

# COMPSCI 70 Self-Study Lecture Notes

Dun-Ming Huang

# Contents

<b>0</b>	<b>Review of Sets and Mathematical Notation</b>	<b>6</b>
0.1	Set Theory . . . . .	6
0.2	Mathematical Notation . . . . .	8
0.3	Personal Regard on this Topic . . . . .	10
<b>1</b>	<b>Propositional Logic</b>	<b>11</b>
1.1	Propositional Logic . . . . .	11
1.2	Quantifiers . . . . .	13
1.3	Negation . . . . .	15
1.4	Personal Thoughts on This Topic . . . . .	17
<b>2</b>	<b>Forms of Proofs</b>	<b>18</b>
2.1	Introduction to Proofs . . . . .	18
2.2	Proof Technique . . . . .	19
2.3	Personal Regard on This Topic . . . . .	24
<b>3</b>	<b>Mathematical Induction</b>	<b>25</b>
3.1	Introduction to Mathematical Induction . . . . .	25
3.2	Simple Induction and Strong Induction . . . . .	28
3.3	Build-Up Error . . . . .	29
3.4	Personal Regard on This Topic . . . . .	30
<b>4</b>	<b>The Stable Matching Problem</b>	<b>31</b>
4.1	Description of Stable Matching Problem . . . . .	31
4.2	The Properties of Propose-And-Reject Algorithm . . . . .	31
4.3	Review Questions from Midterm . . . . .	35
4.4	Personal Regard on This Topic . . . . .	36
<b>5</b>	<b>Graph Theory</b>	<b>37</b>
5.1	An Introduction to Graph Theory . . . . .	37
5.2	Koenigsberg's Seven Bridges and Eulerian Tour . . . . .	40

5.3	Planarity, Euler's Formula, Coloring . . . . .	42
5.4	Classes of Graphs . . . . .	45
5.5	Personal Regard on This Topic . . . . .	47
<b>6</b>	<b>Modular Arithmetics</b>	<b>48</b>
6.1	Introduction to Modular Arithmetics . . . . .	48
6.2	Exponentiation . . . . .	49
6.3	Bijections . . . . .	50
6.4	Inverse Operations . . . . .	51
6.5	Chinese Remainder Theorem . . . . .	54
6.6	Personal Regard on This Topic . . . . .	56
<b>7</b>	<b>Public Key Cryptography</b>	<b>57</b>
7.1	Introduction to Cryptography . . . . .	57
7.2	Cryptography in Modular Arithmetics Theorems . . . . .	58
7.3	My Personal Regard on This Topic . . . . .	60
<b>8</b>	<b>Polynomials</b>	<b>61</b>
8.1	Re-Introduction to Polynomial . . . . .	61
8.2	Polynomial Interpolation . . . . .	61
8.3	Re-Introducing Polynomial Division . . . . .	62
8.4	Interlude: Finite Fields . . . . .	63
8.5	Secret Sharing . . . . .	64
8.6	Personal Regard on This Topic . . . . .	64
<b>9</b>	<b>Error Correcting Codes</b>	<b>65</b>
9.1	Erasure Errors . . . . .	65
9.2	Polynomial Interpolation in Error Correction . . . . .	66
9.3	General Errors . . . . .	66
9.4	Personal Regard on This Topic . . . . .	68
<b>10</b>	<b>Counting</b>	<b>69</b>
10.1	Sampling without Replacement . . . . .	69
10.2	Sampling with Replacement . . . . .	71
10.3	Combinatorial Proofs . . . . .	72
10.4	The Principle of Inclusion-Exclusion . . . . .	75
10.5	Personal Regards on this Topic . . . . .	76
<b>11</b>	<b>Countability</b>	<b>78</b>
11.1	More on Bijections . . . . .	78
11.2	Countability . . . . .	80

11.3 Cantor's Diagonalization . . . . .	81
11.4 Cantor's Set . . . . .	82
11.5 Power Sets: Higher Order Infinity . . . . .	83
<b>12 Computability</b>	<b>84</b>
12.1 The Liar's Paradox . . . . .	84
12.2 Self Replicating Programs . . . . .	84
12.3 The Halting Problem . . . . .	85
12.4 Godel's Incompleteness Theorem . . . . .	87
<b>13 Introduction to Discrete Probability</b>	<b>88</b>
13.1 Random Experiments . . . . .	88
13.2 Defining Probability Space and Events . . . . .	89
13.3 Examples after Examples of Discrete Probability . . . . .	89
13.4 Summary . . . . .	91
<b>14 Conditional Probability, Independence, Combination of Events</b>	<b>92</b>
14.1 Conditional Probability . . . . .	92
14.2 Bayesian Inferenece . . . . .	93
14.3 Bayes Rule and Total Probability Rule . . . . .	94
14.4 Combination of Events . . . . .	96
<b>15 Random Variables: Distribution and Expectation</b>	<b>99</b>
15.1 Introduction to Random Variables . . . . .	99
15.2 Probability Distribution . . . . .	99
15.3 Multiple Random Variables and Independence . . . . .	101
15.4 Expectation . . . . .	102
<b>16 Variance</b>	<b>103</b>
16.1 Functions of Random Variables . . . . .	103
16.2 Introducing and Defining Variance . . . . .	103
16.3 Sum of Independent Random Variables . . . . .	105
16.4 Covariance and Correlation . . . . .	105
<b>17 Concentration Inequalities and the Laws of Large Numbers</b>	<b>107</b>
17.1 Introduction . . . . .	107
17.2 Markov's Inequality . . . . .	108
17.3 Chebyshev's Inequality . . . . .	109
17.4 The Law of Large Numbers . . . . .	109
<b>18 Geometric and Poisson Distributions</b>	<b>111</b>

18.1 Geometric Distribution . . . . .	111
18.2 Poisson Distribution . . . . .	115
<b>19 Prediction!</b>	<b>118</b>
19.1 Joint Distributions, Conditional Expectation . . . . .	118
19.2 Iterated Expectation and Wald's Identity . . . . .	119
19.3 Minimal Mean Square Error? . . . . .	119
19.4 Linear Regression . . . . .	119
<b>20 Continuous Probability Distributions</b>	<b>121</b>
20.1 Continuous Random Variables . . . . .	121
20.2 Exponential Distribution . . . . .	124
20.3 Normal Distribution . . . . .	125
20.4 Central Limit Theorem . . . . .	128

# Chapter 0

## Review of Sets and Mathematical Notation

In this chapter, we provide a brief review of fundamental mathematical notations that will be used throughout the coursework.

This involves set notation, set properties, as well as convenient mathematical notations that aid to simplify, abbreviate, as well as clarify the different conditions and clauses that we used to literally write out in mathematical statements.

### 0.1 Set Theory

This chapter describes the fundamental properties of mathematical objects called sets.

Sets are widely applied throughout and beyond the subfields of mathematics, whose origins and characteristics shared across those applications are described by the mathematical definition of sets.

#### 0.1.1 Fundamental Properties of Sets

As you may have guessed, in set theory, we mainly discuss a mathematical object called "set".

##### Definition 0.1.1. Set

A set is a collection of objects.

The mathematical objects that are members of a set are called elements.

Sets are normally expressed in the format of:

$$\{\text{element 1, element 2, } \dots, \text{element n}\}$$

The symbol representing membership in a set works as follows:

##### Symbol 0.1.1. Membership Symbol

If a mathematical object  $x$  is a member of a set  $A$ , then we may mathematically express that  $x \in A$ .

Otherwise, for a mathematical object  $y$  that does not belong to  $A$ ,  $y \notin A$ .

There are several more properties to sets:

- The size of a set is defined as the number of elements in it. This property of set is called **cardinality**. The cardinality of a set  $A$  can be mathematically denoted as  $|A|$ .
- The equality of a set is held when two sets have exactly same elements; order and repetition does not matter, albeit in computers, many implementations of sets already do not allow repetition. We call a set whose elements can be repeated a “multiset”.
- Sets whose cardinality are 0, also known as empty sets, are denoted as  $\{\}$  as well as  $\emptyset$ .

To denote a set of mathematical objects that all belong to some other set but suffices some property, we may utilize a notation called **set-builder notation**:

#### Symbol 0.1.2. Set-Builder Notation

A set  $A$  whose members all follow expression  $exp$ , and the terms of  $exp$  all fit some list of conditions can be written as:

$$\{exp : \text{condition 1}, \dots, \text{condition n}\}$$

The colon above can be replaced by a vertical bar. Since we will be using vertical bars for other purposes in later chapters, a colon is perhaps the least abusive notation to play with.

But, to conform the formatting of official lecture notes, let us continue with vertical bars.

For example, the list of all rational numbers can be written as:

$$\mathbb{Q} = \left\{ \frac{p}{q} \mid p, q \in \mathbb{R}, q \neq 0 \right\}$$

## 0.1.2 Relationships between Sets

We often compare sets in terms of their sizes and elements they contain. Out of these standards for comparison, we can define the relationship between a set and another larger set that contains all elements of the former as follows:

#### Definition 0.1.2. Subset

If every element of a set  $A$  is also in a set  $B$ , then  $A$  is a subset of  $B$ .

Mathematically, we write it as  $A \subseteq B$ .

Or, stating the equivalent in an opposite direction,  $B \supseteq A$ , and this would state  $B$  as a superset of  $A$ .

And a stricter similar relationship follows:

#### Definition 0.1.3. Strict Subset

If  $A \subseteq B$  but  $A$  excludes at least an element of  $B$ , then we say that  $A$  is a strict subset (proper subset) of  $B$ .

Mathematically denoted,  $A \subset B$ .

Utilizing the definitions above, we may form some observations:

- The empty set is a proper subset of any nonempty set.
- The empty set is a subset for any set, including itself's.
- While every set is not a proper subset of itself, every set is a subset of itself.

While we compare sets based on their members, we can also "add", create sets based on the members of two sets. There are two ways of doing so, being **intersection** and **union**.

**Definition 0.1.4. Intersection**

The intersection of set  $A$  and  $B$ , written as  $A \cap B$ , is the set containing all elements that are both in  $A$  and  $B$ .  
In set builder notation:

$$A \cap B = \{x | x \in A \wedge x \in B\}$$

**Definition 0.1.5. Union**

The union of set  $A$  and  $B$ , written as  $A \cup B$ , is the set of all elements contained in either  $A$  or  $B$ .  
In set builder notation, it pronounces very similarly with the notation for intersections:

$$A \cup B = \{x | x \in A \vee x \in B\}$$

Regarding the way empty sets work in the above arithmetic, for an arbitrary set  $A$ :

- $A \cap \emptyset = \emptyset$
- $A \cup \emptyset = A$

At last, let us regard the notion of “**complements**”, the difference between sets.

**Definition 0.1.6. Complement**

Imagine the arithmetic of sets to work solely upon their members, and let the difference of sets be defined such that:

$$A - B = A \setminus B = \{x \in A | x \notin B\}$$

This set  $A - B$  is called the set difference between  $A$  and  $B$ , or alternatively, the relative complement of  $B$  in  $A$ .

This also indicates that the “subtraction arithmetic” for sets is not commutative. Using empty sets as examples:

- $A \setminus \emptyset = A$
- $\emptyset \setminus A = \emptyset$

And the last arithmetic is the “**multiplication of sets**”: **Cartesian Products**.

**Definition 0.1.7. Cartesian Products**

The Cartesian Product (cross product) of sets  $A$  and  $B$  is defined such that:

$$A \times B = \{(a, b) | a \in A, b \in B\}$$

And last but not least, we have mathematical operations that generate a set of sets (nested) based on other sets:

**Definition 0.1.8. Power Set**

The power set of  $S$  can be written as:

$$\mathcal{P}(S) = \{P | P \subseteq S\}$$

as one of its various denotations, the one on lecture note using a Weierstrass  $\mathcal{P}$ .

The power set of a set is the set of all of its possible subsets.

In addition, if the cardinality of  $S$  in the above definition is  $|S| = k$ , then we can state that  $|\mathcal{P}(S)| = 2^k$ . This will be an explored notion as we discuss a discrete math topic “Counting” in near future.

## 0.2 Mathematical Notation

Here we describe and explore mathematical notations that will abbreviate equations, statements, as well as facilitate us to view mathematical expressions in different perspectives. Mathematical notations introduced in this section will



be very frequently employed in coming courseworks.

### 0.2.1 Summation and Products

The summation of some expression dependent on some other variable can be stated with a summation symbol, such that:

#### Symbol 0.2.1. Summation

Let  $f(x)$  be an expression based on some input  $x$ , then the following expressions are equivalent:

$$f(m) + f(m+1) + \cdots + f(m+n) = \sum_{k=1}^n f(m+k)$$

While the summation expression is highly applicable to numerous theorems, there is another notation similar in nature to the summation expression– the product expression:

#### Symbol 0.2.2. Product

Let  $f(x)$  be an expression based on some input  $x$ , then the following expressions are equivalent:

$$f(m) \times f(m+1) \times \cdots \times f(m+n) = \prod_{k=1}^n f(m+k)$$

### 0.2.2 Quantifiers

Quantifiers are symbols that help us target a range of elements from some set when we need to write a conditional statement for it. There are two specific quantifiers we will use here.

#### Symbol 0.2.3. Universal Quantifier

The universal quantifier  $\forall$ , pronounced as "for all", targets all elements from a set.

For example, the phrase " $\forall x \in A$ " is equivalent to the phrase "for all objects  $x$  that belongs to the set  $A$ ". This can serve as the first clause of a statement, while the second clause would refer a statement for the elements targeted by the quantifier.

Meanwhile, we have a different quantifier that selects a different range from the universal quantifier:

#### Symbol 0.2.4. Existential Quantifier

The existential quantifier  $\exists$ , pronounced as "there exists (at least one)", states that there should exist at least one element from some set that fits a condition stated in the second clause. The syntax, therefore, is as follows:

$$\exists x \in A : \text{statement}$$

This means: there exists at least an object  $x$  belonging to the set  $A$  such that the "statement" holds. There are some more variants to this overall syntax, such as replacing  $\in$  with  $\notin$  etc.

There also exists some variations to the existential quantifier:

- There exists at least one:  $\exists$
- There exists only one:  $\exists!$
- There does not exist:  $\nexists$

### 0.3 Personal Regard on this Topic

Personal log for mathematical topics, because why not! Just don't treat this section as educational content.

Sets are very useful objects for mathematical generalizations, which we'd do a TON in discrete mathematics, hence its practicality. There will, however, come times where statements are hard to write out.

In that case, don't bother with making the statements multi-part! There's a clear trend in homework questions and solutions for that.

Also, my reader didn't grade one of my homework questions and I missed my regrade deadline :v

# Chapter 1

## Propositional Logic

This chapter introduces the fundamentals of propositional logic while enhancing the student's understanding of quantifiers, as well as discussing popular uses of propositional operations as implications and negations.

### 1.1 Propositional Logic

#### 1.1.1 Definition of Proposition

To speak the language of mathematics, we must familiarize ourselves with objects of the mathematical language. One fundamental block of this mathematical language is a **proposition**.

##### Definition 1.1.1. Proposition

A proposition is a logical statement that is either true or false.

In other words, it is a mathematical statement that could be correct, but also could be incorrect.

Propositions need to be either true or false, which means it cannot be statements that can yield some ambiguous result. And because they usually are not known to be true or false, we often need to prove propositions in discrete mathematics.

Exercise Reader! Exercise! 1.1.1: What should a proposition be like?

Which of the following qualifies as a proposition?

- (a)  $\sqrt{3}$  is rational.
- (b)  $2 + 2$
- (c) I often never give up or let down

Option (a) is a statement that has to either hold true or false, because this value  $\sqrt{3}$  is either rational or irrational.

Option (b) is a mathematical expression, and there is no true or false aspect to it. It is not a proposition.

Last but not least, option (c) is just a statement, because the word "often" is not properly defined. If propositions are supposed to clearly hold true or false, it should not contain blurrily defined words.

The answer, therefore, is option (a).

### 1.1.2 Logical Operations of Proposition

There are approximately four methods of using preexistent propositions to make new propositions, and we do so to produce complicated propositions from simpler subparts.

These four ways being:

- **Conjunction:**  $P \wedge Q$ , works similarly as the boolean "and".
- **Disjunction:**  $P \vee Q$ , works similarly as the boolean "or".
- **Negation:**  $\neg P$ , works similarly as the boolean "not".
- **Implies:**  $P \implies Q$ , which states a familiar phrase of "if  $P$ , then  $Q$ ". Here, the proposition  $P$  is called a hypothesis, and  $Q$  a conclusion.

To demonstrate the behaviors of these operations, we may use a **truth table**. A truth table example for the operation "Conjunction" is attached below:

Figure 1.1.1. Truth Table for Conjunction

A truth table is a table that points out the result of an operation on statements  $P$  and  $Q$ , for a set of given boolean values that  $P$  and  $Q$  each result in.

$P$	$Q$	$P \wedge Q$
T	T	T
T	F	F
F	T	F
F	F	F

Now we may discuss some other properties of propositions.

First of all, there is a fundamental principle called **Law of the Excluded Middle**, which states what is as follows:

#### Axiom 1.1.1. Law of Excluded Middle

For any proposition  $P$ ,  $P \vee \neg P$  always holds true, because  $P$  is either true or false (not true).

Here, propositions like  $P \vee \neg P$  that always hold true are called **tautologies**.

Meanwhile, propositions that are always false, such as  $P \wedge \neg P$  (which cannot be true, since  $P$  is either true or false), are called **contradictions**.

Implications also have interesting truth table properties:

Figure 1.1.2. Truth Table for Implication

A truth table for implications is as follows:

$P$	$Q$	$P \implies Q$	$\neg P \vee Q$
T	T	T	T
T	F	F	F
F	T	T	T
F	F	T	T

While this provides that the implication is only false when  $P$  is true and  $Q$  is false, this also provides a not so intuitive insight: An implication is trivially true when the hypothesis is false.

This situation is known as "**vacuously true**".

The reason why this insight seems to hold is because while we cannot confirm the conclusion is false when the hypothesis is false, this is also the nature of how if-then statements are expected to be defined.

Indicated via the above truth table, as  $P \implies Q$  and  $\neg P \vee Q$  have the exact same boolean behavior on truth table, they are logically equivalent. This means they are statements that perform the logical behavior and meaning.

Interestingly, there are also various ways to pronounce implications, such as "if P, then Q", as well as "Q if P", and many more.

Sometimes we also come across words that seem to connect two propositions in a tighter implication, such as "iff".

An **iff** relationship is also standardly pronounced as **if and only if**, which is mathematically expressed as  $P \iff Q$ . This "iff" holds if both  $P \implies Q$  and  $Q \implies P$  are true. Consequentially, and as encountered in previous coursework, such as EECS 16A, the proof of an "iff" statement is at its nature the proof of two smaller statements.

### 1.1.3 Implications from Implications

For an implication  $P \implies Q$ , we may define two other implications out of it:

- **Contrapositive:**  $\neg Q \implies \neg P$ , which is logically equivalent to  $P \implies Q$ . At its implication: if a statement  $P$ 's contrapositive is true, so should the statement  $P$  per se.
- **Converse:**  $Q \implies P$ , which does not always hold true. In some cases, converse statements can hold true, but usually accompanies some extra conditions. Some examples can be observed in the contents of MATH 53.

When two propositions are logically equivalent, we may mathematically denote it with the  $\equiv$  symbol.

An example of usage would be stating an implication the logical equivalent of its contrapositive, in the forms:

$$(P \implies Q) \equiv (\neg Q \implies \neg P)$$

An application of this is a new approach of performing proofs. To prove a prompt that states a theorem as an implication  $P \implies Q$ , by proving its contrapositive, which is in some cases easier than proving the prompt directly, we can manage to indirectly prove the prompt.

This application will be introduced in the next chapter as "proof by contrapositive".

## 1.2 Quantifiers

### 1.2.1 Use of Quantifiers

In the previous chapter, we introduced the usage of quantifiers as an abbreviation and selector of elements in sets to which a condition holds, in the format of:

<quantifier selecting clause> <second clause containing condition to hold for selected values>

The second clause is effectively a proposition, while the entire sentence containing the first and second clause is also effectively a proposition.

There is a more sophisticated way to express the first clauses of these forms of propositions.

The first clause has a quantifier, and throughout a "universe" we are working with (abstractly, a wide variety of mathematical objects following some conventions), the statement is quantified (holding true) for a selected range of objects in this universe.

So syntactically, the quantifier serves to involve this range of quantification for our proposition.

Let me expound on this abstraction of "universe" further. Say we attempt to qualify a proposition on all natural numbers, such as:

$$\forall x \in \mathbb{N}, x \in \mathbb{Q}$$

Then, the total set of objects I'm working through is the set of all natural numbers, and is therefore the **"universe": a collection of all objects to be qualified in this proposition**. In that sense, universe is a set.

In a finite universe, there is also another way of conceptualizing the quantifiers we know:

## Exercise Reader! Exercise! 1.2.1: Propositional Logic in Quantifiers

Let us work with a universe  $\mathbb{U} = U_1, \dots, U_n$ , where  $n$  is some finite natural number and elements of this universe are some mathematical object.

Then:

- **Universal Quantifiers:**  $(\forall x \in \mathbb{U}(P(x))) \equiv (P(U_1) \wedge \dots \wedge P(U_n))$
- **Existential Quantifiers:**  $(\exists x \in \mathbb{U}(P(x))) \equiv (P(U_1) \vee \dots \vee P(U_n))$

To express more complicated propositions, we must expand beyond the previous first-clause-second-clause format for propositions. This is especially when we are selecting two ranges of values to describe a multi-variable proposition with.

## Exercise Reader! Exercise! 1.2.2: How to read complicated propositions?

Say we are dealing with the proposition:

$$(\forall x \in \mathbb{Z})(\exists y \in \mathbb{Z})(x < y)$$

Here, I can observe two quantifying clauses and the last clause that stands for a proposition. In fact, propositions involving quantifiers will always have the last clause as a proposition, due to the syntactical limits and customs of mathematical language.

Now let us try translating each clause into the English language:

1.  $\forall x \in \mathbb{Z}$ : For all objects  $x$  that belongs to  $\mathbb{Z}$ , the set of all integers.
2.  $\exists y \in \mathbb{Z}$ : There exists an object  $y$  that belongs to  $\mathbb{Z}$ , the set of all integers.
3.  $x < y$ : Object  $x$  is smaller than  $y$ .

Piecing the clauses together: For all objects  $x$  that belongs to  $\mathbb{Z}$ , the set of all integers, there exists an object  $y$  that belongs to  $\mathbb{Z}$  such that  $x$  is smaller than  $y$ .

This proposition, in fact, holds true, since we can always find a larger integer.

However, the proposition with a reversed order for quantifier clauses:

$$(\exists y \in \mathbb{Z})(\forall x \in \mathbb{Z})(x < y)$$

has to hold False.

With the above translating process, this statement states that: there is an object  $y$  that belongs to  $\mathbb{Z}$ , the set of all integers, such that for all objects  $x$  that belongs to  $\mathbb{Z}$  such that  $x$  is smaller than  $y$ .

In other words, it assumes the existence of one largest integer, which does not exist!

Then, let us look at an example from the Fall 2022 Midterm 1:

## Practice 1.2.1: Reading Complicated Quantifier Statements

Which logical expression(s) below correspond to the predicate that  $p$  is prime?

- a.  $(\forall b \in \mathbb{N})((b \leq 1) \vee \neg(\exists a \in \mathbb{N})(a \neq 0 \wedge ab \neq p))$
- b.  $\neg((\exists a, b \in \mathbb{N})((a \geq 2) \wedge (a < p) \wedge (p = ab)))$

For the expression (a), it states that:

For all natural numbers  $b$ , either  $b \leq 1$  or there does not exist a natural number  $a$  such that  $a$  is nonzero and  $ab \neq p$ .

Therefore, for any positive integer  $b$  larger than 2, there does not exist another positive integer  $a$  such that  $ab \neq p$ .

However, let's consider the possibility  $p = 2$ , which is prime, then for an integer  $b = 2$ , there exists  $a = 1$  such that  $ab = 2$ .

Therefore, the trueness of expression (a) does not indicate that  $p$  is prime.

Now, consider expression (b), which instead states that:

It is not such that there exists two natural numbers  $a, b$  where simultaneously,  $a \geq 2$ ,  $a < p$ , and  $p = ab$ .

To simplify the above quote, there does not exist a pair of natural numbers such that when one of them is a natural number  $2 \leq a < p$ , the other number multiplied by this number is  $p$ .

Since  $a < p$ , to achieve  $ab = p$ , it is impossible that  $b = 1$ . This successfully circumvents the limitation encountered in expression (a), qualifying expression (b) as the correct solution. This expression does successfully describe the definition of prime number  $p$ :

The only pair of natural number whose product is prime  $p$  is  $(1, p)$

By demonstrating that there exist no pair of natural number whose product is  $p$  if one of the element is exclusively between 1 and  $p$ .

This is an interesting discussion, as the question presented above has in fact been regraded!

## 1.3 Negation

### 1.3.1 Meaning of Negation

If a proposition  $P$  is false, its negation is true.

This inspires another approach for proofs (proof by contradiction), which usually work for simpler prompts since their contradictions would only be simpler to prove, if provided that this method is optimal to begin with.

However, when  $P$  appears to be a more complicated proposition involving other smaller propositions, we should look for some way to make our understanding of it more concise.

A logical law, called the **De Morgan's Laws**, facilitates us with such matter:

#### Axiom 1.3.1. De Morgan's Law

The De Morgan's Law states the two following logical equivalences:

$$\neg(P \wedge Q) \equiv (\neg P \vee \neg Q)$$

$$\neg(P \vee Q) \equiv (\neg P \wedge \neg Q)$$

In fact, negations of propositions involving quantifiers follow analogous laws, due to the similarity of universal quantifier with conjunction and existential quantifier with disjunction:

**Axiom 1.3.2. Extension of De Morgan's Law**

Based on the similarity of quantifiers with operations of propositions:

$$\neg(\forall x P(x)) \equiv (\exists x \neg P(x))$$

$$\neg(\exists x P(x)) \equiv (\forall x \neg P(x))$$

These above equivalences can provide flexibility in coming mathematical proofs.

**1.3.2 Complex Examples of Negation**

Let us discuss the example attached on the lecture notes from Summer 2022.

**Prompt 1: Write a proposition that states, "there are at least three distinct integers  $x$  that satisfies  $P(x)$ ."**

And the proposition our textbook provided was:

$$\exists x \exists y \exists z (x \neq y \neq z \wedge P(x) \wedge P(y) \wedge P(z))$$

Exercise Reader! Exercise! 1.3.1: Proposition for Prompt 1

Let us analyze this proposition via some translation:

$$\exists x \exists y \exists z (x \neq y \wedge x \neq z \wedge y \neq z \wedge P(x) \wedge P(y) \wedge P(z))$$

→ "there exists a  $x, y, z$ " such that

" $x, y, z$  are not equal to each other, and all of  $P(x), P(y)$ , and  $P(z)$  holds".

→ "there exists at least one possible combination of three distinct integers" such that "all of  $P(x), P(y)$ , and  $P(z)$  holds".

→ there are at least three distinct integers  $x$  that satisfies  $P(x)$ .

In the above translation process, we first analyzed each phrases independently. Then, having each clauses translated, we move on to combine them together semantically.

For example, the fact that the three integers  $x, y, z$  are not equal to each other means they are distinct, and so we may summarize the quantifying clause from "there exists a ..." into "there exists at least one combination of three integers such that ...".

This helps us imply that, if there exists at least one combination, then there are at least three integers to satisfy  $P(x)$ .

This logic of inter-language and inter-context translation allows us to convert a proposition from its mathematical form to pure English text. It should be familiarized via practice and experience.

For purposes of practicing, let us analyze another proposition and attempt to translate it:

$$\exists x \exists y \exists z \forall d (P(d) \implies d = x \vee d = y \vee d = z)$$

And in the following portion, let us perform again the same procedure of directly translating the proposition from symbols to English, and then rephrase each English clause into more concise descriptions.

Exercise Reader! Exercise! 1.3.2: Explain this proposition for me, will you?

Translate the proposition to English:  $\exists x \exists y \exists z \forall d (P(d) \implies d = x \vee d = y \vee d = z)$

This proposition first has a quantifying clause that states: "there exists a  $x, y$ , and  $z$  for all values  $d$ ", or that there exists a set of three numbers  $x, y$  and  $z$  for all values  $d$ .

The last clause is an implication stating: "if  $P(d)$  holds true, then  $d = x \vee d = y \vee d = z$ ". That means  $d$  is equal



to either  $x$ ,  $y$ , or  $z$ . So, there exists three numbers out of all possible values in the universe such that if  $P(d)$  holds, then  $d$  is one of the three numbers we observed before.

Meanwhile, we are not provided that  $x$ ,  $y$ , and  $z$  are distinct values, so there are at most three values for  $d$  such that  $P(d)$  holds.

Let us switch a perspective and view the contrapositive of the later clause and verify our previous interpretation.

If  $d$  is equal to neither of  $x$ ,  $y$ ,  $z$ , then  $P(d)$  will not hold.

Assuming they are distinct, this statement would provide that if  $d$  is equal to neither of some (up to three) distinct numbers,  $P(d)$  will not hold.

Therefore, the proposition is stating " **$P(d)$  holds for at most three distinct integers**".

Here is one demonstration about conjunctions.

By performing conjunction for the propositions above: that  $P(d)$  holds for at least and at most 3 integers, we successfully create a new proposition that states " $P(d)$  holds for exactly 3 integers".

Mathematical propositions that come with limits, or quantifying clauses, for some other proposition are useful in helping us locate a range of values for which a condition holds.

If we combine propositions that hold for same conditions but with different quantified limits, we managed to combine the limits together and form some tighter statement providing a proposition with tighter limitations.

## 1.4 Personal Thoughts on This Topic

Propositional logic is one aspect of the universal mathematical language that our community will continue to work, especially in fields as artificial intelligence, where many propositions need to hold for some operation(s) to be optimized and conducted correctly. So, we will indeed come across complicated propositions.

The most intuitive way to deal with it is practice translating them into English, and figure how to structurally present them the best in our process!

For now, we are still in the "Let's have a taste of discrete math" stage, and we will have to study just two or three more topics to look at more specific mathematical objects. However, the skills we build here are worthwhile, so make sure to stick around and motivate yourself to read more propositions!

Propositions are essentially methods to address "the rules of a game" in mathematics; in other words, propositions are rules of mathematics that support its framework. It defines mathematics.

This brings upon an interesting question. What if a proposition is not provable? Does this mean that mathematics is in fact not a complete enough system, such that something cannot be determined as right or wrong by mathematics? Is mathematics "incomplete" then?

As you see, propositions indeed bring creative and destructive consequences for mathematics.

# Chapter 2

## Forms of Proofs

This chapter provides an introduction and demonstration for the lines of thoughts behind numerous examples of proofs. The reader will also be exposed to the variety of proof techniques, and be encouraged to explore them via examples provided in the lecture notes.

In this chapter, we concentrate specifically on mathematical proof for its breadth of application and significant role as the foundational concept and tool of understanding and observing the coursework's contents.

Proofs are known to be more complicated concepts, and as well the first barricade for students when studying discrete mathematics. The creativity and "intuition" for generating proofs usually comes with practice and experience with variety of prompts. Readers are encouraged to explore beyond the examples of this lecture note for that reason.

But before we start, let me preface with this. Proofs have learning curves. My advice for you, sincerely, is to have fun and enjoy tackling challenges.

### 2.1 Introduction to Proofs

#### 2.1.1 Definition of Proofs

In many occasions, we aim to provide some explanation towards the trueness and falseness of mathematical propositions, statements. The process of doing so is called "proving", and produces a **mathematical proof**.

##### Definition 2.1.1. Mathematical Proof

A mathematical proof provides a means for guaranteeing the truthfulness of a statement. Concretely, it is a finite sequence of steps that construct an entire logical deduction, such that this deduction proves the statement for the range of values it is quantified over. In some cases, this means guaranteeing a proposition for infinite cases.

Proofs are widely applied not only on the mathematic fields, but also those of computer science.

For instance, we might want to prove the correctness of a program as in following its contract, or even on the famous "halting problem". In that regard, for the variety of fields the concept of "proof" is applied on, it is not necessarily mathematic. Still, it is completely based on logic.

Proofs that can successfully guarantee the trueness of a statement, or provide conclusive evidence for it, helps us be certain of a proposition. Such proof is welcomedly named a **rigorous** proof.

A rigorous proof has some structural convention. Proofs usually begin with **axiom**, or self-evident propositions, and the proof proceeds with a sequence of simple logical steps from axioms to further statements. In some sense, it is human thinking put on paper without any jumps.

### 2.1.2 Basic Notations

There are some basic mathematical notations we have not introduced in previous chapters, but will prove useful and foundational in discrete mathematics, including the symbols of some universes that we have discussed before. It will be listed in the following table:

Symbol 2.1.1. Table of Basic Notations in Proofs

$\mathbb{N}$	Set of all natural numbers
$\mathbb{Z}$	Set of all integers
$\mathbb{Q}$	Set of all rational numbers
$\mathbb{R}$	Set of all real numbers
$\mathbb{C}$	Set of all complex numbers
$a b$	Number $a$ divides number $b$ ( $b = aq$ )
$a := b$	Definition: number $a$ is defined to have value $b$

## 2.2 Proof Technique

### 2.2.1 Direct Proof

Each techniques of proofs will have to do with guaranteeing some form of implication, say  $P(x) \implies Q(x)$ . The main approach of a Direct Proof user is to prove the legitimacy of  $Q(x)$  via assuming  $P(x)$ .

We assume an **arbitrary**, another word for "general, any", value  $x$  such that  $P(x)$  holds. We then write logical deduction that guarantees  $Q(x)$  assuming that  $P(x)$  is true, therefore proving that if  $P(x)$  holds then  $Q(x)$  works.

That above is a very literal description of direct proof, and in fact, might have sound like it doesn't offer much about the supposedly esoteric arts of mathematical proofs. Direct Proof is known as the foundational technique of proofs, and most of the time (as encountered in other courseworks), while the approach of direct proof is concise and simply comprehensible, finding the mathematical rules to support this approach is often difficult. It requires creativity. To demonstrate the Direct Proof Technique, let us consider the following prompt:

Exercise Reader! Exercise! 2.2.1: Example of Direct Proof

Prove that: for any  $a, b, c \in \mathbb{Z}$ , if  $a|b$  and  $a|c$ , then  $a|(b+c)$

The proposition above is quantified over some three integers  $a, b$  and  $c$ . The hypothesis is that  $a$  divides  $b$  and  $b$  divides  $c$ . Let us write the hypothesis in mathematical notation:

$$b = aq_1, c = aq_2, q_1, q_2 \in \mathbb{Z}$$

Let us observe the value of  $(b+c)$ , and see how it could guide us to proving the prompt:

$$b+c = aq_1 + aq_2 = a(q_1 + q_2), (q_1 + q_2) \in \mathbb{Z}$$

Then, since  $(b+c)$  can be written as a product of integer  $a$  with some other integer  $q_1 + q_2$ ,  $a$  does divide  $(b+c)$ . Therefore,  $a|(b+c)$ .

From the above we can observe that direct proof is a very straightforward approach for proving theorems, especially when the logical process guiding from hypothesis to conclusion is clear. This might also explain why direct proof is a very instinctive and popular approach for proofs in the previous courses. From here on, we will utilize direct proof in many other proof techniques as well.

Derivations are often forms of direct proof as well. This means direct proofs can also be algebraically heavier than

seen above. In the above example, we also did not attempt to translate the proposition into mathematical symbols, for the sake of brevity. Therefore, in the example that follows, we will demonstrate a much fuller, complicated example for executing direct proof on a theorem.

Exercise Reader! Exercise! 2.2.2: Show me another example of direct proof

Prove that: For an integer between 0 and 1000 exclusively, if such integer is divisible by 9, then the sum of its digits is also divisible by 9.

