

Contents

1	How to Count	3
1.1	Basic Counting Principles	3
1.2	Permutations	5
1.2.1	Permutations with Identical Objects	6
1.3	Combinations	6
1.4	Binomial and Multinomial Coefficients	7
2	Axioms of Probability	8
2.1	Sample Space and Events	8
2.2	Probability	8
2.3	Inclusion-Exclusion Principle	10
3	Conditional Probability	12
3.1	Conditional Probability	12
3.2	Bayes's Formula	13
3.3	Indepedent Events	14
4	Random Variables	16
4.1	Random Variables	16
4.2	Discrete Random Variables	16
4.2.1	Expectation of Discrete Random Variables	17
4.2.2	Variance	19
4.2.3	Bernoulli and Binomial Random Variables	21
4.2.4	Poisson Random Variable	23
4.2.5	Geometric Random Variable	25
4.2.6	Negative Binomial Random Variable	25
4.2.7	Hypergeometric Random Variable	26
4.3	Continuous Random Variables	27
4.3.1	Uniform Random Variable	28
4.3.2	Normal Random Variable	29
4.3.3	Exponential Random Variable	30
4.3.4	Gamma Random Variable	32
4.3.5	Beta Random Variable	33
4.4	Jointly Distributed Random Variables	34

4.4.1	Independent Random Variables	36
4.4.2	Sums of Independent Random Variables	38
4.5	Conditional Distribution	39
5	Expectation	42
5.1	Sums of Random Variables	42
5.2	Moments	44
5.3	Covariance and Correlations	46
5.3.1	Covariance	46
5.3.2	Correlation	48
5.4	Conditional Expectation and Variance	49
5.4.1	Conditional Expectation	49
5.4.2	Conditional Variance	49
5.4.3	Prediction	50

How to Count

1.1 Basic Counting Principles

An important motivation to study combinatorics is to count the **number of ways** in which an event may occur. Intuitively, we have two approaches to count.

The first approach is to categorise the event into **non-overlapping cases**. This means that we break an event into mutually exclusive sub-events, after which we can count the number of ways for each sub-event to occur. The aggregate of these counts is the total number of ways for the original event to occur.

Those familiar with basic set theory may consider E to be the set containing all distinct ways for an event to occur. By breaking up the event, we essentially establish a **partition** of E , so that the sum of cardinalities of all the elements in that partition equals the cardinality of E .

This motivates us to write the following principle using set notations.

Theorem 1.1.1 ► Addition Principle (AP)

Let $k \in \mathbb{N}^+$ and let A_1, A_2, \dots, A_k be k finite sets which are pairwise disjoint, i.e. for all i, j such that $1 \leq i, j \leq k$, $A_i \cap A_j = \emptyset$ whenever $i \neq j$, then

$$\left| \bigcup_{i=1}^k A_i \right| = \sum_{i=1}^k |A_i|.$$

Proof. The case where $k = 1$ is trivial.

Suppose that when $k = n$, we have

$$\left| \bigcup_{i=1}^n A_i \right| = \sum_{i=1}^n |A_i|$$

for any n finite sets which are pairwise disjoint. Let A_{n+1} be an arbitrary finite set

which is disjoint with any of the A_i 's from the n sets. So we have:

$$\begin{aligned}
 \left| \bigcup_{i=1}^{n+1} A_i \right| &= \left| \left(\bigcup_{i=1}^n A_i \right) \cup A_{n+1} \right| \\
 &= \left| \bigcup_{i=1}^n A_i \right| + |A_{n+1}| - \left| \left(\bigcup_{i=1}^n A_i \right) \cap A_{n+1} \right| \\
 &= \left(\sum_{i=1}^n |A_i| \right) + |A_{n+1}| - |\emptyset| \\
 &= \sum_{i=1}^{n+1} |A_i|.
 \end{aligned}$$

Therefore, the original statement holds for all $k \in \mathbb{N}^+$. □

In more casual language, this means that if an event E_k has n_k distinct ways to occur, then there is $\sum_{i=1}^k n_k$ ways for at least one of the events E_1, E_2, \dots, E_k to occur, provided that E_i and E_j can never occur concurrently whenever $i \neq j$.

Given an event E , the other approach to count the number of ways for it to occur is to break E up internally into **non-overlapping stages**.

With set notations, we can write the i -th stage for E to occur as e_i , and so a way for E to occur can be represented by an ordered tuple (e_1, e_2, \dots, e_k) , where k is the total number of stages to undergo for E to occur.

Let E_i denote the set of all distinct ways to undergo the i -th stage of E , then it is easy to see that E is just the **Cartesian product** of all the E_i 's. Hence, we derive the following principle:

Theorem 1.1.2 ► Multiplication Principle (MP)

Let $k \in \mathbb{N}^+$ and let A_1, A_2, \dots, A_k be k pairwise disjoint finite sets, then

$$\left| \prod_{i=1}^k A_i \right| = \prod_{i=1}^k |A_i|.$$

Proof. The case where $k = 1$ is trivial.

Suppose that when $k = n$, we have

$$\left| \prod_{i=1}^n A_i \right| = \prod_{i=1}^n |A_i|$$

for any n finite sets which are pairwise disjoint. Let A_{n+1} be an arbitrary finite set which is disjoint with any of the A_i 's from the n sets. Take $a_i, a_j \in A_{n+1}$. Note that for all $\mathbf{a} \in \prod_{i=1}^n A_i$, $(\mathbf{a}, a_i) \neq (\mathbf{a}, a_j)$ whenever $a_i \neq a_j$. This means that

$$\begin{aligned} \left| \prod_{i=1}^{n+1} A_i \right| &= \left| \prod_{i=1}^n A_i \times A_{n+1} \right| \\ &= \left| \prod_{i=1}^n A_i \right| |A_{n+1}| \\ &= \left(\prod_{i=1}^n |A_i| \right) |A_{n+1}| \\ &= \prod_{i=1}^{n+1} |A_i| \end{aligned}$$

Therefore, the original statement holds for all $k \in \mathbb{N}^+$. □

In more casual language, this means that if an event E requires k stages to be undergone before it occurs and the i -th stage has n_i ways to complete, then there is $\prod_{i=1}^k n_i$ ways for E to occur, provided that no two different stages complete concurrently.

1.2 Permutations

A fundamental problem in combinatorics is described as follows: given a set S , how many ways are there to arrange r elements in S , i.e. how many **distinct sequences** can be formed using the elements in S without repetition? The process of selecting elements from S and arranging them as a sequence is known as *permutation*.

Note that forming a sequence using r elements from a set S is an event consisting of r stages, as we need to select an element for each of the r terms of the sequence. Suppose S has n elements. For the first term of the sequence, we can choose any of the elements in S , so there is n ways to do it. For the second term, since we cannot repeat the elements, we are left with $(n - 1)$ choices.

Continue choosing elements in this way, we realise that if we choose the terms sequentially, when we reach the k -th term we will be left with $n - k + 1$ options as the previous $(k - 1)$ terms have taken away $(k - 1)$ elements. By Theorem 1.1.2, we know that the number of sequences which can be formed is given by $\prod_{i=1}^r (n - r + i)$.

Definition 1.2.1 ▶ Permutations

Let A be a finite set such that $|A| = n$, an r -permutation of A is a way to arrange r elements of A , denoted as P_r^n and given by

$$P_r^n = \prod_{i=1}^r (n - r + i) = \frac{n!}{(n - r)!}.$$

1.2.1 Permutations with Identical Objects**Theorem 1.2.2 ▶ Generalised Formula for Permutations**

Let $k \in \mathbb{N}^+$ and let A_1, A_2, \dots, A_k be k distinct objects, where A_i occurs $n_i > 0$ times for $i = 1, 2, \dots, k$, then the number of permutations for these k objects are given by

$$\frac{\left(\sum_{i=1}^k n_i\right)!}{\prod_{i=1}^k (n_i)!}.$$

1.3 Combinations**Definition 1.3.1 ▶ Combinations**

Let A be a finite set such that $|A| = n$, an r -combination of A is a way to choose r elements from A regardless of the order of selection, denoted as C_r^n and given by

$$C_r^n = \frac{P_r^n}{P_r^r} = \frac{n!}{r!(n - r)!} = \binom{n}{r}.$$

Remark. Two obvious results:

1. If $r > n$ or $r < 0$, $C_r^n = 0$;
2. $C_r^n = C_{n-r}^n$.

Theorem 1.3.2 ▶ Pascal's Triangle

Let n be an integer with $n \geq 2$ and let r be an integer with $0 \leq r \leq n$, then

$$C_r^n = C_{r-1}^{n-1} + C_r^{n-1}.$$

1.4 Binomial and Multinomial Coefficients

Consider the expansion of $(x + y)^n$ where $n \in \mathbb{N}$. Note that this expansion is a linear combination of terms in the form of $x^k y^{n-k}$ where $k = 0, 1, 2, \dots, n$.

Thus, fix any k , to determine how many copies of $x^k y^{n-k}$ there are, it suffices to compute C_k^n . Therefore, in the expanded form of $(x + y)^n$, the coefficient is exactly C_r^n .

Theorem 1.4.1 ► Binomial Expansion

Let $n \in \mathbb{N}$, then

$$(x + y)^n = \sum_{k=0}^n \left[\binom{n}{k} x^k y^{n-k} \right].$$

We can extend the idea of binomial coefficients onto multinomial expansions, i.e. expressions in the form of $(\sum_{i=1}^r x_i)^n$.

Note that the binomial coefficient C_r^n is essentially equivalent to dividing n distinct elements into two groups with r and $(n-r)$ members respectively. Now we consider dividing n distinct elements into r groups with n_1, n_2, \dots, n_r members respectively for each group.

Notice that we can simply permute the n distinct elements and assign them sequentially into the r groups, i.e. the first n_1 elements will go into the first group and so on.

