

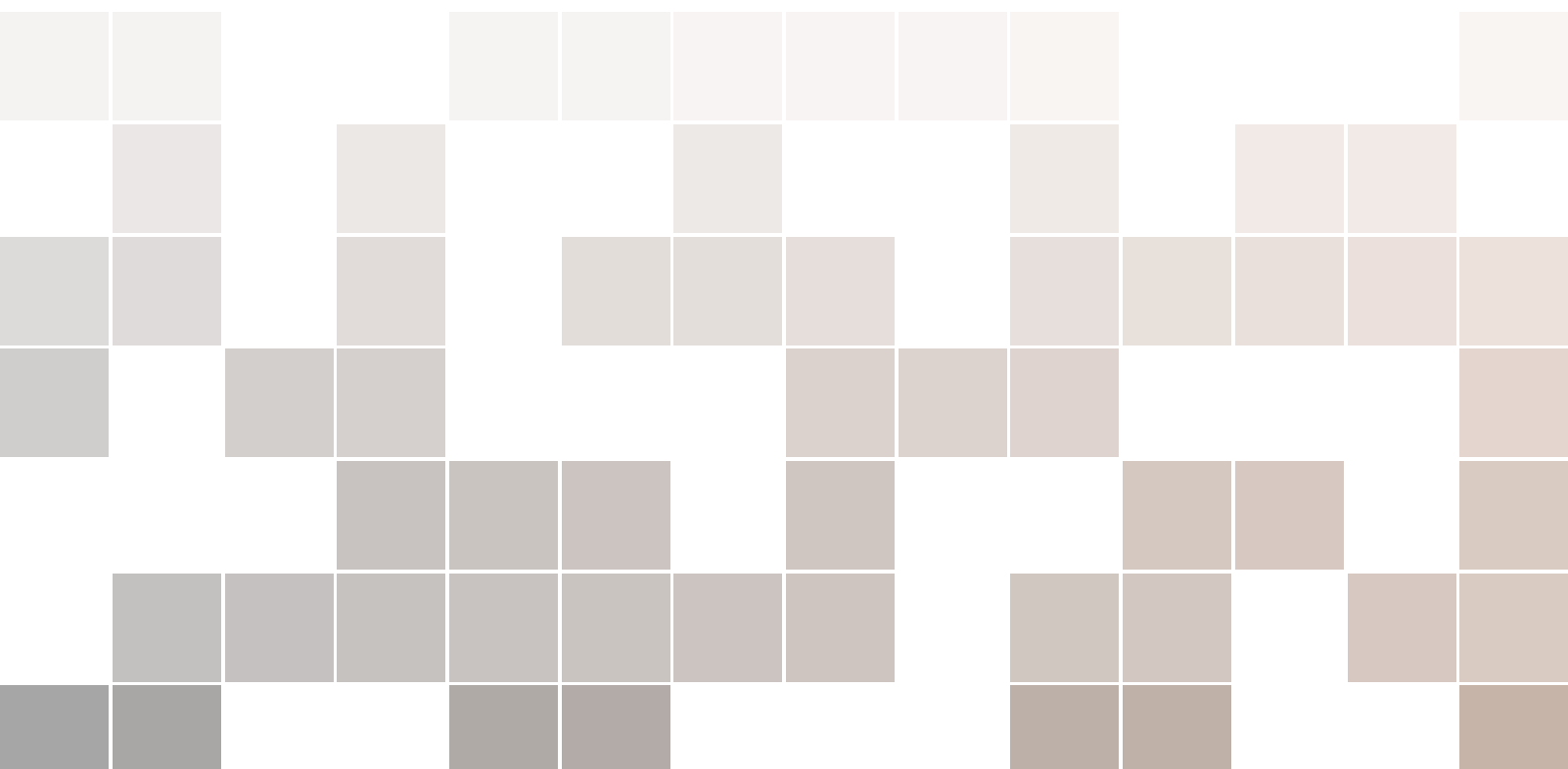


STA247

Probability with Computer Applications

Prof. K. H. Wong

Sinan Li



Copyright © 2022 UofT - STA247

ARTSCI.CALENDAR.UTORONTO.CA/COURSE/STA247H1

Licensed under the Creative Commons Attribution-NonCommercial 3.0 Unported License (the “License”). You may not use this file except in compliance with the License. You may obtain a copy of the License at <http://creativecommons.org/licenses/by-nc/3.0>. Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an “AS IS” BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

First digital, October 2022

Contents

I	Contents	
	Introduction	7
	Textbooks	7
	Course Information	8
	Discussion Board	8
1	Introduction to Probability	11
1.1	Useful Terminologies	11
1.1.1	Random Experiments	11
1.1.2	Operations on Sets	12
1.1.3	Important Laws	13
1.2	Probability Functions	14
1.2.1	Probability Functions	14
1.2.2	Probability and Event Relations	15
2	Counting	19
2.1	Counting and Probability	19
2.2	Fundamental Principle of Counting	19
2.3	Permutations	20
2.4	Combinations	22
3	Conditional Probability	25
3.1	Conditional Probability	25

3.2	Independence	27
3.3	Bayes' Rule	28
4	Discrete Distributions	31
4.1	Random Variables	31
4.1.1	Introduction to Random Variables	31
4.1.2	Characteristics of Random Variables	33
4.2	Cumulative Distribution Function	36
4.2.1	Properties of CDF	37
4.2.2	Chebyshev's Inequality	37
4.3	Common Discrete Distributions	39
4.3.1	Binomial Distribution	39
4.3.2	Poisson Distribution	41
4.3.3	Geometric Distribution	45
4.3.4	Negative Binomial Distribution	45
4.3.5	Hypergeometric Distribution	46
5	Continuous Random Variables	49
5.1	Probability Density Function	49
5.2	Common Continuous Distributions	51
5.2.1	Uniform Distribution	51
5.2.2	Cumulative Distribution Function	52
5.2.3	Exponential Distribution	53
5.2.4	Gamma Distribution	54
5.3	Percentile	57
6	Bivariate Distributions	59
6.1	Multivariate and Joint Distributions	59
	Bibliography	61
	Textbooks	61
	Index	63

Contents

	Introduction	7
	Textbooks	
	Course Information	
	Discussion Board	
1	Introduction to Probability	11
1.1	Useful Terminologies	
1.2	Probability Functions	
2	Counting	19
2.1	Counting and Probability	
2.2	Fundamental Principle of Counting	
2.3	Permutations	
2.4	Combinations	
3	Conditional Probability	25
3.1	Conditional Probability	
3.2	Independence	
3.3	Bayes' Rule	
4	Discrete Distributions	31
4.1	Random Variables	
4.2	Cumulative Distribution Function	
4.3	Common Discrete Distributions	
5	Continuous Random Variables	49
5.1	Probability Density Function	
5.2	Common Continuous Distributions	
5.3	Percentile	
6	Bivariate Distributions	59
6.1	Multivariate and Joint Distributions	
	Bibliography	61
	Textbooks	
	Index	63



Introduction

This is an introductory course to probability, where our main focus will be developing an understanding of probability and the concept of probability distributions, both for discrete and continuous quantities. This includes developing the intuition for how probabilities ‘behave’ and the situations in which it is valid to describe randomness using probability, as well as relying on simulations in R to help us visualize these properties.

By the end of the course, students should be able to...

- Describe random quantities in various ways, such as by: their features, density functions, distribution functions, graphs
- Select an appropriate probability model based on their unique properties to quantify the randomness of random quantities
- Compute and interpret the various features of a random quantity: expected value, variance, standard deviation, correlation, covariance, event probabilities (either exactly, through approximations, or simulation)
- Select the most appropriate model to represent randomness and compute probabilities
- Use simulation in R for estimation purposes
- Explain the relationship of transforming a random variable and its effect on its distribution
- Use bivariate distributions to describe the association between two random quantities

Textbooks

We will be referencing these two textbooks regularly:

1. *Probability with Applications and R*, 2nd ed. by *Wagaman and Dobrow* through the library [here](#) and with student companion site [here](#).
2. *Modern Mathematical Statistics with Applications*, 3rd ed. by *Devore and Berk* available through the library [here](#).

We will be using R throughout the course to help us understand probability distributions, and to simulate probabilities for quantities that are harder to compute by hand. R Markdown files will be provided for you with necessary starter code. You’ll also gain some experience with using \LaTeX to produce documents with well-presented math notation.

For new users to R, you may find chapters 3-7 from R for Data Science to be helpful as reference material in data visualization and data manipulation tools in R.

Course Information

All course-related materials can be found on our Quercus page:

- **Lectures** are held in MC 102 for both sections. We begin 10 minutes after the hour, and no recordings will be provided.
- **Weekly materials** such as slides, suggested problems for practice, and reminders of upcoming due dates are posted each week's page on our course home page.
- **Tutorial materials** will be distributed at the beginning of your tutorial.
- **Announcements**: this is the primary channel to distribute important information to everyone! You're expected to check and read announcements regularly to ensure you don't miss any important communication!
- **Assignments** and documents will be distributed and submitted through Crowdmark.

Discussion Board

Throughout the term, beginning September 14, there will be pinned weekly discussions on our discussion board page. They will remain pinned from Wednesday 8 AM to the following Tuesday 8 PM.

- **General Q&A page**: general questions, clarifications, request for additional explanations, share your thoughts/understanding of topics
- **Grouped practice problem discussions**: post your solutions, thoughts, approaches, questions here for that collection of textbook questions.
 - Each discussion thread (i.e. reply to one original comment) is dedicated to ONE question, labeled in bold in the original comment

You earn your discussion board credits by contributing at least five (5) times during this course to any of the discussions during the pinned period (Wednesday to Tuesday) in any of the following ways:

- Posting a question to a problem you tried, with a clear explanation of your process, and if you got stuck: what you tried to do, and where you need help moving forward
- Responding to a question with a thoughtful explanation to help your peer by sharing your own understanding of the problem
- Posting to the general Q&A page with your own question, about a course topic that is still unclear and being specific in describing what you do not understand (and perhaps what you did understand!)
- Providing a detailed and clear response to a question in the general Q&A page

To ensure boards are easy to reference, posting similar content is discouraged and these would be ineligible for earning credits.

- Only contributions during the pinned period will be counted. The discussions will remain open the rest of term for students who come across new problems or would like to continue the discussions.
- You are encouraged to keep the discussion going, but in terms of credit, it will be capped at 5 points.
- While there is no weekly cap on points, the maximum points you can earn in the last two weeks is 2 points, with no more than 2 points per week (i.e. don't wait to the last minute to participate in the class discussion).
- The discussion boards exist to facilitate peer-to-peer collaboration and learning, while also encouraging regular active engagement with course content. The course offers many

opportunities that most students shouldn't find themselves unable to contribute in a unique way.

Why a Discussion Board?

- There are records and studies that have shown the process of explaining and teaching to others is an effective way to learn, consolidate, and retain what has been taught. See [here](#) and [here](#).
- It's a space for students to come together to work collaboratively, receive and provide peer support.
- It's also a space to get feedback and guidance from TAs and myself.
- It's a good opportunity to self-assess ('how comfortable am I explaining this to another student?', 'how often do I need to refer back to my notes to explain this concept clearly?') – an important component of good study skills!
- It's valuable information to us! Common questions/misunderstandings that pop up in the weekly discussions can be addressed during our weekly lecture meetings.

