

---

# Discrete Probability

---

**Brandon Smith**

## Abstract

Life is filled with many uncertainties and Probability Theory allows us to effectively think and deal with it. In this article, we provide an introduction to the field of Probability, in particular, Discrete Probability. We focus on three specific learning objectives where, after reading, one should be able to explain and define probability models and axioms, such as sample space and events, compute the conditional probability directly using Bayes' rule on real-world problems, recognise and apply discrete distributions to find the probability of various events.

Probability theory is nothing but  
common sense reduced to  
calculation.

---

*Laplace, Pierre Simon (1812)*

## 1 Introduction

Whether we have formally studied Probability Theory or not, we all have in some way, engaged with it unconsciously. Suppose you are sitting in your room, looking into the cloudy grey sky: "Looks like a storm is coming" you thought.

How did you come to this prediction? It most likely, almost certainly, from past experiences which formed a basis to infer the likelihood of a storm. This very thought process, at least conceptually, is the idea of Probability Theory. To apply this in a more robust, structured manner, one must distil these observations into systematic calculations.

Various fields utilise Probability Theory to handle the randomness in this world, ever more so today than ever with the boom of big data. While each field adapts these concepts to their particular needs, they all originate from the basic principles this paper looks into today. It's my hope that these concepts will not only be useful but also inspire to dive deeper into this fascinating subject.

## 2 Preliminaries

Throughout this module, we'll be using several concepts which are not explained here, predominantly relating to Set Theory.

A set is a collection of distinct items, known as elements, while an empty set, is a set without any elements, which is denoted as  $\emptyset$ . Sets can be expressed in various ways. Let's say we have a finite set  $S$  that contains the elements  $a_1, a_2$  up to  $a_n$ . We could list these elements within curly braces like so:

$$S = a_1, a_2, \dots, a_n$$

Subsets on the other hand, which means if every element in one set (we'll call this  $E$ ) is also in another set (we'll call this  $S$ ), then we say  $E$  is a subset of  $S$  and is written as  $E \subseteq S$ . In other words, a subset is a set composed entirely of elements within another set.

Let's say we have a set  $S$  which contains the objects  $A, B$  and  $C$ :

$$S = A, B, C$$

In this case, sets like  $\{A\}$ ,  $\{B\}$ ,  $\{C\}$ ,  $\{A, B\}$ ,  $\{A, C\}$ ,  $\{B, C\}$  and  $\{A, B, C\}$  are all subsets of  $S$ . The empty set  $\emptyset$  is considered a subset of every set. This idea of subsets forms the backbone for many discussions in probability theory, especially when we're talking about event spaces within a particular sample space.

For further reading, I recommend reading *Naïve Set Theory* by Paul R. Halmos. It is a friendly and fun-to-read introduction to Set Theory. We will also be using Venn Diagrams to illustrate some ideas of Probability Theory; I have provided a resource here for which you can use to get yourself familiar.

## 3 Probability Models

Probability models serve as mathematical representations of uncertain situations or experiments, helping us navigate randomness which affect our very lives. They can be thought of as blueprints that guide us through uncertainty and allow us to make informed decisions.

The creation of a probability model involves a few important steps. Firstly, we identify all possible outcomes of an experiment or situation - we establish what's known as a sample space.

Once we've defined our sample space, we establish a probability law. This assigns a likelihood to each individual outcome or groups of outcomes within our sample space, enabling us to understand which results are more likely to transpire. However, these assigned probabilities need to adhere to a set of rules or principles, ensuring logical coherence within our model. These principles form the foundational pillars of probability theory, known as the axioms. While only a few, these axioms serve as the bedrock upon which we construct many important insights and theorems.

In the next few points in this report, we will provide a detailed exploration of these concepts, breaking down their mechanisms and highlighting their significance. We will also shine a light on the practical applications of these probability models, demonstrating their substantial impact in real-world situations.

---

### 3.1 Probability Functions

Probability functions works by assigning probabilities to events, or specific sets of outcomes from a random experiment. The function abides by Kolmogorov's axioms – which we will discuss in detail later.

1. The domain of the function is the event space, a collection of subsets from the sample space.
2. The co-domain of a probability function falls between 0 and 1, inclusive, reflecting the probability of an event's occurrence. It's important to note that this will differ from the functions we'll explore in later discussions on discrete distributions.

Let's assume the probability function is denoted by  $P$ , with the domain being the event space  $E$  and the co-domain falling between 0 and 1 inclusive. This can be represented as:

$$P : E \rightarrow [0, 1]$$

### 3.2 Sample Space and Events

Let us imagine that we're conducting an experiment with an unpredictable outcome. We can't foresee the exact result, but we know all the potential outcomes that could occur. This collection of all conceivable results is what we call the sample space of the experiment and in this module, we will denote it as  $S$ . However, you may find that it is can also be denoted as either  $\Omega$ , or  $U$ .

#### Sample Space

We begin with the idea of a sample space. This is the full set of potential outcomes for a given scenario or experiment. For example, let's consider an experiment where a coin is tossed twice and we record the face that shows after each toss. We can represent this as  $H$  for heads and  $T$  for tails. The sample space of this experiment, denoted  $S$ , would be:

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

In more general terms, an experiment can be thought of as a procedure or action that could be repeated under the same conditions, leading to uncertain outcomes.

Take rolling a pair of dice or flipping a coin or drawing a card from a deck as an example. In each of these experiments, the specific outcomes remain unknown until the experiment is conducted and the collective set of all possible outcomes constitutes the sample space.

#### Event Space

An event is a specific outcome or a collection of outcomes from the sample space, also known as the subset of the sample space. In the coin-tossing example, an event could be getting a head followed by a tail, or a tail followed by a head. This can be denoted like so:

$$E = \{HT, TH\}$$

This event represents the proposition that when the coin is tossed twice, it will show one head and one tail. Events are collections of individual sample points that make up some interesting statement or proposition in the experiment.

## 1. Discrete Probability

---

In Kolmogorov's second remark from Foundations of the theory of probability [1], he states that if an event has no potential outcomes within the sample space, we can safely consider it "impossible" and hence it has a probability of zero.

An example of such an event would be flipping a fair coin and landing on both heads and tails at the same time — this is physically impossible. However, having zero probability doesn't necessarily make an event "impossible". A zero-probability event may seem practically impossible within a single instance or a finite event space, but when considering infinite trials or an infinite event space, such an event can indeed occur. This points to the event's extreme "improbability", not its absolute "impossibility".

### Example:

A survey is conducted among students in SIT192 at Deakin University to determine their pizza preferences. Specifically, students are queried about their liking for Meat Lovers, Hawaiian, both types, or neither.

The survey receives responses from a total of 100 students, leading to the following:

- Total number of students who responded:  $N = 100$
- Number preferring Meat Lovers only:  $M = 35$
- Number preferring Hawaiian only:  $H = 20$
- Number preferring both Meat Lovers and Hawaiian:  $M \cap H = 15$
- Number not preferring either Meat Lovers or Hawaiian:  $N - (M \cup H) = 30$

From these results, the sample space  $S$  representing all possible outcomes of the students is defined as:

$$S = \{\text{Meat Lovers only, Hawaiian only, Both, Neither}\}$$

This way, we can categorise the preferences of the entire sample space into distinct events, providing us a more clear and systematic understanding of the students' pizza preferences.

### Definition 1.0.1: Sample Space & Events

Consider  $S$  to be the sample space for a given experiment. Any subset  $E$  within  $S$ , which can range from the empty set to the entire sample space  $S$ , is defined as an event.

---

### 3.3 Operations on Events

Given an experiment and its associated sample space, we can perform several operations on events, which are subsets of the sample space.

Using Definition 1.0.1, it defines an event as any subset within the sample space, which could range from the empty set to the entire sample space itself. From here, we can dive into several operations on these events within the sample space.

#### Definition 1.0.2: Union of Events

Denoted as  $A \cup B$ , this represents the occurrence of either event  $A$ , event  $B$ , or both.

The union of  $A$  and  $B$  defined in 1.0.2 and visually shown in figure 1.1, is denoted as  $A \cup B$ , is the event representing all students who prefer either Meat Lovers, Hawaiian, or both. This group includes students who enjoy either pizza individually or also those who like both, covering all possibilities related to these two pizza types.

