedure
- **Equal Variances:** Residual plots, Levene's test
- **Additivity:** Residual analysis, Tukey's test for non-additivity

10.12 Summary

10.12.1 Design Selection Guide

Table 10.10: Guide for Selecting Experimental Design

Design	Use When	Controls For
CRD	Homogeneous units	Nothing
RBD	One source of variation	One blocking factor
LSD	Two sources of variation	Two blocking factors
Factorial	Multiple factors of interest	Factor interactions

10.12.2 Key Principles

1. Always randomize treatment assignments
2. Include adequate replication
3. Control known sources of variation through blocking
4. Consider factorial designs when studying multiple factors
5. Check model assumptions before drawing conclusions

10.12.3 Benefits of Proper Design

- Valid and reliable conclusions
- Efficient use of resources
- Ability to detect important effects
- Control of experimental error
- Broader applicability of results

Design of Experiments provides a systematic framework for conducting efficient and reliable experiments. By following the principles of randomization, replication, and blocking, and choosing appropriate designs, researchers can obtain maximum information while minimizing costs and ensuring valid conclusions.

Chapter 11

Correlation and Regression

11.1 Correlation

In statistics, **correlation** between two random variables is a measure of the degree of linear association between them.

For example, imagine you record how many hours each of the six students studies in a week and their corresponding exam scores:

Student	Hours Studied (X)	Exam Score (Y)
1	2	65
2	3	70
3	4	76
4	5	78
5	6	82
6	7	89

Table 11.1: *Study hours and exam scores of six students.*



Figure 11.1: *Scatter plot of study time and exam performance.*

As illustrated in the table and the accompanying scatter plot, exam scores increase as students dedicate more hours to studying. This clear upward pattern reflects a positive correlation: students who invest more time in preparation tend to achieve higher marks. In general, the direction of correlation can be classified as follows:

1. **Positive correlation:** An increase in one variable is accompanied by an increase in the other.

2. **Negative correlation:** An increase in one variable is accompanied by a decrease in the other.
3. **No (zero) correlation:** Changes in one variable show no consistent association with changes in the other.

11.2 Pearson Correlation Coefficient

Pearson's correlation coefficient for a population, commonly denoted by the Greek letter ρ (rho), is also known as the **population correlation coefficient**. Given a pair of random variables X and Y , the population correlation coefficient ρ_{XY} is defined as

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\mathbb{E}[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$

where $\mu_X = \mathbb{E}[X]$ and $\mu_Y = \mathbb{E}[Y]$ are the means of X and Y , σ_X and σ_Y are their standard deviations, and $\text{Cov}(X, Y)$ is the covariance between X and Y .

Like all population parameters, The value of ρ_{XY} is not known to us. We may need to estimate it from the random sample observation pairs (X, Y) . Let

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

be n pairs of observations with respective means \bar{X}, \bar{Y} and variances S_X^2, S_Y^2 respectively. It turns out that a point estimate of $\text{Cov}(X, Y)$ is the sample covariance:

$$S_{XY} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

The point estimate of σ_X is sample standard deviation of X :

$$S_X = \sqrt{\frac{1}{n - 1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

The point estimate of σ_Y is sample standard deviation of Y :

$$S_Y = \sqrt{\frac{1}{n - 1} \sum_{i=1}^n (y_i - \bar{y})^2}$$

Substituting these estimates for their population counterparts, we get the formula for the **sample correlation coefficient** as

$$r_{XY} = \frac{S_{XY}}{S_X \cdot S_Y}$$

Alternatively, r_{XY} can be expressed as

$$r_{XY} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

We can simplify the formula for r_{XY} into a form most convenient for computing from raw data by

using the identities:

$$\begin{aligned}\sum_i (x_i - \bar{x})(y_i - \bar{y}) &= \sum_i x_i y_i - n\bar{x} \cdot \bar{y} = \sum_i x_i y_i - \frac{\sum_i x_i \sum_i y_i}{n} \\ \sum_i (x_i - \bar{x})^2 &= \sum_i x_i^2 - n\bar{x}^2 = \sum_i x_i^2 - \frac{(\sum_i x_i)^2}{n} \\ \sum_i (y_i - \bar{y})^2 &= \sum_i y_i^2 - n\bar{y}^2 = \sum_i y_i^2 - \frac{(\sum_i y_i)^2}{n}\end{aligned}$$

Thus we obtain the formula for r_{XY} :

$$r_{XY} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \cdot \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$

The value of r_{XY} lies within -1 and 1 i.e. $-1 \leq r_{XY} \leq 1$.

- $r > 0$ indicates a **positive correlation**.
- $r < 0$ indicates a **negative correlation**.
- $r = 0$ suggests **no correlation** (but not necessarily independence).
- $r = \pm 1$ indicates a **perfect linear correlation** (+1 for positive and -1 for negative).



Figure 11.2: Scatterplots illustrating positive, negative, zero, and perfect linear correlations.

Example: Suppose we have data on students' hours studied (X) and exam scores (Y) as given in the table 11.1. We need to calculate the Pearson correlation coefficient.

The formula we will use:

$$r_{XY} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \cdot \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$

Student	x_i	y_i	x_i^2	y_i^2	$x_i y_i$
1	2	65	4	4225	130
2	3	70	9	4900	210
3	4	76	16	5776	304
4	5	78	25	6084	390
5	6	82	36	6724	492
6	7	89	49	7921	623
$n = 6$	$\sum x_i = 27$	$\sum y_i = 460$	$\sum x_i^2 = 139$	$\sum y_i^2 = 35630$	$\sum x_i y_i = 2149$

Table 11.2: *Correlation coefficient calculation table.*

Substituting the values:

$$\begin{aligned}
r_{XY} &= \frac{6 \cdot 2149 - 27 \cdot 460}{\sqrt{6 \cdot 139 - 27^2} \cdot \sqrt{6 \cdot 35630 - 460^2}} \\
&= \frac{12894 - 12420}{\sqrt{834 - 729} \cdot \sqrt{213780 - 211600}} \\
&= \frac{474}{\sqrt{105} \cdot \sqrt{2180}} \approx \frac{474}{10.247 \cdot 46.690} \\
&\approx \frac{474}{478.74} \approx 0.9907
\end{aligned}$$

This indicates a very strong positive correlation between hours studied and exam scores.

11.3 Effect of Linear Transformations on Correlation Coefficient

Theorem: Suppose we have two random variables X and Y with correlation coefficient ρ_{XY} . Define new variables

$$U = a + bX, \quad V = c + dY$$

where a, c are constants (shifts) and $b, d \neq 0$ are scaling factors. Then the population correlation coefficient between U and V is

$$\rho_{UV} = \frac{bd}{|b||d|} \cdot \rho_{XY}$$

Proof:

$$\rho_{UV} = \frac{\text{Cov}(U, V)}{\sigma_U \sigma_V} = \frac{\text{Cov}(a + bX, c + dY)}{\sqrt{\text{Var}(a + bX)} \cdot \sqrt{\text{Var}(c + dY)}}$$

Now,

$$\text{Cov}(a + bX, c + dY) = bd \text{Cov}(X, Y)$$

$$\text{Var}(a + bX) = b^2 \text{Var}(X)$$

$$\text{Var}(c + dY) = d^2 \text{Var}(Y)$$

Thus,

$$\rho_{UV} = \frac{\text{Cov}(bX, dY)}{\sqrt{\text{Var}(bX)} \cdot \sqrt{\text{Var}(dY)}} = \frac{bd \text{Cov}(X, Y)}{|b|\sigma_X \cdot |d|\sigma_Y} = \frac{bd}{|b||d|} \cdot \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

Since

$$\frac{bd}{|b||d|} = \text{sgn}(b) \cdot \text{sgn}(d)$$

where $\text{sgn}(\cdot)$ is the sign function¹, it follows that

$$\rho_{UV} = \frac{bd}{|b||d|} \cdot \rho_{XY} = \text{sgn}(b) \cdot \text{sgn}(d) \cdot \rho_{XY}$$

■

This shows that linear transformations preserve the magnitude of the correlation coefficient but may reverse its sign depending on the scaling factors.

- If b and d have the same sign, then $\rho_{UV} = \rho_{XY}$.
- If b and d have opposite signs, then $\rho_{UV} = -\rho_{XY}$.
- Adding constants a and c has no effect on the value of the correlation.

The sample correlation coefficient also satisfies the same transformation relation:

$$r_{UV} = \text{sgn}(b) \cdot \text{sgn}(d) \cdot r_{XY}$$

11.4 Correlation Is Not Causation

It is important to understand that a high correlation between two variables does not necessarily imply that one variable causes the other. Correlation measures the strength and direction of a linear relationship between variables, but it does not provide evidence about causality.

There are several reasons why correlation does not imply causation:

- **Confounding variables:** A third variable may influence both variables under study, creating a spurious correlation.
- **Reverse causality:** The direction of cause and effect may be opposite to what is assumed.
- **Coincidence:** Sometimes, correlations arise purely by chance.

For example, consider a study might that finds a positive correlation between umbrella sales and slipping accidents. This doesn't mean umbrellas cause slips or vice versa. Instead, rainy weather acts as the common factor: rain simultaneously drives people to buy umbrellas and makes sidewalks slippery, creating the illusion of a direct relationship between the two variables.

This example illustrates why it is crucial to use careful experimental design, statistical controls, and domain knowledge before concluding causal relationships from correlated data.

¹The **sign function**, denoted by $\text{sgn}(x)$, is defined as:

$$\text{sgn}(x) = \begin{cases} -1, & \text{if } x < 0, \\ 0, & \text{if } x = 0, \\ 1, & \text{if } x > 0. \end{cases}$$

It extracts the sign of a real number x .

11.5 Regression Analysis

Regression analysis is a statistical method used to explore and model the relationship between a dependent variable and one or more independent variables.

It plays a central role in data analysis, prediction, and inference, particularly when trying to establish a functional relationship between variables.

In its simplest form—**simple linear regression**—we wish to study the relationship between two variables X and Y and use it to predict Y from X . The variable X acts as the **independent variable** (predictor, causal variable) whose values are controlled by the experimenter and Y is the **dependent variable** (response) which is also subjected to unaccountable variations (errors).

For example, a teacher wants to examine whether there's a relationship between how long students study and the scores they achieve in a test. By treating the number of hours studied as the independent variable and the test score as the dependent variable, regression helps us determine whether there is a consistent trend between the two.

11.6 The Simple Linear Regression Model

A simple linear regression model assumes the existence of a linear relationship between X (predictor variable) and Y (response) that is disturbed by a random error ϵ which can be written as an equation of the form:

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

where:

- β_0 : y -intercept of the line,
- β_1 : slope of the line (rate of change in Y per unit increase in X),
- ϵ : random error, accounting for unexplained variation

Given a dataset of n observations, represented as pairs:

$$(x_1, y_1), (x_2, y_2), \dots (x_n, y_n)$$

the objective is to estimate the unknown parameters β_0 and β_1 , and then use these estimates to define a straight line that best fits the data.

The **fitted regression line** (also called the prediction equation) is:

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

Here, $\hat{\beta}_0$ and $\hat{\beta}_1$ are the estimated values of the intercept and slope, respectively, obtained from the sample data.



11.7 Estimating Parameters Using Least Squares

Given n data points $(x_1, y_1), \dots, (x_n, y_n)$, the estimates of the parameters β_0 and β_1 should result in a line that is (in some sense) a “best fit” to the data. To define what we mean by a “best fit” line, consider each data point (x_i, y_i) and its corresponding prediction $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ from the regression line. This predicted value is known as the **fitted value**. The difference between the observed value y_i and the fitted value \hat{y}_i is called the **residual** (error), denoted by ϵ_i :

$$\epsilon_i = y_i - \hat{y}_i$$

The residual ϵ_i represents the vertical distance between an observed data point and the regression line. A positive residual indicates that the point lies above the regression line, while a negative residual indicates that it lies below. The closer the residuals are to zero, the better the fitted values approximate the observed data. Therefore, the estimates of the parameters β_0 and β_1 should be such that these residuals (errors) are as small as possible. However, minimizing the simple sum of residuals is not appropriate, because the positive and negative errors can cancel each other out, even when individual errors are large. A more effective approach is the **method of least squares**, which involves minimizing the sum of the squares of the residuals.



The **least-squares line** is the line that minimizes the **residual sum of squares (RSS)**²:

²In some texts, the residual sum of squares (RSS) is also called the **sum of squares of errors (SSE)**

$$RSS = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

To derive the expressions for $\hat{\beta}_0$ and $\hat{\beta}_1$, we have to take the partial derivatives of RSS with respect to $\hat{\beta}_0$ and $\hat{\beta}_1$ and set them to zero.

$$\begin{aligned} \frac{\partial}{\partial \hat{\beta}_0}(RSS) &= -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0, \\ \frac{\partial}{\partial \hat{\beta}_1}(RSS) &= -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0 \end{aligned}$$

Simplifying these two equations yields

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i, \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n y_i x_i \end{aligned}$$

This set of two linear equations in the unknowns β_1 and β_2 , called the **normal equations**, which provides the best fit to a given set of paired data in accordance with the criterion of least squares. Multiplying the first equation by $\sum_{i=1}^n x_i$ and second equation by n and then subtracting second from the first yields

$$\left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) - n \sum_{i=1}^n x_i y_i = \hat{\beta}_1 \left(\sum_{i=1}^n x_i \right)^2 - n \hat{\beta}_1 \sum_{i=1}^n x_i^2$$

Thus we get the expression for $\hat{\beta}_1$

$$\begin{aligned} \hat{\beta}_1 &= \left[n \sum_{i=1}^n x_i y_i - \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \right] / \left[n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2 \right] \\ &= \left[\sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \right] / \left[\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right] \end{aligned}$$

It is convenient to introduce some notation for the sums of squared deviation from mean and sums of cross-products of deviation.

$$\begin{aligned} S_{xx} &= \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \\ S_{yy} &= \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n y_i^2 - \frac{1}{n} \left(\sum_{i=1}^n y_i \right)^2 \\ S_{xy} &= \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \end{aligned}$$

Using these notation we can write,

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$

The expression for $\hat{\beta}_0$ is calculated as

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

11.7.1 Calculation of RSS

We can now calculate the value of RSS based on the value of $\hat{\beta}_0$ and $\hat{\beta}_1$ found using the method of least square:

$$\begin{aligned} \text{RSS} &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \sum_{i=1}^n \left[y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \right]^2 \quad \text{Since } \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \\ &= \sum_{i=1}^n \left[y_i - (\bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_i) \right]^2 \quad \text{Since } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \\ &= \sum_{i=1}^n \left[(y_i - \bar{y}) - \hat{\beta}_1 (x_i - \bar{x}) \right]^2 \\ &= \sum_{i=1}^n (y_i - \bar{y})^2 - 2\hat{\beta}_1 \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) + \hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= S_{yy} - 2\hat{\beta}_1 S_{xy} + \hat{\beta}_1^2 S_{xx} \\ &= S_{yy} - 2 \left(\frac{S_{xy}}{S_{xx}} \right) S_{xy} + \left(\frac{S_{xy}}{S_{xx}} \right)^2 S_{xx} \quad \text{Since } \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \\ &= S_{yy} - \frac{2S_{xy}^2}{S_{xx}} + \frac{S_{xy}^2}{S_{xx}} \\ &= S_{yy} - \frac{S_{xy}^2}{S_{xx}} \\ &= S_{yy} - \hat{\beta}_1 S_{xy} \end{aligned}$$

$$\text{RSS} = S_{yy} - \hat{\beta}_1 S_{xy}$$

11.7.2 Example

We aim to model the relationship between the number of hours studied (X) and the corresponding test score (Y) using simple linear regression. The goal is to estimate a linear equation that best describes this relationship based on observed data from five students. Once the model is established, we will use it to predict the expected test score for a student who studies for six hours.

