

# Thinking like a Frequentist

Tin Vu

January 2026

*This page is intentionally left blank.*

# Preface

*"Are you watching closely?"*

Alfred Borden, *The Prestige* (2006)

When I was learning Probability and Statistics (ProbStat for short), I did not really understand why things were distributed *normally*. Were they *normal* because they were, or were *we just assuming "Everything is normal for simplicity"*? Sometimes, I felt we were overusing this term and in many situations, we were separating our theory from the actual data due to the initial assumption of *normal* distribution.

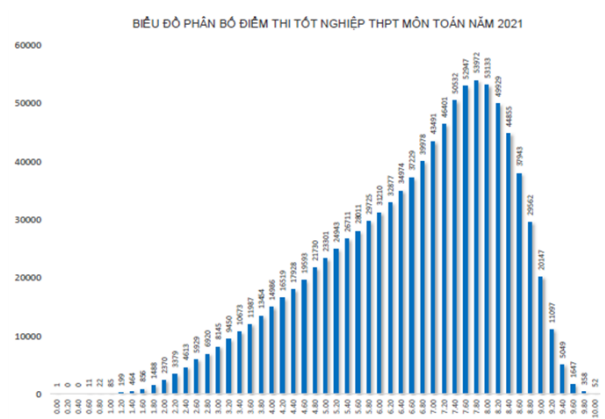


Figure 1: Distribution of Math scores in the 2021 National Entrance Exam

As you can see here, this graph **does not really resemble a bell-shaped curve**; and clearly, this data is not distributed normally. In fact, not much real-world statistical data fits a bell-shaped curve. But why are we still using it?

In my opinion, I think the main reason is due to **Central Limit Theorem** (CLT for short), every sample in the same population (regardless of how the data is distributed) has a **mean** that converges to a *normal distribution*, or **bell-shaped curve**. This topic will be explained in detail in **Chapter 7: Fundamentals of Statistics**.

This book approaches the fundamentals of probability and statistics in a very mathematically rigorous way to ensure maximum accuracy. I encourage readers to prove all of theorems, corollaries; and you should create your own examples for each theorem to gain a deeper understanding of their origins.

# Acknowledgements

First of all, I want to thank you for reading this book. Although I am not a native English speaker, but I truly enjoy writing in English. I have tried my best to express my ideas in English, but minor grammatical or spelling errors are unavoidable. If you find any of them, I would be very happy to receive your feedback via my email address: [tinvu1309@gmail.com](mailto:tinvu1309@gmail.com)

Secondly, I am extremely grateful and would like to express my thanks to the authors of the textbook "Probability and Statistics for Engineers and Scientists", 9th edition. This book truly saved my student life.

Finally, the idea of writing this book was inspired by Prof. Steve Brunton lectures on YouTube. You should check out his videos too!

—Tin Vu—

# Contents

|          |  |           |
|----------|--|-----------|
| <b>0</b> | <b>Introduction to Probability and Statistics</b>                        | <b>9</b>  |
| 0.1      | Coin Tossing . . . . .   | 9         |
| 0.2      | Weather Forecasting . . . . .  | 10        |
| 0.3      | Relationship between Probability and Statistics . . . . .                | 10        |
| <b>1</b> | <b>Fundamentals of Probability</b>                                       | <b>12</b> |
| 1.1      | Sample Space and Events . . . . .  | 12        |
| 1.1.1    | Sample Space . . . . .   | 12        |
| 1.1.2    | Events . . . . .   | 12        |
| 1.2      | Counting Sample Points . . . . .   | 13        |
| 1.2.1    | Rule of Product . . . . .  | 14        |
| 1.2.2    | Permutations . . . . .   | 14        |
| 1.2.3    | Combinations . . . . .   | 15        |
| 1.3      | Probability of an Event . . . . .  | 16        |
| 1.4      | Conditional Probability . . . . .  | 19        |
| 1.4.1    | Conditional Probability . . . . .  | 19        |
| 1.4.2    | Independent Events . . . . .   | 20        |
| 1.5      | Total Probability and Bayes' rule . . . . .                              | 21        |
| 1.5.1    | Total Probability . . . . .  | 21        |
| 1.5.2    | Bayes' rule . . . . .  | 21        |
| <b>2</b> | <b>Random Variables and Probability Distributions</b>                    | <b>23</b> |
| 2.1      | Definition of Random Variables . . . . .                                 | 23        |
| 2.2      | Discrete Probability Distributions . . . . .                             | 24        |
| 2.3      | Continuous Probability Distributions . . . . .                           | 25        |
| 2.4      | Joint Probability Distributions . . . . .                                | 26        |
| 2.4.1    | Case of Discrete Random Variables . . . . .                              | 26        |
| 2.4.2    | Case of Continuous Random Variables . . . . .                            | 29        |
| <b>3</b> | <b>Mathematical Expectation</b>  | <b>31</b> |
| 3.1      | Mean of a Random Variable . . . . .                                      | 31        |
| 3.2      | Variance and Covariance of Random Variables . . . . .                    | 35        |
| 3.2.1    | Variance of Random Variables . . . . .                                   | 35        |
| 3.2.2    | Covariance of Random Variables . . . . .                                 | 36        |
| 3.3      | Means and Variances of Linear Combinations of Random Variables . . . . . | 38        |
| 3.4      | Markov's and Chebyshev's Inequalities . . . . .                          | 40        |
| 3.4.1    | Markov's Inequality . . . . .  | 40        |
| 3.4.2    | Chebyshev's Inequality . . . . .   | 40        |

|          |  |           |
|----------|--|-----------|
| <b>4</b> | <b>Some Discrete Probability Distributions</b>                       | <b>42</b> |
| 4.1      | Bernoulli, Binomial and Poisson Distributions . . . . .              | 42        |
| 4.1.1    | Bernoulli Distribution . . . . .                                     | 42        |
| 4.1.2    | Binomial Distribution . . . . .                                      | 44        |
| 4.1.3    | Poisson Distribution . . . . .                                       | 46        |
| 4.2      | Negative Binomial and Geometric Distributions . . . . .              | 51        |
| 4.2.1    | Negative Binomial Distribution . . . . .                             | 51        |
| 4.2.2    | Geometric Distribution . . . . .                                     | 51        |
| 4.3      | Hypergeometric Distribution . . . . .                                | 52        |
| <b>5</b> | <b>Some Continuous Probability Distributions</b>                     | <b>53</b> |
| 5.1      | Uniform Distribution . . . . .                                       | 53        |
| 5.2      | Normal Distribution . . . . .  | 54        |
| 5.2.1    | The Idea behind the Normal Distribution . . . . .                    | 54        |
| 5.2.2    | Standard Normal Distribution . . . . .                               | 58        |
| 5.2.3    | Normal Distribution . . . . .  | 59        |
| 5.3      | Exponential, Gamma and Chi-Squared Distributions . . . . .           | 61        |
| 5.3.1    | Exponential Distribution . . . . .                                   | 61        |
| 5.3.2    | Gamma Distribution . . . . .   | 63        |
| 5.3.3    | Chi-Squared Distribution . . . . .                                   | 66        |
| <b>6</b> | <b>Functions of Random Variables</b>                                 | <b>67</b> |
| 6.1      | Transformations of Variables . . . . .                               | 67        |
| 6.1.1    | Linear Transformations . . . . .                                     | 67        |
| 6.1.2    | Non-linear Transformations . . . . .                                 | 71        |
| 6.2      | Moment-Generating Functions . . . . .                                | 73        |
| 6.2.1    | Definition of Moment-Generating Functions . . . . .                  | 73        |
| 6.2.2    | Some Useful Moment-Generating Functions . . . . .                    | 74        |
| 6.2.3    | Linear Combinations of Random Variables . . . . .                    | 75        |
| <b>7</b> | <b>Fundamentals of Statistics</b>                                    | <b>77</b> |
| 7.1      | The Big Picture of Statistics . . . . .                              | 78        |
| 7.1.1    | Populations and Samples . . . . .                                    | 78        |
| 7.1.2    | Sample Mean and Sample Variance . . . . .                            | 79        |
| 7.2      | Sampling Distribution of Means . . . . .                             | 80        |
| 7.2.1    | Central Limit Theorem . . . . .                                      | 80        |
| 7.2.2    | t-Distribution . . . . .   | 82        |
| 7.3      | Sampling Distribution of Variances . . . . .                         | 85        |
| 7.4      | Case Study: Seed Germination Time . . . . .                          | 86        |
| 7.4.1    | Estimating the Mean using t-test . . . . .                           | 87        |
| 7.4.2    | Estimating the Variance using Chi-test . . . . .                     | 88        |
| <b>8</b> | <b>Classical Methods of Estimation</b>                               | <b>90</b> |
| 8.1      | Definition of Unbiased Estimator . . . . .                           | 90        |
| 8.2      | Determining Statistical Intervals Like A Pro: Step By Step . . . . . | 92        |
| 8.2.1    | Selecting Your Test Statistic . . . . .                              | 92        |
| 8.2.2    | Establishing The Confidence Interval . . . . .                       | 93        |

# List of Figures

|     |   |    |
|-----|---|----|
| 1   | Distribution of Math scores in the 2021 National Entrance Exam . . . . .                      | 2  |
| 1   | General model of simple ProbStat problems . . . . .   | 11 |
| 1.1 | Tree diagram for rule of product . . . . .  | 14 |
| 1.2 | Visualizing definition of probability . . . . .   | 17 |
| 1.3 | Visualizing how conditional probability is calculated . . . . .                               | 20 |
| 1.4 | Total probability . . . . .   | 21 |
| 2.1 | What is the probability that a number will fall within this range? . . . . .                  | 26 |
| 3.1 | Pdf of fair coin tossing experiment with its mean . . . . .                                   | 31 |
| 3.2 | Pdf of unfair coin tossing experiment with its mean . . . . .                                 | 32 |
| 3.3 | The mean value of "luck level" . . . . .  | 33 |
| 3.4 | How is the data distributed around the <b>mean</b> ? . . . . .                                | 35 |
| 3.5 | Example of a valid pdf graph . . . . .  | 41 |
| 4.1 | The pdf of Bernoulli distribution . . . . .   | 43 |
| 4.2 | The pdf of Binomial distribution . . . . .  | 46 |
| 4.3 | $\mathcal{B}(x; 100, 0.6)$ graph . . . . .  | 46 |
| 4.4 | $\mathcal{B}(x; 100, 0.05)$ graph . . . . .   | 47 |
| 4.5 | The pdf of Poisson distribution . . . . .   | 50 |
| 4.6 | A $4 \times 4$ grid . . . . .   | 50 |
| 4.7 | Illustration of Hypergeometric distribution . . . . .   | 52 |
| 5.1 | The pdf of Uniform distribution . . . . .   | 54 |
| 5.2 | A circular courtyard with raindrops falling inside . . . . .                                  | 55 |
| 5.3 | The probability of raindrops falling around the center is higher than anywhere else . . . . . | 57 |
| 5.4 | The pdf of Standard Normal distribution . . . . .   | 58 |
| 5.5 | Visualizing how the general normal distribution formula is formed . . . . .                   | 59 |
| 5.6 | Normal approximation to the binomial . . . . .  | 60 |
| 5.7 | Illustration of Exponential distribution . . . . .  | 62 |
| 5.8 | The pdf of Exponential distribution and its memoryless property . . . . .                     | 63 |
| 6.1 | Changing the domain region . . . . .  | 71 |
| 7.1 | Why was Mendel so certain about 3 : 1 ratio? . . . . .  | 77 |
| 7.2 | Sampling process . . . . .  | 78 |
| 7.3 | Inferring from $(\bar{x}, s^2) \rightarrow (\mu, \sigma^2)$ . . . . .                         | 80 |
| 7.4 | t-distribution and Standard Normal distribution . . . . .                                     | 85 |
| 7.5 | Experiment's measured data . . . . .  | 87 |
| 7.6 | t-distribution with 95% confidence interval . . . . .   | 87 |
| 7.7 | Chi-Squared distribution with 95% confidence interval . . . . .                               | 89 |

|     |                                       |    |
|-----|---------------------------------------|----|
| 8.1 | Two-sided confidence bounds . . . . . | 93 |
|-----|---------------------------------------|----|



# List of notations

Since much of my work is handwritten, so I have modified some commonly used notations for probability distribution functions using curved script for convenience. You should notice that my conventions are not the international standards.

1.  $\mathcal{I}$ : Bernoulli distribution
2.  $\mathcal{B}$ : Binomial distribution
3.  $\mathcal{B}^*$ : Negative binomial distribution
4.  $\mathcal{H}$ : Hypergeometric distribution
5.  $\mathcal{G}^*$ : Geometric distribution
6.  $\mathcal{P}$ : Poisson distribution
7.  $\mathcal{U}$ : Uniform distribution
8.  $\mathcal{N}$ : Normal distribution
9.  $\mathcal{G}$ : Gamma distribution
10.  $\mathcal{E}$ : Exponential distribution
11.  $\mathcal{C}$ : Chi-squared distributon
12.  $\mathcal{T}$ : t-distribution (or Student distribution)

# Chapter 0

## Introduction to Probability and Statistics

### 0.1 Coin Tossing

Imagine you have a coin, like this one. It is an ordinary coin that can be found everywhere. Everyone knows that the chance of heads (H) appearing after each toss is 50%, no one even doubts that. So, the **probability** of an **event** that heads appearing is 50%. But because you are a very curious person, you do not easily accept the truth like others, so you will find a way to test it. The fact "The probability of heads appearing is 50%" will be tested, and it can be called a **hypothesis**.

The simplest way to test your hypothesis is tossing a coin many times. You might toss a coin  $N = 100$  times and see heads appear  $n = 40$  times, so your actual probability is:

$$\hat{p} = \frac{n}{N} = 0.4$$

But your previous assumption (or hypothesis) states that the ideal probability is:

$$p = 0.5$$

Is there anything wrong here? Dissatisfied, you continue your experiment. This time, you toss a coin  $N = 500$  times and see heads appear  $n = 270$  times, so the new actual probability is:

$$\hat{p} = \frac{n}{N} = 0.54$$

This time the result is closer with initial hypothesis  $p$ , but your hand now must be very tired after tossing a coin 600 times. So do you think you have **enough** evidence to conclude that if  $N \rightarrow +\infty$ , then  $\hat{p} \rightarrow p$ ? If so, your assumption is correct because *can not be refuted*; and if not, our common sense might be wrong because the evidence *strongly refutes it*.

Another strategy to toss a coin is fixing and dividing  $N$ . Instead of tossing  $N = 600$  times and only getting the final  $p$  value, you can divide large number of coin tosses  $N = 600$  to smaller tosses  $N_1 = N_2 = \dots = N_6 = 100$  (but do not too small, at least each  $N_i > 30$ ), and get 6 values of  $\hat{p}_i$ . Suppose that you would get:

| $\hat{p}_1$ | $\hat{p}_2$ | $\hat{p}_3$ | $\hat{p}_4$ | $\hat{p}_5$ | $\hat{p}_6$ |
|-------------|-------------|-------------|-------------|-------------|-------------|
| 0.43        | 0.46        | 0.56        | 0.53        | 0.49        | 0.51        |

You might see that the value  $p = 0.5$  too idealistic to occur in our experiment, and  $\hat{p}$  constantly changes. Therefore instead of determining the **exact value** of  $p$ , we **predict** that the  $p$  value lies inside a **closed interval** with a **certainty** of  $(1 - \alpha)100\%$ . For example, we can confidently conclude that the value of  $p$  lies within the range  $(0.48, 0.51)$  with 95% accuracy.

## 0.2 Weather Forecasting

You do not need to know anything about geography to predict what the weather will be tomorrow. Everything you need is just knowledge about Poisson, exponential distributions and the **average number** of rainy days per month where you are living.

For example, in December there are **average** 2 rainy days. What are the probabilities of:

1. There will be 3 rainy days this month.
2. Tomorrow will be a rainy day, if today is the 10th and you have not seen any rain since the beginning of the month.

This problem will be covered in detail in **Chapter 4: Some Discrete Probability Distributions** and **Chapter 5: Continuous Probability Distributions**, but if you do have experience with random variables, you can try solving it!

The answer for the first question is:

$$P(X = 3) = \frac{e^{-2}2^3}{3!} = 0.1804 = 18.04\%$$

The answer for the second question is:

$$P(X < 11|X > 10) = 1 - P(X > 11|X > 10) = 1 - e^{\frac{-2}{31}} = 0.0624 = 6.24\%$$

Since 6.24% is quite low, so tomorrow you do not have to bring an umbrella.

## 0.3 Relationship between Probability and Statistics

The formal definitions of **probability** and **statistics** will be discussed later in **Chapter 1: Fundamentals of Probability** and **Chapter 7: Fundamentals of Statistics**, but now let's focus on the general model below:

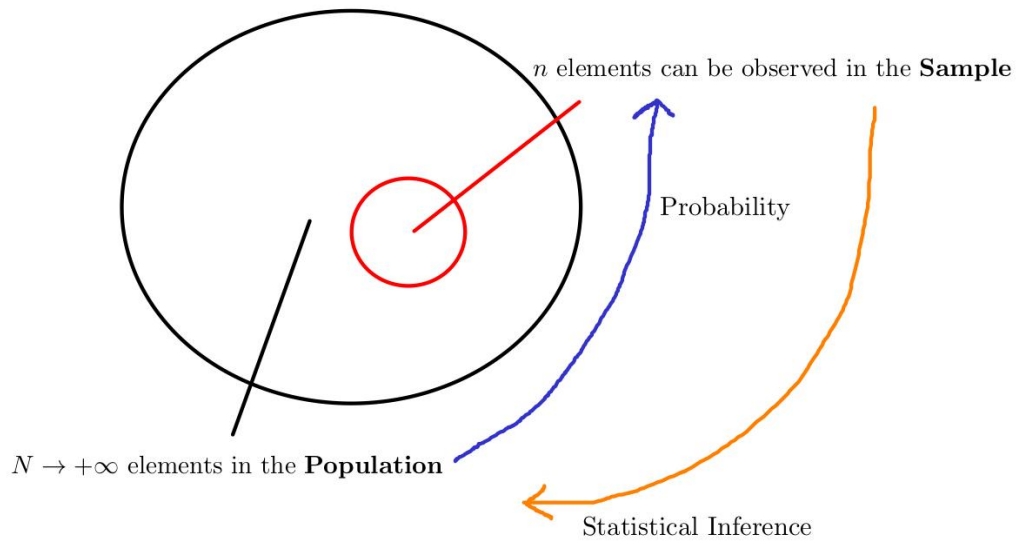


Figure 1: General model of simple ProbStat problems

You want to know some attributes of a very large **population** such as the probabilities, means,  $\dots$  of the quantities. But there is no way you can observe the entire population (in many cases,  $N \rightarrow +\infty$ ), so you have to take some **samples** ( $n$  elements) and use **probability rules** to process them. After that, you can try applying **statistical inference rules** to draw conclusions about the original population. That is how it works!

Referring to the previous example, you want to check if the probability of heads appearing is  $p = 50\%$ , so you toss a coin multiple times. But since you can only toss a coin for 600 times consecutively, so you conclude based on your final result  $\hat{p} = 0.54$  (or closed interval with a certainty) that might be if  $N \rightarrow +\infty$  then  $\hat{p} \rightarrow p = 0.5$ . You took the **samples** ( $N = 600$  or  $N_i = 100$ ), used **probability rule** to calculate the actual  $p$  value, and then applied **statistical inference rules** (hypothesis or closed interval) to draw conclusions respectively.

# Chapter 1

## Fundamentals of Probability

### 1.1 Sample Space and Events

#### 1.1.1 Sample Space

If you toss a coin, you will see that there are 2 possible outcomes: heads and tails, denoted by capital letters H and T; or if you roll a die, you will see that there are 6 possible outcomes, from 1 to 6. Tossing a coin, or rolling a die are typical examples of **experiments**.

**Definition 1.1.1.** An ***experiment*** is the process that generates a set of outcomes (or data).

All of possible outcomes are collected into a single set.

$$S_1 = \{H, T\}$$

$$S_2 = \{1, 2, 3, 4, 5, 6\}$$

**Definition 1.1.2.** The ***sample space*** is the set of all possible outcomes of an ***experiment***.

The sample space is usually represented by the symbol  $S$  or  $\Omega$ . You can easily see that  $S$  is not always a countable or finite set. For instance,  $S_3$  is the set of all random numbers you can choose within the range  $(0, 1)$ :

$$S_3 = \{x \mid 0 < x < 1\}$$

$S_4$  is the set of all number of coin tosses until first heads appear:

$$S_4 = \{1, 2, 3, \dots\}$$

It is possible that you toss a coin forever, and heads never appear.

**Definition 1.1.3.** A ***sample point*** is a single outcome of the sample space  $S$ .

$S_1, S_2$  have 2 and 6 sample points respectively, while  $S_3, S_4$  have infinite sample points.

#### 1.1.2 Events

For any given experiment, we are often interested in the occurrence of certain **events** rather than a specific element (or **sample point**) in the sample space. For example, in the die roll experiment, you may want to know when the outcome is an even number. This will happen if the result is an element of the subset  $E_2$  of the sample space  $S_2$ :

$$E_2 = \{2, 4, 6\}$$

**Definition 1.1.4.** An **event** is the subset of a sample space.

Events are always denoted by capital letters like  $A, B, C, \dots$ . Similarly, an event is not always a countable or finite subset.

The **complement** of an event  $E_2$  with respect to  $S_2$  is the set of all *odd outcomes*, and can be represented as follows:

$$\overline{E}_2 = \{1, 3, 5\}$$

There are many ways to denote the **complement** of an event:  $\overline{E}, E', E^c$ , but in this book I choose overline notation for convenience and make it easy to relate with Boolean algebra.

**Definition 1.1.5.** The **complement** of an event  $E$  with respect to  $S$  is the subset of all elements of  $S$  that are not in  $E$ , and can be denoted by the symbol  $\overline{E}$ .

Applying set theory, we can perform many set operations like joint, disjoint, union,  $\dots$

**Definition 1.1.6.** The **intersection** of two events  $A$  and  $B$ , denoted by the symbol  $A \cap B$  is the event containing all elements that are common to  $A$  and  $B$ .

For example, if  $A$  and  $B$  are the subset of  $S_3$  and defined as:

$$\begin{aligned} A &= \{x \mid 0 < x < 0.7\} \\ B &= \{x \mid 0.2 < x < 0.9\} \\ \Rightarrow A \cap B &= \{x \mid 0.2 < x < 0.7\} \end{aligned}$$

**Definition 1.1.7.** Two events  $A$  and  $B$  are **mutually exclusive**, or **disjoint**, if  $A \cap B = \emptyset$ , that is  $A$  and  $B$  have nothing in common.

For example,  $E_2$  and  $\overline{E}_2$  are mutually exclusive.

**Definition 1.1.8.** The **union** of the two events  $A$  and  $B$ , denoted by the symbol  $A \cup B$  is the event containing all of the elements that belong to  $A$  or  $B$  or both.

For example,  $A \cup B = \{x \mid 0 < x < 0.9\}$ , and  $(A \cup B) \subset S_3$ . Note that  $E_2 \cup \overline{E}_2 = S_2$  is an useful result and can be generalized to corollary below:

**Corollary 1.1.0.1.** If  $A$  is an event with respect to sample space  $S$ , then  $A \cup \overline{A} = S$

De Morgan's laws can also be applied to set theory.

**Corollary 1.1.0.2.** If  $A$  and  $B$  are events with respect to sample space  $S$ , then

$$\begin{aligned} \overline{A \cap B} &= \overline{A} \cup \overline{B} \\ \overline{A \cup B} &= \overline{A} \cap \overline{B} \end{aligned}$$

These results above are really useful in many cases, but you do not have prove them since they are pretty simple. Proving them rigorously is not the focus of Probability and Statistics book.

## 1.2 Counting Sample Points

In this section, we develop some *counting techniques* to count the number of points in the sample space  $S$  and its event subset  $E$  without actually listing each elements. These techniques play an important role in solving some simple probability problems. Notice that these techniques can only be applied when your sample space is finite and countable.

Obviously, you can not count the number of elements of  $S_3$  and  $S_4$ , since they are infinite and uncountable sets.

### 1.2.1 Rule of Product

You toss a pair of coins. How many sample points are there in the sample space? First you can try to list all of the possible outcomes:

$$S = \{HH, HT, TH, TT\}$$

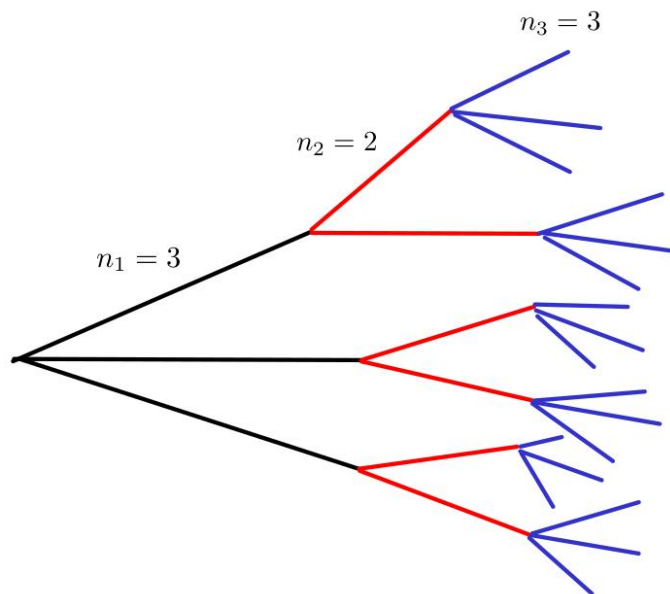
There are total 4 sample points in  $S$ . But listing all of the elements might be not so clever idea, so you use **rule of product**.

The first coin can land heads or tails, so can the second coin. We multiply the number of possible outcomes of the two coins:

$$n_1 n_2 = 2 \cdot 2 = 4 \text{ (possible outcomes)}$$

**Definition 1.2.1.** *If an operation can be performed in  $n_1$  ways, and if for each of these a second operation can be performed in  $n_2$  ways, and for each of the first two a third operation can be performed in  $n_3$  ways, and so fourth, the the sequence of  $k$  operations can be performed in:*

$$\prod_{i=1}^k n_i \text{ (ways)}$$



There are a total  $n_1 n_2 n_3 = 18$  ways to perform this operation

Figure 1.1: Tree diagram for rule of product

### 1.2.2 Permutations

**Definition 1.2.2.** A **permutation** is an arrangement of **all** or **part** of a set of objects.

For instance, the number of permutations of  $n$  distinct objects can be counted by using rule of product:

$$n(n-1)(n-2) \cdots 3 \cdot 2 \cdot 1$$

We introduce a new notation for such a number.

**Definition 1.2.3.** For any non-negative integer  $n$ ,  $n!$ , called " $n$  factorial", is defined as:

$$n! = n(n-1)(n-2) \cdots 3 \cdot 2 \cdot 1$$

with special case  $0! = 1$

**Theorem 1.2.1.** The number of **permutations** of  $n$  distinct objects is  $n!$

In general, the number of permutations of  $n$  distinct objects taken  $r$  at a time can also be counted by using rule of product:

$$n(n-1)(n-2) \cdots (n-r+2)(n-r+1)$$

This product can be represented by the new symbol  $nPr$ .

**Theorem 1.2.2.** The number of **permutations** of  $n$  distinct objects taken  $r$  at a time is:

$$nPr = \frac{n!}{(n-r)!}$$

But how about our  $n$  objects are not distinct? Assume that  $n_1$  are of one kind,  $n_2$  are of a second kind,  $\cdots$  and  $n_k$  of a  $k$ th kind.

**Theorem 1.2.3.** The number of **permutations** of  $n$  things of which  $n_1$  are of one kind,  $n_2$  of a second kind, and so forth is:

$$\frac{n!}{n_1!n_2! \cdots n_k!}$$

with  $\sum_{i=1}^k n_i = n$ .

For example, from the digits 1 to 9, we can form  $9!$  numbers made up of 9 distinct digits, or we can form  $9P5$  five-digit numbers such that all the digits are different. If we allow repetition, five-digits numbers are now formed by 2 digits 1, 2 digits 2 and 1 digit 3, then the number of permutation that satisfy is:

$$\frac{5!}{2!2!1!} = 30$$

### 1.2.3 Combinations

Consider another problem, now you have a set of  $n$  elements. How many ways you can partition it into  $k$  cells with  $n_1$  elements in the first cell,  $n_2$  elements in the second cell, and so forth? Coincidentally, the equation in **Theorem 1.2.3** appears again.

**Theorem 1.2.4.** The number of ways of partitioning a set of  $n$  objects into  $k$  cells with  $n_1$  elements in the first cell,  $n_2$  elements in the second, and so forth, is:

$$\binom{n}{n_1, n_2, \dots, n_k} = \frac{n!}{n_1!n_2! \cdots n_k!}$$

with  $\sum_{i=1}^k n_i = n$ .

In many problems, we are interested in the number of ways of selecting  $k$  objects from  $n$  without regard to order. These selections are called **combinations**. It is not too hard to realize that a **combination** is just a partition with 2 cells, one cell containing the  $k$  objects and the other containing the  $(n-k)$  objects.

$$\binom{n}{k, n-k} \text{ is often shortened to } \binom{n}{k}$$



**Definition 1.2.4.** A **combination** is a selection of items from a set that has distinct elements, such that **the order of selection does not matter**.

**Theorem 1.2.5.** The number of combinations of  $n$  distinct objects taken  $k$  at a time is:

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

For example, the number of subsets of 3 elements that can be obtained from an original set of 10 elements is:

$$\binom{10}{3} = 120$$

For the rest of this book, we mainly focus on **combinations**, and derive many probability distribution functions based on them. Like set theory, we do not go deeply into **counting techniques** since they are not the main concern of ProbStat.