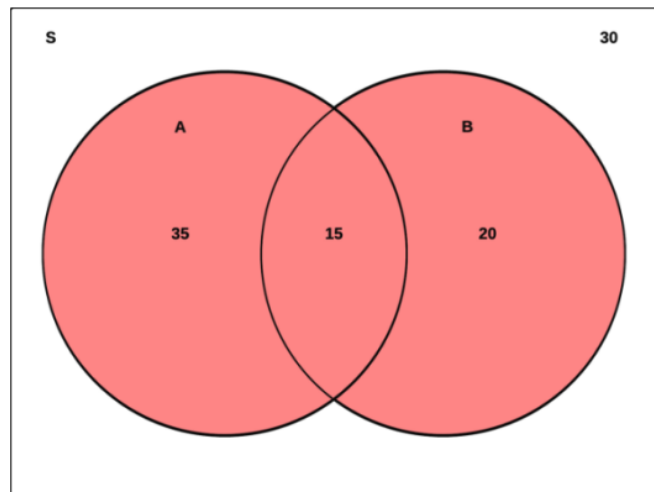


Figure 1.1: Union of  $A$  and  $B$  inside the sample space  $S$ .

#### Definition 1.0.3: Intersection of Events

Represented as  $A \cap B$ , this signifies that events  $A$  and  $B$  occur simultaneously.

Using definition 1.0.3, the intersection of events  $A$  and  $B$  visually shown in ,  $A \cap B$ , represents students who enjoy both Meat Lovers and Hawaiian pizzas.

We can see that in figure 1.2, aside from the 15 students who prefer both types of pizzas (which is  $A \cap B$ ), 35 students prefer Meat Lovers only and 20 students prefer Hawaiian only.

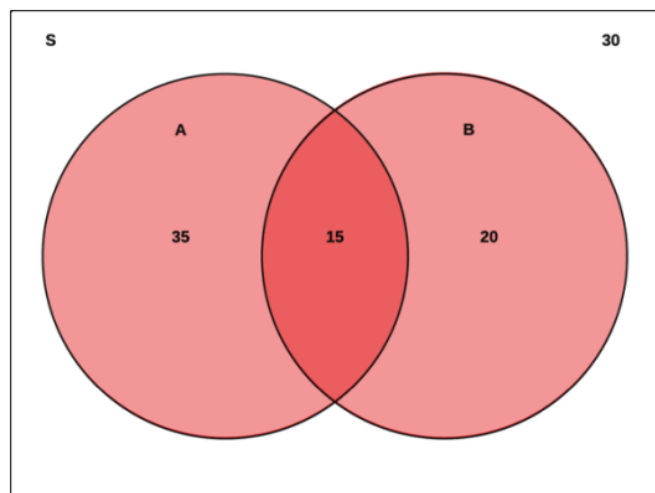


Figure 1.2: Intersection of  $A$  and  $B$  inside the sample space  $S$ .

The area outside both events  $A$  and  $B$  but within the sample space  $S$  represents the 30 students who do not prefer either Meat Lovers or Hawaiian pizza.

**Definition 1.0.4: Complement of an Event**

Denoted as  $A'$  or  $\neg A$ , this signifies the non-occurrence of event  $A$ .

Figure 1.3 illustrates the concept of complement in the context of our pizza survey. Within the sample space  $S$ , event  $A$  represents students who prefer Meat Lovers pizza, indicated by the unshaded circle labelled  $A$ .

The shaded region represents  $A'$ , the complement of  $A$ , which includes all students who do not prefer Meat Lovers pizza.

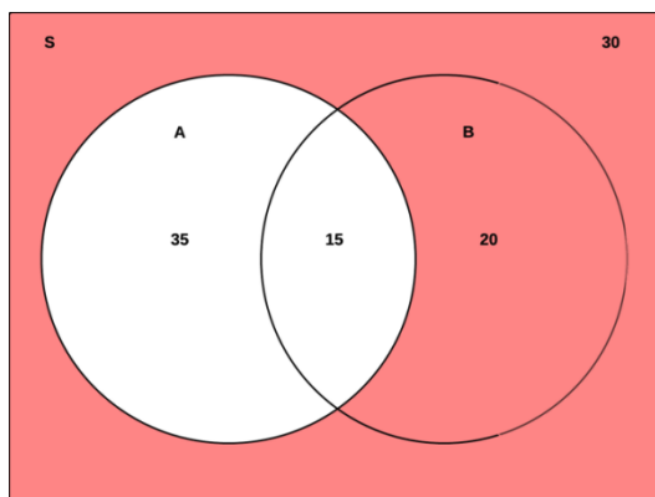


Figure 1.3: The complement of  $A$  ( $\neg A$ ) inside the sample space  $S$ .

---

### 3.4 Product and Sum rule

This section is dedicated to defining the product and sum rules, which are widely accepted mathematical principles that form the basis for the subsequent sections of this module.

#### Definition 1.0.5: Product Rule

For independent events  $A$  and  $B$ , the product rule determines the probability of both events happening concurrently, denoted as  $P(A \cap B)$ . It is defined as:

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

#### Definition 1.0.6: Sum Rule

The sum rule calculates the probability of either or both events  $A$  and  $B$  occurring, denoted as  $P(A \cup B)$ , as:

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

It corrects for the double counting of  $A$  and  $B$  occurring together in the individual probabilities  $P(A)$  and  $P(B)$ .

### 3.5 Kolmogorov Axioms

In this section, we will introduce Kolmogorov's axioms. The axioms described by Kolmogorov in his publication *Foundations of the Theory of Probability* [1] are as follows:

"The theory of probability, as a mathematical discipline, can and should be developed from axioms in exactly the same way as Geometry and Algebra."

What this means is that, after we've identified the basic elements and their relationships and laid out the rules (axioms) guiding these relationships, everything that follows should be based strictly on these rules. This approach suggests that we should set aside real-world interpretations for a bit and focus on the logic and reasoning that these rules provide.

But why do we take such an approach? The purpose is not merely abstraction for abstraction's sake but to strive for generality in our understanding and application of the theory. This point resonates with Kolmogorov's assertion that an abstract theory admits an unlimited number of concrete interpretations beyond those from which it was originally derived.

This means that by focusing on the abstract structure and logic provided by the axioms, we are creating a framework that can be adapted to a wide list of different scenarios and disciplines.

Set Theory	Random Events
1. $A$ and $B$ are disjoint sets, i.e., $A \cap B = \emptyset$ .	1. Events $A$ and $B$ are mutually exclusive.
2. The intersection of sets $A, B, \dots, N$ is represented as $AB \dots N$ and if $AB \dots N = \emptyset$ , it means that sets $A, B, \dots, N$ are mutually exclusive.	2. If events $A, B, \dots, N$ are mutually exclusive, it means that they cannot occur simultaneously.
3. If $AB \dots N = X$ , it signifies that $X$ represents the intersection of sets $A, B, \dots, N$ .	3. Event $X$ happens if and only if events $A, B, \dots, N$ all occur simultaneously or concurrently.
4. $A + B + \dots + N = X$ means that $X$ represents the union or combination of $A, B, \dots, N$ . It indicates that $X$ occurs if at least one of $A, B, \dots, N$ occurs.	4. Event $X$ is defined as the occurrence of at least one of the events $A, B, \dots, N$ . It means that $X$ happens if any of the events occur.
5. The complementary set of $A$ , denoted as $\bar{A}$ . Which means the opposite of $A$ .	5. The event $A'$ , representing the non-occurrence of event $A$ .
6. $A = \emptyset$ .	6. Event $A$ is impossible within the given sample space.
7. $A = E$	7. Event $A$ represents the entire sample space $E$ , which indicates that all outcomes in the sample space are certain to occur.
8. $B$ is a subset of $A$ (which means included in $A$ ): $B \subseteq A$ .	8. If event $B$ happens, it guarantees that event $A$ will happen without a doubt.



---

### The Axioms

**Axiom 1:** The probability of any event  $E$  within the sample space  $S$  can never be negative.

$$P(E) \geq 0 \quad \forall E \in S$$

**Axiom 2:** The probability of at least one event occurring within the sample space is certain, hence equal to 1.

$$P(S) = 1$$

**Axiom 3:** If events are mutually exclusive, then the probability of the union of these events equals the sum of their individual probabilities.

$$P(A_1 \cup A_2 \cup A_3 \cup \dots) = P(A_1) + P(A_2) + P(A_3) + \dots$$

Kolmogorov's set theory terminology gives us some essential definitions for understanding his axioms. In particular, he uses the notation  $A = \emptyset$  to indicate that event A is impossible within the given sample space. This means that the probability of an impossible event is  $P(\emptyset) = 0$ .

However, this fact is not explicitly mentioned in Kolmogorov's axioms (and it doesn't need to be). That is because it is a consequence that can be derived directly from the axioms.