Let us organize the proposition into mathematical symbol:

$$(\forall n \in \mathbb{N})(n < 1000) \implies (9|n \implies 9|(\text{sum of digits of } n))$$

And now, let us initiate another logical deduction, assuming that the hypothesis,  $9|n$ , is true:

$$9|n \implies n = 9k, k \in \mathbb{Z}$$

(Let the numeric expression of  $n$  be  $abc$ , then:)

$$\implies 9k = 100a + 10b + c$$

$$\implies 9k = 99a + 9b + (a + b + c)$$

$$\implies (a + b + c) = 9k - 99a - 9b$$

$$\implies (a + b + c) = 9(k - 11a - b)$$

(Since  $k, a, b$  are all integers:)

$$\implies (a + b + c) = 9l, l = k - 11a - b \in \mathbb{Z} \implies 9|(a + b + c)$$

## 2.2.2 Proof by Contraposition

This is the second proof technique we will discuss, called **Proof by Contraposition**. We should first recap on what exactly is a contraposition.

The contraposition to an implication  $P \implies Q$  is another logically equivalent proposition to the original implication, in the form of  $\neg Q \implies \neg P$ . Since the implication is logically equivalent to its contrapositive, by proving the contrapositive statement of the prompt, we have indirectly proved the prompt.

Exercise Reader! Exercise! 2.2.3: Demonstrating the Proof by Contraposition Technique

Prove the theorem: Let  $n$  be a positive integer and  $d|n$ . If  $n$  is odd, then  $d$  is odd.

We will use proof by contraposition to prove this theorem.

Stating the prompt in mathematical symbol:

$$(\forall n \in \mathbb{Z})(d|n) \implies (n \pmod{2} = 1 \implies d \pmod{2} = 1)$$

The contrapositive of this prompt assuming the provided information for  $n$  and  $d$  would be:

$$(\forall n \in \mathbb{Z})(d|n) \implies (d \pmod{2} \neq 1 \implies n \pmod{2} \neq 1)$$

Provided so, in the contrapositive proposition, our hypothesis is  $2|d$ . Meanwhile, by prompt,  $d|n$ , thus for some  $k \in \mathbb{Z}$ ,  $n = dk$ .

Combining the above knowledge:

$$\begin{aligned} n &= d \times k = 2\left(\frac{d}{2} \times k\right) \\ 2|d &\implies \frac{d}{2} \in \mathbb{Z} \\ \left(\frac{d}{2} \times k\right) \in \mathbb{Z} &\implies 2|n \end{aligned}$$

In the above proof, we first attempted to find a contrapositive proposition to the original prompt. After so, we directly apply the Direct Proof technique in our current technique.

In other words, proof by contraposition is just finding the contrapositive proposition to the prompt and using direct proof on it: it is a variant of direct proof.

### 2.2.3 Proof by Contradiction

Proof by Contradiction proves a proposition  $P$  by first assuming its negation,  $\neg P$ . After so, attempt to directly prove that this situation,  $\neg P$ , can lead to both results of  $R$  and  $\neg R$ . In other words, we prove the implications  $\neg P \implies (\neg R \wedge R)$  simultaneously.

This is contradictory.

Furthermore, applying the definition of implications, this provides the insight  $\neg P \implies (\neg R \wedge R) \equiv \text{False}$  (referring to the Law of Excluded Middle). The contrapositive of such proven proposition is  $\text{True} \implies P$ , proving that proposition  $P$  is true.

This process is slightly similar to "disproving": via showing that the negation of the prompt to be proven is a ridiculous, impossible situation to be real.

Let us begin with an example then:

Exercise Reader! Exercise! 2.2.4: Now, an example of Proof by Contraposition Technique

Prove that: There are infinitely many prime numbers.

Now, let us consider the negation of the prompt: there are finitely many prime numbers.

If so, then let us also enumerate these finite amount of (let us say  $k$ ) prime numbers:  $p_1 < p_2 < \dots < p_k$ .

**Step 1: For the sake of Contradiction...**

Let us then have a number:

$$q := p_1 p_2 \dots p_k + 1$$

This number cannot be prime, for the largest prime number is  $p_k < q$ . Since  $q$  is not prime, it then has a prime divisor  $p$  that is among the enumerated list of prime numbers.

**Step 2: However, if we follow the assumption...**

Provided  $p|p_1 p_2 \dots p_k$  and supposedly  $p|p_1 p_2 \dots p_k + 1$ , we reach a proposition that would prove fatal:

$$p|((p_1 p_2 \dots p_k + 1) - (p_1 p_2 \dots p_k))$$

(This comes from another proof in the lecture, stating that  $(a|b \wedge a|c \wedge (b > c)) \implies (a|(b - c))$ )

The above proposition tells us that  $p|1$ , and this requires that  $p \leq 1$ . Since there does not exist such prime divisor  $p$ ,  $q$  has no prime divisor and is therefore prime itself.

We have created a contradiction, stating that if there are finitely many prime numbers, then such defined number  $q$  is prime and not prime.

**Step 3: We have formed a contradiction:**

From this contradiction, we clearly see that there cannot just be finitely many prime numbers.

Rather, as the prompt says, there must be infinitely many prime numbers!

Let us continue the series of demonstration with another example to enhance our understanding towards the process of disproving:

Exercise Reader! Exercise! 2.2.5: Another example for Proof by Contraposition Technique

Prove that:  $\sqrt{2}$  is irrational.

To prove by contradiction, we should first present a contradiction in the proposition " $\sqrt{2}$  is rational". If that above negation of the prompt works true, then the following proposition holds by the definition of rational numbers:

$$\exists p \exists q \left( \sqrt{2} = \frac{p}{q} \wedge p \text{ and } q \text{ are coprime} \right)$$

And therefore, squaring both sides of the later clause for this proposition:

$$2 = \frac{p^2}{q^2}$$

And therefore,  $2q^2 = p^2$ .

Now, let us enter a smaller proof: "For an integer  $a$ , if  $a^2$  is even, so is  $a$  even".

Let us consider the contrapositive of that minor prompt, which would be: "For an integer  $a$ , if  $a$  is odd, so is  $a^2$  odd". Any odd integer, meanwhile, can be considered as some number  $a = 2n + 1$  for another integer  $n$ . Squaring such odd integer:

$$\begin{aligned} a^2 &= (2n + 1)^2 \\ &= 4n^2 + 4n + 1 \\ &= 2(2n^2 + 2n) + 1 \\ (2n^2 + 2n) &\in \mathbb{Z} \\ a^2 \% 2 &= 1 \implies a^2 \text{ is odd} \end{aligned}$$

Having proven this smaller property, let us move back to the original demonstration. The minor proof indicates to us that if  $p^2$  is even, so should  $p$ . Therefore,  $p = 2c$  for some integer  $c$ .

This leads us to  $2q^2 = p^2 = 4c^2$ , which shows  $q^2 = 2c^2$ , and provided from there that  $q^2$  is even, so should  $q$ . If  $p$  and  $q$  are both even, then they are not coprime. Yet, by the definition of irrational number, it is defined that  $p$  and  $q$  are coprime.

There exists a contradiction formed from the hypothesis " $\sqrt{2}$  is rational": it is wrong.

Via proof by contradiction, we have disproved the negation of prompt and in turn proved  $\sqrt{2}$  is irrational.

## 2.2.4 Proof by Cases

The technique of proof by cases is slightly more complex.

For some claim with an unknown variable, let us set finite amount of cases for that unknown variable, without knowing which of the possible cases is true. Then, prove for each possible cases that for either one to be a legitimate proposition, the general prompt still holds.

Therefore, in turn, for any possibility of the unknown variable in our prompt, the prompt holds.

Let us demonstrate this underlying line of thought with an example below:

Exercise Reader! Exercise! 2.2.6: An example for Proof by Cases

Prove that:  $((\exists x \notin \mathbb{Q})(\exists y \notin \mathbb{Q}))(x^y \in \mathbb{Q})$

Let us attempt the proof by cases.

Since the prompt's quantifying is involved with existential quantifiers, showing one set of numbers  $x$  and  $y$  would suffice to prove the prompt.

Let us use the set  $x = y = \sqrt{2}$  to proceed with the prompt. Our cases to prove would then be:

1. **Case 1:**  $\sqrt{2}^{\sqrt{2}} \in \mathbb{Q}$

2. **Case 2:**  $\sqrt{2}^{\sqrt{2}} \notin \mathbb{Q}$

#### Case 1

If this case holds, then the prompt is already proved the case's existence.

#### Case 2

If this case holds, then we may apply a new set:  $x = \sqrt{2}^{\sqrt{2}}, y = \sqrt{2}$ .

However:

$$\begin{aligned} (\sqrt{2}^{\sqrt{2}})^{\sqrt{2}} &= (\sqrt{2})^{\sqrt{2} \times \sqrt{2}} \\ &= (\sqrt{2})^2 = 2 \in \mathbb{Q} \end{aligned}$$

Therefore, we still manage to prove that there exists irrational numbers  $x, y$  such that  $x^y$  is rational, via proving every possible cases of our prompt's unknown variable letting our prompt guaranteed to be true.

Albeit the introduction said "finitely many cases", we in reality are expected to deal with just 2 or 3 cases for this technique. If we manage to encounter a prompt that requires more cases for so, revise the categorization of cases or try some other technique of proofs.

Applying the above claim to a broader horizon, every proof has some techniques to prove by, and some are to be more efficient than others. The selection of method is purely context-based.

It is therefore encouraged of mathematicians to explore some different techniques, and encouraged of students to familiarize with each of the four to provide flexibility during problemsolving.

### 2.2.5 General Advice in Proofs, Bullet Points

- Justification can be stated without proof only if it is correct or will be automatically agreed with by any reader.
- If a step cannot be justified and explained clearly, there exists a jump, and some intermediate steps must be added.
- A subsidiary result useful in a more complex proof is called a lemma. It has to be already proved somewhere else. People also break proofs into proofs of lemmas to provide organization and structure.
- The line between lemma and theorems is not clear cut, but theorems are propositions that usually should be exported to the mathematical world, while lemma exists on a more local standing to just facilitate the proof.
- Still, there are some famous lemmas in the world that serve almost the same roles as theorems. This is why the boundary is said to not be that clear cut.

**Side Note from Future:** You might find another form of proof, **Proof by Construction**, useful as we proceed into constructions and considerations of larger mathematical objects, such as those that can span several dimensions or levels. However, they are structurally very similar to mathematical inductions, which is the topic of discussion for the next chapter.

## 2.3 Personal Regard on This Topic

Proofs are difficult, and we all need to practice it. In fact, most of the problems in computer science are resolved or concluded via proofs.

Why are proofs so difficult? Here are some of my own difficulties and my own methods of resolving it:

- Where do I start on a proof? The prompt looks so abstract! It looks like the Minecraft enchantment texts!

The worst thing you can do to your brain during an exam is making it panick. Minecraft enchantment texts prespective will mystify the prompt, and possibly discourage you from approaching it. At these times, try just listing out every single thing the prompt has provided you. Work with these pieces of puzzles, and paraphrase the prompt mathematically.

- Ho-ho. Are you approaching the proof prompt? Then, what's the next step after writing out all the ideas that trivially exist?

Think about equivalent statements. Let's take *Think! 2.2.4* as an example. If there are infinitely many prime numbers, then that means there doesn't exist a smallest or largest prime numbers. So the contradiction would be that.

- Now that I have equivalent statements, what do I head for next, if I just have very insufficient information to work with?

If you don't have anything, try crafting something. In this case, **for the sake of contradiction** saved our lives in 2.2.4 by providing us an imaginary case to work with, where there exists a largest prime number  $p_k$  and a number that is some value more than it and must be non-prime. In my view, everything afterwards really just depends on creative work.

- How do I get better at creative work?

Read more proofs. Read how people approach different problems with different perspectives, internalize homework solutions and online forum proofs, the criticisms on them. You will only broaden your horizon via looking at how other people create proofs (well, at the cost of playing less Smash Bros Ultimate).

- What do you really mean by a perspective?

In COMPSCI 61B, we have approached a topic called "view", which means an abstraction of a complicated system. For example, the view of a social network might be a graph, the view of a playlist might be a queue. In this case, we can manage and switch between views when encountering proofs with different mathematical objects.

View, or perspectives (which is my preferred word choice), usually deals greater effort on mathematical object proofs. Number theory related proofs take more insights and approaches that rely on the properties of numbers, such as division results and rationality of numbers.

And as mentioned before, these efforts can be improved by reading more proofs and learning more about conducting proofs by some individual learning.

In summary, everything starts with making a first step.



## Chapter 3

# Mathematical Induction

In this chapter, we introduce and discuss a new form of deduction called "Mathematical Induction". In the examples offered by the chapter, we discuss the power and practicality of mathematical induction, as well as the different forms of inductions that come with different tips of using.

Meanwhile, as the chapter does not include the lecture note's description on a connection between COMPSCI 61A's "Leap of Faith" and mathematical induction, the readers are encouraged to nurture an understanding towards this connection in their vision.

Hopefully, it will facilitate the understanding towards mathematical induction and provide a firmer ground on recursive functions/programs.

### 3.1 Introduction to Mathematical Induction

#### 3.1.1 Definition of Mathematical Induction

Previously, we introduced techniques and directions we can take to produce mathematical proofs. Here, we will introduce a new powerful tool, a new proof technique, a form of deduction called "mathematical induction". Such form of logic establishes a statement holds for all natural numbers.

This is an important contribution.

First of all, many mathematical objects are closely associated with natural numbers, especially those that we will encounter in COMPSCI 70.

And second of all, when we are asked to prove statements for all natural numbers, there are an infinite number of values to be checked for the statement. It poses the problem: while mathematical proof requires a finite sequence of logic, the prompt requires us to guarantee a proposition for infinite values. This implies how strong a tool mathematical induction can be for mathematicians, and furthermore, computer scientists.

How exactly does mathematical induction work? Or, how does it guarantee the statement to be right for all natural numbers? Here is a structural summary of mathematical induction to answer this doubt:

- **Base Case:** that the statement works for an initial value  $k$ . In the case we are proving for all natural numbers,  $k = 0$  is a popular choice. Therefore, we attempt to prove at this step that the prompt proposition  $P(n)$  works for the value 0.
- **Induction Hypothesis:** We here make the assumption that for an arbitrary value  $n \geq k$ ,  $P(n)$  is true. Building on that, we proceed into the next step:
- **Inductive Step:** With the induction hypothesis, we here attempt to prove that  $P(n + 1)$  is true.

Let us see an example of mathematical induction utilized in a direct proof:

## Exercise Reader! Exercise! 3.1.1: Mathematical Induction in Direct, Algebraic Proof

Prompt:  $\forall n \in \mathbb{N} (\sum_{i=0}^n i = \frac{n(n+1)}{2})$

For brevity, let the proposition  $P(x)$  be the second clause of the prompt:  $\sum_{i=0}^x i = \frac{x(x+1)}{2}$ .

**Base Case:** Prove  $P(0)$ .

$$\sum_{i=0}^0 i = 0 = \frac{0(0+1)}{2}$$

Hereby,  $P(0)$  has been proven.

**Induction Hypothesis:** Assume for an arbitrary  $k$ ,  $P(k)$  holds, such that  $\sum_{i=0}^k i = \frac{k(k+1)}{2}$ .

**Inductive Step:** Assuming  $P(k)$  is true, prove  $P(k+1)$ .

$$\begin{aligned} \sum_{i=0}^k i &= \frac{k(k+1)}{2} \\ \sum_{i=0}^{k+1} i &= \frac{k(k+1)}{2} + (k+1) \\ &= \frac{k^2+k}{2} + \frac{2k+2}{2} \\ &= \frac{k^2+3k+2}{2} = \frac{(k+1)(k+2)}{2} = \frac{(k+1)((k+1)+1)}{2} \end{aligned}$$

As seen in the above proof, we assumed the induction hypothesis to hold when attempting to complete the inductive step. Often, the induction hypothesis even serves as a critical step towards the completion of an induction.

### 3.1.2 Graphical Mathematical Induction

Mathematical Inductions not only work from the algebraic perspective. Since it can prove statements to hold for all natural numbers, these statements ought to not only be propositions for algebraic identities. It can also be geometrical phenomena. Let us view an example from the lecture notes:

## Exercise Reader! Exercise! 3.1.2: Mathematical Induction in Geometric Proof

Define  $n$ -colorable as able to color a map with  $n$  colors such that no adjacent areas share the same color.

Prompt: Let  $P(n)$  denote the statement "Any map where the border of all areas are marked by within  $n$  lines is two-colorable", then prove that  $\forall n \in \mathbb{N} (P(n))$ .

**Base Case:** Consider  $P(0)$ . In this case, a map with 0 lines to mark the border of areas is just one huge area, which is definitely two-colorable.

**Induction Hypothesis:** Assume for some arbitrary  $k \geq 0$ ,  $P(k)$  holds.

**Inductive Step:** Let us assume that the induction hypothesis is true. In the inductive step, our map will provided a new line to mark the borders of newly created regions. These newly created regions of the map must come as a partition of previously whole areas.

Now, let us set the line onto the map. Let us defined the side left or above the new added line (the  $(k+1)^{th}$  line) as the left hand side (LHS), and the other (RHS).

We notice that the current map would not be two-colorable, because the newly created regions must have the same color with their other partition.

Let us save this proof with a small maneuver. By swapping all colors existing on the LHS to their opposite, then while the two-colorable property is preserved under any color swapping, all original whole areas will now hold partitions of different colors, for their partitions are splitted across the LHS (the swapped side) and RHS

(the unswapped side).

Consequently, for the  $k$ -line map is two-colorable as the induction hypothesis states, the  $k + 1$ -line map must also be. Therefore, under such assumption, the inductive step works and  $P(k + 1)$  holds.

Not only are mathematical inductions useful on such graphical proofs, it can further provide strong proofs towards theorems that are considered really hard to prove in algebraic means.

Graphical proofs rely more on geometry and logic than mere algebra, which makes it different in property from the conventional mathematical inductions we see in the lecture notes dealing with number theory and arithmetic manipulations. This comes very clear when we start using abstractions for graphical objects in discrete mathematics.

### 3.1.3 The Strength of Induction Hypothesis

However, mathematical induction is not all-able. If the prompt is too difficult to prove, or if the induction hypothesis is too weak for us to proceed the inductive step with, we will be unable to complete a proof.

Let us visit an example from the lecture notes regarding this problem.

Exercise Reader! Exercise! 3.1.3: When the Induction Hypothesis is SUSpicious

Prompt: for all  $n \geq 1$ , the sum of the first  $n$  odd numbers is a perfect square.

**Base Case:** Prove the prompt for  $n = 1$ . The first odd number is 1, and  $1 = 1^2$  is a perfect square.

**Induction Hypothesis:** Assume that for an arbitrary  $k \geq 1$ , the sum of the first  $k$  odd numbers is a perfect square.

**Inductive Step:** The  $(k + 1)^{th}$  odd number would be  $(2k + 1)$ . Considering the sum of the first  $k$  odd numbers is a perfect square  $m^2$ , the sum of the first  $k + 1$  odd numbers would be  $m^2 + 2k + 1$ .

What next.....? Proof incomplete.

We were not able to complete the above induction because our induction hypothesis was not powerful enough to drive us towards completing the inductive step. This induction hypothesis, or a prior assumption that we may call **priori**, needs to be a different proposition that speaks the same thing as an induction hypothesis does. Or, it has to be a stricter, stronger proposition for us to assume.

Let us try a different induction hypothesis for the same prompt:

Exercise Reader! Exercise! 3.1.4: When the Induction Hypothesis is Reliable

Prompt: for all  $n \geq 1$ , the sum of the first  $n$  odd numbers is  $n^2$ .

**Base Case:** Prove the prompt for  $n = 1$ . The first odd number is 1, and  $1 = 1^2$ .

**Induction Hypothesis:** Assume that for an arbitrary  $k \geq 1$ , the sum of the first  $k$  odd numbers is a  $k^2$ .

**Inductive Step:** The  $(k + 1)^{th}$  odd number would be  $(2k + 1)$ . Considering the sum of the first  $k$  odd numbers is a perfect square  $k^2$ , the sum of the first  $k + 1$  odd numbers would be  $k^2 + 2k + 1 = (k + 1)^2$ .

What next.....? Proof incomplete.

While this prompt allows us to prove the previous, it also has a stricter proposition to prove with, which comes with a more convenient statement to prove with.

In other words, the original prompt was too vague, causing our Induction Hypothesis to be weaker.

As for how we can discover a stronger hypothesis, we can usually manually demonstrate the proposition of prompt and attempt to see the pattern in which it was achieved:

## Exercise Reader! Exercise! 3.1.5: Finding A Stronger Prompt

Prompt: for all  $n \geq 1$ , the sum of the first  $n$  odd numbers is a perfect square.

$$n = 1 : 1 = 1^2$$

$$n = 2 : 1 + 3 = 2^2$$

$$n = 3 : 1 + 3 + 5 = 3^2$$

$$n = 4 : 1 + 3 + 5 + 7 = 4^2$$

We found there seems to be a pattern that: Prompt: for all  $n \geq 1$ , the sum of the first  $n$  odd numbers is  $n^2$ . We then performed mathematical induction on it as an attempt to find a stronger hypothesis to prove the original prompt with, which worked!

## 3.2 Simple Induction and Strong Induction

In this section, we describe strong induction as an idea similar to simple induction, demonstrate the TPO of using either techniques, and will form an analogy of both inductions as tools that perform the same purpose with different equipped power.

We might recognize a strong induction as a TI-84, if simple induction is a TI-83. How so? Let us attempt to answer this by exploring the contents of this section.

### 3.2.1 The Strength of Induction: Simple or Strong

Up until now, we have used the three-part format of mathematical induction known as "simple induction", or "weak induction".

Notably, there is another notion of induction called "strong induction", which expects a slightly different induction hypothesis:

$$\bigwedge_{i=0}^k P(i) \text{ is true}$$

While both forms of induction hold the same purpose and accomplish the same result, the amount of power needed in completing objectives is different across these two tools. Strong induction provides us a stronger hypothesis to work with, thus is more simple to use. Let us demonstrate so with an example:

## Exercise Reader! Exercise! 3.2.1: New Technique: Strong Induction

Prove:  $(\forall n \in \mathbb{N})((n > 12) \implies (\exists x, y \in \mathbb{N}(n = 4x + 5y)))$

**Base Case:**

$$n = 12 = 3 \times 4 + 0 \times 5$$

$$n = 13 = 2 \times 4 + 1 \times 5$$

$$n = 14 = 1 \times 4 + 2 \times 5$$

$$n = 15 = 0 \times 4 + 3 \times 5$$

**Induction Hypothesis:** Assume that the prompt holds for all  $12 \leq n \leq k$  for some arbitrary  $k \geq 15$ .

**Inductive Step:** We need to prove the claim for numbers beyond such limit  $k$ . We will prove for the case  $n = k + 1 \geq 16$ , or with some algebra: the case  $(k + 1) - 4 \geq 12$ .

Since  $12 \leq (k + 1) - 4 \leq k$ , by the induction hypothesis,  $\exists a, b \in \mathbb{N}((k + 1) - 4 = 4a + 5b)$ .

It would then provide:  $(k + 1) = 4a + 4 + 5b = 4(a + 1) + 5b$ . Therefore, the prompt holds for value  $(k + 1)$

Let us provide another example to witness the efficiency in strong induction, as well as to explore an example on the lecture note.

Exercise Reader! Exercise! 3.2.2: New Technique: Strong Induction, ex2

Prove that: every natural number  $n > 1$  can be written as a product of one or more primes.

**Base Case:** The number  $n = 2$  itself is prime, therefore can be written as a product of a prime number (itself).

**Induction Hypothesis:** Let us assume the prompt holds for all numbers  $2 \leq n \leq k$ , for an arbitrary  $k$ .

**Inductive Step:** We face the number  $k + 1$  to prove the prompt for. Here are two cases: either  $k + 1$  is prime, or not prime.

If  $k + 1$  is prime, then the problem is resolved as it can be written as a product of one prime number (itself).

If it is not prime, however, then we would assume it is the product of two natural numbers  $x, y$  such that  $1 < x \leq y < k + 1$ , as multiplication should work.

However, such numbers  $x, y$  should by the induction hypothesis be able to be written as a product of one or more prime numbers. If so, then so can  $k + 1$ . Via proof by case and mathematical induction, we have proved the prompt to hold true over the range of natural numbers it quantifies over.

### 3.3 Build-Up Error

However, it is very easy to make mistakes on mathematical inductions if we blindly assume the induction hypothesis to be non-vacuously true.

In other words, we can dangerously assume for an arbitrary  $k$  and the prompt's proposition  $P$  that  $P(k)$  is true, while it might not be!

Exercise Reader! Exercise! 3.3.1: When Mathematical Induction is SUSpicious

**Prompt:** Every horse has a same color

**Wrong Proof:**

**Base Case:** Let there be only one horse. One horse has one color, prompt is trivially true.

**Induction Hypothesis:** Let a set of  $n$  horses always have the same color. (Hey, this is MEGA-SUS! We'll explain why later.)

**Induction Step:** Let the new set of horses have  $n + 1$  members. Then, any  $n$ -element subset of this set of horses must have the same color, thus all horses must have the same color!

Where did it go wrong? Or, more exactly, why is mathematical induction sus? Well, there suddenly exist counterexamples to the induction hypothesis, and we blindly assumed that there isn't.

This is known as a **build-up error**, where we falsely assume the prompt to work for any natural number input already and thus falsely prove the prompt as right.

Let us take a stricter look at the above proof.

In the induction step, let us first remove an arbitrary horse from the set of  $n + 1$  horses. The induction hypothesis can tell us that this leftover set of  $n$  horses all have a same color.

But when we apply  $n = 1$ , this entire claim becomes suspicious: we, if following the induction step, actually get two disjoint sets (which means without common elements) of horses, each of cardinality 1. So the horses are never in the same color group, and cannot necessarily just have the same color!

The build-up error occurs: it shows some valid reasoning via a reasoning that is fundamentally flawed, and it is up to the prover to recognize the flaw.

### 3.4 Personal Regard on This Topic

Mathematical Induction is difficult to deal with, because induction hypotheses are sometimes hard to seek and in many cases difficult to configure and debug around.

But as always, practicing will help a lot in familiarizing with the content and coming with esoteric and maverick ways to prove diverse prompts!

For more advices regarding proofs, try searching in the previous chapter and your experiences.

## Chapter 4

# The Stable Matching Problem

In this chapter, we introduce the Stable Matching Problem and the Propose-And-Reject Algorithm, and use it as a case study of how proofs can be applied to problems of the computer science field to guarantee the efficiency and legitimacy of an algorithm.

### 4.1 Description of Stable Matching Problem

In this section, readers will be introduced to a frequently occurring problem of institutions, called the "Stable Matching Problem", as well as learn an algorithmic solution of it that we will analyze the inner workings of in the next following section.

#### 4.1.1 Stable Matching Problem and Propose-And-Reject Algorithm

Say we have an employment system, having to match  $n$  jobs to  $n$  candidates. Each job has a list of  $n$  candidates ordered based on preference, while a preference-ordered list of  $n$  jobs also exist for each candidates.

Our problem is, how do we find the best match, defined as in "switching jobs does not benefit each candidate"?

The algorithm that achieves this is called the "Propose-and-Reject algorithm", or more formally known as the Gale-Shapley algorithm. Let us see how this algorithm works in the hypothetical context we presented above, in "phases":

Exercise Reader! Exercise! 4.1.1: How does the Gale-Shapley Algorithm Work?

At phase 1, every job position proposes to the most preferred candidate on its list who has not yet rejected their position.

At phase 2, each candidate collect all offers they receive, and responds "maybe" to their favoirt e positions while "no" to other offers.

At phase 3, each rejected position crosses off the candidate who rejected them.

We repeat the loop from phase 1, as we have completed phase 3. We stop looping when no offers are rejected, as that means each candidate has a job offer of their "maybe". The candidates will then accept their offers.

In the following section, let us analyze this algorithm and answer the essential query: "How on Earth does this simple approach solve the significant problem?"

### 4.2 The Properties of Propose-And-Reject Algorithm

In this section, we will analyze the properties of Gale-Shapley Algorithm, as in what makes it a working algorithm and provide the proof that it works. This section will involve many mathematical induction, and we will finally see

induction being applied to problems of the computer science fields!

### 4.2.1 The Properties of Propose-And-Reject Algorithm

When discussing whether an algorithm is good or not, we mind about two properties: whether it halts, and whether it outputs a stable (good) matching.

The former can be proven true using a proof:

Exercise Reader! Exercise! 4.2.1: Does the Gale-Shapley Halt?

On each day the algorithm doesn't halt, at least one job must eliminate some candidate from its list. Or else, it means that every candidate accepted an offer, and therefore the matching has ended.

Assuming this worst case for lists of  $n$  jobs and  $n$  candidates, this algorithm has a worst case runtime of  $O(n^2)$  for it halts in at most a finite number,  $n^2$ , of iterations.

How do we define a "stable" match then? What is the heuristic, the metric, or the standards for deciding how good a matching is?

In the context we work with, let us define the metric as:

- A matching is unstable if there is a job and a candidate who both prefer working with each other over a current matching.
- A matching is stable if there are no couples like above, called "rogue couples".

First of all, we want to assess whether there could always be stable matchings, since if not, then the algorithm could still lead us astray from the combinations we want.

Let us observe how the algorithm's list of choices develops.

While each jobs begin with its first choice, their choice prime candidates could betray and reject them. Therefore, jobs have worse choices over time. However, candidates will reject any choices that they do not like among all they receive and "maybe", therefore only having better choices over time.

This can be formalized into a lemma:

#### Lemma 4.2.1. Improvement of Choices in GS Algorithm

If job  $J$  makes an offer to candidate  $C$ , then on every subsequent day  $C$  has a job offer in hand which candidate  $C$  likes at least as much as  $J$ .

##### Proof:

Let us perform induction on a day  $i \geq k$  **Base Case** ( $i = k$ ): on day  $k$ , candidate receives at least one offer from job  $J$ . Since the candidate  $C$  chooses whichever of the current favorite and offer from  $J$  is better, in the case  $J$  is not chosen, the offer  $C$  holds is better than  $J$ , and in the case  $J$  is chosen,  $C$  has an offer as good as  $J$  per se.

**Induction Hypothesis:** Let us assume the lemma holds for some arbitrary  $n \geq k$ .

**Inductive Step:** Let the current favorite offer of candidate  $C$  be  $B$ , and the offer received today be from job  $J$ . Since  $C$  must choose the better among  $B$  and  $J$ , if  $B$  is chosen, then  $B$  is at least more favored and liked than  $J$ , while if  $J$  is chosen, the offer that  $C$  holds is as good as  $J$ . Considering both cases, the offer can only be as good as or better than  $J$ .

Or, there is an alternative way to prove this lemma, utilizing another mathematical principle that makes its debut in this chapter.



## Exercise Reader! Exercise! 4.2.2: Proving the Improvement Lemma, Alternative

Before moving on, let us consider:

## Definition 4.2.1. Well-Ordering Principle

If  $S \subseteq \mathbb{N}$  and  $S$  is not empty, then  $S$  has a smallest element. Or, put in mathematical symbols as an implication:

$$S \subseteq \mathbb{N} \wedge S \neq \emptyset \implies \exists x \in S (\forall y \in S (x \leq y))$$

Now, let us assume that there exists an  $n > k$  after our base case in the previous induction proof such that at day  $n$ , a first counterexample to the lemma occurs. The candidate  $C$  has either no offer or an offer from some other job  $H$  less favorite than the current favorite  $J$ .

On day  $n - 1$ , candidate then receives an offer from some job  $K$  and liked it at least as much as  $J$ , which is as the lemma proposes. Therefore, this offer  $K$  still exists on day  $n$ , the exception day, which must make the new favorite be better than  $K$ , thus  $J$ , thus  $H$ .

We see from the above logic that the proposition  $\neg(P(i) \implies P(i+1))$  cannot hold for any values  $i$ . Therefore, by the law of excluded middle,  $(P(i) \implies P(i+1))$  will hold for any arbitrary value  $i$ . A Proof by Contraposition guides to the completion of a mathematical induction.

The above proof is secured by the Well-Ordering Principle because the order of proof in induction that assumes the smaller natural number to be ordered before the larger natural number is required for both the hypothesis of some day  $n$  being a first counterexample and the basic hypothesis of induction on proving the proposition to an incremented number.

Now, we only need to know two more things:

First of all, the propose-and-reject algorithm does halt, and furthermore, terminate with a matching.

Second of all, the matching produced by the algorithm is always stable, which we prove via the improvement lemma:

## Exercise Reader! Exercise! 4.2.3: Does the GS Algorithm always terminate with a matching?

Let us suppose it does not, such that there is a job  $J$  left unpaired when the algorithm terminates. That means the job  $J$  has offered to and been rejected by all possible candidates.

By the Improvement Lemma, this means each candidate must have had a better offer than what job  $J$  provided. Therefore, there in fact exists  $n$  jobs better than  $J$ , leaving us with a list of  $n + 1$  jobs. However, we only assumed there to be  $n$  jobs.

Therefore, via proof by contradiction, we have proved the lemma:

## Lemma 4.2.2. Termination of GS Algorithm

The Gale-Shapley algorithm always terminates with a matching.

Taking this to the next step, we verify the stability of each matching.

## Exercise Reader! Exercise! 4.2.4: The Stability of Matching by Termination

Let us consider an arbitrary couple  $(J, C)$  in the final result of matching.

Let us suppose that this arbitrary job  $J$  prefers some other arbitrary better candidate  $C^*$  to  $C$ . However,  $C^*$  is currently paired with some other job  $J^*$ , and by the Improvement Lemma, this must be because that  $J^*$  is at least as favorable as  $J$ , making  $C^*$  having favored  $K$  more than  $J$ .

Job  $J$  also wouldn't want to switch to a worse candidate, so the best it can have now is  $C$ . The job  $J$  can therefore never be involved in a rogue couple.

We have proved there cannot be any job  $J$  involved in a rogue couple, such that for any candidate  $J$  prefers better, that candidate doesn't want  $J$  more.

**Theorem 4.2.1. The Property of GS Algorithm**

The matching produced by the Gale-Shapley algorithm is always stable.

**4.2.2 Optimality of Propose-And-Reject Algorithm**

We would also like to discuss whether the matches this algorithm provides is optimal. However, before discussing so, to make good mathematical propositions, we must define optimality.

**Definition 4.2.2. Optimality of Match in GS Algorithm**

For a given job  $J$ , its optimal candidate is the highest rank candidate on the preference list of  $J$  that  $J$  could be paired with in any stable matching.

For a given candidate  $C$ , its optimal job is the highest rank job on the preference list of  $C$  that  $C$  could be paired with in any stable matching.

Put more straightforward, optimal is the highest ranked choice possible for the resulting match to be stable.

A matching where each job is paired with its optimal candidate is known as a job optimal matching. On the other hand, a matching where each candidate is paired with their respective optimal job is known as a candidate optimal matching. The opposite of this concept is "pessimal".

And now, let us decide whether the matching output by our GS Algorithm is job-optimal:

**Exercise Reader! Exercise! 4.2.5: Is the GS Algorithm providing a job-optimal matching?**

For the sake of contradiction, let us propose that the algorithm doesn't.

Let there exist a day on which a job has its offer rejected by the optimal candidate, and the first day of this counter be on day  $k$  (mind that we are again using the well-ordered principle here via having a "first counterexample").