Discussion Board Rubric

Points	1 point	0.5 points	0 points
Quality of contribution	<p>Student has made a substantial and unique contribution with detailed explanations and/or clearly outlined process of the approach to a problem.</p> <p>Student was involved in follow-up discussions and worked collaboratively with their peers to develop a better understanding of the concepts involved.</p>	<p>Student has made a contribution to the discussion that is dismissive, lacking in detail, or not completely unique. Unable to further the discussion in a way that fosters a collaborative learning environment.</p> <p>e.g. responses such as 'you just need to integrate this and solve for it' or 'I got the same answer doing... (reiterates OPs process)'</p>	<p>Student has not contributed to the weekly discussion topic thread, or whose posts are off-topic/irrelevant/do not contribute to the thread or is not unique to what has already been discussed in the thread.</p> <p>e.g. 'I got the same answer', 'How did you get that number?'</p>

Tutorial

- A mix of R labs and collaborative pair work
- R labs: These labs have guided exercises to practice the R tools covered in class or learn new R skills. Labs will be TA-guided.
- R labs require individual .rmd and knitted document submissions at the end of tutorial (Note that the labs are guided and meant to be completed within the tutorial time)
- Pair work: In your tutorial section, you will with a partner of your choice work on more challenging but guided problems together. Discussion and sharing your ideas is a great way to learn from one another, and consider different approaches to problem solving! TAs will be there to support and help answer clarification questions.

Habits for Success

- Attend and participate in lecture. Try to work along with the problems presented. Ask questions and interact with your peers during the open work periods!
- Make sure you focus on *understanding* the concepts and how they relate and build upon each other. This one is hard!
- Regularly attempt as many of the suggested problems as you can and work towards being able to work on the problems closed book. Use this to gauge your familiarity with the material – if it takes you an hour to work through two problems, then it's a good indication to seek out advice and support from the teaching team. The earlier, the better!
- Drop by during the office hours or post on discussion board if you get stuck. Work through practice problems with classmates. Take turns *explaining* your thinking and problem solving process.
- **Create a schedule and stick with it.** This course covers many topics to ensure you have a good foundation for latter courses. Many topics build on top of each other so falling behind can quickly snowball and make it difficult to catch up.

Some Suggestions from Previous Students

1. Will the final be cumulative?
Yes.
2. I find that I *really* struggle with keeping my mind focused on the question I'm presenting working on, and not think about how a similar question was solved previously. I thought I was making some progress, but then I realized this wasn't the case at all after today when I spent 30 minutes trying to recall the solution I wrote for A2 (I don't remember even reading the test question..I just caught a few phrases and began regurgitating my assignment solution incorrectly).

I realize the most helpful thing I can do for myself now is to practice regularly, but I can't say I've been doing that well. For the last week, I've been attempting the textbook questions (open book)in preparation for the test. For the final, I plan to solve questions independently, which was something I didn't do until last minute. What else can I do to improve my problem-solving/critical thinking skills?

Definitely do textbook / slide questions closed book. The only thing you should be looking at is the sample aid sheet. Otherwise you will never know if you actually understand the material.

3. You need to know when you don't know something. Meaning that when you see a question, you need to recognize if you actually know how to solve it or not. And if you don't know how to solve it you need to skip it and go to the next question. I'm sure you know that spending 30 min on any question is not an efficient use of the time.
4. Go to office hours. There is only like 2 other people that I've ever seen at office hours.
5. What does data have to look like for an exponential distribution to describe it?

Why do we use different distributions? What's the importance of justifying a decision to use one distribution over another? (the textbook often tells you the distribution so image all of the questions posed to you without that information. Would you still be confident?) Answering these questions thoroughly will make you a lot better at figuring out how to approach a question when you see it.

Understanding the conceptual side in depth is the most important part of studying I would think. Computational skills come with practice, eventually you get it. But usually the problem with stats or courses like this is interpreting the question - is figuring out how to even approach the question in the first place. Knowing how to solve all of the normal distribution questions that confront you is great, but that won't help you if you can't even recognize when and when not to use it.

1. Introduction to Probability

1.1 Useful Terminologies

1.1.1 Random Experiments

Definition 1.1.1 — Random Experiment. A process that allows us to gather data or observations. Experiment can be repeated multiple times under the same conditions. The set of possible outcomes of the experiment are known, but the outcome of a specific experiment is not known.

- **Example 1.1 — Random Experiments.** Below are some examples of random experiments.
- Rolling a die and observing the top-facing number
 - Rolling a pair of dice and observing the sum of top-facing numbers
 - A patient being administered a painkiller and observing the amount of time in minutes before relief is felt

■ **Definition 1.1.2 — Sample Space.** The *sample space* is the set of all possible outcomes / results from a random experiment, usually denoted by Ω or S . The elements in the sample space are determined by the outcome of interest. Elements are often denoted by ω .

- **Example 1.2 — Sample Space.** Below are some examples of the sample space of random experiments.

- Experiment: Rolling a 20-sided die and observing the top-facing result.
 - $\Omega = \{1, 2, 3, \dots, 20\}$ OR
 - $S = \{1 \leq x \leq 20, x \in \mathbb{Z}\}$
- Experiment: Selecting a random student and observing whether they are a CS student
 - $\Omega = \{\text{CS Student}, \text{Non-CS Student}\}$ OR
 - $S = \{0, 1\}$ where $\begin{cases} 0 = \text{Non-CS Student} \\ 1 = \text{CS Student} \end{cases}$
- Experiment: Select a random student and record the amount of liquids consumed that day
 $\Omega = \{L \geq 0, L \in \mathbb{R}\}$

Definition 1.1.3 — Event. An *event* is a subset of the sample space, usually represented by a capital letter near the beginning of the alphabet. A *simple event* has exactly one element of the sample space, while a *compound event* consists of multiple elements.

Definition 1.1.4 — Complement Event. The *complement event* is the set of outcomes in Ω that are not in A . Can be denoted as one of: A^c (preferred), \bar{A} , or A' .

■ **Example 1.3** Consider the following example.

- $A = B^c = \text{roll an even number} = \{2, 4, 6\}$
- $B = A^c = \text{roll an odd number} = \{1, 3, 5\}$

Here A and B are complement events. ■

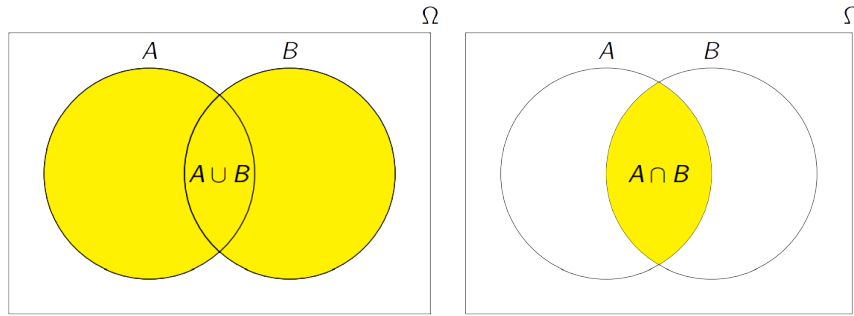
1.1.2 Operations on Sets

Definition 1.1.5 — Union. The *union* of two events A and B is the set of outcomes that are elements of A , B , or both. This is denoted as $A \cup B$.

- Union of events is usually described as A *or* B
- Note that $A \cup A^c = \Omega$

Definition 1.1.6 — Intersection. The *intersection* of two events A and B is the set of outcomes that are common to both A and B . This is denoted as $A \cap B$, or AB in the textbook.

- Intersection of event is usually described as A *and* B .
- Note that $A \cap A^c = \emptyset$



■ **Example 1.4 — Examples of Event Relations.** Let's keep it simple: suppose we roll two differently coloured¹ dice and record the paired outcomes. List out the following events:

a) Outcomes are doubles.

Let D be the event we roll doubles.

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

b) Outcomes sum to 8.

Let E be the event where rolls sum to 8.

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

c) Outcomes where one dice has twice the face value as the other.

Let T be the event where one die outcome is twice the other.

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

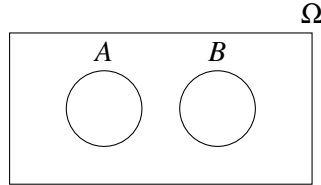
Note that here $E \cap T = \emptyset$ but $E \cup T \neq \Omega$. E and T are *disjoint* or *mutually exclusive*. E is **not** the complement of T . ■

subsectionEvent Relations

¹Note that these two dice are distinguishable, so we would get $\omega = (i, j)$

Definition 1.1.7 — Mutually Exclusive. Two events A and B are *mutually exclusive* if the events cannot both occur or occur simultaneously as an outcome of the experiment. This means that the **intersection** of A and B is **empty** and they have no overlapping elements. A and B are also called *disjoint* events.

In the following Venn Diagram, A and B would be disjoint.



Definition 1.1.8 — Independent. Two events A and B are *independent* if the occurrence of one event does not alter the probability of occurrence of the other in any way.

■ **Example 1.5** Classify the following pairs of events / variables as mutually exclusive, independent, or dependent.

- a) A student surveyed at random:
 $A = \{\text{Studies CS}\}$, $B = \{\text{Studies Stats}\}$
 They are **dependent**.
- b) A biased coin (80% Head) is tossed twice:
 $C = \{\text{Head on first toss}\}$ and $D = \{\text{Tail on second toss}\}$
 They are **independent**.
- c) An individual surveyed at random:
 $E = \{\text{Dislikes hiking}\}$ and $F = \{\text{Likes or indifferent to hiking}\}$
 They are **mutually exclusive**.
- d) A playing card is drawn:
 $G = \{\text{Card is red}\}$ and $H = \{\text{Card is Queen of Hearts}\}$
 They are **dependent**.

■

1.1.3 Important Laws

The following laws are useful relationships between unions and intersections of events that can help re-express events in simpler forms.

Theorem 1.1.1 — Commutative Law.

$$A \cup B = B \cup A$$

Theorem 1.1.2 — Associative Law.

$$(A \cup B) \cup C = A \cup (B \cup C)$$

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

Theorem 1.1.3 — Distributive Law.

$$A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$$

In addition to the commutative, associative, and distributive laws, **DeMorgan's Laws** present some interesting relationships between union and intersection of events.

Theorem 1.1.4 — DeMorgan's Laws. For two events A and B ,

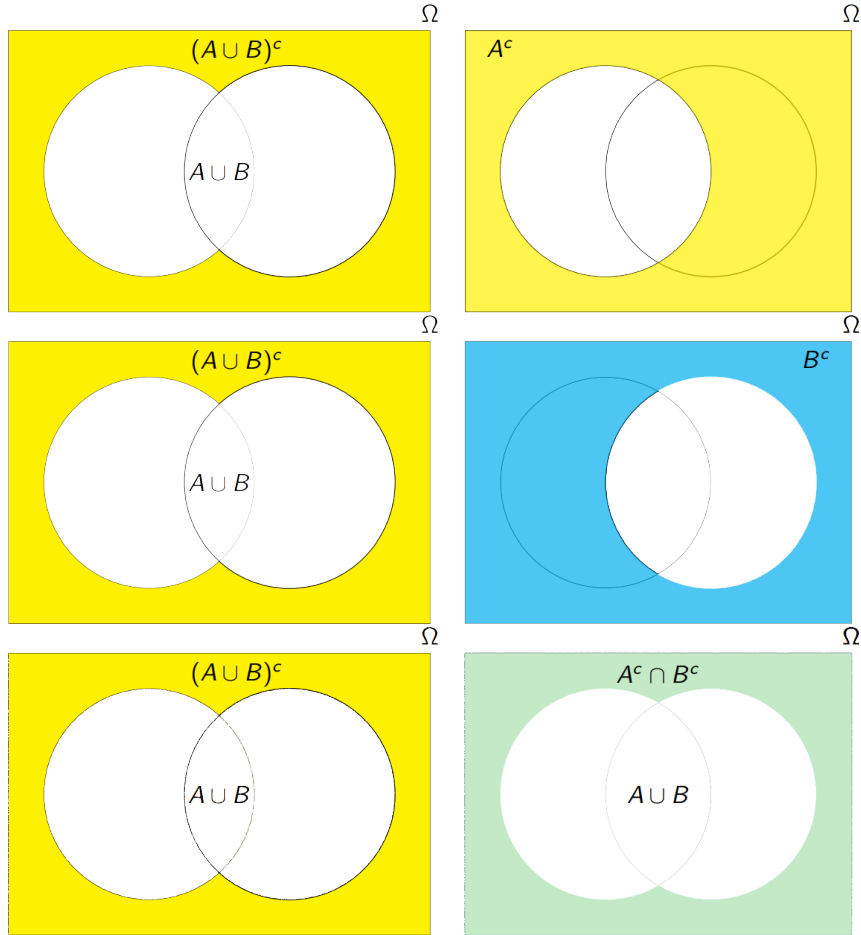
$$(A \cup B)^c = A^c \cap B^c$$

$$(A \cap B)^c = A^c \cup B^c$$

Or more generally, for the set of events $\{A_1, A_2, \dots, A_n\}$,

$$\left(\bigcup_{i=1}^n A_i \right)^c = \bigcap_{i=1}^n A_i^c$$

$$\left(\bigcap_{i=1}^n A_i \right)^c = \bigcup_{i=1}^n A_i^c$$



1.2 Probability Functions

1.2.1 Probability Functions

Definition 1.2.1 — Probability. In a random experiment with sample space Ω , the *probability* of an event A , denoted as $P(A)$ is a function that assigns to event A a **numerical value** ($P(A) \in [0, 1]$) that measures the chance that event A will occur.