#### Proof 1: The Axioms

Let's consider the empty set  $\emptyset$ , an event that represents impossibility as it contains no outcomes. Even though it's not explicitly stated, the probability of the empty set is embedded within these axioms.

By Axiom 3, if we consider two mutually exclusive events which are both the empty set, we have:

$$P(\emptyset \cup \emptyset) = P(\emptyset) + P(\emptyset)$$

But the union of the empty set with itself remains the empty set, so we end up with:

$$P(\emptyset) = P(\emptyset) + P(\emptyset)$$

Which simplifies to (by subtracting  $P(\emptyset)$  from both sides):

$$0 = P(\emptyset)$$

Therefore, by using the axioms, we can prove that  $P(\emptyset) = 0$ .

### 3.6 Probabilities of Unions & Intersections

#### Theorem 1

The probability of the union of two events  $A$  and  $B$ , denoted as  $P(A \cup B)$ , is calculated using the sum rule 1.0.6.

#### Proof 2

We start by defining our events:

1.  $E1 = A \cap B'$
2.  $E2 = A \cap B$
3.  $E3 = A' \cap B$

Here,  $A'$  and  $B'$  denote the complements of  $A$  and  $B$ , respectively, representing all outcomes not in  $A$  or  $B$ . Also,  $E1$ ,  $E2$  and  $E3$  are pairwise mutually exclusive events. This means that if one of these events happens, the others cannot. It's an either/or situation, not both.

- $E1$  and  $E2$  are mutually exclusive because if  $E1$  happens, then  $E2$  cannot.
- $E2$  and  $E3$  are mutually exclusive because if  $E2$  happens, then  $E3$  cannot.
- $E1$  and  $E3$  are mutually exclusive because if  $E1$  happens, then  $E3$  cannot.

Therefore, by applying axiom 3, we have:

$$P(A) = P(E1) + P(E2)$$

$$P(B) = P(E2) + P(E3)$$

Looking at it now, it's clear that  $A \cup B = E1 \cup E2 \cup E3$ . This means by applying axiom 3 again we have:

$$P(A \cup B) = P(E1) + P(E2) + P(E3)$$

Since  $P(A) = P(E1) + P(E2)$  and  $P(B) = P(E2) + P(E3)$ , we have:

$$P(A \cup B) = P(A) + P(B) - P(E2)$$

Since  $E2 = A \cap B$ , we have:

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

Therefore, we have successfully proven Theorem 1.

---

From the theorem, we were also able to derive the following:

1. **Inclusion-Exclusion Principle:** The theorem shows that the probability of either  $A$  or  $B$  occurring isn't just the sum of the probabilities of  $A$  and  $B$ . The intersection is double-counted when we add these probabilities, so we need to subtract it out.
2. **Upper Bound:** The theorem also provides an upper bound on the probability of the union of two events. The probability of either  $A$  or  $B$  happening is at most the sum of the probabilities of  $A$  and  $B$ . This keeps probabilities from going above 1, respecting the boundaries of the probability space.

Now, going back to our pizza party example for students enrolled in SIT192. Remember, we surveyed 100 students about their pizza preferences. Let  $A$  represent the event that a student prefers Meat Lovers and  $B$  represent the event that a student prefers Hawaiian.

In the new survey, we found that:

- 50 students prefer Meat Lovers only.
- 30 students prefer Hawaiian only.
- 20 students enjoyed both.

Using the Probability of Union, we get:

$$\begin{aligned}P(A \cup B) &= P(A) + P(B) - P(A \cap B) \\&= 0.50 + 0.30 - 0.20 \\&= 0.60\end{aligned}$$

So, the probability that a randomly chosen student prefers either Meat Lovers or Hawaiian (or both) is 0.6 or 60

The reason I have chosen this example is that it highlights the usage of the theorem and its proof. From the survey, we were able to systematically understand the students' pizza preferences and make an informed decision about the pizza order for the SIT192 party. Not only have we seen this proof in action, but we also see the consequences of the theorem at play in this example:

When we calculated the probability of a student liking either of the pizzas, we had to subtract the probability of students liking both to avoid double-counting. This was using the inclusion-exclusion principle.

We also see the upper bound consequence, the probability of a student liking either type of pizza (or liking both) is less than or equal to the sum of the individual probabilities of liking Meat Lovers and Hawaiian.

**Probability of Intersection:** The probability of the intersection of two events  $A$  and  $B$ , denoted as  $P(A \cap B)$ , represents the probability that both events  $A$  and  $B$  happen at the same time.

We apply the product rule 1.0.5 when two events are independent of each other – which means two events are independent if the occurrence of one does not affect the probability of the occurrence of the other.

We will investigate this further when I introduce the concept of independent events in the upcoming section.

### 4 Conditional Probability Using Bayes' Rule

Here's probability theory starts to get a little more interesting. However, that is not to say the previous concepts *weren't interesting*, it just gets more... exciting. We're about to dive into areas where probability is not just about counting outcomes and assigning equal chances, but where we consider the context and specific conditions.

In my opinion, applying the fundamentals isn't as intimidating as one might initially believe. Take the example of the SIT192 pizza party example; the application of probability concepts appears straightforward, does it not? However, as we dive deeper, we start to see that the application of these concepts is only the tip of the iceberg. The real challenge lies beneath the surface, in the ability to apply logical thinking and problem-solve. That's where things become less trivial because suddenly, the principles that seemed intuitive don't fit as neatly into the problems we're trying to solve. That, in my opinion, is the true allure of probability.

The core of probability lies in translating a logical problem into a quantitative one, turning abstract ideas into numbers. You will also start to find that it is not just about inserting values into formulas to arrive at an answer. To quote Laplace from 1819, "Probability theory is nothing but common sense reduced to calculation." I found this a nice way to think of probability. I first read this quote in *Probability Theory: The Logic of Science* by E. T. Jaynes, which is a very interesting read and definitely a book I recommend.

Think about it like this, imagine yourself as a data scientist working on predicting a customer's likelihood of buying a particular product. You find that 10% of customers buy this product. Ok, so what if you also find that there are customers who have bought similar products in the past? How does this new information change the model? It's not just about the overall probability anymore, it's about the probability given this new piece of information - given that the customer has a history of buying similar products. The odds of that customer buying this product are likely to be higher than 10%.

#### 4.1 Independent Events

##### Definition 1.0.7: Independence

Events  $A$  and  $B$  are said to be independent if the fact that one event happens does not affect the probability that the other event will happen.

If whether event  $B$  happens does not affect the probability of event  $A$ , then event  $A$  is independent of event  $B$ . When two events are independent, such as  $A$  and  $B$ , we can express it as:

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

Let me introduce 1.0.7 with an entirely new example. Suppose we have a fair six-sided die and a fair coin. We are interested in event  $A$ , which is rolling a 6 and event  $B$ , which is flipping heads. We know that these two events are independent of each other. Rolling a 6 has nothing to do with flipping a head. If you don't believe me, I suggest you try conducting this experiment to see for yourself.

---

Since rolling a 6 is one out of six possible outcomes, the probability of  $A$  is  $1/6$  and since flipping heads is one out of two possible outcomes, the probability of  $B$  is  $1/2$ .

Using the product rule (Rule 2.4.1), we can now figure out the probability of getting a 6 and a head:

$$P(A \cap B) = P(1/6) \cdot P(1/2) = 1/12 \approx 8.3\%$$

Let's think more deeply about this. We know that rolling a 6 and landing heads are independent events. How do we know this? If we refer back to Kolmogorov's notes, he defines an event  $X$  as one that occurs if and only if events  $A$ ,  $B$  and so forth, all occur simultaneously or concurrently. In this case, event  $X$  is the outcome where we roll a 6 and flip heads, while event  $A$  is rolling a 6 and event  $B$  is flipping heads. Event  $X$  cannot happen unless both  $A$  and  $B$  happen.

With this understanding, we can apply the product rule, a direct consequence of the definition of independence, to calculate the probability of the intersection of two independent events.

## 4.2 Conditional Probability

The axioms of probability provided by Kolmogorov form a cornerstone for the concept of conditional probability. It's based on the natural understanding of the revised likelihood of an event given that another event has already occurred.

### Definition 1.0.8: Conditional Probability

If we consider two events  $A$  and  $B$  within a given sample space, with  $P(B) > 0$ , the conditional probability of  $A$  given  $B$  is defined as follows:

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

If whether event  $B$  happens does not affect the probability of event  $A$ , then event  $A$  is independent of event  $B$ . When two events are independent, such as  $A$  and  $B$ , we can express it as:

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

In the definition above,  $B$  acts as our new sample space and we are interested in the likelihood of event  $A$  within this new space. It's important to remind ourselves that  $P(B) > 0$  is required to avoid division by zero.