Suppose on that day,  $J$  was rejected by the optimal candidate  $C^*$  in favor of an offer from  $J^*$ . Then, this couple  $(J, C^*)$  should exist in some stable matching. Therefore, the matching would look something like:

$$T = \{\dots, (J, C^*), \dots, (J^*, C'), \dots\}$$

First of all, if  $C^*$  rejected  $J$  in favor of  $J^*$ , then  $C^*$  must have liked to work on  $J^*$  more. Furthermore,  $J^*$  is assumed to have made an offer to  $C^*$  at some day before, but finally accepts  $C'$  in turn. This means  $C^*$  is at least as good as  $C'$ , and that  $J^*$  would've wanted to work with  $C^*$  instead.

We have come to see that  $(J^*, C^*)$  is in fact a rogue couple, making  $T$  an unstable match rather than a stable one as we assumed with the above logic.

Therefore, via proof by contradiction, we found:

**Theorem 4.2.2. The Optimality of GS Algorithm I**

The matching produced by the Gale-Shapley algorithm is always Job-Optimal.

as well as candidate-optimal:

**Exercise Reader! Exercise! 4.2.6: Is the GS Algorithm providing a candidate-optimal matching?**

Let us utilize the Job-Optimality of GS Algorithm and let  $T$  be an employer-optimal matching:

$$T = \{\dots, (J, C), \dots\}$$

For the sake of contradiction, let us propose that there exists a stable matching:

$$S = \{\dots, (J^*, C), \dots, (J, C'), \dots\}$$

such that job  $J^*$  is ranked lower than job  $J$  on the preference of  $C$ .

However, if so, then  $C$  must have wanted to work with  $J$  more and  $J$ 's current candidate  $C'$  is not preferred more than  $C$  given the stable matching  $T$ .

We managed to show that the matching  $S$  is both stable and unstable, which itself is a contradiction. Therefore, the matching from GS Algorithm cannot be candidate-optimal.

If so, then the employer-optimal pairing is already the worst possible pairing for candidates, and therefore:

**Theorem 4.2.3. The Optimality of GS Algorithm II**

The matching produced by the Gale-Shapley algorithm is always Candidate-Pessimal.

Notably, this algorithm can still be modified in finer details to provide a priori in preference from jobs, as well as converting the algorithm into possessing candidate-optimality instead.

## 4.3 Review Questions from Midterm

There are a couple of tricky questions I came across during practice, of which I will present some on this note for review purpose.

Q1 What is the max amount of job partaking in rogue couple(s) if we permute one job's ranking list of its candidates?

A1 It would be 1, since other jobs and candidates do not change their preferences, and this job would be the only job who starts to prefer some candidate more.

Q2 True or False: Any stable matching is either job-optimal or candidate-optimal

A2 False. This statement would assume that at most 2 matchings exist per stable matching instance, which is not true.

Q3 If a stable matching is optimal for a job  $j$ , it is pessimal for its partner  $c$ .

A3 Let there be a stable matching:

$$T = \{\dots, (J, C), (J', C'), \dots\}$$

which is optimal for job  $J$ .

Here, let  $J'$  be ranked lower than  $J$  on the list of preferences. We will then find another matching where:

$$T' = \{\dots, (J', C), (J, C^*), \dots\}$$

where then  $(J, C)$  becomes a rogue couple.

Therefore,  $C$  cannot be matched to a worse job than  $J$ , whose pairing is defined optimal. In this case, the candidate  $C$  is indeed receiving pessimal matching.

On this chance, let me introduce you to an interesting stable matching instance:

A	1	2	3	1	C	A	B
B	1	2	3	2	B	C	A
C	2	1	3	3	-	-	-

Which is interesting because it completely resolves all following subquestion:

Q4 In the job-propose and candidate-reject algorithm:

1. There is always one job who is never rejected.
2. There is always one candidate who never rejects a job.
3. Let  $(J, C)$  be a proposal made on the last day. If we delete  $J$  and  $C$ , running the propose-and-reject algorithm always produces the same matching for the remaining jobs and candidates.

A4 Let us consider the above stable matching instance:

1. Prove that jobs can get rejected via the above instance.
2. Prove that candidates can have to reject via the above instance.
3. Deleting job  $B$  waives all collisions that was to come due to its similarity with job  $A$ , causing different matchings.

Q5 Prove that the maximum number of rejections in an instance with  $n$  jobs and  $n$  candidates is at most  $(n-1)^2$

A5 One job can be rejected  $n - 1$  times, while all other jobs can be rejected  $n - 2$  times. This causes a sum of  $(n - 1) + (n - 2)(n - 1) = (n - 1)^2$ .

## 4.4 Personal Regard on This Topic

When doing proofs regarding stability, it is always advised to operate with imaginary entities. This almost always offers a more concrete insight to the proof of operation.

This algorithm is an interesting topic that demonstrate the importance of mathematical induction, as well as showcases the discussion and analytics we may perform on an algorithm.

Personally, here are a few topics that I chose to highlight during revision:

- **Improvement Lemma:** Candidates will only hold better or equal offers.
- **Optimality:** Candidate-pessimal and job-optimal for our current settings. However, by switching the places of candidates and jobs, this optimality reverses.
- **Stability:** Definition is customizable across different short answer settings, but the commonality is there should not be rogue couples (matchings), where the second and first element of some different pair wants to be with each other instead.
- **Uniqueness of Stable Match:** Stable matches can be generated via other algorithms. This is just one of them. Therefore, there can exist more stable matchings than we know.

Furthermore:

- No entities are guaranteed favorites as long as everyone gets rejected once.
- At most  $n^2 - n$  possible rogue pairs are possible.
- If every entity has their least favorite corresponding entity, the stable matching is unstable.
- Stable matching can exist under accidental rejections.

Imaginary entities, or denoting entities, is VERY useful. It allows you to form arguments from constructions and produce a proof, which graphs, a mathematical object that we work with in the future, will require you to perform. Please use it wisely.

# Chapter 5

## Graph Theory

In this chapter, we will discuss Graph Theory, a highly applicable set of mathematical theories that help with algorithmic operations via representational abstractions in computational and mathematical notions.

### 5.1 An Introduction to Graph Theory

#### 5.1.1 Motivation

Abstraction. It is the universal theme of many things computer science. We attempt to capture the essence, simple representation of a complex situation, from concepts to functions.

And graphs happen to be a very popular abstraction. Graphs are sets of connections between customizable individuals, working from neural networks (individuals are neurons) to your mom (individuals are cells, connections are the chemical interactions of cells).

This mathematical object, graphs, have been developed by mathematicians and scientists. Eventually, we are provided a framework for the object, and this framework is known as **graph theory**.

##### Definition 5.1.1. Graph Theory

The mathematical theory of properties and applications of Graphs.

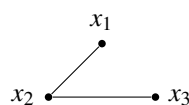
Albeit I realize that grades would be a popular motivation for studying, understanding graphs will provide good fluency across many advanced computational concepts. It is necessary, and is essential to the language of mathematics and computers.

#### 5.1.2 Foundations and Formation of Graphs

Formally, a **graph** is defined by a set of **vertices**  $V$  and a set of **edges**  $E$ . Vertices are the individuals and edges are the connections.

Graphically, each vertex corresponds to the small circles in the following example of graphs, while edges are line segments connecting these vertices:

Figure 5.1.1. Example of a Graph



The *LaTeX* code is shown as follows:

```
\begin{tikzpicture}
  [point/.style = {circle, fill, inner sep = 1pt}]
  \node[point] (1) [label=above:$x_1$] {};
  \node[point] (2) [below left of=1, label=left:$x_2$] {};
  \node[point] (3) [below right of=1, label=right:$x_3$] {};
  \draw (1) -- (2);
  \draw (2) -- (3);
\end{tikzpicture}
```

Exercise Reader! Exercise! 5.1.1: How do you mathematically represent the graph in Figure 5.1.1?

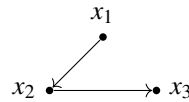
The set of its vertices would be  $V = \{x_1, x_2, x_3\}$ .

Its edge, meanwhile, can be contained in the set  $E = x_1, x_2, x_2, x_3$

Notice that here,  $E$  is a multiset, meaning there can be an element appearing multiple times if an edge is repeated. But we would really like  $E$  to be a set instead.

To let repeated edges that might have had different directions be counted distinctly, let us also define a directed graph, where each edge must have a direction:

Figure 5.1.2. Example of a Graph



The *LaTeX* code is shown as follows:

```
\begin{tikzpicture}
  [point/.style = {circle, fill, inner sep = 1pt}]
  \node[point] (1) [label=above:$x_1$] {};
  \node[point] (2) [below left of=1, label=left:$x_2$] {};
  \node[point] (3) [below right of=1, label=right:$x_3$] {};
  \draw[->] (1) -- (2);
  \draw[->] (2) -- (3);
\end{tikzpicture}
```

The example shown above is known as a **directed graph**, graphs equipped with directed edge.

#### Definition 5.1.2. Directed-ness of Edge

If an edge models a one-way path from one vertex to another, such an edge is known as a **directed edge**. Mathematically, it is represented as an edge with an arrowmark.

For the edge set from the above example,

$$(x_1, x_2) \in E \wedge (x_2, x_1) \notin E$$

If an edge is not arrowmarked, such that it allows travel between two vertices in both directions, then the edge is an **undirected edge**.

We can use a directed edge on multiple abstractions. For example, in a circuit diagram, current flows from one point to another in a single direction.

But, an undirected graph is also useful for modeling, say a metro map, where it is free to commute between stations

for any direction.

### 5.1.3 Edge and Degree

For an undirected edge  $e = (u, v)$  (or in equivalent notation,  $\{u, v\}$ ), such edge  $e$  is said **incident** on vertices  $u$  and  $v$ . Meanwhile, these vertices connected by a same edge would be neighbors, or in more formal terms, **adjacent**. If a vertex is involved in no edge, it is called an **isolated vertex** and would also be disconnected from the graph. A directed graph, however, has one-directional edges. In that case, vertices will still be incident, but the way to calculate visits slightly differs. We will see when we discuss the metric of visit frequency of vertices. We may also use a metric to see how frequently a vertex is involved in the travelings of graphs: **degree**.

#### Definition 5.1.3. Degree

In general, **degree** measures the amount of edges a vertex has been visited by or involved in. Such measure would differ slightly in the context of graph: whether it is directed or undirected. If the graph is undirected, then degree of a vertex  $u$  is calculated as the amount of edges that involve  $u$ :

$$\deg(u) = |\{v \in V : u, v \in E\}|$$

If the graph is directed, there are thwn two measures of degree:

- In-degree of  $u$ : The amount of edges that directs towards  $u$ .
- Out-degree of  $u$ : The amount of edges that directs from  $u$ .

If an edge originates from one vertex to itself, then this edge  $u, u$  is known as a **self-loop**. Again, depending on abstraction, this can either be utterly useless information or very useful information.

Since it is difficult to analyze graphs with self-loops, we will talk about graphs without self-loops in the following sections, and will not discuss the occassion of multiple edges between a pair of vertices (unless their directions are all distinct).

### 5.1.4 Traversing a Graph, Mathematically

We traverse graphs, algorithmically most of the time. How do we mathematically characterize a traverse?

#### Definition 5.1.4. Path

A **path** in a graph  $G$  is a sequence of edges. A path would, then, start from one vertex and end at a vertex. Essentially, it is an itinerary, a projectory of traverse!  
A path is *simple* when the vertices it travels are distinct.

#### Definition 5.1.5. Cycle, Tour

A **cycle** is a sequence of edges that starts and ends at the same vertex. More specifically, all vertices traveled are distinct except the starting and ending vertex, and no edges are repeated.  
If a sequence of edges simply starts and ends at the same vertex, it is known as a **tour**.

Let's make a brief summary:

Figure 5.1.3. Walk vs Path vs Cycle vs Tour

Type	No Repeated Vertices	No repeated edges	Start is End
Walk			
Path	○	○	
Tour			○
Cycle	except start and end	○	○

5.1.5 Connectivity

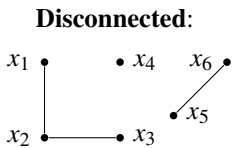
Previously, we discussed the notion of isolated vertices. How do we, then, mathematically describe the connected-ness of a graph?

Definition 5.1.6. Connectivity

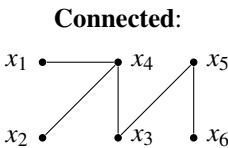
A graph is **connected** if there is a path between any two distinct vertices. This means isolated vertex is in fact a connected component.  
On the contrary is a **disconnected** graph.

Let me demonstrate with some examples:

Figure 5.1.4. Connected vs Disconnected



As you may have seen, vertex 4 is an isolated vertex that no one visits (just like me).  
Meanwhile, the subgraph involving vertex 5 and 6 is disconnected from the rest of the graph, as well as the isolated vertex.



In such a connected graph, there is a path from every vertex to another vertex.

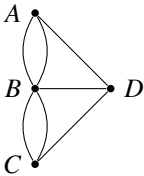
5.2 Koenigsberg’s Seven Bridges and Eulerian Tour

5.2.1 The Seven Bridges of Koenigsberg

Imagine that there are seven bridges connecting four cities, abstractable to the graph as seen below:

Figure 5.2.1. Seven Bridges of Koenigsberg

The Seven Bridges of Koenigsberg can be represented by the graph as follows:



Where the seven edges each are distinct bridges, and the four vertices are distinct cities.



Now, the question is: is there a route that would traverse each seven bridges precisely once and return to the starting point.

### 5.2.2 Euler's Take

Euler, having contributed a lot to graph theory, came with some terminologies regarding itineraries of traversals:

#### Definition 5.2.1. Eulerian Itineraries

**Eulerian Walk** is a walk of graph  $G$  that uses each of its edge exactly once.

**Eulerian Tour** is a closed Eulerian walk where the starting and ending point are the same vertex.

This rephrases the Seven Bridge Problem to: "Is there an Eulerian tour for its representative graph?"  
The answer lies in the Eulerian Theorem:

#### Theorem 5.2.1. Eulerian Theorem

An undirected graph  $G = (V, E)$  has an Eulerian tour iff  $G$  is even degree, and connected (except isolated vertices).

An **even-degreed** graph is a graph where all vertices have even degree.

Its proof is as follows:

#### Exercise Reader! Exercise! 5.2.1: Prove the Eulerian Theorem!

##### Proving the forward direction of theorem:

Assume  $G$  has an Eulerian tour. Since it has to use every single edge, every vertex that has an adjacent edge will be involved in this tour, therefore connected with all other vertices on the tour.

Meanwhile, because everytime a tour enters a vertex along an edge, it exits along a different edge, each vertex will have a pair of edge adjacent to it. The start vertex would also have the same phenomenon, since we begin from and end at the starting vertex.

In that case, the graph would be even-degreed for the pairs of edge that occurs on vertices (which must all be used up), and would be connected.

##### Proving the backward direction of theorem:

Let us first have a subroutine called FIND\_TOUR that finds a tour from a graph  $G$ .

Notice that while the tour is not necessarily Eulerian, it must get stuck at the vertex it started at, because it would then deplete the pairs of edges it can continue its travel with. This does not serve as a formal proof, but hopefully the conceptual summary of lecture notes' claim is good enough.

Then, let's also involve another subroutine SPLICE, which outputs a single tour  $T'$  provided a tour  $T$  that intersects each of the other arguments  $T_1, \dots, T_i$ . Last but not least, these tours are merged together. The input tours are all edge-disjoint.

Last but not least, EULER, which once provided a graph  $G$  and starting vertex  $s$ , will attempt to find a spliced tours  $G_1, \dots, G_i$  from an arbitrary found tour  $T$  of  $EULER(G_i, s)$ , and its edge-disjoint intersected tours.

Let us use a shortened mathematical induction to prove the functionality of EULER:

**Base Case:**

The graph has 0 edges, there is no tour to find.

**Induction Hypothesis:**

EULER outputs an Eulerian Tour for any even degree, connected graph with at most  $m \geq 0$  edges. This is a strong induction!

**Induction Step:**

Suppose  $G$  has  $m + 1$  edges. Removing the edges of  $T$  from  $G$ , since the tour would remove pairs of edges from vertices, we would be left with a remaining even-degree, connected graph with less than  $m$  edges.

This  $T$  intersects each of the edge disjointed but connected subtours  $G_i$ , at some first vertex of intersection.

The induction hypothesis allows us to assume these resultant subtours to have Eulerian tours. Therefore, our end result of EULER is in fact a splice of individual Eulerian tours, into a large tour whose union is all edges of  $G$  and thus also an Eulerian tour.

Such is the general description of this proof.

In this case, since the graph of Seven Bridge is not even-degreed, it cannot have an Eulerian Path. Therefore, the solution is impossible.

## 5.3 Planarity, Euler's Formula, Coloring

### 5.3.1 Trees

Trees are defined as connected acyclic graphs, with some interesting properties upon observation:

- On a tree, any node can be set as a root.
- Any degree-1 node on a tree would be called a “leaf” node.
- A tree is in fact bipartite. We will talk about this later.
- A tree is always planar. We will talk about this later.
- Adding an edge to anywhere of a tree will make it cyclic.

### 5.3.2 Planar Graphs

What is a planar graph?

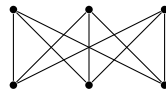
**Definition 5.3.1. Planarity**

A graph is planar if it can be drawn on a plane without crossings.

In that case, all trees are planar.

**Symbol 5.3.1. Three Houses-Three Wells**

The graph  $K_{3,3}$  is a graph where there are two sets of vertices, each of size three, and all edges between the two sets of vertices are present:



Shown in above is an example of a nonplanar graph,  $K_{3,3}$ .

Another example of a nonplanar graph would be a complete graph with five nodes, notated as  $K_5$ . A complete graph is a graph where every possible edge is present.

And from here, let us introduce another useful notion in the discussion of graphs by its edges:

#### Definition 5.3.2. Bipartite Graphs

A bipartite graph,  $G = (V, E)$ , is a graph where vertices are split into two groups and edges only exist between groups.

Mathematically, let the groups of vertices be  $L, R$ , then  $V = L \cup R$  and  $E \subseteq L \times R$ .

### 5.3.3 Euler's Formula

Planar graphs, when drawn on the plane, separates a plane into multiple regions. How do we express this fact mathematically?

#### Definition 5.3.3. Faces

The faces of a graph are the regions which the graph subdivides the plane into.

Faces are distinguishable for a planar graph, but not much so for nonplanar graphs!

Euler has contributed a great formula for planar graphs, which are great generalizations of polyhedra.

#### Theorem 5.3.1. Euler's Formula

Let  $v, f, e$  be respectively the amount of vertices, faces, and edges in a connected planar graph.

Then,  $v + f = e + 2$ .

And let's proceed with a proof:

#### Example Question 5.3.1: Prove Euler's Formula

We can perform mathematical induction on value of  $e$ .

**Base Case**  $e = 0$ :

In this case, since there is no edge, the amount of vertex and face are all 1. The formula holds.

**Induction Hypothesis** For all connected planar graphs, the prompt holds.

**Induction Step** Let's consider the following cases:

- If it is a tree, then the amount of faces is 1 and the amount of edge  $e$  in a tree is equal to  $v - 1$ .
- If it is not a tree, find a cycle of the graph and delete any edge of the cycle to break it. This will reduce both  $e$  and  $f$  by one. By induction hypothesis, formulas work in the smaller graph, and the arithmetics we just performed would prove the formula true in a graph with  $e - 1 + 1 = e$  edges.

For a non-connected planar graph, the situation is slightly different.

#### Theorem 5.3.2. The Sparsity of Graph

Take its connected subpart, a planar graph that has  $f > 1$  faces, and count the number of sides on each face. In

the end, as we double count each side, we double count each edge of the graph, resulting in the conclusion:

$$\sum_{i=1}^f (\text{Number of sides in face } i) = 2e$$

Knowing that there cannot be parallel edges between two same nodes, and assuming that the graph has at least two edges (to secure at least three vertices), then each face has at least three sides.

Number of sides in face  $i \geq 3$

$$\sum_{i=1}^f 3 \leq \sum_{i=1}^f (\text{Number of sides in face } i) = 2e$$

Since we are considering a connected planar graph, let us transform the above inequality according to Euler's Formula:

$$\begin{aligned} 2e &\geq 3f \\ 2e &\geq 3(e + 2 - v) \\ 2e &\geq 3e - 3v + 6 \\ e &\leq 3v - 6 \end{aligned}$$

This implies that planar graphs are sparse, or, they cannot have too many edges. Therefore, specific graphs whose number of edges and vertices do not follow the above inequality can be seen as non-planar:

#### Definition 5.3.4. Planarity

If a graph is a connected planar graph,  $e \leq 3v - 6$ .

Mind that this does not serve as a planarity test, as such property is satisfied by  $K_{3,3}$  as well.

However, its contrapositive can serve as a non-planarity test.

In that case, both  $K_5$  and  $K_{3,3}$  are non-planar. We can explore these graphs further: these are the only non-planar graphs, made precise by the mathematician Kuratowski (where the capital K came from).

#### Theorem 5.3.3. Kuratowski's Theorem

A graph is non-planar iff it contains  $K_5$  or  $K_{3,3}$ .

Let us define contain as where the nodes in the graph  $K_5$  or  $K_{3,3}$  are identifiably connected as their corresponding graph rules via paths, such that no two of those paths share vertices.

In otherwords, let the graphs noted in this theorem be a subgraph of the larger inspected graph.

This mathematical result is sometimes also known as Kuratowski's theorem, and its direction imply that if a graph contains one of these two non-planar graphs, itself must also be non-planar.

### 5.3.4 Duality and Coloring

The Greeks knew everything, from the section above to the next notion:

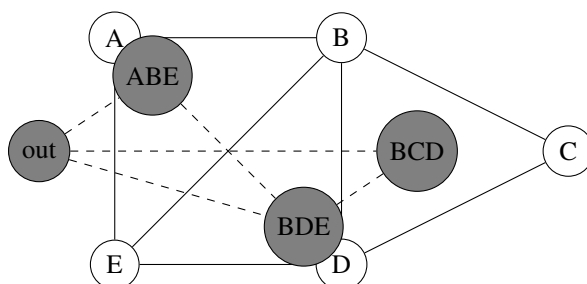
#### Definition 5.3.5. Dual Graph

The dual graph  $G^*$  of a graph  $G$  is a graph formed by converting faces of  $G$  into vertices and drawing edges between adjacent faces.

And the CS70 Lecture Note will ask you to *think about it* ::::D. So do it, maybe search a few Google Images up. I can't LaTeX this. Skill issue.

But Brandon, **I hear your thoughts**. Why did you put a LaTeX-ized dual graph on the screen already

Figure 5.3.1. Dual Graph



and why is Duality important at all?

That's because it reduces area-coloring map into a planar graph. Since a map's faces ought to be distinguishable, its dual graph can serve as a planar graph for us to work with.

Now let's talk about coloring. Essentially, two-coloring a graph is working with a bipartite graph, where one group of vertex has one specific color, the other group the other color.

Or, we may proceed with a breadth-first algorithm that colors any uncolored neighbor with alternative colors. This will let us discover cycles of odd lengths in the event of failing to color.

Notably, we have the four color theorem, which states that you may graph any map with four colors. This implies that any planar graph can be vertex-colored in four or more colors!

Furthermore, there are in fact two types of coloring: vertex coloring and edge coloring.

To color a graph by some component  $X$  means all adjacent components of type  $X$  must not have a same color.

## 5.4 Classes of Graphs

### 5.4.1 Complete Graphs

Starting with a definition,

#### Definition 5.4.1. Complete Graph

A complete graph contains all possible edges between its vertices.

A complete graph on  $n$  vertex is unique and denoted as  $K_n$ , and has  $\frac{n(n-1)}{2}$  edges.

### 5.4.2 Trees

Starting with a definition... well, in fact, many definitions.

#### Definition 5.4.2. Trees

Trees can be defined as any of the following:

- A graph that is connected and acyclic (as spoken in last section).
- A graph that is connected and has  $n - 1$  edges (this guarantees it to be acyclic).
- A graph that is connected, and the removal of an arbitrary edge disconnects the graph.
- An acyclic graph that becomes cyclic upon the addition of edge on any possible position.

Why do we want trees? I don't know, maybe retaking CS61B will provide you some insights.

In a rooted trees, which algorithms are mainly concerned with, there is a designated node called the root at the topmost level of a tree. The bottom-most nodes are called leaves, and the intermediate nodes are called internal nodes.

In a rooted tree, a root should never be a leaf; and throughout types of trees, leaves ought to be degree-1 vertices.

But why does this matter? That's a nice question. James. I have been tested on this when practicing Leetcode. I think that shows something.

### 5.4.3 Hypercubes

We would like a graph to be strongly connected, and a complete graph is such example. But, complete graphs take a lot, a lot of edges. Since graphs are abstractions of real-life objects, this implies huge operational and development costs. Our alternative solution would be a hypercube.

#### Definition 5.4.3. Hypercube

A hypercube is a graph whose vertex set is given by:

$$V = \{0, 1\}^n$$

which is equivalent of the set of all  $n$ -bit strings.

The edge set is then defined as:

$$E = \{x, y : x \text{ and } y \text{ differ in one bit}\}$$

Such a hypercube is also called an  $n$ -dimensional hypercube.

Notably, we may define hypercube recursively.

The vertex set of a  $n - 1$  dimensional hypercube happens to be  $\{0x | x \in V_{n-1}\}$ . The  $n$ -dimensional hypercube is obtained by placing an edge between each pair of vertices in the 0-subcube, which we described via an  $n - 1$  dimensional hypercube, and the 1-subcube where every leftmost digit of such vertex group is 1 instead!

Well, let's take a further investigation into its connectivity, so to make sure that we are working with the alternative solution we have wanted:

#### Lemma 5.4.1. Total Number of Edges in Hypercube

By the definition of hypercube,  $E(n) = 2E(n - 1) + 2^{n-1}$ , where  $E(1) = 1$ .

Using mathematical induction, we may show that  $E(n) = n2^{n-1}$ .

In a simpler approach, since the degree of each vertex is  $n$  and  $n$  bit positions can be flipped in any of the  $2^n$  vertices, there is a total of  $n2^n$  edges (non-distinct) that we can count between cubes.

We counted each edge twice, so we divide that total number by 2 to attain the result that the total number of edges in a  $n$ -dimensional hypercube is still  $n2^{n-1}$ .

#### Theorem 5.4.1. The number of discarded edges from disconnected vertices

Let  $S \subseteq V$  such that  $|S| \leq |V - S|$ , and let  $E_S$  denote the set of edges connecting  $S$  to  $V - S$ :

$$E_S := \{\{u, v\} \in E : u \in S \wedge v \in (V - S)\}$$

Then, it must be that  $|S| \leq |E_S|$ .

#### Example Question 5.4.1: Prove the above theorem

Let us perform mathematical induction on the dimension of hypercube.

**Base Case:**  $n = 1$

We may only separate 1-dimensional hypercubes into one vertex and another. Either way, there will be only one edge connecting  $S$  and  $V - S$ , and the cardinality of such two sets would be 1. Theorem holds.

**Induction Hypothesis:** Assume claim holds for  $1 \leq n \leq k$ ; we perform strong induction.

**Induction Step:** Prove the claim for  $n = k + 1$ .

Let us assume that among the two separated sets, it would be  $|S| \leq 2^k$ .

Let  $S_0$  represent the vertices from 0-subcube in  $S$ , and respectively for  $S_1$ . Then, either  $S$  has a fairly equal intersection size with the two subcubes, or it does not.

*Case 1:*  $|S_0| \leq 2^{k-1} \wedge |S_1| \leq 2^{k-1}$

The induction hypothesis works on both subcubes, which means there are at least  $|S_0|$  edges between  $S_0$  and its complement, and similarly for  $|S_1|$  and  $S_1$ . If so, the total number of edges between  $S$  and  $V - S$  would be at least  $|S_0| + |S_1| = |S|$ , as the theorem attempts to claim.

*Case 2:*  $|S_0| \geq 2^{k-1}$

In this case, we cannot apply the induction hypothesis as wanted before on the 0-subcube, but only on the 1-subcube. But we can apply this onto  $|V_0 - S_0|$  instead and infer the amount of edges between them is at least  $2^k - |S_0|$ . The current total of edges is calculated as  $2^k - |S_0| + |S_1|$ .

The edges that cross between the 0-subcubes and 1-subcubes would be at least  $|S_0| - |S_1|$  because there is an edge between every pair of vertices  $(0x, 1x)$ . Therefore, finally, the number of edges crossing is indeed at least  $2^k \geq |S|$  as desired.

(The  $|S_0| - |S_1|$  was derived from the cardinality of complement from  $\{x : \overline{0x} \in S_0\} - \{x : \overline{1x} \in S_1\}$ ).

## 5.5 Personal Regard on This Topic

This is probably the most abstract topic you'll encounter in CS70, just because everyone is new to graphs. Here is also an organization of results we have proven during discussion and homework:

- Every even-degreed connected graph has an Eulerian Tour, a tour on which every edge is used. (Eulerian Theorem)
- If a graph contains  $K_{3,3}$  or  $K_5$ , it is non-planar. (Kuratowski)
- For a connected planar graph,  $v + f = e + 2$ . (Euler's Formula)
- For a connected planar graph,  $e \leq 3v - 6$ . This inequality is even tighter on disconnected planar graphs. (Proven in Lecture)
- Every planar graph can be vertex-colored with four colors. (Four Color Theorem)
- Hypercubes of  $n$  dimension has  $n2^{n-1}$  edges. (Proven in Lecture)
- The number of vertices in an undirected graph with odd-degree is even. (Discussion 2B)
- An  $n$ -dimensional hypercube can be edge-colored with  $n$  colors. (Discussion 3A)
- Every hypercube is bipartite. (Discussion 3A)
- A Hamiltonian Tour is a tour where each vertex appears only once, each pair of consecutive vertices are connected by an edge, and the first and last vertex of this path are connected by an edge. (Definition)

## Chapter 6

# Modular Arithmetics

In this chapter, we will discuss a very foundational aspect of modular arithmetics that supports a wide majority of computational methods. This aspect of mathematics is modular arithmetics: the arithmetics of remainders and divisions widely researched, boding influence from number theory, and potentially linear algebra!

### 6.1 Introduction to Modular Arithmetics

#### 6.1.1 Motivation

Across several computational notions, we would utilize number lines that wrap around. What do we exactly mean? Think back to Project Enigma, where we increment letters by numbers, and we would like to increment the letter Z by one position. That would provide us A, but because Z is computational represented as number 25, and A 0, we need to wrap a number around back to 0 whenever it exceeds 26.

The mathematical operation that allows it is called a number circle.

On a number line, we proceed from one number to the next, until perhaps infinity; in a number circle, counting works like the hours of a clock: whenever we reach the maximum possible hour (12), we notate the next hour as the minimum possible hour (1). Essentially, the next number to the largest of a circle is the smallest number of a circle.

Mathematically, such arithmetics are known as **Modular Arithmetics**.

#### 6.1.2 Foundation

Let us take alphabets as an example. To computationally represent them, we would use the integers in the following range:

$$\{0, 1, \dots, 24, 25\}$$

There is a one-to-one correspondence of the 26 digits towards 26 alphabets. To consider alphabets whose representation is temporarily larger than the maximum of this range (number circle), 25, since there are a total of 26 consecutive integers from 0 in the range: **the true numerical representation of a letter is the remainder of division from its current representation divided by the size of range**.

##### Symbol 6.1.1. Modulo

The Modulo operator,  $(\text{mod})$ , is used in the manner:

$$x \ (\text{mod} \ m) = \text{remainder of } x \text{ divided by } m$$

For example:

$$26 \ (\text{mod} \ 26) = 0$$



And as a review of the rules of remainder, let  $r = x \pmod{m}$ , it is ruled that  $0 \leq r \leq m-1 \wedge r \in \mathbb{Z}$ .  
Therefore, symbol wise:

$$29 \pmod{26} = -23 \pmod{26} = 3$$

#### Symbol 6.1.2. Properties of $\pmod{\phantom{x}}$

- $(a+b) \pmod{n} = [a \pmod{n} + b \pmod{n}] \pmod{n}$
- $(a-b) \pmod{n} = [a \pmod{n} - b \pmod{n}] \pmod{n}$
- $(a \times b) \pmod{n} = [a \pmod{n} \times b \pmod{n}] \pmod{n}$

### 6.1.3 Set Representation

For a modulo  $k$ , it happens that we can categorize all integers into  $k$  sets  $S_i$ , such that:

$$S_i = \{z : z \pmod{k} = i\}$$

These sets are called residue classes  $\pmod{k}$ . And, to mathematically express their comraderies, we may use the following symbol:

#### Symbol 6.1.3. Congruency

$$x \equiv y \pmod{k} \iff (x \pmod{k} = y \pmod{k})$$

For example,

$$29 \equiv -23 \pmod{26}$$

Let us attempt to prove a relevant theorem:

Exercise Reader! Exercise! 6.1.1: Prove this theorem in the following box

Prove:

$$(a \equiv c \pmod{m} \wedge b \equiv d \pmod{m}) \implies (a+b \equiv c+d \pmod{m} \wedge ab \equiv cd \pmod{m})$$

Let  $c = a + k \cdot m, d = b + l \cdot m$ :

$$\begin{aligned} c+d &= a+b+(k+l) \cdot m \\ &\equiv (a+b) \pmod{m} \end{aligned}$$

$$\begin{aligned} cd &= ab + alm + bkm + klm^2 \\ &= ab + (al + bk + klm) \cdot m \\ &\equiv (ab) \pmod{m} \end{aligned}$$

## 6.2 Exponentiation

How can we compute  $x^y \pmod{m}$ ?

Well, for instance, we can compute the sequence  $x \pmod{m}, x^2 \pmod{m}, \dots$  up to the  $y^{th}$  term, but that takes time exponential to the number of bits in  $y$ , which is less efficient than wanted.

Instead, let us not calculate all items of the sequence, but just some:

**Theorem 6.2.1. Algorithm for Modular Exponentiation**

We solve the value of  $x^y \bmod m$  with the following procedure:

- If  $y$  is 0, then return 1. This is trivial solution!
- Else, we'd use the following procedure:
  - Compute a new variable  $z = x^{\frac{y}{2}} \bmod m$ .
  - If  $y$  is even, then return  $z^2 \bmod m$ .
  - Else, return  $z^2 x \bmod m$

Essentially, we use the mathematical facts that:

$$\begin{cases} x^{2a} = (x^a)^2 \\ x^{2a+1} = x \times (x^a)^2 \end{cases}$$

The concrete analysis of the above algorithm would be performed in CS170, but this algorithm holds essentially an  $O(n)$  runtime, where  $n$  is the number of bits in  $y$ .

This is very efficient, much more efficeint than the exponential-time naive approach.

## 6.3 Bijections

Take this as a side track!

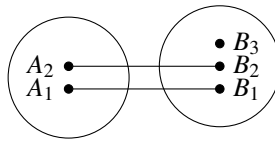
A function is, in its nature, the relationship of sets: it is a relationship of elements from set  $A$  to set  $B$ . Mathematically expressed,

$$(\forall x \in A)(f(x) \in B), f : A \rightarrow B$$

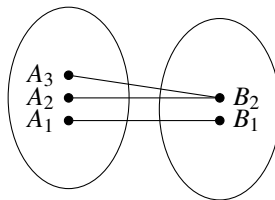
**Definition 6.3.1. Injective, Surjective, Bijection**

Injective, Surjective, and Bijective are all properties of functions, relationship between two sets  $A$  and  $B$ . We call the  $A$  the set of *pre-images*, and  $B$  the set of *images*.

- **Injective:** Each image (output) has at most one pre-image (input). This type of function is thus called a "one-to-one" function.



- **Surjective:** Each image (output) has at least one pre-image (input). This type of function is "onto".



- **Bijective:** A function is Bijective if the function is both injective and surjective. In other words, each image (output) has exactly one pre-image (input) and vice versa.