The observed data are as follows:

Student	Hours Studied (x_i)	Test Score (y_i)
1	2	65
2	3	70
3	5	75
4	7	85
5	9	95

Step 1: Compute Means

$$\bar{x} = \frac{2 + 3 + 5 + 7 + 9}{5} = \frac{26}{5} = 5.2,$$

$$\bar{y} = \frac{65 + 70 + 75 + 85 + 95}{5} = \frac{390}{5} = 78$$

Step 2: Compute S_{xx} and S_{xy}

$$\begin{aligned} S_{xx} &= \sum (x_i - \bar{x})^2 \\ &= (2 - 5.2)^2 + (3 - 5.2)^2 + (5 - 5.2)^2 + (7 - 5.2)^2 + (9 - 5.2)^2 \\ &= 32.8 \end{aligned}$$

$$\begin{aligned} S_{xy} &= \sum (x_i - \bar{x})(y_i - \bar{y}) \\ &= (2 - 5.2)(65 - 78) + (3 - 5.2)(70 - 78) + \dots + (9 - 5.2)(95 - 78) \\ &= 137.0 \end{aligned}$$

Step 3: Estimate Parameters

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{137.0}{32.80} \approx 4.177$$

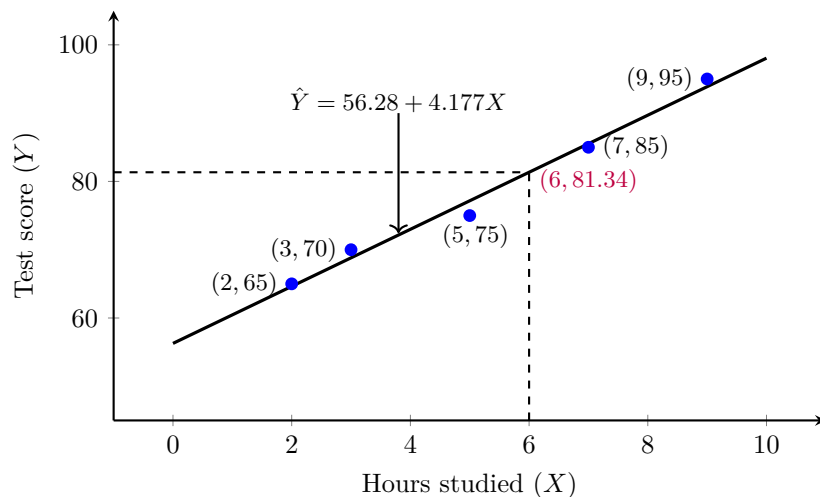
$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \overline{xx} = 78 - 4.177 \times 5.2 \approx 56.28$$

Step 4: Regression Equation

$$\hat{Y} = 56.28 + 4.177X$$

Step 5: Predict Test Score for $X = 6$

$$\hat{Y} = 56.28 + 4.177 \times 6 \approx 81.34$$



11.8 Multiple Regression

Sometimes a response variable Y depends on multiple independent variables X_1, X_2, \dots, X_k . The general form of the multiple regression model is:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon$$

where:

- β_0 is the intercept parameter,
- β_j is the slope parameter for the j -th independent variable,
- ε is the error term.

Based on data we try to obtain good estimates of the regression coefficients $\hat{\beta}_0, \hat{\beta}_1 \dots \hat{\beta}_p$. The prediction equation is given by

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \dots + \hat{\beta}_p X_p + \varepsilon$$

Generally the least square (LS) method is used to estimate the regression coefficients as before. LS method can be applied only if number of samples n is bigger than the number of predictors p . Let x_{ki} be the i th sample of the k th predictor input and y_i be the corresponding response. Then

$$\text{RSS} = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n \left(y_i - \hat{\beta}_0 - \sum_{k=1}^p \hat{\beta}_k x_{ki} \right)^2$$

As before we could take the derivatives with respect to $\hat{\beta}_0, \hat{\beta}_1 \dots \hat{\beta}_p$, set them equal to zero and derive the least square **normal equations** that our parameters $\hat{\beta}_0, \hat{\beta}_1 \dots \hat{\beta}_p$ would have to fulfill.

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_{1i} + \hat{\beta}_2 \sum_{i=1}^n x_{2i} + \dots + \hat{\beta}_p \sum_{i=1}^n x_{pi} &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_{1i} + \hat{\beta}_1 \sum_{i=1}^n x_{1i}^2 + \hat{\beta}_2 \sum_{i=1}^n x_{1i}x_{2i} + \dots + \hat{\beta}_p \sum_{i=1}^n x_{1i}x_{pi} &= \sum_{i=1}^n x_{1i}y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_{2i} + \hat{\beta}_1 \sum_{i=1}^n x_{2i}x_{1i} + \hat{\beta}_2 \sum_{i=1}^n x_{2i}^2 + \dots + \hat{\beta}_p \sum_{i=1}^n x_{2i}x_{pi} &= \sum_{i=1}^n x_{2i}y_i \\ &\vdots \\ \hat{\beta}_0 \sum_{i=1}^n x_{pi} + \hat{\beta}_1 \sum_{i=1}^n x_{pi}x_{1i} + \hat{\beta}_2 \sum_{i=1}^n x_{pi}x_{2i} + \dots + \hat{\beta}_p \sum_{i=1}^n x_{pi}^2 &= \sum_{i=1}^n x_{pi}y_i \end{aligned}$$

This procedure can be easily carried out if we formulate the following prediction equation in matrix

and vector form.

$$\begin{aligned}\hat{y}_1 &= \hat{\beta}_0 + \hat{\beta}_1 x_{11} + \hat{\beta}_2 x_{21} + \cdots + \hat{\beta}_p x_{p1} \\ \hat{y}_2 &= \hat{\beta}_0 + \hat{\beta}_1 x_{12} + \hat{\beta}_2 x_{22} + \cdots + \hat{\beta}_p x_{p2} \\ &\vdots \\ \hat{y}_n &= \hat{\beta}_0 + \hat{\beta}_1 x_{1n} + \hat{\beta}_2 x_{2n} + \cdots + \hat{\beta}_p x_{pn}\end{aligned}$$

In matrix form it is written as

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$$

where,

$$\hat{\mathbf{y}} = \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{21} & \cdots & x_{k1} \\ 1 & x_{12} & x_{22} & \cdots & x_{k2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & x_{2n} & \cdots & x_{kn} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}$$

If \mathbf{y} is the vector for observed variable and $\boldsymbol{\varepsilon}$, then

$$\boldsymbol{\varepsilon} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$$

where,

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

Then the equation for residual sum of square (RSS) in matrix form is given by

$$\begin{aligned}\text{RSS} &= \sum_{i=1}^n \varepsilon_i^2 = \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} \\ &= (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \\ &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{y} + \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}}\end{aligned}$$

To derive the expressions for $\hat{\boldsymbol{\beta}}$, we have to take the partial derivatives of RSS with respect to $\hat{\boldsymbol{\beta}}$ using matrix calculus³ and set it equal to zero. The partial derivative of RSS is

$$\begin{aligned}\frac{\partial}{\partial \hat{\boldsymbol{\beta}}} (\text{RSS}) &= \frac{\partial}{\partial \hat{\boldsymbol{\beta}}} (\mathbf{y}^T \mathbf{y} - 2\hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{y} + \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}}) \\ &= -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} \\ &= 2(\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{X}^T \mathbf{y}) \\ &= 2(\hat{\boldsymbol{\beta}} \mathbf{X}^T - \mathbf{y}^T) \mathbf{X}\end{aligned}$$

³Some matrix calculation formulas:

$$\begin{aligned}\frac{\partial}{\partial \mathbf{x}} (\mathbf{y}^T \mathbf{A} \mathbf{x}) &= \mathbf{y}^T \mathbf{A} \\ \frac{\partial}{\partial \mathbf{y}} (\mathbf{y}^T \mathbf{A} \mathbf{x}) &= \mathbf{x}^T \mathbf{A}^T \\ \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{A} \mathbf{x}) &= (\mathbf{A} + \mathbf{A}^T) \mathbf{x}\end{aligned}$$

To minimize RSS with respect to $\hat{\beta}$, we set

$$\begin{aligned}\frac{\partial}{\partial \hat{\beta}} (\text{RSS}) &= 0 \\ \Rightarrow (\hat{\beta} \mathbf{X}^T - \mathbf{y}^T) \mathbf{X} &= 0 \\ \Rightarrow \hat{\beta} \mathbf{X}^T \mathbf{X} - \mathbf{y}^T \mathbf{X} &= 0 \\ \Rightarrow \mathbf{X}^T \mathbf{X} \hat{\beta} - \mathbf{X}^T \mathbf{y} &= 0\end{aligned}$$

Therefore we can write:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

11.8.1 Assumptions of Multiple Regression

The classical multiple regression model relies on the following assumptions:

1. **Linearity:** The relationship between Y and each X_j is linear
2. **Independence:** Observations are independent of each other
3. **Homoscedasticity:** $\text{Var}(\varepsilon_i) = \sigma^2$ for all i
4. **Normality:** $\varepsilon_i \sim N(0, \sigma^2)$
5. **No perfect multicollinearity:** The matrix $\mathbf{X}^T \mathbf{X}$ is invertible

11.8.2 Least Squares Estimation

The ordinary least squares (OLS) estimator minimizes the sum of squared residuals:

$$S(\beta) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{1i} - \cdots - \beta_k X_{ki})^2 \quad (11.1)$$

In matrix form:

$$S(\beta) = (\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta) \quad (11.2)$$

Taking the derivative with respect to β and setting equal to zero:

$$\frac{\partial S}{\partial \beta} = -2\mathbf{X}^T \mathbf{Y} + 2\mathbf{X}^T \mathbf{X} \beta = \mathbf{0} \quad (11.3)$$

The normal equations are:

$$\mathbf{X}^T \mathbf{X} \beta = \mathbf{X}^T \mathbf{Y} \quad (11.4)$$

The OLS estimator is:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad (11.5)$$

11.8.3 Properties of the OLS Estimator

Under the classical assumptions, the OLS estimator has the following properties:

1. **Unbiasedness:** $E[\hat{\beta}] = \beta$

2. **Variance-Covariance Matrix:**

$$\text{Var}(\hat{\beta}) = \sigma^2(\mathbf{X}^T \mathbf{X})^{-1} \quad (11.6)$$

3. **Best Linear Unbiased Estimator (BLUE):** By the Gauss-Markov theorem

4. **Consistency:** $\hat{\beta} \xrightarrow{p} \beta$ as $n \rightarrow \infty$

5. **Asymptotic Normality:** $\hat{\beta} \sim N(\beta, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$

11.8.4 Estimation of Error Variance

The unbiased estimator of σ^2 is:

$$s^2 = \frac{\text{SSE}}{n - k - 1} = \frac{\sum_{i=1}^n e_i^2}{n - k - 1} \quad (11.7)$$

where $e_i = Y_i - \hat{Y}_i$ are the residuals and $n - k - 1$ represents the degrees of freedom.

11.8.5 Coefficient of Determination (R^2)

The coefficient of determination measures the proportion of variance explained by the regression:

$$R^2 = \frac{\text{SSR}}{\text{SST}} = 1 - \frac{\text{SSE}}{\text{SST}} \quad (11.8)$$

where:

$$\text{SST} = \sum_{i=1}^n (Y_i - \bar{Y})^2 \quad (\text{Total Sum of Squares}) \quad (11.9)$$

$$\text{SSR} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 \quad (\text{Regression Sum of Squares}) \quad (11.10)$$

$$\text{SSE} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (\text{Error Sum of Squares}) \quad (11.11)$$

The adjusted R^2 accounts for the number of predictors:

$$R_{\text{adj}}^2 = 1 - \frac{\text{SSE}/(n - k - 1)}{\text{SST}/(n - 1)} = 1 - \frac{n - 1}{n - k - 1}(1 - R^2) \quad (11.12)$$

11.8.6 Hypothesis Testing

Test for Overall Significance of Regression

To test whether the regression model is significant overall:

$$H_0 : \beta_1 = \beta_2 = \cdots = \beta_k = 0 \quad (11.13)$$

$$H_1 : \text{At least one } \beta_j \neq 0 \quad (11.14)$$

The test statistic follows an F-distribution:

$$F = \frac{\text{MSR}}{\text{MSE}} = \frac{\text{SSR}/k}{\text{SSE}/(n-k-1)} \sim F_{k,n-k-1} \quad (11.15)$$

We reject H_0 if $F > F_{\alpha,k,n-k-1}$ or if the p-value $< \alpha$.