Suppose you are planning to go for a surf tomorrow. There are two possible events in your case:

- Event  $A$ : It storms tomorrow.
- Event  $B$ : You go surfing.

## 1. Discrete Probability

---

You have access to a pretty accurate weather app and feel confident in using it to conclude on a decision.

1. Based on the weather forecast and historical data, you estimate there's a 30% chance it will storm tomorrow. So,  $P(A) = 0.30$ .
2. You're a brave soul who loves to surf storm or shine, so there's a 90% chance you'll surf.  $P(B) = 0.90$ .
3. Based on past experiences, you know that 20% of your surfs have occurred on stormy days. This is the intersection of events  $A$  and  $B$ .  $P(A \cap B) = 0.20$ .

We want to find  $P(A|B)$ , which is the probability it will storm given that you'll surf. By using the formula of conditional probability (Definition 3.2):

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0.20}{0.90} \approx 0.22 \text{ or } 22\%$$

So, given that you decide to surf, there's a 22% chance that it will storm.

Now, let's also calculate  $P(B|A)$ , the probability that you will surf given that it storms:

$$P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{0.20}{0.30} \approx 0.67 \text{ or } 67\%$$

So, if it storms, there's a 67% chance that you will surf.

What does this mean? Well... Given that you decide to surf, there's a 22% chance that it will storm. This tells you the risk of encountering a storm if you go surfing. However, if it storms, there's a 67% chance that you will surf. This tells you how likely you are to stick with your plan to surf, even if it storms.

So really, it depends on your attitude toward risk and how much you dislike surfing in a storm. If you're not deterred by the prospect of a storm, the high probability of surfing (90%) means you're likely to go regardless of the weather. Knowing that there's a 22% chance of a storm shouldn't change this.

On the other hand, if you're considering whether to change your plans when you know there's a storm coming, the 67% figure is relevant. This is pretty high, which means you're fairly committed to surfing, even when a storm is brewing.

Of course, you might also consider other factors, like how severe the storm is likely to be, whether there are safe places to shelter, etc. But from a purely probabilistic standpoint, if you're okay with a 22% chance of surfing in a storm, then you should go for it!

---

### 4.3 Law of Total Probability

The Law of Total Probability is like a recipe that helps you find the overall likelihood of an event by considering all the different ways it could happen. This "recipe" or law is especially handy when you have a situation that can be divided into several distinct or non-overlapping scenarios, called partitions. Each of these partitions represents a different way the event of interest can happen.

The law states that the total probability of an event  $A$  is the sum of the probabilities of  $A$  happening across all possible scenarios. If you think of the event  $A$  as a pizza, then the Law of Total Probability tells us that to find the total "probability pizza," we need to add up all the "probability slices" from each different scenario.

#### Definition 1.0.9: Partition of a Sample Space

A collection of events  $\{B_1, B_2, \dots, B_n\}$  is said to be a partition of the sample space  $S$  if the following conditions are satisfied:

1. The events  $B_i$  are pairwise mutually exclusive and exhaustive – for any pair of unique events in the collection  $B_i$  and  $B_j$ , we have  $B_i \cap B_j = \emptyset$ .
2. The union of all the events in the collection forms the sample space:  $B_1 \cup B_2 \cup \dots \cup B_n = S$ .
3. If  $P(B_i) = 0$  for some  $i$ , then  $B_i$  is an event that never occurs, but this doesn't affect the partition conditions. However, for computing conditional probabilities like  $P(A|B_i)$ , we need  $P(B_i) > 0$  to avoid division by zero.

#### Theorem 2: The Law of Total Probability

Let  $B_1, B_2, \dots, B_n$  be a partition of the sample space  $S$  but together cover all possible outcomes. Let  $A$  be an event. Then the probability of  $A$  can be calculated as:

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

**Proof 3: The Law of Total Probability**

To prove the theorem, we can use the definition of conditional probability defined in Definition 1.0.2, Kolmogorov's third axiom and Definition 1.0.3.

Let us assume that  $\{B_1, B_2, \dots, B_n\}$  forms a partition of the sample space  $S$ . This means that the sets  $B_i$  are mutually exclusive and exhaustive—no two sets share elements and their union is the entire sample space. Therefore, we can represent Event  $A$  as a union of mutually exclusive events and by applying the third axiom:

$$P(A) = P(A \cap B_1) + P(A \cap B_2) + \dots + P(A \cap B_n)$$

Then, by using the definition of conditional probability, we have:

$$P(A \cap B_i) = P(A|B_i) \cdot P(B_i), \text{ for each } i$$

To understand why, let me first illustrate quickly:

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

If we multiply both sides by  $P(B_i)$ , we can then cancel out the like terms to get:

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

Now we apply this definition to  $P(A)$ , which gives:

$$P(A) = P(A|B_1) \cdot P(B_1) + P(A|B_2) \cdot P(B_2) + \dots + P(A|B_n) \cdot P(B_n)$$

Which now can be rewritten as:

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

This completes the proof of Theorem 2.



---

## 4.4 Bayes' Theorem

Bayes' theorem presents an approach for computing the conditional probability of an occurrence  $A$ , given occurrence  $B$ . It operates on the premise that we're already aware of the probabilities of  $A$ ,  $B$  considering  $A$  and  $B$  in the absence of  $A$ .

### Theorem 3: Bayes' Theorem

For events  $A$  and  $B$ , with  $P(A) > 0$  and  $P(B) > 0$ , the conditional probability of  $A$  given  $B$ , denoted as  $P(A|B)$ , is given by:

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

### Proof 4: Baye's Theorem

We will prove Bayes' theorem using the definition of conditional probability and the Law of Total Probability, defined in definition 1.0.3 and Theorem 2.

Firstly, we start by using the definition of conditional probability:

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

and also:

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

While proving Theorem 2, we pointed out that  $P(A \cap B) = P(B|A) \cdot P(A)$ . Therefore, we can replace  $P(A \cap B)$  with this expression:

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

Now, Bayes' theorem is proved. However, let me introduce this by revisiting our previous example in conditional probability. Previously we calculated that:

$$P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{0.30}{0.60} = 0.50$$

But what if now we want to find the conditional probability  $P(A|B)$ , which represents the probability that a student understands Propositional Logic given that they already understand Sets and Functions.

Well, using Bayes' theorem, we can calculate this as:

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} = \frac{0.50 \cdot 0.60}{0.50} = 0.60$$

Also, sometimes, it is much easier to find  $P(B|A)$  than it is  $P(A|B)$ .

Understanding this theorem is important, especially in probability. Naïve Bayes used for classification problems is built on this theorem.

### 5 Discrete Distributions

In this final section of the module, we will explore a set of discrete probability distributions. These distributions are used for describing various types of random phenomena, particularly those which involve discrete outcomes or events.

We will explore five fundamental probability distributions where each type of distribution are fundamental tools for describing various types of random phenomena.

We start with the Bernoulli distribution, a probability distribution that models the outcome of a single binary experiment.

We'll then introduce the binomial distribution, which is an extension of the Bernoulli distribution, but instead of a single trial, we consider multiple independent and identical Bernoulli trials.

Here, the outcome of one trial does not influence the others. But, when we count the number of successes across these trials, we see the binomial distribution presenting a comprehensive depiction of how randomness behaves over repeated trials.

Following this, we examine the geometric and hypergeometric distributions. The geometric distribution describes the number of trials required to achieve the first success, which is helpful in situations like product testing or clinical trials where the focus is on achieving that first breakthrough.

To finish off, we'll encounter the Poisson distribution. Which is used estimate the frequency of events within a specific period.

These discrete distributions and their corresponding Probability Mass Functions (PMFs) make up a nice toolkit for handling randomness. Their applications span a wide variety of fields, a few too many to mention.

#### 5.1 Bernoulli Distribution

Picture this: you're at a roulette table, placing bets on either red or black. Each bet you place, each spin of the wheel, is an independent trial - a Bernoulli trial, to be exact.

So, what exactly is a Bernoulli trial? Well, it's quite simple, yet very important. In a Bernoulli trial, you're dealing with something random, a situation that has exactly two outcomes. It could be red or black on the roulette wheel, heads, or tails on a coin flip, but more formally we would say success or failure.