There are certain axioms that must hold for probability functions.

axiom 1.2.1 — Axiom 1.

$$P(A) \geq 0$$

axiom 1.2.2 — Axiom 2.

$$P(\Omega) = 1$$

axiom 1.2.3 — Axiom 3. For a set of **disjoint** (mutually exclusive) elements A_1, A_2, \dots, A_n in Ω ,

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

Probabilities for outcomes of a random experiment can be represented in many ways. When we have outcomes that can be represented discretely, we can express the associated probability as a function, called a *probability function*. A valid probability function must satisfy all the probability axioms.

Definition 1.2.2 — Probability Function. Suppose the sample space Ω can be represented with a finite number of elements: $\Omega = \{\omega_1, \omega_2, \dots, \omega_n\}$ or a countably infinitely elements $\Omega = \{\omega_1, \omega_2, \dots\}$, then the *probability function* P is a function on Ω with the following properties:

1. $P(\omega) \geq 0$ for all $\omega \in \Omega$.
2. $\sum_{\omega \in \Omega} P(\omega) = 1$
3. For all events $A \subseteq \Omega$, $P(A) = \sum_{\omega \in A} P(\omega)$

1.2.2 Probability and Event Relations

Using the definition of probability, as well as the event relations learned last week, we can relate probabilities of event relations.

Proposition 1.2.4 — Complement Probability. $A \cup A^c = \Omega$ and A is disjoint from A^c . Then,

$$P(\Omega) = P(A \cup A^c)$$

We can further derive that $1 = P(A \cup A^c)$ Axiom 2

$$1 = P(A) + P(A^c) \quad \text{Axiom 3}$$

$$P(A^c) = 1 - P(A)$$

■ **Example 1.6** Suppose we have a box of 20 dice. Eight of the dice are custom made, and 15 of them are 10-sided dice. Assuming that each die belongs in either of these two categories, can we determine how many are custom made and are 10-sided?

Let $n(X)$ be the number of elements in X .

Let C be the custom made die

Let T be the 10-sided die

We know that $n(\Omega) = 20$, $n(C) = 8$, $n(T) = 15$, and $n((C \cup T)^c) = 0$.

$$8 + 15 = 23$$

notice $n(C \cup T)$ was counted twice

excess = over count

$$n(C \cap T) = 3$$

$$n(C \cup T) = n(\Omega)$$

■

■ **Example 1.7 — Union and Intersection.** Notice event B (or A) can be broken down into the following mutually exclusive events:

$$\begin{aligned} B &= (B \cap A) \cup (B \cap A^c) \\ P(B) &= P(A \cap B) + P(A^c \cap B) \quad \text{Axiom 3} \\ P(A^c \cap B) &= P(B) - P(A \cap B) \end{aligned}$$

Note also that $A \cup B$ can also be broken down in a similar way:

$$\begin{aligned} P(A \cup B) &= P(A \cup (A \cap B) \cup (B \cap A^c)) \\ &= P(A \cup (B \cap A^c)) && A \cap B \subseteq A \\ &= P(A) + P(A^c \cap B) && \text{Axiom 3, since disjoint} \\ P(A \cup B) &= P(A) + P(B) - P(A \cap B) && \text{from above} \end{aligned}$$

If A and B are disjoint (mutually exclusive), then $P(A \cup B) = P(A) + P(B)$ since $P(A \cap B) = P(\emptyset) = 0$. ■

■ **Example 1.8** In a class of 50 students, 23 could not roll their tongue, 15 had attached earlobes, and 10 could roll their tongues and had attached earlobes. A student is randomly selected from the class. Let T denote the event that the student can roll their tongue, and E denote the event that they have attached earlobes. Symbolically denote the following events and identify the number of students in each.

- a) The student can roll his or her tongue.

$$\begin{aligned} n(T) &= n(\Omega) - n(T^c) \\ &= 50 - 23 \\ &= 27 \end{aligned}$$

- b) The student can neither roll his or her tongue nor has attached earlobes.

$$\begin{aligned} n((T \cup E)^c) &= n(\Omega) - n(T \cup E) \\ &= n(\Omega) - (n(T) + n(E) - n(T \cap E)) \\ &= 50 - (27 + 15 - 10) \\ &= 18 \end{aligned}$$

In general, for the set of events $\{A_1, A_2, \dots, A_n\}$, we have the **Inclusion-Exclusion Principle**

Theorem 1.2.5 — Inclusion-Exclusion Principle.

$$\begin{aligned} P\left(\bigcup_{i=1}^n A_i\right) &= \sum_{r=1}^n \left((-1)^{r+1} \sum_{i_1 < i_2 < \dots < i_r} P(A_{i_1} \cap A_{i_2} \cap \dots \cap A_{i_r}) \right) \\ &= \sum_{i=1}^n P(A_i) - \sum_{i < j} P(A_i \cap A_j) + \dots + (-1)^{r+1} \sum_{i_1 < i_2 < \dots < i_r} P(A_{i_1} \cap A_{i_2} \cap \dots \cap A_{i_r}) + \dots \\ &\quad + (-1)^{n+1} P(A_{i_1} \cap A_{i_2} \cap \dots \cap A_{i_n}) \end{aligned}$$

If the set of events $\{A_1, A_2, \dots, A_n\}$ are **disjoint**, then the probability of the union is simply the sum of the probabilities of each set.

■ **Example 1.9** For three events A_1, A_2 , and A_3 ,

$$\begin{aligned}
 P(A_1 \cup A_2 \cup A_3) &= \sum_{r=1}^n \left((-1)^{r+1} \sum_{i_1 < i_2 < \dots < i_r} P(A_{i_1} \cap A_{i_2} \cap \dots \cap A_{i_r}) \right) \\
 &= (-1)^{1+1} (P(A_1) + P(A_2) + P(A_3)) + \\
 &\quad (-1)^{2+1} (P(A_1 \cap A_2) + P(A_2 \cap A_3) + P(A_3 \cap A_1)) + (-1)^{3+1} P(A_1 \cap A_2 \cap A_3) \\
 &= P(A_1) + P(A_2) + P(A_3) - P(A_1 \cap A_2) - P(A_2 \cap A_3) - P(A_3 \cap A_1) + \\
 &\quad P(A_1 \cap A_2 \cap A_3)
 \end{aligned}$$

■

■ **Example 1.10 — MMSA ex. 15.** A consulting firm presently has bids out on three projects. Let $A_i = \{\text{Award project } i\}$, for $i = 1, 2, 3$, and suppose that $P(A_1) = 0.22$, $P(A_2) = 0.25$, $P(A_3) = 0.28$, $P(A_1 \cup A_2) = 0.11$, $P(A_1 \cap A_3) = 0.05$, $P(A_2 \cap A_3) = 0.07$, and $P(A_1 \cap A_2 \cap A_3) = 0.01$.

Express in words each of the following events, and compute the probability of each event.

a) $A_1^C \cap A_2^C$

Event where they get neither project 1 nor project 2.

$$\begin{aligned}
 P(A_1^C \cap A_2^C) &= P((A_1 \cup A_2)^C) && \text{DeMorgan's Law} \\
 &= 1 - P(A_1 \cup A_2) \\
 &= 1 - (P(A_1) + P(A_2) - P(A_1 \cap A_2)) \\
 &= 1 - (0.22 + 0.25 - 0.11) \\
 &= 0.64
 \end{aligned}$$

b) $A_1^C \cap A_2^C \cap A_3$

Event where they only win project 3.

$A_1^C \cap A_2^C \cap A_3 = (A_1 \cup A_2)^C \cap A_3$ by DeMorgan's Law.

Note that $A_3 = ((A_1 \cup A_2)^C \cap A_3) \cup ((A_1 \cup A_2) \cap A_3)$.

By Axiom 3, $P(A_3) = P(((A_1 \cup A_2)^C \cap A_3)) + P(((A_1 \cup A_2) \cap A_3))$.

$$\begin{aligned}
 \text{That is, } P((A_1 \cup A_2)^C \cap A_3) &= P(A_3) - P((A_1 \cup A_2) \cap A_3) \\
 &= P(A_3) - P((A_1 \cap A_3) \cup (A_2 \cap A_3)) \\
 &= P(A_3) - (P(A_1 \cap A_3) + P(A_2 \cap A_3) - P(A_1 \cap A_2 \cap A_3)) \\
 &= 0.28 - (0.07 + 0.05 - 0.01) \\
 &= 0.17
 \end{aligned}$$

c) What event can be used to represent the outcome the firm is awarded at least two projects?

Find this probability.

$$(A_1 \cap A_2) \cup (A_2 \cap A_3) \cup (A_3 \cap A_1)$$

Let $X = A_1 \cap A_2$, $Y = A_2 \cap A_3$ and $Z = A_3 \cap A_1$.

Then, by the inclusion-exclusion principle, we have

$$\begin{aligned}
 P(X \cup Y \cup Z) &= P(X) + P(Y) + P(Z) - P(X \cap Y) - P(Y \cap Z) - P(Z \cap X) + P(X \cap Y \cap Z) \\
 &= 0.11 + 0.07 + 0.05 - 0.01 - 0.01 - 0.01 + 0.01 \\
 &= 0.21
 \end{aligned}$$

■

2. Counting

2.1 Counting and Probability

Proposition 2.1.1 — Probability as Relative Frequency. The long-run relative frequency of an event A will approach the probability of event A . In cases where the sample space consists of **equally likely** elements, we can find this probability by calculating the relative frequency of A in Ω by

$$P(A) = \frac{\text{Number of outcomes in } A}{\text{Total number of possible outcomes in the random experiment}} = \frac{n(A)}{n(\Omega)}$$

This is valid **only if** each element in Ω is **equally likely**.

2.2 Fundamental Principle of Counting

Experiments that involve equally likely discrete outcomes make calculating probabilities much easier when we have methods to count outcomes.

- For an experiment with two events of interest, A and B , you can use the Inclusion-Exclusion Principle to count the number of outcomes in event $A \cup B$: $n(A \cup B) = n(A) + n(B) - n(A \cap B)$, where $n(X)$ denotes the number of elements in event X .
- For experiments that involve **multiple ordered** stages, we can use the **Fundamental Principle of Counting** (FPC) to count the number of unique outcomes from this multi-stage experiment.