Since the order of elements within each group does not matter, we need to remove repeated selections by dividing by $\prod_{i=1}^r (n_i!)$. So we have the following definition:

Definition 1.4.2 ► Multinomial Coefficients

The **multinomial coefficient** is defined by

$$\binom{n}{n_1, n_2, \dots, n_k} = \frac{n!}{\prod_{i=1}^k (n_i!)}$$

Theorem 1.4.3 ► Multinomial Expansion

Let $n \in \mathbb{N}$, then

$$\left(\sum_{i=1}^r x_i \right)^n = \sum_{\substack{n_1, n_2, \dots, n_r \in \mathbb{N} \\ \sum_{j=1}^r n_j = n}} \left[\binom{n}{n_1, n_2, \dots, n_r} \prod_{i=1}^r x_i^{n_i} \right]$$

Axioms of Probability

2.1 Sample Space and Events

Definition 2.1.1 ► Sample Space

Consider an experiment whose outcome is **not** predictable, then the set of all possible outcomes of the experiment is called the **sample space** of the experiment, denoted by S .

Remark. Note that $S \neq \emptyset$.

Definition 2.1.2 ► Events

Let S be a sample space, a set $E \subseteq S$ is known as an **event**.

Remark. S itself is known as the **sure event** and \emptyset is known as the **null event**.

Note that since sample spaces and events are sets, we can apply operations onto events precisely in the same way for sets.

By convention, the intersection of two events E and F is preferably written as EF . Two events which are disjoint are called *mutually exclusive*.

2.2 Probability

Definition 2.2.1 ► Probability

Let E be any event of an experiment and let $n(E)$ be the number of occurrences of E in the first n repetitions of the experiment, then the **probability** of E is

$$P(E) = \lim_{n \rightarrow \infty} \frac{n(E)}{n},$$

if the limit exists.

However, notice that from the above, the notion of probability may not be well-defined as $n(E)$ is not a function, which means that the limit is not defined.

To avoid this problem, we shall use an axiomatic definition instead, i.e., we define probability to be such that if it exists and is well-defined, then it satisfies a series of axioms.

Definition 2.2.2 ► Axioms of Probability

Let S be a sample space and let $P(E)$ be a real number defined for every $E \subseteq S$. If

- $0 \leq P(E) \leq 1$,
- $P(S) = 1$, and
- for all mutually exclusive E and F , $P(E \cup F) = P(E) + P(F)$,

then $P(E)$ is the **probability** of E .

With induction, one can easily show that if E_1, E_2, \dots to be any sequence of events in a sample space S , then

$$P\left(\bigcup_{i=1}^{\infty} E_i\right) = \sum_{i=1}^{\infty} P(E_i).$$

We now follow up with proofs for two seemingly intuitive results.

Theorem 2.2.3 ► The Null Event

Consider the null event \emptyset , we have

$$P(\emptyset) = 0.$$

Proof. Let S be a sample space and let E_1, E_2, \dots be a countably infinite sequence of events such that $E_i = \emptyset$ for all $i \in \mathbb{N}^+$. We can write

$$P\left(\bigcup_{i=1}^{\infty} E_i\right) = \sum_{i=1}^{\infty} P(E_i).$$

Note that the countable union of empty sets is empty, so the above is equivalent to

$$P\left(\bigcup_{i=1}^{\infty} \emptyset\right) = P(\emptyset) = \sum_{i=1}^{\infty} P(\emptyset).$$

This means that $P(\emptyset)$ equals the sum of a countably infinite sequence of itself, so

$$P(\emptyset) = 0.$$

□

Theorem 2.2.4 ► Monotonicity of Probability

Let E and F be events such that $E \subseteq F$, then

$$P(F) \geq P(E).$$

Proof. Note that E and $F - E$ are mutually exclusive, so

$$P(F) = P(E \cup (F - E)) = P(E) + P(F - E).$$

Note that $P(F - E) \geq 0$, so $P(E) + P(F - E) \geq P(E)$, which means

$$P(F) \geq P(E).$$

□

2.3 Inclusion-Exclusion Principle

It is easy to compute the probability of a countable union of mutually exclusive events. However, it may get tricky when an event is the union of events which are not mutually exclusive. Intuitively, we can sum up the probabilities of all individual events and subtract the portions which are double-counted. This approach is rigorously summarised as follows:

Theorem 2.3.1 ► Inclusion-Exclusion Principle

Let S be a sample space and let E_1, E_2, \dots, E_n be a sequence of events. In general, we have

$$P\left(\bigcup_{i=1}^n E_i\right) = \sum_{j=1}^n \left[(-1)^{j+1} \left(\sum_{k_1 \leq k_2 \leq \dots \leq k_j} P\left(\bigcap_{h=1}^j E_{k_h}\right) \right) \right].$$

Proof. Define a function $f_S : S \rightarrow \{0, 1\}$ by

$$f_S(x) = \begin{cases} 1 & \text{if } x \in S \\ 0 & \text{if } x \notin S \end{cases}.$$

Let $E = \bigcup_{i=1}^n E_i$. Consider the function $g : S \rightarrow \{0, 1\}$ given by

$$g(x) = \prod_{i=1}^n (f_E(x) - f_{E_i}(x)).$$

For any $x \in S$, if $x \in E$, then $x \in E_k$ for some $k \in \{x \in \mathbb{N} : x \leq n\}$, which means that $f_E(x) - f_{E_k}(x) = 0$; if $x \notin E$, then $f_E(x) = f_{E_i}(x) = 0$ for all $i \in \{x \in \mathbb{N} : x \leq n\}$. In either case, $g(x) = 0$. \square

Theorem 2.3.2 ▶ Boole's Inequality

Let $E_1, E_2, \dots, E_n, \dots$ be a countable sequence of events, then

$$P\left(\bigcup_{i=1}^{\infty} E_i\right) \leq \sum_{i=1}^{\infty} P(E_i).$$

In particular, equality is achieved if and only if the E_i 's are mutually exclusive.

Theorem 2.3.3 ▶ Probability in a Finite Sample Space

Let S be a sample space which is finite and let $E \subseteq S$ be an event, then

$$P(E) = \frac{|E|}{|S|}.$$

Conditional Probability

3.1 Conditional Probability

Given a sample space S , we may wish to find the probability of two events E and F both occurring, $P(EF)$. However, suppose that we already know that event F **has occurred**, then necessarily, the sample space we consider would no longer be S . Essentially, this condition of F having occurred has restricted our sample space to F . Thus, we give the following definition:

Definition 3.1.1 ► Conditional Probability

Let S be a sample space and $E, F \subseteq S$ be two events. If $P(F) \geq 0$, then the **conditional probability** is the probability that E occurs given that F has occurred, denoted by

$$P(E|F) = \frac{P(EF)}{P(F)}.$$

In particular, if $E \subseteq F$, we have $P(E|F) = \frac{P(E)}{P(F)}$.

Remark. Note that $P(E|F) = P(EF|F)$.

It is easy to see that $P(EF) = P(E|F)P(F)$, i.e., the probability of E and F both occurring is the product of the probability of F occurring and the probability of E occurring given the occurrence of F . This complies with our intuition. We can generalise this for a countable number of events:

Proposition 3.1.2 ► Multiplication Rule

Let S be a sample space and let $E_i \subseteq S$ for $i = 1, 2, \dots, n$ be n events, where $n \geq 2$. Suppose that $P\left(\bigcap_{i=1}^{n-1} E_i\right) > 0$, then

$$P\left(\bigcap_{i=1}^n E_i\right) = P(E_1) \prod_{i=2}^n P\left(E_i \left| \bigcap_{j=1}^{i-1} E_j \right.\right).$$

Proof. The case where $n = 2$ is immediate from Definition 3.1.1.

Suppose that there is some $k \in \mathbb{N}$ and $k \geq 2$ such that

$$P\left(\bigcap_{i=1}^k E_i\right) = P(E_1) \prod_{i=2}^k P\left(E_i \left| \bigcap_{j=1}^{i-1} E_j\right.\right),$$

then we consider

$$\begin{aligned} P\left(E_{k+1} \left| \bigcap_{i=1}^k E_i\right.\right) &= \frac{P\left(\bigcap_{i=1}^{k+1} E_i\right)}{P\left(\bigcap_{i=1}^k E_i\right)} \\ &= \frac{P\left(\bigcap_{i=1}^{k+1} E_i\right)}{P(E_1) \prod_{i=2}^k P\left(E_i \left| \bigcap_{j=1}^{i-1} E_j\right.\right)}. \end{aligned}$$

Therefore,

$$\begin{aligned} P\left(\bigcap_{i=1}^{k+1} E_i\right) &= \left[P(E_1) \prod_{i=2}^k P\left(E_i \left| \bigcap_{j=1}^{i-1} E_j\right.\right) \right] P\left(E_{k+1} \left| \bigcap_{i=1}^k E_i\right.\right) \\ &= P(E_1) \prod_{i=2}^{k+1} P\left(E_i \left| \bigcap_{j=1}^{i-1} E_j\right.\right) \end{aligned}$$

□

3.2 Bayes's Formula

Consider a sample space S and two events $E, F \subseteq S$. Suppose that E occurs, then either F has occurred or F has never occurred (i.e. F^c occurred). Therefore, it is easy to see that

$$P(E) = P(EF) + P(EF^c) = P(E|F) + P(E|F^c).$$

We can extend the above argument for more than two events. Suppose that F_1, F_2, \dots, F_n are n mutually exclusive events such that $\bigcup_{i=1}^n F_i = S$, then obviously $\{F_1, F_2, \dots, F_n\}$ is a *partition* of S .