Let us practice this concept with some modular functions:

#### Example Question 6.3.1: Bijectivity of Modular Functions

Show whether the following functions are bijective:

$$\begin{cases} f(x) \equiv x + 1 \pmod{m} \\ g(x) \equiv 2x \pmod{m} \end{cases}$$

Where both functions map  $\{0, \dots, m-1\}$  onto themselves.

The function  $f$  is bijective: every image, which are nonnegative integers from 0 to  $m-1$ , have a corresponding pre-image that is themselves.

The function  $g$ , however, would suffer under an even  $m$  as odd nonnegative images will have no preimages. In that case,  $g$  is only bijective for an odd  $m$ .

## 6.4 Inverse Operations

### 6.4.1 Multiplicative Inverse

We have discussed adding as an inverse operation of subtraction, but what is the inverse operation of multiplication in a modulo sense?

In normal arithmetics, multiplying a number by  $x$  is equivalent of dividing it by  $\frac{1}{x}$ , and for an equation  $xy = 1$ , we would call  $y$  the multiplicative inverse of  $x$  (because multiplying  $y$  reduces the product into 1). In linear algebra, we also had "inverse matrices" that can form multiplicative inverses via yielding a product of identity matrices  $I_n = A^{-1} \cdot A$ .

In a modular sense, then, the multiplicative inverse would be defined as follows:

#### Definition 6.4.1. Multiplicative Inverse of Modular Arithmetics

For  $x, y, m$  such that  $xy \equiv 1 \pmod{m}$ ,  $y$  is the multiplicative inverse of  $x$  modulo  $m$ .

There turns out to be conditions under which the multiplicative inverse of a number  $x$  modulo  $m$  exists:

#### Theorem 6.4.1. Condition of Existence for Multiplicative Inverse in Modular

**Theorem:** The necessary condition of multiplicative inverse of  $x$  modulo  $m$  is that  $\gcd(m, x) = 1$ .

If there exists an inverse  $y$  of  $x$  modular  $m$ , then  $xy \equiv 1 \pmod{m}$ . Therefore,  $xy = km + 1$ .

In that case,  $xy - km = 1$ .

Let  $c$  be a common divisor of  $x$  and  $m$ , then that divisor must be the factor of  $xy$  and  $km$  too, thus the factor of  $xy - km = 1$ . However, the only number that can be a factor of 1 is 1.

Therefore, assuming there exists a multiplicative inverse, then the only possible common divisor of  $x$  and  $m$  is 1. Therefore,  $m$  and  $x$  are also coprime to each other.

And the multiplicative inverse would be unique:

#### Theorem 6.4.2. Uniqueness of Multiplicative Inverse under Modular

**Theorem:** The multiplicative inverse of  $x$  modular  $m$  is unique.

Let us state  $a$  as the multiplicative inverse of  $x$  modular  $m$ , then  $ax \equiv 1 \pmod{m}$ .  
 For the sake of contradiction, let's assume another distinct multiplicative inverse  $b$ , such that  $bx \equiv 1 \pmod{m}$ .  
 In that case,  $bx - ax = x(b - a) \equiv 1 \pmod{m}$ , so if two or more distinct multiplicative inverse of  $x$  modular  $m$  exist, then  $x$  and  $m$  are not coprime to each other. This causes a contradiction with the condition of existence, which states that such inverses can only exist for coprime  $x$  and  $m$ .  
 Therefore, by contradiction, there can only exist one multiplicative inverse of  $x$  modular  $m$ .

While the inverse of  $x$  can be written as  $y = x^{-1} \pmod{m}$ , this is generally considered an abuse of notation, as  $x^{-1}$  can also stand for  $\frac{1}{x}$  under such ambiguous context.

Another way to express the multiplicative inverse can be  $y = (x)_m^{-1}$

But we cannot just settle with knowing that an inverse exist. Let us attempt at finding a way to compute it; namely, let us produce an algorithm that can compute the multiplicative inverse of  $x$  modulo  $m$ .

Let us begin our passage by thinking about greatest common divisors. Suppose that for any pair of numbers  $x, y$ , then the greatest common divisors can be expressed as:

$$d = \gcd(x, y) = ax + by$$

Let's express the relationship of numbers  $m$  and  $x$  when there exists a multiplicative inverse of  $x$  modulo  $m$ , in the above mathematical format:

$$1 = \gcd(m, x) = am + bx$$

This would imply that  $bx \equiv am + bx \equiv 1 \pmod{m}$ , such that  $b$  is a multiplicative inverse of  $x$  modulo  $m$ .

Euclid's Algorithm, which is for computing greatest common divisors, helps us to find the integers  $a$  and  $b$  in the above equation. Thus, let's explore Euclid's Algorithm as a method of computing modulo inverses!

### 6.4.2 Euclid's Algorithm

The Euclid's Algorithm, which is more of a folk algorithm by ancient engineers, relies on the theorem below:

#### Theorem 6.4.3. Reduction of GCD Computation via Modular Arithmetics

**Theorem:** Let  $x \geq y > 0$ , then  $\gcd(x, y) = \gcd(y, x \pmod{y})$

Let us express  $x = qy + r$ , where  $q \in \mathbb{Z}$  and  $r = x \pmod{y}$ . Note its an equal sign, not congruence.

If their greatest common divisor  $\gcd(x, y) = d$  divides  $x$  and  $y$ , then it also divides  $x$  and  $qy$ , thus also for  $r = x - qy$ .

Therefore,  $\gcd(x, y) = \gcd(y, x \pmod{y})$ .

The Euclid's Algorithm is the process to keep applying  $\gcd(x, y) = \gcd(y, x \pmod{y})$  until the second argument becomes 0. At termination, the first argument is the greatest common divisor of original inputs  $x$  and  $y$ .

Now, let's verify whether it works:

#### Theorem 6.4.4. Correctness of Euclid's Algorithm

Prove that Euclid's Algorithm correctly computes  $\gcd(x, y)$ .

Let us perform strong induction on the second argument (which should be the smaller of two inputs).

Let the proposition for functionality of Euclid's Algorithm be that:

$P(n)$  : The algorithm computes  $\gcd(x, n)$  for all  $x$  correctly, and  $x \geq y > 0$

**Base Case:**  $n = 0$

In this case,  $\gcd(x, 0) = x$ , which the algorithm correctly computes; meanwhile, the inequality stated in proposition regarding input size is maintained.

**Induction Hypothesis:** Assume  $P$  holds for all values  $k < y$ .

**Induction Step:** Now's the highlight, prove that  $P(y)$  holds.

According to the steps of Euclid's Algorithm,  $\gcd(x, y) = \gcd(y, x \bmod y)$ . This equality is proven in an above cube, and it is also guaranteed that  $y > x \bmod y$  by the nature of *mod*. Therefore,  $P(y)$  holds!

The runtime of Euclid's Algorithm is  $O(n)$ , where  $n$  is the total number of bits in input  $(x, y)$ . This is because for every two recursion call, the first initial input is reduced by at least one bit!

### 6.4.3 Extended Euclid's Algorithm

The extended Euclid's Algorithm provides us the coefficients  $a$  and  $b$  such that  $d = \gcd(x, y) = ax + by$ :

Exercise Reader! Exercise! 6.4.1: Extended Euclid's Algorithm

```
def extended_gcd(x, y):
    if y == 0:
        return (x, 1, 0)
    d, a, b = extended_gcd(y, x % y)
    return (d, b, a - x // y * b)
```

Let us now consider the inner workings of it: how the return for non-base-case functions?

Exercise Reader! Exercise! 6.4.2: Inner Working of EEA

To establish some foundation,  $d = \gcd(x, y) = ay + b(x \bmod y) = Ax + By$ .

We will need to find the expressions for new coefficients  $A, B$ , in terms of previous coefficients  $a, b$  throughout the calls of recursions.

$$\begin{aligned} d &= ay + b(x \bmod y) \\ &= ay + b(x - \lfloor \frac{x}{y} \rfloor y) \\ &= bx + (a - \lfloor \frac{x}{y} \rfloor)y \\ A &= b \\ B &= a - b \cdot \lfloor \frac{x}{y} \rfloor \end{aligned}$$

Or, instead of having to compute the weird arithmetics of new second coefficient, just alternate the secondary coordinate  $b$  into the primary of the next traced case (as algorithm demonstrates), and autonomously find a value of new  $b$  that works to satisfy  $ax + by = d$ .

This usually works for exams because  $d$  is usually 1. If it is not and the situation is complicated, stick to the algorithm.

### 6.4.4 Brief Notes: Division via Inverse

We have learned that to find the multiplicative inverse of  $x$  modulo  $m$ , we would just need to solve  $b$  from the equation  $1 = ak + bn$ , where  $k = \max(x, m)$  and  $n$  the other.

Let us attempt to apply this onto solving a congruence:

Example Question 6.4.1: Solve the following congruence using modular inverses

Solve:  $8x \equiv 9 \pmod{15}$

The inverse of 8 modulo 15 is 2. How?  $1 = -1 \times 15 + 2 \times 8$

So, let us multiply both sides of the congruence by  $8^{-1} \pmod{15}$ :

$$x \equiv 9 \times 2 \equiv 3 \pmod{15}$$

## 6.5 Chinese Remainder Theorem

The Chinese Remainder Theorem is another craftsman algorithm developed by Sunzi for counting the number of soldiers. The simplest case of the theorem is as follows:

Definition 6.5.1. Chinese Remainder Theorem, Claim

For  $m, n$  with  $\gcd(m, n) = 1$ , there is exactly one  $x \pmod{mn}$  such that:

$$x \equiv a \pmod{n} \wedge x \equiv b \pmod{m}$$

– Sunzi, *The Art of War*

Example Question 6.5.1: Prove the Simplest Case of Chinese Remainder Theorem

The proof follows from the existence of multiplicative inverse for  $n, m$ , respectively modulo  $m$ , modulo  $n$ . By prior discussions, this holds when  $\gcd(n, m) = 1$ .

Let us follow from this fact and declare some variables:

$$\begin{cases} u \equiv m(m^{-1} \pmod{n}) \\ v \equiv n(n^{-1} \pmod{m}) \end{cases}$$

Here, for analysis on  $u$ :

$$\begin{aligned} m(m^{-1} \pmod{n}) &\equiv mm^{-1} \equiv 1 \pmod{n} \\ u &\equiv 0 \pmod{m} \text{ (because it is a multiple of } m) \end{aligned}$$

Applying a similar logic onto the value of  $n(n^{-1} \pmod{m})$ , we will attain the following:

$$\begin{cases} u \equiv 1 \pmod{n} \wedge u \equiv 0 \pmod{m} \\ v \equiv 0 \pmod{n} \wedge v \equiv 1 \pmod{m} \end{cases}$$

And therefore, let  $x = au + bv$ :

$$\begin{cases} x = au + bv \equiv a(1) + b(0) \equiv a \pmod{n} \\ x = au + bv \equiv a(0) + b(1) \equiv b \pmod{m} \end{cases}$$

Such a solution can be argued unique:

Exercise Reader! Exercise! 6.5.1: Might the solutions of CRT be unique?

Consider for the sake of contradiction a second distinct solution  $y$  exists, then due to the modular properties,  $x - y \equiv 0 \pmod{m}$  and  $x - y \equiv 0 \pmod{n}$ .

Since  $\gcd(m, n) = 1$ , and  $x - y$  is some multiple of  $m, n$ , and  $x \neq y$ , then either  $|x - y| = 0mn$  or  $|x - y| \geq mn$ . The former implies  $x = y$ , while the latter suggests that one of  $x$  and  $y$  is not in  $\{0, \dots, mn - 1\}$ .

And therefore, it's not exactly that the solution is unique. It's that distinct solutions must have a difference that is a multiple of  $mn$ . In other words, it potentially has infinite solutions.

We have now figured that CRT helps to solve systems of congruences, and now we are able to solve systems of 2 congruences. What is left to us is to solve systems of more congruences.

And so, let us now move onto the more complex cases of CRT:

#### Theorem 6.5.1. Chinese Remainder Theorem

Let  $n_1, n_2, \dots, n_k$  be coprime natural numbers, then for any sequence of integers  $a$ , there is a unique integer  $0 < x < N = \prod_{i=1}^k n_i$  such that:

$$\begin{cases} x \equiv a_1 \pmod{n_1} \\ \vdots \\ x \equiv a_i \pmod{n_i} \\ \vdots \\ x \equiv a_k \pmod{n_k} \end{cases}$$

where

$$b_i = \frac{N}{n_i} \left( \frac{N}{n_i} \right)^{-1}_{n_i}, \quad x \equiv \left( \sum_{i=1}^k a_i b_i \right) \pmod{N}$$

– Sunzi, *The Art of War*

Exercise Reader! Exercise! 6.5.2: Prove the Existence of Solution for CRT

The proof of solution existence follows a very similar framework to part (b).

Let us first abbreviate an expression in terms of a function:

$$f(b) = \prod_{j=0, j \neq b}^k n_j = \frac{N}{n_b}$$

Let there be values:

$$x_i = a_i f(i) \cdot f(i)_{n_i}^{-1}$$

Then, following the similar logic of part (b), we can dictate that the solution would be:

$$x = \sum_i x_i$$

As for why the above works:

For each congruence under modulo  $n_i$  we deal with, the only term left to interpretation due to not being a multiple of  $n_i$  would be  $x_i$ ; however, due to how the term  $x_i$  is constructed via involving a number and its multiplicative inverse under  $n_i$ :

$$x_i \equiv a_i \pmod{n_i}$$

The uniqueness of this solution follows the fact that, for the sake of proof, let  $y$  be a distinct solution to the system from  $x$ :

$$\forall n_i (x - y \equiv 0 \pmod{n_i})$$

This follows from the observation that they are both equivalent to a same value under any modulo  $n_i$  to qualify as a solution.

But, the above claim would hint that  $x - y \equiv 0 \pmod{N}$ , stating that each solution is of distance  $N$  from each other, and thus there would be a unique solution if under modulo  $N$ .

Now, a linear algebra perspective:

#### Exercise Reader! Exercise! 6.5.3: The set of solutions is a Vector Space

Here, each  $b_i$  is congruent to 1  $\pmod{n_i}$  and 0  $\pmod{n_j}$  by the theorem, so a multiple of  $b_i$  can be used to satisfy congruence  $\pmod{n_i}$  while not interfering with other congruences of other modulo.

In that sense, we can construct multiple unit vectors from  $b_i$  across the indices of congruence to construct linear combinations that become a solution by itself. The set of all solutions is then that span.

This might be a convenient point of vectorization for computation.

The uniqueness condition is very similar: distinct solutions must have a distance that is a multiple of  $N$  as defined in theorem.

## 6.6 Personal Regard on This Topic

Just like my friend who plays Smash Ultimate says when I first tried to play against a tournament player who plays an overpowered character:

My advice for you, is to have fun.

No, just kidding. Here are the advises I meant to put instead:

- Be aware of the timings at which you use  $\equiv$  and  $=$ . These are different. If you are using *mod* as an operator, note it clearly.
- The proof for CRT is useful, and the reason why its solution works associates extremely closely to its proof. This is applicable for polynomials, surprisingly.
- Extended Euclid Algorithm will save you in the next chapter.
- If you are like most of the readers and this author might be, the original name of CRT in its native language is in case you want to yell some move names during your midterm at Pimentel 1 to cheer yourself up when using CRT.



## Chapter 7

# Public Key Cryptography

In this chapter, we discuss public key cryptography: a convenient means of confidential conversation, and explore it as an application of modular arithmetics, which was the explained subject of previous chapter.

### 7.1 Introduction to Cryptography

Cryptography is the study that develops the most secure ways of delivering messages.

The basic setting for cryptography is typically described over a plot involving three characters:

- Crewmate A: Trying to communicate confidentially over a link.
- Crewmate B: Trying to communicate confidentially over a link.
- Sus person: Trying to listen in crewmates' conversations to do sussy things.

Or, to formalize what is a 'sus person', feel free to use the word 'spy' or 'evesdropper' too.

To let the communication be confidential and secure between crewmates, each crewmate will apply an Encryption function  $E$  to a message  $x$ , and send  $E(x)$  to the other crewmate. Upon the receipt, the other crewmate applies a decryption function  $D$  onto  $E(x)$  such that  $D(E(x))$  provides information resembling the original message  $x$ . In most cases,  $D(E(x)) = x$ .

These functions  $E$  and  $D$  are not necessarily single-variable. We will touch on this soon. But, without knowing anything about the inner workings of  $E$  and  $D$ , the 'sus person' cannot know what the original function  $x$  at all.

For many centuries, humans study **private-key protocols**, which also requires a codebook, which contains all future correspondence between encrypted and decrypted messages. This requires a codebook, to begin with.

On the other hand, **public-key schemes** allow crewmates to communicate, encrypt, and decrypt without the need of codebook. This seems off. If crewmate  $B$  and the sus person both have equal information on the encrypted message  $E(x)$ , are they not equally capable of solving the encryption?

The essence of **public-key schemes** hides in implementing a digital lock, which only the decrypter can have a key for. Therefore, the decrypter send the encrypter a method of sending secured message that the decrypter has the key for. While the sus person can attempt finding the key, it should be difficult to do so, so arduous that it is almost impossible.

That's the point of cryptography.

Now, since we also need a key to decrypt the message, the function  $D$  is indeed multivariate rather than single-variable on the encrypted message  $E(x)$ .

## 7.2 Cryptography in Modular Arithmetics Theorems

### 7.2.1 The RSA Scheme of Cryptography

The RSA scheme is based heavily on modular arithmetics.

Let  $p$  and  $q$  be two large primes (typically with 512 bits, so larger than  $2^{511}$ ).

Let  $N = pq$ . Messages to the decrypter are numbers modulo  $N$ , excluding 0 and 1.

Then, let  $e$  be any number relatively prime to  $(p-1)(q-1)$ .

The public key of decryption is the pair of numbers  $(N, e)$ . While this key is published to the whole world, the numbers  $p$  and  $q$  are not public.

The private key of decryption, meanwhile, is  $d = e^{-1}_{(p-1)(q-1)}$ , which would exist by the definition of  $e$ .

Now, let us explore the encryption and decryption functions:

#### Definition 7.2.1. Encryption and Decryption functions for RSA Scheme

Provided the values:

$p, q$  are large prime numbers

$$N = pq$$

$e$  is coprime with  $(p-1)(q-1)$

$$d = e^{-1}_{(p-1)(q-1)}$$

Among these values, only  $N$  and  $e$  are known.

**Encryption:**  $E(x) \equiv x^e \pmod{N}$

**Decryption:**  $D(y) \equiv y^d \pmod{N}$

#### Example Question 7.2.1: Demonstration of RSA Scheme

Let  $p = 5$ ,  $q = 11$ , so  $N = pq = 55$ . In this case,  $(p-1)(q-1) = 40$ . For a number coprime with it, let's choose  $e = 3$ .

The public key of this scheme is  $(N, e) = (55, 3)$ , and the private key is  $d \equiv 3^{-1}_{40} \equiv 27$ .

For any message  $x$  that the encrypter will send, the encryption of  $x$  is  $y \equiv x^e \pmod{55}$ . Meanwhile, the decryption of  $y$  is  $x \equiv y^d \pmod{55}$ .

Therefore, for a message  $x = 13$ , the encryption is  $y = 13^3 \equiv 52 \pmod{55}$ , and the decryption of  $y$  is  $x = 52^{27} \equiv 13 \pmod{55}$ .

### 7.2.2 The Supporting Pillars: Theorems

Let us first introduce a significant theorem of modular arithmetics:

#### Theorem 7.2.1. Fermat's Little Theorem

**Theorem:** For any prime  $p$  and any  $a \in \{1, 2, \dots, p-1\}$  (as long as  $a$  is not a multiple of  $p$ ; in other words, coprime with the prime  $p$ ),  $a^{p-1} \equiv 1 \pmod{p}$ .

Let  $S$  denote the set of possible results from mod  $p$ .

Provided that  $p$  and  $a$  are coprime,  $\gcd(p, a) = 1$ .

Consider the set of numbers  $\{a, 2a, \dots, (p-1)a\}$ . Because  $a$  and  $p$  are coprime, for two numbers of the prior set to be congruent under  $\pmod{p}$ , it must be their multiple is the same.

Consequently, the following sets are equal:

$$\{1, 2, \dots, p-1\} = \{a \pmod{p}, 2a \pmod{p}, \dots, (p-1)a \pmod{p}\}$$

Following from that equality of sets:

$$1 \times 2 \times \dots \times (p-1) \equiv (p-1)! \pmod{p}$$

$$a \times 2a \times \dots \times (p-1)a \equiv a^{p-1}(p-1)! \pmod{p}$$

$$1 \times 2 \times \dots \times (p-1) = a \times 2a \times \dots \times (p-1)a$$

$$a^{p-1}(p-1)! \equiv (p-1)! \pmod{p}$$

Here, multiply both sides by  $(p-1)!^{-1}_p$

$$a^{p-1} \equiv 1 \pmod{p}$$

As a side note, the value  $(p-1)!^{-1}_p$  exists because  $p$  is prime and must therefore be relatively prime to the components of  $(p-1)!$ .

Now, let us ensure that the message will get encrypted and decrypted correctly:

#### Theorem 7.2.2. Correctness of Encryption and Decryption

**Theorem:** Under the above definitions of the encryption and decryption functions  $E$  and  $D$ ,

$$\forall x \in \{0, 1, \dots, N-1\} (D(E(x)) \equiv x \pmod{N})$$

Mathematically expressing the prompt:

$$\forall x \in \{0, 1, \dots, N-1\} (D(E(x)) \equiv (x^e)^d \equiv x \pmod{N})$$

By the definition of  $d = e^{-1}_{(p-1)(q-1)}$ ,  $ed \equiv 1 \pmod{(p-1)(q-1)}$ . Let us then express  $ed = 1 + k(p-1)(q-1)$ , leading to:

$$x^{ed} - x = x^{1+k(p-1)(q-1)} - x = x(x^{k(p-1)(q-1)} - 1)$$

Let us inspect whether  $pq|x^{ed} - x$ , or in other words,  $p|x^{ed} - x \wedge q|x^{ed} - x$ .

For the argument of divisibility with  $p$ , we consider two cases: the trivial case where  $p|x$  and the harder case where  $p \nmid x$ .

If  $x$  is not a multiple of  $p$ , then since  $p$  is prime, by Fermat's Little Theorem:

$$x^{p-1} \equiv 1 \pmod{p}$$

$$x^{p-1k(q-1)} \equiv 1^{k(q-1)} \pmod{p}$$

$$x^{k(p-1)(q-1)} - 1 \equiv 0 \pmod{p}$$

$$p|x^{k(p-1)(q-1)} - 1$$

If  $p|x$ , then the case is trivial. We are just subtracting two multiples of  $p$ .

We can prove under similar logic that  $q|x^{ed} - x$ , so in the end,  $pq|x^{ed} - x$ .

Consequently:

$$x^{ed} \equiv x \pmod{pq}$$

Chinese Remainder Theorem can also be used to provide a similar proof, where we construct a system of congruences under modulo  $p$  and  $q$  exploiting  $x^{ed} \equiv x$  under these modulus, and finally find  $x \pmod{pq}$  to be the only viable solution by the uniqueness property of CRT.

So, why is RSA secure? It is based on the assumption:

Given  $N$ ,  $e$ , and  $y \equiv x^e \pmod{N}$ , there is no efficient algorithm for determining  $x$ .

This is because to guess  $x$ , we would have to try on the order of  $N$  values; however  $N$  is the product of two 512-bit prime numbers, which makes this unrealistic.

Or, even if we factor  $N$  into its prime factors, this is still a majorly impossible problem to deal with computationally. So in summary, the security of RSA depends on the difficulty to solve a message from it.

### 7.2.3 Prime Number Theorem

Now, all that is left is for the crewmates to:

1. Find huge prime numbers  $p$  and  $q$ .
2. Compute exponentials mod  $N$ .

Both of which are not so shrimple.

But, to find large prime numbers, we can attempt to find efficient algorithms of determining whether a number is prime, and fortunately, a large portion of numbers are prime.

#### Theorem 7.2.3. Prime Number Theorem

**Theorem:** Let  $\pi(n)$  denote the number of prime numbers that are less than or equal to  $n$ .

$$\forall n \geq 17 \left( \pi(n) > \frac{n}{\ln(n)} \right)$$

For exponentiation, the exponentiation algorithm provided in Chapter 6 is actually good enough. So RSA is quite plausible for implementation and usage!

## 7.3 My Personal Regard on This Topic

Takeaway: Fermat's Little Theorem is goated.

Otherwise, make sure you understand that in RSA Scheme...:

- The known values are  $N$  and  $e$ , where  $e$  is coprime with  $(p-1)(q-1)$  and  $N = pq$ .
- The reason your algorithm makes the eavesdropper mald is because  $p$  and  $q$  are hard to find.
- If it works properly,  $D(E(x)) \equiv x \pmod{N}$ . You use the public key to encrypt and decrypt, it's the private key that depends on  $(p-1)$  and  $(q-1)$ .
- This algorithm can work if you find two or more prime factors to form  $N$ , and you'd just adjust the formula for  $e$  and  $d$  accordingly. But, if you just have a prime number  $N$  as your public key, your private key is very findable because it's just a multiplicative inverse of a known number under a known modulo.

If you're interested in cryptography, the mathematics of making hackers mald, make sure to tune in for cryptography classes!

# Chapter 8

## Polynomials

In this chapter, we discuss the application of polynomials and their mathematical properties, as well as a concise instance of its performance in cryptography.

### 8.1 Re-Introduction to Polynomial

Polynomial is a category of functions that are easy in their algebra expression while also widely applicable to various mathematical techniques. While we have already encountered polynomials in high school mathematics, let us define polynomial in a formal fashion that we might not have done before:

#### Definition 8.1.1. Polynomial

A polynomial of single variable is a single-variable expression with an associated function:

$$p(x) = a_dx^d + a_{d-1}x^{d-1} + \cdots + a_1x + a_0$$

Here, while the variable and coefficients are usually real numbers, the exponents need be integer.

An  $n$ -degree polynomial is a polynomial whose largest exponent is  $n$ .

The root of a polynomial is the value  $x$  for a polynomial function  $p$  such that  $p(x) = 0$ .

And following the existence of “roots”, we now discuss some properties of roots and polynomials.

#### Symbol 8.1.1. Properties of Polynomials

1. A non-zero  $n$ -degree polynomial has at most  $n$  roots.
2. Given  $n+1$  pairs  $(x_1, y_1), \dots, (x_{d+1}, y_{d+1})$  with all  $x_i$  distinct, there is a unique polynomial  $p(x)$  of degree  $n$  such that  $p(x_i) = y_i$  for  $1 \leq i \leq d+1$ .

The implication of property 2 would be that a line is define-able by two points.

### 8.2 Polynomial Interpolation

Polynomial Interpolation is the property that helps locate the  $n$ -degree polynomial described in property 2, and the name of the algorithmic method for it is Lagrange Interpolation:

## Exercise Reader! Exercise! 8.2.1: A Startpoint

Suppose that we will construct a polynomial  $p(x)$  where  $y_1 = k$ ,  $y_j = 0$ , where  $2 \leq j \leq d+1$ . Let us begin with an arbitrary  $d$ -degree polynomial:

$$q(x) = (x - x_2)(x - x_3) \dots (x - x_{d+1})$$

Here, since  $(x - x_j)$  is involved as a factor of  $q(x)$ , we can confirm  $q(x_j) = 0$ .

Meanwhile, as  $q(x_1) = (x_1 - x_2)(x_1 - x_3) \dots (x_1 - x_{d+1})$ , we may define  $p(x)$  as follows:

$$p(x) = \frac{kq(x)}{q(x_1)} \rightarrow p(x_1) = \frac{kq(x_1)}{q(x_1)} = k$$

Let us generalize the above example into the construction of polynomial  $\Delta_i(x)$ , the  $d$ -degree polynomial where  $y_i = k$ ,  $y_j = 0$ :

$$\Delta_i(x) = \frac{\prod_{j \neq i} (x - x_j)}{\prod_{j \neq i} (x_i - x_j)}$$

So, for a set of  $d+1$  points  $\{(x_1, y_1), \dots, (x_{d+1}, y_{d+1})\}$ , we will respectively produce  $\Delta_1(x), \dots, \Delta_{d+1}(x)$ . We can acquire the polynomial  $p$ :

$$p(x) = \sum_{i=1}^{d+1} y_i \Delta_i(x)$$

We have put our coefficients  $y_i$  back at this point, so the definition of  $\Delta(x)$  is easier to handle.

Let us discuss why such a  $p$  work. First of all,  $p$  is still  $d$ -degree, since we are just adding multiple  $d$ -degree polynomials together. Second of all, when evaluated at  $x_i$ , the  $\Delta_i$  polynomials that contain  $(x - x_i)$  will not contribute to the sum due to being zero, and will let  $\Delta_i(x)$  alone contribute!

The discussion of second property above will probably remind you of constructing solutions for Chinese Remainder Theorem!

## Theorem 8.2.1. Property 2 of Polynomials from Section 8-1

**Theorem:** Given  $d+1$  pairs  $(x_1, y_1), \dots, (x_{d+1}, y_{d+1})$ , with all  $x_i$  distinct, there is a unique polynomial  $p(x)$  of degree at most  $d$ , such that  $\forall i (p(x_i) = y_i)$ .

The existence of such a polynomial is secured by Lagrange Interpolation. We now just perform a separate proof for uniqueness.

For the sake of contradiction, assume the existence of  $q(x)$  such that  $q(x_i) = y_i$  for all listed pairs.

Then, consider a polynomial  $r(x) = p(x) - q(x)$ . Since  $p$  and  $q$  are, for the sake of contradiction as aforementioned, considered different polynomials,  $r(x)$  must be a non-zero polynomial of degree at most  $d$ . Property 1 suggests it can only have at most  $d$  roots.

However,  $\forall i \in 1, 2, \dots, d+1 (r(x_i) = p(x_i) - q(x_i) = 0)$ . In this case,  $r$  has at least  $d+1$  roots, producing a contradiction.

Therefore, there can only exist one polynomial to satisfy the listed points.

### 8.3 Re-Introducing Polynomial Division

For a polynomial  $p(x)$  of degree  $d$ , we may divide it by another polynomial  $q(x)$  of degree  $\leq d$  via long division:

$$p(x) = q'(x)q(x) + r(x)$$

This technique has been employed in Algebra 2 as well as integration in MATH 1A (or 1B, somewhere in the Calculus BC range to be sure).

We would define the quotient of such division  $q'(x)$ , and the remainder  $r(x)$ . The degree of  $r(x)$  must be smaller than

that of  $q(x)$ .

Instead of taking this space to re-introduce long division, let's leave that as an exercise for the reader to review high-school mathematics, and move on to the proof of Property 1 as listed in Section 8-1:

**Theorem 8.3.1. Property 1 of Polynomials from Section 8-1**

**Theorem:** A non-zero polynomial of degree  $d$  has at most  $d$  roots.

This proof can be constructed via the proof of two smaller claims.

**Claim 1:** If  $a$  is a root of  $p(x)$  with degree  $\geq 1$ , then  $p(x) = (x - a)q(x)$ , where  $q(x)$  is a  $d - 1$ -degree polynomial.

To prove this claim, we can use the property of remainders in polynomial division. If we divide  $p(x)$  by  $(x - a)$ , then the remainder of this division must be smaller than the degree of  $(x - a)$ , which would make the remainder a constant term  $c$ .

Meanwhile, since  $a$  is a root of  $p(x)$ , and  $p(x) = (x - a)q(x) + c$ ,  $p(a) = c = 0$ .

Therefore, there doesn't really exist a remainder.

**Claim 2:** A polynomial  $p(x)$  of degree  $d$  with distinct roots  $a_1, \dots, a_d$  can be written as  $p(x) = c(x - a_1) \dots (x - a_d)$ , where  $c \in \mathbb{R}$ .

Time for mathematical induction to be back!

Base Case:  $d = 0$ .

If  $p(x)$  is degree-0, then it is a constant  $c$  to be written in the above form.

Induction Hypothesis: The prompt holds for some  $d \geq 0$ .

Induction Step: prove the case for  $d + 1$ .

Let  $p(x)$  be a polynomial of degree  $d + 1$  with distinct roots, then Claim 1 supports that  $p(x) = (x - a_{d+1})q(x)$  for some degree- $d$   $q(x)$ .

For all  $i \neq d + 1$ ,  $p(a_i) = (x - a_{d+1})q(a_i) = 0$ . It will be helpful if we reference the induction hypothesis, notice that  $q(x)$  has  $d$  distinct roots, and thus is a polynomial of degree  $d$ .

Therefore,  $q(x) = c(x - a_1) \dots (x - a_d)$ , and following the product form of  $p$ :

$$p(x) = c(x - a_1) \dots (x - a_{d+1})$$

Claim 2 implies Property 1 if we can show that the only roots of  $p(x)$  are  $a_1, \dots, a_d$ , but that also can be confirmed from the factorization of such polynomial  $p(x)$ .

## 8.4 Interlude: Finite Fields

Property 1 and 2 of polynomial expression as noted above would both work when coefficients and bases are chosen complex or rational, other than the real cases we have assumed along the way. This is because the proofs of properties rely on legal arithmetics that receive arguments of one set (rational, complex... etc.) and return arguments of a same set.

This would imply that if we have chosen to conduct the proof of properties 1 and 2 via, say natural numbers and integers, then since we cannot secure the property of receiving and returning arguments of same set, the proofs cannot function.

Is there a workaround for natural numbers then? If we work with numbers modulo some prime number  $m$ , then legal arithmetics will still return us any number from  $\{1, \dots, m - 1\}$ , so properties 1 and 2 would still hold! So interestingly, properties 1 and 2 hold under a finite set like that of numbers modulo  $m$  as well.

We describe such situation as we are working over a finite field.

### 8.4.1 Counting

How many possible polynomials are there modulo  $m$ ? Well, for a  $n$ -degree polynomial who has  $m$  values to choose from per  $n$  coefficient, it would be  $m^n$ .

How many possible polynomials are there that satisfy  $k$  points modulo  $m$ ?

Well, for each point we neglect, we gain  $m$  possibilities for that empty coordinate that can bode any value of  $y$ . Therefore, considering  $n$  points, the total number of possible polynomials modulo  $m$  would be  $m^{k-n}$ .

## 8.5 Secret Sharing

Let us attempt to devise a secret sharing scheme such that, for the condition of secret decryption to be having a  $k$ -member faction:

1. Any group of  $k$  of all members can pool their information to configure the secret.
2. No group of  $k - 1$  or fewer members have any information about the secret, no matter how they pool their knowledge.

So now, suppose there are  $n$  officials indexed from 1 to  $n$ , and the secret is a natural number  $s$ . Let  $q$  be a prime number larger than  $n$  and  $s$ , let us work over the finite field of modulo  $q$ .

Pick a random polynomial  $P(x)$  of degree  $k - 1$  such that  $P(0) = s$ , and provide  $P(i)$  to each official indexed  $i$ .

Now, since you need  $k$  points to figure out the explicit expression of a  $k - 1$  degree polynomial  $P(x)$ , you must need at least  $k$  officials who hold one point of information to contribute.

## 8.6 Personal Regard on This Topic

The Lagrange interpolation method is, once again, not the only method of interpolation. However, it is important to note the mechanical aspects of it, such as converting fractions into integers using multiplicative inverse under some modular finite field (GF).

Sometime, you will also encounter questions that do not have a GF. In that case, the question might be a bit more quirky, a bit difficult to solve for the constraints it does not have. Still, polynomial's close experience during high school should aid you on the way.

Please always also remember that, a degree  $d$  polynomial is interpolable from  $d + 1$  points, not otherwise. For example, two points form a line, and three points form a quadratic curve, while one point forms one horizontal line.

When asked to develop a secret-sharing mechanism via polynomial, you'd like to pay attention to how many pairs/roots to give to each individuals for interpolation.

This can also associate to the notion of how much power someone has in a voting system to decrypt a secret to a specific party.