■ **Example 2.1 — Toy Example.** A new sandwich shop seems to only have limited customization options. They offer three types of greens (lettuce, spinach, mixed greens), five types of deli meat, and four types of cheese. Sandwiches are built by layering with greens, followed by deli meat, and topping it off with cheese. How many unique sandwiches can be created if a customer randomly chooses one of each item to include in their sandwich?

This experiment involves 3 ordered stages:

- Stage 1: Pick a green - 3 choices
- Stage 2: Pick a deli meat - 5 choices
- Stage 3: Pick a cheese - 4 choices

Each option is equally likely to be chosen since we're considering all possible combinations
There are in total $3 \times 5 \times 4 = 60$ unique combinations. ■

When counting the number of (ordered) outcomes from a multistage experiment, we can use the Fundamental Principle of Counting.

Theorem 2.2.1 If an experiment consists of m (**ordered**) stages with n_1 possible outcomes in stage 1, n_2 possible outcomes in stage 2, ..., n_m possible outcomes in stage m , then the total number of possible outcomes is

$$\prod_{i=1}^m n_i$$

2.3 Permutations

■ **Example 2.2** FPC counts specifically **ordered stages**. In the toy example, the number of sandwiches only include those that are layered in a specific order: with greens on the bottom, then deli meat, then cheese on top. However, customers might have a preference for how the ingredients are layered.

- How many different ways can the ingredients (greens, deli meat, cheese) be layered / permuted?
 - Stage 1: 3 choices
 - Stage 2: 2 choices
 - Stage 3: 1 choice

There are in total $3 \times 2 \times 1 = 6$ ways.

- If the choice and order of ingredients each result in a 'different' sandwich, how many sandwich choices are there?
6 choices to order ingredients \times 60 ingredient combinations = 360 'different' sandwiches. ■

Definition 2.3.1 — Permutations - ${}_nP_n$. The number of ways to order n **distinct** item is

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

Definition 2.3.2 — Permutations - ${}_nP_k$. The number of ways to select **ordered subset** of k elements from a group of n **distinct** items is

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

The intuition behind this formula is to count all possible arrangements ($n!$) and group together all arrangements that have the same objects in the first k stages. The resulting number of 'groups' is the number of unique ordered subset. The number of elements in each group is equivalent to the number of ways to arrange the remaining $(n-k)$ objects.

■ **Example 2.3** Suppose you are selecting numbers for Lotto 649. You must pick 6 numbers between 1 and 49.

- How many possible winning numbers can be generated if the rules specify that each number can be selected more than once and sequence matters?

$$n(\Omega) = 49 \times 49 \times 49 \times 49 \times 49 \times 49$$

$$= 49^6$$

$$= 13,841,287,201$$

- b) How many possible winning numbers can be generated if the rules specify that each number can only be selected once and sequence matters?

$$\begin{aligned}
 n(\Omega) &= {}_{49}P_6 \\
 &= \frac{49!}{(49-6)!} \\
 &= 49 \times 48 \times 47 \times 46 \times 45 \times 44 \\
 &= 10,068,247,500
 \end{aligned}$$

- c) You win the lottery jackpot if you match all six winning numbers. Under the rules of (b), how likely are you to win the lottery if the six winning numbers must appear in the correct order? In any order?

- If only in correct order, $P(\text{win}) = \frac{n(\text{win})}{n(\Omega)}$

$$\begin{aligned}
 &= \frac{1}{10,068,247,500} \\
 &\approx 9.93 \times 10^{-9}\%
 \end{aligned}$$
- In in any order, $P(\text{win}) = \frac{n(\text{win})}{n(\Omega)}$

$$\begin{aligned}
 &= \frac{6!}{10,068,247,500} \\
 &\approx 7.15 \times 10^{-6}\%
 \end{aligned}$$

■ **Example 2.4** A new student union made up of 5 representatives with different roles is to be established in the next school year. There are 30 candidates applying to become a representative. If each candidate is equally qualified and likely to be elected, how many different groups of representatives can be created?

$$\begin{aligned}
 n(\text{student union}) &= {}_{30}P_5 \\
 &= \frac{30!}{(30-5)!} \\
 &= 30 \times 29 \times 28 \times 27 \times 26 \\
 &= 17,100,770
 \end{aligned}$$

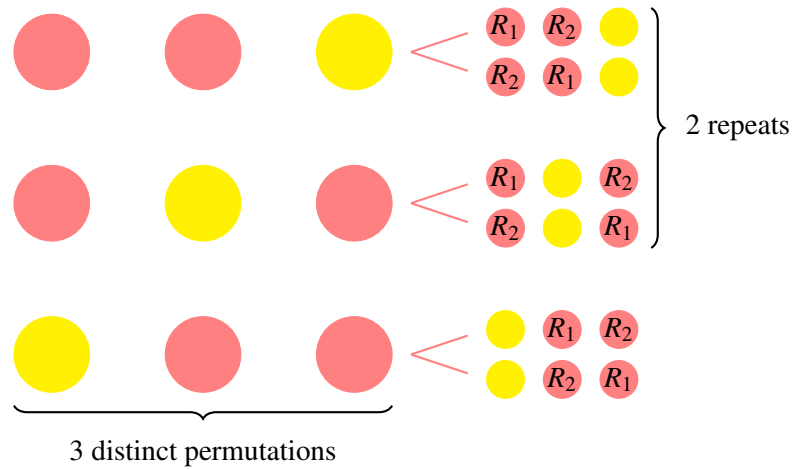
■ **Example 2.5** At a party there are 10 guests. In how many ways can at least two guests have the same birthday? (Assume no leap year births, so 365 days of the year).

$$\begin{aligned}
 P(\text{at least 2}) &= \frac{n(\text{at least 2})}{n(\Omega)} \\
 &= \frac{n(\Omega) - n(\text{distinct birthdays})}{n(\Omega)} \\
 &= \frac{365^{10} - {}_{365}P_{10}}{365^{10}} \\
 &\approx 11.69\%
 \end{aligned}$$

■ **Example 2.6** An urn contains 4 red balls, 1 yellow ball, 3 green balls, and 2 blue balls. How many different ways can you arrange all 10 of these balls?

- We are arranging **all 10 non-distinct** balls. We can use the idea of ${}_nP_k$ to count the number of distinct orderings by taking into account the repeat coloured balls.
- But there's a problem! consider
 (R1) (R2) (Y1) (G1) (G2) (G3) (R3) (R4) (B1) (B2) is the same ordering as
 (R2) (R1) (Y1) (G3) (G2) (G1) (R3) (R4) (B1) (B2)

- How might we adjust for these repeat orderings?
We count the number of ways to rearrange identical balls.



The red balls have 2 repeats, so there are $2!$ ways to rearrange the identical red balls. This gives a total of $\frac{3!}{2!} = 3$ ways to arrange the three balls. ■

2.4 Combinations

Definition 2.4.1 — Combinations - ${}_nC_k$. The number of ways to select an *unordered* subset of k items from a group of n *distinct* items without replacement is

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

Similar to the intuition for permutations, we divide by $(n-k)!$ to remove all the different ways of ordering the remaining $(n-k)$ items. However, for every ${}_nP_k$ ordering of distinct objects, there exists $k!$ orderings of the same collection of k objects. Thus, divide $k!$ to get the number of unique groupings of k objects.

■ **Example 2.7** A new student union made up of 5 representatives with equal roles is to be established in the next school year. There are 30 candidates applying to become a representative. 20 of the candidates are senior students, while 10 are junior students. If each candidate is equally qualified and likely to be elected, how many different student unions can be formed? How likely is it that the union comprises of all senior candidates?

$$\begin{aligned} \text{We are selecting 5 students where order is irrelevant, so } n(\text{student union}) &= {}_{30}C_5 \\ &= \frac{30!}{(30-5)! \cdot 5!} \\ &= 142,510 \end{aligned}$$

$$\text{The possibility that all the candidates are seniors, } P(\text{all seniors}) = \frac{n(\text{all seniors})}{n(\Omega)}$$

$$\begin{aligned} &= \frac{{}_{20}C_5}{{}_{30}C_5} \\ &= \frac{15,504}{142,506} \\ &\approx 10.88\% \end{aligned}$$

■ **Example 2.8** Consider a standard deck of 52 cards.

- a) How many different hand of 5 cards can be drawn from the 52 cards?

Order is irrelevant, so $n(\text{hand of 5}) = {}_{52}C_5$

$$\begin{aligned} &= \frac{52!}{(52-5)! \cdot 5!} \\ &= 2,598,960 \end{aligned}$$

- b) How many different hands will have 4 face cards and 1 numeric card?

- Stage 1: choose 4 from the 12 face cards, ${}_{12}C_4$
- Stage 2: choose 1 from the 36 numeric cards, ${}_{36}C_1$

In total, by FPC, we have ${}_{12}C_4 \times {}_{36}C_1 = 17,820$ different hands.

■

■ **Example 2.9** A two-pair in poker is a five card hand consisting of a two pairs of two distinct ranks, and a single card of third rank. Recall that a standard deck of 52 card has 13 ranks: A, 2, 3, ..., 10, Jack, Queen, King. Each rank comes in four suits: ♣, ♠, ♦, ♥. An example of a two-pair hand is $9\clubsuit 9\heartsuit 8\spadesuit 2\diamondsuit 2\clubsuit$. How many distinct two pair poker hands are there?

- Pick two ranks, ${}_{13}C_2$
- Pick two suits for each of the ranks, ${}_4C_2$ each
- Pick the last card, ${}_{44}C_1$

In total, by FPC, we have ${}_{13}C_2 \times ({}_4C_2)^2 \times {}_{44}C_1 = 123,552$ two pair hands.

■

■ **Example 2.10** Simple comparison sort algorithms used to sort a list of integers do so by comparing two integers, checking which is larger, and swapping the elements according to size. For example in bubble sort, the algorithm iterates through a list of integers repeatedly, compares two consecutive integers at a time, swaps their positions if necessary. It repeats this until all integers are in order. Suppose we want to sort a list from least to greatest:

- a) You might have learned that the number of comparisons in required in the worst case will grow proportionally to $n \log_{10}(n)$ as the number of inputs n increases. For a given list of 8 distinct integers, what would the worst case sorting outcome look like?

There are $8! = 40,320$ ways for 8 integers to be arranged.

- In worst case, each time we compare 2 inputs, we are left with many permutations to sort through
 - We are left with some permutations that has the most number of comparisons remaining.
- b) Of the number of possibilities in a), only one results in the desired sorted list. In the worst case, how many comparisons do we need at least to sort n elements?

Let C be the maximum number of comparisons to sort the worst case.

C is the height of the decision tree, where each node have two branches (for $x < y$): Yes or No.

That is,

$$2^C \geq n!$$

$$2^C \geq 8!$$

$$\log_2(2^C) \geq \log_2(8!)$$

$$C \geq \log_2(8!) \approx 15.29$$

■

3. Conditional Probability

3.1 Conditional Probability

■ **Example 3.1 — Thought Exercise.** For the following, assume that the probability of having a child of either sex (male or female) is 50%.

- a) A family has two children. What are the chances this family has two boys?

$$\Omega = \{BB, BG, GB, GG\}$$

$$P(BB) = \frac{1}{4}$$

- b) A family has two children, and you know that one of the children is a boy. What are the chances this family has two boys?

$$\Omega = \{BB, BG, GB\}$$

$$P(2 \text{ boys if one boy}) = \frac{1}{3}$$

■

■ **Example 3.2** Below is a contingency table of counts in a fictional study of colourblindness among the two sexes. C denotes the event that a surveyed individual is colourblind, and M denotes the event that a surveyed individual is male.