Consider any event E and let $e \in E$. Clearly, e must be in one and only one of F_1, F_2, \dots, F_n . It then follows that $\{E \cap F_1, E \cap F_2, \dots, E \cap F_n\}$ is a partition of E . Generalising this further to a countably infinite number of mutually exclusive events F_1, F_2, \dots such that $\bigcup_{i=1}^{\infty} F_i = S$,

we arrive at the following formula:

$$P(E) = \sum_{i=1}^{\infty} P(E|F_i)P(F_i).$$

This leads to the *Bayes's Formula*:

Theorem 3.2.1 ► Bayes's Formula

Let F_1, F_2, \dots be a countably infinite sequence of events over a sample space S such that $\bigcup_{i=1}^{\infty} F_i = S$. For any event $E \subseteq S$, we have

$$P(F_j|E) = \frac{P(E|F_j)P(F_j)}{\sum_{i=1}^{\infty} P(E|F_i)P(F_i)}.$$

3.3 Independent Events

Note that in general, for two events E and F , $P(E|F) \neq P(E)$, i.e., the occurrence of F may affect the occurrence of E . However, in some cases, we notice that the occurrence of E is *independent* of F , and so we introduce the following definition:

Definition 3.3.1 ► Independent Events

Let S be a sample space and let $E, F \subseteq S$ be two events. We say that E and F are **independent** if $P(EF) = P(E)P(F)$, and **dependent** otherwise.

Remark. The following results are immediate:

1. If $P(E) = 0$ or $P(F) = 0$, then E and F are independent.
2. If $P(E) > 0$ (respectively, $P(F) > 0$), then E and F are independent if and only if $P(F|E) = P(F)$ (respectively, $P(E|F) = P(E)$).

Intuitively, given independent events E and F , we may believe that if the occurrence of E does not affect the occurrence of F , then naturally the occurrence of E should also not affect the “not-occurring” of F , i.e., the following is true:

Proposition 3.3.2

E and F are independent events if and only if E and F^c are independent events.

Proof. Since $F = (F^c)^c$, it suffices to prove for one direction.

Notice that $EF \cup EF^c = E(F \cup F^c) = E$, so

$$P(E) = P(EF) + P(EF^c) = P(E)P(F) + P(EF^c).$$

Therefore,

$$P(EF^c) = P(E) - P(E)P(F) = P(E)(1 - P(F)) = P(E)P(F^c),$$

and so E and F^c are independent. □

Random Variables

4.1 Random Variables

In many contexts, we might wish to generalise a formula to compute the probability of the occurrence of a certain event. However, in cases where the events are abstract or unquantifiable (e.g. the event “tomorrow is rainy”), it becomes hard to formulate a well-defined mapping from a sample space to $[0, 1]$. Thus, to model all events easily using functions and mappings, we introduce the notion of *random variables*.

Definition 4.1.1 ► Random Variable

Let Ω be a sample space, the **random variable**

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

is a real-valued function such that for any event $E \subseteq \Omega$,

$$P(E) = P(X[E]) = P(X \in X[E]),$$

where

$$X[E] := \{X(\omega) : \omega \in E\}$$

is the image of the event E under X .

4.2 Discrete Random Variables

Intuitively, there are certain events whose outcomes are finite or can be enumerated. In such cases, we may associate these events with a *discrete random variable*.

Definition 4.2.1 ► Discrete Random Variable (DRV)

Let X be a random variable over a sample space S , if $\text{ran}(X)$ is countable, then X is called a **discrete random variable**.

Definition 4.2.2 ▶ Probability Mass Function (PMF)

Let X be a discrete random variable over a sample space S , the function

$$p_X : X[S] \rightarrow [0, 1]$$

where $p_X(a) = P(X = a)$ is known as the **probability mass function** of X .

Remark. $\sum_{a \in X[S]} p_X(a) = 1$.

For any discrete random variable X , $p_X(a) = 0$ if and only if $\{s \in S : X(s) = a\} = \emptyset$. This essentially means that if p_X evaluates to 0, then the corresponding event is an impossible event.

Note that p_X essentially gives the probability of **singleton** events in a sample space. Naturally, we can represent the probability of a union of singleton events as a linear combination of the values of p_X .

Definition 4.2.3 ▶ Cumulative Distribution Function (CDF)

Let X be a discrete random variable over a sample space S with PMF p_X , the function

$$F_X : X[S] \rightarrow [0, 1]$$

where $F_X(a) = \sum_{x \leq a} p_X(x)$ is known as the **cumulative distribution function** of X .

Remark. Suppose that $X(s_i) = a_i$ for all $s_i \in S$ such that $a_i < a_j$ whenever $i < j$, then F_X is a non-decreasing **step function**, i.e., for all a such that $a_i \leq a < a_{i+1}$,

$$F_X(a) = \sum_{x \leq a} p_X(x) = \sum_{k=1}^i p_X(a_k).$$

4.2.1 Expectation of Discrete Random Variables

Suppose X is a discrete random variable with range $\{x_1, x_2, \dots, x_m\}$ and PMF p_X . By repeating an experiment of X for n times, we can approximate the total number of occurrences of $X = x_i$ by $np_X(x_i)$. Therefore, the average value of X can be approximated by

$$\frac{\sum_{i=1}^m nx_i p_X(x_i)}{n}.$$

For large n , we have

$$\lim_{n \rightarrow \infty} \frac{\sum_{i=1}^m nx_i p_X(x_i)x_i}{n} = \lim_{n \rightarrow \infty} \sum_{i=1}^m x_i p_X(x_i) = \sum_{i=1}^m x_i p_X(x_i).$$

Similarly, if the range of X is countably infinite, replacing from the above $\sum_{i=1}^m nx_i p_X(x_i)x_i$ with $\sum_{i=1}^{\infty} nx_i p_X(x_i)x_i$ will yield the same limit.

Intuitively, the above limit represents the *expected value* of X when a large number of experiments are conducted, which leads to the following definition:

Definition 4.2.4 ► Expectation of Discrete Random Variables

Let X be a discrete random variable. The **expectation** (or **mean**, **expected value**) of X is defined to be

$$E[X] = \sum_{i=1}^m [x_i p_X(x_i)] = \sum_{i=1}^m [x_i P(X = x_i)]$$

if $|\text{ran}(X)| = m$, and

$$E[X] = \sum_{i=1}^{\infty} [x_i p_X(x_i)] = \sum_{i=1}^{\infty} [x_i P(X = x_i)]$$

if $\text{ran}(X)$ is countably infinite.

By convention, we use μ to denote expectation, so $E[X]$ can be written as μ_X .

For a discrete random variable X , we can define a function $g : \text{ran}(X) \rightarrow \mathbb{R}$. It is easy to see that $g(X)$ is also a discrete random variable. Therefore, we may have the following result:

Theorem 4.2.5 ► Expectation of Functions

Let X be a discrete random variable and define $Y = g(X)$, then

$$E[Y] = E[g(X)] = \sum_x [g(x)P(X = x)].$$

Proof. Note that for each $y \in \text{ran}(Y)$, $g(x) = y$ for some $x \in \text{ran}(X)$, and so

$$P(Y = y) = \sum_{g(x)=y} P(X = x).$$

Therefore,

$$\begin{aligned}
 E[Y] &= \sum_y [yP(Y = y)] \\
 &= \sum_y \left[y \sum_{g(x)=y} P(X = x) \right] \\
 &= \sum_y \left[\sum_{g(x)=y} g(x)P(X = x) \right] \\
 &= \sum_x [g(x)P(X = x)].
 \end{aligned}$$

□

Two simple corollaries to the above theorem are:

$$\begin{aligned}
 E[aX + b] &= aE[X] + b \\
 E[X + Y] &= E[X] + E[Y].
 \end{aligned}$$

In later sections, we will prove that the same rule applies to continuous random variables as well. The above theorem gives rise to the following notion of *moments*:

Definition 4.2.6 ► Moment

Let X be a random variable. $E[X^n]$ is called the n -th **moment** of X .

Following Theorem 4.2.5, it is easy to see that if X is discrete, then

$$E[X^n] = \sum_{i=1}^{\infty} x_i^n p_X(x_i).$$

4.2.2 Variance

Note that given two different discrete random variables X and Y , their probability mass functions can be different but they can still have identical expectations. For example, consider p_X to be identically 0 and p_Y to be such that $p_Y(0) = 1$ and $p_Y(y) = 0$ for all $y \neq 0$.

This motivates us to find other properties to classify and characterise random variables. One of these properties is the **spread** of the possible values taken by a random variable with respect to its mean, i.e., consider the random variable X with $E[X] = \mu$, we wish to determine $E[|X - \mu|]$ or equivalently $E[(X - \mu)^2]$. This spread is known as the *variance* of a random variable.

Definition 4.2.7 ▶ Variance

Let X be a random variable with $E[X] = \mu$, the **variance** of X is defined to be

$$\text{Var}(X) = E[(X - \mu)^2] = E[X^2] - (E[X])^2.$$

By convention, we use σ^2 to denote variance, so $\text{Var}(X)$ can be written as σ_X^2 .

The formula for $\text{Var}(X)$ can be derived via Theorem 4.2.5:

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

Another term we hear often is *standard deviation*, which is defined as follows:

Definition 4.2.8 ▶ Standard Deviation

Let X be a random variable, the **standard deviation** of X is defined to be

$$\text{SD}(X) = \sqrt{\text{Var}(X)} = \sigma_X.$$

Note that we have computed the general formula for any linear combination of discrete random variables. We shall do the same for variance.

Proposition 4.2.9 ▶ Variance of Linear Combinations of Random Variables

Let X be a random variable, then

$$\begin{aligned}\text{Var}(aX + b) &= a^2 \text{Var}(X) \\ \text{SD}(aX + b) &= |a| \text{SD}(X)\end{aligned}$$

for all $a, b \in \mathbb{R}$.

Proof. By using Theorem 4.2.5, we have

$$\begin{aligned}\text{Var}(aX + b) &= E[(aX + b)^2] - (E[aX + b])^2 \\ &= a^2 E[X^2] + 2abE[X] + b^2 - [a(E[X])^2 + 2abE[X] + b^2] \\ &= a^2 [E[X^2] - (E[X])^2] \\ &= a^2 \text{Var}(X).\end{aligned}$$

Therefore,

$$\text{SD}(aX + b) = \sqrt{\text{Var}(aX + b)} = |a|\text{SD}(X).$$

□

4.2.3 Bernoulli and Binomial Random Variables

Suppose we conduct an experiment. In the most simplistic view, only two outcomes are considered, namely **success** and **failure**. We can model such experiments using a discrete random variable whose range has a cardinality of 2.