Key properties of a Bernoulli trial:

1. Independence: The outcomes of each trial are independent events, meaning the occurrence of one outcome does not influence the occurrence of another outcome.
2. Two possible outcomes: The sample space  $S$  contains exactly two distinct outcomes,  $\{0, 1\}$ , representing failure and success respectively.
3. Constant probabilities: The probabilities of success and failure,  $p$  and  $q$ , remain fixed throughout the experiment. These probabilities are not dependent on previous trials or outcomes.
4. Probability distribution: Where there are two outcomes, the probabilities can be defined as  $P(0) = q$  and  $P(1) = p$ , satisfying the conditions  $P(0) + P(1) = q + p = 1$ .

---

Now that we understand what a Bernoulli trial is, what exactly is the Bernoulli distribution? Well, it is the probability distribution of a random variable which takes a binary, Boolean output: 1 (success) or 0 (failure).

Well, that is because it is. Essentially, a Bernoulli distribution represents the outcome of a single Bernoulli trial and we can model this by using the Probability Mass Function (PMF) of the Bernoulli distribution:

$$P(X = k) = p^k(1 - p)^{1-k} \text{ for } k \in \{0, 1\}$$

Suppose there are 150 students in a course. Based on previous data, we know that 85% of students pass the course while the rest fail 15

$$\text{When } k = 1 : P(X = 1) = 0.85^1(1 - 0.85)^{1-1} = 0.85 \text{ or } 85\%$$

$$\text{When } k = 0 : P(X = 0) = 0.85^0(1 - 0.85)^{1-0} = 0.15 \text{ or } 15\%$$

Visually, we can imagine the distribution looking like so:

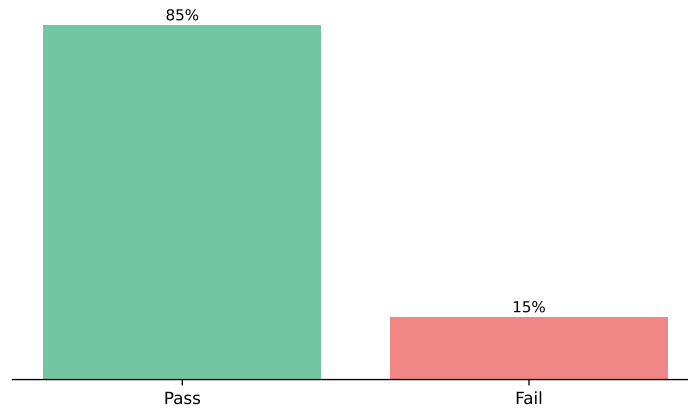


Figure 1.4: Representation of Bernoulli Distribution where  $p = 0.85$  and there are 150 samples (see code 1).

## 5.2 Probability Mass Functions (PMF)

Here we discuss what a probability mass function (PMF) is, which was first introduced prior when introducing the Bernoulli distribution. The reason it was structured this way, is because it helps us prepare for a broader understanding of PMF. By understanding binary outcomes and their associated probabilities, we can more intuitively understand the concept of a PMF, which can allow variables to take on more than just two outcomes.

However, before I provide a more general definition of this, we need to understand what a random variable is.

A random variable, which we typically denote as  $X$ , is a variable representing the numeric outcomes of a random phenomenon. Random variables are either discrete or continuous. However, we will only be focusing on discrete probabilities.

However, what is a discrete random variable? Well, it is a random variable which has the distinct feature of assuming a set of specific, countable values.

## 1. Discrete Probability

Imagine you're at a basketball game and we're observing Stephen Curry who's known for his three pointers. We could define a variable  $Y$  to represent the number of successful shots he makes in a game. Here,  $Y$  isn't just a typical variable. Instead, its value relies on the unpredictable performance of the player, thereby making  $Y$  a random variable.

Considering the random characteristic of  $Y$ , it wouldn't be logical to ask outright if  $Y = 3$ , for example, but it's entirely reasonable to ask about the probability that  $Y = 3$ , or if  $Y < 2$ . This is because  $Y = 3$  and  $Y < 2$  correlate to events within our sample space, set out by the player's shooting outcomes. The likelihoods of these specific outcomes, or the  $Y = 3$  and  $Y < 2$  events, are formally characterised by the Probability Mass Function (PMF).

### Definition 1.0.10: Probability Mass Function (PMF)

If  $Y$  is a discrete random variable with a possible range of values represented as  $\{y_1, y_2, y_3, \dots\}$ , then the Probability Mass Function of  $Y$ , denoted as  $P(Y = y_i)$  for all  $i$  in the index set, assigns probabilities to each of these feasible values of  $Y$ .

One important property of the PMF is that the sum of probabilities for all possible values of the discrete random variable equals 1. This is because the probabilities for each outcome collectively represent all possible events within the sample space. While I used  $Y$  in this example, it's more common to see  $X$  used to represent a random variable in PMF definitions. However, this is a matter of convention and any letter can be used to denote the random variable.

## 5.3 Binomial Distribution

The Binomial distribution models the number of successes in a fixed number of independent Bernoulli trials. Recall from our previous discussions on Bernoulli trials that they consist of only two outcomes, labelled as "success" and "failure," with probabilities  $p$  and  $q = 1 - p$ , respectively. The distribution of such trials is governed by the Bernoulli distribution.

Now suppose we conduct  $n$  independent Bernoulli trials, each with a success probability  $p$ . Let's denote the number of successes as  $k$ . We are interested in finding the probability of obtaining exactly  $k$  successes, written as:

$$P(X = k)$$

To derive this, we'll use the concept of the Bernoulli distribution which was explained in the section prior and expand it to  $n$  trials. For any sequence of  $k$  successes and  $n - k$  failures, the probability will be the product of individual probabilities. This is due to the property of independence in Bernoulli trials. Therefore, for a specific sequence, we have:

$$p^k(1 - p)^{n-k}$$

However, there are multiple ways of obtaining  $k$  successes in  $n$  trials. The number of such combinations is given by the binomial coefficient  $\binom{n}{k}$ , representing the number of ways to choose  $k$  successes from  $n$  trials.

The total probability  $P(X = k)$  for  $k$  successes in  $n$  trials is given by multiplying the probability of a single sequence by the number of such sequences. This gives us the Probability Mass Function (PMF) of the Binomial distribution:

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

---

This shows that the Binomial distribution is indeed a generalisation of the Bernoulli distribution to a fixed number of  $n$  independent trials. It also gives a comprehensive model for representing the number of successes in  $n$  trials.

Using the same pizza example from the section of this module, let's consider each student's preference for Hawaiian as a Bernoulli trial. Each student either prefers Hawaiian (success) or doesn't (failure).

Suppose we take a random sample of 10 students out of 100. What is the probability that exactly  $k$  students prefer Hawaiian? We can start by using the Binomial distribution PMF:

$$P(X = k) = \binom{10}{k} \cdot 0.35^k (1 - 0.35)^{10-k}$$

If  $k = 3$ , what would the probability be?

$$P(X = 3) = \binom{10}{3} \cdot 0.35^3 \cdot (1 - 0.35)^{10-3} \approx 0.252 \text{ or } 25.2\%$$

This means the probability of finding exactly 3 students in a random sample of 10 – meaning they were randomly chosen from the 100 students, would be approximately 25.2%. How about for other  $k$  values in the sample of 10?

Well we can use a visualisation to show us each:

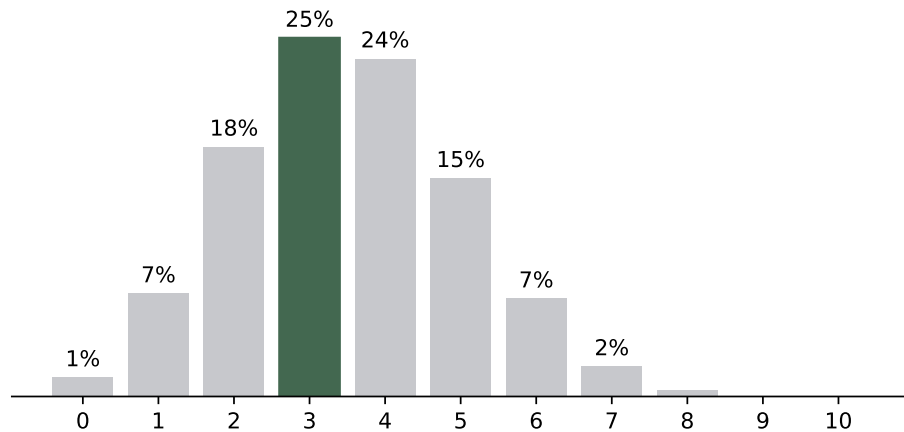


Figure 1.5: Representation of Binomial Distribution where  $p = 0.35$  and there are 10 students from 100. (see code 1)

How about we make this a little more interesting and we include the visualisation of the number of students picking Meat Lovers as well? The math stays the same, the only thing that changes is  $p$ :

$$P(X = k) = \binom{10}{k} \cdot 0.50^k \cdot (1 - 0.50)^{10-k}$$

## 1. Discrete Probability

Now what is the probability we find exactly 3 students picking Meat Lovers? Well, we input  $k = 3$  and...

$$P(X = 3) = \binom{10}{3} \cdot 0.503 \cdot (1 - 0.50)^{10-3} \approx 0.117 \text{ or } 11.7\%$$

What is interesting is that even though the probability of finding a student who likes Meat Lovers is significantly higher, the probability of finding exactly 3 out of 10 students who prefer Meat Lovers is significantly lower than finding 3 who prefer Hawaiian. Why is that?