	C	C^c	Row Totals
M	106	1175	1281
M^c	7	1212	1219
Column Totals	113	2397	2500

- a) What is the estimated probability that an individual is male and colourblind?

$$\begin{aligned} P(M \cap C) &= \frac{n(M \cap C)}{n(\Omega)} \\ &= \frac{106}{2500} \\ &\approx 4.52\% \end{aligned}$$

- b) What is the estimated probability that a male individual is colourblind? That a female individual is colourblind?

$$\begin{aligned}
 P(C \text{ if } M) &= \frac{n(C \cap M)}{n(M)} \\
 &= \frac{106}{1281} \\
 &\approx 8.27\%
 \end{aligned}$$

\therefore 8.27% of males surveyed reported to be colourblind.

$$\begin{aligned}
 P(C \text{ if } M^c) &= \frac{n(C \cap M^c)}{n(M^c)} \\
 &= \frac{7}{1219} \\
 &\approx 0.57\%
 \end{aligned}$$

\therefore 0.57% of females surveyed reported to be colourblind. ■

Definition 3.1.1 — Conditional Probability. The notation $P(A|B)$ denotes the probability that of event A under the condition that event B is known to have occurred.

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

Rearranging the above, we have the following relations we can use depending on available information

$$P(A \cup B) = P(A|B) \cdot P(B) \quad \text{or} \quad P(A \cup B) = P(B|A) \cdot P(A)$$

Conditional probabilities provide additional information when we know partially the outcome of a random experiment. Conditional probabilities are probability distributions on a **restricted sample space**, and follow the same probability axioms.

axiom 3.1.1 Consider a random experiment with sample space Ω . Let B be an event ($B \subseteq \Omega$) with $P(B) > 0$. Let b denote the elements of event B . Then,

1. $P(b|B) \geq 0$ for all $b \in B$.
2. $\sum_{b \in B} P(b|B) = 1$
3. For $A \subseteq B$, $P(A|B) = \sum_{b \in A} P(b|B)$

Some examples include

- determining the probability distribution of disease status when someone tests negative for a disease
- applications in Bayesian statistics, Bayes classification, etc.

■ **Example 3.3** You pick a card at random from a standard deck of cards. Define events Q where a queen of hearts is drawn and R where a red card is drawn. Describe the events below and find their probabilities.

a) $Q|R$

$$\begin{aligned}
 P(Q|R) &= \frac{P(Q \cap R)}{P(R)} \\
 &= \frac{1/52}{26/52} \\
 &= \frac{1}{26}
 \end{aligned}$$

b) $R|Q$

$$\begin{aligned}
 P(R|Q) &= \frac{P(R \cap Q)}{P(Q)} \\
 &= \frac{1/52}{1/52} \\
 &= 1
 \end{aligned}$$

$$\begin{aligned}
 \text{c) } Q^c|R \\
 P(Q^c|R) &= 1 - P(Q|R) \\
 &= 1 - \frac{1}{26} \\
 &= \frac{25}{26}
 \end{aligned}$$

■

3.2 Independence

Recall that two events are independent if the occurrence of one (A) does not alter the chances of the other event (B). Formally, we have the following.

Definition 3.2.1 — Independent Events. Two events A and B are *independent* if

$$P(A|B) = P(A) \quad \text{provided that } P(B) > 0$$

$$P(B|A) = P(B) \quad \text{provided that } P(A) > 0$$

Using this, we can show that event A is independent of event B , then $P(A \cap B) = P(A) \cdot P(B)$. Otherwise, the two events are dependent.

Definition 3.2.2 — Mutually Exclusive. Two events are *mutually exclusive* if the occurrence of one (A) excludes the occurrence of the other (B). Event-wise, the two sets are disjoint ($A \cap B = \emptyset$) and $P(A \cap B) = 0$. This implies that the events are dependent.

For example, given $P(A) > 0$ and $P(B) > 0$, A and B are mutually exclusive when $P(A|B) = P(B|A) = 0$.

If events A and B are independent, then so are their complements A^c and B^c .

For a collection of n events, A_1, A_2, \dots, A_n ,

- If all n events are independent, then

$$P(A_1 \cap A_2 \cap \dots \cap A_n) = P(A_1) \times P(A_2) \times \dots \times P(A_n)$$

- A_1, \dots, A_n are **mutually independent** if for any subset of k events, $k = 2, 3, \dots, n$

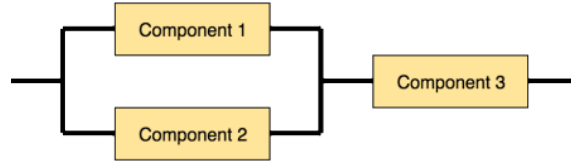
$$P(A_{i_1} \cap A_{i_2} \cap \dots \cap A_{i_k}) = P(A_{i_1}) \times P(A_{i_2}) \times \dots \times P(A_{i_k})$$

■ **Example 3.4** Two events E and F have the properties $P(E) = 0.44$, $P(F) = 0.6$, and $P(E \cap F) = 0.35$. Use the information to determine whether events E and F are independent.

$$\begin{aligned}
 P(E|F) &= \frac{P(E \cap F)}{P(F)} & P(E) \times P(F) &= 0.44 \times 0.6 \\
 &= \frac{0.35}{0.6} & &= 0.264 \\
 &\approx 0.5933
 \end{aligned}$$

Clearly, $P(E|F) \neq P(E) \times P(F)$. E and F are NOT independent. ■

■ **Example 3.5** A system below is made of independent components. The probability that the first component works is 0.9, 0.95 for the second component, and 0.99 for the third component. The signal can travel from left to right if there is a circuit made of working components. Find the probability that the signal is blocked.



Let B be the event that the signal is blocked. Let C_1 , C_2 and C_3 be the events that components 1, 2, and 3 are working.

$$\begin{aligned}
 P(B) &= P(C_3^c \cup (C_1^c \cap C_2^c)) \\
 &= P(C_3^c) + P(C_1^c \cap C_2^c) - P(C_1^c \cap C_2^c \cap C_3^c) \\
 &= P(C_3^c) + P(C_1^c) \times P(C_2^c) - P(C_1^c) \times P(C_2^c) \times P(C_3^c) \\
 &= 0.01 + 0.1 \times 0.05 - 0.1 \times 0.05 \times 0.01 \\
 &= 1.495\%
 \end{aligned}$$

Alternatively, we can use the indirect method.

$$\begin{aligned}
 P(B) &= 1 - P(B^c) \\
 &= 1 - P(C_3 \cap (C_1 \cup C_2)) \\
 &= 1 - P((C_3 \cap C_1) \cup (C_3 \cap C_2)) && \text{distribution law} \\
 &= 1 - (P(C_1 \cap C_3) + P(C_2 \cap C_3) - P(C_1 \cap C_2 \cap C_3)) \\
 &= 1 - (P(C_1) \times P(C_3) + P(C_2) \times P(C_3) - P(C_1) \times P(C_2) \times P(C_3)) \\
 &= 1 - (0.9 \times 0.99 + 0.95 \times 0.99 - 0.9 \times 0.95 \times 0.99) \\
 &= 1 - 0.98505 \\
 &= 1.495\%
 \end{aligned}$$

3.3 Bayes' Rule

Given two events, we know that their conditional probability is defined as

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

which can be rearranged to

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

Suppose now we have a sample space consisting **only** of events A, B_1, B_2, \dots, B_k where the B_i 's **partition the sample space**. That is, the B_i 's are disjoint and $\bigcup_{i=1}^k B_i = \Omega$. From Axiom 3, we can show that

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

Combining this with our knowledge of conditional probability, we get the **Law of Total Probability**.

Theorem 3.3.1 — Law of Total Probability. If B_1, B_2, \dots, B_k is a collection of mutually exclusive (disjoint) and exhaustive events that partition the sample space, then for any event A ,

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

Putting together the Law of Total Probability and definition of conditional probability, we get **Bayes' Rule**.

Theorem 3.3.2 — Bayes' Rule. Let B_1, B_2, \dots, B_k form a partition of the sample space and let A be an event in Ω . Then

$$\begin{aligned} P(B_i|A) &= \frac{P(A \cap B_i)}{P(A)} \\ &= \frac{P(A|B_i) \cdot P(B_i)}{\sum_{i=1}^k P(A|B_i) \cdot P(B_i)} \end{aligned}$$

■ **Example 3.6** A ball is drawn at random from an urn containing one red and one white ball. If the white ball is drawn, it is put back into the urn. If the red ball is drawn, it is returned to the urn together with two more red balls. Then a second draw is made.

- a) What is the probability a red ball was drawn on both the first and the second draws?

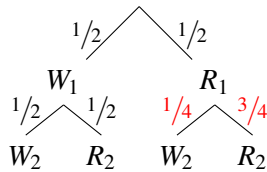
Let W_i be the event where a white ball is drawn on the i^{th} draw.

Let R_i be the event where a red ball is drawn on the i^{th} draw.

$$P(R_1 \cap R_2) = P(R_1) \times P(R_2|R_1)$$

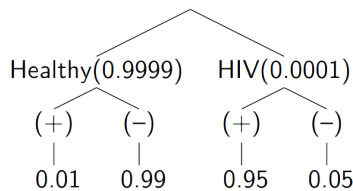
$$\begin{aligned} &= \frac{1}{2} \times \frac{3}{4} \\ &= \frac{3}{8} \end{aligned}$$

- b) What is the probability that a red ball was drawn first if the second ball drawn was red?



$$\begin{aligned} P(R_1|R_2) &= \frac{P(R_1 \cap R_2)}{P(R_2)} \\ &= \frac{3/8}{P(R_1 \cap W_1) + P(R_2 \cap R_1)} \\ &= \frac{3/8}{P(R_2|W_1) \times P(W_1) + P(R_2|R_1) \times P(R_1)} \\ &= \frac{3/8}{1/2 \times 1/2 + 3/4 \times 1/2} \\ &= \frac{3}{5} \\ &= 60\% \end{aligned}$$

■ **Example 3.7** Consider HIV testing. An HIV test will correctly test positive 95% of the time (sensitivity), and incorrectly test positive 1% (false positive rate) of the time. Suppose we know that 99.99% of the population is HIV-free. Under these conditions, what is the probability that a patient who tested positive is actually HIV positive?



$$\begin{aligned} P(\text{HIV}|+) &= \frac{P(\text{HIV} \cap +)}{P(+)} \\ &= \frac{P(+ \cap \text{HIV})}{P(+ \cap \text{Healthy}) + P(+ \cap \text{HIV})} \\ &= \frac{0.95 \times 0.0001}{0.01 \times 0.9999 + 0.95 \times 0.0001} \\ &= 0.94\% \end{aligned}$$

∴ If we were randomly pick a person and test them for HIV, there chance of them having HIV with a positive test is 0.94%. ■

4. Discrete Distributions

4.1 Random Variables

4.1.1 Introduction to Random Variables

Definition 4.1.1 — Random Variable. A *random variable* is a real-valued function that assigns a numerical value to each event in the sample space Ω arising from a random experiment. A random variable X is a **real-valued function** $X : \Omega \rightarrow \mathbb{R}$ such that for every $\omega \in \Omega$, $X(\omega) = x \in \mathbb{R}$. It is a mapping from the sample space to the real numbers.

■ **Example 4.1** Consider the random experiment of tossing a coin.

- $\Omega = \{H, T\}$
- Let X be the RV denoting the outcome of a toss. We can define X such that $X(H) = 1$, $X(T) = 0$ essentially converting each outcome into a number.
- By convention, we will denote random variables with capital letters, and a particular (unknown value) of a random variable with its lower case equivalent. i.e. for a random variable X , a particular value of this RV would be denoted by x .