Definition 4.2.10 ► Bernoulli Random Variable

A random variable X is a **Bernoulli random variable** if

$$p_X(x) = \begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0 \\ 0, & \text{otherwise} \end{cases}$$

for some $p \in [0, 1]$.

Now, consider n **independent** trials of an experiment with a probability for success of p . Let X be the number of successes among these n trials, then clearly,

$$p_X(x) = P(X = x) = \binom{n}{x} p^x (1 - p)^{n-x}.$$

Definition 4.2.11 ► Binomial Random Variable

A random variable X is a **binomial random variable** if

$$p_X(x) = P(X = x) = \binom{n}{x} p^x (1 - p)^{n-x}$$

for some $p \in [0, 1]$. X is said to have a **binomial distribution** with parameters (n, p) , denoted by $X \sim B(n, p)$.

Remark. In particular, if X is a Bernoulli random variable, then $X \sim B(1, p)$.

Suppose $X \sim B(n, p)$. Let N be the average number of successes in the n trials, it is expected

that $p \approx \frac{N}{n}$. Therefore, we may conjecture that $N \approx np$.

Theorem 4.2.12 ► Expectation and Variance of Binomial Distribution

Let $X \sim B(n, p)$, then $E[X] = np$ and $\text{Var}(X) = np(1 - p)$.

Proof. Note that $iC_i^n = nC_{i-1}^{n-1}$, so

$$\begin{aligned} E[X] &= \sum_{i=0}^n \left[i \binom{n}{i} p^i (1-p)^{n-i} \right] \\ &= \sum_{i=1}^n \left[n \binom{n-1}{i-1} p^i (1-p)^{n-i} \right] \\ &= n \sum_{j=0}^{n-1} \left[\binom{n-1}{j} p^{j+1} (1-p)^{n-1-j} \right] \\ &= np \sum_{j=0}^{n-1} \left[\binom{n-1}{j} p^j (1-p)^{n-1-j} \right] \\ &= np, \end{aligned}$$

$$\begin{aligned} \text{Var}(X) &= E[X^2] - (E[X])^2 \\ &= \sum_{i=0}^n \left[i^2 \binom{n}{i} p^i (1-p)^{n-i} \right] - n^2 p^2 \\ &= n \sum_{i=1}^n \left[i \binom{n-1}{i-1} p^i (1-p)^{n-i} \right] - n^2 p^2 \\ &= n \sum_{j=0}^{n-1} \left[(j+1) \binom{n-1}{j} p^{j+1} (1-p)^{n-1-j} \right] - n^2 p^2 \\ &= np \left\{ \sum_{j=0}^{n-1} \left[j \binom{n-1}{j} p^j (1-p)^{n-1-j} \right] + 1 \right\} - n^2 p^2 \\ &= np[(n-1)p + 1] - n^2 p^2 \\ &= np - np^2 \\ &= np(1 - p). \end{aligned}$$

□

Let $X \sim B(n, p)$, consider

$$\begin{aligned}\frac{p_X(i+1)}{p_X(i)} &= \frac{\frac{n!}{(i+1)!(n-i-1)!} p^{i+1} (1-p)^{n-i-1}}{\frac{n!}{i!(n-i)!} p^i (1-p)^{n-i}} \\ &= \frac{\frac{1}{i+1} p}{\frac{1}{n-i} (1-p)} \\ &= \frac{(n-i)p}{(i+1)(1-p)}.\end{aligned}$$

Suppose $p_X(i+1) < p_X(i)$, then $(n-i)p < (i+1)(1-p)$. This implies that $i > (n+1)p - 1$, which means that

- $p_X(i)$ is monotonically increasing on $[0, (n+1)p - 1]$.
- $p_X(i)$ maximises when $i = \lceil (n+1)p - 1 \rceil = \lfloor (n+1)p \rfloor$.
- $p_X(i)$ is monotonically decreasing on $((n+1)p - 1, n]$.

4.2.4 Poisson Random Variable

Suppose that $X \sim B(n, p)$ such that n is large and p is small. Let $\lambda = np$, then

$$\begin{aligned}p_X(i) &= \binom{n}{i} p^i (1-p)^{n-i} \\ &= \frac{\prod_{j=0}^{i-1} (n-j)}{i!} \left(\frac{\lambda}{n}\right)^i \left(1 - \frac{\lambda}{n}\right)^{n-i} \\ &= \frac{\prod_{j=0}^{i-1} (n-j)}{n^i} \cdot \frac{\lambda^i}{i!} \cdot \frac{\left(1 - \frac{\lambda}{n}\right)^n}{\left(1 - \frac{\lambda}{n}\right)^i}\end{aligned}$$

Therefore,

$$\lim_{n \rightarrow \infty} p_X(i) = e^{-\lambda} \frac{\lambda^i}{i!}.$$

Note that this means that we can use $e^{-\lambda} \frac{\lambda^i}{i!}$ as a good approximation for $p_X(i)$ when n is large and p is small! In this case, λ is the expected frequency of occurrences of the event corresponding to $X = 1$ within a unit interval.

Definition 4.2.13 ► Poisson Random Variable

A random variable X is a **Poisson random variable** if

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

for some $\lambda > 0$, denoted as $X \sim \text{Po}(\lambda)$.

Note that for $X \sim \text{Po}(\lambda)$, we can find some $Y \sim B(n, p)$ where n is large and p is small such that $np = \lambda$. Therefore, it is expected that

$$\begin{aligned} E[X] &\approx E[Y] = \lambda, \\ \text{Var}(X) &\approx \text{Var}(Y) = np(1 - p) \approx \lambda. \end{aligned}$$

Theorem 4.2.14 ► Expectation and Variance of Poisson Random Variables

If $X \sim \text{Po}(\lambda)$ where $\lambda > 0$, then $E[X] = \text{Var}(X) = \lambda$.

Definition 4.2.15 ► Weakly Dependent

Let E and F be two events. If $P(E) \approx P(E \mid F)$, we say that E and F are **weakly dependent**.

Let $i = 1, 2, 3, \dots, n$ and p_i be the probability of event i occurring. If the i 's are independent or weakly dependent, then we can approximate for large n that the rate of occurrences of these events is $\sum_{i=1}^n p_i$. Let X be the number of events which occur, then

$$X \sim \text{Po}\left(\sum_{i=1}^n p_i\right).$$

Theorem 4.2.16 ► Poisson Process

Let E be an event which occurs randomly. Assume that

1. there are λ occurrences per unit interval;
2. no two occurrences happen at the same point;
3. numbers of occurrences in disjoint intervals are independent.

Let $N(t)$ be the number of occurrences of E in an interval of length t , then $N(t) \sim \text{Po}(\lambda t)$.

4.2.5 Geometric Random Variable

Suppose we perform some experiment with a probability of success of p . Let X be the number of failures before the first success occurs, then clearly,

$$P(X = x) = (1 - p)^x p.$$

Additionally, let Y be the number of trials needed to reach the first success, then

$$P(Y = y) = (1 - p)^{y-1} p.$$

Note that both $(P(X = x))$ and $(P(Y = y))$ form geometric sequences.

Definition 4.2.17 ► Geometric Random Variable

A random variable X is called a **geometric random variable** with parameter $p \in (0, 1)$, denoted by $X \sim \text{Geo}(p)$, if

$$p_X(n) = (1 - p)^{n-1} p.$$

Theorem 4.2.18 ► Expectation and Variance of Geometric Random Variables

If $X \sim \text{Geo}(p)$, then $E[X] = \frac{1}{p}$ and $\text{Var}(X) = \frac{1-p}{p^2}$.

4.2.6 Negative Binomial Random Variable

negBinDRV Suppose we perform some experiment with a probability of success of p . Let X be the number of trials needed to achieve the r -th success, then clearly, for $X = n$, we need $(r - 1)$ successes (i.e., $(n - r)$ failures) in the first $(n - 1)$ trials and the r -th trial to be a success. Therefore,

$$P(X = n) = C_{r-1}^{n-1} p^{r-1} (1 - p)^{n-r} p = C_{r-1}^{n-1} p^r (1 - p)^{n-r}.$$

Definition 4.2.19 ► Negative Binomial Random Variable

A random variable X is called a **negative binomial random variable** if

$$p_X(n) = \binom{n-1}{r-1} p^r (1 - p)^{n-r},$$

where $0 < p < 1$ and $n \geq r$, denoted as $X \sim \text{NB}(r, p)$.

Theorem 4.2.20 ► Expectation and Variance of Negative Binomial Variables

Let $X \sim \text{NB}(r, p)$, then

$$E[X] = \frac{r}{p}, \quad \text{Var}(X) = \frac{r(1-p)}{p^2}.$$

4.2.7 Hypergeometric Random Variable

Suppose a collection contains N objects, m of which are of type A. If n objects are selected randomly without replacement and let X be the number of objects of type A selected, then

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

Definition 4.2.21 ► Hypergeometric Random Variable

A random variable X is called a **hypergeometric random variable** with parameters (n, N, m) if

$$p_X(x) = \frac{\binom{m}{x} \binom{N-m}{n-x}}{\binom{N}{n}},$$

where $0 \leq m, n \leq N$.

Theorem 4.2.22 ► Expectation and Variance of Hypergeometric Random Variables

Let X be a hypergeometric random variable with parameters (n, N, m) , then

$$E[X] = np, \quad \text{Var}(X) = np(1-p) \left(1 - \frac{n-1}{N-1}\right).$$

4.3 Continuous Random Variables

In real life, the outcomes of certain events are infinitely many, and so they cannot be enumerated as discrete cases. Thus, we will need to use *continuous random variables* to model these events.