It is because the binomial distribution considers not just the total number of trials and the number of successes, but also considers the underlying probability of success. In our pizza example, since more students prefer Meat Lovers, we're more likely to see a higher number of students preferring Meat Lovers in a random sample of 10 students.

Conversely, fewer students prefer Hawaiian, so it's more probable to see a lower number of students preferring Hawaiian in our sample. To finish off, I've provided a visualisation of the binomial distributions for both Meat Lovers and Hawaiian preferences below.

This gives us some interesting insights into how different population proportions can influence the probabilities of different outcomes in a binomial distribution.

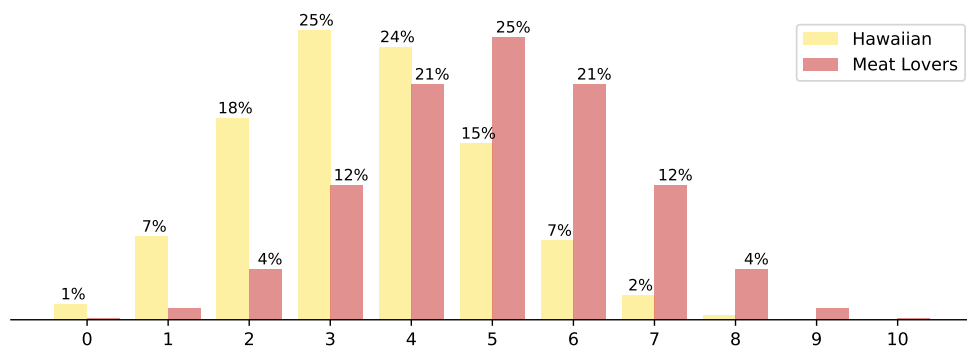


Figure 1.6: Representation of Binomial Distributions for Student Preferences, where one distribution corresponds to a probability of 0.35 which is Hawaiian and the other to a probability of 0.5, Meat lovers (see code 1).

---

## 5.4 Geometric Distribution

The geometric distribution is a probability distribution that models the number of trials needed to get the first success in repeated, independent Bernoulli trials. Here's an easy way to remember it: the geometric distribution effectively answers the question, "How many times do I need to try before I succeed?"

Just like the Binomial distribution, we can prove that the PMF of the geometric distribution using the concept of independent Bernoulli trials.

Suppose we denote the number of trials required as  $k$ . The event of interest is that the first success happens on the  $k$ -th trial.

Like in the binomial example, the Bernoulli trials are independent and each trial has a success probability of  $p$ . Therefore, any specific sequence of  $k - 1$  failures followed by a success has a probability of:

$$(1 - p)^{k-1} \cdot p$$

This represents the probability of getting  $k - 1$  failures in a row and  $p$  represents the probability of getting a success on the  $k$ -th trial.

However, unlike the binomial example, there's only one "ordering" of outcomes that results in the first success occurring on the  $k$ -th trial.

That is, you must observe  $k - 1$  failures followed by a success. This means that there are no combinations to consider and the probability of the first success occurring on the  $k$ -th trial is just the probability of this single sequence of outcomes.

Therefore, we have the PMF of the geometric distribution as follows:

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

Now, consider ice hockey player Connor McDavid who has a faceoff win percentage of 53.7%. Let's suppose you and a mate make a wager about whether Connor will succeed on their first faceoff after initially missing the first two.

Considering the player's success rate, you may feel confident about the bet, but let's determine the actual probability of success on the third attempt.

$$P(X = 3) = (1 - 0.537)^{3-1} \cdot 0.537 \approx 0.115 \text{ or } 11.5\%$$

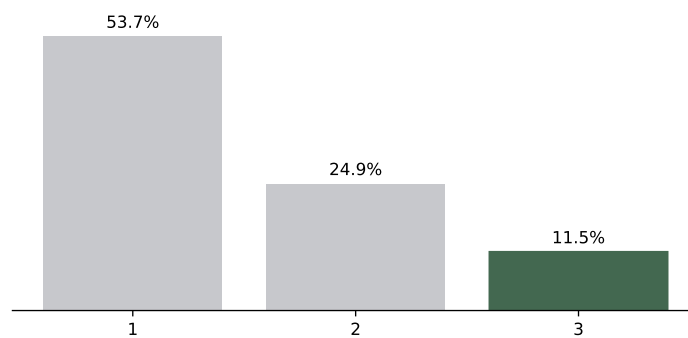


Figure 1.7: Representation of the Geometric Distribution for the Probability of Success on the  $k$ -th Trial (see code 1).

## 5.5 Hyper-geometric Distribution

Hyper-geometric distribution differs from the first two distributions (Bernoulli and geometric) because the trials are no longer independent, and the probabilities do not remain constant. This is because we are dealing with a finite population and selections are made without replacement which changes the probabilities with each trial.

Suppose we have a finite population of  $N$  items, where  $K$  items are considered "successes" and  $N - K$  items are "failures". We draw  $n$  items without replacement from this population. Any specific sequence of  $k$  successes and  $n - k$  failures has a certain likelihood. However, unlike in the Bernoulli or the geometric distributions, the probability is not simply:

$$p^k(1 - p)^{n-k}$$

This is because the events are not independent and the probability changes with each draw. However, like the Bernoulli and geometric scenarios, there are many ways that  $k$  successes can occur in  $n$  draws. We need to consider all possible orderings of these successes and failures.

- The number of ways that we can choose  $k$  successes from the  $K$  available is given by the binomial coefficient  $\binom{K}{k}$  – which is equivalent to  $\binom{n}{k}$  which we used prior.
- The number of ways we can choose  $n - k$  failures from the  $N - K$  failures available is given by  $\binom{N-K}{n-k}$ .

From this, we know that the total number of ways we can obtain a sequence of  $k$  successes and  $n - k$  failures is by  $n$  draws in the given product:

$$\binom{K}{k} \cdot \binom{N-K}{n-k}$$

The total number of ways one can draw  $n$  items from the population  $N$  is given by the binomial coefficient:

$$\binom{N}{n}$$

Therefore, with this understanding we can conclude that the hyper-geometric distribution is obtained by taking the ratio of the number of favourable outcomes to the total number of outcomes and can be given by:

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

Let me introduce this with an example. Suppose that you are a consultant, and a high school has hired you to help them with a problem they are currently facing. They have around 2500 students, and they have recently found that around 150 of them are secretly using ChatGPT to solve their schoolwork. They know this only because of the similarities between each other work and because of the incorrect answers which the model produces.

They have an exam coming up which will provide a scholarship to a top University in Mount Sasquatch. There are 10 students who were selected for this exam, and they want to know the probability of finding exactly 3 students cheating in the exam.



---

For this problem we can define the following:

$$N = 2500$$

$$K = 150$$

$$n = 10$$

$$k = 3$$

Putting this into the PMF function for Hypergeometric distributions:

$$P(X = 3) = \frac{\binom{150}{3} \cdot \binom{2350}{7}}{\binom{2500}{10}} \approx 0.0166$$

You find that the result of finding exactly 3 cheaters in the exam is approximately 0.0166 or 1.66%. You go back to the high school and provide them with the results. To assist them further you decided to compute the hyper-geometric distribution for them, so they can view the different probabilities in the distribution of 10 students.