For instance, if person  $A$  has more pairs of the polynomial than person  $B$  has, then person  $A$  has a larger power in this voting system because he has more portions of the decryption key.



# Chapter 9

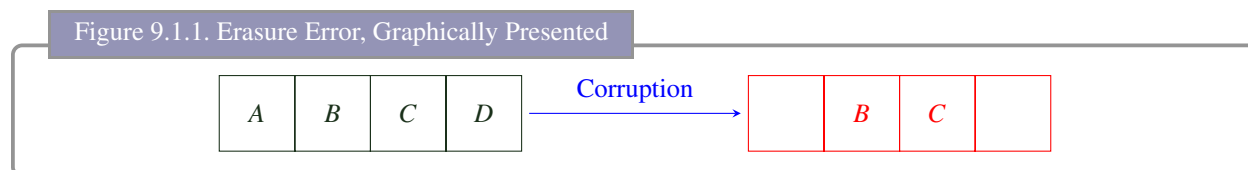
## Error Correcting Codes

We can deliver messages all the way we want, but what if the message is lost, dropped, or corrupted?

In this chapter, we learn about error correcting codes: a major study of mathematics and EECS with many underlying theories and applications on media playback devices. If you have guessed from the combination of disciplines involved, these applications contain significant linear algebra based on finite fields!

### 9.1 Erasure Errors

Let us consider a corrupted message as follows:



These errors are erasure errors, referring to the above situation(s) where a message with  $n$  packets has lost at most  $k$  packets during transmission.

To restore such messages, we can use some mathematical procedure.

Let us assume the contents of each packet is a number modulo  $q$ , where  $q$  is prime. The properties over  $GF(q)$  will then resolve the erasure error via polynomials. We are now applying note 8 contents!

Let us denote the message to be sent as  $m_1, \dots, m_n$ , where  $m_i$  is a number in  $GF(q)$ . Then:

1. There exists a unique polynomial  $P(x)$  of degree  $n - 1$  such that  $P(i) = m_i$ . Relevant properties were the focus subject of discussion in previous chapter.
2. For each message to be sent,  $m_i = P(i)$ . We may then generate additional packets via evaluating this polynomial  $P$  at points  $n + k$ . The interpolation would work fine as long as  $n + k \leq q$ , and thus we have attained a trasmitted codeword, which would have more packets than messages to account for packet loss.
3. This polynomial  $P(x)$  can be reconstructed from its values at any  $n$  distinct points, since it is degree  $n - 1$ . Therefore, as long as  $n$  packets are transmitted from the codeword, we can recover the original message.

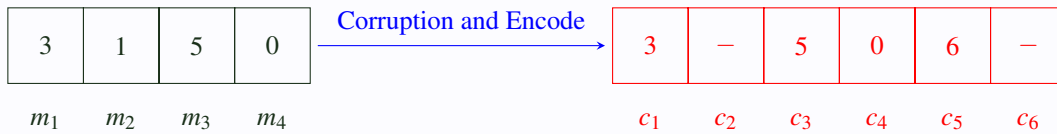
So in summary, we'd transmit a  $n + k$  packet codeword, which would have at most  $k$  packages lost. However, the codeword is based on the original message, which has  $n$  packets and thus is degree  $n + 1$ , and since we would have at least  $n$  packets left after any errorneous transmission or disk-biting, disk-crushing, using my DVD as a frisbee... you get the point: at least  $n$  packets are left for the polynomial interpolation to occur.

And the polynomial is characteristic of the message. As long as this polynomial can be retained from the successfully transmitted packets, we know the original message's individual packets via recovering the polynomial.

The uniqueness of our polynomial is thus CRUCIAL to the error correcting code!!!!!!

### Exercise Reader! Exercise! 9.1.1: Example of Recovering Erasure Errors

Let us review the example from lecture notes. Suppose our final result of transmission was something like this:



The original message, notated as terms of  $m$ , will form a polynomial interpolation  $P(x) = x^3 + 4x^2 + 5$ . Therefore, the codeword would actually be a rectangular array of:

3	1	5	0	6	1
$c_1$	$c_2$	$c_3$	$c_4$	$c_5$	$c_6$

But, we are able to interpolate that exact polynomial  $P(x)$  even if we are left with just  $c_1, c_3, c_4, c_5$ , since we have these  $n$  points to interpolate a  $n - 1$  degree polynomial. Therefore, computing  $P(2)$  will recover  $m_2$ .

## 9.2 Polynomial Interpolation in Error Correction

There are ways to perform polynomial interpolation other than the Lagrange Interpolation, which will be useful in the next section.

The goal of polynomial interpolation, again, is to take input of  $d + 1$  pairs (points) and return a polynomial  $p(x)$  such that it contains every point inputted.

### Theorem 9.2.1. Polynomial Interpolation for General Error

Let us discuss Polynomial Interpolation in this alternative method.

We may first write a system of  $d + 1$  linear equations in  $d + 1$  variables, where the variables are the coefficients of the polynomial.

Then, each equation is obtained by fixing  $x$  to be one of  $d + 1$  values (first item of per inputted tuple), with the rightside of equation (which is constant for a system solving context) be the second item of corresponding tuple.

This polynomial interpolation method is fundamentally linear-algebra based, and only requires very fundamental skills (Gaussian Elimination) that we might just make it a new EECS 16A question.

## 9.3 General Errors

The context of general error is, in summary, the general case of erasure errors. In this case, erased packets are replaced by corrupted packets, and our problem setting is once again, to work with a length- $n$  message with at most  $k$  packets corrupted.

By the similarity of context, you may probably have guessed that our settings for solution once again follows these rules:

- Each character is a number modulo  $q$  for some prime number  $q$ .
- We can describe a message as a polynomial  $P(x)$  of degree  $n - 1$  over  $GF(q)$ , such that each pair of message  $(i, m_i)$  is contained by  $P(x)$ .
- To cope, with general errors specifically, we will transmit additional characters (so, sending a codeword, not a message).

In this case, we transmit  $2k$  additional characters. Therefore, at least  $n + k$  of the received packets in a codeword is non-corrupted.

Now that we have the settings of our riddle, we must discuss: does Bob have sufficient information for recovering, or interpolating, the characteristic polynomial of original message?

First of all, for any  $n + k$  packets of the codeword, there should exist a unique polynomial that fits all  $n + k$  packets. So it is at least possible to recover the polynomial.

The proof of its uniqueness follows: Among these  $n + k$  points there can be at most  $k$  errors. Let  $P'(x)$  be another polynomial of degree  $n - 1$  that goes through these  $n + k$  packets, then at least  $n$  of the packets would be contained by both  $P$  and  $P'$ . However, our previous chapter has already revealed that a polynomial of degree  $n - 1$  is defined by its values at  $n$  points, so  $P$  and  $P'$  are in fact same polynomials.

Now that we know it is possible to find a polynomial, what is an efficient method of finding such polynomial, amidst the possible unlucky case of encountering some of the  $k$  errors?

#### Theorem 9.3.1. Error-Locator Polynomial

An error locator polynomial is designed to locate errors:

$$E(x) = (x - e_1)(x - e_2) \dots (x - e_k)$$

We do not know this polynomial explicitly, since we don't know where errors are.

However, we have found that:

$$P(i)E(i) = r_i E(i)$$

where  $r_i$  is the value of packet contained at index  $i$  of the corrupted message.

Since the value of  $i$  has the constraint  $1 \leq i \leq n + 2k$ , we can indeed produce a system of  $n + 2k$  linear equations with  $n + 2k$  unknowns.

Let us reflect on the approach.

Define  $Q(x) := P(x)E(x)$ , which has degree  $n + k - 1$  and therefore is described by  $n + k$  coefficients.

Meanwhile,  $E(x)$  is described by  $k + 1$  coefficients as a  $k$ -degree polynomial.

Therefore, we interpolate  $n + k$  coefficients from  $Q(x) = r_x E(x)$  and  $k$  coefficients from  $E(x)$  (which is a polynomial of  $(x - e_i)$ , and we find that  $e_i$  and there are  $k$  of these).

This forms  $n + 2k$  linear equations for us to configure the coefficients of  $Q(x)$  and  $E(x)$ .

At last, we perform long division to acquire  $P(x)$ :

$$P(x) = \frac{Q(x)}{E(x)}$$

#### Theorem 9.3.2. Correctness of Error-Locator Approach

Let us for the sake of proof assume the existence of alternative solutions  $Q', E'$ , such that:

$$P(x) = \frac{Q(x)}{E(x)} = \frac{Q'(x)}{E'(x)}$$

Well, since we can follow from the definition of  $Q$  to determine that  $Q(x) = r_x E(x)$ , we may conclude ahead:

$$\begin{cases} Q'(x) = r_x E'(x) \\ Q(x) = r_x E(x) \end{cases}$$

Following the above equations:

$$\begin{aligned} Q'(x)E(x) &= r_x E(x)E'(x) \\ Q(x)E'(x) &= r_x E(x)E'(x) \\ Q'(x)E(x) &= Q(x)E'(x) \\ \frac{Q(x)}{E(x)} &= \frac{Q'(x)}{E'(x)} \end{aligned}$$

Therefore, any possible set of polynomials would function to find us the right polynomial!

## 9.4 Personal Regard on This Topic

It really probably is just about reading the notes a few more times and knowing how to perform polynomial interpolation :p

For now, there's nothing I'd like to add on top of the notes for this section.

# Chapter 10

## Counting

The “counting” we address here in a discrete mathematical sense is not about counting from number to number in an arithmetic sequence. It’s about estimating the number of possibilities.

Consider a deck of 40 cards. What is the probability from which you draw all 5 parts of the inexorable Exodia? What is the total number of cases about the combination of 5 cards you draw from your deck, and from there within, what are all possible cases of drawing five cards of some same category (number 8, club, diamond...)?

Most importantly, what is the probability of it?

The techniques we employ to estimate and learn the total number of possible cases according to some conditions is known as **counting**.

In this chapter, let us begin from counting in simpler circumstances, like counting sequences.

### 10.1 Sampling without Replacement

#### 10.1.1 Counting Sequences

Let there be a set  $S = \{1, 2, \dots, n\}$ .

We may pick  $k \leq n$  elements from this set  $S$ , one at a time, while removing what we have sampled (selected) from  $S$ . This is known as **sampling without replacement**: once we select something from a set, it is removed from the set.

When counting the number of different ways to do this, either order matters, or order doesn’t. In the case that order does matter, we are counting ordered sequences.

The most daily-life example would be Poker.

When we deal a card, that card disappears is removed from the deck. The card is being sampled without replacement of it back into the card. Meanwhile, when the sequence is ordered, the sequence of dealing card A then card B is different from the sequence of dealing card B then card A: in a sequence, order matters.

Therefore, at the first card, we have 52 distinct possibilities for the card we now deal.

The second card, however, depends on the first card, since it cannot coincide with the first card. Your hand cannot have two exactly same card. That is cheating, and COMPSCI 70 does not condone cheating.

This means the second card has 51 possible choices, whichever card we picked from the first. Notice again that albeit the number of possibilities is the same, the second card you can get comes from a different set of cards depending on your first card.

Therefore, the total amount of possibilities for drawing 5 cards from a 52-card poker deck is in fact:

$$52 \times (52 - 1) \times (52 - 2) \times (52 - 3) \times (52 - 4)$$

This brings us to discuss the **First Rule of Counting**:

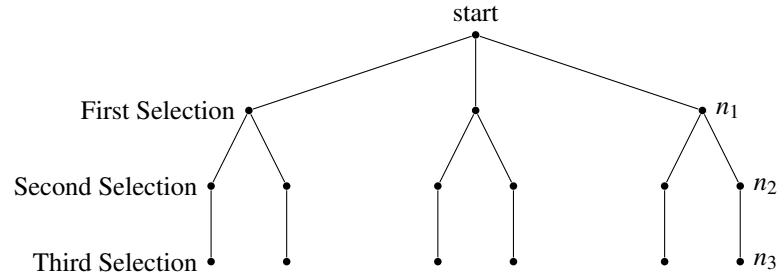
**Axiom 10.1.1. The First Rule of Counting**

If an object can be made by a succession of  $k$  choices, where there are  $n_i$  ways to make the  $i^{th}$  choice, then the total number of distinct objects that can be made in this way is:

$$\prod_{i=1}^k n_i$$

Let us use this opportunity to introduce another instinct:

Figure 10.1.1. The Tree of Counting



Say we are looking for the number of ordered sequences produced from sampling from a three-card deck. Then, we will once again have 3 possibilities for the first card, 2 for the second card, and 1 for the third card (provided it's the last card chosen).

In the above figure, each leaf node is a termination of the sampling process: a possibility. Therefore, we just count the number of leaf nodes that exist in this tree: in this case:

$$n_1 \times n_2 \times n_3 = 3 \times 2 \times 1 = 6$$

## 10.1.2 Counting Sets

Sets are slightly different stories from sequences, since there is no order in set.

Put in a deck drawing context as discussed in prior subsection, we would then be counting the number of distinct hands: where hands do not all have same members (cards).

To count the distinct subsets of a set  $S$  (in this case our deck) with 5 cards, we can let each sequence get placed into some bin corresponding to the set of 5 elements in the sequence. Each bin must contain  $5!$  ways of ordering the 5 elements in the set.

This means we just have to find the number of bins (distinct sets), which is essentially dividing the number of ordered sequences with  $5!$  for this context:

$$\frac{52!}{(52-5)!5!}$$

This quantity,  $\frac{n!}{(n-k)!k!}$ , stands for the number of  $k$ -element subsets from an  $n$ -element set  $S$ , the result of our counting. The alternative but more popular expression of this is  $\binom{n}{k}$ :

**Theorem 10.1.1. Second Rule of Counting**

Assume an object is made by a succession of choices, and the order in which the choices are made does not matter.

Let  $A$  be the set of ordered objects,  $B$  the set of unordered objects. If there exists an  $m$ -to-1 function  $f : A \rightarrow B$ , then we can count the number of ordered objects and divide by  $m$  to obtain the number of unordered objects.

So in the hand-counting context of this subsection, we had a  $5!$ -to-1 function from set of all ordered sequences to set of all distinct hands (which are the unordered sequences we count for).

## 10.2 Sampling with Replacement

What are the mathematical properties of sampling with replacement?

Since no elements are removed from  $S$ , the number of possible choices across each selection will be constant, provided that the number of elements in  $S$  is always constant and does not decrease.

And, since we have  $n$  choices in each trial, across  $k$  selections in the process of creating a  $k$ -length sequence,

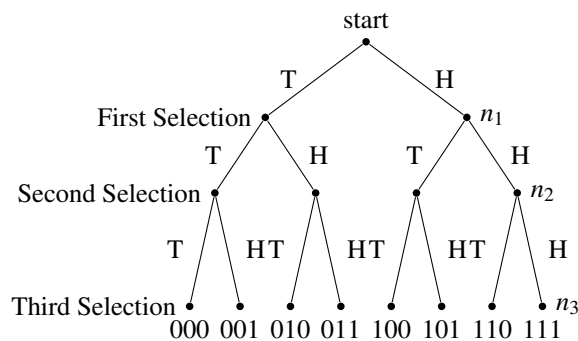
$$n_1 = n_2 = \cdots = n_k = n$$

Therefore, the total number of possible ordered  $k$ -length sequences from an  $n$ -element set  $S$  under sample with replacement is  $n^k$ .

### 10.2.1 Here is an Example

We may visit coin tosses first. The total possible outcomes of tossing a coin  $k = 3$  times can be portrayed by the following tree:

Figure 10.2.1. The Tree of Counting on Coins



In the above graph, each vertex corresponds to a possible result, where each bit stands for whether the toss resulted in a head (1) or a tail (0).

### 10.2.2 But Sometimes, Order Doesn't Matter

Suppose we would like to select 5 cards from a set of three cards  $S = \{1, 2, 3\}$ , where ordering doesn't matter.

In this context, Second Rule of Counting does not help, because we do not have a  $m$ -to-1 relationship between ordered sequences and distinct sets. We would need to apply a different perspective. Let us discuss below:

Let us first generalize back to the original setting. We have a set  $S = \{1, 2, \dots, n\}$  and would like to know the number of ways to choose multisets (sets with repetition) of size  $k$ .

Assume we have one bin for each element of  $S$ : providing  $n$  bins in total. We then count the number of ways to fill these bins with  $k$  elements, without minding the order of the bins themselves but just how many of the  $k$  elements they each contain.

Each allocation to the bin can be written as a bit-string, where 1 represents the change of bin and 0 represents the existence of one element.

For example, provided 5 bins and 5 elements where the leftmost bin has 1 element, the middle bin has 3, and the last has 1, the representation is

011000110

If there exists a space between the 1s, it would imply that the bin contains no elements.

The length of our binary string, provided  $n$  elements and  $k$  bins, is  $n + k - 1$ . We now choose which  $n$  locations should contain 0s in this arbitrary zeros, or alternatively, which of the remaining  $k - 1$  containing 1s.

A position corresponds to one specific value: once it is used, you may not reuse the position for another value and the position will be removed from the set of all possible options. The ordering also does not matter. Therefore, we are essentially sampling  $k$  locations from  $k + n - 1$  without replacement where order does not matter, which provides the total number of multisets as:

#### Theorem 10.2.1. Balls and Urns

Generally, the number of ways to place  $n$  indistinguishable balls into  $k$  distinguishable bins is:

$$\binom{n+k-1}{n} = \binom{n+k-1}{k-1}$$

Let us also describe the intuition we have used in generating the balls and urns approach in producing a match between bitstrings and ball-urn allocations:

#### Theorem 10.2.2. Zeroth Rule of Counting

If a set  $A$  can be placed into a one-to-one correspondence with a set  $B$  (in other words, find a bijection  $f : A \rightarrow B$ ), then  $|A| = |B|$ .

Upon which, let's attempt an elaboration on the knowledge of balls and urns:

#### Exercise Reader! Exercise! 10.2.1: Balls and Urns with No Empty Urns

Q: How many ways are there to put  $n$  balls into  $k$  urns such that there are no empty urns?

This is a variant of the balls and urns derivation, where instead of using 1 to notate the boundary of bins, we use 01 instead to ensure all bins on the left of some boundary receive at least one ball. Meanwhile, a character 0 is put as the rightmost bit of the ball-urn representing bitstring.

Now, how long is this bitstring, if we consider bits to be represented by 0 and 01 (instead of 1)?

We will realize that while  $n - (k - 1)$  bits would be 0,  $k - 1$  bits would be 01. Furthermore, we are only permitted  $n - 1$  positions to freely choose on the insertion of 0 (since  $k - 1$  of 0s are used in bin boundaries, and another 1 is put as the rightmost bit of the bitstring).

Therefore, the number of ways balls and urns operate with no empty bins would be:

$$\binom{n-1}{k-1}$$

## 10.3 Combinatorial Proofs

In short, proofs on combinatorials. They are mostly identity proofs. Let us explore a few examples/categories from below sections:

### 10.3.1 Binomial Theorem

#### Theorem 10.3.1. Binomial Theorem

$$\forall n \in \mathbb{N}, (a+b)^n = \sum_{k=0}^n \binom{n}{k} a^k b^{n-k}$$



**Proof:**

The LHS is a product of  $n$  terms  $(a + b)$ .

Let these factors be labeled sets  $f_1, \dots, f_n$  for generality.

For multiplying factors  $f_1, f_2$ , a total of 4 two-term combinations is made via matching each term of  $f_1$  to each term of  $f_2$ :

$$f_1 \times f_2 = \{(a_{f_1}, a_{f_2}), (a_{f_1}, b_{f_2}), (b_{f_1}, a_{f_2}), (b_{f_1}, b_{f_2})\}$$

$$(a + b)^2 = a^2 + ab + ba + b^2$$

Every term in such polynomial is in some form  $a^k b^{n-k}$ . It would be reasonable to see how this term contributes to the over sums of polynomial.

We produce terms  $a^k b^{n-k}$  when  $k$  copies of  $a$  and  $n$  copies of  $b$  are multiplied, and there are  $\binom{n}{k}$  ways to form such terms.

Therefore, for each index  $k$  possible (which spans from 0 to  $n$ , standing for the number of  $a$  apparent in the monomial), the sum of such monomials is the product of its value and number of such pairing:

$$\binom{n}{k} a^k b^{n-k}$$

Sum this product over all possible  $k$ , we will get the value of  $(a + b)^n$ , which is the sum of all possible monomials (products) achieved:

$$\sum_{k=0}^n \binom{n}{k} a^k b^{n-k}$$

The Binomial Theorem would then contribute to an interesting identity. Let us set  $a = 1, b = -1$  in binomial theorem. Then:

$$(-1 + 1)^n = \sum_{k=0}^n \binom{n}{k} (-1)^k (1)^{n-k}$$

$$= \sum_{k=0}^n \binom{n}{k} (-1)^k = 0$$

Huh, what! Did you just use binomial theorem to prove that  $-1 + 1 = 0$  to the  $n^{\text{th}}$  power is still 0?

Yeah, we did. But the last line of the proof, which is the combinatorial identity we have produced, is arguably more valuable.

### 10.3.2 Combinatorial Identities

Here are two interesting combinatorial identities to capitalize off in several proofs.

For now, let's talk about the first of two:

#### Theorem 10.3.2. Hockey-Stick Identity

The Hockey-Stick Identity states that:

$$\binom{n}{k+1} = \sum_{i=k}^{n-1} \binom{i}{k}$$

**Combinatorial Proof:** The LHS of this identity is the number of ways to choose a  $(k + 1)$ -length subset from an  $n$ -element set  $S$ .

For generality, let us denote  $S$  such that:

$$S = \{S_1, \dots, S_n\}, \quad S_1 < \dots < S_n$$

Make the following observation on how we can produce subsets:

- If we let the lowest-value element of the subset to be  $S_1$ , then we will have  $k$  elements to choose from the rest of  $n - 1$  elements from  $S$ . Therefore, it contributes  $\binom{n-1}{k}$  to the total number of combinations.
- If we let the lowest-value element of the subset to be  $S_i$ , then we will have  $k$  elements to choose from the rest of  $n - i$  elements from  $S$  (because  $i$  elements of  $S$  are smaller than or equal to the smallest element of subset). Therefore, it contributes  $\binom{n-i}{k}$  to the total number of combinations.

This sum will continue until we let the lowest-value element of subset to be  $S_{n-k}$ , where we would be left with  $n - (n - k) = k$  elements to choose from: the bare minimum.

Totaling all above contribution leads to the Hockey-Stick Identity.

The second identity I'd like to introduce to you is:

**Theorem 10.3.3. Sum of combinations from  $n$  items**

$$\sum_{k=0}^n \binom{n}{k} = 2^n$$

This identity can be proved via setting  $a = 1, b = 1$  from the binomial theorem:

$$\begin{aligned} 2^n &= (1 + 1)^n \\ &= \sum_{k=0}^n \binom{n}{k} 1^k 1^{n-k} \\ &= \sum_{k=0}^n \binom{n}{k} 1^k = \sum_{k=0}^n \binom{n}{k} \end{aligned}$$

To discuss this via combinatorial proof (which is via interpretation of both sides of equation), suppose we have a set  $S$  with  $n$  distinct elements.

The LHS of this identity counts the number of ways of choosing an  $i$ -lengthed subset.

The RHS, meanwhile, counts the total number of possible subsets (see the Chapter for Set Theory in this document).

Therefore, the LHS and RHS of this identity indeed denote the same value!

### 10.3.3 Permutation and Derangements

Assume that we collect the namecard of  $n$  people, and redistribute them into the  $n$  people. On average, how many people receive their own namecard?

Let us look at this in a combinatorics perspective.

Label the namecards as  $1, 2, \dots, n$  and let  $\pi_i$  denote the namecard that is returned to the  $i^{\text{th}}$  student.

Then, the vector  $\vec{\pi} = [\pi_1 \ \dots \ \pi_n]^T$ : a permutation of the set  $\{1, \dots, n\}$ . There are then  $n!$  distinct permutations of  $\{1, \dots, n\}$  since it is not a multiset.

In this setting, the situation at which a person receives his own namecard is when  $\pi_i = i$ .

Therefore, the number of people who receive their own namecard is the amount of index at which  $\pi_i = i$ . These points are called “fixed points”.

A permutation with no fixed points is called a **derangement**.

There is a formula for the number of derangements provided a set of  $n$  distinct elements:

**Theorem 10.3.4. Number of Derangements**

**Theorem:** For an arbitrary positive integer  $n \geq 3$ , the number  $D_n$  of derangements of  $\{1, \dots, n\}$  satisfies:

$$D_n = (n-1)(D_{n-1} + D_{n-2})$$

**Proof:**

In a derangement  $(\pi_1, \dots, \pi_n)$ , suppose  $\pi_n = j \in \{1, \dots, n-1\}$ . Then, there are  $n-1$  choices for  $j$ . We also know that  $j \neq n$  because this would then not produce a derangement.

Now, if  $\pi_j = n$ , then we recognize that  $j$  and  $n$  get swapped, and the number of possible derangements with the remaining  $n-2$  numbers is  $D_{n-2}$ .

Or, if  $\pi_j \neq n$ , then  $n$  satisfies the same constraint as  $j$  (in not being able to be the value  $j$ ). Therefore, we have  $n-1$  values left to form derangements with, which provides  $D_{n-1}$  possibilities.

This completes the derivation of the above formula.

Notably, with the boundary conditions  $D_1 = 0$  and  $D_2 = 1$ , we can then determine  $D_n$  computationally via recursion and rather efficiently with the use of memoization.

The derivation is not shown here, but the recursion would then yield that:

$$D_n = n! \sum_{k=0}^n \frac{(-1)^k}{k!}$$

## 10.4 The Principle of Inclusion-Exclusion

Overcounting is a big problem in combinatorics, but here is a technique to prevent most of such situations.

Let  $A_1$  and  $A_2$  be two subsets of the same finite set  $A$  and we want to count the number of elements in  $A_1 \cup A_2$ .

If  $A_1$  and  $A_2$  are disjoint, it is very lucky:

$$|A_1 \cup A_2| = |A_1| + |A_2|$$

If  $A_1$  and  $A_2$  are not disjoint, however, then we need to subtract the number of elements from their intersections:

$$|A_1 \cup A_2| = |A_1| + |A_2| - |A_1 \cap A_2|$$

This brings us to the generalization of this logic over an arbitrary number of subsets:

**Theorem 10.4.1. Principle of Inclusion-Exclusion**

**Theorem:** Let  $A_1, \dots, A_n$  be arbitrary subsets of the same finite set  $A$ :

$$|A_1 \cup \dots \cup A_n| = \sum_{k=1}^n (-1)^{k-1} \sum_{S \subseteq \{1, \dots, n\}: |S|=k} |\cap_{i \in S} A_i|$$

**Proof:**

Note that if  $a$ , an element, is not included by the union of all subsets of  $A$  considered in this scenario, then it is not possible for  $a$  to be counted anyways.

For  $a$  to be counted, it must belong to some subset  $A_i$ , such that  $i \in \{1, \dots, n\}$ .

For the sake of clarifying the proof, define an index set  $M \subseteq \{1, \dots, n\}$  as:

$$M = \{i \in \{1, \dots, n\} : a \in A_i\}$$

applying the prior logic,  $a \in A_i$  iff  $i \in M$ . Here, if the subset  $S$  we find is not contained by  $M$ , then  $a \notin \cap_{i \in S} A_i$ .

Therefore, for  $a$  to be counted,  $a \in \cap_{i \in S} A_i$  for any non-empty subset  $S \subseteq M$ .

The number of subsets that have counted  $a$  on RHS is thus:

$$\sum_{k=1}^n (-1)^{k-1} \sum_{S \subseteq \{1, \dots, n\}: |S|=k} 1 = \sum_{k=0}^m \binom{m}{k} (-1)^{k-1} = 1$$

On the perspective of some venn diagram of events, inclusion-exclusion principle serves to calculate the net area of venn diagram by repeating the process of subtracting overlapped regions and adding back over-subtracted regions. It is notable that, however, inclusion-exclusion used in probabilistic contexts provide diminishing marginal utility due to computational cost, and usually perform a good enough job at the third iteration of adding/subtracting areas.

And applying this to the derangement problem:

#### Example Question 10.4.1: Count Derangements

Let  $A_i$  denote the set of all permutations with  $i$  being a fixed point.

The number of permutations with at least one fixed point is  $|A_1 \cup \dots \cup A_n|$ . Therefore,

$$D_n = n! - |A_1 \cup \dots \cup A_n|$$

Meanwhile, the cardinality of  $A_i$  would be  $(n-1)!$  because it just has to permute  $n-1$  other elements in an arbitrary way.

The intersection of  $A_i$  and  $A_j$  would thus need to permute  $n-2$  arbitrary elements, and has cardinality  $(n-2)!$ .

Let us follow the clue: for a subset  $S \in \{1, \dots, n\}$  such that  $|S| = k$ ,

$$|\cap_{i \in S} A_i| = (n-k)!$$

The rest of algebra follows:

$$\begin{aligned} |A_1 \cup \dots \cup A_n| &= \sum_{k=1}^n (-1)^{k-1} \binom{n}{k} (n-k)! \\ &= \sum_{k=1}^n (-1)^{k-1} \frac{n!}{k!} \\ &= n! \sum_{k=1}^n (-1)^{k-1} \frac{1}{k!} \\ D_n &= n! - |A_1 \cup \dots \cup A_n| \\ &= n! \left( 1 - \sum_{k=1}^n (-1)^{k-1} \frac{1}{k!} \right) \\ &= n! \sum_{k=1}^n \frac{(-1)^k}{k!} \end{aligned}$$

## 10.5 Personal Regards on this Topic

The biggest takeaways from this chapter would be the ability to answer the following questions:

- In the counting problem I work at, does order matter?
- In the counting problem I work at, is sampling with or without replacement?
- How do I perform combinatorial proofs by creating stories (interpretation of combinatorial expression) for both sides of the equations?
- How do I setup counting problems and see what object do I exactly need to count?

And combinatorial identities can prove useful if you bring them along on cheat sheets; as you see, combinatorial arguments of this chapter have hugely relied on the power of many combinatorial identities.

# Chapter 11

## Countability

### 11.1 More on Bijections

Recall upon the Zeroth Rule of Counting:

Two finite sets have the same size iff their elements can be paired up, where each element of one set has a unique partner in another.

This is oddly familiar. It is, in fact, formalized via the concept of bijection, which we will re-introduce here after its brief appearance when discussing modular arithmetics.

#### Definition 11.1.1. Bijections

Consider a function  $f : A \rightarrow B$ , which since it is a function, would satisfy that  $\forall x \in A, f(x) \in B$ .

Then  $f$  is a bijection if such relation holds true in a reverse direction. Therefore, for a bijection  $f : A \rightarrow B$ :

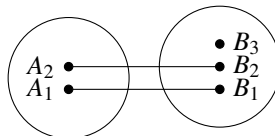
$$[\forall x \in A (f(x) \in B)] \wedge [\forall x \in B (f(x) \in A)]$$

In this case, as discussed before,  $f$  would satisfy two other significant properties simultaneously: injective and surjective.

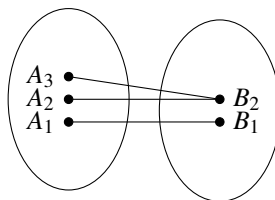
#### Definition 11.1.2. Injective, Surjective, Bijection

We call the  $A$  the set of *pre-images*, and  $B$  the set of *images*.

- **Injective:** Each image (output) has at most one pre-image (input). This type of function is thus called a "one-to-one" function.



- **Surjective:** Each image (output) has at least one pre-image (input). This type of function is "onto".



And once again, the inverse of a bijective function must also be bijective. While this is listed as an exercise to the reader, you may imagine bijections as relationship between sets  $A$  and  $B$ ; then, notice that inverse functions of  $f : A \rightarrow B$  are in fact relationships  $f^{-1} : B \rightarrow A$ .

Since  $f^{-1}$  is the inverse of  $f$ , it carries the exact same pairings of  $f$  and must also be bijective as  $f$  is.

### 11.1.1 Cardinality and Bijection

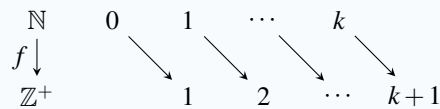
Cardinality is the size of a set.

To show that two sets have the same cardinality, we can attempt to prove the existence of a bijection between them. For example, let us consider the following problem of two sets:

#### Example Question 11.1.1: The cardinality of nonnegative integers and positive integers

For this course,  $0 \in \mathbb{N}$ . Essentially, we ask if the cardinality of  $\mathbb{N}$  is equal to that of  $\mathbb{Z}^+$ .

To prove the equality, let us attempt to form a bijection between these sets:



We can see that there is indeed a bijection  $f : \mathbb{N} \rightarrow \mathbb{Z}^+$ , since every value of each set is paired with another value from the other set. We speak of this as: having **enumerated** each element from set  $\mathbb{Z}^+$ .

Meanwhile, we can notice that the relationship is one-to-one because no two natural numbers have the same image.

A very similar case can be made for  $f : \mathbb{N} \rightarrow 2\mathbb{N}$ , and can be found on the lecture note.

We are left with another interesting case where bijection can be certified, between  $\mathbb{N}$  and  $\mathbb{Z}$ .

Our initial thought here would have to do with the properties of both sets. For each pair of number  $(-k, k)$  in  $\mathbb{Z}$ , where  $k$  is a natural number, we also need to take two respective numbers from  $\mathbb{N}$  to match against them.

If so, then let us attempt with the function:

$$f(x) = \begin{cases} \frac{x}{2} & \text{if } x \text{ is even} \\ -\frac{x+1}{2} & \text{if } x \text{ is odd} \end{cases}$$

And let us see how this pairing works out:



Let us proceed with a proof of its bijection by demonstrating its injectiveness and surjectiveness:

#### Example Question 11.1.2: Prove the One-to-One and Onto of above function

Remember that to prove a relationship  $f$  as bijection, we will be proving its two substatements:

- **Injectiveness:**  $\forall x, y \in \mathbb{U}, f(x) = f(y) \implies x = y$
- **Surjectiveness:**  $\forall y, \exists x, f(x) = y$

#### Injectiveness (One-to-One):

Suppose there exists two arbitrary values  $x, y$  such that  $f(x) = f(y)$ .

Then, both values must have the same sign, so either  $f(x) = \frac{x}{2} = f(y) = \frac{y}{2}$  or  $f(x) = -\frac{x+1}{2} = f(y) = -\frac{y+1}{2}$ .

And in either case, we obtain that  $x = y$ , thus obtaining that:

$$f(x) = f(y) \implies x = y$$

Therefore,  $f$  is injective.

**Surjectiveness (Onto):**

Suppose that  $y$  is non-negative, then  $f(2y) = y$ . Suppose that  $y$  is negative, then  $f(-(2y + 1)) = y$  as well.

Since every image has a pre-image,  $f$  is onto.

Therefore, while  $\mathbb{N}$  has the same size with  $\mathbb{Z}^+$ , it also has the same size with  $\mathbb{Z}$ . Interesting!

## 11.2 Countability

Now we come to an important question. What does it mean for a set to be countable?

### Definition 11.2.1. Countability

A set  $S$  is countable if it is finite, or if there is a bijection between  $S$  and  $N \subseteq \mathbb{N}$ .

Essentially, if we can enumerate all members of a set  $S$ , then the set  $S$  is countable. Vice versa.

Well, just knowing countability by definition would yield fun, but what would yield more fun? Explaining if  $|\mathbb{Q}| = |\mathbb{N}|$ !

Exercise Reader! Exercise! 11.2.1: Proof that  $|\mathbb{N}| = |\mathbb{Q}|$

To do so, let us first observe that  $|\mathbb{N}| \leq |\mathbb{Q}|$ , since we would see that every natural number is a rational number while the reverse does not necessarily hold true.

To continue off this proof, we now need to recognize if also  $|\mathbb{Q}| \leq |\mathbb{N}|$ .