Definition 4.3.1 ► Continuous Random Variable

A random variable X is a **continuous random variable** if there exists some non-negative function f_X such that for all $B \subseteq \mathbb{R}$,

$$P(X \in B) = \int_B f_X(x) dx.$$

The function f_X is known as the **probability density function** of X . The function F_X with $0 \leq F_X(x) \leq 1$ and $F'_X(x) = f_X(x)$ is known as the **cumulative distribution function** of X .

An interesting property of a continuous random variable X is that

$$P(X = x) = \int_x^x f_X(x) dx = 0,$$

which means that the probability of any single outcome of an event is 0, but this does not mean that it is impossible to occur! In particular, it is more meaningful to consider the probability of the occurrence of a range of outcomes. We have

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

$$P(X \leq a) = P(X < a) = \int_{-\infty}^a f_X(x) dx.$$

It is not surprising that a function of a continuous random variable is still a continuous random variable.

Theorem 4.3.2 ► Function of Continuous Random Variables

Let X be a continuous random variable and $Y = g(X)$. If g is strictly monotonic and differentiable, then

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|.$$

Let X be a continuous random variable with probability density function f_X such that $f_X(x) = 0$ for all $x \in \mathbb{R} - [a, b]$. We divide $[a, b]$ into n intervals $[x_{i-1}, x_i]$ for $i = 1, 2, 3, \dots, n$

with equal length $\Delta x = \frac{b-a}{n}$. Thus,

$$P(x_{i-1} < X < x_i) \approx \Delta x f_X(x_i).$$

Let Y be a discrete random variable with $P(Y = x_i) = \Delta x f_X(x_i)$, then

$$E[X] \approx E[Y] = \sum_{i=1}^n x_i \Delta x f_X(x_i).$$

When $n \rightarrow \infty$, i.e., $\Delta x \rightarrow 0$, we have

$$\lim_{\Delta x \rightarrow 0} E[Y] = \int_a^b x f_X(x) dx.$$

By letting $a \rightarrow -\infty$ and $b \rightarrow \infty$, we have arrived at the following definition:

Definition 4.3.3 ► Expectation of Continuous Random Variables

Let X be a continuous random variable with probability density function f_X , then

$$E[X] = \int_{-\infty}^{\infty} x f_X(x) dx.$$

Let Y be a continuous random variable with probability density function f , consider

$$\begin{aligned} \int_0^{\infty} P(Y > y) dy &= \int_0^{\infty} \int_y^{\infty} f(x) dx dy \\ &= \int_0^{\infty} \int_0^x f(x) dy dx \\ &= \int_0^{\infty} x f(x) dy dx \\ &= E[Y]. \end{aligned}$$

Therefore, set $Y = g(X)$, then similar to discrete random variables, we have

$$E[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx.$$

4.3.1 Uniform Random Variable

Intuitively, we may call a random variable X “uniformly” distributed in (a, b) if $P(X = x)$ is a constant for all $x \in (a, b)$.

Definition 4.3.4 ► Uniform Random Variable

A continuous random variable X is a **uniform random variable** if

$$f_X(x) = \begin{cases} \frac{1}{b-a}, & \text{if } a < x < b \\ 0, & \text{otherwise} \end{cases},$$

denoted by $X \sim U(a, b)$.

Let $X \sim U(a, b)$, then the cumulative density function is

$$F_X(x) = \begin{cases} 0, & \text{if } x \leq a \\ \frac{x-a}{b-a}, & \text{if } a < x < b \\ 1, & \text{otherwise} \end{cases}.$$

Theorem 4.3.5 ► Expectation and Variance of Uniform Random Variables

Let $X \sim U(a, b)$, then $E[X] = \frac{a+b}{2}$ and $\text{Var}(X) = \frac{(b-a)^2}{12}$.

4.3.2 Normal Random Variable**Definition 4.3.6 ► Normal Random Variable**

A continuous random variable Z with probability density function ϕ is a **normal random variable** if

$$\phi(z) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(z-\mu)^2}{2\sigma^2}},$$

denoted as $Z \sim \mathcal{N}(\mu, \sigma^2)$.

In particular, $Z \sim \mathcal{N}(0, 1)$ is known as the *standard normal random variable*. Let Φ be the cumulative density function for Z , then

$$\Phi(z) = P(Z < z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{t^2}{2}} dt.$$

Let $X \sim \mathcal{N}(\mu, \sigma^2)$, then $Z = \frac{X-\mu}{\sigma}$, so

$$\Phi_X(x) = P(X < x) = P\left(\frac{X-\mu}{\sigma} < \frac{x-\mu}{\sigma}\right) = \Phi\left(\frac{x-\mu}{\sigma}\right)$$

Theorem 4.3.7 ▶ Expectation and Variance of Normal Random Variables

Let $Z \sim \mathcal{N}(\mu, \sigma^2)$, then $E[Z] = \mu$ and $\text{Var}(Z) = \sigma^2$.

4.3.3 Exponential Random Variable

Let $N(t) \sim \text{Po}(t\lambda)$ be the number of occurrences of an event in an interval of length t . Suppose X is the time before the first occurrence of the event, then

$$P(X > t) = P(N(t) = 0) = e^{-\lambda t}.$$

In other words, if F_X is the cumulative distribution function of X , then $F_X(x) = 1 - e^{-\lambda x}$.

Definition 4.3.8 ▶ Exponential Random Variable

A continuous random variable X is an **exponential random variable** if

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x}, & \text{if } x \geq 0 \\ 0, & \text{otherwise} \end{cases},$$

where $\lambda > 0$, denoted as $X \sim \text{Exp}(\lambda)$.

Theorem 4.3.9 ▶ Expectation and Variance of Exponential Random Variables

Let $X \sim \text{Exp}(\lambda)$, then $E[X] = \frac{1}{\lambda}$ and $\text{Var}(X) = \frac{1}{\lambda^2}$.

Informally, an exponential random variable models the **waiting time** before an event occurs. Suppose we have already waited for s unit of time for the occurrence, we may wish to know the probability of us having to wait for another t unit of time. To solve such questions, we need to understand the *memoryless* property.

Definition 4.3.10 ▶ Memoryless Property

Let X be a random variable, we say that X is **memoryless** if

$$P(X > s + t \mid X > t) = P(X > s).$$

In particular, if $X \sim \text{Exp}(\lambda)$, consider

$$\begin{aligned} P(X > s + t \mid X > t) &= \frac{e^{-\lambda(s+t)}}{e^{-\lambda t}} \\ &= e^{-\lambda s} \\ &= P(X > s). \end{aligned}$$

Therefore, exponential random variables are memoryless. One may also prove that geometric random variables are also memoryless.

Now we introduce another random variable which is closely related to the exponential random variable.

Definition 4.3.11 ► Double Exponential Random Variable

A continuous random variable X is a **double exponential variable** if

$$f_X(x) = \frac{1}{2}\lambda x^{-\lambda|x|}$$

for some $\lambda > 0$.

Consider $Y = |X|$ where X is a double exponential random variable. The double exponential random variable is so named because for any $y \geq 0$,

$$\begin{aligned} P(Y > y) &= P(X > y) + P(X < -y) \\ &= 2P(X > y) \\ &= 2 \int_y^\infty \frac{1}{2}\lambda e^{-\lambda x} dx \\ &= e^{-\lambda y}. \end{aligned}$$

Thus, $Y = |X| \sim \text{Exp}(\lambda)$.

A common application of exponential random variables is to determine the *hazard rate*. Suppose X is the survival time of some object and that the object has already survived for a time t . Consider ϵ to be a small interval, then the probability that the object cannot survive past this small interval is approximately

$$\begin{aligned} P(X < t + \epsilon \mid X > t) &= \frac{P(t < X < t + \epsilon)}{P(X > t)} \\ &\approx \frac{\epsilon f_X(t)}{1 - F_X(t)}. \end{aligned}$$

In general, we have the following definition:

Definition 4.3.12 ► Hazard Rate Function

Let X be a positive continuous random variable and define $\overline{F_X}(x) = 1 - F_X(x)$, then the function

$$\lambda(x) = \frac{f_X(x)}{\overline{F_X}(x)}$$

is known as the **hazard rate function** of X .

In particular, if $X \sim \text{Exp}(\lambda)$, then its hazard rate function is just $\lambda(x) = \lambda$, which is also known as the *rate* of X .

4.3.4 Gamma Random Variable

We have seen that the exponential random variable can be used to model the waiting time between two consecutive occurrences of an event modelled by a Poisson random variable. We would also like to know the waiting time till the n -th occurrence of an event.

Let $N(t)$ be a Poisson process with rate λ , so $P(N(t) = n) = \frac{e^{-\lambda t}(\lambda t)^n}{n!}$. Let T_n be the waiting time till the n -th event with f_k being the probability density function, then

$$\begin{aligned} f_k(t) &= \frac{d}{dt}(1 - P(T_n > t)) \\ &= \frac{d}{dt}(1 - P(N(t) < k)) \\ &= \frac{d}{dt}\left(1 - e^{-\lambda t} \sum_{i=0}^{k-1} \frac{(\lambda t)^i}{i!}\right) \\ &= \frac{\lambda e^{-\lambda t} (\lambda t)^{k-1}}{(k-1)!}. \end{aligned}$$

Definition 4.3.13 ► Gamma Random Variable

A continuous random variable X is called a **gamma** random variable with parameters (n, λ) where $\lambda > 0$ if

$$f(t) = \frac{\lambda e^{-\lambda t} (\lambda t)^{n-1}}{(n-1)!}$$

for all $t \geq 0$.

Definition 4.3.14 ▶ Gamma Function

Let $\alpha > 0$, the **gamma function** is defined as

$$\begin{aligned}\Gamma(\alpha) &= \int_0^{\infty} \lambda e^{-\lambda t} (\lambda t)^{\alpha-1} dt \\ &= \int_0^{\infty} \lambda e^{-x} x^{\alpha-1} dx.\end{aligned}$$

In particular, if X is a continuous random variable with probability density function

$$f(t) = \frac{\lambda e^{-\lambda t} (\lambda t)^{\alpha-1}}{\Gamma(\alpha)},$$

where $\lambda > 0$ and $t \geq 0$, then X is a gamma random variable with parameters (α, λ) .