Test for Individual Regression Coefficients

To test the significance of an individual coefficient β_j :

$$H_0 : \beta_j = 0 \quad (11.16)$$

$$H_1 : \beta_j \neq 0 \quad (11.17)$$

The test statistic is:

$$t_j = \frac{\hat{\beta}_j}{\text{SE}(\hat{\beta}_j)} = \frac{\hat{\beta}_j}{s\sqrt{c_{jj}}} \sim t_{n-k-1} \quad (11.18)$$

where c_{jj} is the $(j+1, j+1)$ element of $(\mathbf{X}^T \mathbf{X})^{-1}$ and $\text{SE}(\hat{\beta}_j) = s\sqrt{c_{jj}}$.

We reject H_0 if $|t_j| > t_{\alpha/2, n-k-1}$ or if the p-value $< \alpha$.

Test for a Group of Predictors (Partial F-Test)

To test the significance of a subset of predictors, we compare a full model with a reduced model.

Full model: $Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k + \varepsilon$

Reduced model: $Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_r X_r + \varepsilon$ (where $r < k$)

Hypotheses:

$$H_0 : \beta_{r+1} = \beta_{r+2} = \cdots = \beta_k = 0 \quad (11.19)$$

$$H_1 : \text{At least one of } \beta_{r+1}, \dots, \beta_k \neq 0 \quad (11.20)$$

The test statistic is:

$$F = \frac{(\text{SSE}_R - \text{SSE}_F)/(k-r)}{\text{SSE}_F/(n-k-1)} \sim F_{k-r, n-k-1} \quad (11.21)$$

where SSE_R and SSE_F are the error sums of squares for the reduced and full models, respectively.

11.8.7 Confidence Intervals

Confidence Interval for Individual Coefficients

A $(1-\alpha)100\%$ confidence interval for β_j is:

$$\hat{\beta}_j \pm t_{\alpha/2, n-k-1} \cdot \text{SE}(\hat{\beta}_j) \quad (11.22)$$

Confidence Interval for Mean Response

For a given set of predictor values $\mathbf{x}_0 = (1, x_{10}, x_{20}, \dots, x_{k0})^T$, the predicted mean response is:

$$\hat{Y}_0 = \mathbf{x}_0^T \hat{\boldsymbol{\beta}} \quad (11.23)$$

The $(1-\alpha)100\%$ confidence interval for the mean response is:

$$\hat{Y}_0 \pm t_{\alpha/2, n-k-1} \cdot s\sqrt{\mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0} \quad (11.24)$$

Prediction Interval for Individual Response

The $(1 - \alpha)100\%$ prediction interval for a new observation is:

$$\hat{Y}_0 \pm t_{\alpha/2, n-k-1} \cdot s \sqrt{1 + \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0} \quad (11.25)$$

11.8.8 Example: Housing Price Prediction

Consider predicting house prices based on size (sq ft), number of bedrooms, and age of the house.

Model: $\text{Price} = \beta_0 + \beta_1 \text{Size} + \beta_2 \text{Bedrooms} + \beta_3 \text{Age} + \varepsilon$

Suppose we have the following results from $n = 100$ observations:

Variable	Coefficient	SE	t-statistic	p-value
Intercept	50,000	15,000	3.33	0.001
Size	120	10	12.00	< 0.001
Bedrooms	8,000	2,500	3.20	0.002
Age	-500	200	-2.50	0.014

Additional statistics: $R^2 = 0.85$, $s = 20,000$, $F = 182.4$ (p-value < 0.001)

Interpretation:

- For each additional square foot, price increases by \$120 on average, holding other variables constant
- Each additional bedroom increases price by \$8,000 on average
- Each additional year of age decreases price by \$500 on average
- The model explains 85% of the variance in house prices
- All coefficients are statistically significant at $\alpha = 0.05$
- The overall model is highly significant ($F = 182.4$, p-value < 0.001)

11.8.9 ANOVA Table for Multiple Regression

Source	DF	SS	MS	F	p-value
Regression	k	SSR	$\text{MSR} = \text{SSR}/k$	MSR/MSE	
Error	$n - k - 1$	SSE	$\text{MSE} = \text{SSE}/(n - k - 1)$		
Total	$n - 1$	SST			

11.8.10 Diagnostic Checks

After fitting a multiple regression model, it's essential to check the assumptions:

1. **Residual plots:** Plot residuals vs. fitted values to check linearity and homoscedasticity
2. **Normal Q-Q plot:** Check normality of residuals
3. **Leverage and influence:** Identify outliers and influential observations
4. **Multicollinearity:** Check variance inflation factors (VIF)

The variance inflation factor for the j -th predictor is:

$$\text{VIF}_j = \frac{1}{1 - R_j^2} \quad (11.26)$$

where R_j^2 is the coefficient of determination from regressing X_j on all other predictors.

Chapter 12

Theory of Errors

12.1 Introduction

The measurement of any physical quantity can never be made with perfect accuracy. There will always be some error or uncertainty present in every experimental observation. For any measurement, there exists an infinite number of factors that can cause the experimentally obtained value to deviate from the true (theoretical) value.

Let us denote the true value of a quantity as X_{true} and the measured value as X_{measured} . The absolute error ΔX is defined as:

$$\Delta X = X_{\text{measured}} - X_{\text{true}}$$

The relative error is given by:

$$\epsilon = \frac{\Delta X}{X_{\text{true}}} = \frac{X_{\text{measured}} - X_{\text{true}}}{X_{\text{true}}}$$

When experimental results are reported, they must be accompanied by an estimate of the experimental uncertainty, which indicates the reliability of the measurement.

Experimental errors can be systematically classified into three fundamental categories:

1. Random errors
2. Systematic errors
3. Gross errors

12.2 Random Errors

Random errors are fluctuations in measurements that vary unpredictably from one observation to another. These errors follow statistical laws and can be analyzed using probability theory.

If we perform n measurements of a quantity X , obtaining values x_1, x_2, \dots, x_n , the random errors ϵ_i are defined as:

$$\epsilon_i = x_i - \mu$$

where μ is the true mean value. Random errors satisfy the following conditions:

$$\begin{aligned}\langle \epsilon_i \rangle &= 0 \quad (\text{zero mean}) \\ \langle \epsilon_i \epsilon_j \rangle &= 0 \quad \text{for } i \neq j \quad (\text{uncorrelated})\end{aligned}$$

Statistical Analysis

For a set of n measurements $\{x_1, x_2, \dots, x_n\}$, the sample mean (best estimate) is:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (12.1)$$

The sample variance is:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (12.2)$$

The standard deviation is:

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \quad (12.3)$$

Derivation of Standard Error of the Mean:

The standard error of the mean $\sigma_{\bar{x}}$ represents the uncertainty in our estimate of the true mean. For independent measurements with individual standard deviation σ :

$$\sigma_{\bar{x}}^2 = \text{Var} \left(\frac{1}{n} \sum_{i=1}^n x_i \right) \quad (12.4)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(x_i) \quad (12.5)$$

$$= \frac{1}{n^2} \cdot n\sigma^2 \quad (12.6)$$

$$= \frac{\sigma^2}{n} \quad (12.7)$$

Therefore:

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{s}{\sqrt{n}} \quad (12.8)$$

Normal Distribution

For large n , random errors typically follow a normal (Gaussian) distribution:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left(-\frac{(x - \mu)^2}{2\sigma^2} \right) \quad (12.9)$$

This leads to important statistical properties:

- 68.27% of measurements fall within $\pm 1\sigma$ of the mean

- 95.45% of measurements fall within $\pm 2\sigma$ of the mean
- 99.73% of measurements fall within $\pm 3\sigma$ of the mean

Example: Repeated Measurements of Length

Consider measuring the length of a rod 10 times with a ruler:

Measurement	Length (cm)
1	15.23
2	15.25
3	15.21
4	15.24
5	15.22
6	15.26
7	15.20
8	15.25
9	15.23
10	15.24

Calculations:

$$\bar{x} = \frac{152.33}{10} = 15.233 \text{ cm} \quad (12.10)$$

$$s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = 0.019 \text{ cm} \quad (12.11)$$

$$\sigma_{\bar{x}} = \frac{s}{\sqrt{n}} = \frac{0.019}{\sqrt{10}} = 0.006 \text{ cm} \quad (12.12)$$

Result: $L = 15.233 \pm 0.006 \text{ cm}$

12.3 Systematic Errors

Systematic errors produce consistent deviations from the true value in the same direction. Unlike random errors, systematic errors do not average to zero when measurements are repeated.

If δ represents a systematic error, then all measurements are affected by the same bias:

$$x_i = X_{\text{true}} + \delta + \epsilon_i$$

where ϵ_i are random errors. The mean of measurements becomes:

$$\bar{x} = X_{\text{true}} + \delta + \frac{1}{n} \sum_{i=1}^n \epsilon_i \approx X_{\text{true}} + \delta$$

Sources of Systematic Errors

1. **Instrumental Errors:** Calibration errors, zero drift, non-linearity
2. **Environmental Errors:** Temperature effects, pressure variations
3. **Method Errors:** Theoretical approximations, neglected effects
4. **Personal Errors:** Consistent bias in reading instruments

Detection and Correction

Method 1: Calibration with Standards

Use known reference standards to determine systematic bias:

$$\delta = \bar{x}_{\text{standard}} - X_{\text{known}} \quad (12.13)$$

Method 2: Comparison with Independent Methods

If two independent methods give results \bar{x}_1 and \bar{x}_2 , the systematic difference is:

$$\Delta_{\text{sys}} = \bar{x}_1 - \bar{x}_2 \quad (12.14)$$

Example: Voltmeter Calibration

A digital voltmeter is tested against a precision standard:

Standard (V)	Measured (V)	Error (V)
1.000	1.003	+0.003
2.000	2.006	+0.006
5.000	5.015	+0.015
10.000	10.030	+0.030

The systematic error shows a linear relationship: $\delta = 0.003 \times V_{\text{reading}}$

Correction formula: $V_{\text{true}} = V_{\text{measured}} \times (1 - 0.003)$

12.4 Gross Errors (Blunders)

Gross errors are large, obvious mistakes that occur due to human error, equipment malfunction, or procedural failures. These errors are typically much larger than random or systematic errors and can severely distort results.

A measurement x_i is considered a potential outlier if:

$$|x_i - \bar{x}| > k \cdot s$$

where k is typically 2 or 3, depending on the desired confidence level.

Statistical Tests for Outliers

Chauvenet's Criterion:

A data point should be rejected if the probability of obtaining such a deviation is less than $\frac{1}{2n}$, where n is the number of measurements.

For a normal distribution, reject x_i if:

$$\left| \frac{x_i - \bar{x}}{s} \right| > t_{\text{Chauvenet}} \quad (12.15)$$

where $t_{\text{Chauvenet}}$ depends on n :

n	$t_{\text{Chauvenet}}$
5	1.65
10	1.96
15	2.13
25	2.33
50	2.57

Grubbs' Test:

The test statistic is:

$$G = \frac{\max |x_i - \bar{x}|}{s} \quad (12.16)$$

Compare with critical value $G_{\text{critical}}(\alpha, n)$ from Grubbs' table.

Example: Detection of Blunder

Measurements of gravitational acceleration (m/s²): 9.78, 9.82, 9.79, 9.81, 9.85, 10.15, 9.80

Calculations:

$$\bar{x} = 9.86 \text{ m/s}^2 \quad (12.17)$$

$$s = 0.13 \text{ m/s}^2 \quad (12.18)$$

$$G = \frac{|10.15 - 9.86|}{0.13} = 2.23 \quad (12.19)$$

For $n = 7$ and $\alpha = 0.05$, $G_{\text{critical}} = 2.02$

Since $G > G_{\text{critical}}$, the value 10.15 is identified as an outlier and should be investigated.

12.5 Error Propagation

When quantities with uncertainties are combined in calculations, the uncertainties propagate according to specific mathematical rules.

12.5.1 General Formula

For a function $f(x, y, z, \dots)$ where x, y, z have uncertainties $\delta x, \delta y, \delta z$:

$$(\delta f)^2 = \left(\frac{\partial f}{\partial x}\right)^2 (\delta x)^2 + \left(\frac{\partial f}{\partial y}\right)^2 (\delta y)^2 + \left(\frac{\partial f}{\partial z}\right)^2 (\delta z)^2 + \dots \quad (12.20)$$

12.5.2 Specific Cases

Addition/Subtraction: $f = x \pm y$

$$\delta f = \sqrt{(\delta x)^2 + (\delta y)^2} \quad (12.21)$$

Multiplication/Division: $f = \frac{xy}{z}$

$$\frac{\delta f}{f} = \sqrt{\left(\frac{\delta x}{x}\right)^2 + \left(\frac{\delta y}{y}\right)^2 + \left(\frac{\delta z}{z}\right)^2} \quad (12.22)$$

Powers: $f = x^n$

$$\frac{\delta f}{f} = |n| \frac{\delta x}{x} \quad (12.23)$$

12.5.3 Derivation Example: Product of Two Variables

For $f = xy$:

$$\frac{\partial f}{\partial x} = y \quad (12.24)$$

$$\frac{\partial f}{\partial y} = x \quad (12.25)$$

Therefore:

$$(\delta f)^2 = y^2(\delta x)^2 + x^2(\delta y)^2 \quad (12.26)$$

$$\frac{(\delta f)^2}{f^2} = \frac{y^2(\delta x)^2 + x^2(\delta y)^2}{x^2y^2} \quad (12.27)$$

$$= \left(\frac{\delta x}{x}\right)^2 + \left(\frac{\delta y}{y}\right)^2 \quad (12.28)$$

12.6 Practical Example: Pendulum Experiment

Consider measuring the acceleration due to gravity using a simple pendulum:

$$g = \frac{4\pi^2 L}{T^2} \quad (12.29)$$

Measurements:

$$L = 1.000 \pm 0.001 \text{ m} \quad (12.30)$$

$$T = 2.006 \pm 0.003 \text{ s} \quad (12.31)$$

Error Propagation:

$$\frac{\delta g}{g} = \sqrt{\left(\frac{\delta L}{L}\right)^2 + \left(2\frac{\delta T}{T}\right)^2} \quad (12.32)$$

Calculations:

$$\frac{\delta L}{L} = \frac{0.001}{1.000} = 0.001 \quad (12.33)$$

$$\frac{\delta T}{T} = \frac{0.003}{2.006} = 0.00150 \quad (12.34)$$

$$\frac{\delta g}{g} = \sqrt{(0.001)^2 + (2 \times 0.00150)^2} = 0.00316 \quad (12.35)$$

Result:

$$g = \frac{4\pi^2 \times 1.000}{(2.006)^2} = 9.79 \text{ m/s}^2 \quad (12.36)$$

$$\delta g = 9.79 \times 0.00316 = 0.03 \text{ m/s}^2 \quad (12.37)$$

Final answer: $g = 9.79 \pm 0.03 \text{ m/s}^2$

12.7 Conclusion

Understanding and properly handling experimental errors is fundamental to scientific measurement. The three types of errors—random, systematic, and gross errors—each require different approaches:

- Random errors are reduced by statistical methods and multiple measurements
- Systematic errors require careful calibration and method validation
- Gross errors must be identified and eliminated through quality control

Proper error analysis ensures that experimental results are reliable, reproducible, and scientifically meaningful. The mathematical framework presented here provides the tools necessary for rigorous uncertainty analysis in experimental physics and engineering.