## 1.3 Probability of an Event

Everyone knows the basic idea of probability. If I toss a **fair** coin once, and I want to know the probability of getting heads; how can I determine it? Strictly, I have to define  $S$  is the sample space of this experiment and  $A$  is the event "Getting heads after one toss".

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

$$A = \{H\}$$

So, the probability of event  $A$  occurring is:

$$P(A) = \frac{\text{number of sample points inside } A}{\text{number of sample points inside } S} = \frac{1}{2}$$

Very easy and intuitive. But how about tossing 4 **fair** coins once, and determining the probability of getting 2 heads? Now sample space  $S$  and its event set can be represented as:

$$S = \{HHHH, HHHT, HHTH, \dots\}$$

$$A = \{TTHH, THTH, \dots\}$$

Now we change our strategy using **counting techniques**:

$$P(A) = \frac{\text{number of sample points inside } A}{\text{number of sample points inside } S} = \frac{\binom{4}{2}}{2^4} = \frac{3}{8} = 0.375$$

In fact, the chance of getting heads is slightly greater than tails (you know, because coins are asymmetrical). Assume that the probability of heads appearing is 60%, and tails appearing is 40%. Since the role of **sample points** inside set  $S$  and subset  $A$  are not **equal** anymore, so we can not use **counting techniques** blindly.

$$P(A) = \binom{4}{2} 0.6^2 0.4^2 = 0.3456$$

Coin tossing is a typical example of **Bernoulli trial**, and the experiment tossing unfair coins is a **Bernoulli process**. The probability  $P(A)$  can be calculated by using **Binomial distribution** formula. You do not have to worry about these terms, we will cover them in **Chapter 4: Some Discrete Probability Distributions** very carefully.

Consider another problem, what is the probability of getting 0.7 when randomly choosing a number (assume that the role of every numbers are equal) inside the interval  $(0, 1)$ ?

$$S = \{x \mid 0 < x < 1\}$$

$$A = \{0.7\}$$

$$P(A) = \frac{\text{number of sample points inside } A}{\text{number of sample points inside } S} = \frac{1}{+\infty} = 0 \text{ (?????)}$$

Since  $P(A) = 0$ , so may we conclude 0.7 will never be chosen? Absolutely incorrect, even in common thinking. From previous examples, now we see the **limitation** of our "common definition" of probability in real life. Mathematically, our definition is just a very special case of the formal one.

**Definition 1.3.1.** The **probability** of an event  $A$  is the sum of the weights (or probabilities) of all sample points in  $A$ . Therefore:

$$0 \leq P(A) \leq 1, \quad P(\emptyset) = 0, \quad P(S) = 1$$

If  $A$  and  $B$  are **mutually exclusive (or disjoint)**, then  $P(A \cup B) = P(A) + P(B)$

This definition can also be called Kolmogorov's axioms. An intuitive way to understand it is sketching a sample space with some events inside.

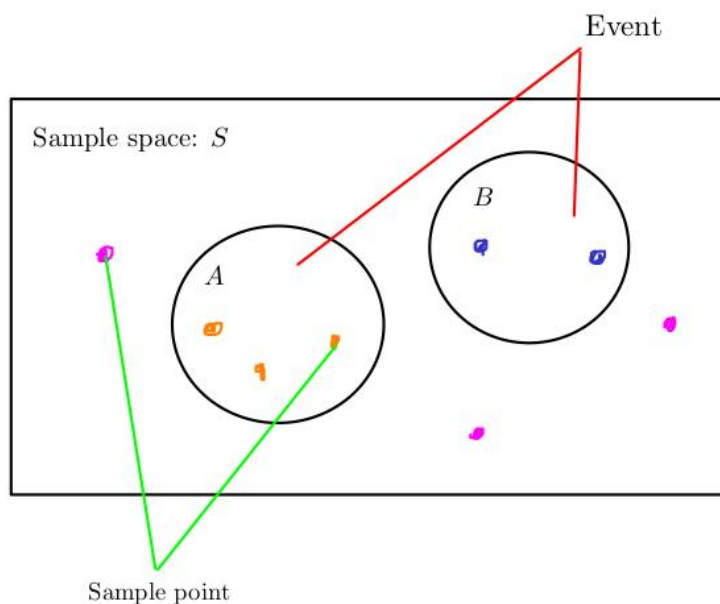


Figure 1.2: Visualizing definition of probability

Literally, probability is just a number that describes how likely an event can occur. I often relate it with "weight" quantity. As you can see in the sample space, I can assign arbitrarily the "weights" (or probabilities) of sample points as follows:

- Each blue point is 0.1
- Each orange point is 0.2

- Each pink point is 0.066

Using **Definition 1.3.1**, now we can obtain these results:

$$P(A) = 3 \cdot 0.2 = 0.6$$

$$P(B) = 2 \cdot 0.1 = 0.2$$

$$P(A \cup B) = P(A) + P(B) = 0.6 + 0.2 = 0.8$$

$$P(\overline{A \cup B}) = 3 \cdot 0.066 = 0.2$$

$$P(S) = 2 \cdot 0.1 + 3 \cdot 0.2 + 3 \cdot 0.066 = 1$$

You should verify that our "special definition" of probability above satisfies the axioms of probability, and can be formalized by a theorem.

**Theorem 1.3.1.** *If an experiment can result in any one of  $N$  different **equally** likely outcomes, and if exactly  $n$  of these outcomes correspond to event  $A$ , then the probability of event  $A$  is:*

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

Logically, you can view probability as a mapping from a set to a closed interval  $[0, 1]$ , then you can freely define the mapping  $P$  by yourself as long as it satisfies Kolmogorov's axioms.

$$P : \text{Set} \rightarrow [0, 1]$$

$$A \xrightarrow{P} P(A)$$

By applying set theory, we can derive several extremely useful theorems and corollaries:

**Theorem 1.3.2.** *If  $A$  and  $B$  are two events, then:*

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

**Theorem 1.3.3.** *If  $A$  is an event of sample space  $S$ , then:*

$$P(S) = P(A \cup \overline{A}) = P(A) + P(\overline{A}) = 1$$

**Corollary 1.3.3.1.** *If  $A$  is an event of sample space  $S$ , then:*

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

Again, De Morgan's laws can also be applied:

**Corollary 1.3.3.2.** *If  $A$  and  $B$  are events of sample space  $S$ , then:*

$$P(\overline{A \cup B}) = P(\overline{A} \cap \overline{B})$$

$$P(\overline{A \cap B}) = P(\overline{A} \cup \overline{B})$$

## 1.4 Conditional Probability

### 1.4.1 Conditional Probability

Imagine you now have a perfectly fair coin. Obviously if you define events  $A$  and  $B$  are heads and tails appearing after one toss, respectively, you will conclude:

$$P(A) = P(B) = 0.5$$

But I might ask you "If event  $A$  **did** occur, might event  $B$  would have any chance to occur?" . Since there is no way heads and tails can simultaneously appear, so your answer must be "No!". Now we can form an equation to represent the probability of a "special event" (or formally, conditional event) that "If  $A$  occurred, then  $B$  would occur.":

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

Now we shift our attention to a more complex problem. I give you a die, and you roll it. Events  $E_1$  and  $E_2$  can be defined as: "Landing a number that greater than 3" and "Landing a number that divisible by 2":

$$P(E_1) = P(E_2) = \frac{1}{3}$$

Sample points of 2 events can be listed:

$$E_1 = \{4, 5, 6\}$$

$$E_2 = \{2, 4, 6\}$$

If a number greater than 3 landed, what is the probability that it could be divisible by 2? You can count the number of sample points inside  $E_1$  and draw a result:

$$P(E_2|E_1) = \frac{2}{3}$$

Conversely, if a number divisible by 2 landed, what is the probability that it could be greater than 3?

$$P(E_1|E_2) = \frac{2}{3}$$

As you can see here, if one event occurs before another event, probability will be completely **changed**. So we call the pairs of events  $A$  and  $B$ ,  $E_1$  and  $E_2$  **dependent events** because they depend on each other. The "special" probabilities  $P(B|A)$ ,  $P(E_2|E_1)$ ,  $P(E_1|E_2)$  are called **conditional probability**.

**Definition 1.4.1.** The **conditional probability** of  $B$ , given  $A$ , denoted by  $P(B|A)$ , is defined:

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

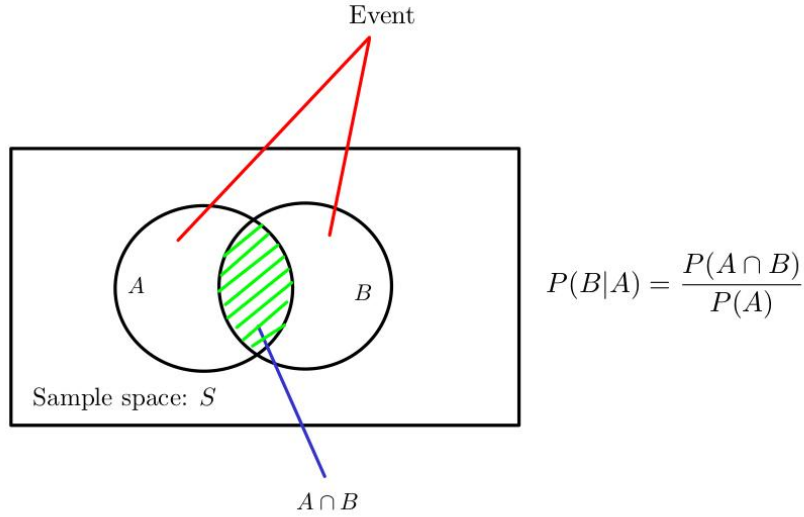


Figure 1.3: Visualizing how conditional probability is calculated

Now back to our previous problems, these conditional probabilities can easily be determined using **Definition 1.4.1** above:

$$P(E_1|E_2) = \frac{P(E_1 \cap E_2)}{P(E_2)} = \frac{2/6}{1/2} = \frac{2}{3}$$

$$P(E_2|E_1) = \frac{P(E_1 \cap E_2)}{P(E_1)} = \frac{2/6}{1/2} = \frac{2}{3}$$

### 1.4.2 Independent Events

Intuitively, we can see that if events  $A$  and  $B$  do not influence each other, then:

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

$$P(B|A) = P(B) \Leftrightarrow \frac{P(B \cap A)}{P(A)} = P(B) \Leftrightarrow P(A \cap B) = P(B)(A) \text{ (if } P(A) > 0)$$

**Definition 1.4.2.** Two events  $A$  and  $B$  are independent if and only if:

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

assuming  $P(A) > 0$

**Theorem 1.4.1.** Two events  $A$  and  $B$  are independent if and only if

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

There are many examples of independent events, like if we define 2 events  $E_1$ : "Getting heads on the first toss." and  $E_2$ : "Getting heads on the second toss.". Intuitively you can see that  $E_1$  and  $E_2$  are unrelated, so they are **independent events**. You can also verify this fact:

$$P(E_1 \cap E_2) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} = P(E_1)P(E_2)$$

It is very important to note that determining the independence of events is completely **unrelated** to their mutual exclusion. Events are considered independent if and only if they satisfy **Theorem 1.4.1**.

## 1.5 Total Probability and Bayes' rule

### 1.5.1 Total Probability

In many situations, we do not know directly the information of  $P(A)$  in **Definition 1.4.1** (or denominator part); so in this section, we will develop a simple formula to handle this problem.

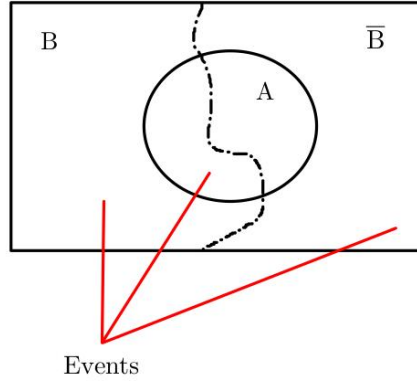


Figure 1.4: Total probability

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

Because  $A \cap B$  and  $A \cap \bar{B}$  are independent events, so:

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

**Theorem 1.5.1.** *If  $A$  and  $B$  are two events of sample space  $S$ , then:*

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

Total probability is an useful formula, especially when you are conducting surveys in practice. For example, in my university, there are 2 types of students: those who "studied the whole semester" and those who "studied only one night before the exam". After the final exam, I asked everyone in my class about their scores and totaled the results. Event  $A$ : "Got  $A+$ " and event  $\bar{A}$ : "Did not get  $A+$ "; event  $B$ : "Studied the whole semester" and event  $\bar{B}$ : "Studied only one night".

$$P(A|B) = 0.8, P(A|\bar{B}) = 0.3, P(B) = 0.4$$

Using total probability formula, I obtained the probability of  $A$ :

$$P(A) = P(A|B)P(B) + P(A|\bar{B})P(\bar{B}) = 0.8 \cdot 0.4 + 0.3 \cdot (1 - 0.4) = 0.5$$

Wow, that was an impressive ratio. Half of a class received perfect score. Was this course too easy?

### 1.5.2 Bayes' rule

If I studied hard, I would get  $A+$ . But how about me, who was not keen on studying boring courses like Computer Architecture but *still survived after final test and even got "A+"*? I could questioned myself "Was I too lucky?". To answer myself, I had to calculate  $P(\bar{B}|A)$ , if the result is not so high, perhaps I was lucky. Bayes' rule was what I needed.

**Theorem 1.5.2.** *If  $A$  and  $B$  are two events of sample space  $S$ , then:*

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

Bayes' rule seems very simple. Indeed an average high school student has no difficulty finding it, but its idea is brilliant. Before Bayes, we typically only reasoned about problems in terms of cause first, then effect. But Bayes' rule opens up a completely new way of thinking for us: knowing the effect beforehand, then understanding how cause influences it.

Back to my previous problem, I applied Bayes' rule:

$$P(B|A) = \frac{P(A|B)P(B)}{P(A|B)P(B) + P(A|\bar{B})P(\bar{B})} = \frac{0.8 \cdot 0.4}{0.5} = 0.64 \Rightarrow P(\bar{B}|A) = 1 - P(B|A) = 0.36$$

Since 0.36 was not a high number, indeed I was incredibly lucky and the course was not as easy as I thought.

# Chapter 2

## Random Variables and Probability Distributions

### 2.1 Definition of Random Variables

Back to our coin tossing game, now you toss a fair coin three times. This experiment has sample space  $S$ :

$$S = \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\}$$

Since writing all of the sample points is quite lengthy and time-consuming, sometimes unnecessarily, so how about we assign each point with a **numerical value**? Because the assignment of values is entirely based on our own conventions, there are no constraints whatsoever. But for this reason, we should choose the **smartest** and **most convenient** way to assign values to suit our concern. For example, if I define an event  $A$ : "Getting 2 heads", then the cleverest way is assigning each sample point with its number of heads.

$$\begin{aligned} HHH &\rightarrow 3 \\ HHT, HTH, THH &\rightarrow 2 \\ HTT, TTH, THT &\rightarrow 1 \\ TTT &\rightarrow 0 \end{aligned}$$

Or event  $B$ : "Getting both tails and heads":

$$\begin{aligned} HHT, HTH, HTT, THH, THT, TTH, TTT &\rightarrow 1 \quad (\text{Yes}) \\ HHH, TTT &\rightarrow 0 \quad (\text{No}) \end{aligned}$$

These values may be viewed as values assumed by the **random variable**  $X$ ,  $X$  can be "number of heads appearing" or "getting both tails and heads state", but *not both*.

**Definition 2.1.1.** A **random variable** is a function that associates a real number with each element in the sample space.

We shall use a capital letter  $X$ , to denote the random variable and its corresponding smaller letter  $x$ , for **one of its values**. A typical mistake is forgetting to define the meaning of random variable  $X$ , so *please do not forget it in your exam*. Let's continue by looking at some examples of random variables and sample spaces.

Roll the die and observe the landing value.  $X$  is the random variable defined as the value we observed:

$$S_1 = \{1, 2, 3, 4, 5, 6\}$$



Roll the die repeatedly until the sixth appears 6 times.  $Y$  is the random variable defined as the number of rolling:

$$S_2 = \{1, 2, 3, \dots\}$$

Use **exponential distribution** to predict if tomorrow will be a rainy day.  $Z$  is the random variable defined as the probability of the event "Tomorrow will be a rainy day":

$$S_3 = \{z \mid 0 < z < 1\}$$

The random variable  $X$  can take one of the values  $x_1 = 1, x_2 = 2, \dots, x_6 = 6$ , and similarly with  $Y$  and  $Z$  can take one of their own values in their sample spaces  $S_2, S_3$  respectively. Now we are interested in classifying 2 types of sample spaces and their random variables.

**Definition 2.1.2.** *If a sample space contains a finite number of possibilities or an unending sequence with as many elements as there are whole numbers, it is called a **discrete sample space**.*

**Definition 2.1.3.** *If a sample space contains an infinite number of possibilities equal to the number of points on a line segment, it is called a **continuous sample space**.*

So  $S_1$  and  $S_2$  are **discrete sample spaces** and  $X, Y$  are called **discrete random variables**;  $S_3$  is **continuous sample space** and  $Z$  is called **continuous random variable**.

## 2.2 Discrete Probability Distributions

Consider coin tossing experiment, now if we assign random variable  $X$  as the number of heads appearing, I can obtain some useful results:

$$\begin{aligned} P(X = 3) &= P(\{HHH\}) = \frac{1}{8} = 0.125 \\ P(X = 2) &= P(\{HHT, HTH, THH\}) = \frac{3}{8} = 0.375 \\ P(X = 1) &= P(\{TTH, THT, HTT\}) = \frac{3}{8} = 0.375 \\ P(X = 0) &= P(\{TTT\}) = \frac{1}{8} = 0.125 \end{aligned}$$

Or in table form:

| $x$        | 0     | 1     | 2     | 3     |
|------------|-------|-------|-------|-------|
| $P(X = x)$ | 0.125 | 0.375 | 0.375 | 0.125 |

Frequently, it is much convenient to represent all of the probabilities of a random variable  $X$  by a **formula**.

**Definition 2.2.1.** *The function  $f(x)$  is a **probability density function** (pdf) of the discrete random variable  $X$  if, for each possible outcome  $X$ :*

$$\begin{cases} f(x) \geq 0 \\ \sum_x f(x) = 1 \\ P(X = x) = f(x) \end{cases}$$

You can derive that the pdf of coin tossing experiment above is:

$$f(x) = P(X = x) = \frac{\binom{3}{x}}{2^3} \quad (x = 0, 1, 2, 3)$$

In many cases, we want to know the probability of  $(X \leq x)$  (you will see it clearly in the next section). For instance, I want to know the probability of heads appearing a maximum of 2 times.

$$P(X \leq 2) = P(X = 0) + P(X = 1) + P(X = 2) = 0.125 + 0.375 + 0.375 = 0.875$$

In general, we define a new function  $F(x)$  to handle these cases as follows:

**Definition 2.2.2.** The **cumulative distribution function** (cdf)  $F(x)$  of a discrete random variable  $X$  with probability distribution  $f(x)$  is:

$$F(x) = P(X \leq x) = \sum_{t \leq x} f(t) \quad (-\infty < x < +\infty)$$

For example, in the above experiment:

$$F(x) = \begin{cases} 0 & (x < 0) \\ 0.125 & (0 \leq x < 1) \\ 0.5 & (1 \leq x < 2) \\ 0.875 & (2 \leq x < 3) \\ 1 & (x \geq 3) \end{cases}$$

## 2.3 Continuous Probability Distributions

Choosing randomly a number within the range  $(0, 1)$ .  $X$  is the random variable, defined as the chosen one. In the previous chapter, we have discussed that the probability of getting a **single number** in the range  $(0, 1)$  is 0.

$$P(X = 0.7) = 0$$

So 0.7 will never be chosen since  $P(X = 0.7) = 0$ ? Now think carefully about the reasons why the probability of an event might be zero. Recall the **Theorem 1.3.1**:

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

There are 2 main reasons that could explain why  $P(A)$  can be zero; the first one is  $n = 0$  **and**  $N$  **is a finite number**, and the second one is  $n$  **is a finite number and**  $N$  **is an infinite number**. This might be the big misconception, since people always claim that the only reason for  $P(A) = 0$  is the first one, and forgot the second. But now you can clearly see that if  $n > 0$ , event will always have a chance of happening. So now we conclude certainly: " $P(A) = 0$  does not mean event  $A$  will never happen."

If sample space  $S$  is *continuous sample space*, which contains *an infinite number of possibilities equal to the number of points on a line segment*, we do not care about the probability of **a single sample point** occurring (because it is always equal 0). We shift our attention to the probability of **the interval that our concern sample point may be fallen inside** occurring.



Figure 2.1: What is the probability that a number will fall within this range?

Intuitively you can conclude that:

$$P(0.6 < X < 0.78) = 0.78 - 0.6 = 0.18$$

Similarly with the discrete random variables, we can also define:

**Definition 2.3.1.** The function  $f(x)$  is a **probability distribution (or density) function** (pdf) for the continuous random variable  $X$ , defined over the set of real numbers, if:

$$\begin{cases} f(x) \geq 0, \text{ for all } x \in R \\ \int_{-\infty}^{+\infty} f(x)dx = 1 \\ P(a < X < b) = \int_a^b f(x)dx \end{cases}$$

Verifying yourself that the pdf of number choosing experiment is:

$$f(x) = \begin{cases} 1 & (0 < x < 1) \\ 0 & (\text{elsewhere}) \end{cases}$$

This pdf is the simplest case of **uniform distribution function**. As I mentioned before, in the continuous sample space case, concerning on the probability of  $(X = x)$  does not make any sense since it is equal 0 and we can not use any information about it. Instead, we turn our focus to the probability of  $(X < x)$  and define the **cumulative distribution function**  $F(x)$  as follows:

**Definition 2.3.2.** The **cumulative distribution function** (cdf)  $F(x)$  of a continuous random variable  $X$  with pdf  $f(x)$  is:

$$F(x) = P(X < x) = \int_{-\infty}^x f(t)dt \quad (-\infty < x < +\infty)$$

The cdf of number choosing experiment is:

$$F(x) = \begin{cases} 0 & (x < 0) \\ x & (0 \leq x < 1) \\ 1 & (x \geq 1) \end{cases}$$

## 2.4 Joint Probability Distributions

### 2.4.1 Case of Discrete Random Variables

In the previous sections, we have considered **single random variable** and its **probability distribution function** case, and restricted ourselves to **one-dimension sample space**. But now we are interested in observing the **simultaneous outcomes** of **several random**

**variables.** For example, you toss 2 fair coins simultaneously and observe their appearing faces. You define 2 random variables,  $X$  for the first coin face, and  $Y$  for the second one. The heads is assigned value 1, and the tails is assigned value 0. These are some results obtained from this experiment:

$$P(X = 0, Y = 0) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} = 0.25$$

$$P(X = 0, Y = 1) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} = 0.25$$

$$P(X = 1, Y = 0) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} = 0.25$$

$$P(X = 1, Y = 1) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} = 0.25$$

These results can be written in table form:

| $P(X = x, Y = y)$ | $Y = 0$ | $Y = 1$ |
|-------------------|---------|---------|
| $X = 0$           | 0.25    | 0.25    |
| $X = 1$           | 0.25    | 0.25    |

Another classic example of **discrete joint probability distribution** is the problem of picking balls from a basket. Now you have a basket with many colorful balls inside; there are 3 red, 4 green and 5 blue balls. You choose randomly 3 balls from the basket, and you define 2 random variables  $X$  and  $Y$ ;  $X$  is the number of red balls and  $Y$  is the number of green balls. Using counting techniques and combinations, you can write the **joint pdf**:

$$f(x, y) = P(X = x, Y = y) = \frac{\binom{3}{x} \binom{4}{y} \binom{5}{3-x-y}}{\binom{12}{3}}$$

Or in table form:

| $P(X = x, Y = y)$ | $Y = 0$         | $Y = 1$        | $Y = 2$         | $Y = 3$        |
|-------------------|-----------------|----------------|-----------------|----------------|
| $X = 0$           | $\frac{1}{22}$  | $\frac{2}{11}$ | $\frac{3}{22}$  | $\frac{1}{55}$ |
| $X = 1$           | $\frac{3}{22}$  | $\frac{3}{11}$ | $\frac{9}{110}$ | -              |
| $X = 2$           | $\frac{3}{44}$  | $\frac{3}{55}$ | -               | -              |
| $X = 3$           | $\frac{1}{220}$ | -              | -               | -              |

**Definition 2.4.1.** The function  $f(x, y)$  is a **joint pdf** of the **discrete random variables**  $X$  and  $Y$  if:

$$\begin{cases} f(x, y) \geq 0 \\ \sum_x \sum_y f(x, y) = 1 \\ P(X = x, Y = y) = f(x, y) \end{cases}$$

**Corollary 2.4.0.1.** For any region  $A$  in the  $xy$  plane:

$$P[(X, Y) \in A] = \sum \sum_A f(x, y)$$

After considering the case where both variables  $X$  and  $Y$  are both varying, what if we **fix** one variable and vary the other? You can see this idea appearing very naturally in the process of finding the values of the above table using a calculator. We introduce the new functions called **marginal distributions** of  $X$  and  $Y$  alone.

**Definition 2.4.2.** The **marginal distributions** of  $X$  alone and of  $Y$  alone for the **discrete case** are:

$$g(x) = \sum_y f(x, y); \quad h(y) = \sum_x f(x, y)$$

Since the general form of  $g(x)$  and  $h(y)$  are not easy to be generalized, and the range of discrete random variables  $X$  and  $Y$  is very narrow, so we should reuse the table above to obtain values of these functions.

| $P(X = x, Y = y)$ | $Y = 0$                | $Y = 1$                | $Y = 2$                | $Y = 3$               | $g(x)$                  |
|-------------------|------------------------|------------------------|------------------------|-----------------------|-------------------------|
| $X = 0$           | $\frac{1}{22}$         | $\frac{2}{11}$         | $\frac{3}{22}$         | $\frac{1}{55}$        | $g(0) = \frac{21}{55}$  |
| $X = 1$           | $\frac{3}{22}$         | $\frac{1}{11}$         | $\frac{9}{110}$        | -                     | $g(1) = \frac{27}{55}$  |
| $X = 2$           | $\frac{3}{44}$         | $\frac{3}{55}$         | -                      | -                     | $g(2) = \frac{27}{220}$ |
| $X = 3$           | $\frac{1}{220}$        | -                      | -                      | -                     | $g(3) = \frac{1}{220}$  |
| $h(y)$            | $h(0) = \frac{14}{55}$ | $h(1) = \frac{28}{55}$ | $h(2) = \frac{12}{55}$ | $h(3) = \frac{1}{55}$ | 1                       |

**Corollary 2.4.0.2.** If  $g(x)$  and  $h(y)$  are marginal distributions of  $X$  alone and  $Y$  alone for the **discrete case**, then:

$$\sum_x g(x) = \sum_y h(y) = 1$$

**Corollary 2.4.0.2** is directly derived from the **Definition 2.4.1**, and you can verify them by using the probability distribution table.

After defining the **marginal distribution**, now we can see clearly the connection between them and regular pdfs are:

$$P(X = x) = g(x); \quad P(Y = y) = h(y)$$

Now if I choose 3 balls from the basket, and I know two of them are green, what is the probability that the remaining ball is red?

$$P(X = 1|Y = 2) = \frac{P(X = 1, Y = 2)}{P(Y = 2)} = \frac{\frac{9}{110}}{\frac{12}{55}} = \frac{3}{8} = 0.375$$

In general, it is not hard to deduce these formulas:

$$f(x|y) = P(X = x|Y = y) = \frac{P(X = x, Y = y)}{P(Y = y)} = \frac{f(x, y)}{h(y)}$$

$$f(y|x) = P(Y = y|X = x) = \frac{P(Y = y, X = x)}{P(X = x)} = \frac{f(x, y)}{g(x)}$$

**Definition 2.4.3.** Let  $X$  and  $Y$  be two **discrete random variables**. The **conditional distribution** of the variable  $Y$  given that  $(X = x)$  is:

$$f(y|x) = \frac{f(x, y)}{g(x)}$$

Similarly for the reverse case:

$$f(x|y) = \frac{f(x, y)}{h(y)}$$

From the previous section **Chapter 1.4: Conditional Probability**, you already knew how to identify two independent events, by checking this condition:

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

Now we use it again, with random variables  $X$  and  $Y$ :

$$P(X = x, Y = y) = P(X = x)P(Y = y) \Leftrightarrow f(x, y) = g(x)h(y)$$

**Definition 2.4.4.** Let  $X$  and  $Y$  be two **discrete random variables**, with joint pdf  $f(x, y)$  and marginal distributions  $g(x)$ ,  $h(y)$ , respectively. They are said to be **statistically independent** if and only if:

$$f(x, y) = g(x)h(y)$$

If we check the case  $(X = 1, Y = 1)$ :

$$\left(f(1, 1) = \frac{3}{11}\right) \neq \left(g(1)h(1) = \frac{27}{55} \cdot \frac{28}{55}\right)$$

So  $X$  and  $Y$  are **not** statistically independent, they are interdependent. But how interdependent are they? Are they highly or minimally interdependent? This question will be answered at the end of the next chapter.

## 2.4.2 Case of Continuous Random Variables

Similarly, the joint pdf of **continuous random variable** can also be defined as follows:

**Definition 2.4.5.** The function  $f(x, y)$  is a **joint pdf** of the **continuous random variables**  $X$  and  $Y$  if:

$$\begin{cases} f(x, y) \geq 0 \\ \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f(x, y) dx dy = 1 \\ P(a < X < b, c < Y < d) = \int_c^d \int_a^b f(x, y) dx dy \end{cases}$$

**Corollary 2.4.0.3.** For any region  $A$  in the  $xy$  plane:

$$P[(X, Y) \in A] = \iint_A f(x, y) dx dy$$

Almost all the formulas in the above section are redefined exactly the same, with a slight difference in the notation for the integral (you can read more about Riemann sum to see the relationship between integrals and infinite discrete sums).