■

Definition 4.1.2 — Discrete Random Variable. A *discrete* random variable X is one that can take on only a finite number or a countably infinite number of possible values x . A random variable X is *continuous* if its domain is an interval of real numbers.

Definition 4.1.3 — Probability Mass Function. A *probability mass function* (PMF) of a discrete random variable is one that assigns a probability to each value $x \in C$ such that

- $0 \leq P(X = x) \leq 1$
- $\sum_{x \in C} P(X = x) = 1$

■ **Example 4.2** Below are some examples of random variables.

Discrete RV Examples

- The number of defects in a day's production of car parts
- The number of new arrivals in a queue

- The status of your internet service: online or offline
- The number of students online at a particular time

Continuous RV Examples

- The weight of a randomly selected individual
- The time it takes to load a video
- The temperature in the morning of a random day

■ **Example 4.3** Determine the value of k such that $f(x) = \frac{kx^2 - x + 2}{4}$ will be a valid probability mass function for $X = \{0, 1, 2, 3, 4\}$.

X	0	1	2	3	4
$P(X = x)$	$\frac{2}{4}$	$\frac{k+1}{4}$	$\frac{4k}{4}$	$\frac{9k-1}{4}$	$\frac{16k-2}{4}$

We need $\sum_{x=0}^4 P(X = x) = 1$. That is, $\frac{2 + (k+1) + (4k) + (9k-1) + (16k-2)}{4} = 1$

$$30k = 4$$

$$k = \frac{2}{15}$$

Thus, k must be $\frac{2}{15}$ for $\sum_{x \in X} P(X = x) = 1$.

X	0	1	2	3	4
$P(X = x)$	$\frac{1}{2}$	$\frac{17}{60}$	$\frac{2}{15}$	$\frac{1}{20}$	$\frac{1}{30}$

■ **Example 4.4** A factory producing computer parts sends out a shipment of 10 parts of which 3 are defective. Find the probability mass function for the number of defectives a customer will get if the first customer randomly purchases 4 computer parts.

Let D be the number of defectives purchased.

$$D = \{0, 1, 2, 3\}.$$

- $P(D = 0) = \frac{{}^7C_4}{{}^{10}C_4}$
- $P(D = 1) = \frac{{}^3C_1 \times {}^7C_3}{{}^{10}C_4}$
- $P(D = 1) = \frac{{}^3C_2 \times {}^7C_2}{{}^{10}C_4}$
- $P(D = 1) = \frac{{}^3C_3 \times {}^7C_1}{{}^{10}C_4}$

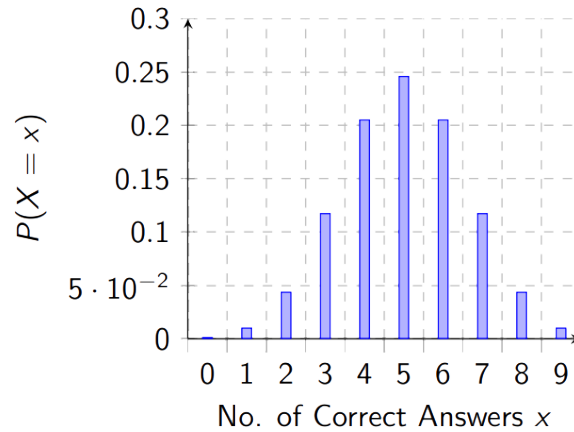
We may conclude that $P(D = d) = \frac{{}^3C_d \times {}^7C_{4-d}}{{}^{10}C_4}$.

■ **Example 4.5** A quiz consists of 10 true / false problems. A student takes the quiz by randomly selecting the answers. Examine the graph of the probability mass function and describe the behaviour of the random number of correct answers X .

For discrete distribution, we use bar charts. The bins are $0, 1, 2, 3, \dots, 10$, which are discrete. Histograms, on the other hand, has bins that span a continuous interval, for example, $[0, 1), [1, 2), [2, 3), \dots$

- 4, 5, 6 have the highest probability mass, meaning they are the most *likely*;
- 0, 1, 2, 8, 9, 10 have the lowest probability mass, meaning that they are the most *unlikely*;

- 3 and 7 lie between *likely* and *unlikely*.



■

4.1.2 Characteristics of Random Variables

Both the PMF and CDF describe the exact distribution of a random quantity. Distributions have two useful characteristics that are often used in statistics:

- **Expected Value:** The long run / theoretical average. If a random experiment were to be conducted n times, then as $n \rightarrow \infty$ then the average of outcomes converges to the expected value. This is often denoted as μ .
- **Variance** or **Standard Deviation:** Are measures of the spread and variability of a random variable. Standard deviation is the square root of variance. Variance is often denoted as σ^2 and SD as σ .

Definition 4.1.4 — Expected Value. The *expected value* of a discrete random quantity X is defined to be

$$\mu = E(X) = \sum_{x \in X} x \cdot P(X = x)$$

For a given transformation of X , $g(X)$, $E(g(x))$ can be found by

$$E(g(x)) = \sum_{x \in X} g(x) \cdot P(X = x)$$

Note that $E(g(x)) \neq g(E(X))$ except when $g(X)$ is a linear transformation.

Review of Summation Rules

- $\sum_{i=1}^n c = n \cdot c$ for a constant c
- $\sum_{i=1}^n c \cdot x_i = c \cdot \left(\sum_{i=1}^n x_i \right)$
- $\sum_{i=1}^n i = \frac{n \cdot (n+1)}{2}$

■ **Example 4.6** A mining company needs a type of drill bit for a project. It is known through historical data that these drill bits for similar projects will last 2, 4, or 7 hours with probabilities 0.1, 0.7, and 0.2. How long do they expect each drill bit to last, on average?

Let L be the longevity of a random drill bit.

l	2	4	7
$P(L=l)$	0.1	0.7	0.2

$$\begin{aligned}
E(L) &= \sum_{l \in L} l \cdot P(L = l) \\
&= 2 \cdot 0.1 + 4 \cdot 0.7 + 7 \cdot 0.2 \\
&= 4.4 \text{ hours}
\end{aligned}$$

∴ on average, the drill bits will last 4.4 hours. ■

Properties of Expectation

For any constants a , b , and discrete variables X , the following are true.

- $E(a) = 0$
- $E(X + a) = E(X) + a = \mu + a$.
Increase in $x \in X$ will shift the centre / average by the same amount.
- $E(aX) = a \cdot E(X) = a \cdot \mu$

$$\begin{aligned}
\text{Proof. Take } g(X) &= aX. \text{ Then, } E(g(x)) = \sum_{x \in X} g(x) \cdot P(X = x) \\
&= \sum_{x \in X} ax \cdot P(X = x) \\
&= a \sum_{x \in X} x \cdot P(X = x) \\
&= a \cdot E(x) \\
&= a \cdot \mu
\end{aligned}$$

That is, $E(aX) = a \cdot E(X) = a \cdot \mu$. ■

- $E(aX + b) = a \cdot E(X) + b = a \cdot \mu + b$

Proof. Take $g(X) = aX + b$. $g(x)$ is a linear transformation of X .
 $E(g(x)) = g(E(x))$ if $g(X)$ is a linear transformation of X .
 Thus, $E(aX + b) = a \cdot E(X) + b = a \cdot \mu + b$. ■

- $E(X + Y) = E(X) + E(Y)$.
- $E(XY) \neq E(x) \cdot E(Y)$ **unless** X and Y are *independent*.

Definition 4.1.5 — Variance. For a discrete variable X , the *variance* of X is defined to be

$$\sigma^2 = V(X) = E((X - \mu)^2) = \sum_{x \in X} (x - \mu)^2 \cdot P(X = x)$$

Variance captures the spread in *units*². The standard deviation, $\sigma = \sqrt{\text{variance}}$ is a measure of spread in the same units as the random variable X .

■ **Example 4.7** A mining company needs a type of drill bit for a project. It is known through historical data that these drill bits for similar projects will last 2, 4, or 7 hours with probabilities 0.1, 0.7, and 0.2. Find the variance and standard deviation in the longevity of this type of drill bit. Interpret the values

Let L be the longevity of a random drill bit.

We also found $\mu = 4.4$

l	2	4	7
$P(L = l)$	0.1	0.7	0.2
$(l - \mu)^2$	$(2 - 4.4)^2 = 5.76$	$(4 - 4.4)^2 = 0.16$	$(7 - 4.4)^2 = 6.76$

$$\begin{aligned}
V(L) &= \sigma_L^2 \\
&= \sum_{l \in L} (l - \mu)^2 \cdot P(L = l) \\
&= 5.76 \cdot 0.1 + 0.16 \cdot 0.7 + 6.76 \cdot 0.2 \\
&= 2.04 \text{ hours}^2
\end{aligned}$$

$$\begin{aligned}
\text{Then, } SD(L) &= \sqrt{V(L)} = \sigma_L \\
&= \sqrt{2.04} \\
&\approx 1.4283 \text{ hours}
\end{aligned}$$

On average, the longevity of the drill bits will vary by about 1.43 hours from the average. Typically, we expect the longevity to be between $(4.4 - 1.4283, 4.4 + 1.4283) = (2.97, 5.83)$ hours.

■

Properties of Variance

For any constants a , b , and discrete variables X , the following are true.

- $V(a) = 0$
Constants do not vary.
- $V(X + a) = V(X) = \sigma^2$
If the spread of X is $V(X)$, increasing each $x \in X$ will not change how spread out the random variable is.
- $V(aX) = a^2 \cdot V(X) = a^2 \cdot \sigma^2$

$$\begin{aligned}
\text{Proof. } V(aX) &= E((aX - E(aX))^2) \\
&= E((ax - a\mu)^2) \\
&= E((a \cdot (x - \mu))^2) \\
&= E(a^2 \cdot (x - \mu)^2) \\
&= a^2 \cdot E((x - \mu)^2) \\
&= a^2 \cdot V(X) = a^2 \cdot \sigma^2
\end{aligned}$$

That is, $V(aX) = a^2 \cdot V(X) = a^2 \cdot \sigma^2$. ■

- $V(aX + b) = a^2 \cdot V(X) = a^2 \cdot \sigma^2$
- $V(X + Y) \neq V(X) + V(Y)$ **unless** X and Y are *independent*.

■ **Example 4.8** A mining company needs a type of drill bit for a project. It is known through historical data that these drill bits for similar projects will last 2, 4, or 7 hours with probabilities 0.1, 0.7, and 0.2.

- a) If they ordered 10 drill bits of the same type for replacement once one drill bit fails, how long can they expect these drill bits to last for this project?

Let L be the longevity of drill bits.

$$\begin{aligned}
E(L_1 + L_2 + \cdots + L_{10}) &= \sum_{i=1}^{10} E(L_i) \\
&= 4.4 + 4.4 + \cdots + 4.4 \\
&= 44 \text{ hours}
\end{aligned}$$

- b) Find the variance and standard deviation in the longevity for the 10 drill bits that were ordered.

All drill bits are independent, so $V(L_1 + L_2 + \cdots + L_{10}) = \sum_{i=1}^{10} V(L_i)$

$$= 2.04 + 2.04 + \cdots + 2.04$$

$$= 20.4 \text{ hours}^2$$

Then, $SD(L_1 + L_2 + \cdots + L_{10}) = \sqrt{20.4}$

$$\approx 4.52 \text{ hours}$$

On average, they can expect the drill bits to last 44 ± 4.52 hours. That is, from 39.48 hours to 48.52 hours.

■

An alternative method for calculating the variance of a discrete random variable X with PMF $f(x)$ can be derived as follows:

$$\begin{aligned} E((X - \mu)^2) &= E(X^2 - 2X\mu + \mu^2) \\ &= E(X^2) - 2\mu \cdot E(X) + \mu^2 \\ &= E(X^2) - 2E(X) \cdot E(X) + E(X)^2 \\ &= E(X^2) - E(X)^2 \end{aligned}$$

This breakdown is “allowed” since $\mu = E(X)$ is an unknown **constant**.