Remark. One can prove that

1. $\Gamma(\alpha + 1) = \alpha \Gamma(\alpha)$ for all $\alpha > 0$.
2. $\Gamma(n) = (n - 1)!$ for all $n \in \mathbb{Z}^+$.

Theorem 4.3.15 ▶ Expectation and Variance of Gamma Random Variables

Let X be a gamma random variable with parameters (α, λ) , then

$$E[X] = \frac{\alpha}{\lambda}, \quad \text{Var}(X) = \frac{\alpha}{\lambda^2}.$$

4.3.5 Beta Random Variable

In real life, sometimes we may not know the exact probability distribution of a random variable and wish to deduce its probability distribution based on experiments. In other words, suppose there are a successes and b failures of an experiment, we wish to know what is the **most likely** probability of success.

Definition 4.3.16 ▶ Beta Random Variable

A continuous random variable X is a **beta** random variable with parameters (a, b) with $a, b > 0$ if

$$f_X(x) = \frac{1}{B(a, b)} x^{a-1} (1-x)^{b-1},$$

where

$$B(a, b) = \int_0^1 x^{a-1} (1-x)^{b-1} dx = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}.$$

If $a = b = 1$, we see that

$$f_X(x) = \frac{\Gamma(2)}{(\Gamma(1))^2},$$

so X is uniform on $(0, 1)$. In particular, we also see that

$$\frac{B(a+1, b)}{B(a, b)} = \frac{a}{a+b}.$$

Theorem 4.3.17 ► Expectation and Variance of Beta Random Variables

Let X be a beta random variable with parameters (a, b) , then

$$E[X] = \frac{a}{a+b}, \quad \text{Var}(X) = \frac{ab}{(a+b)^2(a+b+1)}.$$

4.4 Jointly Distributed Random Variables

Sometimes, the outcomes of the events we wish to study cannot be expressed with a single random variable. In general, if X_1, X_2, \dots, X_n are random variables, we may be interested to know

$$P((X_1, X_2, \dots, X_n) \in C), \quad C \subseteq \mathbb{R}^n.$$

Take $C = \prod_{i=1}^n (-\infty, x_i]$, then $(X_1, X_2, \dots, X_n) \in C$ if and only if $X_i \leq x_i$ for $i = 1, 2, \dots, n$.

Definition 4.4.1 ► Joint Cumulative Distribution Function

Let X_1, X_2, \dots, X_n be random variables, then their **joint cumulative distribution function** is defined as

$$F_X(x_1, x_2, \dots, x_n) = P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n).$$

One may be tempted to think that $P(X \leq x, Y \leq y) = P(X \leq x)P(Y \leq y)$, but in general this is not true!

Remark. $P(X \leq x, Y \leq y) = P(X \leq x)P(Y \leq y)$ if and only if X and Y are independent random variables (Definition 4.4.6).

Consider two random variables X and Y jointly distributed, observe that

$$\begin{aligned} P(x_1 \leq X \leq x_2, Y \leq y) &= P(X \leq x_2, Y \leq y) - P(X \leq x_1, Y \leq y) \\ &= F_{X,Y}(x_2, y) - F_{X,Y}(x_1, y). \end{aligned}$$

Similarly,

$$P(X \leq x, y_1 \leq Y \leq y_2) = F_{X,Y}(x, y_2) - F_{X,Y}(x, y_1).$$

Combining the two equations we have

$$\begin{aligned} P(x_1 \leq X \leq x_2, y_1 \leq Y \leq y_2) &= P(x \leq x_2, y_1 \leq Y \leq y_2) - P(x \leq x_1, y_1 \leq Y \leq y_2) \\ &= F_{X,Y}(x_2, y_2) - F_{X,Y}(x_2, y_1) - F_{X,Y}(x_1, y_2) + F_{X,Y}(x_1, y_1). \end{aligned}$$

Note that with this identity we can compute $P((X, Y) \in C)$ for all Borel sets C . A *Borel set* is a set generated by countable unions and countable intersections of intervals.

Consider jointly distributed discrete random variables. We first introduce the following intuitive definition:

Definition 4.4.2 ► Joint Probability Mass Function

Let X_1, X_2, \dots, X_n be discrete random variables. Their **joint probability mass function** is defined to be

$$p_X(x_1, x_2, \dots, x_n) = P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n).$$

In the continuous case, we have a similar definition:

Definition 4.4.3 ► Joint Continuity

Let X_1, X_2, \dots, X_n be continuous random variables. They are said to be **jointly continuous** if there exists a **joint probability density function** f such that

$$P((X_1, X_2, \dots, X_n) \in C) = \int \cdots \int_C f_X(x_1, x_2, \dots, x_n) dx_1 \cdots dx_n.$$

Let X, Y be discrete random variables. It is easy to see that

$$\begin{aligned} P(X = x) &= \sum_{i=1}^{\infty} P(X = x, Y = y_i), \\ P(Y = y) &= \sum_{i=1}^{\infty} P(X = x_i, Y = y). \end{aligned}$$

On the other hand, if X, Y are continuous random variables, we have

$$\begin{aligned} P(X \in A) &= P(X \in A, Y \in \mathbb{R}) \\ &= \int_A \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy dx. \end{aligned}$$

These are known as *marginal probability density functions*.

Definition 4.4.4 ► Marginal Probability Density Function

Let X, Y be discrete random variables. Their **marginal probability density functions** are defined as

$$p_X(x) = \sum_{i=1}^{\infty} p_{X,Y}(x, y_i),$$

$$p_Y(y) = \sum_{i=1}^{\infty} p_{X,Y}(x_i, y).$$

If X, Y are continuous random variables, then

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

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

Lastly, we state the notion of *joint distribution of functions*.

Theorem 4.4.5 ► Joint Distribution of Functions

Let X_1 and X_2 be jointly continuous random variables. Let $Y_1 = g_1(X_1, X_2)$ and $Y_2 = g_2(X_1, X_2)$. If for all y_1, y_2 , the linear system

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} g_1(x_1, x_2) \\ g_2(x_1, x_2) \end{bmatrix}$$

has a unique solution, and the Jacobian

$$J(x_1, x_2) = \frac{\partial y_1}{\partial x_1} \cdot \frac{\partial y_2}{\partial x_2} - \frac{\partial y_1}{\partial x_2} \cdot \frac{\partial y_2}{\partial x_1}$$

is continuous and non-zero, then

$$f_{Y_1, Y_2}(y_1, y_2) = f_{X_1, X_2}(x_1, x_2) |J(x_1, x_2)|^{-1}.$$

4.4.1 Independent Random Variables

Recall that we have previously defined the notion of independent events (Definition 3.3.1). Observe that an event can be precisely expressed as $X \in A$ for some random variable X and set A , which motivates us to define independence of random variables.

Definition 4.4.6 ▶ Independent Random Variables

Two random variables X and Y are **independent** if for any sets A and B ,

$$P(X \in A, Y \in B) = P(X \in A)P(Y \in B).$$

Remark. Less rigorously, notice that for Borel sets A and B , since they can be expressed as countable unions and intersections of intervals, we can say that X and Y are independent if and only if $F_{X,Y}(x, y) = F_X(x)F_Y(y)$.

Since the cumulative distribution functions are closely related to the probability mass and probability density, we can then check for independence using them instead. We first consider the discrete case.

Theorem 4.4.7 ▶ Independence of Discrete Random Variables

Let X and Y be discrete random variables with probability mass functions p_X and p_Y . X and Y are independent if and only if $p_{X,Y}(x, y) = p_X(x)p_Y(y)$.

Based on Theorem 4.4.7, one may check that if X and Y are independent discrete random variables, then

$$E[XY] = E[X]E[Y], \quad \text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y).$$

In the continuous case, we have a similar conclusion.

Theorem 4.4.8 ▶ Independence of Continuous Random Variables

Let X and Y be continuous random variables with probability density functions f_X and f_Y . X and Y are independent if and only if $f_{X,Y}(x, y) = f_X(x)f_Y(y)$.

The above theorems and definitions can be easily extended to a countable number of independent random variables.

4.4.2 Sums of Independent Random Variables

Suppose X and Y are independent **integer-valued** discrete random variables, then clearly

$$\begin{aligned}
 p_{X+Y}(n) &= P(X + Y = n) \\
 &= \sum_i P(X = i, Y = n - i) \\
 &= \sum_i p_X(i) p_Y(n - i) \\
 &= \sum_{i+j=n} p_X(i) p_Y(j).
 \end{aligned}$$

On the other hand, if X and Y are independent continuous random variables, we have

$$\begin{aligned}
 F_{X+Y}(n) &= \iint_{x+y \leq n} f_{X+Y}(x+y) dx dy \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{n-y} f_X(x) f_Y(y) dx dy \\
 &= \int_{-\infty}^{\infty} F_X(n-y) f_Y(y) dy \\
 &= \int_{-\infty}^{\infty} F_X(n-y) dF_Y(y).
 \end{aligned}$$

We say that F_{X+Y} is the *convolution* $F_X * F_Y$.

Now, we proceed to discussing some special cases.