**Definition 2.4.6.** The **marginal distributions** of  $X$  alone and of  $Y$  alone for the **continuous case** are:

$$g(x) = \int_{-\infty}^{+\infty} f(x, y) dy; \quad h(y) = \int_{-\infty}^{+\infty} f(x, y) dx$$

**Corollary 2.4.0.4.** If  $g(x)$  and  $h(y)$  are marginal distributions of  $X$  alone and  $Y$  alone for the **continuous case**, then:

$$\int_{-\infty}^{+\infty} g(x) dx = \int_{-\infty}^{+\infty} h(y) dy = 1$$

**Definition 2.4.7.** Let  $X$  and  $Y$  be two **continuous random variables**. The **conditional distribution** of the variable  $Y$  given that  $(X = x)$  is:

$$f(y|x) = \frac{f(x, y)}{g(x)}$$

Similarly for the reverse case:

$$f(x|y) = \frac{f(x, y)}{h(y)}$$

But notice that if  $X$  and  $Y$  are continuous random variables, then  $P(X = x|Y = y)$  is always equal zero! The correct way to obtain conditional probability value is shown below:

**Corollary 2.4.0.5.** *If  $X$  and  $Y$  be two **continuous random variables**, then:*

$$P(a < X < b|Y = y) = \int_a^b f(x|y)dx$$

**Definition 2.4.8.** *Let  $X$  and  $Y$  be two **continuous random variables**, with joint pdf  $f(x, y)$  and marginal distributions  $g(x)$ ,  $h(y)$ , respectively. They are said to be **statistically independent** if and only if:*

$$f(x, y) = g(x)h(y)$$

# Chapter 3

## Mathematical Expectation

Due to the relatively high degree of similarity between the formulas in the case of discrete and continuous random variables, this chapter approaches them **in parallel**, rather than separating them like previous chapter.

### 3.1 Mean of a Random Variable

After conducting experiment, now it is time to process your obtained results. Toss a fair coin for 3 times, and define the discrete random variable  $X$  as the number of heads appearing. The pdf of this experiment is:

$$f(x) = P(X = x) = \frac{\binom{3}{x}}{2^3} \quad (x = 0, 1, 2, 3)$$

Or in table form:

| $x$        | 0     | 1     | 2     | 3     |
|------------|-------|-------|-------|-------|
| $P(X = x)$ | 0.125 | 0.375 | 0.375 | 0.125 |

What is the **mean** or **expected value** of  $X$ ? Intuitively, you will use the formula:

$$\mu_X = E(X) = \sum_x x f(x) = 0 \cdot 0.125 + 1 \cdot 0.375 + 2 \cdot 0.375 + 3 \cdot 0.125 = 1.5$$

By illustrating the above results with a graph, the position of  $\mu$  is shown.

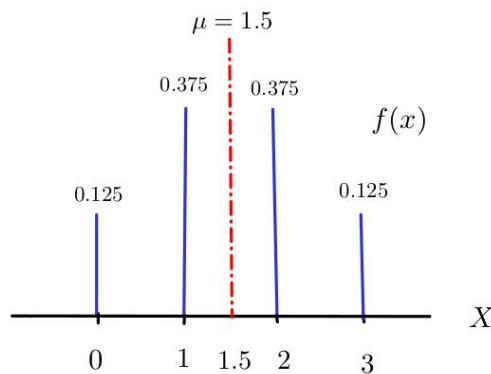


Figure 3.1: Pdf of fair coin tossing experiment with its mean



If our coin is unfair, assume that the probability of getting heads is  $p = 0.3$ , then the probability of getting tails is  $q = 1 - p = 0.7$ . Now the pdf is:

$$f(x) = P(X = x) = \binom{3}{x} p^x q^{3-x} = \binom{3}{x} 0.3^x 0.7^{3-x} \quad (x = 0, 1, 2, 3)$$

Or in table form:

| $x$        | 0     | 1     | 2     | 3     |
|------------|-------|-------|-------|-------|
| $P(X = x)$ | 0.343 | 0.441 | 0.189 | 0.027 |

In this experiment, we **expect** to get the **average value** (or mean) of  $X$ :

$$\mu_X = E(X) = \sum_x x f(x) = 0.0.343 + 1.0.441 + 2.0.189 + 3.0.027 = 0.9$$

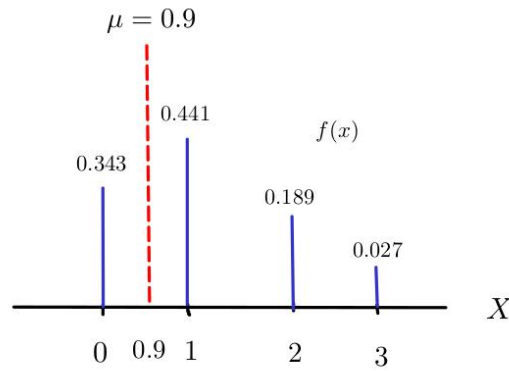


Figure 3.2: Pdf of unfair coin tossing experiment with its mean

**Definition 3.1.1.** Let  $X$  be a random variable with pdf  $f(x)$ . The **expected value** or **mean** of  $X$  is:

$$\mu_X = E(X) = \sum_x x f(x) \quad (\text{if } X \text{ is discrete})$$

$$\mu_X = E(X) = \int_{-\infty}^{+\infty} x f(x) dx \quad (\text{if } X \text{ is continuous})$$

Now imagine you are playing coin tossing game (do not treat it like an experiment). As I mention the rule above, you have 3 turns to toss a coin. If you see 3 heads appear, you are extremely lucky today; but if you do not see any, do not be sad because life is long. A very natural thought occurred to me that we should quantify a player's "luck level" with a quantitative value. Random variable  $Y$  takes this role and can be defined as:

$$Y = \frac{X}{3}$$

Now we are interested in average "luck level" of coin tossing game, so reuse the previous table that we have created:

| $x$               | 0     | 1     | 2     | 3     |
|-------------------|-------|-------|-------|-------|
| $y = \frac{x}{3}$ | 0     | 0.333 | 0.666 | 1     |
| $P(X = x)$        | 0.343 | 0.441 | 0.189 | 0.027 |

The **expected value** or **mean** of "luck level" is:

$$\mu_Y = E(Y) = \sum_y yf(x) = 0.0.343 + 0.333.0.441 + 0.666.0.189 + 1.0.027 = 0.3$$

Since 0.3 is not so high value, so the **mean** of "luck level" is pretty small and players should not look forward to their chances in this game. Furthermore, you should notice the subtle connection between two random variables  $X$  and  $Y$  as:

$$Y = \frac{X}{3} \Leftrightarrow E(Y) = \frac{E(X)}{3}$$

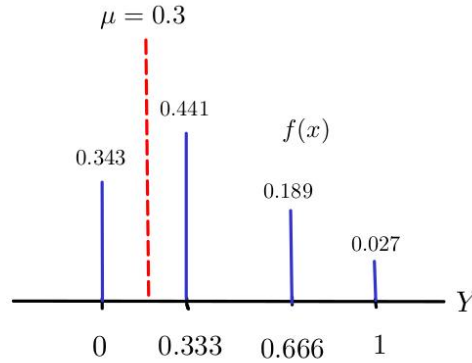


Figure 3.3: The mean value of "luck level"

If  $Y$  is defined as **function of random variable**  $X$  or ( $Y = g(X)$ ), then the **expected value** of it can be calculated using the theorem below:

**Theorem 3.1.1.** *Let  $X$  be a random variable with pdf  $f(x)$ , and the **expected value** of the random variable  $g(X)$  is:*

$$\mu_{g(X)} = E(g(X)) = \sum_x g(x)f(x) \quad (\text{if } X \text{ is discrete})$$

$$\mu_{g(X)} = E(g(X)) = \int_{-\infty}^{+\infty} g(x)f(x)dx \quad (\text{if } X \text{ is continuous})$$

Now I want to make our game funnier because just tossing a coin is too boring, so I came up with an idea. How about toss 3 unfair coins and the bottle cap **at the same time**? The random variable  $X$  is defined as the number of heads appearing with the probability of getting heads is  $p = 0.3$ , and getting tails is  $q = 1 - p = 0.7$ , so the pdf is:

$$g(x) = P(X = x) = \binom{3}{x} p^x q^{3-x} = \binom{3}{x} 0.3^x 0.7^{3-x} \quad (x = 0, 1, 2, 3)$$

A 'PEPSI' bottle cap has just been founded, so I define a new random variable  $Y$  as its face after tossing. The 'PEPSI' face is assigned the value 1, and the other is assigned the value 0. Assume that:

$$\begin{aligned} P(Y = 1) &= p' = 0.8 \\ P(Y = 0) &= q' = 1 - p' = 0.2 \end{aligned}$$

So the pdf of bottle cap tossing experiment is:

$$h(y) = P(Y = y) = (p')^y(q')^{1-y} = 0.8^y 0.2^{1-y} \quad (y = 0, 1)$$

Intuitively, you can certainly conclude that  $X$  and  $Y$  are **statistically independent** since coins and bottle cap tossing process do not affect each other. So the **joint pdf** of them is:

$$f(x, y) = g(x)h(y) = \binom{3}{x} 0.3^x 0.7^{3-x} 0.8^y 0.2^{1-y} \quad (x = 0, 1, 2, 3; y = 0, 1)$$

Or in table form:

| $P(X = x, Y = y)$ | $X = 0$        | $X = 1$        | $X = 2$        | $X = 3$        | $h(y)$       |
|-------------------|----------------|----------------|----------------|----------------|--------------|
| $Y = 0$           | 0.0686         | 0.0882         | 0.0378         | 0.0054         | $h(0) = 0.2$ |
| $Y = 1$           | 0.2744         | 0.3528         | 0.1512         | 0.0216         | $h(1) = 0.8$ |
| $g(x)$            | $g(0) = 0.343$ | $g(1) = 0.441$ | $g(2) = 0.189$ | $g(3) = 0.027$ | 1            |

Now "luck level" can be defined as the value of  $(X + Y)$ . Thinking according to the same logic, the **mean** value of  $(X + Y)$  is:

$$\begin{aligned}
\mu_{(X+Y)} &= E(X + Y) = \sum_x \sum_y (x + y) f(x, y) = \sum_x \left( \sum_y (x + y) f(x, y) \right) \\
&= \sum_x (x f(x, 0) + x f(x, 1) + 0 f(x, 0) + 1 f(x, 1)) \\
&= \sum_x (x f(x, 0) + x f(x, 1) + f(x, 1)) \\
&= 0 f(0, 0) + 1 f(1, 0) + 2 f(2, 0) + 3 f(3, 0) + 0 f(0, 1) + 1 f(1, 1) \\
&\quad + 2 f(2, 1) + 3 f(3, 1) + f(0, 1) + f(1, 1) + f(2, 1) + f(3, 1) \\
&= 0.0882 + 2.0.0378 + 3.0.0054 + 0.3528 + 2.0.1512 \\
&\quad + 3.0.00216 + 0.2744 + 0.3528 + 0.1512 + 0.0216 \\
&= 1.7
\end{aligned}$$

So now, the players can confidently look forward for their opportunity in this game!

**Definition 3.1.2.** Let  $X$  and  $Y$  be random variables with **joint pdf**  $f(x, y)$ . The **mean** or **expected value** of the random variable function  $g(X, Y)$  is:

$$\mu_{g(X,Y)} = E(g(X, Y)) = \sum_x \sum_y g(x, y) f(x, y) \quad (\text{if } X \text{ and } Y \text{ are discrete})$$

$$\mu_{g(X,Y)} = E(g(X, Y)) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} g(x, y) f(x, y) dx dy \quad (\text{if } X \text{ and } Y \text{ are continuous})$$

If  $g(X, Y) = X$ , then we obtain some useful results and can be represented as corollaries:

**Corollary 3.1.1.1.** Let  $X$  and  $Y$  be discrete random variables with **joint pdf**  $f(x, y)$ , the **mean** of  $X$  is:

$$\mu_X = E(X) = \sum_x \sum_y x f(x, y) = \sum_x x \left( \sum_y f(x, y) \right) = \sum_x x g(x)$$

Similarly, the mean of  $Y$  is:

$$\mu_Y = E(Y) = \sum_y y h(y)$$

**Corollary 3.1.1.2.** Let  $X$  and  $Y$  be continuous random variables with **joint pdf**  $f(x, y)$ , the **mean** of  $X$  is:

$$\mu_X = E(X) = \int_{-\infty}^{+\infty} xg(x)dx$$

Similarly, the mean of  $Y$  is:

$$\mu_Y = E(Y) = \int_{-\infty}^{+\infty} yh(y)dy$$

## 3.2 Variance and Covariance of Random Variables

### 3.2.1 Variance of Random Variables

The concept of **mean** or **expected value** is very important in ProbStat. In any probability distribution graph, we use the **mean value** as **main reference point** (although the mean is not always in the center of the graph). We are very interested in how the data is distributed **around the mean value**. Is it distributed far or close to it? How does it spread?

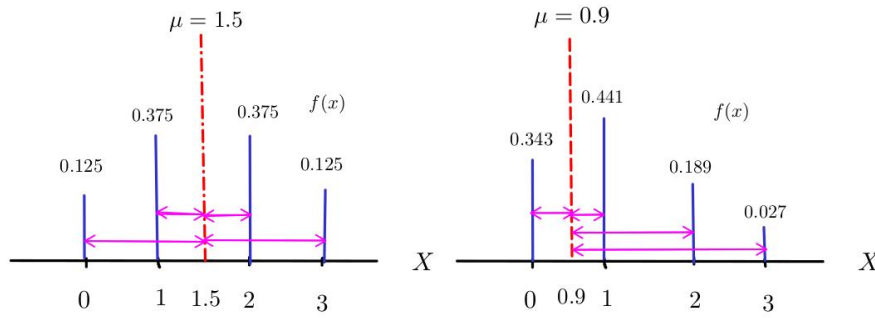


Figure 3.4: How is the data distributed around the **mean**?

The sums of total squares of the distances from the mean are:

$$\begin{aligned} \sum &= (0 - 1.5)^2 + (1 - 1.5)^2 + (2 - 1.5)^2 + (3 - 1.5)^2 = 5 \text{ (fair coin case)} \\ \sum &= (0 - 0.9)^2 + (1 - 0.9)^2 + (2 - 0.9)^2 + (3 - 0.9)^2 = 6.44 \text{ (unfair coin case)} \end{aligned}$$

Because  $6.44 > 5$ , and you can also see on the graph, the data from the unfair coin toss experiment is more **widely distributed** compare to the fair coin case. For simplicity, we often take **the average of sums** and call it **variance**.

**Definition 3.2.1.** Let  $X$  be a random variable with pdf  $f(x)$  and mean  $\mu$ . The **variance** of  $X$  is:

$$\sigma_X^2 = E((X - \mu_X)^2) = \sum_x (x - \mu_X)^2 f(x) \quad (\text{if } X \text{ is discrete})$$

$$\sigma_X^2 = E((X - \mu_X)^2) = \int_{-\infty}^{+\infty} (x - \mu_X)^2 f(x) dx \quad (\text{if } X \text{ is continuous})$$

The positive square root of the variance,  $\sigma_X$ , is called the **standard deviation** of  $X$ .

**Theorem 3.2.1.** The **variance** of a random variable  $X$  is:

$$\sigma_X^2 = E(X^2) - \mu_X^2$$

*Proof.* We only prove for the discrete case because proving for the continuous case using exactly the same logic:

$$\begin{aligned}\sigma_X^2 &= E((X - \mu_X)^2) = \sum_x (x - \mu_X)^2 f(x) = \sum_x x^2 f(x) - 2\mu_X \sum_x x f(x) + \mu_X^2 \sum_x f(x) \\ &= E(X^2) - 2\mu_X^2 + \mu_X^2 = E(X^2) - \mu_X^2\end{aligned}$$

□

At this point you may wonder "Why do not we care about the sum of total **distances** itself?" Well, there are 3 main reasons; the first is absolutely avoiding **negative values** from miscalculation, the second is working with  $\sigma_X^2$  is much easier than  $\sigma_X$  and you will understand the final reason after reading **Chapter 10: Simple Linear Regression**. Using a completely similar line of thinking as before, we can also develop the following formulas:

**Theorem 3.2.2.** *Let  $X$  be a random variable with pdf  $f(X)$ . The **variance** of the random variable  $g(X)$  is:*

$$\sigma_{g(X)}^2 = E((g(X) - \mu_{g(X)})^2) = \sum_x (g(x) - \mu_{g(X)})^2 f(x) \quad (\text{if } X \text{ is discrete})$$

$$\sigma_{g(X)}^2 = E((g(X) - \mu_{g(X)})^2) = \int_{-\infty}^{+\infty} (g(x) - \mu_{g(X)})^2 f(x) dx \quad (\text{if } X \text{ is continuous})$$

**Corollary 3.2.2.1.** *The **variance** of random variable  $g(X)$  is:*

$$\sigma_{g(X)}^2 = E(g(X)^2) - \mu_{g(X)}^2$$

*Proof.* We only prove for the discrete case:

$$\begin{aligned}\sigma_{g(X)}^2 &= E((g(X) - \mu_{g(X)})^2) = \sum_x (g(x) - \mu_{g(X)})^2 f(x) \\ &= \sum_x g(x)^2 f(x) - 2\mu_{g(X)} \sum_x g(x) f(x) + \mu_{g(X)}^2 \sum_x f(x) \\ &= E(g(X)^2) - 2\mu_{g(X)}^2 + \mu_{g(X)}^2 = E(g(X)^2) - \mu_{g(X)}^2\end{aligned}$$

□

### 3.2.2 Covariance of Random Variables

In this subsection, we will answer the question: "How are 2 random variables  $X$  and  $Y$  interdependent?" after checking the **statistically independent criterion**:

$$f(x, y) = g(x)h(y)$$

We define a new concept called **covariance**, which is a **measure of the nature of the association between the two**.

**Definition 3.2.2.** *Let  $X$  and  $Y$  be random variables with joint pdf  $f(x, y)$ . The covariance of  $X$  and  $Y$  is:*

$$\sigma_{XY} = E((X - \mu_X)(Y - \mu_Y)) = \sum_x \sum_y (x - \mu_X)(y - \mu_Y) f(x, y) \quad (\text{if } X \text{ and } Y \text{ are discrete})$$

$$\sigma_{XY} = E((X - \mu_X)(Y - \mu_Y)) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} (x - \mu_X)(y - \mu_Y) f(x, y) dx dy \quad (\text{if } X \text{ and } Y \text{ are continuous})$$

Applying directly **Definition 3.2.2** is not convenient in many situations, so we should rewrite its formula:

**Corollary 3.2.2.2.** *The **covariance** of random variable  $X$  and  $Y$  is:*

$$\sigma_{XY} = E(XY) - \mu_X \mu_Y$$

*Proof.* We only prove for the discrete case:

$$\begin{aligned} \sigma_{XY} &= \sum_x \sum_y (x - \mu_X)(y - \mu_Y) f(x, y) \\ &= \sum_x \sum_y xy f(x, y) - \mu_Y \sum_x \sum_y x f(x, y) - \mu_X \sum_x \sum_y y f(x, y) + \mu_X \mu_Y \sum_x \sum_y f(x, y) \\ &= E(XY) - \mu_Y \sum_x x g(x) - \mu_X \sum_y y h(y) + \mu_X \mu_Y \\ &= E(XY) - \mu_Y \mu_X - \mu_X \mu_Y + \mu_X \mu_Y \\ &= E(XY) - \mu_X \mu_Y \end{aligned}$$

□

Although the **covariance** between two random variables does provide information regarding the nature of the relationship, the magnitude of  $\sigma_{XY}$  does not indicate anything regarding the strength of the relationship, since  $\sigma_{XY}$  is not scale-free quantity. There is a scale-free version of the covariance called the **correlation coefficient** that is widely used in ProbStat.

**Definition 3.2.3.** *Let  $X$  and  $Y$  be random variables with **covariance**  $\sigma_{XY}$  and **standard deviations**  $\sigma_X$  and  $\sigma_Y$ , respectively. The **correlation coefficient** of  $X$  and  $Y$  is:*

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

**Theorem 3.2.3.** *The **absolute value of correlation coefficient**  $\rho_{XY}$  is always smaller or equal 1.*

$$|\rho_{XY}| \leq 1$$

*Proof.* Proving this theorem above is quiet difficult and requires knowledge of **Linear Algebra**. I recommend skipping this proof and moving on to the next section to learn some **useful theorems** first, and then back here later.

In the **random variables vector space** sharing the same joint pdf, we now define the **inner product** is:

$$\langle X, Y \rangle = E(XY)$$

With  $X$ ,  $Y$  and  $W$  are three vectors of this vector space and for any constant  $c \in \mathbb{R}$ , now we check if our **definition** satisfies 5 axioms of inner product.

$$\begin{aligned} \langle X, Y \rangle &= \langle Y, X \rangle \leftrightarrow E(XY) = E(YX) \\ c \langle X, Y \rangle &= \langle cX, Y \rangle \leftrightarrow cE(XY) = E(cXY) \\ \langle X, Y + W \rangle &= \langle X, Y \rangle + \langle X, W \rangle \leftrightarrow E(X(Y + W)) = E(XY) + E(XW) \\ \langle X, X \rangle &\geq 0 \leftrightarrow E(X^2) \geq 0 \\ \langle X, X \rangle &= 0 \Leftrightarrow X = 0 \Leftrightarrow E(X^2) = 0 \Leftrightarrow X = 0 \end{aligned}$$

The original inequality is equivalent to:

$$\begin{aligned} |\rho_{XY}| \leq 1 &\Leftrightarrow \sigma_{XY}^2 \leq (\sigma_X \sigma_Y)^2 \\ &\Leftrightarrow E^2((X - \mu_X)(Y - \mu_Y)) \leq E((X - \mu_X)^2)E((Y - \mu_Y)^2) \end{aligned}$$

Define 2 new random variables  $A$  and  $B$  (or vectors):

$$A = X - \mu_X$$

$$B = Y - \mu_Y$$

So we have to prove:

$$E^2(AB) \leq E(A^2)E(B^2) \Leftrightarrow |E(AB)| \leq \sqrt{E(A^2)E(B^2)}$$

Rewrite in vector form, this is exactly Cauchy-Schwarz inequality:

$$|\langle A, B \rangle| \leq \|A\| \cdot \|B\|$$

Now **Theorem 3.2.3** has been completely proven. The equality of the inequality occurs if and only if  $B = kA$ , where  $k$  is any real constant.

$$B = kA \Leftrightarrow Y - \mu_Y = k(X - \mu_X) \Leftrightarrow Y = kX + (\mu_Y - k\mu_X)$$

Since  $k$  is an arbitrary value in  $\mathbb{R}$ , so we can conclude that if and only if  $|\rho_{XY}| = 1$ , then  $Y = aX + b$  ( $a, b \in \mathbb{R}$ ) or in other words, they have a linear relationship.  $\square$

**Corollary 3.2.3.1.** *If  $\rho_{XY} = 0$ , then  $X$  and  $Y$  are **statistically independent**.*

*Proof.*  $\rho_{XY} = 0 \Leftrightarrow \sigma_{XY} = E(XY) - \mu_X\mu_Y = 0 \Leftrightarrow E(XY) = \mu_X\mu_Y$ . Or in equivalent:

$$\begin{aligned} E(XY) &= \mu_X\mu_Y \\ \Leftrightarrow \sum_x \sum_y xyf(x, y) &= \sum_x xg(x) \sum_y yh(y) \\ \Leftrightarrow \sum_x \sum_y xyf(x, y) &= \sum_x \sum_y xyg(x)h(y) \end{aligned}$$

Finally,  $f(x, y) = g(x)h(y)$ , and we conclude  $X$  and  $Y$  are statistically independent.  $\square$

From the equality condition of the Cauchy-Schwarz inequality, we can deduce:

**Corollary 3.2.3.2.** *If  $\rho_{XY} = 1$ , then:*

$$Y = aX + b \quad (a > 0)$$

*And if  $\rho_{XY} = -1$ , then:*

$$Y = aX + b \quad (a < 0)$$

### 3.3 Means and Variances of Linear Combinations of Random Variables

In this section, we will derive some extremely useful theorems and corollaries. You should note that all results obtained from this section will be reused multiple times throughout the rest of this book.

**Theorem 3.3.1.** *If  $a, b$  are 2 constants, then:  $E(aX + b) = aE(X) + b$*

*Proof.* By definition of the **mean**:

$$E(aX + b) = \sum_x (ax + b)f(x) = a \sum_x xf(x) + b \sum_x f(x) = aE(X) + b$$

□

**Corollary 3.3.1.1.** *Setting  $a = 0$ , we see that  $E(b) = b$*

**Corollary 3.3.1.2.** *Setting  $b = 0$ , we see that  $E(aX) = 0$*

**Theorem 3.3.2.** *If  $g(X)$  and  $h(X)$  are two functions of random variable  $X$ , then:*

$$E(g(X) + h(X)) = E(g(X)) + E(h(X))$$

*Proof.* Also by the definition of the **expected value**:

$$E(g(X) + h(X)) = \sum_x (g(x) + h(x))f(x) = \sum_x g(x)f(x) + \sum_x h(x)f(x) = E(g(X)) + E(h(X))$$

□

**Theorem 3.3.3.** *If  $g(X, Y)$  and  $h(X, Y)$  are two functions of random variable  $(X, Y)$ , then:*

$$E(g(X, Y) + h(X, Y)) = E(g(X, Y)) + E(h(X, Y))$$

*Proof.* Applying directly the definition of the **expected value**, we have:

$$\begin{aligned} E(g(X, Y) + h(X, Y)) &= \sum_x \sum_y (g(x, y) + h(x, y))f(x, y) \\ &= \sum_x \sum_y g(x, y)f(x, y) + \sum_x \sum_y h(x, y)f(x, y) \\ &= E(g(X, Y)) + E(h(X, Y)) \end{aligned}$$

□

**Corollary 3.3.3.1.** *Setting  $g(X, Y) = g(X)$  and  $h(X, Y) = h(Y)$ , then:*

$$E(g(X) + h(Y)) = E(g(X)) + E(h(Y))$$

**Corollary 3.3.3.2.** *Setting  $g(X, Y) = X$  and  $h(X, Y) = Y$ , then:*

$$E(X + Y) = E(X) + E(Y)$$

**Theorem 3.3.4.** *If  $X$  and  $Y$  are statistically independent, then:*

$$E(XY) = \mu_X \mu_Y$$

*Proof.* We already have  $f(x, y) = g(x)h(y)$ .

$$E(XY) = \sum_x \sum_y xyf(x, y) = \sum_x \sum_y xyg(x)h(y) = \left( \sum_x xg(x) \right) \left( \sum_y yh(y) \right) = \mu_X \mu_Y$$

□

You should note that  $X$  and  $Y$  are independent random variables  $\Leftrightarrow E(XY) = \mu_X \mu_Y$



**Theorem 3.3.5.** If  $X$  and  $Y$  are two **independent random variables** with the joint pdf  $f(x, y)$  with three constants  $a, b, c$ , then:

$$\sigma_{aX+bY+c}^2 = a^2\sigma_X^2 + b^2\sigma_Y^2$$

*Proof.* Using definition, we have:

$$\begin{aligned}\sigma_{aX+bY+c}^2 &= E((aX + bY + c - \mu_{aX+bY+c})^2) \\ &= E((aX + bY + c - a\mu_X - b\mu_Y - c)^2) \\ &= E((a(X - \mu_X) + b(Y - \mu_Y))^2) \\ &= E(a^2(X - \mu_X)^2) + 2abE(X - \mu_X)(Y - \mu_Y) + E(b^2(Y - \mu_Y)^2) \\ &= a^2\sigma_X^2 + 2ab\sigma_{XY} + b^2\sigma_Y^2 \\ &= a^2\sigma_X^2 + b^2\sigma_Y^2 \quad (\sigma_{XY} = 0)\end{aligned}$$

□

**Corollary 3.3.5.1.** Setting  $b = 0$ , we see that:  $\sigma_{aX+c}^2 = a^2\sigma_X^2$

**Corollary 3.3.5.2.** Setting both  $b = c = 0$ , we see that:  $\sigma_{aX}^2 = a^2\sigma_X^2$

**Corollary 3.3.5.3.** Setting  $c = 0$ , we see that:

$$\sigma_{aX+bY}^2 = a^2\sigma_X^2 + b^2\sigma_Y^2$$

This corollary is only true when  $X$  and  $Y$  are **statistically independent**.

## 3.4 Markov's and Chebyshev's Inequalities

### 3.4.1 Markov's Inequality

**Theorem 3.4.1.** If  $X$  is a **non-negative** random variable with pdf  $f(x)$ , then:

$$P(X \geq a) \leq \frac{E(X)}{a}$$

*Proof.* From the initial condition  $X \geq 0$ :

$$E(X) = \sum_x xf(x) \geq \sum_{x \geq a} xf(x) \geq \sum_{x \geq a} af(x) = aP(X \geq a) \Rightarrow P(X \geq a) \leq \frac{E(x)}{a}$$

□

### 3.4.2 Chebyshev's Inequality

**Theorem 3.4.2.** If  $X$  is any random variables with pdf  $f(x)$ , then:

$$P(\mu_X - k\sigma_X < X < \mu_X + k\sigma_X) \geq 1 - \frac{1}{k^2}$$

*Proof.* Applying Markov's inequality with  $(X - \mu_X)^2$  as non-negative random variable, now we obtain:

$$\begin{aligned}P((X - \mu_X)^2 \geq a^2) &\leq \frac{E((X - \mu_X)^2)}{a^2} = \frac{\sigma_X^2}{a^2} \\ \Rightarrow P((X - \mu_X)^2 < a^2) &\geq 1 - \frac{\sigma_X^2}{a^2}\end{aligned}$$

Using the positive scale factor  $k$ :  $a = k\sigma_X$  from the **standard deviation**, now we obtain:

$$P(\mu_X - k\sigma_X < X < \mu_X + k\sigma_X) \geq 1 - \frac{1}{k^2}$$

□

Chebyshev's Inequality can also be written as:

$$P(|X - \mu_X| \geq k\sigma_X) \leq \frac{1}{k^2}$$

This result is very important, especially in **case of continuous random variable**, the range of  $X$  usually extends to infinity; it shows us that the first and last ends of any pdf are **bounded**. In general, the shape of a valid pdf graph is always **flattened** at both ends according to the inequality above.