Chapter 13

Statistical Quality Control

13.1 Introduction to Quality and Statistical Quality Control

13.1.1 Meaning of Quality

Quality is about how well a product or service meets the needs of the customer. For example, if you buy a mobile phone and it works smoothly, has good battery life, and doesn't break easily, you would say it is a high-quality product. In both factories and service industries, keeping quality consistent is important to keep customers happy, reduce waste, and stay ahead of the competition.

Quality can be understood in different ways:

- **Conformance to specifications:** Does the product meet the standards set by the manufacturer? For example, if a bolt is supposed to be 10 mm long, and it actually is 10 mm, it meets the specification.
- **Fitness for use:** Does the product do what the customer needs? A raincoat that keeps you dry during heavy rain is fit for use.
- **Value for money:** Is the customer getting good performance for the price they pay? A budget phone that works reliably for daily tasks could be good value for money.

13.1.2 Statistical Quality Control (SQC)

Quality control is the process of monitoring performance and comparing it with a standard. If deviations are found, corrective actions are taken to bring the process back into control.

Statistical Quality Control (SQC) applies statistical techniques to monitor, control, and improve product and process quality. Since inspecting every item is often impractical, a smaller sample is often taken and used to draw conclusions about the entire production lot.

For example, instead of inspecting every item—such as checking every loaf of bread for color and shape—SQC uses sampling. A factory may inspect 10 out of every 1,000 bolts and use control charts to assess whether the process remains stable.

In Statistical Quality Control (SQC), variation in a process is classified into two types: common causes and special causes.

- **Common causes** are natural, random variations that are always present in a stable process, such as slight differences in dough thickness in a bakery.
- In contrast, **special causes** are unexpected variations due to specific, identifiable factors like equipment failure, human error, or defective raw materials.

The aim of SQC is to quickly detect special causes of variation, take corrective action, and maintain a stable, high-quality process.

13.1.3 Core Activities in SQC

SQC typically involves three main steps:

1. **Systematic data collection and graphical representation:** Gathering accurate quality data and visualizing it using charts and graphs.
2. **Data analysis:** Using statistical tools to detect trends, patterns, or deviations from specifications.
3. **Corrective action:** Taking practical engineering or managerial decisions based on analysis results.

These steps help prevent defects, reduce waste, and promote a culture of quality control throughout the production system.

13.1.4 Benefits of SQC

Some key benefits of implementing SQC in manufacturing and service industries include:

1. **Early detection of problems:** SQC helps in identifying defects and unusual variations in the process at an early stage, enabling timely corrective actions and preventing large-scale production issues.
2. **Improved product quality:** Continuous monitoring ensures that the products conform to quality standards, leading to higher consistency and improved customer satisfaction.
3. **Reduced waste and cost:** By minimizing defective output, SQC reduces the need for rework and scrap, thereby lowering production costs and improving efficiency.
4. **Data-driven decision making:** SQC uses statistical evidence for process evaluation, eliminating guesswork and personal bias, and promoting informed decision-making.

13.1.5 Key Statistical Tools in SQC

Several statistical tools are commonly used in SQC for controlling and improving quality:

- **Frequency distribution:** A frequency distribution is a summary of how often each value of a quality characteristic appears in the sample. For example, if you measure the diameter of 100 bolts, a frequency distribution shows how many fall within specific size intervals.
- **Control charts:** Control charts are graphs used to monitor a quality characteristic over time. They help determine whether a process is operating within statistically acceptable limits.
For example, an \bar{X} chart may track the average weight of packed sugar bags. If the average weight stays within control limits, the process is considered stable. Points outside the limits or unusual patterns indicate special causes of variation.
There are two main categories of control charts, those that display variables (continuous variables) data, and those that display attribute data.
- **Acceptance sampling:** Acceptance sampling involves inspecting a sample from a production lot to decide whether to accept or reject the entire lot. It is widely used when 100% inspection is costly or impractical.

- **Data analysis techniques:** SQC also includes advanced data analysis techniques for problem-solving and process improvement, such as:
 - Tolerance analysis
 - Correlation and regression analysis
 - Analysis of variance (ANOVA)
 - Root cause analysis

These tools help engineers and managers understand sources of variation, identify potential improvements, and make informed decisions.

13.2 Control Charts for Variables

Control charts for variables are specialized graphical tools in Statistical Process Control (SPC) designed to track quantitative measurements—such as dimensions, weights, temperatures or pressures—over time. By carefully scrutinizing the chart, a quality control engineer can identify any potential problems with the production process. A control chart consists of three parallel lines:

- **Center line (CL)**
- **Upper Control Limit (UCL)**
- **Lower Control Limit (LCL)**

At the center of a control chart lies the target line, representing the process mean. Positioned symmetrically above and below this line are the Upper Control Limit (UCL) and Lower Control Limit (LCL), typically set at three standard deviations (σ) from the mean¹. The 3σ limits are based on the fact that for a normal distribution, 99.73% of values lie within three standard deviations of the mean.

If the process is in statistical control, the sample means—taken at regular intervals and plotted on the control chart—will fall randomly within the control limits (LCL and UCL) with a probability of approximately 0.9973, assuming a normal distribution and 3-sigma limits.

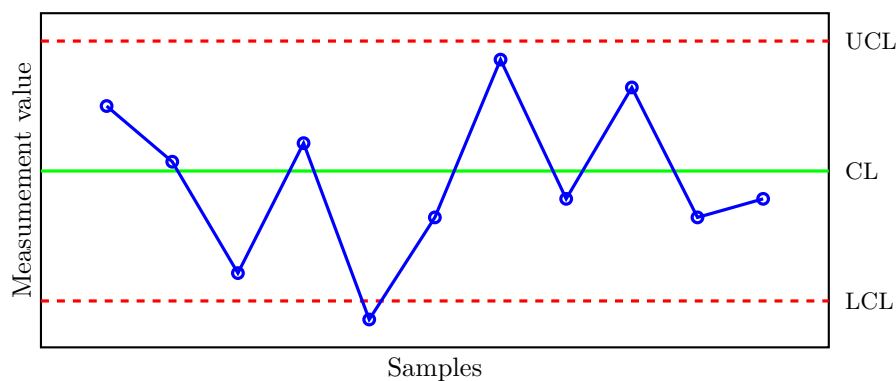


Figure 13.1: *Control Chart.*

Common types of control charts for variables include:

- The \bar{X} (mean) and the R (range) charts.
- The \bar{X} (mean) and the S (standard deviation) charts.
- Individual observations and the R (range) charts.

¹Additionally, some control charts include intermediate lines known as **warning limits**, placed at one and two standard deviations above and below the mean.

13.2.1 \bar{X} and R Charts

For monitoring the process mean and variability simultaneously, we use \bar{X} (sample mean) and R (range) charts.

- **\bar{X} Chart**

This chart tracks the **average value** of each sample. For a sample of size n , the sample mean is:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

And the range is

$$R = \max_i X_i - \min_i X_i$$

After collecting many such samples, we calculate the **grand mean** $\bar{\bar{X}}$, which is the average of all sample means, and the average range \bar{R} , which is the average of all ranges.

The control limits for the \bar{X} chart are:

$$\begin{aligned} UCL &= \bar{\bar{X}} + A_2 \bar{R} \\ CL &= \bar{\bar{X}} \\ LCL &= \bar{\bar{X}} - A_2 \bar{R} \end{aligned}$$

Here, A_2 is a constant that depends on the sample size n (available from statistical control chart tables 13.1).

- **R Chart**

This chart tracks the **range** of each sample, which reflects how much variability exists within each sample. The range for a sample is:

$$R = \max_i X_i - \min_i X_i$$

After collecting several samples, we compute the average range \bar{R} and use it to set the control limits:

$$\begin{aligned} UCL &= D_4 \bar{R} \\ CL &= \bar{R} \\ LCL &= D_3 \bar{R} \end{aligned}$$

The constants D_3 and D_4 also depend on the sample size n , and are provided in control chart reference table 13.1.

n	A_2	D_3	D_4
2	1.880	0	3.267
3	1.023	0	2.574
4	0.729	0	2.282
5	0.577	0	2.114
6	0.483	0	2.004
7	0.419	0.076	1.924
8	0.373	0.136	1.864
9	0.337	0.184	1.816
10	0.308	0.223	1.777

Table 13.1: Control chart constants for \bar{X} and R charts.

Example

A manufacturing process produces metal rods. Five samples are taken every hour, and the diameter is measured. Calculate control limits for \bar{X} and R charts.

Given data for 10 subgroups:

Table 13.2: Sample Data for Rod Diameter (mm)

Subgroup	X_1	X_2	X_3	X_4	X_5	\bar{X}	R
1	10.2	10.1	10.3	10.0	10.1	10.14	0.3
2	10.0	10.2	10.1	10.3	10.2	10.16	0.3
3	10.1	10.0	10.2	10.1	10.0	10.08	0.2
4	10.3	10.1	10.2	10.2	10.1	10.18	0.2
5	10.0	10.1	10.0	10.2	10.1	10.08	0.2

Calculations:

$$\bar{\bar{X}} = \frac{10.14 + 10.16 + 10.08 + 10.18 + 10.08}{5} = 10.13$$

$$\bar{R} = \frac{0.3 + 0.3 + 0.2 + 0.2 + 0.2}{5} = 0.24$$

For $n = 5$: $A_2 = 0.577$, $D_3 = 0$, $D_4 = 2.114$

\bar{X} Chart limits:

$$UCL = 10.13 + 0.577 \times 0.24 = 10.27$$

$$CL = 10.13$$

$$LCL = 10.13 - 0.577 \times 0.24 = 9.99$$

R Chart limits:

$$UCL = 2.114 \times 0.24 = 0.51$$

$$CL = 0.24$$

$$LCL = 0 \times 0.24 = 0$$

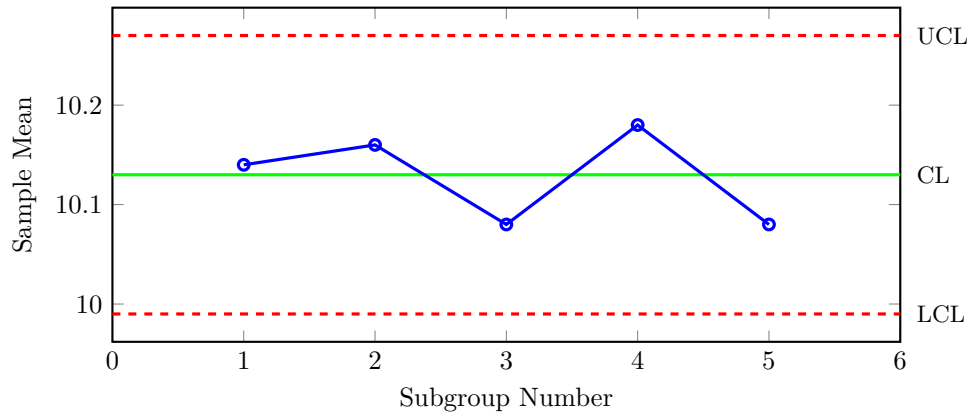


Figure 13.2: \bar{X} control chart for rod diameter.

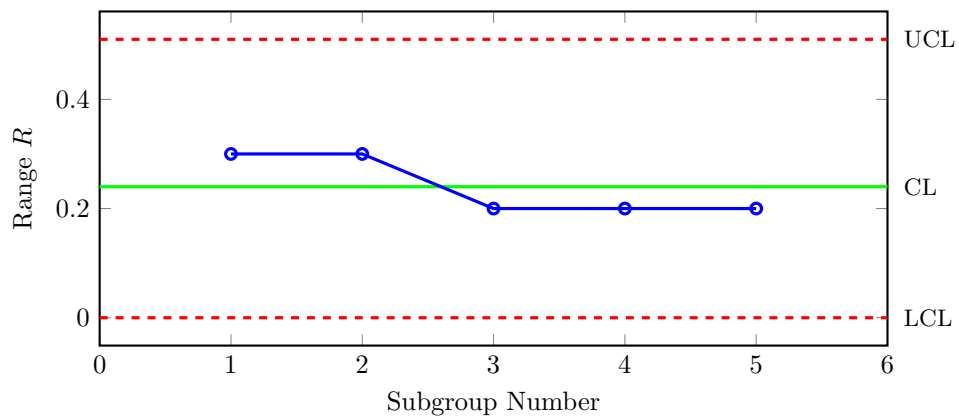


Figure 13.3: R control chart for rod diameter.

13.2.2 \bar{X} and S Charts

When your subgroup size is fairly large (usually $n \geq 10$), the sample standard deviation S gives a more reliable measure of variability than the range R . We therefore pair the \bar{X} chart with an S chart.

- **\bar{X} Chart (Mean Chart)**

This chart still monitors the average of each sample. For a subgroup of size n ,

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

You need to compute the standard deviation in each sample as well.

$$S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2}$$

After you collect many such samples, you need to compute the grand mean $\bar{\bar{X}}$ and average standard deviation \bar{S} over all the samples collected.

The control limits for the \bar{X} chart are then

$$\begin{aligned}
UCL &= \bar{\bar{X}} + A_3 \bar{S} \\
CL &= \bar{\bar{X}} \\
LCL &= \bar{\bar{X}} - A_3 \bar{S}
\end{aligned}$$

where A_3 is a constant depending on n (see Table 13.3).