To do so, we will exhibit an injection  $f: \mathbb{Q} \rightarrow \mathbb{N}$ .

Now, let us consider a counterclockwise spiral that rotates its direction of travel 90 degrees counterclockwise. In such a figure, findable from the lecture notes (Around page 4), we observe that each rational number is represented by the point  $(x, y)$  on the path of such spiral, which is placed on the two-dimensional space  $\mathbb{Z}^2$ .

If we map each lattice point on the spiral to the numeric order of its precedence on the spiral, we then thus prove that every rational number (characterized as a lattice point, again) must correspond to one natural number.

Therefore, knowing that all points on the spiral must be presentable as a natural number, and knowing that each lattice point corresponds to exactly one natural number:  $|\mathbb{Q}| \leq |\mathbb{N}|$ .

Last but not least:

$$|\mathbb{N}| \leq |\mathbb{Q}| \wedge |\mathbb{Q}| \leq |\mathbb{N}| \implies |\mathbb{N}| = |\mathbb{Q}|$$

Exercise Reader! Exercise! 11.2.2: Proof that the set of all polynomials is countable

The set of all binary strings is countable: each binary string has a representative number, and if we must consider the empty string, then we would assign it to the value 0 and assign other bit-strings to the number they represent in binary base added by 1.

This exact reasoning would suffice for ternary strings, quad-ary strings, hexadecimal strings... the list goes on, but this holds a significant implication.

If we can manage to map mathematical objects, like polynomials and graphs (representable by adjacency matrices, for example), then sets of mathematical objects can indeed also be proven as countable.

The lecture note proposes an example proving the set of all polynomials to be countable, via formatting a



polynomial as a ternary string, where each coefficients and exponents are formatted in their binary bit-string form and separated by a character “2”. The former logic has stated that the set of all ternary strings is countable, which would then, via convenience offered from the definition of countability, make the set of all polynomials be as countable as the set of all ternary strings are.

Now, on a similar note, let's prove a much more significant proposition in scale:

Exercise Reader! Exercise! 11.2.3: Proof that the Cartesian product of countable sets is countable

Let countable sets  $S$  and  $T$  be examined for the proof, where we prove that  $S \times T$  is countable. We should first observe that  $S$  and  $T$  themselves must, for their countability, have injections:

$$\begin{aligned} f : S &\rightarrow \mathbb{N} \\ g : T &\rightarrow \mathbb{N} \end{aligned}$$

Therefore, there must also exist some injection:

$$h : S \times T \rightarrow \mathbb{N} \times \mathbb{N}$$

where  $S \times T$  is a set of tuples of enumerated items, so representable as a set of tuples of natural numbers. Then, let us investigate the countability of  $\mathbb{N}^2$ .

We recognize that each natural number larger than 1 has a unique prime factorization. Therefore, let us enumerate all ordered pairs  $(x, y) \in \mathbb{N}^2$  with the number  $2^x 3^y$  (2 and 3 here are replaceable by other prime numbers).

Then, we will have built an injection with the above rules:

$$\begin{aligned} F : \mathbb{N}^2 &\rightarrow \mathbb{N} \\ F(x, y) &:= 2^x 3^y \end{aligned}$$

where the relation  $F$  is indeed injective because each number can only have one unique prime factorization, such that:

$$\forall m, n, r, s \in \mathbb{N}, (2^m 3^n = 2^r 3^s \implies m = r \wedge n = s)$$

Now, such injection suggests that  $\mathbb{N}^2$  is countable; therefore,  $S \times T$ , which has an injection towards  $\mathbb{N}^2$ , must also be countable.

## 11.3 Cantor's Diagonalization

In this chapter, we deal with two great questions.

First of all, rational numbers are so dense? The set of rational numbers is dense, as between any two rational numbers, we can find another rational number. In fact, between any two real numbers there is always a rational number (which we used as an example in Chapter 2 when demonstrating proofs). It's Real Dense.

And second of all, is there a bijection between the rationals and the reals, provided the above fact of the density of set (that real numbers and rational numbers are all dense)? We will resolve this question via an argument from Cantor, known as “Cantor's diagonalization”.

But, before presenting the proof, let us recall that any real number in the interval  $[0, 1]$  can be written as an infinite decimal with no trailing zeros, and rational numbers can be represented as recurring decimals, while irrational ones non-recurring decimals.

Most importantly, all decimal expressions of irrational numbers are all unique and infinitely long. For numbers like 0.5, they will be represented as  $0.4\overline{9}$ .

### Theorem 11.3.1. Uncountability of Real Numbers

**Theorem:** The real interval  $\mathbb{R}[0, 1]$  is uncountable.

**Proof:** For the sake of contradiction, let us consider there is a bijection  $f : \mathbb{N} \rightarrow \mathbb{R}[0, 1]$ . We may then enumerate real numbers in an infinite list  $f(0), \dots$  as follows:

$$\begin{aligned} f(0) &= 0.5214\dots \\ f(1) &= 0.1416\dots \\ f(2) &= 0.9478\dots \\ f(3) &= 0.5309\dots \\ &\vdots \end{aligned}$$

The basic rule of this construction is that for the number of leftup-rightdown diagonal (in this case 0.5479 $\dots$ )  $r$ , its  $(i+1)^{th}$  digit after decimal point must be same with that of  $f(i)$  for any  $i \in \mathbb{N}$ .

However, we can see that a value, such as  $r$  whose digit is incremented by 2 under modulo 10, cannot fit into this enumeration because of its  $(i+1)^{th}$  digit not being equal to that corresponding digit of  $r$ .

Therefore, we see at least one item not counted in the “bijection” that should’ve existed between  $\mathbb{N}$  and  $\mathbb{R}[0, 1]$ , thus deeming  $\mathbb{R}[0, 1]$  uncountable.

The way this technique worked ultimately roots in how we defined a real number as an infinite decimal expansion that can be either recurring or non-recurring. Say today we attempt to apply this argument on  $\mathbb{Q}$ , which is a set of recurring decimal expansions, the proof wouldn’t gain progress because we cannot confirm whether the number of diagonal is indeed a recurring decimal representation.

In other words, Cantor’s Diagonalization relies on the infinite length of expressions from the set it attempts to prove uncountable.

## 11.4 Cantor’s Set

The Cantor’s Set then argues about the mathematical properties of interval  $\mathbb{R}[0, 1]$ . This set is defined by repeatedly removing the middle thirds of line segments infinitely many times.

### Definition 11.4.1. Cantor’s Set, Defined by Construction

A Cantor’s Set can be constructed iteratively with infinite steps.

In each step, we remove the middle one-third of line segments we currently have. For example,

$$[0, 1] \rightarrow \left[0, \frac{1}{3}\right] \cup \left[\frac{2}{3}, 1\right]$$

And then we proceed to remove the middle third of each line segments we deal with:

$$\left[0, \frac{1}{3}\right] \cup \left[\frac{2}{3}, 1\right] \rightarrow \left[0, \frac{1}{9}\right] \cup \left[\frac{2}{9}, \frac{1}{3}\right] \cup \left[\frac{2}{3}, \frac{7}{9}\right] \cup \left[\frac{8}{9}, 1\right]$$

You would think that the Cantor’s Set faces an inevitable emptiness (just like most of human life), but it is merely that the measure of the set is zero, as in it now contains many almost trivial intervals with many isolated points across the interval and “looks mathematically like being empty”. It is the realistic equivalent of looking at a white paper with many atomic black points, which we just cannot perceive with human eyes.

But the Cantor’s Set is uncountable. This is proven by forming a bijection between the Cantor’s Set and an uncountable set:

**Theorem 11.4.1. The Uncountability of Cantor's Set**

**Theorem:** The Cantor's Set is uncountable.

**Proof:** Each element in the Cantor's Set can in fact be expressed as a ternary string, at last containing only those whose ternary representations have 0s and 2s (because those that involve 1s must have been in the middle third of some interval and thus deserves removal).

We will also proceed to write any finite decimal expressions into its infinite form (for example,  $0.1 \rightarrow 0.0\bar{2}$ ).

This leaves us with a set builder representation of the Cantor's Set:

$$\left\{ C = \{x \in [0, 1] : x \text{ has a ternary representation of only 0s and 2s} \} \right\}$$

Which allows us to map the set  $C$  onto the interval  $\mathbb{R}[0, 1]$  via substituting every trit 2 with value 1 and forcing the ternary string into a binary presentation that corresponds to some value in  $\mathbb{R}[0, 1]$ .

This way, we formed a bijection between  $C$  and an uncountable set, so using the Zeroth Rule of Counting,  $C$  the Cantor's Set is also uncountable.

Note that the difference between the above proof and the proof of countability for sets of all binary string, is that in the Cantor's Set proof, leading 0-valued bits exist, as supposed to them not existing in the enumerable set of bitstrings.

## 11.5 Power Sets: Higher Order Infinity

Here is a review on Power Sets:

**Definition 11.5.1. Power Set**

The power set of  $S$  can be written as:

$$\mathbb{P}(S) = \{P | P \subseteq S\}$$

as one of its various denotations, the one on lecture note using a Weierstrass  $\mathbb{P}$ ,  $\wp$ .

The power set of a set is the set of all of its possible subsets.

In addition, if the cardinality of  $S$  in the above definition is  $|S| = k$ , then we can state that  $|\mathbb{P}(S)| = 2^k$ .

This can now be combinatorically understood: every power set is a  $k$ -bit string representation for the absence of each of  $k$  elements.

Now, let us discuss if the power set of  $\mathbb{N}$  can be countable, and from there, reveal a thing or two about infinity.

**Theorem 11.5.1. Uncountability of  $\wp(\mathbb{N})$** 

For the sake of contradiction, let us assume there is a bijection  $f : \mathbb{N} \rightarrow \wp(\mathbb{N})$  defined exactly like Cantor's Diagonalization Argument, where each row of a table describes a subset via a bit-string, and each column stands for whether an element of  $\mathbb{N}$  is present or not.

Let us now present a modified version of the diagonal number (again, leftup-rightdown) that is the flipped version of diagonal (which means 0s and 1s are flipped). That digit then must not exist in the enumeration, therefore proving the bijection malfunctioning, and moreover,

$$|\wp(\mathbb{N})| > |\mathbb{N}|$$

The cardinality of  $\mathbb{N}$  is denoted in the symbol  $\aleph_0$  (aleph-null), and using that symbol,  $|\wp(\mathbb{N})| = 2^{\aleph_0}$ .

Such cardinality of  $\wp(\mathbb{N})$  is the same as the cardinality of  $\mathbb{R}$ , which is known as  $\mathfrak{c}$ .

Even larger infinite cardinalities are denoted as  $\aleph_k$  where  $k > 0$ , and they can be defined using set theory. They, interestingly, obey the rules of arithmetics, and in turn provide assistance towards famous mathematical hypotheses.

# Chapter 12

## Computability

### 12.1 The Liar's Paradox

Have you ever played the hit game Among Us? If not, here is a fundamental rule of the game Among Us:

Let  $A$  be the set of players in an Among Us game, then there exists one player in  $A$  such that the player is an impostor.

But here is a problematic proposition:

An impostor says: "All impostors are liars."

Why do we have so many Among Us jokes? That I do not know, and memes are the DNA of souls so it's only natural to use memes in mathematical contexts.

But the key inquiry we'd like to make is: how exactly is this statement problematic?

Well, if the statement is true, then all impostors are liars. However, since this sentence is said by an impostor, it is false, and all impostors are in fact not liars. But if impostors are not liars, then this statement is true.

The logical loop goes on due to the self-referencing nature of this sentence.

A same inquiry is reflected on this famous proposition in discrete mathematics:

In the village with just one barber, some men shave themselves or go to that barber.

The barber proclaims: I shall shave all and only those men who do not shave themselves.

Does the barber shave himself?

If a barber's statement is true, then he shaves only those who do not shave themselves. He would shave himself if he does not shave himself. If he does not shave himself, however, he does shave himself.

Such is the wonder of **liar's paradox**: you cannot get any effective results from interpreting these statements due to their self-referencing and contradictory nature.

### 12.2 Self Replicating Programs

Let us view self-reference as a computational concept: can we use self-reference to design a program that outputs itself?

Let us consider how we can do this if we write the program in pseudocode, where thus we get:

```
print the following sentence twice, the second time in quotes and indented:
    "print the following sentence twice, the second time in quotes and indented:"
```

If we execute this instruction above, then we will get the output be exactly what we executed. This self-replicating program is known as a quine (which you might have encountered in COMPSCI 61A as an extra credit). To construct such a quine in general, we can use:

#### Theorem 12.2.1. Recursion Theorem

Given any program  $P(x, y)$ , we can always convert it into another program  $Q(x)$  such that  $Q(x) = P(x, Q)$ .  $Q$  behaves exactly as  $P$  would if its second input is the description of program  $Q$ . Therefore,  $Q$  has access to its own description, and is essentially a self-aware version of  $P$ .

The Recursion Theorem will eventually lend great strength to our efforts in proving and disproving existence of some program.

## 12.3 The Halting Problem

Now we see self-reference working in computation (albeit still not in human logic), let's then come to ask: is there anything a computer cannot do?

Yes.

In 1936, Alan Turing managed to show that there is no program that can answer for itself:

“does your execution lead to an infinite loop?”

The proof has two ingredients: self-reference, and how inseparable programs and data are to each other. But before inspecting the proof, let's inspect and understand the problem more.

Let us write a program that behaves as follows:

$$\text{TestHalt}(P, x) = \begin{cases} \text{"yes"}, & \text{if program } P \text{ halts on input } x \\ \text{"no"}, & \text{if program } P \text{ loops on input } x \end{cases}$$

And now, let us access the proof:

#### Theorem 12.3.1. The Uncomputability of Halting Problem

**Theorem:** The Halting Problem is uncomputable.

**Proof:** For the sake of contradiction, let us construct the following program:

```
Turing (P)
    if TestHalt(P, P) is ``yes``:
        loop forever
    else:
        halt
```

Let us then look at the behavior of  $\text{Turing}(\text{Turing})$ , which either halts or not.

- If it halts, then it must be that  $\text{TestHalt}(\text{Turing}, \text{Turing})$  returned “no”, which means Turing loops on Turing.
- If it does not halt, then it loops on forever, which implies that Turing should've halted on Turing.

In both cases, we find a contradiction. The program  $\text{TestHalt}$  must've been dysfunctional.

The above proof is actually a diagonalization argument:

### Theorem 12.3.2. Continued, The Uncomputability of Halting Problem

**Proof:**

Let us then construct a table of program enumeration vs program enumeration ( $P_i$  vs.  $P_i$ ).

On this table, the  $(i, j)^{th}$  entry in the table is  $H$  if  $P_i$  halts on  $P_j$  (halts when  $P_j$  is an input), and  $L$  if otherwise (which means it loops).

Supposedly, the program  $P_n = \text{Turing}$  exists. Yet, it cannot be enumerated, as outlined with the following proof by case:

- If the  $(n, n)$  entry is  $H$ , then  $P_n$  halts on  $P_n$  and Turing will thus loop forever.
- If the  $(n, n)$  entry is  $L$ , then Turing halts on  $P_n$  instead, but contradicts the original assumption.

which the logic follows the prior proof performed.

Then, since Turing is a simple modification of TestHalt, if Turing does not get enumerated, neither can TestHalt.

Therefore, the Halting problem is unsolvable.

## 12.3.1 Easy Halting Problem

Brandon. Who in the right mind would feed a program back to the original program?

To that answer, I respond “*computer scientists*”, which you are one, for the record of this fact.

But, computer scientists still proposed a new, more practical problem according to your concerns, complaints, insights:

### Theorem 12.3.3. The Uncomputability of Easy Halting Problem

**Theorem:** The Easy Halting Problem states that:

There does not exist a program that can decide whether any given program  $P$  will halt on 0.

**Proof:** Let there be a program TestEasyHalt that resolves this problem (for the sake of contradiction):

$$\text{TestEasyHalt}(P, x) = \begin{cases} \text{“yes”}, & \text{if program } P \text{ halts on input } 0 \\ \text{“no”}, & \text{if program } P \text{ loops on input } 0 \end{cases}$$

Then following the Recursion Theorem, we can use this as a subroutine to build the following algorithm:

```
Halt (P, x)
  def P' (x) :
    if x == 0:
      return P(x)
    return TestEasyHalt (P')
```

Well,  $P'(0)$  would only halt iff  $P(x)$  does, so if such a program TestEasyHalt exists, then Halt effectively solves the Halting Problem because it can tell whether an arbitrary  $P$  holds or not on an arbitrary  $x$ . However, we learned that the Halting Problem is impossible to solve!

Therefore, there cannot be such a program Halt, and thus cannot be a program TestEasyHalt.

## 12.4 Godel's Incompleteness Theorem

In 1900, David Hilbert posed two questions about the foundation of logic in mathematics:

- Is arithmetic consistent?
- Is arithmetic complete?

Such questions made interesting contributions to many things in the world, one being mathematics, and the other being an overall existential crisis across the human species and Pythagoreanists.

To understand the problems he posted, we should remember that mathematics is a formal system built on axioms and rules of inferences, hence all the theorems we have relied on. Axioms are the cornerstones of mathematics.

So, the first question above asks whether it is possible to prove a proposition  $P$  and its negation  $\neg P$  (which we attempted to not do in the Law of Excluded Middle). The case in which this is proven will cause arithmetic to be inconsistent; otherwise, consistent.

If arithmetics is inconsistent, many currently wrong statements can later be proven right.

The second question asks whether every true statement in arithmetic can be mathematically proved.

Let's talk about history. There was indeed an unproven theorem (Fermat's Last Theorem) until around 1990s.

Meanwhile, as these questions were proposed in 1928 as "Entscheidungsproblem", which we hope the answer is 'yes' to, they were immediately proven "no" by Godel in 1930.

This proof exploited the deep connection between proofs and arithmetics, as well as proofs and computation. We will, in the following sections, discuss the sketch and essence of Godel's Proof on the incompleteness of mathematics.

### 12.4.1 Sketch of Godel's Proof

Suppose we have a formal system  $F$  and assume it is sufficiently expressive that we can use it to express all of arithmetic. Suppose we can write the following statement:

$$S(F) = \text{"This statement is not provable in system } F\text{"}$$

Two things can happen:

- Case 1:  $S(F)$  is provable. Then the statement is true, but by inspecting the content of the statement itself, we attempt to see that this implies  $S(F)$  is not provable. This makes  $F$  inconsistent.
- Case 2:  $S(F)$  is not provable, which means  $S(F)$  is true, and  $F$  is therefore already incomplete.

To prove that mathematics is incomplete, we would just need to construct such a statement  $S(F)$ , which is the reduction of the larger problem for proving mathematics wrong.

### 12.4.2 Proof via Halting Problem

Suppose arithmetic is both consistent and complete, we will now use this statement to perform a proof by contradiction. Let the proposition  $S_{P,x}$  depict "Program  $P$  halts on input  $x$ ", and we can thus phrase it as an arithmetic:

$$\exists z(z \text{ encodes a valid halting exception sequence of } P \text{ on input } x)$$

Now let us assume that arithmetic is both consistent and complete. This means the above statement, which represents the halting problem, can be proven either true or false due to the completeness of arithmetics.

However, this then in turn implies the halting problem is prove-able via searching through every single possible proof enumerable (which is countably many). The problem, however, is that the Halting Problem is uncomputable, not allowing such a program for searching proofs to exist.

Therefore, the initial assumption is false. Arithmetics is inconsistent and incomplete.

# Chapter 13

## Introduction to Discrete Probability

Occasionally, the voice in your brain will protest:

**STOP DOING PROBABILITY!** Years of calculating yet no real-life usage on finding an affordable housing in Berkeley!

But, probability is in fact useful across many other disciplines and purposes. One of the most applicable example in real life is gambling, gacha, as well as several applications of computer science where we work with statistical and probabilistic frameworks to produce efficient models.

What do statements concerning probability look like? It may be of the following forms:

- The average time between system failure is 5 days.
- In average, 1 in 1000 students from UC Berkeley lives in an affordable housing.
- The probability at which the game character Mr. Game and Watch rolls a 9 on a fair 8-face dice with values 2 to 9 is precisely  $\frac{1}{8}$ .

And so on.

Implicitly, all statements concern with a **probability space**: the result of random experiments that we constructed or simulated. Therefore, we would greatly benefit from understanding what random experiments are as we come to discuss *probability space* and *sample space*, two “frames” upon which the theories of probability depend on.

### 13.1 Random Experiments

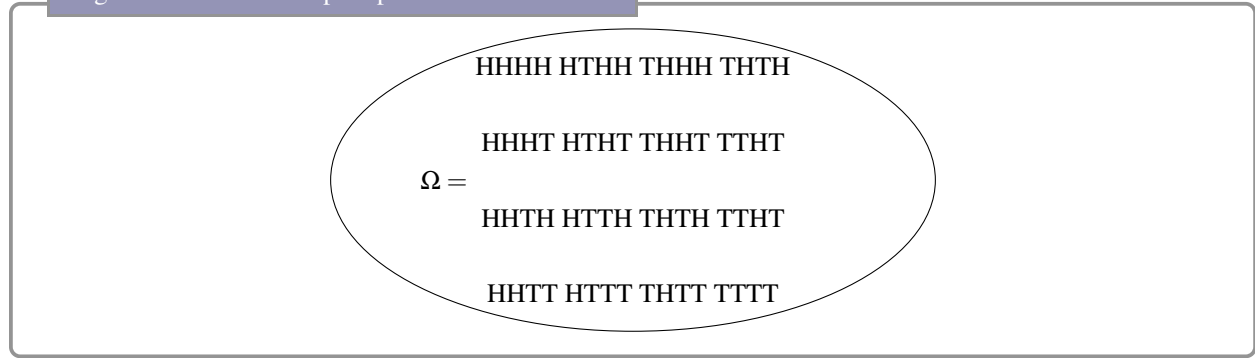
In general, a **random experiment** consists of sampling  $k$  elements from some set  $S$  such that  $|S| = n \geq k$ . The possible outcomes of such experiment are the objects we count.

Here, the four categories of counting situation (ordered or orderless, with or without replacement) also exist in random experiments. The outcome of random experiment is a **sample point**, where the **sample space** ( $\Omega$ ) is the set of all possible outcomes of an experiment.

Let us run through an example:



Figure 13.1.1. The Sample Space of 4 Coin-Tosses



Determining the chance of each particular outcome in the sample space would then have to do with the probability for each sample point.

## 13.2 Defining Probability Space and Events

A **probability space** is a sample space  $\Omega$  with a probability  $\mathbb{P}[\omega]$  for each sample point  $\omega$ , while following the two properties of:

- **Non-negativity:**  $\forall \omega \in \Omega, 0 \leq \mathbb{P}[\omega] \leq 1$
- **Total one:**  $\sum_{\omega \in \Omega} \mathbb{P}[\omega] = 1$

A uniform distribution is some probability space where:

$$\text{for } |\Omega| = N, \text{ then } \forall \omega \in \Omega, \mathbb{P}[\omega] = \frac{1}{N}.$$

When looking at probability spaces, we would also like to look at the chance at which a specific condition occurs, rather than that of when a sample point becomes the outcome. For example, “exactly one head has occurred”, as well as “at least one head has occurred” in one trial. Such conditions are called **events**.

An event is a subset of the sample space.

For example, the event where exactly one head appeared in the span of four coin tosses would be a subset:

$$\{HTTT, THTT, TTHT, TTTH\} \subset \Omega$$

The probability of an event  $A$  would thus be expressed as:

$$\mathbb{P}[A] = \sum_{\omega \in A} \mathbb{P}[\omega]$$

## 13.3 Examples after Examples of Discrete Probability

Let us look at a selection of notable examples across which discrete probability can be exercised practically and helpfully to expand our knowledge of listed contexts:

### 13.3.1 Coin Tosses

The probability of tossing a head in a coin toss is  $\mathbb{P}[H] = \frac{1}{2}$ , and since the result of a toss is binary (it can only be either heads or tails), the probability of tossing a tail  $\mathbb{P}[T] = 1 - \mathbb{P}[H] = \frac{1}{2}$  is complementary to that of tossing a head.

The probability at which we attain  $k$  tosses as head, and  $n - k$  as tails, can be thought of calculating the number of combinations from  $n$  tosses where  $k$  are selected to be head, with each of them having the probability of  $\mathbb{P}[H]^k \mathbb{P}[T]^{n-k}$ .

The result of this thought experiment upon random experiments that yield binary results is otherwise known as:

**Definition 13.3.1. Binomial Probability**

For a random experiment with binary results where the sample points are  $A$  or  $B$ , the probability at which exactly  $k$  out of  $n$  trials have result  $A$  is:

$$\binom{n}{k} \mathbb{P}[A]^k \mathbb{P}[B]^{n-k}$$

The application of binary probability works beyond coin tosses; say, the probability at which a player succeeds in 5 out of 10 baseball throws.

**13.3.2 Balls and Bins**

In this experiment, we allocate  $n$  balls into  $k$  bins, where each ball is equally likely to land in any urn.

In this case, the sample space  $\Omega$  would be a set of  $n$ -element tuple at which each element shows which of the bin (from number 1 to  $k$ ) does the ball of corresponding position in tuple drop into. The cardinality of such a sample space would then be  $k^n$ .

Note that consequentially, the balls and bins in this context are all distinguishable. Therefore, the balls-and-urns context as discussed in Counting portion of COMPSCI 70 would not apply.

**13.3.3 Birthday Paradox**

For an  $n$ -people group, the sample space  $\Omega$  has a cardinality of  $365^n$  because each of the  $n$  people have 365 different possible birthdays (let us disregard leap years).

Here, let us calculate the probability at which at least a pair of people have the same birthday, noted as event  $A$ . It would be much simpler to determine  $1 - \mathbb{P}[\bar{A}]$ : which is the probability that the complement of event  $A$  does not occur. The amount of sample points at which  $\bar{A}$  occurs would be  $365 \times \cdots \times (365 - n + 1)$  (since we have one less available birthday per progression of people).

Then,

$$\mathbb{P}[A] = 1 - \mathbb{P}[\bar{A}] = 1 - \frac{|\bar{A}|}{|\Omega|} = \frac{365 \times \cdots \times (365 - n + 1)}{365^n}$$

And interestingly, with  $n = 23$  people, such probability would be larger than 50%, and at  $n = 60$ ,  $\mathbb{P}[A]$  reaches over 99%!

**13.3.4 Monty Hall, the TV Show Goat-Car Scenario**

You have probably visited this scenario on mathematics-related YouTube videos, or movies like 21.

If not, don't worry, here is a description of the Monty Hall Problem:

You are now on a TV Show, where three gates are presented in front of you. There are also two goats and one car, each arranged to be behind one door.

You are invited to select one initial door, behind which will be your prize.

Then, the game show host, who knows which gate hides the car, will reveal the other gate that is neither the gate you choose nor the gate the car is behind.

The game show host now asks, "Will you switch the gate?" Do you switch the gate?

Let us mark this as a discrete probability problem. Characterize each sample point as a triplet:

(Door of prize, Initial chosen door, Door opened by host)

The only constraint for the triplet would then be  $i \neq k$  (since the game show host would not open a door with car behind it). In that case, let us further categorize the possible sample points:

$$\begin{cases} (i, j, k), & \text{Type A: Each of the doors are distinct, the initial pick was not a car} \\ (i, i, k), & \text{Type B: The initial pick was the car} \end{cases}$$

Here:

For Type A sample points, there are 6 possible such situations, and the probability of encountering such a sample point is  $\frac{1}{3} \times \frac{1}{3} \times \frac{1}{1}$ .

For Type B sample points, meanwhile, there are also 6 possible such situations, and the probability of encountering such a sample point is  $\frac{1}{3} \times \frac{1}{3} \times \frac{1}{2}$ .

Such a sample space does satisfy the non-negative and total-one properties, and therefore is applicable for use. In this case, the probability of winning by switching a door would be the probability at which the sample point is type A, such that switching leads to  $i = j$ . The resulting probability of winning by switching is then  $\frac{6}{9} = \frac{2}{3}$ .

Alternatively, let us consider a total summary of Monty Hall contexts:

#### Practice 13.3.1: Solve for the General Monty Hall Solution

For  $k$  revealed doors and  $c$  cars (each behind a door), the probability that we have picked the current prize door, resulting in a futile switch, is  $\frac{c}{n}$ .

The probability that a switch works can be computed as:

$$\frac{c \cdot (c - 1) + (n - c) \cdot c}{n \cdot (n - 1 - k)} = \frac{c \cdot (n - 1)}{n \cdot (n - 1 - k)}$$

By counting all situations where a switch works either the initial choice was correct or not, and offered that the second round of choice only has  $n - 1 - k$  possible options.

Here, we see that setting  $c = 1$  and  $k = 1$  reflects the original Monty Hall choice problem, where the probability of winning by switch is as well computed to be  $\frac{2}{3}$ .

## 13.4 Summary

When considering discrete probability problems, here are some guidelines to generate solutions along:

- What is the sample space, namely all the possible sample points?
- What is the probability of each sample point?
- What is the event that the probability should estimate for?
- Does the way we aggregate (sum/multiply) probability up make sense?

## Chapter 14

# Conditional Probability, Independence, Combination of Events

### 14.1 Conditional Probability

Sometimes, the sample space available to us at the moment depends on the events that has occurred before it. For example, if we have sampled without replacement, the probability of getting a particular outcome in some random experiment would change due to the change in sample space. Mathematicians have organized a uniform way to describe such phenomenon where the probability of an event changes based on the event that transpires before it:

#### Definition 14.1.1. Conditional Probability

Let  $A$  denote the event which we'd like to know the probability of.

Then, let  $B$  be the event that has already happened and influences the probability of event  $A$ : for example, some urn is empty provided a ball-and-urn situation, or a ball of some color has been priorly sampled without replacement from a finite pool of balls.

The **conditional probability** of  $A$  given  $B$  is then the probability of event  $A$  provided that event  $B$  has occurred, denoted mathematically as:

$$\mathbb{P}[A|B]$$

The computation of such probability can be portrayed with a venn diagram.

If event  $B$  has occurred, then the sample points we count towards our conditional probability would only concern sample points that signify event  $B$ . Among which, the sample points that signify both event  $A$  and  $B$  is in the sample space  $A \cap B$ .

Therefore:

#### Definition 14.1.2. Computation of Conditional Probability

For events  $A, B \subset \Omega$  in the same probability space such that  $\mathbb{P}[B] > 0$ , the conditional probability of  $A$  given  $B$  is

$$\mathbb{P}[A|B] = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[B]}$$

Now, if you wouldn't mind an example to demonstrate the above:

## Example Question 14.1.1: Conditional Probability of Card Dealing Instance

What is the conditional probability that, when dealing 2 cards from a normal deck of poker cards, given the first card is an ace, the second card is also an ace?

In this case,  $\mathbb{P}[B]$  is the probability of the first card being an ace:

$$\frac{4}{52}$$

The probability of  $\mathbb{P}[A \cap B]$  is the probability of both cards being an ace:

$$\frac{4}{52} \frac{3}{51}$$

Therefore:

$$\mathbb{P}[A|B] = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[B]} = \frac{3}{51}$$

Let us do a sanity check for the above example. If we are provided that one ace is already drawn, then it should become harder for the next drawn card to also be ace.

This intuition is coherent with the results of our computation: the original probability of drawing an ace is  $\frac{4}{52}$ , which is higher than the conditional probability we computed.

## 14.2 Bayesian Inference

Conditional probability offers us the opportunity to explore a probability provided some update to the current state of system. For example, in the above card dealing example, by computing the conditional probability of a situation, we computed the probability of something happening provided an update in the state (contents of deck) of the system (deck).

Such mathematical device lends hand to an important mathematical subject, called **Bayesian Inference**: regarding how to update knowledge after making an observation. In such interpretation, we can treat  $\mathbb{P}[A]$  in the computation of  $\mathbb{P}[A|B]$  as a “priori”, and the conditional probability itself a “posterior”.

Let us inspect the following famous application of conditional probability to get a better idea of Bayesian Inference:

## Example Question 14.2.1: Conditional Probability: Accuracy of Clinical Tests

Provided the following results of a clinical trial:

- When applied to an affected person, the test is positive in 90% of cases, and negative in 10% of cases.
- When applied to a negative person, the test is positive in 20% of cases, and negative in 80% of cases.

Suppose 5% of the US population is truly affected (which is our priori), what is the probability that a random person in the US Population is affected, provided the person tests positive?

First of all, the answer is NOT just 5%. Those who are affected can also erroneously be tested as negative. To consider this question in its entire scope, we can construct the following table:

	Affected	Healthy
Positive	AP	HP
Negative	AN	HN

(Note that we would usually say AP is True Positive, HN is True Negative, HP is False Positive... etc.)

In this case, the probability that a tested person is affected given the person tests positive is in fact the condi-

tional probability:

$$\mathbb{P}[AP|Positive] = \frac{\mathbb{P}[AP]}{\mathbb{P}[AP] + \mathbb{P}[HP]}$$

Now, using the priori that *5% of the US population is truly affected*, let's fill in the table above with concrete probability values:

	Affected	Healthy
Positive	$0.05 \times 0.9$	$0.95 \times 0.2$
Negative	$0.05 \times 0.1$	$0.95 \times 0.8$

Substituting the value, we calculate the conditional probability to be:

$$\mathbb{P}[AP|Positive] = \frac{\mathbb{P}[AP]}{\mathbb{P}[AP] + \mathbb{P}[HP]} = \frac{0.05 \times 0.9}{0.05 \times 0.9 + 0.95 \times 0.8} = 0.1914$$

Let's explore another method of performing the above computation.

#### Example Question 14.2.2: An Alternative Solution

Note that for events  $A$  and  $B$ , there are two ways to express the probability  $\mathbb{P}[A \cap B]$ :

$$\mathbb{P}[A \cap B] = \mathbb{P}[A|B]\mathbb{P}[B] = \mathbb{P}[B|A]\mathbb{P}[A]$$

$$\mathbb{P}[B|A] = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[A]}$$

Combining the above with the conditional probability formula, we learned that:

$$\mathbb{P}[A|B] = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[B]} = \frac{\mathbb{P}[B|A]\mathbb{P}[A]}{\mathbb{P}[B]}$$

Since an event  $A$  either occurs or not, we can also conclude that:

$$\begin{aligned} \mathbb{P}[B] &= \mathbb{P}[A \cap B] + \mathbb{P}[\bar{A} \cap B] = \mathbb{P}[B|A]\mathbb{P}[A] + \mathbb{P}[B|\bar{A}]\mathbb{P}[\bar{A}] \\ &= \mathbb{P}[B|A]\mathbb{P}[A] + \mathbb{P}[B|\bar{A}](1 - \mathbb{P}[A]) \end{aligned}$$

Combining all above equations, we obtain that:

$$\mathbb{P}[A|B] = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[B]} = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[B|A]\mathbb{P}[A] + \mathbb{P}[B|\bar{A}](1 - \mathbb{P}[A])}$$

Substituting event  $A$  with being an affected individual and event  $B$  with the probability of testing positive, we would eventually reach the same computational result as outlined before!

## 14.3 Bayes Rule and Total Probability Rule

Turns out our prior efforts have already been found by prior mathematicians, who then named them as follows:

## Definition 14.3.1. Bayes Rule and Total Probability Rule

**Bayes Rule:**

$$\mathbb{P}[B|A] = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[A]} = \frac{\mathbb{P}[A|B]\mathbb{P}[B]}{\mathbb{P}[A]}$$

**Total Probability Rule:**

$$\mathbb{P}[B] = \mathbb{P}[B|A]\mathbb{P}[A] + \mathbb{P}[B|\bar{A}](1 - \mathbb{P}[A])$$

Let's demonstrate the above principles with some examples:

**14.3.1 Example 1: Tennis Match**

The current player is now fighting for the place of top 3 in a tournament, and there are two possible opponents for the current player:  $X$  and  $Y$ .