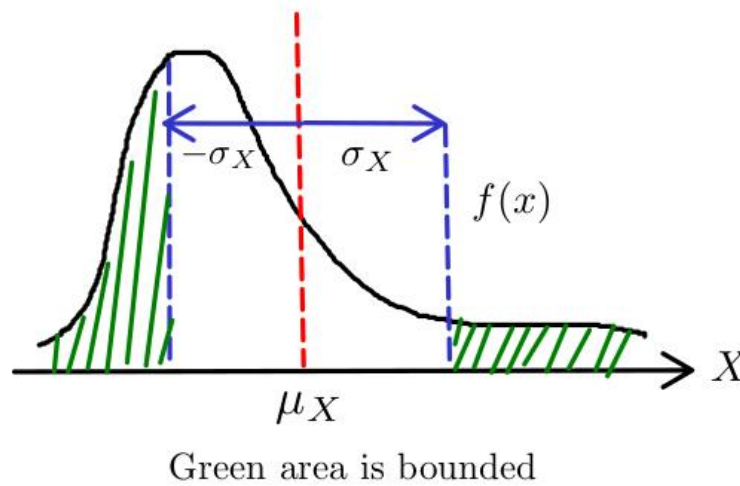


Figure 3.5: Example of a valid pdf graph

# Chapter 4

## Some Discrete Probability Distributions

From this chapter to the end of this book, if  $X$  is a random variable with pdf  $f(x)$ , we can denote it as:

$$X \sim f(x)$$

### 4.1 Bernoulli, Binomial and Poisson Distributions

#### 4.1.1 Bernoulli Distribution

Because tossing a coin many times might be so boring, so in this chapter we will begin with more vivid examples. Let's plan to plant some mung bean plants (or maybe just imagine it)! Now you have to buy a packet of mung bean seeds from the agriculture store. All seed packets clearly state the germination rate of the seeds. Assume that the germination rate is  $p = 0.8$ . Now if we define the random variable  $X$  as the germination state of **a single seed**, then the pdf of  $X$  can be represented as follows, where  $(X = 1)$  describes the germinating state of the seed, and  $(X = 0)$  describes the opposite state:

|            |     |     |
|------------|-----|-----|
| $x$        | 0   | 1   |
| $P(X = x)$ | 0.2 | 0.8 |

Or in equation form:

$$f(x) = P(X = x) = 0.8^x 0.2^{1-x} \quad (x = 0, 1)$$

The experiment to test the germination state of **a single mung bean seed** is a typical example of **Bernoulli trial**. Formally, a **Bernoulli trial** is a random experiment with exactly **two possible outcomes**: "success" and "failure". The success rate is always denoted by the letter  $p$ , and failure rate is  $q$ . Because they complement each other, so:  $q = 1 - p$ .

**Definition 4.1.1.** The **Bernoulli trial** is a random experiment with exactly **two possible outcomes**: "success" and "failure".

**Theorem 4.1.1.** The pdf of a **Bernoulli trial** is:

$$\mathcal{J}(x; 1, p) = p^x q^{1-x} \quad (x = 0, 1)$$

**Corollary 4.1.1.1.** If  $X \sim \mathcal{J}(x; 1, p)$ , then:

$$\begin{aligned}\mu_X &= p \\ \sigma_X^2 &= pq\end{aligned}$$

*Proof.* For the mean:

$$\mu_X = \sum_x x f(x) = p^1 q^0 = p$$

For the variance:

$$\sigma_X^2 = E(X^2) - \mu_X^2 = \sum_x x^2 f(x) - p^2 = p - p^2 = pq$$

□

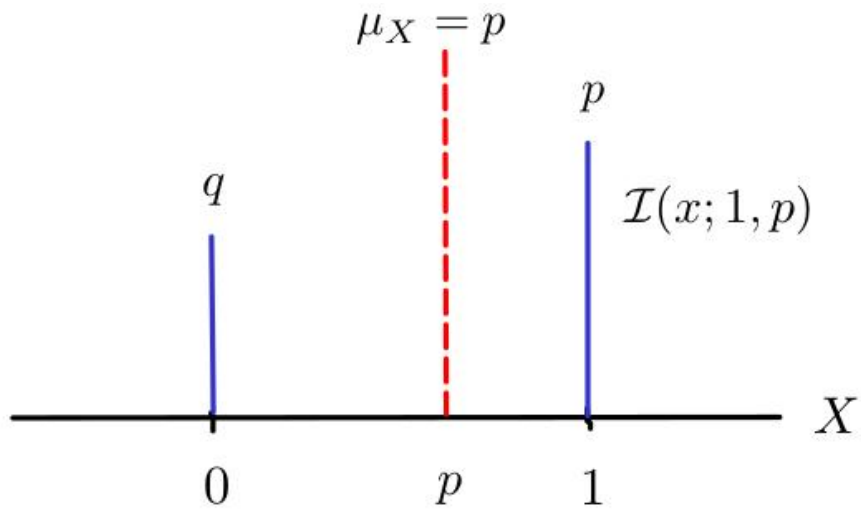


Figure 4.1: The pdf of Bernoulli distribution

### 4.1.2 Binomial Distribution

We know that the germination rate is  $p = 0.8$ . Now you pick randomly 10 mung bean seeds from the packet and put them on the wet tissue.

Perhaps after waiting 3 days, you will find out that only 6 seeds have germinated, while the rest will start to smell. You might ask "But why did only 6 seeds germinated? I thought it would be 8?". Now let me explain.

Firstly I change our previous definition of the random variable  $X$ .  $X$  is now defined as the number of seeds that successfully germinated. The probability of 6 seeds germinating is:

$$P(X = 6) = \binom{10}{6} 0.8^6 0.2^4 = 0.088$$

And the probability of 8 seeds germinating is:

$$P(X = 8) = \binom{10}{8} 0.8^8 0.2^2 = 0.301$$

Although  $0.301 > 0.088$ , it does not mean the event "6 seeds germinated" will never happen (because  $0.088 > 0$ ); but  $1 > 0.301$ , so the event "8 seeds germinated" will not always happen. This is the paradox of ProbStat, we evaluate everything through the question "**How likely will it occur?**", and our certainty is never absolute.

The pdf of our experiment is:

$$f(x) = P(X = x) = \binom{10}{x} 0.8^x 0.2^{10-x} \quad (x \in \{0, 1, 2, \dots, 10\})$$

Or in table form:

| $x$        | 0         | 1                 | 2                   | 3                   | 4                   |
|------------|-----------|-------------------|---------------------|---------------------|---------------------|
| $P(X = x)$ | $10^{-7}$ | $4 \cdot 10^{-6}$ | $7.3 \cdot 10^{-5}$ | $7.8 \cdot 10^{-4}$ | $5.5 \cdot 10^{-3}$ |

| $x$        | 5     | 6     | 7     | 8     | 9     | 10    |
|------------|-------|-------|-------|-------|-------|-------|
| $P(X = x)$ | 0.026 | 0.088 | 0.201 | 0.301 | 0.268 | 0.107 |

From the data table above, you can easily see that even if  $P(X = 8)$  reaches its maximum value, the probability of occurrence is only about 30%. Testing the germination ability of 10 seeds experiment is a classic example of the **Bernoulli process**, which consists of **repeated Bernoulli trials**.

**Definition 4.1.2.** The **Bernoulli process** is the process that consists **repeated Bernoulli trials**.

**Theorem 4.1.2.** A single Bernoulli trial can result is a sucess with probability  $p$  and a failure with probability  $q = 1 - p$ . Then the pdf of **binomial random variable**  $X$ , the **number of successes** in  $n$  independent trials is:

$$\mathcal{B}(x; n, p) = \binom{n}{x} p^x q^{n-x} \quad (x = 0, 1, 2, \dots, n)$$

**Corollary 4.1.2.1.** If  $X \sim \mathcal{B}(x; n, p)$ , then:

$$\begin{aligned} \mu_X &= np \\ \sigma_X^2 &= npq \end{aligned}$$

*Proof.* For the mean:

$$\begin{aligned}
\mu_X &= \sum_x x f(x) = \sum_{x=0}^n x \binom{n}{x} p^x q^{n-x} = \sum_{x=1}^n \frac{n!}{(x-1)!(n-x)!} p^x q^{n-x} \\
&= np \sum_{x=1}^n \frac{(n-1)!}{(x-1)!(n-x)!} p^{x-1} q^{n-x} = np \sum_{x=1}^n \binom{n-1}{x-1} p^{x-1} q^{n-x} \\
&= np(p+q)^{n-1} = np
\end{aligned}$$

Finding the variance directly is not easy, so we have to perform small trick here:

$$\begin{aligned}
E(X(X-1)) &= \sum_x x(x-1) f(x) = \sum_{x=1}^n x(x-1) \binom{n}{x} p^x q^{n-x} \\
&= p^2 n(n-1) \sum_{x=2}^n \frac{(n-2)!}{(n-x)!(x-2)!} p^{x-2} q^{n-x} \\
&= p^2 n(n-1) \sum_{x=2}^n \binom{n-2}{x-2} p^{x-2} q^{n-x} \\
&= p^2 n(n-1)(p+q)^{n-2} = p^2 n(n-1)
\end{aligned}$$

Now we apply the result above to the definition of variance:

$$\sigma_X^2 = E(X^2) - \mu_X^2 = E(X(X-1)) + E(X) - \mu_X^2 = p^2 n(n-1) + np - n^2 p^2 = npq$$

□

*Proof.* There is another subtle way to prove **Corollary 4.1.2.1**, if  $X \sim \mathcal{B}(x; n, p)$  then:

$$X = I_1 + I_2 + \cdots + I_n$$

where each  $I_i \sim \mathcal{J}(x; 1, p)$ . Since they are **independent random variables**, we can apply some useful results:

$$\begin{aligned}
\mu_X &= E(I_1) + E(I_2) + \cdots + E(I_n) = np \\
\sigma_X^2 &= \sigma_{I_1}^2 + \sigma_{I_2}^2 + \cdots + \sigma_{I_n}^2 = npq
\end{aligned}$$

□

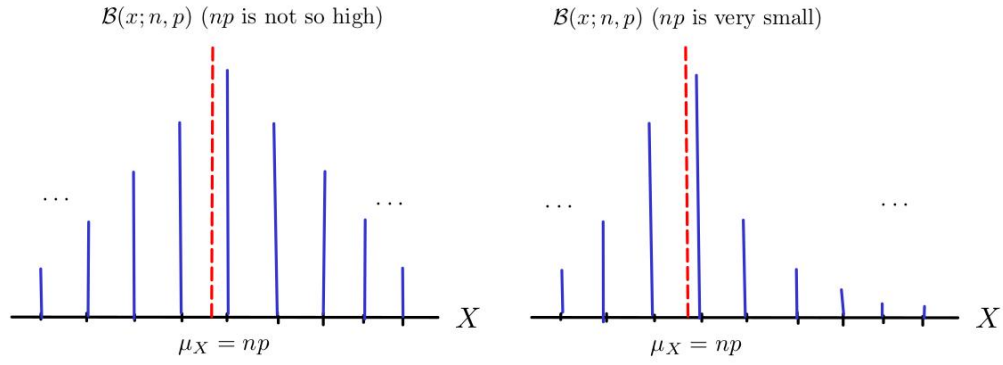


Figure 4.2: The pdf of Binomial distribution

You should note that the position of the **mean value** is very close to the value of the random variable  $X$  at which  $P(X = x)$  reaches its **highest value**. Or in other words,  $\mu_X$  is very close to the **peak** of the pdf graph.

### 4.1.3 Poisson Distribution

#### Using Bell-shaped Curve to approx the Binomial Distribution

Let's begin with a specific example:  $X \sim \mathcal{B}(x; 100, 0.6)$ . Now we can see that the mean value  $\mu_X = np = 60$  is relatively close to the center  $X = 50$ .

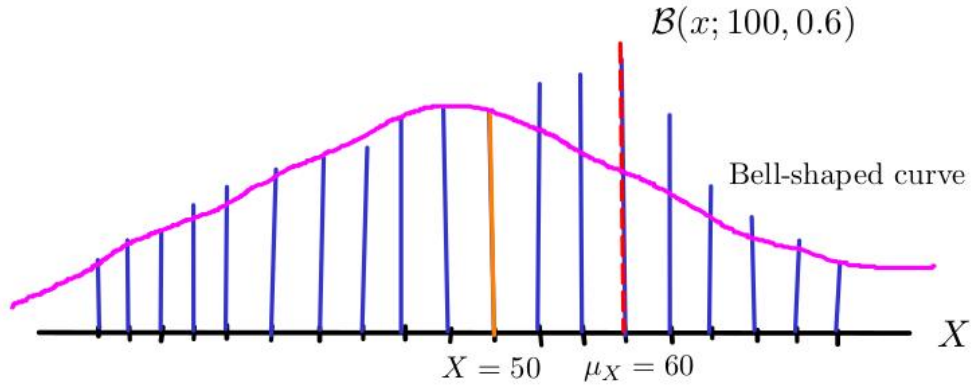


Figure 4.3:  $\mathcal{B}(x; 100, 0.6)$  graph

Because  $np = 60$  value is medium, neither too high nor low compared to the center value  $X = 50$ ; and plotting  $n = 100$  points is high enough for us to connect all of them to obtain a relatively smooth curve. This smooth curve **can be approximated** by the **bell-shaped curve**, so we will perform operations on it instead of the original curve.

In this subsection, we will focus on mathematical concepts rather than delving into calculation methods (you can try it yourself if you want). We will return to specific calculation methods in the next chapter.

**Theorem 4.1.3.** If  $X \sim \mathcal{B}(x; n, p)$  with  $np$  value is medium compared to the center value, then:

$$\text{Original Curve} \approx \text{Bell-shaped Curve: } \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

where:  $\mu = np$  and  $\sigma^2 = npq$

Now let's compare 2 methods: calculating directly and using approximation.

$$P(X \leq 70) = 1 - P(X \geq 71) = 1 - \sum_{x=71}^{100} \binom{100}{x} 0.6^x 0.4^{100-x} = 0.985$$

$$P(X \leq 70) = P\left(Z < \frac{70 - 100 \cdot 0.6 + 0.5}{\sqrt{100 \cdot 0.6 \cdot 0.4}}\right) = P(Z < 2.143) = 0.983$$

Or we can also try:

$$\begin{aligned} P(50 \leq X \leq 65) &= P(X \leq 65) - P(X \leq 49) \\ &= P\left(Z < \frac{65 - 100 \cdot 0.6 + 0.5}{\sqrt{100 \cdot 0.6 \cdot 0.4}}\right) - P\left(Z < \frac{49 - 100 \cdot 0.6 + 0.5}{\sqrt{100 \cdot 0.6 \cdot 0.4}}\right) \\ &= P(Z < 1.122) - P(Z < -2.143) \\ &= P(Z < 1.122) - 1 + P(Z < 2.143) \\ &= 0.868 - 1 + 0.983 = 0.851 \end{aligned}$$

$$P(50 \leq X \leq 65) = \sum_{x=50}^{65} \binom{100}{x} 0.6^x 0.4^{100-x} = 0.852$$

The approximation method works very well and yields results very close to calculating directly. But how about the  $X \sim \mathcal{B}(x; 100, 0.05)$  case? Can we still use the **bell-shaped curve**? Let's check it!

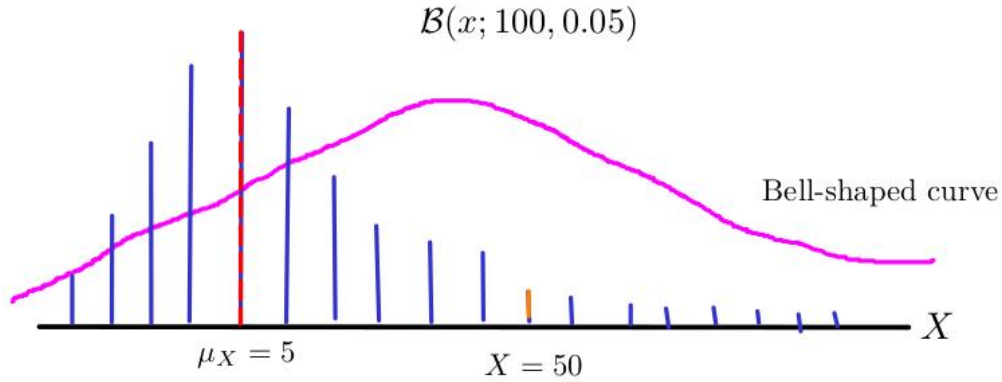


Figure 4.4:  $\mathcal{B}(x; 100, 0.05)$  graph

For example:

$$P(X \leq 7) = \sum_{x=0}^7 \binom{100}{x} 0.05^x 0.95^{100-x} = 0.872$$

$$P(X \leq 7) = P\left(Z < \frac{7 - 100 \cdot 0.05 + 0.5}{\sqrt{100 \cdot 0.05 \cdot 0.95}}\right) = 0.9292$$

The approximate result differs from the actual result nearly 6%. Therefore, for cases where  $\mu_X = np$  is located very far from the center, the bell-shaped curve approximation method is no longer suitable.



## Poisson Distribution

How do we handle extreme value cases of  $np$ ? We must find better way to approximate original curve. Because  $n \gg p$ , so we can perform some equivalent transformations with assumptions  $n \rightarrow +\infty$  and  $p \rightarrow 0$ . We begin with  $X \sim \mathcal{B}(x; n, p)$ :

$$\begin{aligned} P(X = x) &= \mathcal{B}(x; n, p) = \binom{n}{x} p^x q^{n-x} = \frac{n!}{(n-x)!x!} p^x q^{n-x} \\ &= \frac{n(n-1)(n-2) \cdots (n-x+1)}{n^x} \frac{(np)^x}{x!} (1-p)^{n-x} \end{aligned}$$

Because  $n \rightarrow +\infty$  and  $p \rightarrow 0$ :

$$\lim_{n \rightarrow +\infty} \frac{n(n-1)(n-2) \cdots (n-x+1)}{n^x} = 1$$

Our approximation curve will be more accurate if we only care about **the first end** values, it means  $n \gg x$ :

$$\lim_{n \rightarrow +\infty} (1-p)^{n-x} = (1-p)^n$$

By the definition of the constant  $e$ :

$$\lim_{p \rightarrow 0} (1-p)^{\frac{-1}{p}} = e \Rightarrow \lim_{p \rightarrow 0} (1-p) = e^{-p} \Rightarrow (1-p) \approx e^{-p}$$

Now we obtain:

$$(1-p)^n \approx e^{-np}$$

Finally we have:

$$P(X = x) = \frac{(np)^x}{x!} e^{-np}$$

Or in shortened form with  $\mu = np$ :

$$P(X = x) = \frac{\mu^x e^{-\mu}}{x!}$$

Test the new function with the case  $X \sim \mathcal{B}(x; 100, 0.05)$ :

$$P(X \leq 7) = \sum_{x=0}^7 \frac{(100 \cdot 0.05)^x e^{-100 \cdot 0.05}}{x!} = 0.866$$

Obviously, the approximated result is very close to the actual result 0.872.

How about the case  $X \sim \mathcal{B}(x; 100, 0.9)$ ? Since  $p$  is clearly close to 1, not zero, so we must **inverse the role** of  $p$  and  $q$ . We define the new random variable  $Y$  as the **number of failures** to satisfy the condition of valid approximation.

$$\begin{aligned} P(X \geq 90) &= \sum_{x=90}^{100} \binom{100}{x} 0.9^x 0.1^{100-x} = 0.583 \\ P(X \geq 90) &= P(Y \leq 10) = \sum_{y=0}^{10} \frac{(100 \cdot 0.1)^y e^{-100 \cdot 0.1}}{y!} = 0.583 \end{aligned}$$

So relatively speaking, we conclude the following, if  $\mu = np$  position is very **far** from the **left side** of the center, then the Binomial distribution can be approximated by the formula:

$$P(X = x) = \mathcal{B}(x; n, p) \approx \frac{\mu^x e^{-\mu}}{x!}$$

Because of two main constraints  $n \gg p$  and  $p$  is very close to 0, the probability of a success event is relatively **low**. You can think about  $p$  as the **average number** of outcomes per unit time, distance or volume; and  $n$  as a given **time interval** or **specified region**. In practice, to avoid being confused with Binomial distribution and emphasize the fact that  $p \approx 0$ , we usually change our notations to:

$$p \leftrightarrow \lambda$$

$$n \leftrightarrow t$$

**Definition 4.1.3.** The **Poisson process** is the **Bernoulli process** where  $n \gg p$  and  $p \approx 0$ .

**Definition 4.1.4.** A **Poisson event** is an event with a **low** probability of occurring.

**Theorem 4.1.4.** The pdf of the **Poisson random variable**  $X$ , representing the number of outcomes occurring in a **given time interval** or **specified region** denoted by  $t$ , is:

$$\mathcal{P}(x; \mu) = \frac{\mu^x e^{-\mu}}{x!} \quad (x = 0, 1, 2, \dots)$$

where  $\mu = \lambda t$ ,  $\lambda$  is the **average** number of outcomes per unit time, distance area or volume.

**Theorem 4.1.5.** If  $X \sim \mathcal{B}(x; n, p)$ , when  $n \rightarrow +\infty$ ,  $p \rightarrow 0$  and  $np \rightarrow \mu$  remains constant:

$$\mathcal{B}(x; n, p) \rightarrow \mathcal{P}(x; \mu)$$

**Corollary 4.1.5.1.** If  $X \sim \mathcal{P}(x; \mu)$ , then:

$$\mu_X = \mu$$

$$\sigma_X^2 = \mu$$

*Proof.* For the mean:

$$\mu_X = \sum_x x f(x) = \sum_{x=0}^n x \frac{\mu^x e^{-\mu}}{x!} = \sum_{x=1}^n \frac{\mu^x e^{-\mu}}{(x-1)!} = \mu e^{-\mu} \sum_{x=1}^n \frac{\mu^{x-1}}{(x-1)!} = \mu e^{-\mu} e^{\mu} = \mu$$

For the variance, perform the same trick:

$$\begin{aligned} E(X(X-1)) &= \sum_x x(x-1) f(x) = \sum_{x=1}^n x(x-1) \frac{\mu^x e^{-\mu}}{x!} \\ &= \mu^2 e^{-\mu} \sum_{x=2}^n \frac{\mu^{x-2}}{(x-2)!} = \mu^2 e^{-\mu} e^{\mu} \\ &= \mu^2 \end{aligned}$$

Apply directly our previous result:

$$\sigma_X^2 = E(X^2) - \mu_X^2 = E(X(X-1)) + E(X) - \mu_X^2 = \mu^2 + \mu - \mu^2 = \mu$$

□

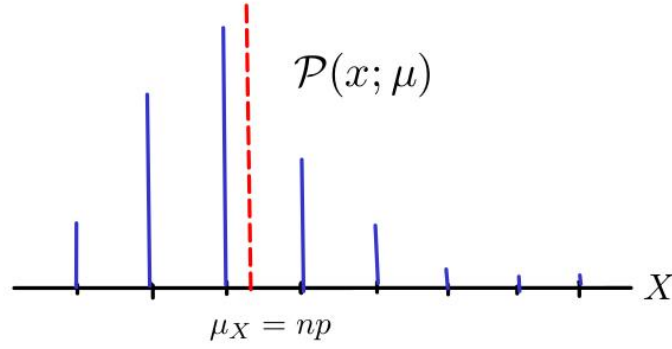


Figure 4.5: The pdf of Poisson distribution

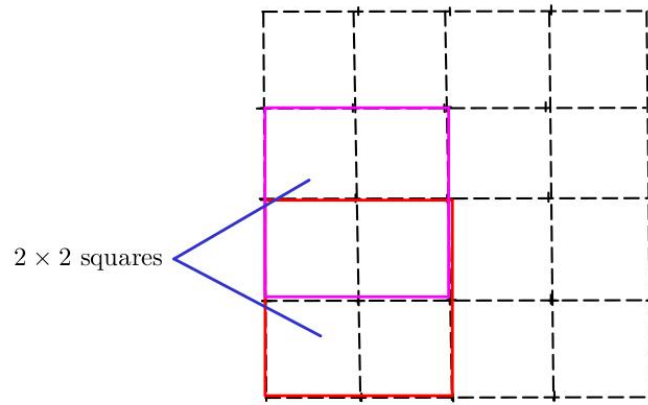


Figure 4.6: A  $4 \times 4$  grid

Back to our mung bean seeds, now you will see how Poisson distribution can be applied in many situations in real life. If you take some seeds from the packet and throw them on a  $4 \times 4$  grid and you observe that there are 10 seeds in the  $2 \times 2$  square closet to you, what are the probabilities of:

1. There are a total of 40 seeds on that  $4 \times 4$  grid.
2. Suppose you knew the exact number of seeds that you had thrown is 40. Evaluate the chance of observing 10 seeds in the  $2 \times 2$  square.

The random variable  $X$  is defined as the number of seeds on a  $4 \times 4$  grid. Using the formula above, we have:

$$\begin{aligned}\mu &= \lambda t = \frac{10}{4} \cdot 16 = 40 \\ \Rightarrow P(X = 40) &= \frac{e^{-\mu} \mu^x}{x!} = \frac{e^{-40} 40^{40}}{40!} = 0.0629\end{aligned}$$

The random variable  $Y$  is defined as the number of seeds in the  $2 \times 2$  square. Similarly, we have:

$$\begin{aligned}\mu &= \lambda t = \frac{40}{16} \cdot 4 = 10 \\ \Rightarrow P(Y = 10) &= \frac{e^{-10} 10^{10}}{10!} = 0.125\end{aligned}$$

So the chance of observing 10 seeds in the  $2 \times 2$  square if you throw 40 seeds on a  $4 \times 4$  grid is pretty low, just 12.5%.

## 4.2 Negative Binomial and Geometric Distributions

### 4.2.1 Negative Binomial Distribution

What is the probability of tossing a coin  $x$  times **until** getting  $k$  heads, if we know that the probability of heads appearing is  $p$  and  $X$  is a random variable, defined as the number of tosses? This question is very easy so I will write down the answer:

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

The function above is a **Negative Binomial Distribution** formula. Formally, we state that:

**Theorem 4.2.1.** *If repeated **Bernoulli trials** can result in a success with probability  $p$  and a failure with probability  $q = 1 - p$ , then the pdf of random variable  $X$ , defined as the **number of the trial** on which the  **$k$ th success occurs** is:*

$$\mathcal{B}^*(x; k, p) = \binom{x-1}{k-1} p^k q^{x-k} \quad (x = k, k+1, k+2, \dots)$$

### 4.2.2 Geometric Distribution

If we only care about the probability of tossing a coin  $x$  times until getting **first** heads, it means  $k = 1$  now and we can rewrite our formula as:

$$P(X = x) = pq^{x-1}$$

The  $k = 1$  case of **Negative Binomial Distribution** can also be called as **Geometric Distribution**.

**Theorem 4.2.2.** *If repeated **Bernoulli trials** can result in a success with probability  $p$  and a failure with probability  $q = 1 - p$ , then the pdf of random variable  $X$ , defined as the **number of trials** on which the **first success occurs** is:*

$$\mathcal{G}^*(x; p) = pq^{x-1} \quad (x = 1, 2, 3, \dots)$$

**Corollary 4.2.2.1.** *If  $X \sim \mathcal{G}^*(x; p)$ , then:*

$$\mu_X = \frac{1}{p}$$
$$\sigma_X^2 = \frac{q}{p^2}$$

*Proof.* For the mean:

$$\mu_X = \sum_x xf(x) = p \sum_{x=1}^{+\infty} xq^{x-1} = p \sum_{x=1}^{+\infty} x(1-p)^{x-1}$$

Now we use the Geometric series and take the derivative once with respect to  $p$ :

$$\begin{aligned} \sum_{x=0}^{+\infty} (1-p)^x &= \frac{1}{1-(1-p)} = \frac{1}{p} \\ \Leftrightarrow \left( \sum_{x=0}^{+\infty} (1-p)^x \right)' &= \left( \frac{1}{p} \right)' \\ \Leftrightarrow - \sum_{x=0}^{+\infty} x(1-p)^{x-1} &= \frac{-1}{p^2} \Rightarrow \sum_{x=1}^{+\infty} x(1-p)^{x-1} = \frac{1}{p^2} \end{aligned}$$

After substituting our result to the previous equation, we obtain:

$$\mu_X = p \sum_{x=1}^{+\infty} x(1-p)^{x-1} = p \cdot \frac{1}{p^2} = \frac{1}{p}$$

Take the second derivative with respect to  $p$  from the previous equation:

$$\left( \sum_{x=1}^{+\infty} x(1-p)^{x-1} \right)' = \left( \frac{1}{p^2} \right)' \Leftrightarrow \sum_{x=2}^{+\infty} x(x-1)(1-p)^{x-2} = \frac{2}{p^3}$$

We want to determine:

$$E(X(X-1)) = \sum_x x(x-1)f(x) = p \sum_{x=1}^{+\infty} x(x-1)(1-p)^{x-1} = \frac{2p(1-p)}{p^3} = \frac{2q}{p^2}$$

For the variance:

$$\sigma_X^2 = E(X^2) - \mu_X^2 = E(X(X-1)) + E(X) - \frac{1}{p^2} = \frac{2q}{p^2} + \frac{1}{p} - \frac{1}{p^2} = \frac{1-p}{p^2} = \frac{q}{p^2}$$

□

### 4.3 Hypergeometric Distribution

Suppose you now have  $N$  mung bean seeds in the packet, of which you know  $N$  of them are spoiled. If you randomly select a handful of  $k$  seeds, what is the probability of getting  $x$  spoiled seeds?