Let X and Y be independent uniform random variables on (a, b) . Note that $f_Y(y) = \frac{1}{b-a}$ for all $y \in (a, b)$, so

$$\begin{aligned}
 f_{X+Y}(n) &= \int_a^b f_X(n-y) f_Y(y) dy \\
 &= \frac{1}{b-a} \int_a^b f_X(n-y) dy.
 \end{aligned}$$

Note that $f_X(n-y) = \frac{1}{b-a}$ if and only if $a < n-y < b$, i.e., $n-b < y < n-a$. Otherwise, $f_X(n-y) = 0$. If $2a < n < a+b$, we have $(a, b) \cap (n-b, n-a) = (a, n-a)$, so

$$\begin{aligned}
 \frac{1}{b-a} \int_a^b f_X(n-y) dy &= \frac{1}{b-a} \int_a^{n-a} \frac{1}{b-a} dy \\
 &= \frac{n-2a}{(b-a)^2}.
 \end{aligned}$$

If $a + b < n < 2b$, then $(a, b) \cap (n - b, n - a) = (n - b, b)$, so

$$\begin{aligned} \frac{1}{b-a} \int_a^b f_X(n-y) dy &= \frac{1}{b-a} \int_{n-b}^b \frac{1}{b-a} dy \\ &= \frac{2b-n}{(b-a)^2}. \end{aligned}$$

Therefore,

$$f_{X+Y}(n) = \begin{cases} \frac{n-2a}{(b-a)^2} & \text{if } 2a < n \leq a+b \\ \frac{2b-n}{(b-a)^2} & \text{if } a+b < n < 2b \\ 0 & \text{otherwise} \end{cases}$$

This is known as a *triangular distribution*.

Let X and Y be independent gamma random variables with parameters (α, λ) and (β, λ) . One may check with some computation that

$$\begin{aligned} f_{X+Y}(n) &= \frac{B(\alpha, \beta)}{\Gamma(\alpha)\Gamma(\beta)} \lambda e^{-\lambda n} (\lambda n)^{\alpha+\beta-1} \\ &= \frac{1}{\Gamma(\alpha + \beta)} \lambda e^{-\lambda n} (\lambda n)^{\alpha+\beta-1}. \end{aligned}$$

Therefore, $X + Y$ is a gamma random variable with parameters $(\alpha + \beta, \lambda)$.

Lastly, we shall state the following theorem without proof:

Theorem 4.4.9 ► Sum of Normal Random Variables Is Normal

Let $X \sim N(\mu_X, \sigma_X^2)$ and $Y \sim N(\mu_Y, \sigma_Y^2)$, then

$$X + Y \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2).$$

4.5 Conditional Distribution

Recall that previously, we have defined the conditional probability

$$P(E | F) = \frac{P(EF)}{P(F)}.$$

Let X and Y be discrete random variables, we can similarly see that

$$\begin{aligned} P(X = x | Y = y) &= \frac{P(X = x, Y = y)}{P(Y = y)} \\ &= \frac{p_{X,Y}(x, y)}{p_Y(y)}. \end{aligned}$$

Definition 4.5.1 ► Conditional Probability Mass Function

Let X and Y be discrete random variables with probability mass functions p_X and p_Y respectively, the **conditional probability mass function** of X given $Y = y$ is defined as

$$p_{X|Y}(x | y) = \frac{p_{X,Y}(x, y)}{p_Y(y)}.$$

Remark. In particular, we have

$$P(X = x | X \in A) = \begin{cases} \frac{P(X=x)}{P(X \in A)}, & \text{if } x \in A \\ 0, & \text{otherwise} \end{cases}.$$

We can similar define for the continuous case:

Definition 4.5.2 ► Conditional Probability Density Function

Let X and Y be jointly continuous random variables with probability density functions f_X and f_Y respectively, the **conditional probability density function** of X given $Y = y$ is defined as

$$f_{X|Y}(x | y) = \frac{f_{X,Y}(x, y)}{f_Y(y)}.$$

Remark. In particular, we have

$$f_{X|X \in A}(x) = \frac{f_X(x)}{\int_A f_X(x) dx}.$$

Now, we can compute the conditional probability for jointly distributed continuous random variables with

$$P(X \in A | Y = y) = \int_A f_{X|Y}(x | y) dx.$$

In particular, we could define the *conditional cumulative distribution function* as

$$\begin{aligned} F_{X|Y}(x | y) &= P(X \leq x | Y = y) \\ &= \int_{-\infty}^x f_{X|Y}(x | y) dx \end{aligned}$$

What if X is continuous but Y is discrete? In this case, we have

$$f_{X|Y}(x | y) = \frac{P(Y = y | X = x)}{P(Y = y)} f_X(x).$$

Expectation

5.1 Sums of Random Variables

In general, we have:

Theorem 5.1.1 ► Expectation of Functions of Jointly Distributed Random Variables

Let X and Y be jointly distributed random variables. If X and Y are discrete, then

$$E[g(X, Y)] = \sum_x \sum_y g(x, y) p_{X,Y}(x, y).$$

If X and Y are continuous, then

$$E[g(X, Y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f_{X,Y}(x, y) dx dy.$$

Based on the above, it is easy to see that if X_1, X_2, \dots, X_n are random variables,

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

Definition 5.1.2 ► Sample Mean

Let X_1, X_2, \dots, X_n be independent random variables. The **sample mean** is defined as

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

Clearly, if $X_i = \mu$ for all $i = 1, 2, \dots, n$,

$$E[\bar{X}] = \frac{1}{n} \sum_{i=1}^n E[X_i] = \mu.$$

This means that given a random sample, its sample mean is an *unbiased estimate* of its expectation.

Using the expectation of the sum of random variables, we can also prove the following inequality:

Theorem 5.1.3 ► Boole's Inequality

Let A_1, A_2, \dots, A_n be events, then

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

Proof. Define

$$X_i = \begin{cases} 1, & \text{if } A_i \text{ occurs} \\ 0, & \text{otherwise} \end{cases}.$$

Let $X = \sum_{i=1}^n X_i$ and $Y = \max\{X_1, X_2, \dots, X_n\}$. Note that $Y = 1$ if and only if at least one of the A_i 's occurs, so

$$E[Y] = P\left(\bigcup_{i=1}^n A_i\right).$$

Note that $E[X] = \sum_{i=1}^n P(A_i)$ and that $X \geq Y$, so

$$\sum_{i=1}^n P(A_i) = E[X] \geq E[Y] = P\left(\bigcup_{i=1}^n A_i\right).$$

□

So far we have been discussing the sum of finitely many random variables. It turns out that in the infinite case, the following applies:

Theorem 5.1.4 ► Expectation of Infinite Sum of Random Variables

Let X_1, X_2, \dots be infinitely many random variables, if

- $X_i \geq 0$ for all $i \in \mathbb{N}^+$, or
- the series $\sum_{i=1}^{\infty} E[|X_i|]$ converges,

then

$$E\left[\sum_{i=1}^{\infty} X_i\right] = \sum_{i=1}^{\infty} E[X_i].$$

5.2 Moments

Consider n events A_1, A_2, \dots, A_n . Define

$$X_i = \begin{cases} 1, & \text{if } A_i \text{ occurs} \\ 0, & \text{otherwise} \end{cases},$$

then $X = \sum_{i=1}^n X_i$ is the number of events which have occurred. Observe that $E[X] = \sum_{i=1}^n P(A_i)$. Define the event

$$E_{m,k} = \bigcap_{i=1}^k A_{m_i}$$

where $1 \leq m_i \leq n$ for any $i = 1, 2, \dots, k$, then the number of such $E_{m,k}$'s which have occurred is given by $\binom{X}{k}$. Note that the event $\bigcap_{i=1}^k A_{m_i}$ occurs if and only if $\prod_{i=1}^k X_{m_i} = 1$, so

$$\binom{X}{k} = \sum_{m_1 < m_2 < \dots < m_k} \left(\prod_{i=1}^k X_{m_i} \right).$$

Therefore,

$$\begin{aligned} E \left[\binom{X}{k} \right] &= E \left[\sum_{m_1 < m_2 < \dots < m_k} \left(\prod_{i=1}^k X_{m_i} \right) \right] \\ &= \sum_{m_1 < m_2 < \dots < m_k} P(E_{m,k}). \end{aligned}$$

Notice that $E \left[\binom{X}{k} \right]$ is a linear combination of the first k -th moments of X . We would like to find a way to quickly compute the n -th moment of a random variable.

Let X be a discrete random variable, then

$$E[X^n] = \sum_x x^n p_X(x).$$

Note that by Maclaurin Series, we have

$$g(t) = \sum_{n=0}^{\infty} \frac{g^{(n)}(0)}{n!} t^n$$

for any function g . A motivation is to construct g such that $g^{(n)}(0) = E[X^n]$. Therefore,

$$\begin{aligned}
 g(t) &= \sum_{n=0}^{\infty} \frac{E[X^n]}{n!} t^n \\
 &= \sum_{n=0}^{\infty} \frac{\sum_x x^n p_X(x)}{n!} t^n \\
 &= \sum_x \left(p_X(x) \sum_{n=0}^{\infty} \frac{(tx)^n}{n!} \right) \\
 &= \sum_x e^{tx} p_X(x) \\
 &= E[e^{tX}].
 \end{aligned}$$

Definition 5.2.1 ► Moment Generating Function

Let X be a random variable, the **moment generating function** of X is defined as

$$M_X(t) = E[e^{tX}]$$

such that

$$E[X^n] = M_X^{(n)}(0) \quad \text{for all } n \in \mathbb{N}.$$

It can be proven that for any random variable X , M_X is unique. In other words, if two random variables have the same moment generating function, they must be the same random variable.

Theorem 5.2.2 ► Uniqueness of Moment Generating Function

Let X and Y be random variables. If

$$\lim_{t \rightarrow 0} M_X(t) = \lim_{t \rightarrow 0} M_Y(t),$$

then X and Y have the same distribution.

By properties of exponential functions, the following theorem can also be easily proven:

Proposition 5.2.3 ► Moment Generating Function of Sum of Independent Random Variables

Let X_1, X_2, \dots, X_n be independent random variables, then

$$M_{\sum_{i=1}^n X_i}(t) = \prod_{i=1}^n M_{X_i}(t).$$

However, note that the converse of Proposition 5.2.3 is not true in general.

Definition 5.2.4 ► Joint Moment Generating Function

Let X and Y be random variables. The **joint moment generating function** of X and Y is defined as

$$M_{X,Y}(s, t) = E[e^{sX+tY}].$$

Observe that in the joint case, $M_X(s) = M_{X,Y}(s, 0)$ and $M_Y(t) = M_{X,Y}(0, t)$. X and Y are independent if and only if $M_{X,Y}(s, t) = M_X(s)M_Y(t)$.