If person  $X$  is the opponent, then the current player has a 70% chance of winning; otherwise, if the opponent is  $Y$ , the chance of winning becomes 30%.

The chance of getting  $X$  as the opponent is 60%.

Provided the above, what is the probability that the current player wins?

## Example Question 14.3.1: Example 1: Tennis Match

Let us first label each possibilities we compute the probability for as an event:

- Let event  $A$  be that the current player wins.
- Let event  $B_X$  be that  $X$  is the opponent.
- Let event  $B_Y$  be that  $Y$  is the opponent.

We now calculate  $\mathbb{P}[A]$ , and by the Total Probability Rule, we may provide that:

$$\begin{aligned}\mathbb{P}[A] &= \mathbb{P}[A|B_X]\mathbb{P}[B_X] + \mathbb{P}[A|B_Y]\mathbb{P}[B_Y] \\ &= 0.7 \times 0.6 + 0.3 \times 0.4 = 0.54\end{aligned}$$

**14.3.2 Example 2: Balls and Bins**

Suppose we are provided two bins, each filled with different contents of balls:

- The first bin,  $B_1$ , has two white, three black balls.
- The second bin,  $B_2$ , has one white, one black ball.

We choose a random bin to sample a ball from, both bins having the same probability of being picked. Then, what is the probability that we picked  $B_1$  provided that a white ball was drawn?

## Example Question 14.3.2: Example 2: Balls and Bins

Let us, once again, first label each assessing possibilities as events:

- Let event  $A$  be that  $B_1$  was chosen.

Based on observation, we may also have concluded that:

- The probability of drawing a white ball and having picked  $B_1$  was  $\mathbb{P}[w \cap A] = \frac{2}{5} \times \frac{2}{1} = \frac{1}{5}$ .
- The probability of having picked  $B_1$  was  $\mathbb{P}[A] = \frac{1}{2}$ .
- The probability of having drawn a white ball was  $\mathbb{P}[w] = \mathbb{P}[A] \times \frac{2}{5} + \mathbb{P}[\bar{A}] \times \frac{1}{2} = \frac{9}{20}$ .

Using Bayes' Rule:

$$\begin{aligned}\mathbb{P}[A|w] &= \frac{\mathbb{P}[A \cap w]}{\mathbb{P}[w]} \\ &= \frac{\frac{2}{10}}{\frac{9}{20}} = \frac{4}{9}\end{aligned}$$

### 14.3.3 Generalization of Bayes Rule and Total Probability Rule

To consider these principles of probability in a more general concept across an event being divisible into two aspects, let us first discuss how events can be partitioned.

In other words, rather than only considering that an event only has two cases (positive or negative, head or tails), let us discuss how an event can be divided into having  $n \geq 2$  cases (for example, picking from a pool of balls in  $n$  colors):

#### Definition 14.3.2. Partition of Event

An event  $A$  is partitioned into  $n$  (sub)events  $A_1, \dots, A_n$  if:

- $A = A_1 \cup \dots \cup A_n$
- All such subevents (partitions) are pairwise disjoint. That is,

$$\forall (i, j) \in \{(x, y) : x \neq y \wedge 1 \leq x \leq y \leq n\}, A_i \cap A_j = \emptyset$$

Therefore, each outcome in  $A$  belongs to exactly one of the subevents  $A_1, \dots, A_n$ .

Upon this view of partitions, we may generalize as such:

#### Definition 14.3.3. Generalization of Bayes Rule and Total Probability Rule

Let  $A_1, \dots, A_n$  be a partition of sample space  $\Omega$ .

The Total Probability Rule of event  $B$  is:

$$\mathbb{P}[B] = \sum_{i=1}^n \mathbb{P}[B \cap A_i] = \sum_{i=1}^n \mathbb{P}[B|A_i]\mathbb{P}[A_i]$$

The Bayes' Rule of event  $A_i$  given  $B$  is:

$$\mathbb{P}[A_i|B] = \frac{\mathbb{P}[B|A_i]\mathbb{P}[A_i]}{\mathbb{P}[B]} = \frac{\mathbb{P}[B|A_i]\mathbb{P}[A_i]}{\sum_{j=1}^n \mathbb{P}[B|A_j]\mathbb{P}[A_j]}$$

## 14.4 Combination of Events

In Computer Science, we frequently discuss the probability of some combination of events. This has sizeable implication on the operations of Data Science as well as many other significant contributions and subfields of computer science.

Therefore, it is worthy to discuss the probability theories of combination of events with this section. We will first introduce independent events, then the intersections of events.

### 14.4.1 Independent Events

Let us first define the term "independent events":



**Definition 14.4.1. Independent Event**

Two events  $A, B$  in the same probability space are independent if  $\mathbb{P}[A \cap B] = \mathbb{P}[A] \times \mathbb{P}[B]$ .  
In turn, the following three conditions all show independence of  $A$  and  $B$ :

- $\mathbb{P}[A \cap B] = \mathbb{P}[A] \times \mathbb{P}[B]$
- $\mathbb{P}[A|B] = \mathbb{P}[A]$
- $\mathbb{P}[B|A] = \mathbb{P}[B]$

Then, let us discuss independent events on some wider scale:

**Definition 14.4.2. Mutual Independence**

Events  $A_1, \dots, A_n$  are said to be mutually independent if for every subset  $I \subseteq \{1, \dots, n\}$  with cardinality greater than 1:

$$\mathbb{P}[\cap_{i \in I} A_i] = \prod_{i \in I} \mathbb{P}[A_i]$$

Equivalently, we can express such expression in an alternative form: for all  $B_i \in \{A_i, \overline{A_i}\}$ :

$$\mathbb{P}[B_1 \cap \dots \cap B_n] = \prod_{i=1}^n \mathbb{P}[B_i]$$

Using a pair-by-pair perspective, we are also able to find that, for mutually independent events  $A_1, \dots, A_n$ :

$$\mathbb{P}[A_i | \cap_{j \in I} A_j] = \mathbb{P}[A_i]$$

However, it is important to note that pair-wise independence does not imply mutual independence, since an event can have dependence on a combination of event rather than a singular one.

**14.4.2 Intersection of Events**

Let us first consider what we have learned from conditional probability:

**Theorem 14.4.1. Product Rule (Chain Rule)**

For any events  $A, B$ , we may conclude:

$$\mathbb{P}[A \cap B] = \mathbb{P}[A] \mathbb{P}[B|A] = \mathbb{P}[B] \mathbb{P}[A|B]$$

On a general scale, for events  $A_1, \dots, A_n$ :

$$\mathbb{P}[\cap_{i=1}^n A_i] = \mathbb{P}[A_1] \times \mathbb{P}[A_2|A_1] \times \mathbb{P}[A_3|A_2 \cap A_1] \times \dots \times \mathbb{P}[A_n | \cap_{i=1}^{n-1} A_i]$$

**Reasoning:** Let us perform a mathematical induction on  $n$ .

At base case, we have one event to consider, and the trivial case holds.

At an inductive case, for an arbitrary  $n > 1$ , let the induction hypothesis be that:

$$\mathbb{P}[\cap_{i=1}^{n-1} A_i] = \mathbb{P}[A_1] \times \mathbb{P}[A_2|A_1] \times \mathbb{P}[A_3|A_2 \cap A_1] \times \dots \times \mathbb{P}[A_{n-1} | \cap_{i=1}^{n-2} A_i]$$

We may then apply the definition of conditional probability, where:

$$\mathbb{P}[\cap_{i=1}^{n-1} A_i \cap A_n] = \mathbb{P}[\cap_{i=1}^{n-1} A_i] \mathbb{P}[A_n | \cap_{i=1}^{n-1} A_i]$$

Which leads to the exact formula outlined by the Product Rule.

### 14.4.3 Union of Events

The tricky part about calculating for the union of events is that: if we directly sum up the probability of members of unions, we make the calculation extremely frequently prone to overcount, if not always. However, we have learned a principle that helps us to eliminate overcountings when counting over unions of sets: Inclusion-Exclusion.

#### Definition 14.4.3. Inclusion-Exclusion Principle in Probability

For events  $A_1, \dots, A_n$  in some probability space, where  $n \geq 2$ :

$$\mathbb{P}[A_1 \cup \dots \cup A_n] = \sum_{k=1}^n \left( (-1)^{k-1} \sum_{S \subseteq \{1, \dots, n\}, |S|=k} \mathbb{P}[\cap_{i \in S} A_i] \right)$$

We see that this reflects to the original inclusion-exclusion principle:

$$|A_1 \cup \dots \cup A_n| = \sum_{k=1}^n (-1)^{k-1} \sum_{S \subseteq \{1, \dots, n\}, |S|=k} |\cap_{i \in S} A_i|$$

But when  $n$  is large, ain't nobody got time for inclusion-exclusion!

So here are a bit of tips to deal with union of events when there are too many cases to consider:

1. If events  $A_1, \dots, A_n$  that we consider are naturally Mutually Exclusive Events, then

$$\mathbb{P}[A_1 \cup \dots \cup A_n] = \sum_{i=1}^n \mathbb{P}[A_i]$$

2. Following the above point, we recognize an opportunity for sanity check in any calculation for unions of events:

$$\mathbb{P}[A_1 \cup \dots \cup A_n] \leq \sum_{i=1}^n \mathbb{P}[A_i]$$

## Chapter 15

# Random Variables: Distribution and Expectation

### 15.1 Introduction to Random Variables

Experiments need measurements. Without measurements, the results of experiments may not be quantified. In a random experiment, our measurements would naturally be associated with the notion of probability.

We can characterize the result of an experiment by a property of it.

For example, in the random experiment of flipping 4 coins, let  $h$  stand for the number of heads that appeared in one run of such experiment. Such value, while random, quantifies the result of an experiment.

In this case, the probability of  $h$  being a value in natural numbers from 0 to 4 can be calculated via binomial probability. The random outcome of this probabilistic experiment is called a **random variable**. We will come across a rigorous definition soon; for now, we will briefly introduce the properties of a random variable.

As spoken before, each random variable can have a probability distribution over the set of possible values it can be. It is random, in the perspective that its value occurs randomly with probability. The property of randomness is practical in this regard: we can thus quantify the likelihood and average of an experiment's outcome via calculations along a known probability distribution.

Then, from here on, let us discuss random variables in the context of some random experiment.

#### Definition 15.1.1. Random Variable

A random variable  $X$  on a sample space  $\Omega$  is a function  $X : \Omega \rightarrow \mathbb{R}$  that assigns to each sample point  $\omega \in \Omega$  a real number  $X(\omega)$ .

But in that regard, we see that a random variable is in fact a function, mapping sample points onto some real number!

### 15.2 Probability Distribution

A probability space comes along two properties:

- The sample space  $\Omega$  consisting of all possible outcomes of the experiment.
- The probability of each of the sample points.

We may find similar features in random variables.

- The set of possible values.
- The probability that it is of some specific value.

Why not, then, characterize the set of all outcomes for which a random variable takes on a specific value, as an event (or alternatively, a partition of the sample space)?

$$\{\omega \in \Omega : X(\omega) = a\}$$

In this case, the probability by which the random variable is a specific value can be spoken of as the probability of an event.

The collection of these probabilities is known as the **distribution** of the random variable  $X$ :

#### Definition 15.2.1. Distribution

The distribution of a discrete random variable  $X$  is the collection of values:

$$\{(a, \mathbb{P}[X = a]) : a \in \mathcal{A}\}$$

where  $\mathcal{A}$  denotes the set of all possible values taken by  $X$ .

Any distribution, as the sum of probabilities of all events available, should thus have an area under curve of 1.

## 15.2.1 Types of Distributions

There are several types of probability distributions that prove to be generally useful. We will go through some of these distributions below, and some in future lecture notes as we progress through CS70.

#### Definition 15.2.2. Bernoulli Distribution

A Bernoulli Distribution of a random variable that takes value in  $\{0, 1\}$  (which makes the random variable binary, as well as making it an indicator variable).

The probability of an outcome is:

$$\mathbb{P}[X = i] = \begin{cases} p, & \text{if } i = 1 \\ 1 - p, & \text{if } i = 0 \end{cases}$$

This is essentially a histogram of two bars.

#### Definition 15.2.3. Binomial Distribution

The binomial distribution outlines the probability of a random experiment whose outcome is essentially binary. For example, coin flipping yields a binary result (heads or tails).

In this case, the probability by which  $k$  results with probability  $p$  occurs in a total of  $n$  runs is, once again:

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

The binomial distribution using a random variable to represent the number of results that occurs with probability  $p$  would then be denoted as:

$$X \sim \text{Bin}(n, p)$$

#### Definition 15.2.4. Hypergeometric Distribution

Let the results of a random example be binary, but in a total of  $N$  sample points, there are  $A$  sample points of one specific result.

The probability at which we may acquire  $k$  results of  $A$  out of  $n$  sample points sampled without replacement is

slightly complicated, but can be generalized as:

$$\mathbb{P}[Y = k] = \binom{n}{k} \frac{\frac{B!}{(B-k)!} \frac{(N-B)!}{(N-B-(n-k))!}}{\frac{N!}{(N-n)!}} = \frac{\binom{B}{k} \binom{N-B}{n-k}}{\binom{N}{n}}$$

This forms what we call a hypergeometric distribution with parameters  $N, B, n$ , and we denote as such:

$$Y \sim \text{Hypergeometric}(N, B, n)$$

## 15.3 Multiple Random Variables and Independence

A sample space can also contain multiple random variables. For example, in the coinflip random experiment, I can also define random variables for:

- The number of tails.
- The number of consecutives.
- The number of nonconsecutives.
- The number of coins that landed on the United States of America as it is being tossed from the X Æ A-12 planet.

That last item didn't cause a LaTeX compilation error.

In any case, we see that some random variables above are independent of each other, while some are not. Meanwhile, we may also see that we can form distributions of multiple two discrete random variables:

### Definition 15.3.1. Joint Distribution

The joint distribution for two discrete random variables  $X$  and  $Y$  is the collection of values:

$$\{(a, b), \mathbb{P}[X = a, Y = b] : a \in \mathcal{A}, b \in \mathcal{B}\}$$

where  $\mathcal{A}$  and  $\mathcal{B}$  are respectively the set of all possible values for random variables  $X$  and  $Y$ .

The joint distribution of random variables  $X_1, \dots, X_n$  is denoted as:

$$\mathbb{P}[X_1 = a_1, \dots, X_n = a_n]$$

We see that we may have potentially formed a two-dimensional histogram with joint distribution! Following that idea, we can thus see the distribution of one random variable when the other stays at a constant value.

### Definition 15.3.2. Marginal Distribution

Readdressing the above idea: when given a joint distribution of RVs  $X$  and  $Y$ , the distribution  $\mathbb{P}[X = a]$  for  $X$  is called the **marginal distribution** of  $X$ , found via summing over the values of  $Y$ :

$$\mathbb{P}[X = a] = \sum_{b \in \mathcal{B}} \mathbb{P}[X = a, Y = b]$$

Furthermore, as foreshadowed, random variables can be dependent or independent upon each other:

### Definition 15.3.3. Independence of Random Variables

Random variables  $X$  and  $Y$  on the same probability space are independent of each other if:

$$\forall a, b (\mathbb{P}[X = a, Y = b] = \mathbb{P}[X = a] \mathbb{P}[Y = b])$$

## 15.4 Expectation

To calculate the center of a distribution, we can, instead of computing over every single sample point, calculate the expectation of a distribution:

### Definition 15.4.1. Expectation

The expectation of a discrete random variable  $X$  is defined as:  $\mathbb{E}[X] = \sum_{a \in \mathcal{A}} a\mathbb{P}[X = a]$  or verbally, the sum over the product of all possible values taken by the random value and their respective probabilities of occurring.

### 15.4.1 Linearity of Expectation

Due to the nature of expectation as a weighted average, it has several linear properties.

For one:

#### Theorem 15.4.1. Linearity of Expectation with Respect to RV

For any two random variables  $X$  and  $Y$  on the same probability space,

$$E[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y] \quad (15.1)$$

$$\mathbb{E}[cX] = c\mathbb{E}[X] \quad (15.2)$$

# Chapter 16

## Variance

### 16.1 Functions of Random Variables

A random variable on a space can be a sum of other random variables. In this case, we see one random variable being a function of some other random variables. So, what other functional relationship might exist outside this?

Formally, let  $X$  be a random variable on a sample space  $\Omega$  with probability distribution  $\mathbb{P}_X$ , then for a function  $f(\cdot)$  on the range of  $X$ ,  $f(X)$  is a random variable  $Y$  on the same sample space where  $Y(\omega) = f(X(\omega))$ ,  $\omega \in \Omega$ .

In terms of sample spaces, the event  $Y = y$  is essentially the event  $X \in f^{-1}(y)$ , where we take  $f^{-1}(y) = \{x | f(x) = y\}$ . With this view, the distribution of  $Y$  can thus be derived from the distribution of  $X$  as:

$$\mathbb{P}_Y[Y = y] = \sum_{x: f(x)=y} \mathbb{P}_X[X = x]$$

The expectation of such variable  $Y$  thus becomes:

$$\mathbb{E}[Y] = \sum_y y \mathbb{P}_Y[Y = y] = \sum_x f(x) \mathbb{P}_X[X = x]$$

Which therefore translates into:

**Theorem 16.1.1. Law of the Unconscious Statistician**

$$\mathbb{E}[f(X)] = \sum_x f(x) \mathbb{P}_X[X = x]$$

For example:

$$\mathbb{E}[X^2] = \sum_x x^2 \mathbb{P}_X[X = x]$$

### 16.2 Introducing and Defining Variance

Now that we have certified a way to express what is the centroid of a distribution, we may concern with a new subject: what is the deviation of points in a distribution from its centroid, and how may we quantify the overall deviation of a distribution?

Let us consider a simple context.

You are being pursued by Oski.

In this case, you may flip a coin to decide whether you move from your current location forward 1 step (by flipping a head), which we'll call the origin, or move backward 1 step.

To avoid being caught by Oski and captured into the terrifyin sanctuary of CS61D, Discrete Memematics and Probably

Theory, you must figure first what is the expectation of your distance from the origin.

Let us call the number of steps you take at last  $S$ , a random variable, as a sum of other random variables  $X_i$  that record the number of steps you perform at the  $i^{\text{th}}$  flip:

$$\mathbb{E}[S] = \mathbb{E}[X_1 + \cdots + X_n] = \mathbb{E}[X_1] + \cdots + \mathbb{E}[X_n] = n \times \left(\frac{1}{2} \times 1 + \frac{1}{2} \times -1\right) = 0$$

LMAO, have fun being caught by Oski.

But here arises a different question. What is the expected distance of the player from the origin? What is the expected value of  $|S|$ ?

Since the absolute value is difficult to work with, alternatively, let us consider what is the expected value of  $S^2$ .

But the essential quality we derive for is, in fact,  $\mathbb{E}[(S - \mu)^2]$ , which captures not just the value of distance from origin but its deviation from the centroid of distribution!

In this case, we may derive the expectation of  $S^2$  as follows:

$$\begin{aligned} \mathbb{E}[S^2] &= \mathbb{E}[(X_1 + \cdots + X_n)^2] \\ &= \mathbb{E}\left[\sum_{i=1}^n X_i^2 + 2 \sum_{i < j} X_i X_j\right] \\ &= \sum_{i=1}^n \mathbb{E}[X_i^2] + 2 \sum_{i < j} \mathbb{E}[X_i X_j] \\ \mathbb{P}[X_i X_j = 1] &= \frac{1}{2} \text{ (either both heads or both tails)} \\ \mathbb{P}[X_i X_j = -1] &= \frac{1}{2} \text{ (one head one tail)} \\ \forall i, j \mathbb{E}[X_i X_j] &= 0 \\ \mathbb{E}[S^2] &= \sum_{i=1}^n \mathbb{E}[X_i^2] = n \end{aligned}$$

Let us now introduce a proper name for this quality:

#### Definition 16.2.1. Variance

For a random variable  $X$  with expectation  $\mathbb{E}[X] = \mu$ , the variance of  $X$  is defined as:

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

And the square root  $\sigma(X) := \sqrt{\text{Var}(X)}$  is called the standard deviation of  $X$ .

We can, in fact, develop an easier expression for variance:

#### Theorem 16.2.1. Formula of Variance

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

Furthermore, the property of variance provides that:

$$\text{Var}(cX) = c^2 \text{Var}(X)$$



## 16.3 Sum of Independent Random Variables

If a random variable  $X$  is the sum of independent random variables:

$$X = X_1 + \cdots + X_n$$

then its variance is the sum of variances of individual random variables.

In particular, for independent random variables  $X_1, \dots, X_n$ :

$$\text{Var}(X_1 + \cdots + X_n) = \sum_{i=1}^n \text{Var}(X_i)$$

Let us explore a similar rule on the product of independent random variables:

### Theorem 16.3.1. Expectation of Product of Random Variables

For independent random variables  $X, Y$ :

$$\begin{aligned} \mathbb{E}[XY] &= \sum_a \sum_b ab \mathbb{P}[X = a, Y = b] \\ &= \sum_a \sum_b ab \mathbb{P}[X = a] \mathbb{P}[Y = b] \\ &= \sum_a a \mathbb{P}[X = a] \times \sum_b b \mathbb{P}[Y = b] = \mathbb{E}[X] \mathbb{E}[Y] \end{aligned}$$

### Theorem 16.3.2. Variance of Sum of Independent Random Variables

For independent random variables  $X, Y$ :

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

## 16.4 Covariance and Correlation

The expression we see for the above proof,  $\mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$ , is a measure of association between variables  $X$  and  $Y$ :

### Definition 16.4.1. Covariance

Covariance of random variables  $X$  and  $Y$ , denoted as  $\text{Cov}(X, Y)$ , is defined as:

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

Here are some interesting properties:

- If  $X$  and  $Y$  are independent, then  $\text{Cov}(X, Y) = 0$ . However, the converse is not true.
- $\text{Cov}(X, X) = \text{Var}(X)$ .
- Covariance is bilinear, which means:

$$\text{Cov}\left(\sum_{i=1}^n a_i X_i, \sum_{j=1}^m b_j Y_j\right) = \sum_{i=1}^n \sum_{j=1}^m a_i b_j \text{Cov}(X_i, Y_j)$$

To measure the correlation of random variables, whose standard deviations are both nonzero, we may use a different statistics:

Definition 16.4.2. Correlation

Suppose  $X$  and  $Y$  are random variables with  $\sigma(X) > 0$  and  $\sigma(Y) > 0$ , then the correlation of  $X$  and  $Y$  is defined as:

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

And this value must be such that:

$$-1 \leq \text{Corr}(X, Y) \leq 1$$

It is worthy to note that while all independent random variables have zero covariance, zero-covariance itself does not imply independence.

## Chapter 17

# Concentration Inequalities and the Laws of Large Numbers

### 17.1 Introduction

How do we estimate the bias of a binary random experiment like coin-flipping? How do we measure the overall deviation of an expected hypothetical result and what was really observed? How do we know if the coin is a cheated coin?

Such a quantity for a random experiment of  $n$  trials can be expressed as  $\hat{p} = \frac{1}{n}S_n$ , where  $S_n$  is the number of success in a binary experiment.

And, let  $p$  denote the true bias of the coin, which is unknown, and we are offered that  $\mathbb{E}[S_n] = np$ , we see that the estimator  $\hat{p}$  works to align with  $p$ .

Therefore,

This means when  $n$  is sufficiently large, the estimator of bias is close to the true value of bias.

We will discuss this rule again later.

But, how large should  $n$  be, so to provide an upper bound to the error between the estimator and bias itself?

We do not truly know, because the estimator of bias itself is also a random variable.

Let us analyze the estimator in the perspectives of random variable:

$$\hat{p} = \frac{1}{n}S_n = \frac{1}{n} \sum_{i=1}^n X_i$$

where  $X_i$  is an indicator random variable that suggests whether the  $i^{th}$  trial of experiment was noted as successful. This means the estimator  $\hat{p}$  is an average of indicator random variables, therefore also a random variable itself.

It is somehow reasonable that there is no concrete theorem for an upperbound of  $n$  such that the estimator estimates the real bias of some binary experiment at extreme accuracy. Or else I wouldn't be malding over gacha games and hack with mathematics.

But, rather than trying to guarantee that the error of estimation has some upperbound, we can demand it to have an upper bound with some confidence  $1 - \delta$ , therefore:

$$\mathbb{P}[|\hat{p} - p| \leq \epsilon] \geq 1 - \delta$$

Spoiler alert: to guarantee and so, we require that:

$$n \geq \frac{1}{4\epsilon^2\delta}$$

We will also discuss this later within the section of *Chebyshev's Inequality*.

## 17.2 Markov's Inequality

Let us discuss the upper bound of error for experiments on nonnegative random variables:

### Theorem 17.2.1. Markov's Inequality

For a nonnegative random variable  $X$  with finite mean, for any positive constant  $c$ :

$$\mathbb{P}[X \geq c] \leq \frac{\mathbb{E}[X]}{c}$$

#### Proof:

Let  $\mathcal{A}$  denote the range of values that random variable  $X$  is in. Then:

$$\begin{aligned} \mathbb{E}[X] &= \sum_{a \in \mathcal{A}} a \times \mathbb{P}[X = a] \\ &\geq \sum_{a \geq c} a \times \mathbb{P}[X = a] \\ &\geq \sum_{a \geq c} c \times \mathbb{P}[X = a] \\ &= c \sum_{a \geq c} \mathbb{P}[X = a] \\ &= c \mathbb{P}[X \geq c] \end{aligned}$$

Here, say we let  $c = 2\mathbb{E}[X]$ , then we see that at most half of the individuals can reach some score that is double the mean in a game.

We can further generalize Markov's Inequality onto random variables based on absolute values:

### Theorem 17.2.2. Generalized Markov's Inequality

Let  $Y$  be an arbitrary random variable with finite mean. Then, for any positive constants  $c, r$ ,

$$\mathbb{P}[|Y| \geq c] \leq \frac{\mathbb{E}[|Y|^r]}{c^r}$$

#### Proof:

For  $c > 0, r > 0$ , let  $I$  be an indicator function:

$$|Y|^r \geq |Y|^r I(|Y| \geq c) \geq c^r I(|Y| \geq c)$$

Taking expectations of both sides then offers that:

$$\mathbb{E}[|Y|^r] \geq c^r \mathbb{E}[I(|Y| \geq c)] = c^r \mathbb{P}[|Y| \geq c]$$

Markov's inequality also concerns the center of its probability distribution:  $\mu = \mathbb{E}[X]$ , as the fulcrum of distributional mass.

To let the mass on left of and right of  $\mu$  be equivalent, up until some point at the right  $k\mu$ , the mass must be 0, while to keep balance, the mass at the left must then be  $(k-1)$  times of the mass at the right:

$$(k-1)mass_{right} \times 1 = mass_{right} \times (k-1)$$

This forms a known phenomenon called “Markov's Bound” with  $\alpha = k\mu$ , where  $\mathbb{P}[X > k\mu] = 0$  can be noted as some constraint of a distribution with  $k$  being a value as regulated by the fulcrum-determined balance addressed above.

## Example Question 17.2.1: Coin Tosses and Markov

Consider fair coin toss, and let  $X$  denote the number of heads observed. What is the probability of observing more than  $\frac{3}{4}n$  heads? Let's apply Markov inequality to obtain an upper bound on  $\mathbb{P}[X \geq \frac{3}{4}n]$ , and since  $\mathbb{E}[X] = \frac{1}{2}n$ :

$$\mathbb{P}[X \geq \frac{3}{4}n] \leq \frac{\mathbb{E}[X]}{c} = \frac{2}{3}$$

## 17.3 Chebyshev's Inequality

### Theorem 17.3.1. Chebyshev's Inequality

For a random variable  $X$  with finite expectation  $\mathbb{E}[X] = \mu$ , for any positive constant  $c$ :

$$\mathbb{P}[|X - \mu| \geq c] \leq \frac{\text{Var}(X)}{c^2}$$

**Proof:** Let  $Y = (X - \mu)^2$ , then we discover that  $\mathbb{E}[Y] = \text{Var}(X)$ . Meanwhile, this presents us that we are looking for the probability where  $Y \geq c^2$ , which we may apply Markov's Inequality from, to derive the following result:

$$\mathbb{P}[|X - \mu| \geq c] = \mathbb{P}[Y \geq c^2] \leq \frac{\mathbb{E}[Y]}{c^2} = \frac{\text{Var}(X)}{c^2}$$

In fact, since variance is  $\sigma^2$ , we may understand the Chebyshev's inequality as the following Corollary:

For any random variable  $X$  with finite expectation  $\mu$  and finite standard deviation  $\sigma$ :

$$\mathbb{P}[|X - \mu| \geq k\sigma] \leq \frac{\sigma^2}{k^2\sigma^2} = \frac{1}{k^2}$$

## 17.4 The Law of Large Numbers

### Theorem 17.4.1. Law of Large Numbers

Let  $X_1, \dots, X_n$  be a sequence of Independent and Identically Distributed random variables with common finite expectation  $\mathbb{E}[X_i] = \mu$  for all  $i$ , then their partial sum:  $S_n = \sum_{i=1}^n X_i$ :

$$\lim_{n \rightarrow \infty} \mathbb{P}\left[\left|\frac{1}{n}S_n - \mu\right| < \varepsilon\right] = 1$$

for any positive constant  $\varepsilon$ .

**Proof:** Let  $\text{Var}(X_i) = \sigma^2$  be the common variance of random variables, we assume that  $\sigma^2$  is finite. Then, we can obtain from known means and variances that:

$$\mathbb{E}\left[\frac{1}{n}S_n\right] = \mu$$

$$\text{Var}\left(\frac{1}{n}S_n\right) = \frac{1}{n^2}\text{Var}(S_n) = \frac{1}{n^2}\sum_{i=1}^n \text{Var}(X_i) = \frac{\sigma^2}{n}$$

Therefore, by Chebyshev's inequality, we have:

$$\mathbb{P}\left[\left|\frac{1}{n}S_n - \mu\right| \geq \varepsilon \leq \frac{\frac{\sigma^2}{n}}{\varepsilon^2} = \frac{\sigma^2}{n\varepsilon^2}\right]$$

Which we can see that the limit of such value approaches 0 when  $n$  approaches  $\infty$ .

Therefore, its complement would approach 0 when  $n$  approaches  $\infty$ .

## Chapter 18

# Geometric and Poisson Distributions

### 18.1 Geometric Distribution

The geometric distribution frequently occurs in situations where we consider “how long we have to wait” for something to occur, such as “how many run until error”, or “how many rolls until SSR”.

The **geometric distribution** is then mathematically defined as below:

#### Definition 18.1.1. Geometric Distribution

A random variable  $X$  such that:

$$\mathbb{P}[X = i] = (1 - p)^{i-1} p, i \in \mathbb{Z}^+$$

is said to have a geometric distribution with parameter  $p$ , mathematically expressable as:

$$X \sim \text{Geometric}(p)$$

This composition of the distribution seems reasonable!

Say, the probability at which you get an ultra rare item from a game is 0.8%. Then, the probability at which you get this item at the  $k^{\text{th}}$  attempt (each attempt is one attempt at the 0.8% probability, as implied), is:

$$\mathbb{P}[X = k] = (1 - 0.008)^{k-1} 0.008$$

as the geometric distribution implies.

All probabilities calculated here are products of powers of nonnegative numbers; thus, will never produce negative probabilities. In that case, we only need the one-total property to hold true for this distribution to appear mathematically sound.

Does such distribution sum to 1? Allow us to verify:

$$\begin{aligned} \sum_{i=1}^{\infty} \mathbb{P}[X = i] &= \sum_{i=1}^{\infty} (1 - p)^{i-1} p \\ &= p \sum_{i=1}^{\infty} (1 - p)^{i-1} \\ &= p \lim_{n \rightarrow \infty} \frac{1 - (1 - p)^n}{1 - (1 - p)} \\ &= p \frac{1}{1 - 1 + p} = 1 \end{aligned}$$

The distribution of  $X$  would thus be some curve that decreases by a factor of  $(1 - p)$  at each step; in other words, an exponential decay. We will see that it is, in fact, an analog to what we will later introduce as the exponential

distribution, in the last chapters of these notes.

And from this above fact, we would know that the amplitude of a geometric distribution's curve is dependent on  $p$ , while the speed at which the curve's height descend across  $x$  depends on  $(1 - p)$ .

### 18.1.1 Mean and Variance of a Geometric RV

Let us first discuss the expectation of a geometric random variable.

If directly computed, we would obtain some result like what follows:

$$\mathbb{E}[X] = \sum_{i=1}^{\infty} i\mathbb{P}[X = i] = p \sum_{i=1}^{\infty} i(1-p)^{i-1}$$

However, there is a way to simplify this above complication:

#### Theorem 18.1.1. Tail Sum Formula

**Theorem:**

$$\mathbb{E}[X] = \sum_{i=1}^{\infty} \mathbb{P}[X \geq i]$$

**Proof:**

$$\begin{aligned} \mathbb{E}[X] &= \sum_{i=1}^{\infty} i\mathbb{P}[X = i] \\ &= p \sum_{i=1}^{\infty} i(1-p)^{i-1} \\ &= 0p_0 + 1p_1 + 2p_2 + \dots \\ &= (p_1 + p_2 + \dots) + (p_2 + \dots) + \dots \\ &= \mathbb{P}[X \geq 1] + \mathbb{P}[X \geq 2] + \dots \end{aligned}$$

What is the significance of this?

#### Theorem 18.1.2. The Expectation of a Geometric RV

$$\begin{aligned} \mathbb{P}[X \geq i] &= (1-p)(1-p)^{i-1} + p(1-p)^{i-1} \\ &= (1-p)^{i-1} \\ \mathbb{E}[X] &= \sum_{i=1}^{\infty} \mathbb{P}[X \geq i] \\ &= \sum_{i=1}^{\infty} (1-p)^{i-1} = \frac{1}{1-(1-p)} = \frac{1}{p} \end{aligned}$$

Upon that, let's also discuss and derive such identity:

#### Theorem 18.1.3. Expectation of $X(X-1)$

**Theorem:**

$$\mathbb{E}[X(X-1)] = \frac{2(1-p)}{p^2}$$



**Proof:** Preliminarily,

$$\begin{aligned}\mathbb{E}[X(X-1)] &= \sum_{i=1}^{\infty} i(i-1) \times \mathbb{P}[X=i] \\ &= \sum_{i=1}^{\infty} i(i-1) \times p(1-p)^{i-1} \\ &= p(1-p) \sum_{i=1}^{\infty} i(i-1)(1-p)^{i-2}\end{aligned}$$

With some calculus and an identity we just learned,

$$\begin{aligned}\sum_{i=1}^{\infty} (1-p)^i &= \frac{1}{p} \\ \frac{d}{dp} \left( \frac{1}{p} \right) &= \sum_{i=1}^{\infty} -i(1-p)^{i-1} \\ &= -\frac{1}{p^2} \\ \frac{d^2}{d^2} \left( \frac{1}{p} \right) &= \sum_{i=1}^{\infty} i(i-1)(1-p)^{i-2} \\ &= \frac{2}{p^3}\end{aligned}$$

Using the result above:

$$\begin{aligned}\mathbb{E}[X(X-1)] &= p(1-p) \sum_{i=1}^{\infty} i(i-1)(1-p)^{i-2} \\ &= p(1-p) \frac{2}{p^3} \\ &= \frac{2(1-p)}{p^2}\end{aligned}$$

Using the linearity of expectation, we can elegantly calculate the variance of a geometric distribution:

**Theorem 18.1.4. The Variance of a Geometric RV**

**Theorem:**

$$\text{Var}(X) = \frac{1-p}{p^2}$$

**Proof:**

$$\begin{aligned}\text{Var}(X) &= \mathbb{E}[X^2] - (\mathbb{E}[X])^2 \\ &= \mathbb{E}[X^2 - X] + \mathbb{E}[X] - (\mathbb{E}[X])^2 \\ &= \frac{2(1-p)}{p^2} + \frac{1}{p} - \frac{1}{p^2} \\ &= \frac{2-2p+p-1}{p^2} = \frac{1-p}{p^2}\end{aligned}$$

In addition, geometric distributions have a property called “memoryless”:

**Definition 18.1.2. Memoryless**

In memoryless distributions, prior trials do not affect the number of additional trials needed. Mathematically:

$$\mathbb{P}[X > n + m | X > m] = \mathbb{P}[X > n]$$

For example, say we have already tossed the coin  $m$  times, what is the probability that we need more than  $n$  additional tosses before getting our first head?

$$\begin{aligned} \mathbb{P}[X > n + m | X > m] &= \frac{\mathbb{P}[X > n + m]}{\mathbb{P}[X > m]} \\ &= \frac{(1-p)^{n+m}}{(1-p)^m} \\ &= (1-p)^n = \mathbb{P}[X > n] \end{aligned}$$

In this sense, “memorylessness” refers to how the distribution, when calculating a tail, will not remember the given condition ( $X > m$ ) since it doesn’t contribute to the end result of a tail even given that condition.