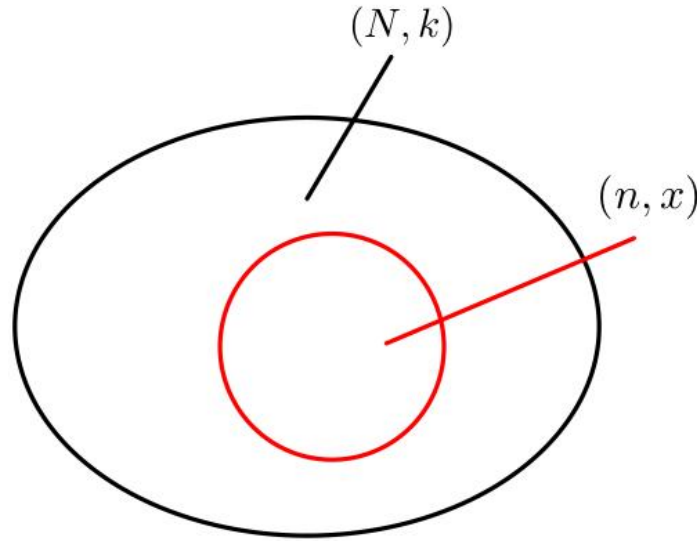


Figure 4.7: Illustration of Hypergeometric distribution

Intuitively, we can see that:

$$P(X = x) = \frac{\binom{k}{x} \binom{N-k}{n-x}}{\binom{N}{n}}$$

**Theorem 4.3.1.** *The pdf of the **hypergeometric** random variable  $X$ , the number of successes in a random sample of size  $n$  selected from  $N$  items of which  $k$  are labeled **success** and  $N - k$  labeled **failure** is:*

$$\mathcal{H}(x; N, k, n) = \frac{\binom{k}{x} \binom{N-k}{n-x}}{\binom{N}{n}} \quad (x = 0, 1, 2, \dots, k)$$

# Chapter 5

## Some Continuous Probability Distributions

Unlike discrete quantities, which are only formed through human counting, continuous quantities appear naturally. You can easily see that most of physical quantities in reality are continuous, such as mass, length, time,  $\dots$  Therefore, this chapter plays crucial foundational role for the rest of this book.

### 5.1 Uniform Distribution

As we discussed above, the pdf of choosing randomly a number within the range  $(0, 1)$  experiment is:

$$f(x) = \begin{cases} 1 & (0 < x < 1) \\ 0 & (\text{elsewhere}) \end{cases}$$

The probability of our selected number within  $(0.3, 0.7)$  interval is:

$$P(0.3 < x < 0.7) = \int_{-\infty}^{+\infty} f(x)dx = \int_{0.3}^{0.7} dx = 0.4$$

Suppose normal seeds will typically germinate within 3 – 5 days of being placed on a wet towel. They will never germinate before 3 days or after 5 days. Now we form the pdf of this experiment:

$$f(x) = \begin{cases} \frac{1}{2} & (3 < x < 5) \\ 0 & (\text{elsewhere}) \end{cases}$$

We can obtain some useful result such as:

$$P(X < x) = F(x) = \int_{-\infty}^x f(t)dt = \begin{cases} 0 & (x \leq 3) \\ \frac{x}{2} & (3 < x < 5) \\ 1 & (x \geq 5) \end{cases}$$

But why do  $f(x) = \frac{1}{2}$  with  $x \in (3, 5)$ ? Well because every pdf must satisfy these criteria:

$$\begin{cases} f(x) \geq 0 \text{ (for all } x \in \mathbb{R}) \\ \int_{-\infty}^{+\infty} f(x)dx = 1 \\ P(a < x < b) = \int_a^b f(x)dx \end{cases}$$

Verifying yourself that the pdf of germinating seeds experiment satisfies all of them.

**Definition 5.1.1.** The pdf of the continuous **uniform random variable**  $X$  on the interval  $(A, B)$  is:

$$\mathcal{U}(x; a, b) = \begin{cases} \frac{1}{b-a} & (a < x < b) \\ 0 & (\text{elsewhere}) \end{cases}$$

**Corollary 5.1.0.1.** If  $X \sim \mathcal{U}(x; a, b)$  then:

$$\mu_X = \frac{b+a}{2}$$

$$\sigma_X^2 = \frac{(b-a)^2}{12}$$

*Proof.* For the mean:

$$\mu_X = \int_{-\infty}^{+\infty} x f(x) dx = \frac{1}{b-a} \int_a^b x dx = \frac{1}{b-a} \cdot \frac{b^2 - a^2}{2} = \frac{b+a}{2}$$

For the variance:

$$\begin{aligned} \sigma_X^2 &= E(X^2) - \mu_X^2 = \int_{-\infty}^{+\infty} x^2 f(x) dx - \frac{(b+a)^2}{4} = \frac{1}{b-a} \int_a^b x^2 dx - \frac{(b+a)^2}{4} \\ &= \frac{1}{b-a} \frac{b^3 - a^3}{3} - \frac{(b+a)^2}{4} = \frac{(b-a)^2}{12} \end{aligned}$$

□

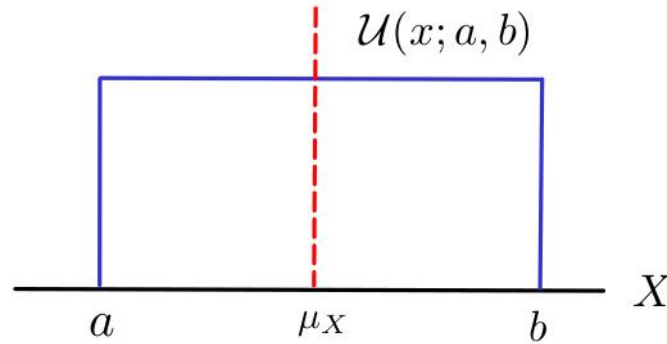


Figure 5.1: The pdf of Uniform distribution

## 5.2 Normal Distribution

### 5.2.1 The Idea behind the Normal Distribution

This section is one of the most important part of this book. Now we shall 'watch' very closely to the question "Is everything normal because they are?" and partially answer my question from the **Preface**. Let's temporarily forget everything we were taught about those concepts and now rebuild them from scratch.

You are watching the rain outside. You observe that the raindrops are all falling **randomly** onto a circular courtyard. You might wonder where a raindrop is most likely to fall? Perhaps at the **center** of the circle? Or at the **edge**? Or neither? Intuitively, you might think the raindrop would concentrate at the center, but are you sure about that when we can not blindly trust intuition? Remember that ProbStat is a paradoxical subject. It tricks our brains with the fact **that not everything can be absolutely certain, even this statement**; while human brain is always naturally seeking for the certainty. The only thing we can trust is the use of mathematical language to describe the uncertainty.

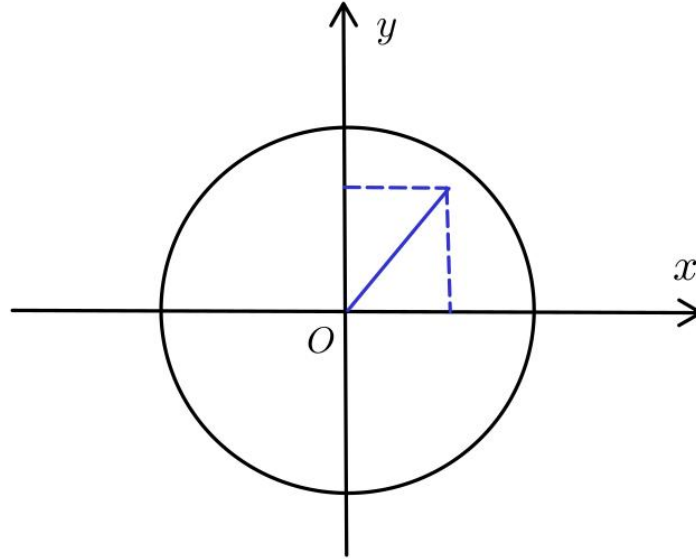


Figure 5.2: A circular courtyard with raindrops falling inside

You define the pdf of an event "A single raindrop falls at the position with coordinate  $(x, y)$ " as  $f(x, y)$ . Since raindrops fall completely randomly and you surely know:

$$f(x|y) = f(x)$$

$$f(y|x) = f(y)$$

Since the events "A raindrop falls on a point with coordinate  $x$ " and "A raindrop falls on a point with coordinate  $y$ " are **independent**, we can conclude:

$$f(x, y) = f(x)f(y)$$

Applying Pythagorean theorem shows that:

$$f(x, y) = f(\sqrt{x^2 + y^2})$$

Substituting to the previous formula yields:

$$f(\sqrt{x^2 + y^2}) = f(x)f(y)$$

Now we have to find the solution of this functional equation. Firstly, we define a new auxiliary function  $g(x)$  as follows:

$$f(x) = g(x^2)$$

Our functional equation can be rewritten as:

$$g(x^2 + y^2) = g(x^2)g(y^2)$$



Setting  $y^2 = x^2$  and  $t = x^2$  yields:

$$g(2t) = g^2(t)$$

We can observe the pattern:

$$g(3t) = g(2t + t) = g(2t)g(t) = g^3(t)$$

Now we want to prove this result by using induction:

$$g(nt) = g^n(t) \quad (n \in \mathbb{N}^*)$$

With special cases  $n = 1, 2, 3$ , our equation is correct, now we assume it can be still correct until  $n = k$  with  $k > 3$ :

$$g(kt) = g^k(t) \quad (k > 3)$$

For the case  $n = k + 1$ , our equation is still correct:

$$g((k + 1)t) = g(kt)g(t) = g^{k+1}(t)$$

By induction we can conclude:

$$g(nt) = g^n(t) \quad (n \in \mathbb{N}^*)$$

Setting  $t = 1$  and  $g(1)$  as a constant  $C$  yields:

$$g(n) = g^n(1) = C^n \quad (n \in \mathbb{N}^*)$$

Constant  $C$  can also be written in an exponential form:

$$g(n) = C^n = (e^c)^n \quad (n \in \mathbb{N}^*)$$

Now we guess the form of  $g(x)$  might be:

$$g(x) = e^{cx} \quad (\forall x \in \mathbb{R})$$

Or equivalent to:

$$f(x) = g(x^2) = e^{cx^2} \quad (\forall x \in \mathbb{R})$$

Test our assumption with an arbitrary constant  $c$ :

$$f(\sqrt{x^2 + y^2}) = e^{c(x^2 + y^2)} = e^{cx^2} e^{cy^2} = f(x)f(y)$$

So the solution of our functional equation is:

$$f(x) = e^{cx^2} \quad (\forall x, c \in \mathbb{R})$$

Due to **Chebyshev's Inequality**, both ends of a valid pdf must be **flattened**, so we choose our solution is:

$$f(x) = e^{-cx^2} \quad (c > 0)$$

The final step is choosing an appropriate constant  $c$  to satisfy the criteria of pdf:

$$\int_{-\infty}^{+\infty} f(x)dx = 1$$

To avoid the constant  $c$  ( $c$  is just a scale factor), we consider the special case where  $c = 1$  and calculate its integral value:

$$I = \int_{-\infty}^{+\infty} e^{-x^2} dx$$

Since  $e^{-x^2}$  is not an elementary function, so its antiderivative can not be found. Again, we must perform a brilliant trick here:

$$I^2 = \int_{-\infty}^{+\infty} e^{-x^2} dx \int_{-\infty}^{+\infty} e^{-y^2} dy = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} e^{-(x^2+y^2)} dx dy = \int_0^{2\pi} \int_0^{+\infty} e^{-r^2} r dr d\theta = \pi$$

I do not go deeply into the changing variables (from Cartesian to Polar coordinates) and using Jacobian transformation steps since they are very easy; but you can see here, the idea of squaring  $I$  is so amazing! Now we obtain:

$$I = \int_{-\infty}^{+\infty} e^{-x^2} dx = \sqrt{\pi}$$

Back to our pdf, now we can change the variable and perform integration:

$$\int_{-\infty}^{+\infty} e^{-cx^2} dx = \frac{1}{\sqrt{c}} \int_{-\infty}^{+\infty} e^{-(\sqrt{c}x)^2} d(\sqrt{c}x) = \sqrt{\frac{\pi}{c}}$$

After choosing  $c = \pi$ , finally our pdf is:

$$f(x) = e^{-\pi x^2}$$

Sketching the pdf, and as you can see here, the chance of raindrops falling around the **center** is **higher** than the **edge** or anywhere esle.

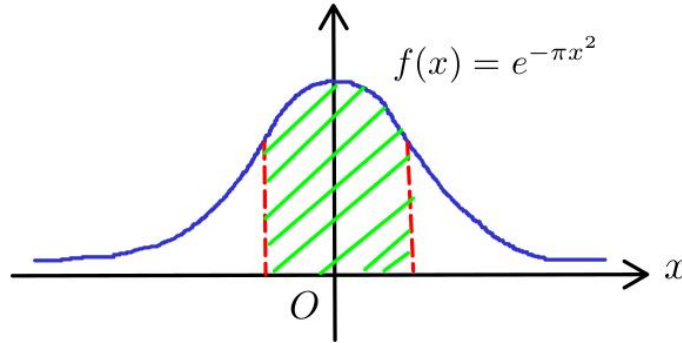


Figure 5.3: The probability of raindrops falling around the center is higher than anywhere else

### 5.2.2 Standard Normal Distribution

For convenience, we usually choose  $c = \frac{1}{2}$  to be standard scale factor and define:

**Definition 5.2.1.** A **standard normal distribution function** is defined as:

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

In the case of **standard normal distribution**, we often denote  $Z$  as our random variable. This is a common convention, and you will understand why later.

**Definition 5.2.2.** If  $Z \sim \mathcal{N}(z; 0, 1)$ , then:

$$Z \sim f(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right)$$

**Corollary 5.2.0.1.** If  $Z \sim \mathcal{N}(z; 0, 1)$ , then:

$$\int_{-\infty}^{+\infty} f(z) dz = 1$$

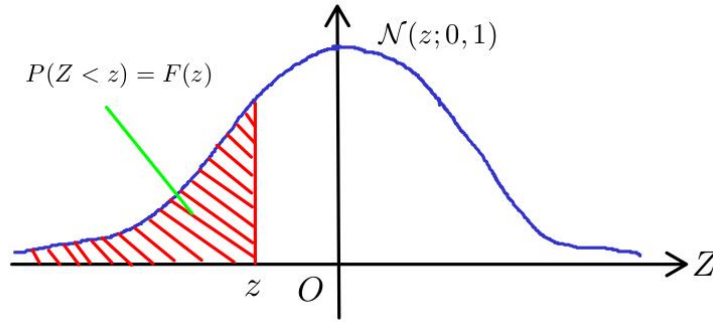


Figure 5.4: The pdf of Standard Normal distribution

We are very interested in determining cdf  $F(z)$  is defined as:

$$F(z) = P(Z < z) = \int_{-\infty}^z f(t) dt$$

But because  $f(z)$  is **not an elementary function**, so the antiderivative  $F(z)$  can not be found exactly. An alternative way is using Taylor-Maclaurin expansion to obtain the **approximate** version of  $F(z)$ :

$$\begin{aligned} F(z) = P(Z < z) &= \int_{-\infty}^z f(t) dt = \int_{-\infty}^0 f(t) dt + \int_0^z f(t) dt = \frac{1}{2} + \frac{1}{\sqrt{2\pi}} \int_0^z \exp\left(-\frac{t^2}{2}\right) dt \\ &= \frac{1}{2} + \frac{1}{\sqrt{2\pi}} \int_0^z \sum_{k=0}^{+\infty} \frac{(-1)^k t^{2k}}{2^k k!} dt = \frac{1}{2} + \frac{1}{\sqrt{2\pi}} \sum_{k=0}^{+\infty} \frac{(-1)^k z^{2k+1}}{(2k+1) 2^k k!} \end{aligned}$$

In practice, we usually do not use the approximation of  $F(z)$ ; instead, we always look for our desired value in the  $Z$  score table. In this book,  $Z$  score table is the **Appendix A**.

Let's try comparing some values of  $F(z)$  when using the table and using the approximation function:

$$F(Z < 0.76) \approx \frac{1}{2} + \frac{1}{\sqrt{2\pi}} \sum_{k=0}^{10} \frac{(-1)^k (0.76)^{2k+1}}{(2k+1) 2^k k!} = 0.7763$$

$$F(Z < 1.65) \approx \frac{1}{2} + \frac{1}{\sqrt{2\pi}} \sum_{k=0}^{10} \frac{(-1)^k (1.65)^{2k+1}}{(2k+1)2^k k!} = 0.9505$$

Our results are perfectly matched with results in the  $Z$  score table. There are many ways to approximate  $F(z)$ , and using Taylor-Maclaurin series is the simplest one; although it is accurate and acceptable, but the complexity is **not optimized**. The topic of optimizing approximation of  $F(z)$  is beyond the scope of this book, and as I said in practice we always look for our desired value in the  $Z$  score table; but I hope now we have a better understanding of the mystery behind how the  $Z$  score table is formed, rather than just simply looking up at it.

### 5.2.3 Normal Distribution

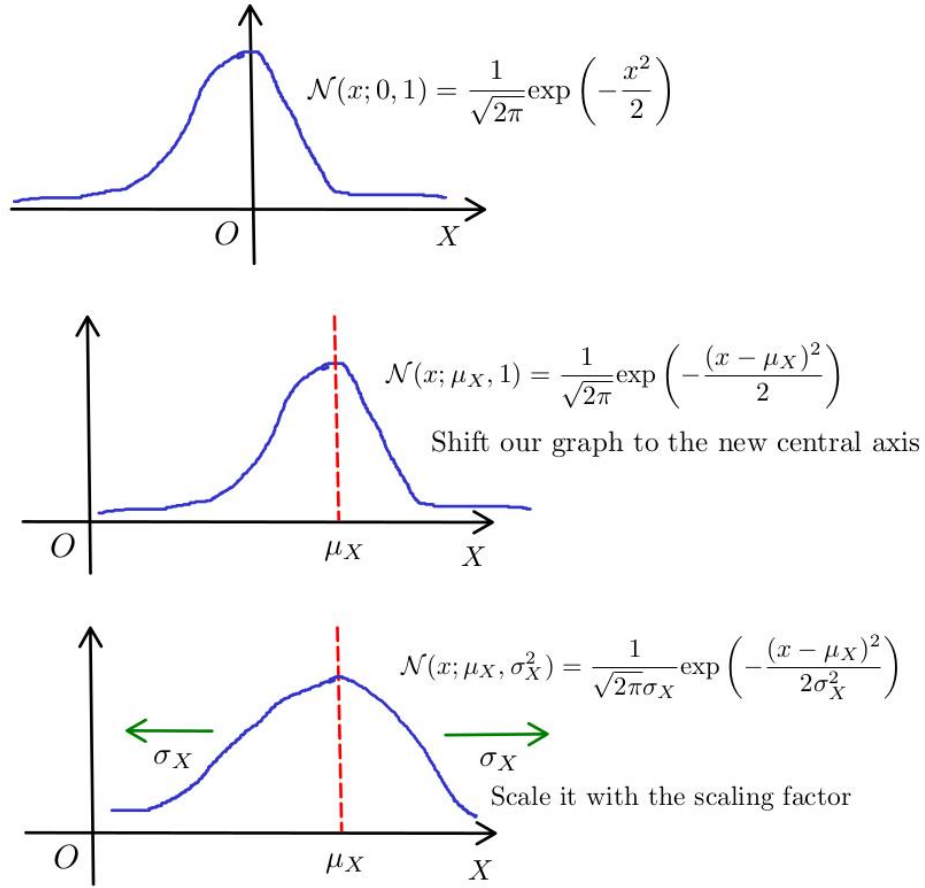


Figure 5.5: Visualizing how the general normal distribution formula is formed

As you can see the general normal distribution function can be formed in two steps: shifting the standard normal distribution curve and scaling it. Now the roles of the **mean**  $\mu_X$  and the **standard deviation**  $\sigma_X$  are very clear. The **mean** acts as the **axis of symmetry**, while the **standard deviation** acts as the scale factor of the general normal distribution graph. Formally, we define:

**Definition 5.2.3.** The pdf of the **normal random variable**  $X$ , with **mean**  $\mu_X$  and **variance**  $\sigma_X^2$  is:

$$\mathcal{N}(x; \mu_X, \sigma_X^2) = \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x - \mu_X)^2}{2\sigma_X^2}\right)$$

**Theorem 5.2.1.** If  $X \sim \mathcal{N}(x; \mu_X, \sigma_X^2)$  and a new random variable  $Z$  (or  $Z$  score) is defined as:

$$Z = \frac{X - \mu_X}{\sigma_X}$$

then:

$$Z \sim \mathcal{N}(z; 0, 1)$$

This is an extremely useful result. We will prove it rigorously later in the next chapter, but you can see the relationship between  $X$  and  $Z$  through graph transformations is very clear and intuitive. These random variables are **linearly mapped**. Now instead of performing directly operations on  $X$ , now we do the smarter choice: transforming from  $X$  to  $Z$  score first, then processing  $Z$  later; since we want to avoid 2 parameters  $\mu_X$  and  $\sigma_X$  in our calculation and we already knew the  $Z$  score table.

For example, if  $X \sim \mathcal{N}(x; 3, 0.5^2)$ , then:

$$\begin{aligned} P(3.5 < X < 4) &= P\left(\frac{3.5 - 3}{0.5} < Z < \frac{4 - 3}{0.5}\right) \\ &= P(1 < Z < 2) \\ &= P(Z < 2) - P(Z < 1) \\ &= 0.9773 - 0.8413 = 0.135 \end{aligned}$$

Finally, this section will conclude with one of the most powerful theorems, which I introduced in the previous chapter.

**Theorem 5.2.2.** Let  $X \sim \mathcal{B}(x; n, p)$ , if  $np$  is medium compared to the center value, then:

$$P(X \leq x) \approx P\left(Z < \frac{x - np + 0.5}{\sqrt{npq}}\right)$$

where  $Z$  is a **standard normal** random variable.

What is the meaning of 0.5? Why do we have to pad 0.5 unit? I think the reason is padding 0.5 to increase the **resolution** from **discrete** to **continuous** version of our approximate line.

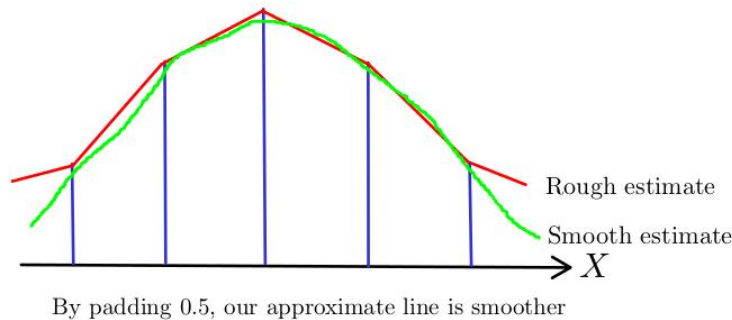


Figure 5.6: Normal approximation to the binomial

Now you should review **Chapter 4: Some Discrete Probability Distributions** and try applying this theorem.

## 5.3 Exponential, Gamma and Chi-Squared Distributions

### 5.3.1 Exponential Distribution

Back to our weather forecasting problem from **Chapter 0: Introduction to Probability and Statistics**; suppose in December, there are average 2 rainy days annually. Applying directly Poisson distribution, you can calculate the probabilities of only 1, or 2, 3, or more rainy days in December.

$$P(X = 1) = \frac{e^{-2}2^1}{1!} = 0.2706$$

$$P(X = 2) = \frac{e^{-2}2^2}{2!} = 0.2706$$

$$P(X = 3) = \frac{e^{-2}2^3}{3!} = 0.1804$$

But our concern is not just the number of rainy days in December, we actually want to know when it will rain (and also plan for the picnic!). For example, if today is the 3rd and you have not seen a drop of rain since the beginning of the month, what is the probability of it raining within the next 4 days? Or what is the probability that the first 10 days of this month will have no rain? Or there will be at least one rainy day? A lot of information can be gathered just from the data "On average, there are 2 rainy days in December."

Firstly, let's find the answer for the easiest question: "What is the probability that the first 10 days of this month will have no rain?". Again, you can just applying directly Poisson distribution without hesitation. The random variable  $Y$  is defined as the number of rainy days within the first 10 days:

$$\begin{aligned}\mu &= \lambda t = \frac{2}{31} \cdot 10 = \frac{20}{31} \\ \Rightarrow P(Y = 0) &= \frac{e^{-\mu}\mu^0}{0!} = e^{-\mu} = 0.5245\end{aligned}$$

Looking from a different perspective, you define the random variable  $X$  as **the waiting time** for the **first** rain:

$$P(Y = 0) = P(X > 10) = e^{-\mu} = 0.5245$$

Now the probability of at least one rainy day in the first 10 days is:

$$P(X < 10) = 1 - P(X > 10) = 1 - 0.5245 = 0.4755$$

Generally, if  $x$  denote the number of days:

$$\begin{aligned}P(Y = 0) &= P(X > x) = e^{-\mu} = e^{-\lambda x} \\ \Rightarrow P(X < x) &= 1 - P(X > x) = 1 - e^{-\lambda x}\end{aligned}$$

From our previous results, now we can derive the pdf of random variable  $X$ , is defined as the waiting time for the **first** rain:

$$F(x) = P(X < x) = 1 - e^{-\lambda x} \Rightarrow f(x) = \frac{dF(x)}{dx} = \lambda e^{-\lambda x}$$

Rainy day is an example of a **Poisson event** because its probability of occuring is **very low**. Since  $f(x)$  is an exponential function, so  $X$  can also be called an **exponential random variable**.

**Definition 5.3.1.** The random variable  $X$  representing the **distance** (in most cases, it is time) between **Poisson events** has pdf:

$$\mathcal{E}(x; \lambda x) = \begin{cases} \lambda e^{-\lambda x} & (x \geq 0) \\ 0 & (\text{elsewhere}) \end{cases}$$

**Corollary 5.3.0.1.** If  $X \sim \mathcal{E}(x; \lambda x)$ , then:

$$\mu_X = \frac{1}{\lambda}$$

$$\sigma_X^2 = \frac{1}{\lambda^2}$$

*Proof.* For the mean:

$$\begin{aligned} \mu_X &= \int_{-\infty}^{+\infty} x f(x) dx = \frac{1}{\lambda} \int_0^{+\infty} (\lambda x) e^{-\lambda x} d(\lambda x) = \frac{1}{\lambda} \int_0^{+\infty} t e^{-t} dt = \frac{-1}{\lambda} \int_0^{+\infty} t d(e^{-t}) \\ &= \frac{-1}{\lambda} \left( t e^{-t} \Big|_0^{+\infty} - \int_0^{+\infty} e^{-t} dt \right) = \frac{1}{\lambda} \end{aligned}$$

For the variance:

$$\begin{aligned} \sigma_X^2 &= E(X^2) - \mu_X^2 = \int_{-\infty}^{+\infty} x^2 f(x) dx - \frac{1}{\lambda^2} = \int_0^{+\infty} x^2 \lambda e^{-\lambda x} dx - \frac{1}{\lambda^2} = \frac{1}{\lambda^2} \int_0^{+\infty} (\lambda x)^2 e^{-\lambda x} d(\lambda x) - \frac{1}{\lambda^2} \\ &= \frac{1}{\lambda^2} \int_0^{+\infty} t^2 e^{-t} dt - \frac{1}{\lambda^2} = \frac{-1}{\lambda^2} \int_0^{+\infty} t^2 d(e^{-t}) - \frac{1}{\lambda^2} = \frac{-1}{\lambda^2} \left( t^2 e^{-t} \Big|_0^{+\infty} - 2 \int_0^{+\infty} t e^{-t} dt \right) - \frac{1}{\lambda^2} \\ &= \frac{2}{\lambda^2} - \frac{1}{\lambda^2} = \frac{1}{\lambda^2} \end{aligned}$$

□

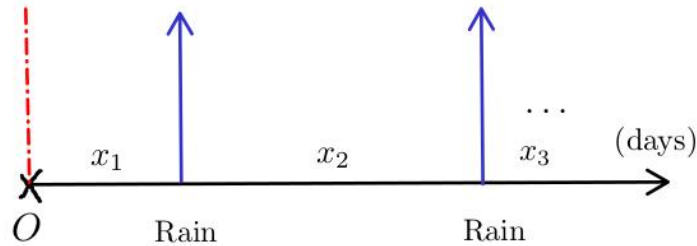


Figure 5.7: Illustration of Exponential distribution

From the illustration, now we can see that from 0:00 on December 1st (denoted  $O$ ), we have to wait  $x_1$  days until the first rain. The random variable  $X_1$  is defined as **the waiting time** for the first rain **since the origin**:

$$X_1 \sim \mathcal{E}(x_1; \lambda x_1)$$

Right after the first rain, we have to wait  $x_2$  days until the second rain. The random variable  $X_2$  is defined as **the waiting time** for the second rain **since the first rain**. Similarly, we see that:

$$X_2 \sim \mathcal{E}(x_2, \lambda x_2)$$

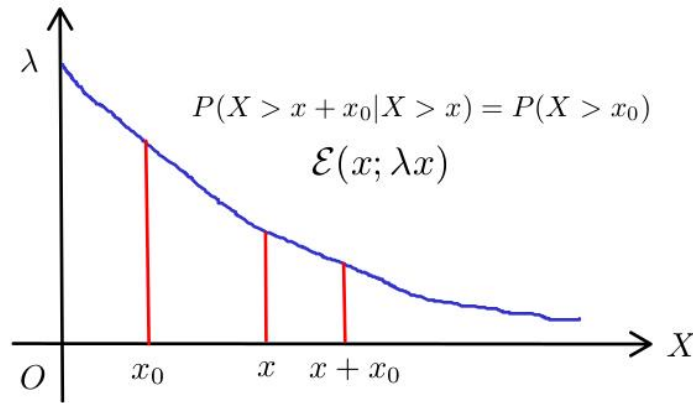
The random variables  $X_1, X_2, X_3, \dots$  are **statistically independent** since they do not have any relationship to the others.