- **S Chart (Standard Deviation Chart)**

This chart monitors the spread within each subgroup by plotting its standard deviation

$$S = \sqrt{\frac{1}{n-1} \sum_i^n (X_i - \bar{X})^2}$$

After collecting several samples we compute the the average \bar{S} and use it to set the control limits.

$$\begin{aligned}
UCL &= B_4 \bar{S} \\
CL &= \bar{S} \\
LCL &= B_3 \bar{S}
\end{aligned}$$

where B_3 and B_4 depend on n (see Table 13.3).

Note: Always check the S chart first. If the process variability is out of control, then any signals on the \bar{X} chart may be misleading.

n	A_3	B_3	B_4
10	0.975	0.284	1.716
11	0.927	0.321	1.679
12	0.886	0.354	1.646
13	0.850	0.382	1.618
14	0.817	0.406	1.594
15	0.789	0.428	1.572
16	0.763	0.448	1.552
17	0.739	0.466	1.534
18	0.718	0.482	1.518
19	0.698	0.497	1.503
20	0.680	0.510	1.490
\vdots	\vdots	\vdots	\vdots

Table 13.3: Control chart constants for \bar{X} and S charts.

Example

A manufacturing process produces metal rods. Ten measurements are taken every hour, and the diameter is recorded. Calculate control limits for \bar{X} and S charts.

Given data for 5 subgroups:

Subgroup	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9	X_{10}	\bar{X}	S
1	10.1	10.0	10.2	10.1	10.3	10.2	10.1	10.0	10.2	10.1	10.13	0.10
2	10.0	10.2	10.1	10.3	10.2	10.1	10.0	10.2	10.1	10.3	10.15	0.11
3	10.1	10.0	10.2	10.1	10.0	10.1	10.0	10.1	10.2	10.1	10.09	0.07
4	10.3	10.1	10.2	10.2	10.1	10.2	10.3	10.1	10.2	10.2	10.19	0.07
5	10.0	10.1	10.0	10.2	10.1	10.0	10.1	10.0	10.2	10.1	10.08	0.07

Table 13.4: *Sample data for rod diameter (mm).*

Calculations:

$$\bar{\bar{X}} = \frac{10.13 + 10.15 + 10.09 + 10.19 + 10.08}{5} = 10.13$$

$$\bar{S} = \frac{0.10 + 0.11 + 0.07 + 0.07 + 0.07}{5} = 0.084$$

For $n = 10$: $A_3 = 0.266$, $B_3 = 0.284$, $B_4 = 1.716$

\bar{X} Chart limits:

$$UCL = 10.13 + 0.266 \times 0.084 \approx 10.15,$$

$$CL = 10.13,$$

$$LCL = 10.13 - 0.266 \times 0.084 \approx 10.11.$$

S Chart limits:

$$UCL = 1.716 \times 0.084 \approx 0.14,$$

$$CL = 0.084,$$

$$LCL = 0.284 \times 0.084 \approx 0.024.$$

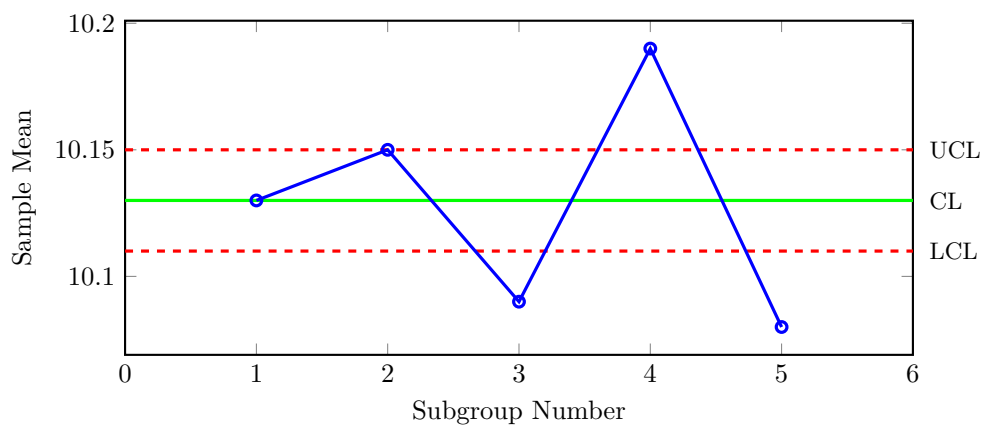


Figure 13.4: \bar{X} control chart for rod diameter.

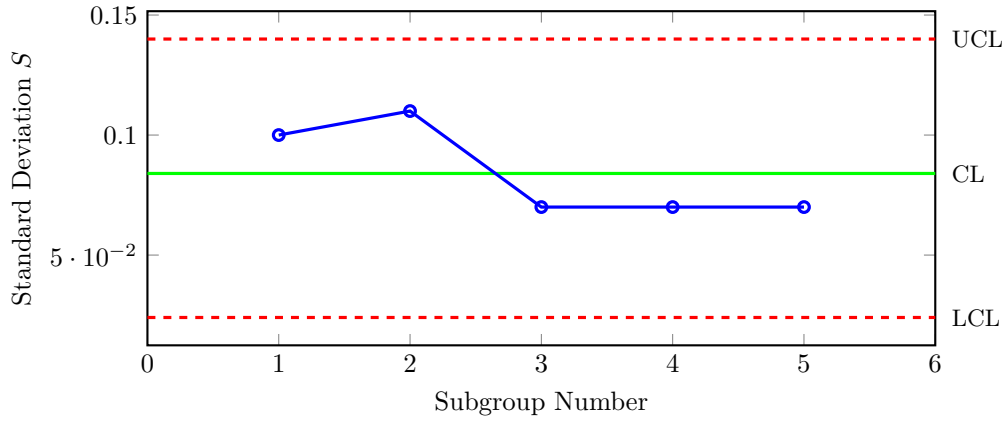


Figure 13.5: S control chart for rod diameter.

13.3 Control Charts for Attributes

Attribute data represents characteristics that can be counted or classified (defective/non-defective, pass/fail).

13.3.1 p Chart (Fraction Defective)

The p chart monitors the fraction of defective items in a sample.

Let D = number of defective items in a sample of size n

$$p = \frac{D}{n}$$

Control limits:

$$UCL = \bar{p} + 3\sqrt{\frac{\bar{p}(1 - \bar{p})}{n}} \quad (13.1)$$

$$CL = \bar{p} \quad (13.2)$$

$$LCL = \bar{p} - 3\sqrt{\frac{\bar{p}(1 - \bar{p})}{n}} \quad (13.3)$$

where \bar{p} is the average fraction defective.

Example: Defective Products

A company inspects 100 items daily. Data for 10 days:

Table 13.5: Daily Inspection Data

Day	Sample Size	Defective Items	p
1	100	5	0.05
2	100	3	0.03
3	100	7	0.07
4	100	4	0.04
5	100	6	0.06
6	100	2	0.02
7	100	8	0.08
8	100	3	0.03
9	100	5	0.05
10	100	7	0.07

$$\bar{p} = \frac{5 + 3 + 7 + 4 + 6 + 2 + 8 + 3 + 5 + 7}{10} = \frac{50}{1000} = 0.05$$

Control limits:

$$UCL = 0.05 + 3\sqrt{\frac{0.05(1 - 0.05)}{100}} = 0.05 + 3\sqrt{0.000475} = 0.115 \quad (13.4)$$

$$CL = 0.05 \quad (13.5)$$

$$LCL = 0.05 - 3\sqrt{\frac{0.05(1 - 0.05)}{100}} = -0.015 \approx 0 \quad (13.6)$$

13.3.2 c Chart (Number of Defects)

The c chart monitors the number of defects per unit when the sample size is constant.

For a Poisson distribution with parameter λ :

$$P(c) = \frac{e^{-\lambda} \lambda^c}{c!}$$

Control limits:

$$UCL = \bar{c} + 3\sqrt{\bar{c}} \quad (13.7)$$

$$CL = \bar{c} \quad (13.8)$$

$$LCL = \bar{c} - 3\sqrt{\bar{c}} \quad (13.9)$$

13.4 Pareto Chart in Statistical Quality Control

A **Pareto chart** is a bar graph that displays the frequency or impact of problems or defects in descending order, alongside a cumulative percentage line. It is based on the **Pareto principle** (also known as the 80/20 rule), which suggests that roughly 80% of the problems arise from 20% of the causes.

In SQC, Pareto charts help prioritize improvement efforts by identifying the most significant sources of defects or errors.

Example: Suppose a factory inspects 1000 products and records the number of occurrences for various types of defects. The data is shown below:

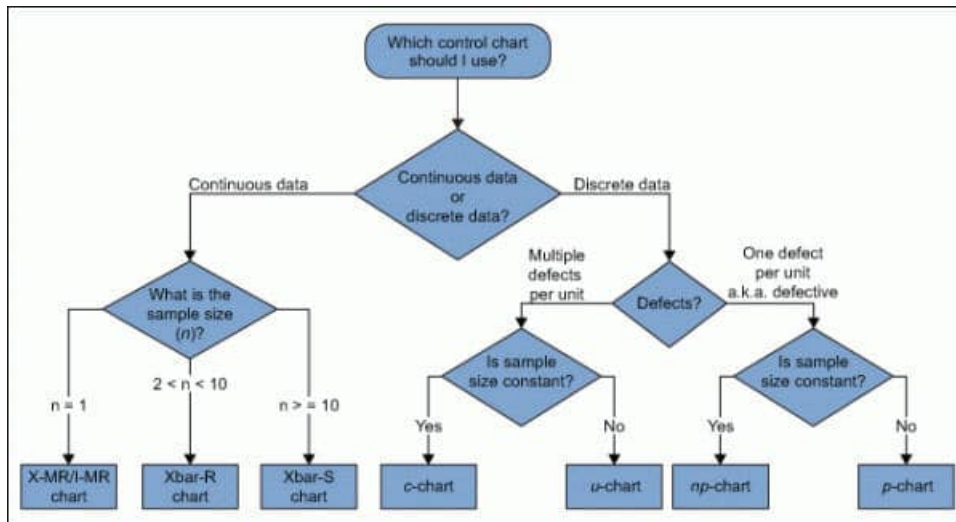


Figure 13.6: Control chart decision tree.

Defect Type	Frequency	Percentage (%)	Cumulative %
Surface Scratch	300	30.0	30.0
Paint Defect	250	25.0	55.0
Loose Screw	150	15.0	70.0
Crack	100	10.0	80.0
Missing Label	80	8.0	88.0
Bent Part	70	7.0	95.0
Other	50	5.0	100.0

Table 13.6: Frequency of Defect Types in Production.

The Pareto chart for the above data is shown below:

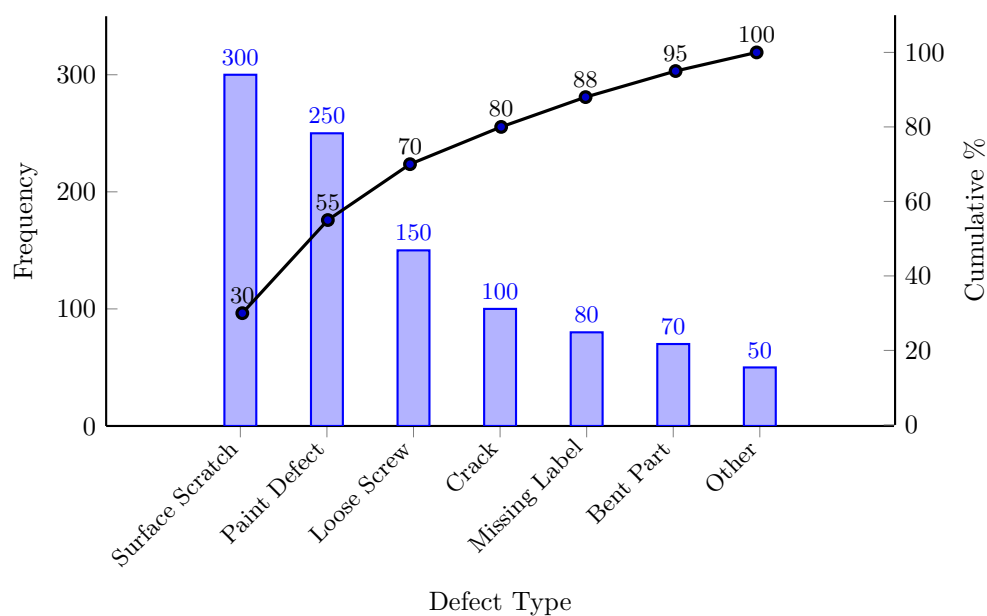


Figure 13.7: Pareto chart for defect types.

Interpretation

From the Pareto chart, it is evident that the majority of the defects (70%) come from the first three types: surface scratches, paint defects, and loose screws. Therefore, quality improvement efforts should focus on these issues first to achieve the greatest impact with minimal effort.

13.5 Process Capability Analysis

The **process capability** allow us to measure or quantify the capability of a process to outputs that meet our design specification.

The **Process Capability Index** is the specific measure we use to quantify process capability.

Every manufacturing or production process has some natural variability, and the goal is to ensure that this variability stays within the acceptable range set by customer or engineering requirements. The Process capability index helps us understand whether the process is capable of producing results that consistently fall within these limits.

For example, consider a process that manufactures plastic bottles with an acceptable weight range of 95 to 105 grams. If the process consistently produces bottles between 96 and 104 grams, it is considered capable. But if the weights vary from 92 to 108 grams, many units may fall outside the limits, indicating poor capability. The Process Capability Index summarizes this relationship between process variation and specification limits in a single number.

13.5.1 Statistical Control and Its Importance

A process is said to be in **statistical control** when it operates consistently over time, with variation coming only from natural, random factors known as **common causes**. These are small, unavoidable fluctuations that are inherent to the system.

In contrast, there are **special causes** which are unusual variations caused by specific issues like operator changes, equipment drift, raw material shifts, and environmental conditions. A process in statistical control shows no unexpected patterns or points outside control limits on a control chart. This is an essential condition or using capability indices like C_k and C_{pk} reliably which we will discuss in coming subsections.

13.5.2 Understanding Short-term vs. Long-term Variation

To interpret process capability indices meaningfully, we must first distinguish between two types of process variation: **short-term** and **long-term** variation. These two reflect different time horizons and include different sources of variability.

Short-term Variation (Within-subgroup)

Short-term variation is measured over a limited time frame—such as within a shift or over a few hours. It represents the natural, inherent variability in the process when no major changes or external disturbances occur. This variation is quantified using the **within-subgroup standard deviation** (σ_{within}) and assumes that the process is stable and under statistical control. It captures only common-cause (natural) variation.