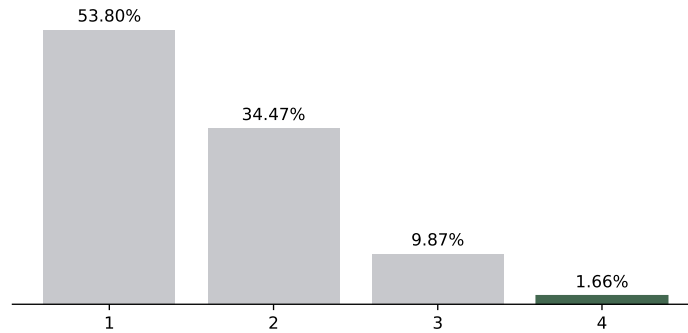


Figure 1.8: Representation of the Hypergeometric Distribution for the Probability of Success on the  $k$ -th Trial (see code 1).

### 5.6 Poisson Distribution

A very interesting distribution and quite different from the previous distributions in terms of the type of events it models and the parameter it uses. The Poisson distribution expresses the probability of a given number of events occurring in a fixed interval of time or space, given a fixed average rate of occurrence. These events must be independent, meaning they don't affect each other.

The PMF of a Poisson distribution is given by:

$$P(k; \lambda) = \frac{\lambda^k \cdot e^{-\lambda}}{k!}$$

Where we have:

- $P(k; \lambda)$  is the Poisson probability, which asks the question "what is the probability of  $k$  events happening in an interval?"
- $\lambda$  is the rate parameter (also known as the rate or mean number of events).
- $k$  is the actual number of events that result.

Now let's say you got another job, remember you're a consultant... the job is to assist Student Connect at Deakin University by providing insights into the number of phone calls they receive per hour. Suppose that they receive an average of 12 calls per hour. The Poisson distribution can then be used to calculate the probability of a certain number of calls (let's say 15 calls) happening in the next hour.

---

So here, we define the following:

$$P(k; \lambda) = P(15; 12)$$

$$\lambda = 12$$

$$k = 15$$

You then plug in the numbers into the PMF:

$$P(15; 12) = \frac{12^{15} \cdot e^{-12}}{15!} \approx 0.0724 \text{ or } 7.24\%$$

After performing the calculations, you find that there is approximately a 7.24% chance that Deakin Connect will receive exactly 15 calls per hour on any given day. They were happy with your result, and you also provided the visualisation of the Poisson distribution of 15 calls an hour with an average rate of 12 per hour.

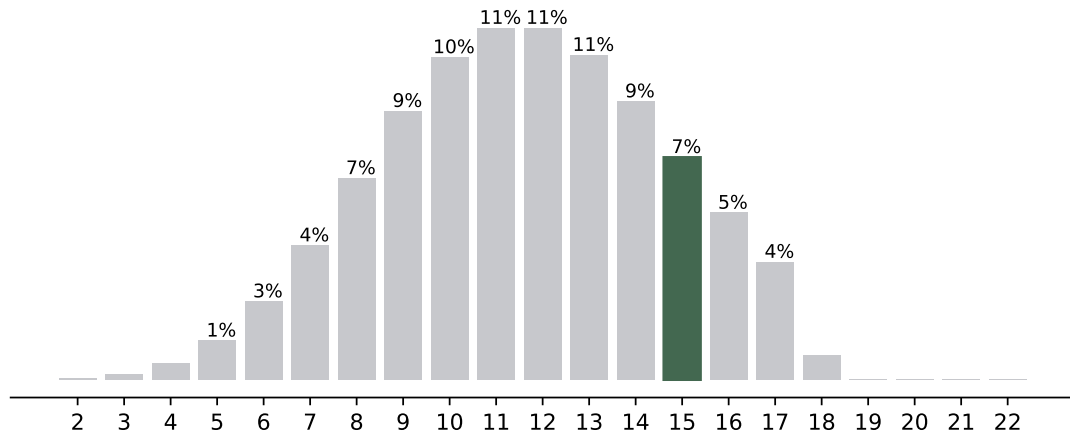


Figure 1.9: Representation of the Poisson Distribution with an average of 12 calls per hour (see code 1).

### 6 Recap

We have finally reached the end of this module, to quickly recap. You were introduced to some core concepts in probability theory and their applications. It started off with probability models, which are mathematical representations that help us navigate and understand uncertain scenarios. We then explored conditional probability using Bayes' Rule, where we examined how prior knowledge or conditions can affect our probability calculations. Then finally, we looked at discrete distributions, which provide a way to model the probability of outcomes for discrete random variables.

We learnt that probability is more than just counting outcomes and assigning chances. It's about understanding the context, the conditions and making sense of it all through logic and quantitative reasoning. From the basic axioms of probability to complex concepts like Bayes' theorem and probability mass functions, each topic builds upon the previous, leading to a deeper understanding of this fascinating field.

Probability theory is an incredibly powerful tool. Which underpins many practical applications, from data science and predictive modelling to everyday technologies. Understanding these fundamental concepts doesn't just deepen our mathematical knowledge; it helps you see the world in a different way. As a quote mentioned prior by Laplace, "Probability theory is nothing but common sense reduced to calculation". But understanding the theory is only half the journey. The other half is applying it. Whether you're predicting customer behaviour, modelling outcomes of a basketball game, or navigating the odds in a game of chance, the principles and methods learned in this module are valuable tools.

I do recommend to check out E. T. Jaynes's "Probability Theory: The Logic of Science" and definitely recommend to look through Foundations of the theory of probability by Kolmogorov, A. N. Remember, probability isn't just a branch of mathematics, it's a way of thinking, a way of making sense of the world around us, one uncertainty at a time.

### About the Author



Brandon Smith is a second year undergraduate student completing his Bachelor of Data Science. He has a keen interest in research, especially in the field of information retrieval and natural language processing. He believes that learning can be simple, we just need to find the right way to learn that suits and dedicate the necessary time to pursue it.

### Acknowledgements

A big thank you to Julien Ugon for not only being a fantastic teacher but also a great mentor and role model. The unit SIT192 really changed the way I learned and greatly motivated me, sparking a passion for learning new things especially in mathematics! This experience has really influenced my passion for research and I have found a path to pursue in my life with full drive—a path which might not have been possible if I had never crossed paths with Julien. Thank you!

---

## References

- [1] A.N. Kolmogorov. *Foundations of the Theory of Probability*.

## Appendix

In this section, I'll showcase examples of concepts I find interesting, particularly by utilising code to demonstrate probability.

Just a quick remark, there are libraries out there that you can use to easily compute probability distributions, matrix multiplications, and other types of mathematical functions. However, my strong recommendation is to try it the first time yourself purely through code. This is how you build a strong intuition rather than relying solely on encapsulated functions provided by these libraries.

This holds true for all things that can become quite complicated later on. We have a lot of software these days that can help us get past not understanding what is lying under the hood. However, as an aspiring researcher or technical professional, knowing the underlying functionality will help with future problems. It allows you to be more innovative and create new solutions.

My rule is start small, then understand the small stuff then move forward. If you are interested in Deep Learning but skip past all the fundamental models such as Linear Regression, Decision Trees and Clustering, you end up missing out on a lot of intuition and possibly better solutions than Deep Learning.

For all distributions I have introduced and provided examples for, you can encapsulate much of the coding by using SciPy, a Python library.

### Law of Large Numbers

We begin with a simple, yet interesting concept the Law of Large Numbers also known as Bernoulli's theorem. This theorem, introduced in nearly all introductory probability textbooks, states that as the number of independent samples increases, their average result should converge to the true mean value.

There are proofs for both small and large numbers, however, sometimes it is much easier to visualise this working in action through a simulation, which can be done if we produce code for it.

In the two visualisations referenced in 1.10, it can be seen that the probability of rolling a one or flipping heads, calculated as 0.16 and 0.5, respectively, does not stabilise until approximately 200 trials for both simulations. Which is not a \*proof\* by any means but just a way to help us conceptualise things.

We can even derive insights from the visualisations, such as we can notice that the number of experiments increases, the observed outcomes approach the expected value, as mentioned in the theorem of the Law of Large Numbers.

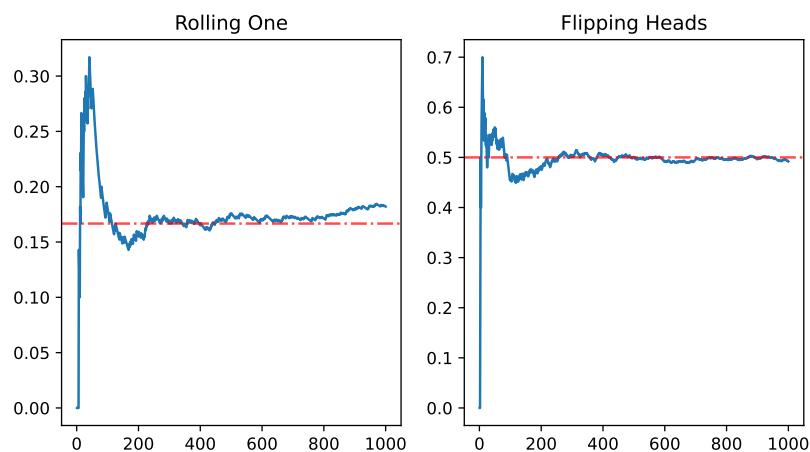


Figure 1.10: Simulation of rolling ones and flipping heads.