Now to answer our second question: "If the first 3 days are not rainy, what is the probability of it raining within the next 4 days?", you should note that this is a classic conditional probability problem:

$$P(X < 7 | X > 3) = 1 - P(X > 7 | X > 3) = 1 - \frac{P(X > 7)}{P(X > 3)} = 1 - \frac{e^{-7 \cdot \frac{2}{31}}}{e^{-3 \cdot \frac{2}{31}}} = 1 - e^{-4 \cdot \frac{2}{31}} = 0.2274$$

Generally, we can see the **memoryless property of exponential distribution**:

$$P(X > x + x_0 | X > x) = \frac{P(X > x + x_0)}{P(X > x)} = \frac{e^{-\lambda(x+x_0)}}{e^{-\lambda x}} = e^{-\lambda x_0} = P(X > x_0)$$



The observation point  $x$  does not affect the final result

Figure 5.8: The pdf of Exponential distribution and its memoryless property

**Theorem 5.3.1.** *If  $X \sim \mathcal{E}(x; \lambda x)$ , then it has **memoryless property**:*

$$P(X > x + x_0 | X > x) = P(X > x_0)$$

## 5.3.2 Gamma Distribution

### Gamma Function

Before exploring the **Gamma Distribution**, let me introduce you to an odd but very powerful function:

**Definition 5.3.2.** *For every  $\alpha > 0$ , the **gamma function** is defined by:*

$$\Gamma(\alpha) = \int_0^{+\infty} x^{\alpha-1} e^{-x} dx$$

From now on, the **gamma function** will appear almost everywhere in ProbStat, so you should learn it carefully.



**Theorem 5.3.2.** For every  $\alpha > 0$ , the **gamma function** can be determined through a **recursion**:

$$\Gamma(\alpha) = (\alpha - 1)\Gamma(\alpha - 1)$$

*Proof.*

$$\begin{aligned}\Gamma(\alpha) &= \int_0^{+\infty} x^{\alpha-1} e^{-x} dx = - \int_0^{+\infty} x^{\alpha-1} d(e^{-x}) = - \left( x^{\alpha-1} e^{-x} \Big|_0^{+\infty} - (\alpha - 1) \int_0^{+\infty} x^{\alpha-2} e^{-x} dx \right) \\ &= (\alpha - 1) \int_0^{+\infty} x^{\alpha-2} e^{-x} dx = (\alpha - 1)\Gamma(\alpha - 1)\end{aligned}$$

□

**Corollary 5.3.2.1.**

$$\Gamma(1) = 1$$

*Proof.*

$$\Gamma(1) = \int_0^{+\infty} x^0 e^{-x} dx = \int_0^{+\infty} e^{-x} dx = 1$$

□

**Corollary 5.3.2.2.** For any positive integer  $n$ :

$$\Gamma(n) = (n - 1)!$$

*Proof.* For any positive integer  $n$ :

$$\Gamma(n) = (n - 1)\Gamma(n - 1) = (n - 1)(n - 2)\Gamma(n - 2) = (n - 1)(n - 2) \cdots 2\Gamma(1) = (n - 1)!$$

□

**Corollary 5.3.2.3.**

$$\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$$

*Proof.*

$$\Gamma\left(\frac{1}{2}\right) = \int_0^{+\infty} \frac{e^{-x}}{\sqrt{x}} dx$$

Change variable by setting  $t = \sqrt{x} \Rightarrow 2t dt = dx$ :

$$\Gamma\left(\frac{1}{2}\right) = \int_0^{+\infty} \frac{e^{-t^2}}{t} 2t dt = 2 \int_0^{+\infty} e^{-t^2} dt = \sqrt{\pi}$$

□

Using recursion property of Gamma function, we can deduce several useful results:

$$\Gamma\left(\frac{3}{2}\right) = \frac{1}{2}\Gamma\left(\frac{1}{2}\right) = \frac{\sqrt{\pi}}{2}$$

$$\Gamma\left(\frac{5}{2}\right) = \frac{3}{2}\Gamma\left(\frac{3}{2}\right) = \frac{3\sqrt{\pi}}{4}$$

## Gamma Distribution

The pdf of Gamma distribution is a little bit messy, so we focus on understanding how it can be applied rather than deriving from scratch.

**Definition 5.3.3.** *The continuous random variable  $X$  has a **gamma distribution**, with 2 parameters  $\alpha$  and  $\beta$  are both positive, if its pdf is given by:*

$$\mathcal{G}(x; \alpha, \beta) = \begin{cases} \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-x/\beta} & (x > 0) \\ 0 & (\text{elsewhere}) \end{cases}$$

**Corollary 5.3.2.4.** *If  $X \sim \mathcal{G}(x; \alpha, \beta)$ , then:*

$$\begin{aligned} \mu_X &= \alpha\beta \\ \sigma_X^2 &= \alpha\beta^2 \end{aligned}$$

*Proof.* For the mean:

$$\begin{aligned} \mu_X &= \int_{-\infty}^{+\infty} x f(x) dx = \int_0^{+\infty} \frac{1}{\Gamma(\alpha)\beta^\alpha} x^\alpha e^{-x/\beta} dx = \alpha\beta \int_0^{+\infty} \frac{1}{\Gamma(\alpha+1)\beta^{\alpha+1}} x^\alpha e^{-x/\beta} dx \\ &= \alpha\beta \int_0^{+\infty} \mathcal{G}(x; \alpha+1, \beta) dx = \alpha\beta \end{aligned}$$

For the variance:

$$\begin{aligned} \sigma_X^2 &= E(X^2) - \mu_X^2 = \int_{-\infty}^{+\infty} x^2 f(x) dx - \mu_X^2 = \int_0^{+\infty} \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha+1} e^{-x/\beta} dx - (\alpha\beta)^2 \\ &= (\alpha+1)\alpha\beta^2 \int_0^{+\infty} \frac{1}{\Gamma(\alpha+2)\beta^{\alpha+2}} x^{\alpha+1} e^{-x/\beta} dx - (\alpha\beta)^2 \\ &= (\alpha+1)\alpha\beta^2 \int_0^{+\infty} \mathcal{G}(x; \alpha+2, \beta) dx - (\alpha\beta)^2 \\ &= (\alpha+1)\alpha\beta^2 - \alpha^2\beta^2 = \alpha\beta^2 \end{aligned}$$

□

Exponential distribution is just a **special case** of Gamma distribution, you should notice that:

$$\mathcal{E}(x; \lambda) = \mathcal{G}\left(x; 1, \frac{1}{\lambda}\right)$$

So what does it mean? The relationship between them can be described through 2 parameters  $\alpha$  and  $\beta$ , with  $\alpha$  as the **number of Poisson events** that you want to observe in a time interval,  $\beta$  is just the **inverse** of  $\lambda$ . For example, if  $X$  is the random variable, defined as the **waiting time until the second rain** from the **beginning of December** (on average, there are 2 rainy days in December), then:

$$X \sim \mathcal{G}\left(x; 2, \frac{31}{2}\right)$$

The probability of there are **at least** 2 rainy days within the first 10 days in December is:

$$P(X < 10) = \int_0^{10} \frac{1}{\Gamma(2) (31/2)^2} x e^{-2x/31} dx = 0.1369$$

### 5.3.3 Chi-Squared Distribution

Chi-Squared distribution is a special case of Gamma distribution, where  $v$  is a positive integer:

$$\alpha = \frac{v}{2}, \beta = 2$$

You will understand its vital role in **Statistics** in the next chapters.

**Definition 5.3.4.** *The continuous random variable  $X$  has a **chi-squared distribution**, with  $v$  **degrees of freedom**, if its pdf is given by:*

$$\mathcal{C}(x; v) = \begin{cases} \frac{1}{2^{v/2}\Gamma(v/2)} x^{v/2-1} e^{-x/2} & (x > 0) \\ 0 & (\text{elsewhere}) \end{cases}$$

where  $v$  is a positive integer.

**Corollary 5.3.2.5.** *If  $X \sim \mathcal{C}(x; v)$ , then:*

$$\begin{aligned} \mu_X &= v \\ \sigma_X^2 &= 2v \end{aligned}$$

*Proof.* For the mean:

$$\mu_X = \alpha\beta = v$$

For the variance:

$$\sigma_X^2 = \alpha\beta^2 = 2v$$

□

# Chapter 6

## Functions of Random Variables

From **Chapter 3: Mathematical Expectation**, we discovered many cool properties of **mean** and **variance** of random variables. For example, given 2 **independent** random variables  $X_1 \sim \mathcal{P}(x_1; \mu_1)$  and  $X_2 \sim \mathcal{P}(x_2; \mu_2)$ , if we define new random variable  $Y = X_1 + X_2$ , we can confidently conclude:

$$\mu_Y = \mu_{X_1} + \mu_{X_2} = \mu_1 + \mu_2$$

$$\sigma_Y^2 = \sigma_{X_1}^2 + \sigma_{X_2}^2 = \mu_1^2 + \mu_2^2$$

But we have not answered the question "How is the random variable  $Y$  distributed?". Is  $Y$  distributed **normally**, or **exponentially**, or neither? The only information we have known so far is just **expected value** and **variance** of  $Y$ . In this chapter, we will find an answer for our question; and I think it should not be skipped (like many other textbooks) just because there is too much Calculus. All of results derived here plays vital roles in **Statistics**.

### 6.1 Transformations of Variables

#### 6.1.1 Linear Transformations

Suppose that the discrete random variable  $X \sim f(x)$ , and we want to find the pdf of  $Y \sim g(y)$  if we knew both  $X$  and  $Y$  have a linear relationship:  $Y = u(X)$ . How can we do that? Firstly, we find the **inverse function** of  $u(X)$ :

$$X = w(Y)$$

Now from the definition of pdf, we have:

$$g(y) = P(Y = y) = P(X = w(y)) = f(w(y))$$

**Theorem 6.1.1.** Suppose that  $X$  is a **discrete** random variable with pdf  $f(x)$ . Let  $Y = u(X)$  define a **linear transformation** so that its **inverse function**  $X = w(Y)$  can be found; then the pdf of  $Y$  is:

$$g(y) = f(w(y))$$

If  $X$  is a **continuous** random variable with pdf  $f(x)$ , then:

$$g(y) = f(w(y))|w'(y)|$$

A very common mistake is forgetting the  $|w'(y)|$  term. To avoid it, you can treat  $|w'(y)|$  role like  $dx$  as the differential part of integration and remember it is always **positive**.

For example, if  $X \sim \mathcal{B}(x; 10, 0.5)$  and  $Y = u(X) = X + 3$ , then the inverse function of  $u(X)$  is:

$$X = w(Y) = Y - 3$$

Finally, we can conclude:

$$g(y) = f(w(y)) = \mathcal{B}(y - 3; 10, 0.5) = \binom{10}{y - 3} 0.5^{y-3} 0.5^{13-y} \quad (y = 3, 4, 5, \dots, 13)$$

Here is a harder example, if  $X \sim \mathcal{U}(x; 0, 1)$  and  $Y = u(X) = -2 \ln X$ , then the inverse function of  $u(X)$  is:

$$X = w(Y) = e^{-Y/2}$$

Applying directly formula yields:

$$g(y) = f(w(y))|w'(y)| = \mathcal{U}(e^{-y/2}; 0, 1) \left| \frac{-1}{2} e^{-y/2} \right| = \frac{1}{2} e^{-y/2} \mathcal{U}(e^{-y/2}; 0, 1)$$

From the definition of **uniform distribution**:

$$\mathcal{U}(e^{-y/2}; 0, 1) = \begin{cases} 1 & (0 < e^{-y/2} < 1) \\ 0 & (\text{elsewhere}) \end{cases}$$

The interval  $(0 < e^{-y/2} < 1)$  is equivalent to  $y > 0$ , we can conclude:

$$g(y) = \begin{cases} \frac{1}{2} e^{-y/2} & (y > 0) \\ 0 & (\text{elsewhere}) \end{cases} = \mathcal{C}(y; 2)$$

Now we have a powerful tool to prove rigorously the relationship between  $X$  and  $Z$ .

**Corollary 6.1.1.1.** *If  $X \sim \mathcal{N}(x; \mu_X, \sigma_X^2)$  and the random variable  $Z$  is defined as:*

$$Z = \frac{X - \mu_X}{\sigma_X}$$

*then:*

$$Z \sim \mathcal{N}(z; 0, 1)$$

*Proof.*

$$Z = u(X) = \frac{X - \mu_X}{\sigma_X} \Rightarrow X = w(Z) = \sigma_X Z + \mu_X$$

$$g(z) = f(w(z))|w'(z)| = \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{z^2}{2}\right) \sigma_X = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) = \mathcal{N}(z; 0, 1)$$

□

In the case of **joint pdf**, we can also use the theorem as follows:

**Theorem 6.1.2.** *Suppose that  $X_1$  and  $X_2$  are **discrete** random variable with joint pdf  $f(x_1, x_2)$ . Let  $Y_1 = u_1(X_1, X_2)$  and  $Y_2 = u_2(X_1, X_2)$  define a **linear transformation** so that their **inverse functions** can be found:*

$$X_1 = w_1(Y_1, Y_2); \quad X_2 = w_2(Y_1, Y_2)$$

*then the joint pdf of  $(Y_1, Y_2)$  is:*

$$g(y_1, y_2) = f(w_1(y_1, y_2), w_2(y_1, y_2))$$

If  $(X_1, X_2)$  are **continuous** random variables with joint pdf  $f(x_1, x_2)$ , then:

$$g(y_1, y_2) = f(w_1(y_1, y_2), w_2(y_1, y_2))|J|$$

where the Jacobian is the  $2 \times 2$  determinant:

$$J = \begin{vmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} \end{vmatrix}$$

Jacobian determinants knowledge belongs to Calculus 2, so I do not want to explain in detail here. You can think the role of the Jacobian simply as a transformation between coordinates.

**Theorem 6.1.2** can yield several significant results.

For example, let's find out the answer for our question from the beginning of this chapter: if  $X_1 \sim \mathcal{P}(x_1; \mu_1)$  and  $X_2 \sim \mathcal{P}(x_2; \mu_2)$  are 2 independent random variables, what is the pdf of  $Y = X_1 + X_2$ ? We define 2 new functions  $u_1(X_1, X_2)$ ,  $u_2(X_1, X_2)$  and determine their inverse functions:

$$Y_1 = u_1(X_1, X_2) = X_1 + X_2; Y_2 = u_2(X_1, X_2) = X_2$$

Like random variables, the functions  $u(X_1, X_2)$  can be arbitrarily chosen, as long as they are **linear**; and we should choose the **smartest** and **most convenient** choices. The inverse functions of  $u(X_1, X_2)$  are:

$$X_1 = w_1(Y_1, Y_2) = Y_1 - Y_2; X_2 = w_2(Y_1, Y_2) = Y_2$$

Since  $X_1$  and  $X_2$  are **independent**, we can conclude:

$$f(x_1, x_2) = f(x_1)f(x_2) = \frac{e^{-\mu_1}\mu_1^{x_1}}{x_1!} \cdot \frac{e^{-\mu_2}\mu_2^{x_2}}{x_2!} = \frac{e^{-(\mu_1+\mu_2)}\mu_1^{x_1}\mu_2^{x_2}}{x_1!x_2!}$$

Now we obtain the joint pdf of  $g(y_1, y_2)$ :

$$g(y_1, y_2) = f(w_1(y_1, y_2), w_2(y_1, y_2)) = f(y_1 - y_2, y_2) = \frac{e^{-(\mu_1+\mu_2)}\mu_1^{y_1-y_2}\mu_2^{y_2}}{(y_1 - y_2)!y_2!}$$

The pdf of  $Y_1$  is a **marginal distribution function** of  $g(y_1, y_2)$ :

$$\begin{aligned} h(y_1) &= \sum_{y_2} g(y_1, y_2) = \frac{e^{-(\mu_1+\mu_2)}}{y_1!} \sum_{y_2=0}^{y_1} \frac{y_1!}{(y_1 - y_2)!y_2!} \mu_1^{y_1-y_2} \mu_2^{y_2} = \frac{e^{-(\mu_1+\mu_2)}}{y_1!} \sum_{y_2=0}^{y_1} \binom{y_1}{y_2} \mu_1^{y_1-y_2} \mu_2^{y_2} \\ &= \frac{e^{-(\mu_1+\mu_2)}}{y_1!} (\mu_1 + \mu_2)^{y_1} = \mathcal{P}(y_1; \mu_1 + \mu_2) \end{aligned}$$

Our final conclusion is:

$$(X_1 + X_2) \sim \mathcal{P}(x_1 + x_2; \mu_1 + \mu_2)$$

. Now if  $X_1 \sim \mathcal{E}(x_1, 1)$ ,  $X_2 \sim \mathcal{E}(x_2, 1)$  are 2 independent random variables, what are the pdf of  $Y_1$  and  $Y_2$  if they are given by:

$$\begin{aligned} Y_1 &= X_1 + X_2 \\ Y_2 &= \frac{X_1}{X_1 + X_2} \end{aligned}$$

Similarly, we define 2 new functions:

$$Y_1 = u_1(X_1, X_2) = X_1 + X_2$$

$$Y_2 = u_2(X_1, X_2) = \frac{X_1}{X_1 + X_2}$$

The inverse of them are:

$$\begin{aligned} X_1 &= w_1(Y_1, Y_2) = Y_1 Y_2 \\ X_2 &= w_2(Y_1, Y_2) = Y_1 - Y_1 Y_2 \end{aligned}$$

Jacobian matrix and its determinant is:

$$J = \begin{vmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} \end{vmatrix} = \begin{vmatrix} y_2 & y_1 \\ 1 - y_2 & -y_1 \end{vmatrix} = -y_1$$

Since  $X_1$  and  $X_2$  are **statistically independent**, the joint pdf of them is:

$$f(x_1, x_2) = f(x_1)f(x_2) = e^{-(x_1+x_2)}$$

The joint pdf  $g(y_1, y_2)$  is:

$$g(y_1, y_2) = f(w_1(y_1, y_2), w_2(y_1, y_2))|J| = f(y_1 y_2, y_1 - y_1 y_2)y_1 = e^{-y_1}y_1$$

Notice that  $y_1 > 0$  and  $1 > y_2 > 0$ , the pdf  $Y_1$  and  $Y_2$  are **marginal distribution function** of  $g(y_1, y_2)$  with respect to  $y_1$  and  $y_2$ :

$$\begin{aligned} h_1(y_1) &= \int_0^1 g(y_1, y_2) dy_2 = \int_0^1 e^{-y_1} y_1 dy_2 = y_1 e^{-y_1} = \mathcal{G}(y_1; 2, 1) \\ h_2(y_2) &= \int_0^{+\infty} g(y_1, y_2) dy_1 = \int_0^{+\infty} e^{-y_1} y_1 dy_1 = 1 = \mathcal{U}(y_2; 0, 1) \end{aligned}$$

Our final results are:

$$\begin{aligned} Y_1 &\sim \mathcal{G}(y_1; 2, 1) \\ Y_2 &\sim \mathcal{U}(y_2; 0, 1) \end{aligned}$$

In this final example, you must **be careful** when determining the interval for random variables. Suppose  $X_1$  and  $X_2$  are both continuous random variables and their pdf is not a special case like our previous examples.

$$f(x_1, x_2) = \begin{cases} 24x_1x_2 & (0 \leq x_1, x_2 \leq 1; x_1 + x_2 \leq 1) \\ 0 & (\text{elsewhere}) \end{cases}$$

$Y_1$  and  $Y_2$  are defined as:

$$Y_1 = u_1(X_1, X_2) = X_1 + X_2; \quad Y_2 = u_2(X_1, X_2) = X_2$$

Their inverse functions are:

$$X_1 = w_1(Y_1, Y_2) = Y_1 - Y_2; \quad X_2 = w_2(Y_1, Y_2) = Y_2$$

Jacobian matrix and its determinant is:

$$J = \begin{vmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} \end{vmatrix} = \begin{vmatrix} 1 & -1 \\ 0 & 1 \end{vmatrix} = 1$$

The joint pdf  $g(y_1, y_2)$  is:

$$g(y_1, y_2) = f(w_1(y_1, y_2), w_2(y_1, y_2))|J| = f(y_1 - y_2, y_2) \cdot 1 = 24(y_1 - y_2)y_2$$

By replacing  $x_1 = y_1 - y_2$  and  $x_2 = y_2$ , we obtain new region:

$$(0 \leq x_1, x_2 \leq 1, x_1 + x_2 \leq 1) \rightarrow \begin{cases} y_1 - 1 \leq y_2 \leq y_1 \\ 0 \leq y_2 \leq 1 \\ y_1 \leq 1 \end{cases}$$

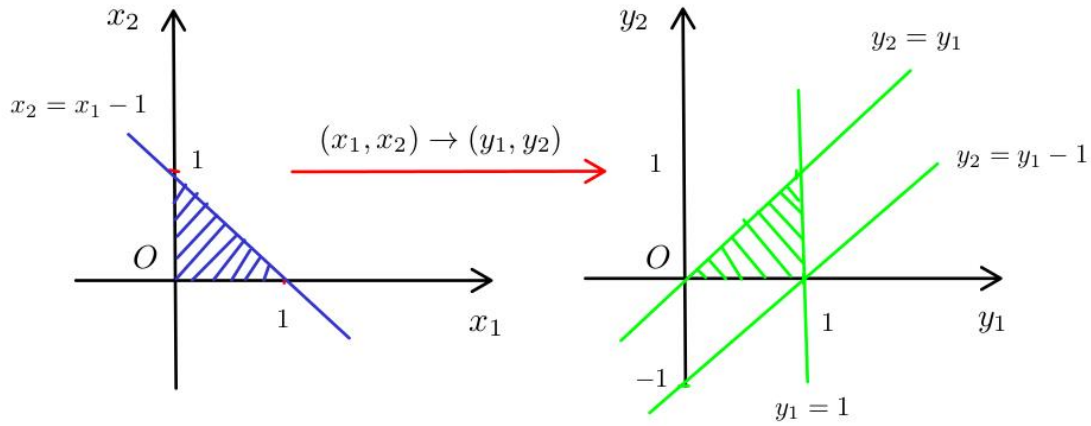


Figure 6.1: Changing the domain region

By looking at the graph, now you can see that  $0 \leq y_2 \leq y_1$ , and it is not so easy to see this bounded area if you do not plot it. Now the pdf of  $Y_1$  is:

$$h(y_1) = \int_{-\infty}^{+\infty} g(y_1, y_2) dy_2 = \int_0^{y_1} 24(y_1 - y_2)y_2 dy_2 = 4y_1^3$$

Since  $0 \leq y_1 \leq 1$ , now you can conclude the pdf of  $Y_1 = X_1 + X_2$  is:

$$Y_1 \sim h(y_1) = \begin{cases} 4y_1^3 & (0 \leq y_1 \leq 1) \\ 0 & (\text{elsewhere}) \end{cases}$$

### 6.1.2 Non-linear Transformations

**Theorem 6.1.3.** Suppose that  $X$  is a **continuous** random variable with pdf  $f(x)$ . Let  $Y = u(X)$  define a **non-linear** transformation. If the interval over which  $X$  is defined can be partitioned into  $k$  mutually disjoint sets such that each of the **inverse** functions:

$$X_1 = w_1(Y), X_2 = w_2(Y), \dots, X_k = w_k(Y)$$

of  $Y = u(X)$  defines a **linear** correspondence, then the pdf of  $Y$  is:

$$g(y) = \sum_{i=1}^k f(w_i(y)) |w'_i(y)|$$

Wow, this theorem seems a bit difficult at first, but its idea is very clear. Let me explain through an example:

Assume that you know the pdf of  $X$  is:

$$X \sim f(x) = \begin{cases} \frac{2(x+1)}{9} & (-1 < x < 2) \\ 0 & (\text{elsewhere}) \end{cases}$$

Now you want to determine the pdf of  $Y = X^2$ . Obviously you know the relationship between  $X$  and  $Y$  is **non-linear**; and using the same thinking line, you define the function  $Y = u(X)$  and find its inverse function:

$$Y = u(X) = X^2$$



But there is a problem, there are **more than one** existing inverse functions  $w(y)$  at the same time:

$$X = w_1(Y) = \sqrt{Y}, \quad X = w_2(Y) = -\sqrt{Y}$$

By partitioning the interval of  $X$ , you can see that:

$$X = \begin{cases} \sqrt{Y} & (0 < X < 2) \\ -\sqrt{Y} & (-1 < X < 0) \end{cases}$$

Since  $Y = X^2$ :

$$\begin{aligned} (0 < X < 2) &\Rightarrow (0 < Y < 4) \\ (-1 < X < 0) &\Rightarrow (0 < Y < 1) \end{aligned}$$

Inside the  $Y$  interval:

$$X = \begin{cases} \pm\sqrt{Y} & (0 < Y < 1) \\ \sqrt{Y} & (1 < Y < 4) \end{cases}$$

After splitting  $Y$  interval into 2 parts, now we apply our theorem above:

$$g(y) = \sum_{i=1}^k f(w_i(y)) |w'_i(y)| = \frac{1}{2\sqrt{y}} (f(\sqrt{y}) + f(-\sqrt{y})) = \frac{2}{9\sqrt{y}} \quad (0 < y < 1)$$

$$g(y) = \sum_{i=1}^k f(w_i(y)) |w'_i(y)| = \frac{1}{2\sqrt{y}} f(\sqrt{y}) = \frac{\sqrt{y} + 1}{9\sqrt{y}} \quad (1 < y < 4)$$

Rewrite our  $g(y)$  function in piecewise form:

$$g(y) = \begin{cases} \frac{2}{9\sqrt{y}} & (0 < y < 1) \\ \frac{\sqrt{y} + 1}{9\sqrt{y}} & (1 < y < 4) \\ 0 & (\text{elsewhere}) \end{cases}$$

**Corollary 6.1.3.1.** *If  $Z \sim \mathcal{N}(z; 0, 1)$  and  $Y = Z^2$ , then:*

$$Y \sim \mathcal{C}(y; 1)$$

*Proof.* This corollary is literally the backbone of **Statistics**. Now we define our function  $u(Z)$  and find its inverse functions  $w(Y)$ :

$$Y = u(Z) = Z^2 \Rightarrow Z = w_1(Y) = \sqrt{Y}, \quad Z = w_2(Y) = -\sqrt{Y}$$

Since  $Y \in (0, +\infty)$  and  $Z \in (-\infty, +\infty)$ , similarly, we have:

$$Z = \begin{cases} \pm\sqrt{Y} & (Y > 0) \\ \text{undefined} & (Y < 0) \end{cases}$$

Apply our theorem:

$$\begin{aligned} g(y) &= \sum_{i=1}^k f(w_i(y)) |w'_i(y)| = \frac{1}{2\sqrt{y}} (f(\sqrt{y}) + f(-\sqrt{y})) \\ &= \frac{1}{2\sqrt{y}} \cdot \frac{1}{\sqrt{2\pi}} \left( \exp\left(-\frac{y}{2}\right) + \exp\left(-\frac{y}{2}\right) \right) \\ &= \frac{1}{y^{1/2}} \cdot \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y}{2}\right) \\ &= \frac{1}{2^{1/2}\Gamma(1/2)} y^{1/2-1} e^{-y/2} \\ &= \mathcal{C}(y; 1) \end{aligned}$$

□

## 6.2 Moment-Generating Functions

### 6.2.1 Definition of Moment-Generating Functions

Every pdf have their own **mean** and **variance**:

$$\mu_X = E(X)$$

$$\sigma_X^2 = E(X^2) - \mu_X^2$$

These functions are all characterized by  $E(X^r)$  functions, you can think them as human fingerprints. 2 pdf are the same if and only if all of their  $E(X^r)$  functions are **matched**. But directly determining  $E(X^r)$  is not so simple, especially with  $r \geq 3$ , therefore we take advantage of the derivative property of exponential function to develop a powerful tool that can yield desired results almost immediately.