Long-term Variation (Overall)

Long-term variation is observed over extended periods such as days, weeks, or months. Unlike short-term variation, it incorporates additional sources of variability that occur between subgroups or over

time. This variation is quantified using the **overall standard deviation** (σ_{overall}). It includes both common-cause and special-cause variation

Example: Coffee Machine

Imagine a coffee machine programmed to dispense exactly 8 oz of coffee per cup.

- **Short-term:** If you record 20 cups dispensed within 10 minutes, the weights might range from 7.9 to 8.1 oz with $\sigma = 0.05$ oz. This shows the machine's precision in a stable, controlled environment.
- **Long-term:** Over three months of operation, including different shifts, bean types, and maintenance cycles, the weights might range from 7.7 to 8.3 oz with $\sigma_{\text{overall}} = 0.15$ oz. This reflects the process's real-world variation.

13.5.3 C_p Index

The C_p index measures the **potential capability of a process** to produce output within the given specification limits, assuming the process is stable (in statistical control) and perfectly centered (i.e., the mean lies exactly midway between the upper and lower specification limits):

$$C_p = \frac{USL - LSL}{6\sigma_{\text{within}}}$$

where,

- USL = Upper specification limit,
- LSL = Lower specification limit,
- σ_{within} is the short-term process standard deviation.

Importantly, C_p is calculated using the **within-subgroup standard deviation** σ_{within} , which reflects **short-term variation**—the natural, inherent fluctuation observed when the process is operating without any external disturbances or special causes. As a result, it reflects the potential capability or maximum capability the process could achieve—if all special causes were eliminated and only inherent variation remained.

The denominator $6\sigma_{\text{within}}$ represents the total spread of the process assuming that it follows a normal distribution (covering approximately 99.73% of the data). The numerator $USL - LSL$ is the allowable spread as specified by the customer or design requirements.

Interpretation

- $C_p = 1$: The process spread exactly fits the specification limits.
- $C_p > 1$: The process variation is smaller than the specification limits—suggesting the process has potential to perform well.
- $C_p < 1$: The process variation is wider than the specification limits—indicating that a significant portion of output may fall outside the acceptable range.

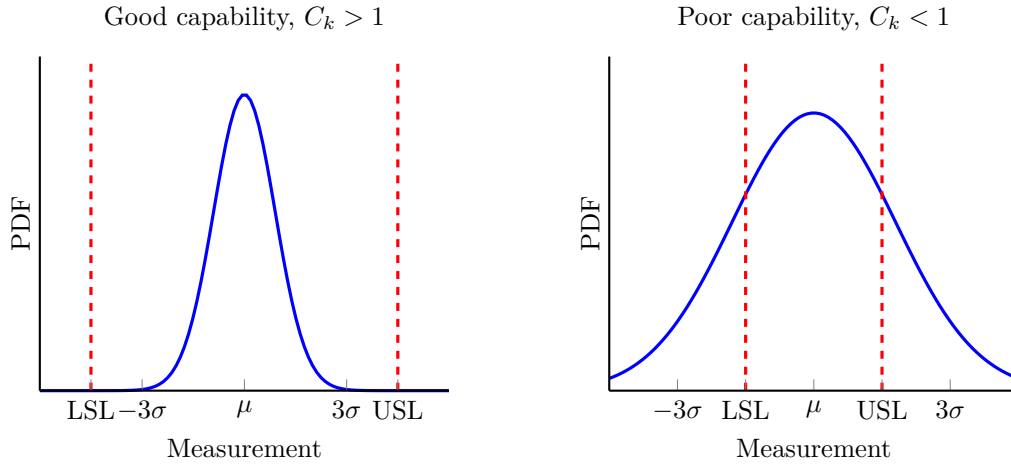


Figure 13.8: Normal distributions illustrating process capability index C_k relative to specification limits (LSL and USL) and process variation ($\pm 3\sigma_{within}$).

C_p Value	Process Capability
$C_p < 1.0$	Inadequate
$1.0 \leq C_p < 1.33$	Marginal
$1.33 \leq C_p < 2.0$	Adequate
$C_p \geq 2.0$	Excellent

Table 13.7: Interpretation of process capability index C_p .

It is important to remember that C_p does not account for how well the process is centered. A process with a high C_p can still produce defects if the mean is not aligned with the target. To evaluate both spread and centering, other indices like C_{pk} are used.

Example

Suppose a manufacturing process has a lower specification limit of 95 mm and an upper specification limit of 105 mm. The short-term standard deviation of the process is 1 mm. Then,

$$C_p = \frac{105 - 95}{6 \times 1} = \frac{10}{6} \approx 1.67$$

This value suggests that the process has a good potential to meet specifications, provided it is properly centered.

One-Sided C_p :

Sometimes, only one-sided limit matters—either just an upper or just a lower limit. In such cases, we use a one-sided version of C_p .

If there is only an upper specification limit (USL), the formula becomes:

$$C_{p,upper} = \frac{USL - \mu_w}{3\sigma_{within}}$$

where μ_w is the process mean within subgroup.

Likewise, if only a lower specification limit (LSL) is applicable, the formula becomes:

$$C_{p,lower} = \frac{\mu_w - LSL}{3\sigma_{within}}$$

These one-sided indices are useful in scenarios where deviations in only one direction are undesirable—for example, when there is a maximum allowable concentration of a chemical, or a minimum required tensile strength in materials.

13.5.4 C_{pk} Index

The C_{pk} index measures the potential capability of a process, considering that the process mean may not be centered between the specification limits.

$$C_{pk} = \min \left(\frac{USL - \mu_w}{3\sigma_{\text{within}}}, \frac{\mu_w - LSL}{3\sigma_{\text{within}}} \right) = \min (C_{p,\text{upper}}, C_{p,\text{lower}})$$

where μ_w is the process mean within subgroup.

The formula takes the minimum of the two distances (in terms of standard deviations) from the process mean to each specification limit. This ensures that any shift in the process mean away from the center reduces the C_{pk} value, even if C_p remains high.

Interpretation

A high C_{pk} value means the process is both low in variability and well-centered. A low C_{pk} indicates either high variation or that the process mean is deviated from the center, increasing the risk of defects.

C_{pk} Value	Process Capability
$C_{pk} < 1.0$	Inadequate
$1.0 \leq C_{pk} < 1.33$	Marginal
$1.33 \leq C_{pk} < 2.0$	Adequate
$C_{pk} \geq 2.0$	Excellent

Table 13.8: Interpretation of process capability index C_{pk} .

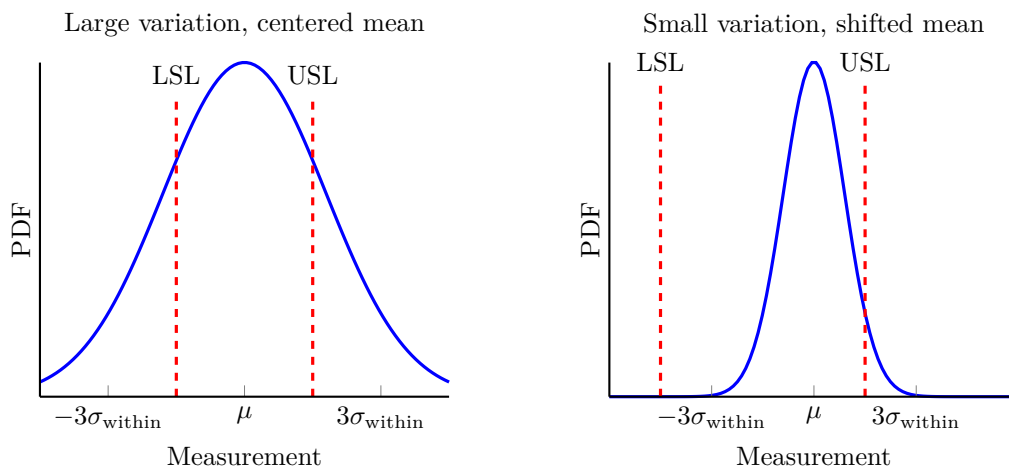


Figure 13.9: Illustration of poor C_{pk} scenarios: (Left) large variation with centered mean and (Right) small variation with shifted mean causing poor C_{pk} values.

Example

Suppose a manufacturing process has a lower specification limit of 95 mm and an upper specification limit of 105 mm. The process mean is 102 mm, and the standard deviation is 1 mm (within subgroup). Then,

$$\frac{USL - \mu_w}{3\sigma_{\text{within}}} = \frac{105 - 102}{3 \times 1} = \frac{3}{3} = 1, \quad \frac{\mu_w - LSL}{3\sigma_{\text{within}}} = \frac{102 - 95}{3 \times 1} = \frac{7}{3} \approx 2.33$$
$$C_{pk} = \min(1, 2.33) = 1.0$$

Although the process has a relatively low variation, the mean is not centered within the specification limits. The value $C_{pk} = 1.0$ indicates that the process is only marginally capable and may produce defects near the upper limit.

Relationship Between C_{pk} and Defect Rate

Assuming the process output is normally distributed, the process capability index can be translated into the proportion of defective items produced. For example, a $C_{pk} = 1.00$ corresponds to 3-sigma capability, implying about 0.27% of the output lies outside the specification limits, or 2,700 defective parts per million (ppm).

Let the specification limits be USL and LSL , and the process mean and standard deviation be μ and σ , respectively. Then, the defect rate can be estimated using the tails of the normal distribution:

$$\begin{aligned} P(\text{Defect}) &= P(X < LSL) + P(X > USL) \\ &= 1 + \Phi\left(\frac{LSL - \mu_w}{\sigma_{\text{within}}}\right) - \Phi\left(\frac{USL - \mu_w}{\sigma_{\text{within}}}\right) \\ &= 1 + \underbrace{\Phi(-3C_{p,\text{lower}})}_{1 - \Phi(3C_{p,\text{lower}})} - \Phi(3C_{p,\text{upper}}) \\ &= 2 - \Phi(3C_{p,\text{lower}}) - \Phi(3C_{p,\text{upper}}) \end{aligned}$$

where $\Phi(\cdot)$ is the CDF of the standard normal distribution. When the process is centered $C_{p,\text{upper}} = C_{p,\text{lower}} = C_{pk}$,

$$P(\text{Defect}) = 2(1 - \Phi(3C_{pk}))$$

13.5.5 P_p Index

The P_p index measures the **actual overall process performance capability**, using the total variation observed in the process data. It does not consider process centering.

$$P_p = \frac{USL - LSL}{6\sigma_{\text{overall}}}$$

where σ_{overall} is the calculated overall standard deviation of the process using all the data. Note that the only difference between C_p and P_p is how the process variation is calculated. The calculated standard deviation includes all the past data and does not take into account if the process is in statistical control. If the process is not in statistical control, the value of the calculated standard deviation will often be inflated.

Interpretation

A higher P_p means the process variation is small relative to the specification width. However, since P_p ignores centering, it might overestimate capability if the process mean is off-center.

Interpreting C_p and P_p together helps evaluate both the potential and actual performance of a process.

- When $C_p \approx P_p$, the process is likely **in statistical control**, meaning the actual performance closely matches the potential capability. If both values are low, the process likely suffers from high variation and requires improvement in design or control. If both values are high, it indicates the process is capable and performing well.
- When $C_p > P_p$, the process has **good potential capability**, but is currently **unstable** due to the presence of special cause variation. In this case, efforts should be directed toward identifying and eliminating those special causes to stabilize the process.
- The condition $C_p < P_p$ is a **rare and unexpected situation** that may suggest an error in calculation or a misestimation of the standard deviation. It is advisable to recheck the data, assumptions, and analysis methods used.

Example

Suppose the process has specification limits 95 mm and 105 mm, and the sample standard deviation $\sigma_{\text{overall}} = 1.2$ mm. Then,

$$P_p = \frac{105 - 95}{6 \times 1.2} = \frac{10}{7.2} \approx 1.39$$

This suggests the process has adequate performance, but this does not consider whether the process mean is centered.

13.5.6 P_{pk} Index

The P_{pk} index evaluates the actual process performance considering both variation and centering, using overall process data.

$$P_{pk} = \min \left(\frac{USL - \mu_o}{3\sigma_{\text{overall}}}, \frac{\mu_o - LSL}{3\sigma_{\text{overall}}} \right)$$

where μ_o and σ_{overall} are the values of the overall process mean and the overall process standard deviation respectively.

Interpretation

A high P_{pk} means the process performs well, with low variation and good centering. A low P_{pk} indicates either high variation or poor centering, increasing defect risk.

Interpreting C_{pk} and P_{pk} together provides insight into both the potential and actual performance of a process with respect to centering.

- When $C_{pk} \approx P_{pk}$, the process is likely **in statistical control** and stable. When both values are low, the process is not well-centered or has high variability. If both are high and close, the process is performing excellently and is well-centered.
- When $C_{pk} > P_{pk}$, the process **has the potential to be well-centered and low in defects**, but is currently affected by special causes of variation.
- The condition $C_{pk} < P_{pk}$ is **very rare** which may indicate calculation error or incorrect standard deviation used.

Example

Suppose a process has $LSL = 95$ mm, $USL = 105$ mm, sample mean $\bar{x} = 102$ mm, and sample standard deviation $\sigma_{\text{overall}} = 1$ mm. Then,

$$\frac{USL - \mu_o}{3\sigma_{\text{overall}}} = \frac{105 - 102}{3 \times 1} = 1, \quad \frac{\mu_o - LSL}{3s} = \frac{102 - 95}{3 \times 1} = \frac{7}{3} \approx 2.33$$

$$P_{pk} = \min(1, 2.33) = 1.0$$

This indicates marginal performance capability, similar to the C_{pk} case but based on sample estimates rather than true parameters.

13.5.7 Recommended Actions for Process Capability

Cp/Pp	Cpk/Ppk	Action Required
≥ 1.33	≥ 1.33	Continue monitoring
	$1.0 - 1.33$	Improve centering
	< 1.0	Improve centering immediately
$1.0 - 1.33$	≥ 1.33	Reduce variation
	$1.0 - 1.33$	Reduce variation & improve centering
	< 1.0	Reduce variation & improve centering immediately
< 1.0	≥ 1.33	Major process redesign
	$1.0 - 1.33$	Major process redesign
	< 1.0	Complete process overhaul

Table 13.9: Recommended actions based on capability indices.