---

### Code for Law of Large Numbers Simulation.

---

```
import numpy as np
import matplotlib.pyplot as plt

def roll_die(n):
    o = np.random.randint(1, 7, size=n)
    c = np.cumsum(o == 1)
    p = c / np.arange(1, n + 1)
    return p

def flip_coin(n):
    o = np.random.randint(0, 2, size=n)
    c = np.cumsum(o == 1)
    p = c / np.arange(1, n + 1)
    return p

def plt_convergence(n):
    fig, axs = plt.subplots(1, 2, figsize=(7, 4))
    axs[0].plot(range(1, n + 1), roll_die(n))
    axs[1].plot(range(1, n + 1), flip_coin(n))

    axs[0].axhline(y=1/6, color='r', linestyle='-.', alpha=0.7)
    axs[1].axhline(y=1/2, color='r', linestyle='-.', alpha=0.7)

    axs[0].set_title('Rolling One')
    axs[1].set_title('Flipping Heads')

    plt.tight_layout()
    plt.show()
```

---

### Probability Distributions

Now I will provide the code for each distribution problem we introduced in the prior sections.

#### Code for simulating Bernoulli Distribution (see figure 1.4).

---

```
p = {'Pass': 0.85, 'Fail': 0.15}

fig, ax = plt.subplots(1,1, figsize=(6, 4))
ax.bar(
    p.keys(),
    p.values(),
    color=['#36AE7C', '#EB5353'],
    alpha=0.7
)

for i, probs in enumerate(p.values()):
    if probs > 0.01:
        ax.annotate(
            f"{np.round(probs*100):.0f}%",
            (k[i], probs),
            xytext=(0, 1),
            textcoords="offset points",
            ha='center',
            va='bottom',
            fontsize=9
        )

ax.spines['top'].set_visible(False)
ax.spines['left'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.set_yticklabels([])
ax.set_yticks([])

ax.set_ylim(0, 1)
plt.tight_layout()
plt.show()
```

---



---

### Code for simulating Binomial Distribution 1 (see figure 1.5).

---

```
import matplotlib.pyplot as plt
import numpy as np
import math

# https://en.wikipedia.org/wiki/Binomial\_coefficient
def nCk(n, k):
    result = math.factorial(n) / (math.factorial(k) * math.factorial(n - k))
    return result

# https://en.wikipedia.org/wiki/Binomial\_distribution
def binomial_distribution(n, k, p):
    result = nCk(n, k) * (p ** k) * ((1 - p) ** (n - k))
    return result

n = 10
k = list(range(11))
p = [binomial_distribution(n, i, 0.35) for i in range(n + 1)]

fig, ax = plt.subplots(1, 1, figsize=(6, 3))
bar = plt.bar(
    k,
    p,
    color='#C7C8CC',
)

bar[3].set_color('#436850')

for i, probs in enumerate(p):
    if probs > 0.01:
        ax.annotate(
            f"{np.round(probs*100):.0f}%",
            (k[i], p[i]),
            xytext=(0, 3),
            textcoords="offset points",
            ha='center',
            va='bottom',
        )

ax.spines['top'].set_visible(False)
ax.spines['left'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.set_yticklabels([])
ax.set_yticks([])

plt.xticks(k)
plt.tight_layout()
plt.show()
```

---

## 1. Discrete Probability

---

### Code for simulating Binomial Distribution 2 (see figure 1.6).

---

```
fig, ax = plt.subplots(1,1, figsize=(8, 3))

p2 = [binomial_distribution(n, i, 0.5) for i in range(n + 1)]

plt.bar([i - 0.2 for i in k], p, width=0.4, label='Hawaiian', color='#FDE767', alpha=0.6)
plt.bar([i + 0.2 for i in k], p2, width=0.4, label='Meat Lovers', color='#D04848', alpha=0.6)

for i, probs in enumerate(p):
    if probs > 0.01:
        ax.annotate(
            f"{np.round(probs*100):.0f}%",
            (k[i], p[i]),
            xytext=(-8, 1),
            textcoords="offset points",
            ha='center',
            va='bottom',
            fontsize=9
        )

for i, probs in enumerate(p2):
    if probs > 0.01:
        ax.annotate(
            f"{np.round(probs*100):.0f}%",
            (k[i], p2[i]),
            xytext=(12, 1),
            textcoords="offset points",
            ha='center',
            va='bottom',
            fontsize=9
        )

ax.spines['top'].set_visible(False)
ax.spines['left'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.set_yticklabels([])
ax.set_yticks([])

plt.xticks(k)
plt.legend()
plt.tight_layout()
plt.show()
```

---

---

### Code for simulating Geometric Distribution (see figure 1.7).

---

```
geometric_distribtion = lambda p, k: [
    (1 - p)**(i - 1) * p for i in range(1, k+1)
]

get_exact_k = lambda p, k: (1 - p)**(k - 1) * p

p = 0.537
k = 3
x_axis = list(range(1, k+1))

fig, ax = plt.subplots(1,1, figsize=(6, 3))
bar = ax.bar(
    x_axis,
    geometric_distribtion(p,k),
    color='#C7C8CC',
)

bar[geometric_distribtion(p,k).index(get_exact_k(p,k))].set_color('#436850')

for i, probs in enumerate(geometric_distribtion(p,k)):
    if probs > 0.01:
        ax.annotate(
            f"{np.round(probs*100, 1):.1f}%",
            (i+1, probs),
            xytext=(0, 3),
            textcoords="offset points",
            ha='center',
            va='bottom',
        )

ax.spines['top'].set_visible(False)
ax.spines['left'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.set_yticklabels([])
ax.set_yticks([])

ax.set_xticks(x_axis)
plt.tight_layout()
plt.show()
```

---

## 1. Discrete Probability

---

Code for simulating Hypergeometric Distribution (see figure 1.8).

---

```
def nCk(n, k):
    result = math.factorial(n) / (math.factorial(k) * math.factorial(n - k))
    return result

def get_exact_k(N, K, n, k):
    return (nCk(K, k) * nCk(N - K, n - k)) / nCk(N, n)

N = 2500
K = 150
n = 10
k = 3

# create the distribution
hypergeo_distribution = [get_exact_k(N, K, n, k) for k in range(4)]
x_axis = list(range(1, k+2))

fig, ax = plt.subplots(1,1, figsize=(6, 3))
bar = ax.bar(
    x_axis,
    hypergeo_distribution,
    color='#C7C8CC',
)

bar[hypergeo_distribution.index(get_exact_k(N, K, n, k))].set_color('#436850')

for i, probs in enumerate(hypergeo_distribution):
    if probs > 0.01:
        ax.annotate(
            f"{np.round(probs*100, 2):.2f}%",
            (i+1, probs),
            xytext=(0, 3),
            textcoords="offset points",
            ha='center',
            va='bottom',
        )

ax.spines['top'].set_visible(False)
ax.spines['left'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.set_yticklabels([])
ax.set_yticks([])

ax.set_xticks(x_axis)
plt.tight_layout()
plt.show()
```

---

---

### Code for simulating Poisson Distribution (see figure 1.9).

---

```
l = 12
k = 15

def get_extact_k(k, l):
    result = (l**k * math.exp(-l)) / math.factorial(k)
    return result

# create the distribution
x_axis = np.arange(l-10, l+11)
poisson_distribution = [get_extact_k(i, l) for i in x_axis]

fig, ax = plt.subplots(1,1, figsize=(7, 3))
bar = ax.bar(
    x_axis,
    poisson_distribution,
    color='#C7C8CC',
)

for i, probs in enumerate(poisson_distribution):
    if probs > 0.01:
        ax.annotate(
            f"{np.round(probs*100):.0f}%",
            (i+1, probs),
            xytext=(23,0),
            textcoords="offset points",
            ha='center',
            va='bottom',
            fontsize=8.5
        )

bar[poisson_distribution.index(get_extact_k(k, l))].set_color('#436850')
ax.spines['top'].set_visible(False)
ax.spines['left'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.set_yticklabels([])
ax.set_yticks([])

ax.set_xticks(x_axis)
plt.tight_layout()
plt.show()
```

---