**Definition 6.2.1.** *The  $r$ th moment about the origin of the random variable  $X$  is given:*

$$\mu'_r = E(X^r) = \sum_x x^r f(x) \quad (\text{if } X \text{ is discrete})$$

$$\mu'_r = E(X^r) = \int_{-\infty}^{+\infty} x^r f(x) dx \quad (\text{if } X \text{ is continuous})$$

**Definition 6.2.2.** *The **moment-generating function** of the random variable  $X$  is:*

$$M_X(t) = E(e^{tX}) = \sum_x e^{tx} f(x) \quad (\text{if } X \text{ is discrete})$$

$$M_X(t) = E(e^{tX}) = \int_{-\infty}^{+\infty} e^{tx} f(x) dx \quad (\text{if } X \text{ is continuous})$$

**Theorem 6.2.1.** *Let  $X$  be a random variable with **moment-generating function**  $M_X(t)$ , then:*

$$\left. \frac{d^r M_X(t)}{dt^r} \right|_{t=0} = \mu'_r$$

*Proof.* We only prove the **discrete case** because the continuous case is similar:

$$\frac{d^r M_X(t)}{dt^r} = \frac{d^r (\sum_x e^{tx} f(x))}{dt^r} = \sum_x x^r e^{tx} f(x)$$

After setting  $t = 0$ , finally we obtain:

$$\left. \sum_x x^r e^{tx} f(x) \right|_{t=0} = \sum_x x^r f(x) = \mu'_r$$

□

## 6.2.2 Some Useful Moment-Generating Functions

### Discrete Probability Distribution

**Corollary 6.2.1.1.** *If  $X \sim \mathcal{B}(x; n, p)$ , then:*

$$M_X(t) = (pe^t + q)^n$$

*Proof.*

$$M_X(t) = \sum_x e^{tx} f(x) = \sum_{x=0}^n \binom{n}{x} (e^t p)^x q^{n-x} = (pe^t + q)^n$$

□

**Corollary 6.2.1.2.** *If  $X \sim \mathcal{P}(x; \mu)$ , then:*

$$M_X(t) = \exp(\mu(e^t - 1))$$

*Proof.*

$$\begin{aligned} M_X(t) &= \sum_x e^{xt} f(x) = \sum_{x=0}^n e^{xt} \cdot \frac{e^{-\mu} \mu^x}{x!} \\ &= e^{-\mu} \sum_{x=0}^n \frac{(e^t \mu)^x \exp(-e^t \mu)}{x!} \exp(e^t \mu) \\ &= \exp(\mu(e^t - 1)) \sum_{x=0}^n \mathcal{P}(x; \mu e^t) \\ &= \exp(\mu(e^t - 1)) \end{aligned}$$

□

**Corollary 6.2.1.3.** *If  $X \sim \mathcal{G}^*(x; p)$ , then:*

$$M_X(t) = \frac{pe^t}{1 - qe^t}$$

*Proof.*

$$\begin{aligned} M_X(t) &= \sum_x e^{xt} f(x) = \sum_{x=1}^{+\infty} e^{xt} p q^{x-1} = pe^t \sum_{x=1}^{+\infty} (e^t q)^{x-1} \\ &= pe^t \frac{1}{1 - qe^t} = \frac{pe^t}{1 - qe^t} \end{aligned}$$

Notice the convergence condition of geometric sum is:  $|qe^t| < 1$

□

### Continuous Probability Distribution

**Corollary 6.2.1.4.** *If  $X \sim \mathcal{N}(x; \mu_X, \sigma_X^2)$ , then:*

$$M_X(t) = \exp\left(\mu_X t + \frac{1}{2} \sigma_X^2 t^2\right)$$

*Proof.*

$$\begin{aligned}
M_X(t) &= \int_{-\infty}^{+\infty} e^{xt} f(x) dx = \frac{1}{\sqrt{2\pi}\sigma_X} \int_{-\infty}^{+\infty} e^{xt} \exp\left(-\frac{(x-\mu_X)^2}{2\sigma_X^2}\right) dx \\
&= \frac{1}{\sqrt{2\pi}\sigma_X} \int_{-\infty}^{+\infty} \exp\left(xt - \frac{(x-\mu_X)^2}{2\sigma_X^2}\right) dx \\
&= \frac{1}{\sqrt{2\pi}\sigma_X} \int_{-\infty}^{+\infty} \exp\left(\frac{-(x - (\mu_X + t\sigma_X^2))^2 + 2(\mu_X t\sigma_X^2 + t^2\sigma_X^4)}{2\sigma_X^2}\right) dx \\
&= \exp\left(\mu_X t + \frac{1}{2}t^2\sigma_X^2\right) \frac{1}{\sqrt{2\pi}\sigma_X} \int_{-\infty}^{+\infty} \exp\left(-\frac{(x - (\mu_X + t\sigma_X^2))^2}{2\sigma_X^2}\right) dx \\
&= \exp\left(\mu_X t + \frac{1}{2}\sigma_X^2 t^2\right) \int_{-\infty}^{+\infty} \mathcal{N}(x; \mu_X + t\sigma_X, \sigma_X^2) dx \\
&= \exp\left(\mu_X t + \frac{1}{2}\sigma_X^2 t^2\right)
\end{aligned}$$

□

**Corollary 6.2.1.5.** *If  $X \sim \mathcal{C}(x; v)$ , then:*

$$M_X(t) = (1 - 2t)^{-v/2}$$

*Proof.*

$$\begin{aligned}
M_X(t) &= \int_{-\infty}^{+\infty} e^{xt} f(x) dx = \int_{-\infty}^{+\infty} e^{xt} \frac{1}{2^{v/2}\Gamma(v/2)} x^{v/2-1} e^{-x/2} dx \\
&= \frac{1}{2^{v/2}\Gamma(v/2)} \int_{-\infty}^{+\infty} x^{v/2-1} \exp\left(-\frac{x}{2/(1-2t)}\right) dx \\
&= \frac{(1-2t)^{-v/2}}{(2/(1-2t))^{v/2}\Gamma(v/2)} \int_{-\infty}^{+\infty} x^{v/2-1} \exp\left(-\frac{x}{2/(1-2t)}\right) dx \\
&= (1-2t)^{-v/2} \int_{-\infty}^{+\infty} \mathcal{G}\left(x; \frac{v}{2}, \frac{2}{1-2t}\right) dx \\
&= (1-2t)^{-v/2}
\end{aligned}$$

□

### 6.2.3 Linear Combinations of Random Variables

Moment-Generating function is an extremely powerful tool, not only for finding moments, but also can be applied to prove several significant theorems related to linear combinations of random variables, like **Central Limit Theorem**. But in this subsection we only focus on 4 main results, which are directly deduced from moment-generating function; all of them are the foundation for the **Statistics**.

**Theorem 6.2.2.**  *$X$  and  $Y$  are 2 random variables with moment-generating functions  $M_X(t)$  and  $M_Y(t)$ , respectively, if:*

$$M_X(t) = M_Y(t) \quad (\forall t)$$

*then  $X$  and  $Y$  have the same pdf.*

Since proving rigorously is beyond the scope of this book, and perhaps it is not necessary because this theorem is too obvious. Like I said before, if 2 random variables have the same moments ("fingerprints"), certainly they share the same pdf.

**Theorem 6.2.3.** If  $X_1, X_2, \dots, X_n$  are independent random variables with moment-generating functions  $M_{X_1}(t), M_{X_2}(t), \dots, M_{X_n}(t)$ , respectively, and  $Y = X_1 + X_2 + \dots + X_n$ , then:

$$M_Y(t) = M_{X_1}(t)M_{X_2}(t) \cdots M_{X_n}(t)$$

*Proof.* It's very easy to see that:

$$\begin{aligned} M_Y(t) &= E(e^{tY}) = E(e^{t(X_1+X_2+\dots+X_n)}) = E(e^{tX_1}e^{tX_2} \dots e^{tX_n}) \\ &= E(e^{tX_1})E(e^{tX_2}) \dots E(e^{tX_n}) = M_{X_1}(t)M_{X_2}(t) \cdots M_{X_n}(t) \end{aligned}$$

□

**Corollary 6.2.3.1.** If  $X_1, X_2, \dots, X_n$  are independent random variables having **normal distributions** with **means**  $\mu_{X_1}, \mu_{X_2}, \dots, \mu_{X_n}$  and **variances**  $\sigma_{X_1}^2, \sigma_{X_2}^2, \dots, \sigma_{X_n}^2$ , respectively, then the random variable  $Y = X_1 + X_2 + \dots + X_n$  has a **normal distribution**.

*Proof.* The moment-generating function of  $Y$  is:

$$\begin{aligned} M_Y(t) &= M_{X_1}(t)M_{X_2}(t) \cdots M_{X_n}(t) \\ &= \exp\left(\mu_{X_1}t + \frac{1}{2}\sigma_{X_1}^2t^2\right) \exp\left(\mu_{X_2}t + \frac{1}{2}\sigma_{X_2}^2t^2\right) \cdots \exp\left(\mu_{X_n}t + \frac{1}{2}\sigma_{X_n}^2t^2\right) \\ &= \exp\left((\mu_{X_1} + \mu_{X_2} + \dots + \mu_{X_n})t + \frac{1}{2}(\sigma_{X_1}^2 + \sigma_{X_2}^2 + \dots + \sigma_{X_n}^2)t^2\right) \end{aligned}$$

Because  $M_Y(t)$  has a normal distribution form, so we conclude  $Y$  has a **normal distribution**.

□

**Corollary 6.2.3.2.** If  $X_1, X_2, \dots, X_n$  are independent random variables having **chi-squared distributions** with  $v_1, v_2, \dots, v_n$  **degrees of freedom**, then the random variable  $Y = X_1 + X_2 + \dots + X_n$  has a **chi-squared distribution**.

*Proof.* The moment-generating function of  $Y$  is:

$$\begin{aligned} M_Y(t) &= M_{X_1}(t)M_{X_2}(t) \cdots M_{X_n}(t) \\ &= (1 - 2t)^{-v_1/2}(1 - 2t)^{-v_2/2} \cdots (1 - 2t)^{-v_n/2} \\ &= (1 - 2t)^{-(v_1+v_2+\dots+v_n)/2} \end{aligned}$$

Because  $M_Y(t)$  has a chi-squared distribution form, so we conclude  $Y$  has a **chi-squared distribution** with  $v = v_1 + v_2 + \dots + v_n$  degrees of freedom.

□

# Chapter 7

## Fundamentals of Statistics

Going back to the 19th and 20th centuries, all the brilliant minds in the field of Biology, such as George Mendel or Ronald Fisher, were masters of Statistics. This is no coincidence, because they not only had to conduct experiments, but also had to process the result. Have you ever wondered how Mendel could derive his theory of heredity just from growing peas? Of course, he could not observe all the peas on Earth; he could only observe some of the peas in his garden. He had to use **rules of statistical inference** to draw conclusions from the **sample set** (his garden) to the entire **population** (the Earth). But now, what if you repeat his experiment again in your own garden and do not get the same result? Is something wrong? The answer is more complex than that. We have to consider many factors:

- First, are your seeds good and purebred quality?
- Second, are you sure that your seeds were randomly selected? Otherwise, your result would be meaningless.
- Ultimately, what is the probability of that result occurring? If it is less than 5%, you can ignore it as an exception.

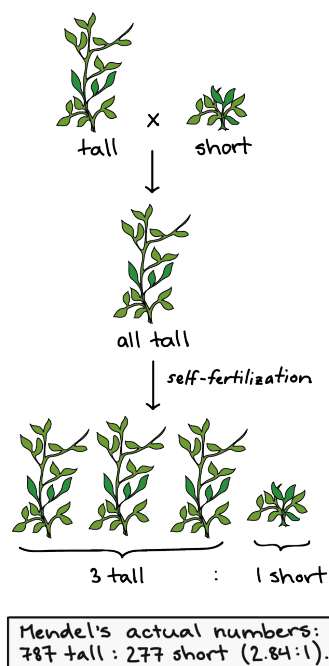


Figure 7.1: Why was Mendel so certain about 3 : 1 ratio?

## 7.1 The Big Picture of Statistics

### 7.1.1 Populations and Samples

The germination rate of mung bean seeds is considered to be  $p = 0.8$ , and germination occurs within 3 – 5 days. How do you know if that is true or not? Like Mendel did, we have to perform an experiment to test it. After buying a packet of mung bean seeds in the agriculture store, you randomly choose 6 **samples**, with each containing exactly 10 seeds. Place each sample on a separate petri dish, lined with wet tissue inside, and avoid placing them in direct sunlight.

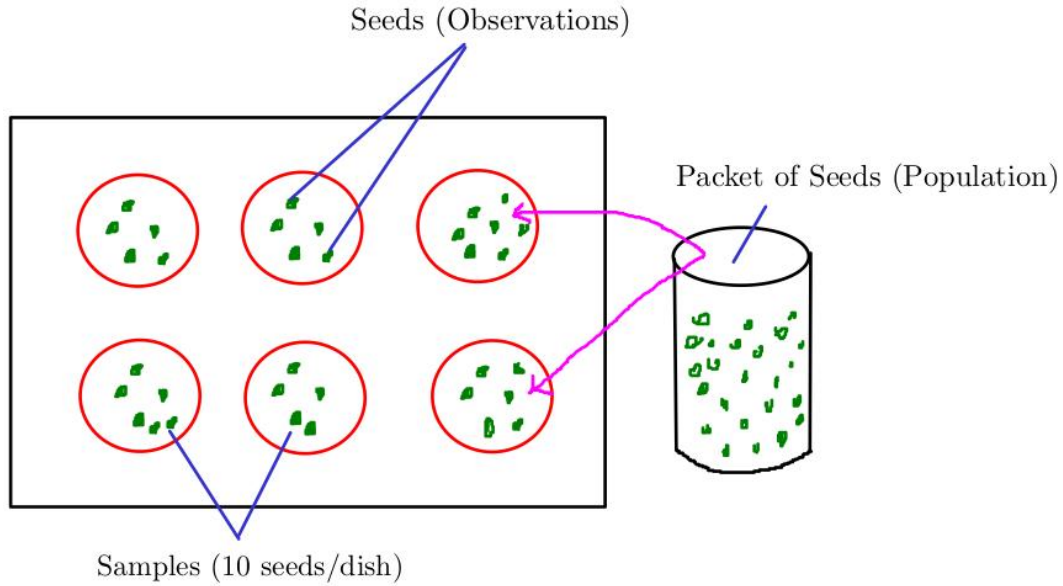


Figure 7.2: Sampling process

Our goal is instead of testing all seeds in the packet (**population**), we observe 6 petri dishes **samples** with each containing 10 seeds (**observations**) and draw conclusions based on our observed data.

**Definition 7.1.1.** A **population** consists of the totality of the **observations** with which we are concerned.

**Definition 7.1.2.** A **sample** is a **subset** of a population.

Conceptually, you can relate these terms with sample space, events and sample points.

$$\text{Sample Points} \subseteq \text{Event} \subseteq \text{Sample Space}$$

$$\text{Observations} \subseteq \text{Sample} \subseteq \text{Population}$$

Now consider a dish and define the random variable  $X$  as germination time of a single seed. There are 10 seeds per dish, so  $X_1, X_2, \dots, X_{10}$  are the germination time of the 1st, 2nd,  $\dots$ , 10th seed. Since they are **independent** random variables, each having the same pdf  $f(x)$ , so we can deduce:

$$f(x_1, x_2, \dots, x_{10}) = f(x_1)f(x_2) \cdots f(x_{10})$$

**Definition 7.1.3.** Let  $X_1, X_2, \dots, X_n$  be  $n$  **independent** random variables, each having the same pdf  $f(x)$ . The **sample** they belong to can also be called a **random sample** of size  $n$  from the **population**  $f(x)$ , and has the joint pdf:

$$f(x_1, x_2, \dots, x_n) = f(x_1)f(x_2) \cdots f(x_n)$$

### 7.1.2 Sample Mean and Sample Variance

After 5 waiting days and with the aid of time-lapsed camera, you can collect 60 values of  $X$ , but now let's take a look at a single dish (**sample**) with 10 values of  $X_i$ .

|       |      |      |      |      |      |      |      |      |     |      |
|-------|------|------|------|------|------|------|------|------|-----|------|
| $i$   | 1    | 2    | 3    | 4    | 5    | 6    | 7    | 8    | 9   | 10   |
| $x_i$ | 4.66 | 4.27 | 4.29 | 4.96 | 3.49 | 3.04 | 2.87 | 3.53 | 5.0 | 2.67 |

Now we define several important quantities, which describe measures of sample such as **sample mean**, **sample variance**,...

**Definition 7.1.4.** The **sample mean** of a sample set containing  $n$  observations is defined as:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

The **mean** of our sample is:

$$\bar{x} = \frac{1}{10} \sum_{i=1}^{10} x_i = 3.878$$

**Definition 7.1.5.** The **sample variance** of a sample set containing  $n$  observations is defined as:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

The **variance** of our sample is:

$$s^2 = \frac{1}{10-1} \sum_{i=1}^{10} (x_i - \bar{x})^2 = 0.757$$

The reason of using  $n-1$  as a divisor rather than  $n$  will be explained in the next chapter. Since calculating the **sample variance** directly is quiet easy to make miscalculations, so in many cases we always apply this theorem belows:

**Theorem 7.1.1.** The **sample variance** of a sample set containing  $n$  observations is:

$$S^2 = \frac{1}{n(n-1)} \left( n \sum_{i=1}^n X_i^2 - \left( \sum_{i=1}^n X_i \right)^2 \right)$$

*Proof.* By definition:

$$\begin{aligned} S^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \left( \sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + n\bar{X}^2 \right) \\ &= \frac{1}{n-1} \left( \sum_{i=1}^n X_i^2 - \frac{1}{n} \left( \sum_{i=1}^n X_i \right)^2 \right) \\ &= \frac{1}{n(n-1)} \left( n \sum_{i=1}^n X_i^2 - \left( \sum_{i=1}^n X_i \right)^2 \right) \end{aligned}$$

□



Applying directly this theorem, now the **variance** of our sample is:

$$s^2 = \frac{1}{9.10} \left( 10 \sum_{i=1}^{10} x_i^2 - \left( \sum_{i=1}^{10} x_i \right)^2 \right) = 0.757$$

You should note that  $\bar{X}$  and  $S$  are 2 random variables constituting from a **random sample**, and they are typical examples of **statistics**.

**Definition 7.1.6.** Any function of the random variables constituting a **random sample** is called a **statistic**.

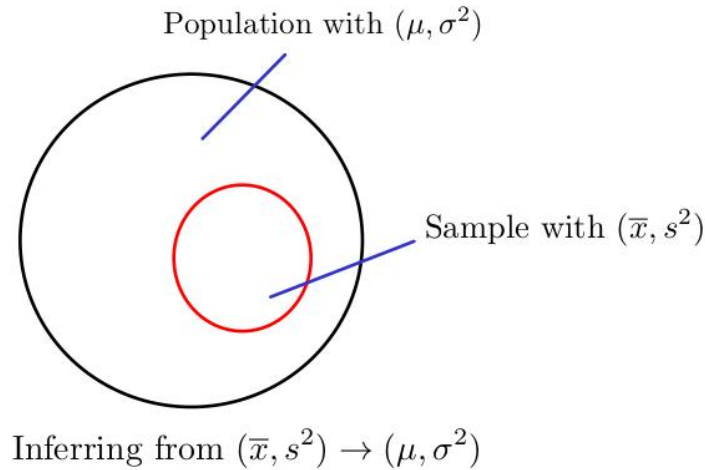


Figure 7.3: Inferring from  $(\bar{x}, s^2) \rightarrow (\mu, \sigma^2)$

After conducting the experiment, you obtain 6 values for  $\bar{x}$  and  $s^2$ . But how do you process them? Do they belong to any probability distribution, or are they simply random numbers? These questions will be answered in the remainder of this chapter!

## 7.2 Sampling Distribution of Means

### 7.2.1 Central Limit Theorem

Before delving into the general case, let's consider a special case of Central Limit Theorem (CLT for short) first.

**Theorem 7.2.1.** (Special case of CLT) If  $\bar{X}$  is the **mean** of a **random sample** of size  $n$  from a **normal population** with mean  $\mu$  and variance  $\sigma^2$ , then the pdf of  $\bar{X}$  is:

$$\bar{X} \sim \mathcal{N} \left( \bar{x}; \mu, \frac{\sigma^2}{n} \right)$$

*Proof.*  $\bar{X}$  is just a **linear combinations** of  $X_1, X_2, \dots, X_n$ :

$$\bar{X} = \frac{1}{n}(X_1 + X_2 + \dots + X_n)$$

Since all of  $X_i$  are normally distributed, from our previous results from **Chapter 3: Mathematical Expectation** and **Chapter 6: Functions of Random Variables**, we know the

random variable  $\bar{X}$  has a **normal distribution** with its **mean** and **variance**:

$$\begin{aligned}\mu_{\bar{X}} &= \frac{1}{n}(\mu + \mu + \dots + \mu) = \frac{1}{n}.n\mu = \mu \\ \sigma_{\bar{X}}^2 &= \frac{1}{n^2}(\sigma^2 + \sigma^2 + \dots + \sigma^2) = \frac{1}{n^2}.n\sigma^2 = \frac{\sigma^2}{n}\end{aligned}$$

Finally, we conclude:

$$\bar{X} \sim \mathcal{N}\left(\bar{x}; \mu, \frac{\sigma^2}{n}\right)$$

□

**Theorem 7.2.2.** (General case of CLT) If  $\bar{X}$  is the **mean** of a **random sample** of size  $n$  from an **arbitrary population** with mean  $\mu$  and variance  $\sigma^2$ , then the pdf of  $\bar{X}$  is:

$$\bar{X} \sim \mathcal{N}\left(\bar{x}; \mu, \frac{\sigma^2}{n}\right)$$

This theorem is so amazing, it states that every pdf will **converge** to a **normal distribution**, regardless of its origin. This is the true power of CLT: the **normal distribution** is the **ultimate goal** of everything and not a magical distribution which drops from the sky.

*Proof.* Proving this theorem is not so easy, but we will go step-by-step. Assume that the random variables  $X_1, X_2, \dots, X_n$  have their moments-generating functions  $M_{X_1}(t), M_{X_2}(t), \dots, M_{X_n}(t)$ . Since  $n\bar{X}$  is just a linear combinations of  $X_i$  and they are in the same **population**:

$$\bar{X} = \frac{1}{n}(X_1 + X_2 + \dots + X_n) \Rightarrow M_{n\bar{X}}(t) = (M_X(t))^n$$

We restrict our constraints  $t \approx 0$  and  $n \rightarrow +\infty$  to use Taylor-Maclaurin approximation:

$$\begin{aligned}(M_X(t))^n &= \left(\sum_{k=0}^{+\infty} \frac{d^k M_X(0)}{k!} t^k\right)^n \\ &= \left(1 + \mu t + \frac{1}{2}\mu'_2 t^2 + \dots\right)^n \\ &\approx \left(1 + \mu t + \frac{\sigma^2 + \mu^2}{2} t^2\right)^n \\ &\approx \left(1 + \mu t + \frac{\sigma^2 t^2}{2}\right)^n \quad \left(\text{because } \mu t \gg \frac{1}{2}\mu^2 t^2\right) \\ &\approx \left(\exp\left(\mu t + \frac{\sigma^2 t^2}{2}\right)\right)^n \\ &= \exp\left(n\mu t + \frac{n\sigma^2 t^2}{2}\right)\end{aligned}$$

From the moments-generating function property:

$$M_{n\bar{X}}(t) = E(e^{nt\bar{X}}) = M_{\bar{X}}(nt)$$

Recall that the moment-generating function of **normal distribution** is:

$$\exp\left(\mu t + \frac{\sigma^2 t^2}{2}\right)$$

If:

$$\bar{X} \sim \mathcal{N}\left(\bar{x}; \mu, \frac{\sigma^2}{n}\right)$$

then we can deduce:

$$M_{\bar{X}}(t) = \exp\left(\mu t + \frac{\sigma^2 t^2}{2n}\right) \Rightarrow M_{\bar{X}}(nt) = \exp\left(n\mu t + \frac{\sigma^2 t^2 n^2}{2n}\right) = \exp\left(n\mu t + \frac{n\sigma^2 t^2}{2}\right) = (M_X(t))^n$$

Finally, we obtain our result is:

$$\bar{X} \sim \mathcal{N}\left(\bar{x}; \mu, \frac{\sigma^2}{n}\right) \quad \square$$

You should note that the condition of our approximation is  $n \rightarrow +\infty$ . In practice, we always choose  $n > 30$  as the criterion for using CLT. You can clearly see that we **can not use CLT** to process our data since each dish (sample) just has 10 seeds (observations) inside and  $n = 10 < 30$ .

**Corollary 7.2.2.1.** *If*

$$\bar{X} \sim \mathcal{N}\left(\bar{x}; \mu, \frac{\sigma^2}{n}\right)$$

*then:*

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(z; 0, 1)$$

## 7.2.2 t-Distribution

We have to change our strategy since our number of observations is pretty low (just 10 seeds per dish). Before developing a new tool, let me introduce some important results which play crucial roles first.

**Corollary 7.2.2.2.** *If  $X_1, X_2, \dots, X_n$  have **normal distribution** with the **same mean and variance**, then the sum:*

$$\sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma}\right)^2$$

*has **chi-squared distribution** with  $n$  degrees of freedom.*

*Proof.* We have already known that if the random variable  $Z$  is defined as:

$$Z_i = \frac{X_i - \mu}{\sigma}$$

then  $Z^2 \sim \mathcal{C}(z; 1)$ . Using the **linear property of chi-squared distribution**, you can conclude the sum:

$$\sum_{i=1}^n Z_i^2 = \sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma}\right)^2$$

also has **chi-squared distribution** with  $n$  degrees of freedom.  $\square$

**Theorem 7.2.3.** *If  $S^2$  is the **variance** of a random sample of size  $n$  taken from a **normal population** having the variance  $\sigma^2$ , then the **statistic**:*

$$\chi^2 = \frac{S^2(n-1)}{\sigma^2} = \sum_{i=1}^n \left(\frac{X_i - \bar{X}}{\sigma}\right)^2 \sim \mathcal{C}(\chi^2; n-1)$$

*Proof.* From the sum:

$$\begin{aligned}
\sum_{i=1}^n (X_i - \mu)^2 &= \sum_{i=1}^n (X_i - \bar{X} + \bar{X} - \mu)^2 \\
&= \sum_{i=1}^n (X_i - \bar{X})^2 + 2 \sum_{i=1}^n (X_i - \bar{X})(\bar{X} - \mu) + \sum_{i=1}^n (\bar{X} - \mu)^2 \\
&= \sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{i=1}^n (\bar{X} - \mu)^2 \\
\Rightarrow \sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2 &= \sum_{i=1}^n \left( \frac{X_i - \bar{X}}{\sigma} \right)^2 + \sum_{i=1}^n \left( \frac{\bar{X} - \mu}{\sigma} \right)^2 \\
&= \sum_{i=1}^n \left( \frac{X_i - \bar{X}}{\sigma} \right)^2 + \left( \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right)^2 \\
&= \chi^2 + \left( \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right)^2
\end{aligned}$$

From the **linear property** of chi-squared distribution, we can conclude the random variable  $\chi^2$  has chi-squared distribution with  $v = n - 1$  degrees of freedom.  $\square$

**Theorem 7.2.4.** Let  $X_1 \sim \mathcal{N}(x_1; 0, 1)$ ,  $X_2 \sim \mathcal{C}(x_2; v)$  and they are both **independent**, then:

$$Y_1 = \frac{X_1}{\sqrt{X_2/v}} \sim h(y) = \frac{1}{\sqrt{\pi v}} \cdot \frac{\Gamma((v+1)/2)}{\Gamma(v/2)} \left( 1 + \frac{y_1^2}{v} \right)^{-(v+1)/2}$$

*Proof.* Since  $X_1$  and  $X_2$  are 2 **statistically independent** random variables, the joint pdf  $f(x_1, x_2)$  is:

$$\begin{aligned}
f(x_1, x_2) &= f(x_1)f(x_2) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x_1^2}{2}\right) \frac{1}{2^{v/2}\Gamma(v/2)} x_2^{v/2-1} \exp\left(-\frac{x_2}{2}\right) \\
&= \frac{1}{\sqrt{2\pi}} \cdot \frac{1}{2^{v/2}\Gamma(v/2)} \exp\left(-\frac{x_1^2 + x_2}{2}\right) x_2^{v/2-1}
\end{aligned}$$

Set 2 auxiliary random variables  $Y_1$  and  $Y_2$  as:

$$Y_1 = u_1(X_1, X_2) = \frac{X_1}{\sqrt{X_2/v}}$$

$$Y_2 = u_2(X_1, X_2) = X_2$$

The inverse functions of  $u(X_1, X_2)$  are:

$$X_1 = w_1(Y_1, Y_2) = \frac{Y_1 \sqrt{Y_2}}{\sqrt{v}}$$

$$X_2 = w_2(Y_1, Y_2) = Y_2$$

The Jacobian of this transformation is:

$$J = \begin{vmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} \end{vmatrix} = \begin{vmatrix} \sqrt{\frac{y_2}{v}} & \frac{y_1}{2\sqrt{v y_2}} \\ 0 & 1 \end{vmatrix} = \sqrt{\frac{y_2}{v}}$$

Now the joint pdf of  $Y_1$  and  $Y_2$  is:

$$\begin{aligned}
g(y_1, y_2) &= f(w_1(y_1, y_2), w_2(y_1, y_2)) |J| \\
&= f\left(\frac{y_1 \sqrt{y_2}}{\sqrt{v}}, y_2\right) \cdot \frac{\sqrt{y_2}}{\sqrt{v}} \\
&= \frac{1}{\sqrt{2\pi}} \cdot \frac{1}{2^{v/2} \Gamma(v/2)} \exp\left(-\frac{1}{2} \cdot \left(\frac{y_1^2 y_2}{v} + y_2\right)\right) y_2^{v/2-1} \frac{\sqrt{y_2}}{\sqrt{v}} \\
&= \frac{1}{\sqrt{2\pi v} 2^{v/2} \Gamma(v/2)} \exp\left(-\frac{y_1^2 y_2 + y_2 v}{2v}\right) \cdot y_2^{((v+1)/2)-1} \\
&= \frac{1}{\sqrt{2\pi v} 2^{v/2} \Gamma(v/2)} \exp\left(-\frac{y_2}{2v/(y_1^2 + v)}\right) \cdot y_2^{((v+1)/2)-1}
\end{aligned}$$

Hence the pdf of  $Y_1$  is the marginal distribution of  $g(y_1, y_2)$ :

$$\begin{aligned}
h(y_1) &= \int_{-\infty}^{+\infty} g(y_1, y_2) dy_2 = \frac{1}{\sqrt{2\pi v} 2^{v/2} \Gamma(v/2)} \int_{-\infty}^{+\infty} \exp\left(-\frac{y_2}{2v/(y_1^2 + v)}\right) \cdot y_2^{((v+1)/2)-1} dy_2 \\
&= \frac{1}{\sqrt{2\pi v} 2^{v/2} \Gamma(v/2)} \int_{-\infty}^{+\infty} \exp\left(-\frac{y_2}{2v/(y_1^2 + v)}\right) \cdot y_2^{((v+1)/2)-1} dy_2 \\
&= \frac{(2v/(y_1^2 + v))^{(v+1)/2} \Gamma((v+1)/2)}{\sqrt{2\pi v} 2^{v/2} \Gamma(v/2)} \\
&= \int_{-\infty}^{+\infty} \frac{1}{(2v/(y_1^2 + v))^{(v+1)/2} \Gamma((v+1)/2)} \exp\left(-\frac{y_2}{2v/(y_1^2 + v)}\right) \cdot y_2^{((v+1)/2)-1} dy_2 \\
&= \frac{(2v/(y_1^2 + v))^{(v+1)/2} \Gamma((v+1)/2)}{\sqrt{2\pi v} 2^{v/2} \Gamma(v/2)} \cdot \int_{-\infty}^{+\infty} \mathcal{G}\left(y_2; \frac{v+1}{2}, \frac{2v}{y_1^2 + v}\right) dy_2 \\
&= \frac{(2v/(y_1^2 + v))^{(v+1)/2} \Gamma((v+1)/2)}{\sqrt{2\pi v} 2^{v/2} \Gamma(v/2)} \\
&= \frac{1}{\sqrt{\pi v}} \cdot \frac{\Gamma((v+1)/2)}{\Gamma(v/2)} \cdot \left(\frac{y_1^2}{v} + 1\right)^{-(v+1)/2}
\end{aligned}$$

□

**Theorem 7.2.4** is completely a mess, but its idea is very clear. We can rephrase this theorem to our convenience.

**Theorem 7.2.5.** Let  $Z$  be standard **normal random variable** and  $V$  is a **chi-squared random variable** with  $v$  degrees of freedom and they are **independent**, then the distribution of the random variable  $T$  is given by:

$$T = \frac{Z}{\sqrt{V/v}} \sim \mathcal{T}(t; v) = \frac{1}{\sqrt{\pi v}} \cdot \frac{\Gamma((v+1)/2)}{\Gamma(v/2)} \cdot \left(\frac{t^2}{v} + 1\right)^{-(v+1)/2} \quad (-\infty < t < +\infty)$$

This pdf is known as the **t-distribution** with  $v$  degrees of freedom.