13.5.8 A Case Study of Process Capability Analysis

Problem Statement

A steel rod manufacturing company produces rods with a target diameter of 25.0 mm. The specification limits are:

- Lower Specification Limit (LSL) = 24.7 mm
- Upper Specification Limit (USL) = 25.3 mm
- Target Value = 25.0 mm

Quality engineers collected 100 measurements over 20 subgroups (5 measurements per subgroup) to assess process capability as follows:

Subgroup	X_1	X_2	X_3	X_4	X_5	\bar{X}	R
1	24.95	25.02	24.98	25.01	24.97	24.986	0.07
2	25.03	24.96	25.00	24.99	25.02	25.000	0.07
3	24.92	25.05	24.98	25.01	24.94	24.980	0.13
4	25.08	24.97	25.03	24.99	25.01	25.016	0.11
5	24.99	25.04	24.96	25.02	25.00	25.002	0.08
6	25.01	24.95	25.07	24.98	25.03	25.008	0.12
7	24.94	25.01	24.97	25.05	24.99	24.992	0.11
8	25.06	24.98	25.02	24.96	25.04	25.012	0.10
9	24.97	25.03	24.99	25.01	24.95	24.990	0.08
10	25.05	24.99	25.01	24.97	25.08	25.020	0.11
11	24.96	25.02	25.04	24.98	25.00	25.000	0.08
12	25.01	24.94	25.06	25.02	24.97	25.000	0.12
13	24.98	25.05	24.99	25.03	25.01	25.012	0.07
14	25.04	24.97	25.00	24.96	25.09	25.012	0.13
15	24.93	25.01	24.98	25.04	24.99	24.990	0.11
16	25.07	24.99	25.02	24.95	25.03	25.012	0.12
17	24.98	25.06	25.00	25.01	24.97	25.004	0.09
18	25.02	24.94	25.05	24.99	25.01	25.002	0.11
19	24.99	25.03	24.96	25.07	25.00	25.010	0.11
20	25.01	24.98	25.04	24.97	25.05	25.010	0.08

Table 13.10: *Steel rod diameter measurements (mm)*

Using the collected data, evaluate the process capability.

Overall Mean

$$\bar{\bar{X}} = \frac{24.986 + 25.000 + 24.980 + \dots}{20} = \frac{500.202}{20} = 25.001 \text{ mm}$$

Average Range

$$\bar{R} = \frac{0.07 + 0.07 + 0.13 + \dots}{20} = \frac{2.01}{20} = 0.1005 \text{ mm}$$

Short-term Variation (Within Subgroups)

For short-term capability indices, we estimate the within-subgroup standard deviation using:

$$\hat{\sigma}_{\text{within}} = \frac{\bar{R}}{d_2}$$

where $d_2 = 2.326$ for subgroup size $n = 5$.

$$\hat{\sigma}_{\text{within}} = \frac{0.1005}{2.326} = 0.0432 \text{ mm}$$

Long-term Variation (Overall Process)

For long-term capability indices, we use the overall standard deviation from all $n = 100$ measurements:

$$\hat{\sigma}_{\text{overall}} = 0.0376 \text{ mm}$$

Process Capability (C_p)

$$C_p = \frac{USL - LSL}{6\hat{\sigma}_{\text{within}}} = \frac{25.3 - 24.7}{6 \times 0.0432} = \frac{0.6}{0.2592} = 2.31$$

Process Capability Index (C_{pk})

$$C_{pk} = \min \left(\frac{USL - \bar{\bar{X}}}{3\hat{\sigma}_{\text{within}}}, \frac{\bar{\bar{X}} - LSL}{3\hat{\sigma}_{\text{within}}} \right)$$

$$\begin{aligned} C_{pk} &= \min \left(\frac{25.3 - 25.001}{3 \times 0.0432}, \frac{25.001 - 24.7}{3 \times 0.0432} \right) \\ &= \min \left(\frac{0.299}{0.1296}, \frac{0.301}{0.1296} \right) = \min(2.31, 2.32) = 2.31 \end{aligned}$$

Process Performance (P_p)

$$P_p = \frac{USL - LSL}{6\hat{\sigma}_{\text{overall}}} = \frac{25.3 - 24.7}{6 \times 0.0376} = \frac{0.6}{0.2256} = 2.66$$

Process Performance Index (P_{pk})

$$P_{pk} = \min \left(\frac{USL - \bar{\bar{X}}}{3\hat{\sigma}_{\text{overall}}}, \frac{\bar{\bar{X}} - LSL}{3\hat{\sigma}_{\text{overall}}} \right)$$

$$\begin{aligned} P_{pk} &= \min \left(\frac{25.3 - 25.001}{3 \times 0.0376}, \frac{25.001 - 24.7}{3 \times 0.0376} \right) \\ &= \min \left(\frac{0.299}{0.1128}, \frac{0.301}{0.1128} \right) = \min(2.65, 2.67) = 2.65 \end{aligned}$$

Interpretation of Results

Index	Value	Interpretation
C_p	2.31	Excellent potential capability (short-term, assuming centered)
C_{pk}	2.31	Excellent potential capability (short-term, actual centering)
P_p	2.66	Excellent overall performance (long-term, assuming centered)
P_{pk}	2.65	Excellent overall performance (long-term, actual centering)

Table 13.11: *Process capability results based on short-term and long-term data.*

Defect Rate Estimation

Based on the capability indices, we can estimate defect rates:

For $C_{pk} = 2.31$:

$$\text{Defect Rate} \approx 2 \times (1 - \Phi(C_{pk})) = 2 \times (1 - \Phi(2.31)) \approx 0.021 \text{ ppm}$$

This means approximately 0.021 parts per million would be outside specifications - essentially zero defects.

Conclusion

The steel rod manufacturing process demonstrates excellent capability with all indices well above acceptable thresholds. The process is:

- Excellent process capability (all indices > 2.0).
- Well-centered on target.
- Stable over time.
- Producing virtually zero defects (less than 0.002 ppm).

The primary focus should be on maintaining this excellent performance. No corrective actions are needed, but continuous monitoring and documentation of best practices are recommended.

13.6 Acceptance Sampling

Acceptance sampling is a statistical method used to determine whether to accept or reject a batch of products based on inspection of a sample.

13.6.1 Single Sampling Plans

A single sampling plan is characterized by (n, c) where:

- n = sample size
- c = acceptance number

Decision rule: Accept the lot if the number of defective items $d \leq c$

13.6.2 Operating Characteristic (OC) Curve

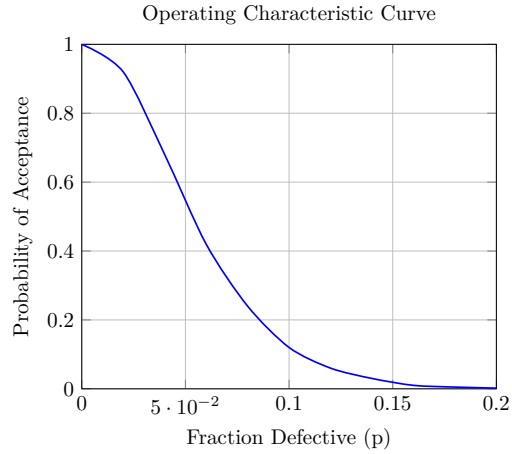
The OC curve shows the probability of accepting a lot as a function of the lot fraction defective p .

For a binomial distribution:

$$P_a(p) = \sum_{d=0}^c \binom{n}{d} p^d (1-p)^{n-d}$$

Example: OC Curve Calculation

For a sampling plan $(n = 50, c = 2)$:



13.6.3 Producer's and Consumer's Risks

- **Producer's Risk (α):** Probability of rejecting a good lot
- **Consumer's Risk (β):** Probability of accepting a bad lot

Typically:

- AQL (Acceptable Quality Level): $\alpha = 0.05$
- LTPD (Lot Tolerance Percent Defective): $\beta = 0.10$

13.7 Reliability and Life Testing

Reliability is the probability that a system performs successfully during a specified time interval under specified conditions.

13.7.1 Reliability Function

The reliability function $R(t)$ is defined as:

$$R(t) = P(T > t) = 1 - F(t)$$

where T is the random variable representing time to failure and $F(t)$ is the cumulative distribution function.

13.7.2 Failure Rate Function

The hazard rate or failure rate function $h(t)$ is:

$$h(t) = \frac{f(t)}{R(t)} = \frac{f(t)}{1 - F(t)}$$

where $f(t)$ is the probability density function.

13.7.3 Common Life Distributions

Exponential Distribution

For constant failure rate λ :

$$f(t) = \lambda e^{-\lambda t}, \quad t \geq 0$$

$$R(t) = e^{-\lambda t}$$

$$h(t) = \lambda$$

Mean Time To Failure (MTTF):

$$MTTF = \frac{1}{\lambda}$$

Weibull Distribution

The Weibull distribution is widely used in reliability analysis:

$$f(t) = \frac{\beta}{\eta} \left(\frac{t}{\eta} \right)^{\beta-1} e^{-(t/\eta)^\beta}$$

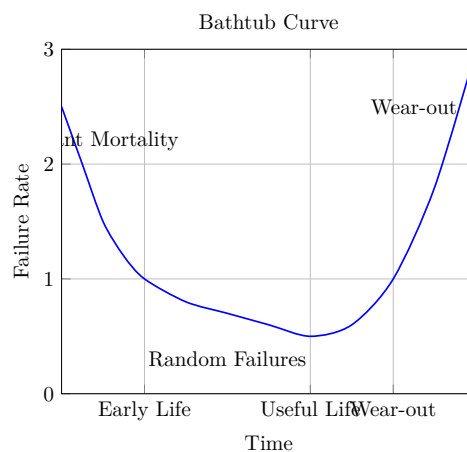
where β is the shape parameter and η is the scale parameter.

$$R(t) = e^{-(t/\eta)^\beta}$$

$$h(t) = \frac{\beta}{\eta} \left(\frac{t}{\eta} \right)^{\beta-1}$$

Bathtub Curve

The bathtub curve describes the failure rate pattern over product lifetime:



13.7.4 Life Testing

Complete Testing

All units are tested until failure. Provides complete information but time-consuming.

Censored Testing

Type I Censoring: Testing stops at predetermined time t_0

Type II Censoring: Testing stops after predetermined number of failures r

Accelerated Life Testing

Testing under elevated stress conditions to accelerate failures.

Arrhenius model:

$$\lambda(T) = Ae^{E_a/(kT)}$$

where T is temperature, E_a is activation energy, k is Boltzmann constant, and A is a constant.

13.7.5 System Reliability

Series System

For a series system with n components:

$$R_s(t) = \prod_{i=1}^n R_i(t)$$

Parallel System

For a parallel system with n components:

$$R_p(t) = 1 - \prod_{i=1}^n [1 - R_i(t)]$$

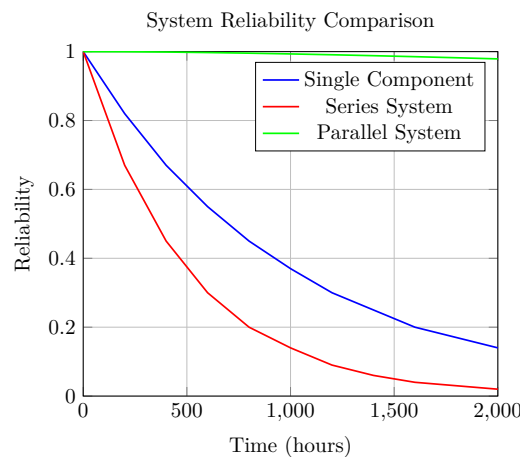
Example: System Reliability Calculation

Consider a system with 3 components in series, each with reliability $R(t) = e^{-0.001t}$:

$$R_s(t) = [e^{-0.001t}]^3 = e^{-0.003t}$$

For the same components in parallel:

$$R_p(t) = 1 - [1 - e^{-0.001t}]^3$$



13.8 Conclusion

Statistical Quality Control provides a comprehensive framework for monitoring, controlling, and improving quality in manufacturing and service processes. The key principles include:

- Use of statistical methods to distinguish between common and special causes of variation
- Implementation of control charts for continuous monitoring
- Application of process capability analysis to ensure specifications are met
- Employment of acceptance sampling for lot-by-lot decisions
- Utilization of statistical tools for quality improvement
- Application of reliability methods for long-term performance assessment

The successful implementation of SQC requires understanding of statistical concepts, proper data collection, and commitment to continuous improvement. These methods, when properly applied, lead to reduced variability, improved quality, and enhanced customer satisfaction.

Chapter 14

Index Numbers

14.1 Introduction

14.2 Introduction to Index Numbers

An **index number** is a statistical measure that shows the relative change in a variable or group of variables with respect to time, geographical location, or other characteristics. Index numbers are expressed as percentages or ratios that compare the value of a variable in a given period to its value in a base period.

To illustrate why we need index number we consider the following example. Suppose a consumer purchases two commodities: rice and cooking oil. The prices and quantities in the base year and current year are given below:

Commodity	Price (Base Year)	Price (Current Year)
Rice (per kg)	40	50
Oil (per litre)	100	120

If we want to understand how much prices have risen overall, it would be misleading to simply average the percentage increases, because it does not account for how much of each commodity is actually purchased. Instead, we fix quantities from the base period and compute a weighted average of the price increases using index numbers.

Assume the consumer bought 30 kg of rice and 10 litres of oil in the base year. Then:

- Cost in base year = $40 \times 30 + 100 \times 10 = 2200$
- Cost in current year = $50 \times 30 + 120 \times 10 = 2700$

The **Price Index Number** is calculated as:

$$\text{Index Number} = \frac{\text{Cost in Current Year}}{\text{Cost in Base Year}} \times 100 = \frac{2700}{2200} \times 100 \approx 122.73$$

This means that, overall, prices have increased by approximately 22.73% compared to the base year.

This example highlights the importance of index numbers. Instead of comparing the price of each item separately, an index number provides a single, meaningful figure that summarizes the overall change in the price level. Such measures are invaluable in economic analysis for tracking inflation, evaluating cost-of-living adjustments, comparing standards of living across regions, or measuring

industrial and agricultural output over time. They help policymakers, businesses, and consumers make informed decisions based on relative changes rather than absolute values.