Note that $E(X^2) \neq E(X)^2$. To compute $E(X^2)$, refer back to the definition of expected value. That is,

$$E(X^2) = \sum_{x \in X} x^2 \cdot f(x)$$

4.2 Cumulative Distribution Function

The probability behaviour of a random variable can be represented in many ways, such as with the probability mass function. Another representation is with the *cumulative distribution function*.

Definition 4.2.1 — Cumulative Distribution Function. The *cumulative distribution function* (CDF) $F(x)$ of a discrete random variable with probability mass function $P(x)$ or $f(x)$ is a function that returns the cumulative (total) probability up to and including $X = x$.

$$F(b) = P(X \leq b) = \sum_{x \in \{x \leq b\}} P(x)$$

The domain of the CDF is always over the set of real numbers! As such, CDFs are often represented as a piecewise function.

■ **Example 4.9** Find the cumulative distribution function for PMF below:

x	0	1	2	3
$P(X = x)$	$\frac{1}{6}$	$\frac{1}{2}$	$\frac{3}{10}$	$\frac{1}{30}$

$$F(x) = \begin{cases} 0 & \text{if } x < 0 \\ \frac{1}{6} & \text{if } 0 \leq x < 1 \\ \frac{2}{3} & \text{if } 1 \leq x < 2 \\ \frac{29}{30} & \text{if } 2 \leq x < 3 \\ 1 & \text{if } x \geq 3 \end{cases}$$

■

4.2.1 Properties of CDF

CDF of a Discrete Random Variable

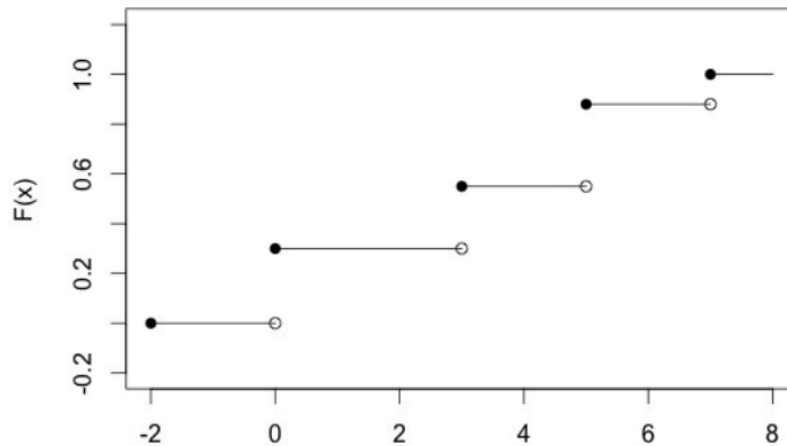
For a discrete random variable X with CDF $F(X)$:

1. The graph of the CDF will be a **non-decreasing step-function**. That is for $a < b$, $F(a) \leq F(b)$.
2. The graph of the CDF is **right continuous**. That is, $\lim_{x \rightarrow c^+} F(x) = F(c)$.
3. $\lim_{x \rightarrow \infty} F(x) = 1$
4. $\lim_{x \rightarrow -\infty} F(x) = -1$

■ **Example 4.10** A discrete random variable X has cumulative distribution function defined by

$$F(X) = \begin{cases} 0 & x < 0 \\ 0.3 & 0 \leq x < 3 \\ 0.55 & 3 \leq x < 5 \\ 0.88 & 5 \leq x < 7 \\ 1 & x \geq 7 \end{cases}$$

a) Plot the CDF.



b) What is the probability that X is less than 5?

$$P(x < 5) = 0.55$$

c) What is the probability of $X = 2$?

$$\begin{aligned} P(X = 2) &= P(X \leq 2) - P(X < 2) \\ &= 0.3 - 0.3 \\ &= 0 \end{aligned}$$

Note that $P(X = 2) = 0$ has no changes to CDF on interval of $[0, 3)$. The only outcomes with probability mass are $X = 0$ and $X = 3$.

■

4.2.2 Chebyshev's Inequality

For a given random variable X , μ_X and σ_X^2 are measures of two features of the distribution of X : it's 'centre' and it's spread. How can we use these two values to better understand the distribution of X , especially in the absence of the exact distribution such as the probability mass function (PMF) or the cumulative distribution function (CDF)?

Theorem 4.2.1 — Markov's Inequality. Let X be a **non-negative** random variable with mean $E(X)$. Then, for some constant $a > 0$,

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

Proof. Since X is a non-negative random variable, $\forall x \in X, x \geq 0$.

$$\begin{aligned} \text{Then, } E(X) &= \sum_{x \in X} x \cdot P(X = x) \\ &= \sum_{x < a} x \cdot P(X = x) + \sum_{x \geq a} x \cdot P(X = x) \\ &\geq \sum_{x \geq a} x \cdot P(X = x) \\ &\geq \sum_{x \geq a} a \cdot P(X = x) && \text{since } x \geq a \\ &= a \cdot \sum_{x \geq a} P(X = x) \\ &= a \cdot P(x \geq a) \end{aligned}$$

That is, $E(X) \geq a \cdot P(X \geq a)$ ■

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

Theorem 4.2.2 — Chebyshev's Inequality. Let X be a random variable with mean (expected value) μ and finite variance σ^2 . Then for any positive k ,

$$P(|x - \mu| < k\sigma) \geq 1 - \frac{1}{k^2}$$

Chebyshev's Inequality applies to all discrete distributions with **finite** $E(X)$ and $V(X)$ for the random variable X .

$$\begin{aligned} \text{Proof. } P(|X - \mu| < k\sigma) &= P((X - \mu)^2 < k^2 \sigma^2) && \text{since RV's are non-negative and } k > 0 \\ &= 1 - (P((X - \mu)^2 \geq k^2 \sigma^2)) \end{aligned}$$

By Markov's Inequality, for a non-negative random variable X and a positive constant a , $P(X \geq a) \leq \frac{E(x)}{a}$, $a > 0$. Consider $(X - \mu)^2 > 0$ as x and $k^2 \sigma^2 > 0$ as a in Markov's Inequality,

$$\begin{aligned} \text{we have } P((X - \mu)^2 \geq k^2 \sigma^2) &\leq \frac{E((X - \mu)^2)}{k^2 \sigma^2} = \frac{\sigma^2}{k^2 \sigma^2} \\ &= \frac{1}{k^2} \end{aligned}$$

$$\text{That is, } P(|X - \mu| < k\sigma) \geq 1 - \frac{1}{k^2}. \quad \blacksquare$$

■ **Example 4.11** Based on past data, the average daily number of tech support requests at a local call centre is 115 with a standard deviation of 10 calls.

a) What can be said about the fraction of days on which the number of calls received is between 100 and 130?

- Distribution info: missing
- We are given: $\mu = 115, r = 10$

Let C be the random number of daily calls.

$$\begin{aligned}
P(100 \leq C \leq 130) &= P(-15 \leq C - 115 \leq 15) \\
&= P(-15 \leq C - \mu \leq 15) \\
&= P(|C - \mu| \leq 15) \\
&= P(|C - \mu| < 16) \\
&= P(|C - \mu| < 1.6\sigma) \\
&\geq 1 - \frac{1}{1.6} = 0.6094
\end{aligned}$$

\therefore At least 60.94% of the time they will have between 100 to 130 calls a day.

- b) What number of calls can they expect to receive at least 90% of the time?

We want to find the number of σ that correspond to an interval that has at least 90% chance of occurring.

$$P(|X - \mu| < k\sigma) = 1 - \frac{1}{k^2} = 90\%$$

$$\text{That is, } 1 - \frac{1}{k^2} = 0.90$$

$$0.1 = \frac{1}{k^2}$$

$$k^2 = 10$$

$$k = \sqrt{10}$$

For all distributions, more than 90% of the outcomes lie within $\sqrt{10} \approx 3.16$ standard deviation of its mean.

$$(115 - k\sigma, 115 + k\sigma) = (115 - \sqrt{10} \cdot 10, 115 + \sqrt{10} \cdot 10) \approx (83.38, 146.62)$$

\therefore They can expect to receive [83, 147] number of calls at least 90% of the time.

■

4.3 Common Discrete Distributions

4.3.1 Binomial Distribution

Definition 4.3.1 — Bernoulli Trials. A *Bernoulli trial* is a random experiment consisting of exactly one trial involving two possible outcomes, often called a *success* or a *failure*. Let X be the outcome of a Bernoulli trial where

$X = 0$ if the outcome is a failure

$X = 1$ if the outcome is a success

We define p to be the probability of **success**, and $q = 1 - p$ to be the probability of **failure**. The *probability mass function* is then

$$f(x) = p^x \cdot (1 - p)^{1-x}$$

■ **Example 4.12 — Bernoulli Trials.** Below are some examples of Bernoulli trials.

- Whether a randomly selected part is defective
- Whether there is an error in a line of code
- Whether a randomly selected individual is taller than 5'7"
- Whether a switch is in the on or off proposition
- Whether your lotto ticket is the winning number

■

■ **Example 4.13** Consider a multiple choice quiz 4 questions. A student selects an answer at random for each question, and each question is a Bernoulli experiment: the student either guesses

correctly (1) or incorrectly (0). The sum of these four Bernoulli experiments will then be the random number of correct answers for a student that completes a similar quiz in this way. Our sample space is:

$$\Omega = \{\checkmark\checkmark\checkmark\checkmark, \checkmark\checkmark\checkmark\text{X}, \checkmark\checkmark\text{X}\checkmark, \checkmark\text{X}\checkmark\checkmark, \text{X}\checkmark\checkmark\checkmark, \checkmark\checkmark\text{XX}, \checkmark\text{X}\checkmark\text{X}, \text{X}\checkmark\checkmark\text{X}, \text{X}\checkmark\text{X}\checkmark, \text{XX}\checkmark\checkmark, \checkmark\text{XX}\checkmark, \\ \text{XXXX}\checkmark, \text{XX}\checkmark\text{X}, \text{X}\checkmark\text{XX}, \checkmark\text{XXX}, \text{XXXX}\}$$

Suppose these MCQs have four options each, and each question only has one correct option. Find the following probabilities.

- a) Guessing each question correctly.

$$P(\checkmark\checkmark\checkmark\checkmark) = \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{1}{4} = \frac{1}{4^4}$$

- b) Guessing each question incorrectly.

$$P(\text{XXXX}) = \frac{3}{4} \cdot \frac{3}{4} \cdot \frac{3}{4} \cdot \frac{3}{4} = \frac{3^4}{4^4}$$

- c) Guessing exactly 2 questions correctly.

$$P(\checkmark\checkmark\text{XX}) = \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{3}{4} \cdot \frac{3}{4} = \frac{3^2}{4^4}$$

$$P(2 \text{ correct}) = {}_4C_2 \cdot \frac{3^2}{4^4}$$

■

Often, we are interested in modeling the number of successes among multiple trials instead of the results of a single trial:

Definition 4.3.2 — Binomial Distribution. A *Binomial experiment* consists of n independent and identical Bernoulli trials. The probability of success, p , is fixed for each trial.

Let X be the random variable representing the number of successes among the n trials. Then X can be modeled by the binomial distribution with parameters n and p , denoted as $X \sim \text{Bin}(n, p)$. The binomial distribution has probability mass function:

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

If $X \sim \text{Bin}(n, p)$, we can show that $E(X) = np$ and $V(X) = np(1 - p)$

■ **Example 4.14 — Binomial Experiments.** Below are some examples of Binomial experiments.