Relating to our previous results, we already have:

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(z; 0, 1)$$

$$V = \chi^2 = \frac{S^2(n-1)}{\sigma^2} \sim \mathcal{C}(v; n-1)$$

Rewriting the theorem above yields an interesting result:

$$T = \frac{Z}{\sqrt{V/v}} = \frac{(\bar{X} - \mu)/(\sigma/\sqrt{n})}{S/\sigma} = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim \mathcal{T}(t; n - 1)$$

**Theorem 7.2.6.** *If  $\bar{X}$  is the **mean** of a **random sample** of size  $n$  from a **normal population** with mean  $\mu$ , **unknown variance**  $\sigma^2$ , then the pdf of random variable  $T$  is:*

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim \mathcal{T}(t; n - 1)$$

where  $S^2$  is the **variance** of the sample in which  $\bar{X}$  belongs to.

Compared to CLT, you can see that in t-distribution **deriving process**, we do not even have to touch the constraint  $n \rightarrow +\infty$ , or at least  $n > 30$ ; so t-distribution is a great tool when dealing with small sample sets where  $n < 30$ . But trade-off with small sample sets is the constraint **population must be normally distributed**; so t-distribution should not be used if you are not certain if the population is normal.

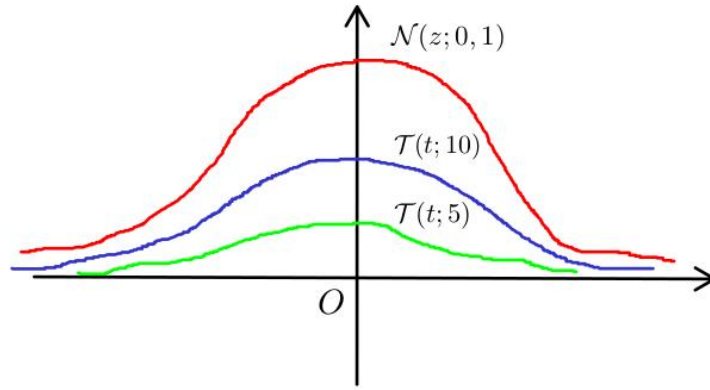


Figure 7.4: t-distribution and Standard Normal distribution

The relationship between t-distribution and standard normal distribution is remarkable. Both of them have **bell-shaped** curves, and logically from their usage (large or small sample), you can describe their relationship as useful approximations:

**Theorem 7.2.7.** *If  $v \rightarrow +\infty$ , then:*

$$\mathcal{T}(t; v) \rightarrow \mathcal{N}(z; 0, 1)$$

In practice, if  $v$  degrees of freedom is greater than 30 (or equivalent to sample size  $n > 31$  because  $v = n - 1$ ), we can use standard normal instead of t-distribution.

**Corollary 7.2.7.1.** *If a random sample has size  $n > 31$ , then:*

$$t \approx z \Rightarrow s \approx \sigma$$

### 7.3 Sampling Distribution of Variances

After gathering 6 values of  $S^2$  from our experiment, now we wish to know what the probability distribution they belong to. Let me recall an useful result which we have just proved in the previous section:

**Corollary 7.3.0.1.** *If  $S^2$  is the **variance** of a random sample of size  $n$  taken from a **normal population** having the variance  $\sigma^2$ , then the statistic:*

$$\chi^2 = \frac{S^2(n-1)}{\sigma^2} \sim \mathcal{C}(\chi^2; n-1)$$

The criterion **normal population** is very important; and our conclusions will not be corrected if it is violated.

## 7.4 Case Study: Seed Germination Time

It is time to summarize all the theorems that we have developed so far and try applying them as tools in a real experiment. There are 3 main tools and they can be called as Z-test, t-test and Chi-test, respectively:

1. CLT (General case): In a **large** ( $n > 30$ ) random sample from an **arbitrary** population with mean  $\mu$  and **known** variance  $\sigma^2$ :

$$\bar{X} \sim \mathcal{N}\left(\bar{x}; \mu, \frac{\sigma^2}{n}\right) \Rightarrow Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(z; 0, 1)$$

2. t-distribution: In a **small** ( $n < 30$ ) random sample from a **normal** population with mean  $\mu$  and **unknown** variance  $\sigma^2$ :

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim \mathcal{T}(t; n-1)$$

3. Chi-distribution: In a random sample (size does not matter) from a **normal** population:

$$\chi^2 = \frac{(n-1)S^2}{\sigma^2} \sim \mathcal{C}(\chi^2, n-1)$$

The first step of all statistical problems is always **choosing our test** (as tool). Now which one should we choose?

- Since we know nothing about the distribution of our population  $f(x)$ , so ensure the accuracy, Z-test is the best option; but our sample size is very small (just  $n = 10$  seeds per dish), and the population variance  $\sigma^2$  is a mystery.
- t-test is a potential choice, since it does not require  $\sigma^2$  value; but the trade-off is the population must be **normally distributed**.
- Chi-test is a very powerful tool too; but the trade-off is still remaining, the population must be **normally distributed**.

We hit the dead end. All of our testing tools have their own pros and cons; and we can not be absolutely certain for our conclusions.

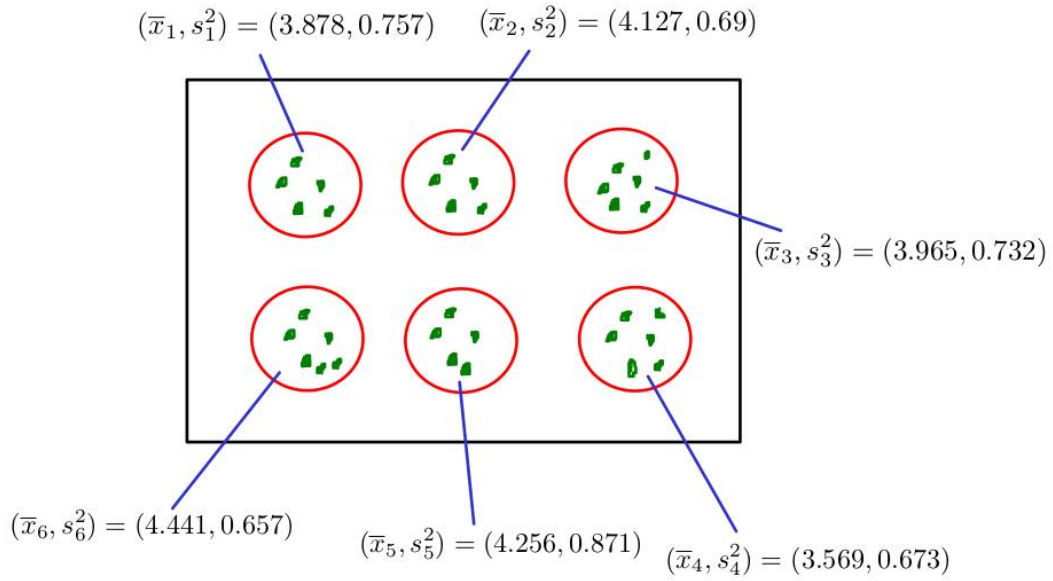


Figure 7.5: Experiment's measured data

Here is our measured data after 5 days conducting experiment:

| $i$         | 1     | 2     | 3     | 4     | 5     | 6     |
|-------------|-------|-------|-------|-------|-------|-------|
| $\bar{x}_i$ | 3.878 | 4.127 | 3.965 | 3.569 | 4.256 | 4.441 |
| $s_i^2$     | 0.757 | 0.69  | 0.732 | 0.673 | 0.871 | 0.657 |

Of course ensuring everything correct is very **impossible**, so **trading-off** our unknown distribution  $f(x)$  with normal assumption is a promising idea; and we can determine the **interval** of  $\mu$  or  $\sigma^2$  may be fallen in with up to 95% certainty. As I mentioned before, because of our small sample size, we can not use Z-test; t-test and Chi-test are the better choices in our situation.

#### 7.4.1 Estimating the Mean using t-test

The key idea of using test statistics is we know exactly the **interval** where 95% of the  $Z$ ,  $t$  or  $\chi^2$  values are inside. Let me explain via the graph of t-distribution:

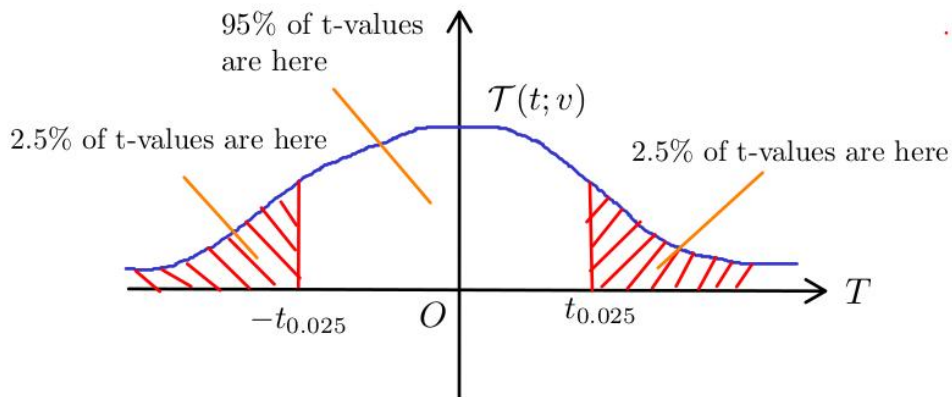


Figure 7.6: t-distribution with 95% confidence interval

Let  $t_\alpha$  represent the  $t$  value above which we find an area of  $\alpha$ . Similarly to others test statistics,



same meaning can also be applied for  $z_\alpha$  and  $\chi_\alpha^2$ . So  $t_{0.025}$  and  $-t_{0.025}$  are the values that we find the **red area** above is 0.025 (because  $\mathcal{T}(t; v)$  graph is symmetrical). The **white area** bounded by two ends is 0.95:

$$-t_{0.025} < T < t_{0.025}$$

Intuitively, you can see that 95% of t-values focusing inside the bounded **white area**. Back to our problem, now you obtain:

$$-t_{0.025} < T < t_{0.025} \Rightarrow -t_{0.025} < \frac{\bar{X} - \mu}{S/\sqrt{n}} < t_{0.025} \Rightarrow \bar{X} - t_{0.025} \frac{S}{\sqrt{n}} < \mu < \bar{X} + t_{0.025} \frac{S}{\sqrt{n}}$$

After substituting the actual values of  $\bar{x}_i$  and  $s_i$ , you can find the interval of  $\mu$  with up to 95% certainty. It means you can confidently state that the **actual**  $\mu$  value lies somewhere inside our interval with a 5% chance of error. The final problem is determining  $t_{0.025}$  value. In practice, instead of using t-distribution formula, we always use t-distribution table for simplicity, which is **Appendix B**. When searching for desired  $t$  value, you must check the parameter  $v$  degrees of freedom first to ensure the accuracy. As you can see here, our sample size is  $n = 10$ , so:

$$v = n - 1 = 10 - 1 = 9 \text{ (degrees of freedom)}$$

Searching for  $t_{0.025}$  value with  $v = 9$  yields  $t_{0.025} = 2.262$ . Finally we conclude:

$$\bar{x}_i - t_{0.025} \frac{s_i}{\sqrt{n}} < \mu < \bar{x}_i + t_{0.025} \frac{s_i}{\sqrt{n}}$$

Substituting each pairs  $(\bar{x}_i, s_i)$  yields 6 intervals:

$$\mu \in (3.68, 4.07)$$

$$\mu \in (3.93, 4.31)$$

$$\mu \in (3.77, 4.15)$$

$$\mu \in (3.38, 3.75)$$

$$\mu \in (4.04, 4.46)$$

$$\mu \in (4.25, 4.62)$$

Finally, we take average of all the starting and ending points of the intervals:

$$\mu \in (3.841, 4.226) \text{ (95\% certainty)}$$

We can conclude that 95% of the seeds germinate within the range of (3.841, 4.226) day.

## 7.4.2 Estimating the Variance using Chi-test

Using similar line of thinking, we have to find the interval that 95% of  $\chi^2$  values fall inside. The only difference is  $\mathcal{C}(\chi^2; v)$  is asymmetrical, so you must search for both of  $\chi_{0.975}^2$  and  $\chi_{0.025}^2$  in chi-squared distribution table, which is **Appendix C**.

$$\chi_{0.975}^2 < \chi^2 < \chi_{0.025}^2 \Rightarrow \chi_{0.975}^2 < \frac{S^2(n-1)}{\sigma^2} < \chi_{0.025}^2 \Rightarrow \frac{S^2(n-1)}{\chi_{0.025}^2} < \sigma^2 < \frac{S^2(n-1)}{\chi_{0.975}^2}$$

With  $v = 9$  degrees of freedom,  $\chi_{0.025}^2 = 19.023$  and  $\chi_{0.975}^2 = 2.7$ . Finally we conclude:

$$\frac{s_i^2(n-1)}{\chi_{0.025}^2} < \sigma^2 < \frac{s_i^2(n-1)}{\chi_{0.975}^2}$$

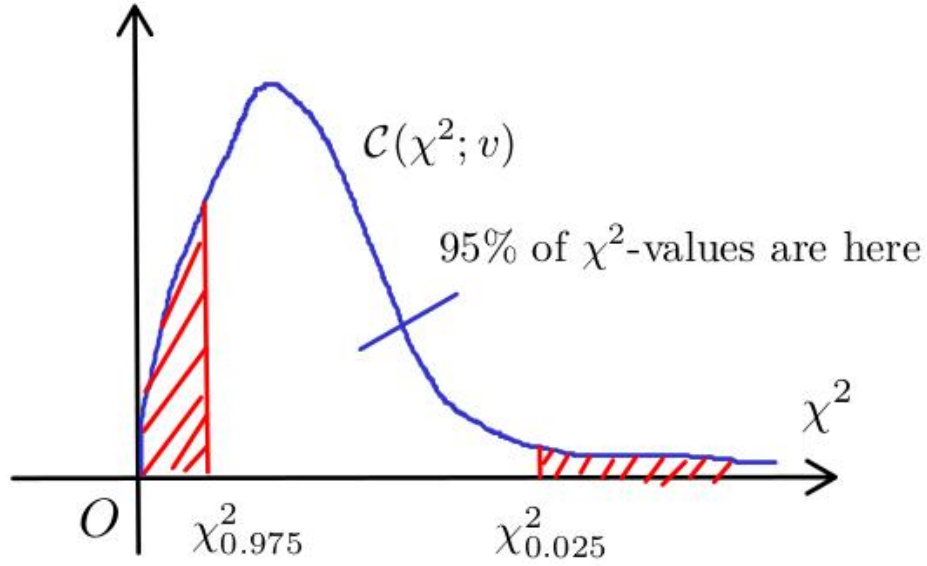


Figure 7.7: Chi-Squared distribution with 95% confidence interval

Substituting each  $s_i$  values yields 6 interval:

$$\sigma^2 \in (0.35, 2.52)$$

$$\sigma^2 \in (0.32, 2.30)$$

$$\sigma^2 \in (0.34, 2.44)$$

$$\sigma^2 \in (0.31, 2.24)$$

$$\sigma^2 \in (0.41, 2.90)$$

$$\sigma^2 \in (0.31, 2.19)$$

We take average of all starting points and ending points of the intervals:

$$\sigma^2 \in (0.34, 2.431) \text{ (95\% certainty)}$$

Finally, we can conclude if  $f(x)$  is **normally distributed**, then:

$$f(x) \sim \mathcal{N}(x; \mu, \sigma^2)$$

with  $\mu \in (3.841, 4.226)$  and  $\sigma^2 \in (0.34, 2.431)$  with 95% certainty.

# Chapter 8

## Classical Methods of Estimation

### 8.1 Definition of Unbiased Estimator

In the previous chapter, we used **sample mean**  $\bar{X}$  to estimate **population mean**  $\mu$ ; and **sample variance**  $S^2$  to estimate **population variance**  $\sigma^2$ . Both  $\bar{X}$  and  $S^2$  are **unbiased estimator**; intuitively, now we can form the definition of **unbiased estimator** as follows:

**Definition 8.1.1.** A statistic  $\hat{\Theta}$  is said to be an **unbiased estimator** of the parameter  $\theta$  if:

$$\mu_{\hat{\Theta}} = E(\hat{\Theta}) = \theta$$

Mathematically, in a sample set, there are many quantites like **sample mean**  $\bar{X}$ , **sample variance**  $S^2$ , **proportion**  $\hat{P}, \dots$  and logically, we can not arbitrarily use  $\bar{X}$  to estimate  $\sigma^2$ . The definition above acts like "key-lock" mechanism to guarantee there will be no ambiguity here.

**Theorem 8.1.1.**  $\bar{X}$  is an **unbiased estimator** of the parameter  $\mu$ .

*Proof.* From CLT, we know:

$$\bar{X} \sim \mathcal{N}\left(\bar{x}; \mu, \frac{\sigma^2}{n}\right)$$

From the definition of **unbiased estimator**, we have to show that a statistic  $\bar{X}$  has its mean equals to  $\mu$ :

$$\mu_{\bar{X}} = E(\bar{X}) = \mu$$

Obviously, it is true. □

**Theorem 8.1.2.**  $S^2$  is an **unbiased estimator** of the parameter  $\sigma^2$ .

*Proof.* From the definition of **sample variance**, we have

$$\begin{aligned} E(S^2) &= \frac{1}{n-1} E\left(\sum_{i=1}^n (X_i - \bar{X})^2\right) \\ &= \frac{1}{n-1} E\left(\sum_{i=1}^n (X_i - \mu + \mu - \bar{X})^2\right) \\ &= \frac{1}{n-1} E\left(\sum_{i=1}^n (X_i - \mu)^2 + 2\sum_{i=1}^n (X_i - \mu)(\mu - \bar{X}) + \sum_{i=1}^n (\mu - \bar{X})^2\right) \\ &= \frac{1}{n-1} E\left(\sum_{i=1}^n (X_i - \mu)^2 + 2(\mu - \bar{X})(n\bar{X} - n\mu) + n(\mu - \bar{X})^2\right) \\ &= \frac{1}{n-1} E\left(\sum_{i=1}^n (X_i - \mu)^2 - n(\mu - \bar{X})^2\right) \end{aligned}$$

Since this sum has chi-squared distribution with  $n$  degrees of freedom, then:

$$E \left( \sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2 \right) = n \Rightarrow E \left( \sum_{i=1}^n (X_i - \mu)^2 \right) = n\sigma^2$$

And this sum has chi-squared distribution with 1 degree of freedom:

$$E \left( \left( \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right)^2 \right) = 1 \Rightarrow E((\bar{X} - \mu)^2) = \frac{\sigma^2}{n}$$

Finally we have:

$$\mu_{S^2} = E(S^2) = \frac{1}{n-1}(n\sigma^2 - \sigma^2) = \sigma^2$$

□

Now you understand **why we divide by  $n - 1$  rather than  $n$**  in sample variance formula to guarantee  $S^2$  is an unbiased estimator. In this chapter, we define a new statistic quantity  $\hat{P}$  (proportion) and prove that it is an unbiased estimator.

**Definition 8.1.2.** A statistic  $\hat{P}$  is defined as:

$$\hat{P} = \frac{X}{n}$$

where  $X$  is a **Bernoulli random variable**.

You can refer to **Chapter 0.1: Coin Tossing** to read a typical example of the quantity  $\hat{P}$  in a practical experiment. In **Statistics**, the terms "probability" and "proportion" can be used interchangeably, depending on the context.

**Theorem 8.1.3.**  $\hat{P}$  is an **unbiased estimator** of the parameter  $p$ .

*Proof.* Because  $X \sim \mathcal{B}(x; n, p)$  we have:

$$E(\hat{P}) = E \left( \frac{X}{n} \right) = \frac{1}{n}E(X) = \frac{np}{n} = p$$

□

It is very easy to indicate that:

$$\sigma_{\hat{P}}^2 = \frac{1}{n^2}\sigma_X^2 = \frac{npq}{n^2} = \frac{pq}{n}$$

**Corollary 8.1.3.1.** The statistic  $\hat{P}$  also has **Bernoulli distribution** with mean and variance:

$$\begin{aligned} E(\hat{P}) &= p \\ \sigma_{\hat{P}}^2 &= \frac{pq}{n} \end{aligned}$$

*Proof.* Obviously the random variable  $\hat{P}$  has Bernoulli distribution since it is just a scaled version of  $X$ . □

**Corollary 8.1.3.2.** If  $n \rightarrow +\infty$ , then:

$$Z = \frac{\hat{P} - p}{\sqrt{pq/n}} \sim \mathcal{N}(z; 0, 1)$$

The most confusing point you need to be aware of in this chapter is the correct use of these notations  $(\hat{\Theta}, \hat{\theta}, \theta)$ :

- $\hat{\Theta}$ : Statistic.
- $\hat{\theta}$ : A single value of Statistic.
- $\theta$ : Parameter of population need to be estimated.

So far we have just found 3 typical examples of these triples  $(\hat{\Theta}, \hat{\theta}, \theta)$  are  $(\bar{X}, \bar{x}, \mu)$ ,  $(S^2, s^2, \sigma^2)$  and  $(\hat{P}, \hat{p}, p)$  are **unbiased** estimator. In the final chapter of this book, we will prove 2 more triples sharing the same characteristics are  $(B_0, b_0, \beta_0)$  and  $(B_1, b_1, \beta_1)$  also **unbiased** too.

## 8.2 Determining Statistical Intervals Like A Pro: Step By Step

### 8.2.1 Selecting Your Test Statistic

The very first step of all statistical problems is **selecting test statistic**. To choose the most appropriate test you must deeply understand their uses and **trade-off**. Remember nothing can be absolutely perfect!

- If your sample size is relatively large ( $n > 30$ ) and you want to estimate the **mean**  $\mu$ , choose Z-test without any doubts.
- If your sample size is small ( $n < 30$ ) and you want to estimate the **mean**  $\mu$ , choose t-test.
- Sample size does not matter when you want to estimate the **variance**  $\sigma^2$ , choose Chi-test.
- If you want to estimate the **proportion**  $p$ , your sample size must be **large**, and choose Z-test.

Since all of our testing tools are based on the **normal population** assumption (except CLT), so you have to be careful when you do not know how your population is distributed. A very common approach is assuming that the entire population follows a **normal distribution**, sacrificing the precision of  $f(x)$  to simplify our problem. From this point onward in this book, **all populations are assumed to be normal**.

After choosing your desired test among 3 tests, write down your testing formula based on the distribution of the quantity that you want to estimate, for example:

$$\begin{aligned}\bar{X} &\sim \mathcal{N}\left(\bar{x}; \mu, \frac{\sigma^2}{n}\right) \Rightarrow Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(z; 0, 1) \\ (\bar{X}_1 - \bar{X}_2) &\sim \mathcal{N}\left(\bar{x}_1 - \bar{x}_2; \mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right) \Rightarrow Z = \frac{\bar{X}_1 - \bar{X}_2 - (\mu_1 - \mu_2)}{\sqrt{(\sigma_1^2/n_1) + (\sigma_2^2/n_2)}} \sim \mathcal{N}(z; 0, 1) \\ \frac{S^2(n-1)}{\sigma^2} &= \chi^2 \sim \mathcal{C}(\chi^2; n-1) \\ \frac{\bar{X} - \mu}{S/\sqrt{n}} &= T \sim \mathcal{T}(t; n-1)\end{aligned}$$

$Z$ ,  $t$  and  $\chi^2$  are 3 new **normalized estimators**  $\hat{\Theta}$ , and we use them to determine confidence interval.

## 8.2.2 Establishing The Confidence Interval

### Two-Sided Confidence Bounds

The idea of establishing the confidence interval is very simple and intuitive, we did it (but still do not give it a name) in the previous chapter. Before learning the methods, let me introduce some general notations first:

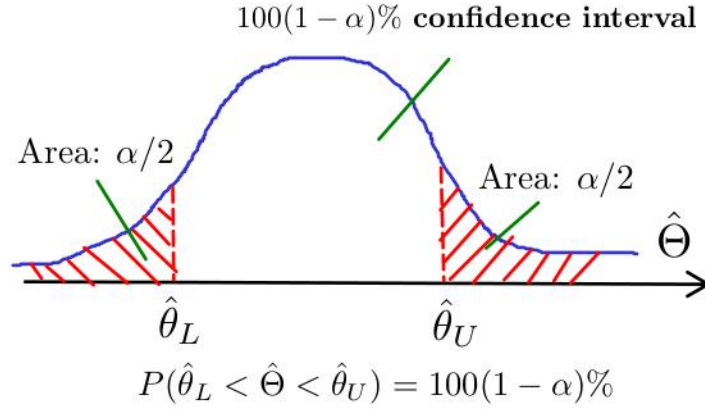


Figure 8.1: Two-sided confidence bounds

$\hat{\theta}_L$  and  $\hat{\theta}_U$  stand for the lower and upper values of our boundaries;  $100(1 - \alpha)\%$  is the **confidence interval**, where **exact**  $100(1 - \alpha)\%$  values of  $\theta$  inside. Obviously you can see that the value of  $\hat{\theta}$  must be chosen to satisfies the white area it bounded is  $1 - \alpha$ . Because standard normal and t-distribution curves are both symmetrical, so we often choose:

$$\hat{\theta}_L = -t_{\alpha/2}, \quad \hat{\theta}_U = t_{\alpha/2} \quad (\text{for t-test})$$

$$\hat{\theta}_L = -z_{\alpha/2}, \quad \hat{\theta}_U = z_{\alpha/2} \quad (\text{for Z-test})$$

But chi-distribution curve is asymmetrical, you have to be careful here:

$$\hat{\theta}_L = \chi^2_{1-\alpha/2}, \quad \hat{\theta}_U = \chi^2_{\alpha/2} \quad (\text{for Chi-test})$$

General interval has been established with  $(1 - \alpha)\%$  certainty:

$$\hat{\theta}_L < \hat{\theta} < \hat{\theta}_U$$

Substituting everything down here yields your final interval, for example:

$$-t_{\alpha/2} < T < t_{\alpha/2} \Rightarrow -t_{\alpha/2} < \frac{\bar{X} - \mu}{s/\sqrt{n}} < t_{\alpha/2} \Rightarrow \bar{X} - t_{\alpha/2} \frac{s}{\sqrt{n}} < \mu < \bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}}$$

From your actual experiment data, set the random variable  $\bar{X}$  be one of your result:

$$\bar{x} - t_{\alpha/2} \frac{s}{\sqrt{n}} < \mu < \bar{x} + t_{\alpha/2} \frac{s}{\sqrt{n}}$$

Using the same line of thinking, you can find every two-sided confidence bounds intervals, with an arbitrary  $\alpha$  value. In practice, we often choose  $\alpha = 0.01, 0.05$  or even  $\alpha = 0.1$ . I do not want to give too many formulas here, since they share the same logic, and you can derive them

yourself easily. Before moving on to the next part, you should pay attention to this subtle formula:

$$-z_{\alpha/2} < Z < z_{\alpha/2} \Rightarrow -z_{\alpha/2} < \frac{\hat{P} - p}{\sqrt{pq/n}} < z_{\alpha/2} \Rightarrow \hat{P} - z_{\alpha/2} \sqrt{\frac{pq}{n}} < p < \hat{P} + z_{\alpha/2} \sqrt{\frac{pq}{n}}$$

Directly estimating  $p$  seems very hard now, because it appears in both sides, but is not impossible (you should try). Since  $n \rightarrow +\infty$  so we can approx  $\hat{p} \approx p$ , then:

$$\hat{P} - z_{\alpha/2} \sqrt{\frac{\hat{p}\hat{q}}{n}} < p < \hat{P} + z_{\alpha/2} \sqrt{\frac{\hat{p}\hat{q}}{n}}$$