Index numbers serve as powerful tools for:

- Measuring changes in economic variables over time
- Comparing economic conditions across different regions
- Making informed business and policy decisions
- Deflating monetary values to real terms

The concept of index numbers originated in the 18th century when economists needed to measure price changes over time. Today, they are extensively used in economics, business, and social sciences.

14.3 Basic Concepts and Terminology

14.3.1 Key Terms

Base Period: The period against which comparisons are made, typically assigned an index value of 100.

Current Period: The period for which the index is being calculated.

Base Year: The reference year for comparison (e.g., 2010 = 100).

Weights: Numerical values representing the relative importance of different items.

Price Relative: The ratio of current price to base price, expressed as a percentage.

Quantity Relative: The ratio of current quantity to base quantity, expressed as a percentage.

14.3.2 Mathematical Foundation

For a single commodity, the price relative is given by:

$$P_r = \frac{P_1}{P_0} \times 100 \quad (14.1)$$

Where:

$$P_r = \text{Price relative} \quad (14.2)$$

$$P_1 = \text{Price in current period} \quad (14.3)$$

$$P_0 = \text{Price in base period} \quad (14.4)$$

Similarly, for quantity relative:

$$Q_r = \frac{Q_1}{Q_0} \times 100 \quad (14.5)$$

14.4 Types of Index Numbers

14.4.1 Classification by Purpose

Price Index Numbers

Price index numbers measure the average change in prices of a basket of goods and services over time.

Table 14.1: Price Data for CPI Calculation

Item	Base Price (P_0)	Current Price (P_1)	Weight (W)	Price Relative (P_r)	Weighted Price Relative ($W \times P_r$)
Food	100	120	40	120	4800
Housing	200	240	30	120	3600
Transport	80	100	20	125	2500
Healthcare	150	180	10	120	1200
Total			100		12100

Example: Consumer Price Index (CPI) calculation

$$\text{Weighted Price Index} = \frac{\sum W \times P_r}{\sum W} = \frac{12100}{100} = 121$$

This indicates a 21% increase in prices from the base period.

Quantity Index Numbers

These measure changes in the physical quantities of goods produced, consumed, or traded.

Value Index Numbers

Value indices measure changes in the total value (price \times quantity) of goods and services.

$$\text{Value Index} = \frac{\sum P_1 Q_1}{\sum P_0 Q_0} \times 100 \quad (14.6)$$

14.5 Methods of Constructing Price Index Numbers

14.5.1 Simple Price Index Numbers

Simple Aggregative Method

$$P_{01} = \frac{\sum P_1}{\sum P_0} \times 100 \quad (14.7)$$

Simple Average of Price Relatives Method

$$P_{01} = \frac{\sum \left(\frac{P_1}{P_0} \times 100 \right)}{n} \quad (14.8)$$

Example: Calculate simple aggregative and average of relatives index

$$\text{Simple Aggregative Index} = \frac{86}{70} \times 100 = 122.86$$

$$\text{Simple Average of Relatives} = \frac{490}{4} = 122.5$$

Table 14.2: Price Data for Simple Index Calculation

Commodity	Base Price (P_0)	Current Price (P_1)	Price Relative $\left(\frac{P_1}{P_0} \times 100\right)$
A	10	12	120
B	20	26	130
C	15	18	120
D	25	30	120
Total	70	86	490

14.5.2 Weighted Price Index Numbers

Weighted Aggregative Method

Laspeyres' Price Index Uses base period quantities as weights:

$$P_{01}^L = \frac{\sum P_1 Q_0}{\sum P_0 Q_0} \times 100 \quad (14.9)$$

Paasche's Price Index Uses current period quantities as weights:

$$P_{01}^P = \frac{\sum P_1 Q_1}{\sum P_0 Q_1} \times 100 \quad (14.10)$$

Fisher's Ideal Price Index Geometric mean of Laspeyres' and Paasche's indices:

$$P_{01}^F = \sqrt{P_{01}^L \times P_{01}^P} \quad (14.11)$$

Comprehensive Example:

Table 14.3: Data for Weighted Index Calculation

Item	P_0	P_1	Q_0	Q_1
Wheat	20	25	100	110
Rice	30	36	80	85
Sugar	40	48	50	45

Table 14.4: Calculations for Weighted Indices

Item	$P_0 Q_0$	$P_1 Q_0$	$P_0 Q_1$	$P_1 Q_1$
Wheat	2000	2500	2200	2750
Rice	2400	2880	2550	3060
Sugar	2000	2400	1800	2160
Total	6400	7780	6550	7970

Calculations:

$$P_{01}^L = \frac{7780}{6400} \times 100 = 121.56 \quad (14.12)$$

$$P_{01}^P = \frac{7970}{6550} \times 100 = 121.68 \quad (14.13)$$

$$P_{01}^F = \sqrt{121.56 \times 121.68} = 121.62 \quad (14.14)$$

14.6 Quantity Index Numbers

14.6.1 Laspeyres' Quantity Index

$$Q_{01}^L = \frac{\sum P_0 Q_1}{\sum P_0 Q_0} \times 100 \quad (14.15)$$

14.6.2 Paasche's Quantity Index

$$Q_{01}^P = \frac{\sum P_1 Q_1}{\sum P_1 Q_0} \times 100 \quad (14.16)$$

14.6.3 Fisher's Quantity Index

$$Q_{01}^F = \sqrt{Q_{01}^L \times Q_{01}^P} \quad (14.17)$$

Using the previous data:

$$Q_{01}^L = \frac{6550}{6400} \times 100 = 102.34 \quad (14.18)$$

$$Q_{01}^P = \frac{7970}{7780} \times 100 = 102.44 \quad (14.19)$$

$$Q_{01}^F = \sqrt{102.34 \times 102.44} = 102.39 \quad (14.20)$$

14.7 Properties and Tests of Index Numbers

14.7.1 Fisher's Tests

Time Reversal Test

An index should satisfy: $P_{01} \times P_{10} = 1$ (or 100 if expressed as percentage)

Factor Reversal Test

Price index \times Quantity index should equal Value index:

$$P_{01} \times Q_{01} = V_{01} \quad (14.21)$$

Where $V_{01} = \frac{\sum P_1 Q_1}{\sum P_0 Q_0} \times 100$

Circular Test

For three periods: $P_{01} \times P_{12} \times P_{20} = 1$

14.7.2 Test Results Summary

14.8 Chain Index Numbers

Chain indices compare each period with the immediately preceding period, then link them together.

Table 14.5: Test Results for Different Index Methods

Index Method	Time Reversal	Factor Reversal	Circular Test
Laspeyres'	Fails	Fails	Fails
Paasche's	Fails	Fails	Fails
Fisher's	Satisfies	Satisfies	Fails

$$\text{Chain Index}_{0n} = \frac{P_1}{P_0} \times \frac{P_2}{P_1} \times \frac{P_3}{P_2} \times \dots \times \frac{P_n}{P_{n-1}} \quad (14.22)$$

Example:

Table 14.6: Chain Index Calculation

Year	Price	Link Relative	Chain Index
2020	100	-	100.0
2021	105	105.0	105.0
2022	110	104.8	110.0
2023	115	104.5	115.0

14.9 Deflating and Real Values

Index numbers are used to convert nominal values to real values by removing the effect of price changes.

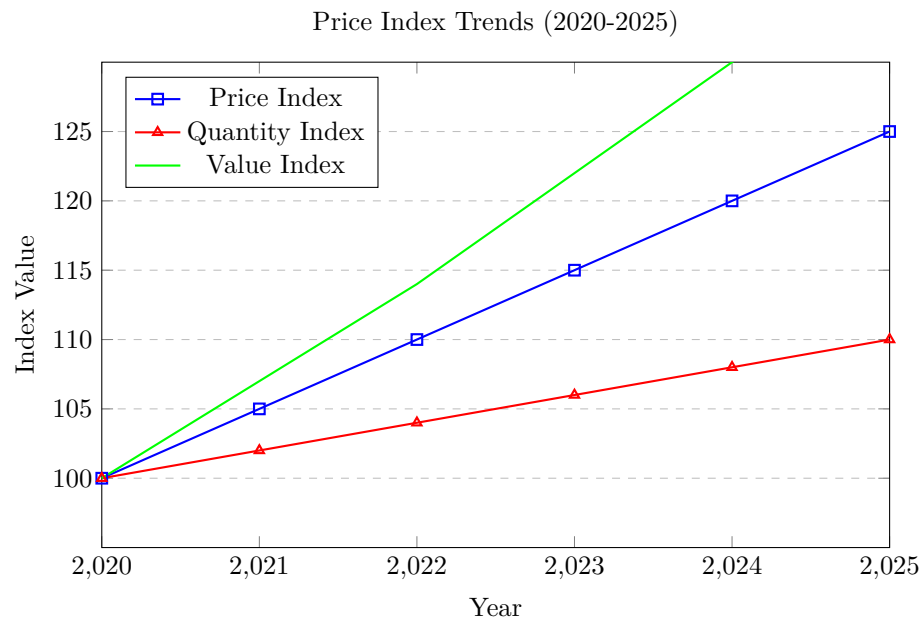
$$\text{Real Value} = \frac{\text{Nominal Value}}{\text{Price Index}} \times 100 \quad (14.23)$$

Example: If nominal wage increased from \$1000 to \$1200, but price index increased from 100 to 125:

$$\text{Real wage change} = \frac{1200}{125} \times 100 = 960$$

This shows that despite nominal increase, real purchasing power decreased.

14.10 Graphical Representation



14.11 Applications of Index Numbers

14.11.1 Economic Applications

- **Consumer Price Index (CPI):** Measures inflation
- **Producer Price Index (PPI):** Measures wholesale price changes
- **GDP Deflator:** Converts nominal GDP to real GDP
- **Stock Market Indices:** Track market performance

14.11.2 Business Applications

- Cost of living adjustments in salaries
- Escalation clauses in contracts
- Performance measurement
- Budgeting and forecasting

14.12 Limitations of Index Numbers

1. **Quality Changes:** Don't account for improvements in product quality
2. **New Products:** Difficulty in incorporating new goods and services
3. **Substitution Bias:** Consumers may substitute expensive items with cheaper alternatives
4. **Base Period Selection:** Choice of base period affects results
5. **Sampling Issues:** Representative sampling challenges
6. **Aggregation Problems:** Combining diverse items into single index

14.13 Worked Examples

14.13.1 Example 1: Consumer Price Index

A household's monthly expenditure pattern and price changes:

Table 14.7: Household Expenditure Data

Category	Weight (%)	2020 Price	2023 Price	Price Relative
Food	35	100	125	125
Housing	25	100	115	115
Transportation	20	100	140	140
Healthcare	10	100	110	110
Education	10	100	105	105

CPI calculation:

$$\text{CPI} = \frac{\sum(\text{Weight} \times \text{Price Relative})}{\sum \text{Weight}} \quad (14.24)$$

$$= \frac{35 \times 125 + 25 \times 115 + 20 \times 140 + 10 \times 110 + 10 \times 105}{100} \quad (14.25)$$

$$= \frac{4375 + 2875 + 2800 + 1100 + 1050}{100} \quad (14.26)$$

$$= \frac{12200}{100} = 122 \quad (14.27)$$

The CPI indicates a 22% increase in cost of living from 2020 to 2023.

14.13.2 Example 2: Fixed Base vs Chain Base

Table 14.8: Comparison of Fixed Base and Chain Base Methods

Year	Price	Fixed Base (2020=100)	Link Relative	Chain Index
2020	50	100.0	-	100.0
2021	55	110.0	110.0	110.0
2022	62	124.0	112.7	124.0
2023	68	136.0	109.7	136.0
2024	75	150.0	110.3	150.0

14.14 Summary of Key Formulas

14.14.1 Basic Relatives

$$\text{Price Relative} = \frac{P_1}{P_0} \times 100 \quad (14.28)$$

$$\text{Quantity Relative} = \frac{Q_1}{Q_0} \times 100 \quad (14.29)$$

$$\text{Value Relative} = \frac{P_1 Q_1}{P_0 Q_0} \times 100 \quad (14.30)$$

14.14.2 Simple Index Numbers

$$\text{Simple Aggregative} = \frac{\sum P_1}{\sum P_0} \times 100 \quad (14.31)$$

$$\text{Simple Average of Relatives} = \frac{\sum \left(\frac{P_1}{P_0} \times 100 \right)}{n} \quad (14.32)$$

14.14.3 Weighted Price Indices

$$\text{Laspeyres' Index} = \frac{\sum P_1 Q_0}{\sum P_0 Q_0} \times 100 \quad (14.33)$$

$$\text{Paasche's Index} = \frac{\sum P_1 Q_1}{\sum P_0 Q_1} \times 100 \quad (14.34)$$

$$\text{Fisher's Index} = \sqrt{P_{01}^L \times P_{01}^P} \quad (14.35)$$

14.14.4 Weighted Quantity Indices

$$\text{Laspeyres' Quantity Index} = \frac{\sum P_0 Q_1}{\sum P_0 Q_0} \times 100 \quad (14.36)$$

$$\text{Paasche's Quantity Index} = \frac{\sum P_1 Q_1}{\sum P_1 Q_0} \times 100 \quad (14.37)$$

$$\text{Fisher's Quantity Index} = \sqrt{Q_{01}^L \times Q_{01}^P} \quad (14.38)$$

14.14.5 Weighted Average of Relatives

$$\text{Price Index} = \frac{\sum W \times P_r}{\sum W} \quad (14.39)$$

$$\text{where } P_r = \frac{P_1}{P_0} \times 100 \quad (14.40)$$

14.14.6 Chain Index

$$\text{Chain Index}_{0n} = \prod_{i=1}^n \frac{P_i}{P_{i-1}} \times 100 \quad (14.41)$$

14.14.7 Tests of Adequacy

$$\text{Time Reversal Test: } P_{01} \times P_{10} = 1 \quad (14.42)$$

$$\text{Factor Reversal Test: } P_{01} \times Q_{01} = V_{01} \quad (14.43)$$

$$\text{Circular Test: } P_{01} \times P_{12} \times P_{20} = 1 \quad (14.44)$$

14.14.8 Deflation

$$\text{Real Value} = \frac{\text{Nominal Value}}{\text{Price Index}} \times 100 \quad (14.45)$$

Note: Fisher's Ideal Index satisfies both Time Reversal and Factor Reversal tests, making it theoretically superior, though computationally more complex than Laspeyres' or Paasche's indices.