5.3 Covariance and Correlations

5.3.1 Covariance

Suppose that X and Y are independent random variables. We have seen that

$$E[XY] = E[X]E[Y].$$

This motivates us to consider the quantity $E[XY] - E[X]E[Y]$.

Definition 5.3.1 ► Covariance

Let X and Y be random variables. The **covariance** of X and Y is defined as

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

Remark. We have seen that if X and Y are independent, then $\text{Cov}(X, Y) = 0$. But the converse is not true in general!

Note that covariance has the following properties:

1. $\text{Cov}(X, X) = \text{Var}(X)$.
2. $\text{Cov}(X, Y) = \text{Cov}(Y, X)$.

3. $\text{Cov}(aX, Y) = a\text{Cov}(X, Y)$.
4. $\text{Cov}(X_1 + X_2, Y) = \text{Cov}(X_1, Y) + \text{Cov}(X_2, Y)$.

Remark. Let V be the set of all random variables, then $\text{Cov}(\cdot, \cdot) : V \times V \rightarrow \mathbb{R}$ is an inner product over V .

Let X and Y be random variables, then by the above properties, we have

$$\begin{aligned}
 \text{Var}(X + Y) &= \text{Cov}(X + Y, X + Y) \\
 &= \text{Cov}(X, X + Y) + \text{Cov}(Y, X + Y) \\
 &= \text{Cov}(X + Y, X) + \text{Cov}(X + Y, Y) \\
 &= \text{Cov}(X, X) + \text{Cov}(Y, X) + \text{Cov}(X, Y) + \text{Cov}(Y, Y) \\
 &= \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y).
 \end{aligned}$$

In general, we have the following formula:

Proposition 5.3.2

Let X_1, X_2, \dots, X_n be random variables, then

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

In particular, if all of the X_i 's are independent, then

$$\text{Var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \text{Var}(X_i).$$

Definition 5.3.3 ▶ Deviation

Let X_1, X_2, \dots, X_n be independent random variables with mean μ and variance σ^2 . The **deviation** is defined as

$$X_i - \bar{X} = X_i - \frac{1}{n} \sum_{i=1}^n X_i.$$

Similar to the sample mean, we would wish the deviation gives an unbiased estimate for the variance. However,

$$E[X_i - \bar{X}] = \frac{n-1}{n} \sigma^2.$$

Therefore, we need to eliminate the bias from the deviation.

Definition 5.3.4 ▶ Sample Variance

Let X_1, X_2, \dots, X_n be independent random variables with mean μ and variance σ^2 . The **sample variance** is defined as

$$S^2 = \sum_{i=1}^n \frac{X_i - \bar{X}}{n-1}.$$

5.3.2 Correlation

Intuitively, we view covariance $\text{Cov}(X, Y)$ as a measure of the degree of spread of the joint distribution of X and Y . Naturally, we can make use of it to measure how related are X and Y .

Definition 5.3.5 ▶ Correlation

Let X and Y be random variables with positive variances. The **correlation** of X and Y is defined as

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}.$$

In particular, X and Y are said to be **uncorrelated** if $\rho(X, Y) = 0$.

It is easy to see that for $c > 0$, $\rho(\pm cX, Y) = \pm \rho(X, Y)$. Suppose X and Y have variances σ_X^2 and σ_Y^2 respectively, we have

$$\rho(X, Y) = \rho\left(\frac{X}{\sigma_X}, \frac{Y}{\sigma_Y}\right) = \text{Cov}\left(\frac{X}{\sigma_X}, \frac{Y}{\sigma_Y}\right).$$

The above property leads to the following result:

Theorem 5.3.6 ▶ Boundedness of Correlation

Let X and Y be random variables, then

$$-1 \leq \rho(X, Y) \leq 1.$$

Note that $\text{Var}(X) = 0$ if and only if X is a constant, i.e., $P(X = c) = 1$ and $P(X = x) = 0$ for all $x \neq c$. Therefore, $\rho(X, Y) = \pm 1$ if and only if $Y = \pm aX + b$ for some constants a, b with $a > 0$.

5.4 Conditional Expectation and Variance

5.4.1 Conditional Expectation

Recall that in Definitions 4.5.1 and 4.5.2, we have defined conditional distribution of a random variable X given $Y = y$. We will study the notion of *conditional expectation* now.

Definition 5.4.1 ► Conditional Expectation

Let X and Y be random variables. The **conditional expectation** of X given $Y = y$ is defined as

$$E[X | Y = y] = \sum_x x p_{X|Y}(x | y)$$

if X and Y are discrete, and

$$E[X | Y = y] = \int_{-\infty}^{\infty} \frac{f_{X,Y}(x, y)}{f_Y(y)} dx$$

if X and Y are jointly continuous.

All properties of expectation are still satisfied by conditional expectation. In particular, notice that $E[X | Y]$ is a function in Y , so it is a random variable by itself. Let $Z = E[X | Y]$, what is the expectation of Z ?

Proposition 5.4.2

Let X and Y be random variables, then

$$E[E[X | Y]] = E[X].$$

5.4.2 Conditional Variance

Similarly, we can define *conditional variance*.

Definition 5.4.3 ► Conditional Variance

Let X and Y be random variables. The **conditional variance** of X given Y is defined as

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

Notice that $\text{Var}(X | Y)$ is also a random variable as it is a function of Y . Since $\text{Var}(X | Y) =$

$E[X^2 | Y] - (E[X | Y])^2$, we have

$$\begin{aligned} E[\text{Var}(X | Y)] &= E[E[X^2 | Y]] - E[(E[X | Y])^2] \\ &= E[X^2] - E[(E[X | Y])^2], \\ \text{Var}(\text{Var}(X | Y)) &= E[(E[X | Y])^2] - (E[E[X | Y]])^2. \end{aligned}$$

With some manipulations we can prove the following proposition:

Proposition 5.4.4

Let X and Y be random variables, then

$$\text{Var}(X) = E[\text{Var}(X | Y)] + \text{Var}(E[X | Y]).$$

5.4.3 Prediction

In many situations, we seek to **predict** the outcome of an event. Specifically, given $X = x$, we wish to find a function $g : \text{Range}(X) \rightarrow \text{Range}(Y)$ such that $g(x) = \hat{y}$ is the prediction for Y . In an ideal case, we would want $g(X)$ to be the **closest** to Y . In other words, we would minimise $g(X) - Y$. To eliminate the inconvenience of a negative difference, we would choose a g such that $E[(Y - g(X))^2]$ is minimised.

Let us consider a simpler case. When $g(x) = c$, to minimise $E[(Y - c)^2]$, consider

$$\begin{aligned} E[(Y - c)^2] &= c^2 - 2cE[Y] + E[Y^2] \\ &= (c - E[Y])^2 + \text{Var}(Y). \end{aligned}$$

Therefore, we need $c = E[Y]$. Since $c = g(x)$, this implies that $g(x) = E[Y | X = x]$. This is summarised into the following theorem:

Theorem 5.4.5 ► Best Predictor

Let X and Y be random variables. Given X , the best predictor of Y is the function

$$g(X) = E[Y | X].$$

Essentially, this implies that for all function f of X , $E[(Y - f(X))^2] \geq E[(Y - E[Y | X])^2]$.

However, in some cases we do not know the exact joint distribution of X and Y and so we cannot find $E[Y | X]$. Therefore, we may attempt to predict Y with a linear function of X which can be easily formulated. Hence, we will find constants a and b such that $E[(Y - a - bX)^2]$ is minimised.

Let $E[X] = \mu_X$, $E[Y] = \mu_Y$, $\text{Var}(X) = \sigma_X^2$, $\text{Var}(Y) = \sigma_Y^2$. Define

$$X' = \frac{X - \mu_X}{\sigma_X}, \quad Y' = \frac{Y - \mu_Y}{\sigma_Y}.$$

Then, $E[X'] = E[Y'] = 0$ and $\text{Var}(X') = \text{Var}(Y') = 1$. Note that

$$Y - a - bX = \sigma_Y \left(Y' - \frac{a + b\mu_X - \mu_Y}{\sigma_Y} - \frac{b\sigma_X}{\sigma_Y} X' \right).$$

Therefore, setting $\frac{a + b\mu_X - \mu_Y}{\sigma_Y} = a'$ and $\frac{b\sigma_X}{\sigma_Y} = b'$, we have

$$\begin{aligned} E[(Y - a - bX)^2] &= \sigma_Y^2 E[(Y' - a' - b'X')^2] \\ &= \sigma_Y^2 E[Y'^2 + a'^2 + b'^2 X'^2 - 2a'Y' + 2a'b'X' - 2b'X'Y'] \\ &= \sigma_Y^2 (1 + a'^2 + b'^2 - 2b'\rho(X, Y)) \\ &= \sigma_Y^2 (a'^2 + (b' - \rho(X, Y))^2 + 1 - \rho(X, Y)^2). \end{aligned}$$

Therefore, we need $a' = 0$ and $b' = \rho(X, Y)$, i.e.,

$$a = \mu_Y - b\mu_X, \quad b = \rho(X, Y) \frac{\sigma_Y}{\sigma_X}.$$

Therefore,

$$\begin{aligned} a + bX &= \mu_Y - b\mu_X + bX \\ &= \mu_Y + b(X - \mu_X) \\ &= \mu_Y + \rho(X, Y) \frac{\sigma_Y}{\sigma_X} (X - \mu_X) \end{aligned}$$

and the minimum of $E[(Y - a - bX)^2]$ is $\sigma_Y^2 (1 - \rho(X, Y)^2)$.

Theorem 5.4.6 ► Best Linear Predictor

Let X and Y be random variables with means μ_X and μ_Y , and variances σ_X^2 and σ_Y^2 respectively. Given X , the best linear predictor of Y is

$$g(X) = E \left[\left(Y - \mu_Y - \rho(X, Y) \frac{\sigma_Y}{\sigma_X} (X - \mu_X) \right) \right].$$