Furthermore, let us perform an interesting mathematical derivation as follows:

**Practice 18.1.1: The Distribution of  $\min(X, Y)$** 

**Question:** What is the distribution of  $\min(X, Y)$ , provided *independent* random variables  $X \sim \text{Geometric}(p)$  and  $Y \sim \text{Geometric}(q)$ ?

**Solution:** Here, let us attempt to find the probability at which  $\min(X, Y) = k$ :

$$\begin{aligned} \mathbb{P}[\min(X, Y) = k] &= \mathbb{P}[X = k, Y \geq k] + \mathbb{P}[Y = k, X > k] \\ &= \mathbb{P}[X = k \cap Y \geq k] + \mathbb{P}[Y = k \cap X > k] \\ &= (1-p)^{k-1} p (1-q)^k + (1-q)^{k-1} q (1-p)^{k-1} \\ &= (1-p)^{k-1} (1-q)^{k-1} (p(1-q) + q) = \text{Geometric}(1 - (1-p)(1-q)) \end{aligned}$$

Furthermore, because

$$X + Y = \min(X, Y) + \max(X, Y)$$

We can use the distributions of  $X \sim \text{Geometric}(p)$ ,  $Y \sim \text{Geometric}(q)$ , and  $\min(X, Y) \sim \text{Geometric}(1 - (1-p)(1-q))$  to infer about  $\max(X, Y)$  as a random variable.

**18.1.2 Example (Application)**

What other greater application can CS70 be made on other than probabilistic problems, and more specifically, gacha games!?

Let us look at the “Coupon Collector Problem”:

Let us attempt to collect a set of  $n$  different cards.

We can get one card by performing one summon, each of the  $n$  cards equivalently likely to appear from that summon.

How many summons do we have to do until we have collected at least one copy of every card?

Let  $S_n$  denote the number of summoning we need to perform to collect all  $n$  cards, then the distribution of  $S_n$  is difficult to compute directly. But, its expected value can be calculated via the gift of linearity.

Let  $X_i$  be a random variable representing the number of boxes needed while trying to get the  $i^{\text{th}}$  new card, starting immediately after the most recent new card’s obtain.

Then,

$$S_n = X_1 + X_2 + \cdots + X_n$$

Therefore,

$$\mathbb{E}[S_n] = \mathbb{E}[X_1 + X_2 + \cdots + X_n] = \mathbb{E}[X_1] + \mathbb{E}[X_2] + \cdots + \mathbb{E}[X_n]$$

The expected value of  $X_1$  would be 1, since any first card is a new card.

Now, how about  $X_i$ ?

Each time we perform a summoning, we will get an new card with probability  $\frac{n-i}{n}$ . We will continue performing summoning, which can be characterized as a coin flip of heads probability  $p = \frac{n-i}{n}$  for simplicity. This makes  $X_i$  the number of tosses until first head.

There,

$$X_i \sim \text{Geometric}\left(\frac{n-i}{n}\right)$$

And consequentially,

$$\mathbb{E}[X_i] = \frac{n}{n-i}$$

Coming back to the original simplification of  $\mathbb{E}[S_n]$ , then:

$$\mathbb{E}[S_n] = \sum_{i=1}^n \frac{n}{n-i} = n \sum_{i=1}^n \frac{1}{i}$$

This can be well approximated as

$$n \sum_{i=1}^n \frac{1}{i} \simeq n(\ln n + \gamma_E)$$

Where  $\gamma_E = 0.5772 \dots$  is known as Euler's constant.

## 18.2 Poisson Distribution

### Definition 18.2.1. Poisson Distribution

A random variable  $X$  for which:

$$\mathbb{P}[X = i] = \frac{\lambda^i}{i!} e^{-\lambda}, i \in \mathbb{Z}^+$$

is said to have the Poisson distribution with parameter  $\lambda$ .

We may mathematically abbreviate such as:

$$X \sim \text{Poisson}(\lambda)$$

This distribution is used to estimate the number of times an event happens (notated by the random variable  $X \sim \text{Poisson}(\lambda)$ ) assuming the average number of occurrence in some periodic unit is  $\lambda$ , where such parameter may be in any division of measurement unit (say clicks occurring per second, but also clicks occurring per minute).

Therefore, while the objective of Poisson is very distinct, it is also a distribution to be carefully utilized along the contexts of a problem. Specifically, depending on the problem's description of length of period (10 seconds? 10 minutes?) on which we model the Poisson distribution, we might need to multiply some given average rate of occurrence with the length of period to obtain the parameter  $\lambda$  (once again, stands for average number of occurrence). Essentially, then,  $\lambda$  is the mean of some distribution!

Let us, once again, verify this distribution by testing whether it sums up to 1:

$$\begin{aligned}\sum_{i=0}^{\infty} \frac{\lambda^i}{i!} e^{-\lambda} &= e^{-\lambda} \sum_{i=0}^{\infty} \frac{\lambda^i}{i!} \\ e^x &= 1 + x + \frac{x^2}{2!} + \dots \\ e^{-\lambda} \sum_{i=0}^{\infty} \frac{\lambda^i}{i!} &= e^{-\lambda} e^{\lambda} = 1\end{aligned}$$

### 18.2.1 Properties of Poisson RV

For both expectation and variance, we can discover a very elegant identity.

#### Theorem 18.2.1. Expectation of Poisson RV

$$\begin{aligned}\mathbb{E}[X] &= \sum_{i=0}^{\infty} i \mathbb{P}[X = i] \\ &= \sum_{i=0}^{\infty} i \frac{\lambda^i}{i!} e^{-\lambda} \\ &= e^{-\lambda} \sum_{i=1}^{\infty} \frac{\lambda^i}{(i-1)!} \\ &= e^{-\lambda} \lambda \sum_{i=1}^{\infty} \frac{\lambda^{i-1}}{(i-1)!} \\ &= e^{-\lambda} \lambda e^{\lambda} = \lambda\end{aligned}$$

And as for variance:

#### Theorem 18.2.2. Variance of Poisson RV

$$\begin{aligned}\mathbb{E}[X(X-1)] &= \sum_{i=0}^{\infty} i(i-1) \mathbb{P}[X = i] \\ &= \sum_{i=0}^{\infty} i(i-1) \frac{\lambda^i}{i!} e^{-\lambda} \\ &= \lambda^2 e^{-\lambda} \sum_{i=2}^{\infty} \frac{\lambda^{i-2}}{(i-2)!} \\ &= \lambda^2 e^{-\lambda} e^{\lambda} = \lambda^2 \\ \text{Var}(X) &= \mathbb{E}[X^2] - (\mathbb{E}[X])^2 \\ &= \mathbb{E}[X(X-1)] + \mathbb{E}[X] - (\mathbb{E}[X])^2 \\ &= \lambda^2 + \lambda - \lambda^2 = \lambda\end{aligned}$$

#### Theorem 18.2.3. Superposition of Poisson Distribution

**Theorem:** Let  $X \sim \text{Poisson}(\lambda)$  and  $Y \sim \text{Poisson}(\mu)$  be independent of each other, then  $X + Y \sim \text{Poisson}(\lambda + \mu)$ .

**Proof:**

$$\begin{aligned}
 \mathbb{P}[X + Y = k] &= \sum_{j=0}^k \mathbb{P}[X = j, Y = k - j] \\
 &= \sum_{j=0}^k \mathbb{P}[X = j] \mathbb{P}[Y = k - j] \\
 &= \sum_{j=0}^k \frac{\lambda^j}{j!} e^{-\lambda} \frac{\mu^{k-j}}{(k-j)!} e^{-\mu} \\
 &= e^{-(\lambda+\mu)} \sum_{j=0}^k \frac{\lambda^j \mu^{k-j}}{j!(k-j)!} \\
 &= e^{-(\lambda+\mu)} \frac{(\lambda + \mu)^k}{k!}
 \end{aligned}$$

### 18.2.2 Binomial Distribution and Poisson Distribution

You may have discovered that binomial distribution and poisson distributions are very similarly used: to estimate the number of times some event occur, provided some metric of occurrence frequency. In fact, this similarity reflects on a mathematical dimension as well:

#### Theorem 18.2.4. Poisson as a Limit of Binomial

**Theorem:** Let  $X \sim \text{Binomial}(n, \frac{\lambda}{n})$  where  $\lambda > 0$  is constant, then

$$\forall i \in \mathbb{N}, \lim_{n \rightarrow \infty} \mathbb{P}[X = i] = \frac{\lambda^i}{i!} e^{-\lambda}$$

**Proof:** First, let's do some algebraic simplification.

$$\begin{aligned}
 \mathbb{P}[X = i] &= \binom{n}{i} \left(\frac{\lambda}{n}\right)^i \left(1 - \frac{\lambda}{n}\right)^{n-i} \\
 &= \frac{\lambda^i}{i!} \cdot \frac{n!}{(n-i)!n^i} \cdot \left(1 - \frac{\lambda}{n}\right)^n \cdot \left(1 - \frac{\lambda}{n}\right)^{-i}
 \end{aligned}$$

Here,

$$\frac{n!}{(n-i)!n^i} = \frac{n(n-1)(n-2) \cdots (n-i+1)}{n^i}$$

This value approaches 1 as  $n$  approaches  $\infty$ .

Meanwhile, from calculus, we may determine that:

$$\lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^n = e^{-\lambda}$$

And last but not least, since  $i$  is a fixed variable,

$$\lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^{-i} = 1$$

Therefore, multiplying these limits together, we obtain that:

$$\lim_{n \rightarrow \infty} \mathbb{P}[X = i] = \frac{\lambda^i}{i!} e^{-\lambda}$$

# Chapter 19

## Prediction!

Prediction? Machine Learning? Uh... I think so?

Next thing you know, I'm reincarnated as a Machine Learning Hypebeast. I estimated the price of your house based on your mom's age, demographics, and your IP address (and the geographic data it encompasses).

However, here's the problem. I overestimated your property. The result that quantifies such an inaccurate property value is known as the **loss**, or, the error.

When dealing with prediction, we attempt to minimize error over each entry of a prediction.

In this case, we will characterize the errors of prediction with Mean Squared Error:

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

This is essentially minimizing a parabola  $Ax^2 + Bx + C$ , from the above equation, where the minimum is at

$$x = \mathbb{E}[X]$$

An important observation is the minimum value of MSE of this estimator is exactly the variance of  $X$ :

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

To summarize, the best MSE estimator for  $X$  is  $\mathbb{E}[X]$ , whose MSE is  $\text{Var}(X)$ .

### 19.1 Joint Distributions, Conditional Expectation

Let's offer a recap of joint distributions:

#### Definition 19.1.1. Joint Distribution

The joint distribution for two discrete random variables  $X$  and  $Y$  is the collection of values:

$$\left\{ \left( (a, b), \mathbb{P}[X = a, Y = b] \right) : a \in \mathcal{A}, b \in \mathcal{B} \right\}$$

where  $\mathcal{A}$  is the set of all possible values taken by  $X$ , and  $\mathcal{B}$  is the set of all possible values taken by  $Y$ .

By such property:

$$\mathbb{P}[X = a] = \sum_{b \in \mathcal{B}} \mathbb{P}[X = a, Y = b]$$

The conditional probability of  $X = x$  given  $Y = y$  would thus also be:

$$\mathbb{P}[X = x|Y = y] = \frac{\mathbb{P}[X = x, Y = y]}{\mathbb{P}[Y = y]}$$

The conditional expectation of  $X$  given  $Y = y$  is defined naturally as follows, where we average along an edge (a margin, an axis) of the distribution's variables:

$$\mathbb{E}[X|Y = y] = \sum_{x \in \mathcal{A}} x \mathbb{P}[X = x|Y = y]$$

## 19.2 Iterated Expectation and Wald's Identity

### Theorem 19.2.1. Law of Iterated Expectations

$$\mathbb{E}[X] = \sum_{y \in \mathcal{B}} \mathbb{E}[X|Y = y] \mathbb{P}[Y = y] = \mathbb{E}[\mathbb{E}[X|Y]]$$

So, let the random variable  $Y = X_1 + \dots + X_N$  be some sum of other random variables that are IID, and  $N$  (the number of IID RVs) is also itself a random variable.

Then, to compute the expectation of  $Y = \sum_{i=1}^N X_i$  where each of  $X_i$  is independent to  $N$ , we derive the following identity:

### Theorem 19.2.2. Wald's Identity

$$\begin{aligned} \mathbb{E}[Y] &= \mathbb{E}[\mathbb{E}[Y|N]] \\ &= \sum_n \mathbb{E}[Y|N = n] \mathbb{P}[N = n] \\ &= \sum_n n \mathbb{E}[X_1] \mathbb{P}[N = n] \\ &= \mathbb{E}[X_1] \sum_n n \mathbb{P}[N = n] = \mathbb{E}[X_1] \mathbb{E}[N] \end{aligned}$$

This can be useful for modelling the total time to serve customers in a time interval, where we model arrivals with Poisson distribution.

We can thus have a Poisson random variable,  $N \sim \text{Poisson}(\lambda)$ , that determines the number of customers and each  $X_i$  is an indicator variable that corresponds to time needed to serve customer  $i$ .

## 19.3 Minimal Mean Square Error?

What estimate should we use for a random variable  $Y$  to minimize mean squared error if one does not have a predictor variable's value  $X$  in hand?

In the context of identifying an estimator  $y$  upon a known predictor variable  $x$ , then we know that the estimator for  $Y$ ,  $y$ , minimizes:

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

We can therefore iterate this computation through all possible  $x$ , where the value is minimized. This means we take the expected value over a joint distribution.

Therefore, we can follow the above reasoning to acquire some estimator that minimizes the mean squared error throughout the entire sample space:

$$\mathbb{E}_{W,H}[(W - \mathbb{E}[W|H])^2]$$

## 19.4 Linear Regression

Let's say we will try to predict  $y$  with a linear function  $\mathcal{L}(x)$ , characterized by some slope  $m$  and intercept  $b$  of a linear equation.

The mean squared error, in this case,  $\mathbb{E}[(Y - \mathcal{L}(x))^2]$ , should be minimized.

Let us first consider the random variables:

$$\bar{X} = X - \mathbb{E}[X], \bar{Y} = Y - \mathbb{E}[Y]$$

In this case, we will note one value of the random variable  $\bar{X}$  as  $\bar{x}$ .

To estimate  $\bar{Y}$ , which is zero-mean now for convenience, with  $x$ , we will consider a function:

$$\bar{f}(x) = m\bar{x} + b$$

Once again, consider that we are attempting to minimize with  $b = 0$ :

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

The quadratic function structure of above equation suggests that, to minimize this error, the slope would be:

$$m = -\frac{-2\mathbb{E}[\bar{X}\bar{Y}]}{2\mathbb{E}[\bar{X}^2]} = \frac{\mathbb{E}[\bar{X}\bar{Y}]}{\mathbb{E}[\bar{X}^2]}$$

However,  $b \neq 0$  yields a more interesting case:

$$\mathbb{E}[\bar{Y}]^2 - 2(m\mathbb{E}[\bar{X}\bar{Y}] + b\mathbb{E}[\bar{Y}]) + (m^2\mathbb{E}[\bar{X}]^2 + 2mb\mathbb{E}[\bar{X}] + b^2)$$

Which thus lead to:

$$2\mathbb{E}[\bar{Y}] + 2m\mathbb{E}[\bar{X}] + 2b = 0$$

Where we see that, since the random variables are purposely zero-meaned,  $b = 0$ .

So indeed  $b = 0$ , which in turn provides a linear function:

$$\mathcal{L}(X) = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}(X - \mathbb{E}[X]) + \mathbb{E}[Y]$$

Let us expand this too:

$$\mathcal{L}(X) = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}X + \mathbb{E}[Y] - \frac{\text{Cov}(X, Y)}{\text{Var}(X)}\mathbb{E}[X]$$

This becomes the optimal regression line equation, as outlined in DATA C100!

Furthermore, simplifying our mean squared error applying the solution that  $b = 0$ , we can see the following derivation process:

$$\begin{aligned} & \mathbb{E}[\bar{Y}^2] - 2m\mathbb{E}[\bar{X}\bar{Y}] + m^2\mathbb{E}[\bar{X}^2] \\ & \mathbb{E}[\bar{Y}^2] - \frac{\text{Cov}(\bar{X}, \bar{Y})^2}{\text{Var}(\bar{X})} \\ & 1 - (\text{Corr}(\bar{X}, \bar{Y}))^2 \end{aligned}$$



## Chapter 20

# Continuous Probability Distributions

Up until now we have focused on discrete sample spaces  $\Omega$ , where there are finite or countably infinite sample points to deal with.

Discrete, as in the sample points themselves are discretely identifiable. However, discrete random variables cannot model everything in the world. Many real life quantities that we wish to model probabilistically are in fact real-valued. For example, position of particle, as well as direction of travel. To aid our future efforts for probability-related studies, let's study continuous random variables.

### 20.1 Continuous Random Variables

Suppose we spin a wheel of fortune and record the position at which some needle points to, in relation to the position along a circumference of length  $l$  (which means the position of needle on wheel of fortune is within the range  $[0, l]$ ). In such case, that position is continuous. There are too many possible points to discretize this situation. But, to resolve this situation, we may as well consider the following scheme:

Consider the probability at which needle falls between positions  $[a, b]$  on all possible positions  $[0, l]$ . Then, the probability of this situation would be:

$$\mathbb{P}[x \in [a, b]] = \frac{b - a}{l}$$

This above probability assumes that each position has a uniform probability of having the needle fall on it. If there are some positions whose probability of being pointed to is higher, then this entire derivation for probability would've been different.

This is how we may relate the concept of event from a discrete point of view to a continuous random variable situation: via its analogy as one section of circumference.

Then, in a continuous concept, random variables  $X$  no longer be probabilistically assessed as “if  $X$  is equal to some value”, since if  $X$  is equal to some value on a continuous space, it would require us to calculate the probability that  $X$  locates on exactly one point (probability that needle stops exactly on some point of the wheel), which is realistically  $\frac{1}{\infty}$ .

Rather, a meaningful question for continuous random variables would be: what is the probability that  $X$  is between some values  $a, b$ ? And to answer this question, we would need to know the distribution of probability across each values between  $a$  and  $b$ :

### Definition 20.1.1. Probability Density Function

A probability density function for a real-valued random variable  $X$  is a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  satisfying the following requirements (of probability distributions!):

- $f$  is non-negative:  $\forall x \in \mathbb{R}, f(x) \geq 0$
- The total integral of  $f$  (Area under curve of distribution) is 1. (All probabilities sum up to 1).

Then the distribution of  $X$  may be provided as:

$$\mathbb{P}[a \leq X \leq b] = \int_a^b f(x)dx$$

for all  $a < b$ .

Which, in other words, must satisfy that:

$$\mathbb{P}[-\infty \leq X \leq \infty] = \int_{-\infty}^{\infty} f(x)dx = 1$$

In this case, the wheel of fortune analogy may receive a piecewise function for its PDF:

$$f(x) = \begin{cases} 0, & \text{for } x < 0 \\ \frac{1}{l}, & \text{for } 0 \leq x \leq l \\ 0, & \text{for } x > l \end{cases}$$

**But remember that PDF itself is not probability.** For example, the above function doesn't necessarily sum up to 1; its integral does, and that is the only thing that matters.

**Its distribution is what exhibits and determines the feature of random variables that we call “probability”.**

## 20.1.1 CDF

For some continuous RV  $X$ , we can discuss the cumulative distribution function of it, which is by nature just some shorthand expression defined as follows:

$$F(x) = \mathbb{P}[x \leq X] = \int_{-\infty}^x f(z)dz$$

But this grants us the opportunity to recognize a probability density function from cumulative distribution function via calculus:

$$f(x) = \frac{d}{dx}F(x)$$

## 20.1.2 Expectation and Variance

### Definition 20.1.2. Expectation of Continuous RV

The expectation of a continuous random variable  $X$  with a PDF  $f$  would then be:

$$\mathbb{E}[X] = \int_a^b xf(x)dx$$

And in turn,

**Definition 20.1.3. Variance of Continuous RV**

The variance of a continuous random variable  $X$  with a PDF  $f$  would then be:

$$\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2] = \int_a^b x^2 f(x) dx - \left( \int_a^b x f(x) dx \right)^2$$

**20.1.3 Joint Distribution**

For some joint distribution considering plural random variables at the same time, the density function would have some slight modifications for rules:

**Definition 20.1.4. Joint Density**

A joint density function for two random variables  $X$  and  $Y$  is a function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  satisfying:

- $f$  is nonnegative (as always, negative density doesn't make sense).
- The total integral of  $f$  is 1.

The consequence of that second rule above would be:

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1$$

So on a distribution's sense:

$$\mathbb{P}[a \leq X \leq b, c \leq Y \leq d] = \int_a^b \int_c^d f(x, y) dx dy$$

Once again, by Fubini's Theorem, the order of  $dx$  and  $dy$  wouldn't matter as long as the bounds of integrations are altered accordingly and still correspond to their variables.

**20.1.4 Independence of Continuous RVs**

For two continuous random variables, their independence is demonstrated by the following rule:

**Definition 20.1.5. Independence of Continuous Random Variables**

Two continuous random variables  $X, Y$  are independent if the events  $a \leq X \leq b$  and  $c \leq Y \leq d$  are independent for all  $a \leq b$  and  $c \leq d$ , such that:

$$\mathbb{P}[a \leq X \leq b, c \leq Y \leq d] = \mathbb{P}[a \leq X \leq b] \mathbb{P}[c \leq Y \leq d]$$

The definition has an implication on the joint density of independent continuous random variables:

$$\begin{aligned} f(x, y) dx dy &\simeq \mathbb{P}[x \leq X \leq x + dx, y \leq Y \leq y + dy] \\ &= \mathbb{P}[x \leq X \leq x + dx] \cdot \mathbb{P}[y \leq Y \leq y + dy] \\ &\simeq f_X(x) dx \times f_Y(y) dy \\ &= f_X(x) f_Y(y) dx dy \end{aligned}$$

Essentially, gives birth to the following theorem:

**Theorem 20.1.1. Joint Density of Independent Continuous RVs**

The joint density of independent RVs  $X$  and  $Y$  is the product of marginal densities:

$$\forall x, y \in \mathbb{R}, f(x, y) = f_X(x)f_Y(y)$$

where marginal densities  $f_N$  are the densities of some specific random variable  $N$  in this joint distribution.

## 20.2 Exponential Distribution

Let us first define this distribution:

**Definition 20.2.1. Exponential Distribution**

For  $\lambda > 0$ , a continuous random variable  $X$  with pdf:

$$f(x) = \begin{cases} \lambda e^{-\lambda x}, & \text{if } x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

is an exponential random variable with parameter  $\lambda$ , which we mathematically express as:  $X \sim \text{Exp}(\lambda)$

This distribution is nonnegative, and sums up to 1:

$$\int_{-\infty}^{\infty} f(x) dx = \int_0^{\infty} \lambda e^{-\lambda x} dx = -e^{-\lambda x} \Big|_0^{\infty} = 1$$

### 20.2.1 Analyses of Exponential Distribution

The expected value can be computed using integration by parts:

**Definition 20.2.2. Expectation of Exponential Distribution**

$$\begin{aligned} \mathbb{E}[X] &= \int_{-\infty}^{\infty} x f(x) dx \\ &= \int_0^{\infty} \lambda x e^{-\lambda x} dx \\ &= -x e^{-\lambda x} \Big|_0^{\infty} + \int_0^{\infty} e^{-\lambda x} dx \\ &= -\frac{e^{-\lambda x}}{\lambda} \Big|_0^{\infty} = \frac{1}{\lambda} \end{aligned}$$

By that logic, the expected value of  $\mathbb{E}[X^2]$  can be evaluated via similar tricks, such that:

$$\mathbb{E}[X^2] = \frac{2}{\lambda^2}$$

The above analyses leads to the conclusion that:

$$\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \frac{2}{\lambda^2} - \frac{1}{\lambda^2} = \frac{1}{\lambda^2}$$

Furthermore, exponential distribution serves as an analog of geometric distribution.

For example, the exponential distribution would describe that:

$$\mathbb{P}[X > t] = \int_t^{\infty} \lambda e^{-\lambda x} dx = -e^{-\lambda x} \Big|_t^{\infty} = e^{-\lambda t}$$

Meanwhile, let us consider that  $\lambda$  can be a fixed rate of success per unit time, such that in a geometric distribution, the rate of success per time is modeled as  $p = \lambda \delta$  where  $\delta$  resembles a small amount of time (small for discrete, near-zero for continuous).

Then, mathematical conclusions point to the fact that geometric distributions  $Geo(\lambda \delta)$  is about equal to the exponential distribution's regional probability marked as  $e^{-\lambda x}$ .

## 20.3 Normal Distribution

Another continuous distribution prevalent in many applications is known as normal, or Gauss distribution.

This distribution uses two parameters,  $\mu$  and  $\sigma^2$ , which respectively are the mean (expectation) and variance of the distribution.

And,

### Definition 20.3.1. Normal Distribution

For any  $\mu \in \mathbb{R}$ ,  $\sigma > 0$ , a continuous random variable  $X$  with PDF:

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

is called a normal random variable with parameters  $\mu$  and  $\sigma^2$ , and mathematically expressed:

$$X \sim N(\mu, \sigma^2)$$

When  $\mu = 0$ ,  $\sigma = 1$ , such random variable is particularly noted as “standard normal distribution”

Here is an interesting manipulation that can prove helpful:

### Theorem 20.3.1. Manipulations towards Standard Normal Distribution

**Theorem:** If  $X \sim N(\mu, \sigma^2)$ , then  $Y = \frac{X-\mu}{\sigma} \sim N(0, 1)$

**Proof:**

$$\mathbb{P}[a \leq Y \leq b] = \mathbb{P}[\sigma a + \mu \leq X \leq \sigma b + \mu] = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{\sigma a + \mu}^{\sigma b + \mu} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

Perform a change of variable here, where  $Y = \frac{X-\mu}{\sigma}$ . Then,

$$\frac{1}{\sqrt{2\pi\sigma^2}} \int_{\sigma a + \mu}^{\sigma b + \mu} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = \frac{1}{\sqrt{2\pi\sigma^2}} \int_a^b e^{-\frac{y^2}{2}} dy$$

And as an extra offering from this note, let's discuss the area under the bell curve.  $f(x) = e^{-x^2}$ , which would rationalize why  $\pi$  becomes involved in a normal distribution.

### Exercise Reader! Exercise! 20.3.1: Area under Bell Curve

Let us first define a variable:

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

However, this expression itself is very difficult to derive for. Let's perform some math hack:

$$\begin{aligned} a &= \int_{-\infty}^{\infty} e^{-x^2} dx = \int_{-\infty}^{\infty} e^{-y^2} dy \\ a^2 &= \int_{-\infty}^{\infty} e^{-x^2} dx \int_{-\infty}^{\infty} e^{-y^2} dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-y^2} e^{-x^2} dy dx = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-(x^2+y^2)} dy dx \end{aligned}$$

Seeing the expression involving some variable  $x^2 + y^2$ , we can attempt some change of coordinate, such that:

$$\int \int_R f(x,y) dy dx = \int \int_R f(r \cos(\theta), r \sin(\theta)) r d\theta dr$$

Therefore,

$$\begin{aligned} a^2 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-(x^2+y^2)} dy dx \\ &= \int_{-\infty}^{\infty} \int_0^{2\pi} e^{-r^2} r d\theta dr \\ &= 2\pi \left. \frac{-e^{-r^2}}{2} \right|_0^{\infty} = \pi \end{aligned}$$

Consequently,

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

### 20.3.1 Analyses of Normal Distribution

Let us prove that the expectation and variance of a normal distribution is indeed as noted by its parameters.

#### Theorem 20.3.2. Expectation of Normal Distribution

We will be exploiting the linearity of expectation values here, by first computing the expectation of  $N(0, 1)$ , then performing a change of variable as noted in the previous theorem.

So, first, for  $X \sim N(0, 1)$ :

$$\begin{aligned} \mathbb{E}[X] &= \int_{-\infty}^{\infty} x f(x) dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x e^{-\frac{x^2}{2}} dx \\ &= \frac{1}{\sqrt{2\pi}} \left( \int_{-\infty}^0 x e^{-\frac{x^2}{2}} dx + \int_0^{\infty} x e^{-\frac{x^2}{2}} dx \right) = 0 \end{aligned}$$

Performing the change of variable:

$$X = \frac{Y - \mu}{\sigma}$$

Such that  $X$  is the standard version of  $Y$ , we may then deduce that:

$$\mathbb{E}[X] = \frac{\mathbb{E}[Y] - \mu}{\sigma} = 0$$

Meaning that:

$$\mathbb{E}[Y] = \sigma \mathbb{E}[X] + \mu = \mu$$

**Theorem 20.3.3. Variance of Normal Distribution**

The variance of standard normal distribution would then also be computable via intergration by parts:

$$\begin{aligned} \text{Var}(X) &= \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^2 e^{-\frac{x^2}{2}} dx \\ &= \frac{1}{\sqrt{2\pi}} \left( -xe^{-\frac{x^2}{2}} \right) \Big|_{-\infty}^{\infty} + \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2}} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2}} dx = 1 \end{aligned}$$

The last line of that derivation was an interesting extra content offered from my semester of MATH 53. Essentially, it relies on a change of variable into the polar coordinates.

Using the change of variable between  $X$  and  $Y$ , then:

$$\text{Var}(X) = \text{Var}\left(\frac{Y - \mu}{\sigma}\right) = 1$$

Therefore:

$$\text{Var}(Y - \mu) = \text{Var}(Y) = \sigma^2$$

Let us discuss the superpositioning properties for independent normal random variables then. First of all:

**Theorem 20.3.4. Sum of Independent Standard Normal Random Variables**

**Theorem:** Let  $X \sim N(0, 1)$  and  $Y \sim N(0, 1)$  be independent standard normal random variables, and suppose  $a, b \in \mathbb{R}$  are constants, then  $Z = aX + bY \sim N(0, a^2 + b^2)$ .

**Proof:**

To begin with,  $X$  and  $Y$  are independent, so the joint density of it can be found as:

$$f(x, y) = f(x)f(y) = \frac{1}{2\pi} e^{-\frac{x^2+y^2}{2}}$$

Such is rotationally symmetric around the origin, where  $f(x, y)$  depends on the distance a point is from the origin.

And upon that, recognize the fact that any rotation of this density function about the origin along the vertical axis would not change how the distribution is, as it is rotationally symmetrical.

Therefore, given some value  $t \in \mathbb{R}$ , we can first recognize that:

$$\mathbb{P}[Z \leq t] = \mathbb{P}[aX + bY \leq t] = \mathbb{P}[(X, Y) \in \text{halfplane containing all points } ax + by \leq t]$$

Meanwhile, notice that this halfplane, now with a diagonal boundary, can be rotated to have a vertical boundary.

Let's call this rotation:

$$T(A) = \{(x, y) | x \leq \frac{t}{\sqrt{a^2 + b^2}}\}$$

Which thus indicates by the unchanging property of density function's rotational symmetry along z-axis, that:

$$\mathbb{P}[Z \leq t] = \mathbb{P}[(X, Y) \in A] = \mathbb{P}[(X, Y) \in T(A)] = \mathbb{P}[X \leq \frac{t}{\sqrt{a^2 + b^2}}] = \mathbb{P}[X\sqrt{a^2 + b^2} \leq t]$$

Here, we will conclude that  $\sqrt{a^2 + b^2}X$  and  $Z$  are identically distributed, for the above equation holds on any value of  $t$ .

The normal distribution  $\sqrt{a^2 + b^2}X$ , a rotation of  $Z$ , has normal distribution with mean 0 and variance  $a^2 + b^2$ , by the work of Theorem 3.1 in this chapter.

Therefore, the normal distribution  $Z = aX + bY$  should have the same parameters. Consequentially,

$$Z = aX + bY \sim N(0, a^2 + b^2)$$

The theorem now guides us to a more general statement:

**Theorem 20.3.5. Addition of Normal Random Variables**

**Theorem:** Let  $X \sim N(\mu_X, \sigma_X^2)$  and  $Y \sim N(\mu_Y, \sigma_Y^2)$  be independent normal random variables. Then, the random variable  $Z = aX + bY$  for real coefficients  $a, b$  is also normally distributed as  $Z \sim N(a\mu_X + b\mu_Y, a^2\sigma_X^2 + b^2\sigma_Y^2)$ .

**Proof:**

Generate standard normal random variables  $Z_1 = \frac{X - \mu_X}{\sigma_X}$  and  $Z_2 = \frac{Y - \mu_Y}{\sigma_Y}$ .

Then, we may use the previous theorem to perform the following derivation:

$$\begin{aligned} Z = aX + bY &= a(\mu_X + \sigma_X Z_1) + b(\mu_Y + \sigma_Y Z_2) \\ &= (a\mu_X + b\mu_Y) + (a\sigma_X Z_1 + b\sigma_Y Z_2) \end{aligned}$$

We know the distribution by the previous theorem that:

$$(a\sigma_X Z_1 + b\sigma_Y Z_2) \sim N(0, (a\sigma_X)^2 + (b\sigma_Y)^2)$$

Therefore, as we add the non-random variable terms of  $\mu_X$  and  $\mu_Y$ , we only shift the distribution by that amount horizontally and finally attain:

$$Z \sim N(a\mu_X + b\mu_Y, a^2\sigma_X^2 + b^2\sigma_Y^2)$$

## 20.4 Central Limit Theorem

**Theorem 20.4.1. Central Limit Theorem**

Let  $X_1, \dots, X_n$  be a sequence of IID random variables with common finite expectation  $\mathbb{E}[X_i] = \mu$  and finite variances  $\text{Var}(X_i) = \sigma_i^2$ .  $S_n$  is the sum  $\sum_{i=1}^n X_i$ .

Then the distribution of  $\frac{S_n - n\mu}{\sigma\sqrt{n}}$  converges to  $N(0, 1)$  as  $n \rightarrow \infty$ .

Therefore, as  $n \rightarrow \infty$ :

$$\mathbb{P}\left[\frac{S_n - n\mu}{\sigma\sqrt{n}} \leq c\right] \rightarrow \frac{1}{\sqrt{2\pi}} \int_{-\infty}^c e^{-\frac{x^2}{2}} dx$$

Which would in fact be the distribution of sample mean.

Then, the Central Limit Theorem mathematically promise such that the above distribution converges onto a standard normal distribution.

By that context, the distribution of observed average of  $n$  observations on some random experiment will be a normal distribution centered at the true value (population mean in statistical context) of that random variable's expectation ( $\mu$ ), with standard deviation  $\frac{\sigma}{\sqrt{n}}$ .