- The number of people who tried the dalgona candy challenge following ‘Squid Game’
- The number number of randomly selected students who started playing Animal Crossing in 2020
- The number of games won out of 7 independent games with the same opponent

■

■ **Example 4.15** While studying by a window, you find yourself noticing many cars at a nearby intersection that fail to fully come to a stop at the stop sign before passing through the intersection. Based on your months of data, you reliably calculate the probability of drivers failing to do a complete stop to be 70%. Assuming the stopping behaviour of each car is independent of all others, find the probability that among 20 randomly observed cars that...

- a) Exactly 5 will come to a complete stop?

Let S be the number of cars that stop.

$$S \sim \text{Bin}(n = 20, p = 0.3).$$

$$\begin{aligned}
 P(S=5) &= \binom{20}{5} \cdot 0.3^5 \cdot (1-0.3)^{20-5} \\
 &= \frac{20!}{(20-5)! \cdot 5!} \cdot 0.3^5 \cdot 0.7^{15} \\
 &\approx 0.1789
 \end{aligned}$$

b) At least 3 will come to a complete stop?

$$\begin{aligned}
 P(S \geq 3) &= \sum_{s=3}^{20} \binom{20}{s} \cdot 0.3^s \cdot (1-0.3)^{20-s} && \text{(direct)} \\
 &= 1 - P(S < 3) && \text{(indirect)} \\
 &= 1 - P(S \leq 2) \\
 &= 1 - \left(\binom{20}{0} \cdot 0.7^{20} + \binom{20}{1} \cdot 0.3^1 \cdot 0.7^{19} + \binom{20}{2} \cdot 0.3^2 \cdot 0.7^{18} \right) \\
 &\approx 0.9645
 \end{aligned}$$

c) At most 3 will come to a complete stop?

$$\begin{aligned}
 P(S \leq 3) &= P(S=0) + P(S=1) + P(S=2) + P(S=3) \\
 &= \binom{20}{0} \cdot 0.7^{20} + \binom{20}{1} \cdot 0.3^1 \cdot 0.7^{19} + \binom{20}{2} \cdot 0.3^2 \cdot 0.7^{18} + \binom{20}{3} \cdot 0.3^3 \cdot 0.7^{17} \\
 &\approx 0.1071
 \end{aligned}$$

■ **Example 4.16** A local hospital has several backup generators to support critical technologies in the event of a power outage or failure. Each backup generator is identical in make, and operate independently of others. Suppose each backup generator has a 20% chance of failing when used. How many generators should be installed so that the system has at least a 99.5% probability of functioning in the event of a power outage?

Let n be the number of generators (a fixed quantity).

Let G be the number of generators that functions.

$G \sim \text{Bin}(n, p = 0.8)$.

We want to find n such that $P(G \geq 1) \geq 99.5\%$.

That is, by the indirect method, we need $P(G = 0) \leq 0.005$

$$\binom{n}{0} \cdot 0.2^n \leq 0.005$$

$$0.2^n \leq 0.005$$

$$n \geq 3.29$$

∴ At least 4 generators should be installed. ■

4.3.2 Poisson Distribution

Consider modeling the number of Shiba, D , spotted at a nearby park over any 1 day with a probability model. (How is this different from a Binomial model if it still models the number of ‘successes’?)



This is **not** a binomial distribution, as trials are *discrete*, while time is *continuous*!

Let's try to formulate this problem so it resembles a Binomial model: first we will arbitrarily divide the 1 day into n equally-sized time interval with the following properties for any one interval:

- $P(D = 1) = p$
- $P(D = 0) = 1 - p$
- $P(D > 1) = 0$ (i.e. the event of Shiba sighting is “rare”)

Let us also assume that each time interval behaves independently, and the average (mean) number of Shiba sightings per day is fixed and denoted by λ .

Based on the construction and assumptions, we have n independent trials with equal probabilities of “success” p . This can be modeled as a binomial distribution where $D \sim \text{Bin}(n, p)$, which has an expected values of $E(D) = np$.

Since the mean number of daily sightings is constant,

- $E(D) = np = \lambda$, and $p = \frac{\lambda}{n}$
- The number of time intervals is arbitrarily decided, neither n nor p are known.
- In order to ensure daily average $\lambda = np$ remains constant, as n increases, p must decrease so that np remains unchanged

We'll get more accurate probabilities of daily sightings in a day if we allow each time interval to shrink to 0 (not too different from using Riemann sums to approximate area under continuous curves!). Let's see how the binomial PMF behaves as $n \rightarrow \infty$ and $p \rightarrow 0$.

The resulting function will model the *probability of D number of Shiba sightings over a continuous period of 1 day*.

Let D be the number of Shiba sighted in “ n ” sub-intervals.

$$D \sim \text{Bin}(n, p = \frac{\lambda}{n}).$$

$$E(D) = np = \lambda.$$

$$\begin{aligned}
 P(D = d) &= \lim_{n \rightarrow \infty} \binom{n}{d} \cdot p^d \cdot (1 - p)^{n-d} \\
 &= \lim_{n \rightarrow \infty} \frac{n!}{(n-d)! \cdot d!} \cdot \left(\frac{\lambda}{n}\right)^d \cdot \left(1 - \frac{\lambda}{n}\right)^{n-d} \\
 &= \frac{\lambda^d}{d!} \lim_{n \rightarrow \infty} \frac{n(n-1)(n-2) \cdots (n-d+1)(\cancel{n-d}!)!}{(\cancel{n-d}!)!} \cdot \frac{1}{n^d} \left(1 - \frac{\lambda}{n}\right)^{n-d} \\
 &= \frac{\lambda^d}{d!} \lim_{n \rightarrow \infty} \frac{n}{n} \times \frac{n-1}{n} \times \frac{n-2}{n} \times \cdots \times \frac{n-d+1}{n} \cdot \left(1 - \frac{\lambda}{n}\right)^{n-d} \\
 &= \frac{\lambda^d}{d!} \lim_{n \rightarrow \infty} (1) \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{d-1}{n}\right) \left(1 - \frac{\lambda}{n}\right)^{n-d} \\
 &= \frac{\lambda^d}{d!} \lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^{n-d} \\
 &= \frac{\lambda^d}{d!} \lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-d} \xrightarrow{d \rightarrow 1} 1 \\
 &= \frac{\lambda^d}{d!} \cdot e^{-\lambda} \quad \text{as } \lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^n = e^{-\lambda} \\
 &= \frac{\lambda^d e^{-\lambda}}{d!}
 \end{aligned}$$

Definition 4.3.3 — Poisson Distribution. A discrete random variable X denoting the number of (sometimes rare) events of interest in an interval, with the mean number of occurrences per unit

interval denoted by λ , is *Poisson Distributed* if it has the probability mass function

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

X has expectation $E(X) = \lambda$ and variance $V(X) = \lambda$.

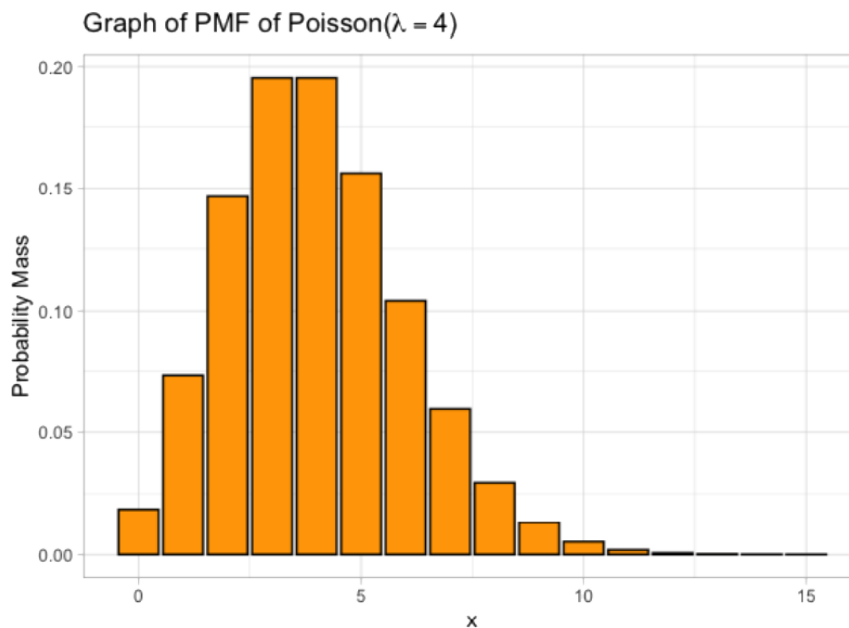
Note that λ is also called the *rare parameter* as it describes the average *rate of occurrence* of the event. λ should always be adjusted for the time interval being considered.

Poisson distribution is appropriate for discrete random variables that count the number of occurrences in a **continuous interval** where

- No more than one occurrence can occur simultaneously
- Non-overlapping intervals have occurrences that behave independently
- Expected number of occurrences in each fixed time interval is constant

As suggested in the construction of the Poisson distribution, this distribution can be used to approximate binomial probabilities where n is large and p is small. The approximation improves as n increases and p decreases. The advantage here is that Poisson distribution can be easier to compute since it doesn't involve the $\binom{n}{x}$ factor in its PMF.

What does a Poisson distribution look like? Generally unimodal and right-skewed, there is greater likelihood in observing values near but less than the mean $\lambda = 4$.



■ **Example 4.17** An area of a forest as on average 6 trees per 100 m².

a) Do the assumptions for a Poisson model seem appropriate for modeling tree distribution in a forest

- Two simulations trees (two trees occupying the same spot) is 0
- Could reasonably assume that the number of trees per area is independent

b) What is the probability of having at least 1 tree in a 100 m² area?

Let T be the number of trees per 100 m².

$T \sim \text{Pois}(\lambda = 6)$.

Using the indirect method, $P(T \geq 1) = 1 - P(T = 0)$

$$= 1 - \frac{6^0 e^{-6}}{0!}$$

$$\approx 0.9975$$

\therefore There is 99.75% change of having at least 1 tree in a 100 m² area. ■

■ **Example 4.18** The number of car repairs R that arrive at a mechanic can be well modeled by a Poisson random variable with a mean of $3t$, where t denotes the time in hours of operating hours. The average profit per repair is given by $Y = 85R^2 - 60t$. Assuming that vehicle arrival occurs independently,

- a) What is the probability that in a 10 hour workday, they will service between 28 to 31 cars?
We can reasonably model $R \sim \text{Pois}(\lambda = 3 \cdot 10 = 30)$.

$$\begin{aligned} \text{Then, } P(28 \leq R \leq 31) &= \sum_{r=28}^{31} P(R=r) \\ &= \frac{30^{28} e^{-30}}{28!} + \frac{30^{29} e^{-30}}{29!} + \frac{30^{30} e^{-30}}{30!} + \frac{30^{31} e^{-30}}{31!} \\ &\approx 0.2858 \end{aligned}$$

\therefore 28.58% of the time they will service between 28 to 31 cars.

- b) What is the corresponding profit earned for the cars serviced in a)?

$$Y = 85R^2 - 60 \cdot 10 = 85R^2 - 600 \text{ with } R \in [28, 31].$$

$$\text{That is, } y \in [85 \cdot 28^2 - 600, 85 \cdot 31^2 - 600] = [66040, 81085]$$

\therefore They are expected to earn from \$66,040 to \$81,085.

- c) What is the expected profit for a typical 8 hour workday?

$$R \sim \text{Pois}(\lambda = 3 \cdot 8 = 24)$$

$$E(Y) = E(85R^2 - 60 \cdot 8)$$

$$= E(85R^2 - 480)$$

$$= 85E(R^2) - 480$$

$$\text{Since } V(R) = E(R^2) - E(R)^2, E(R^2) = V(R) + E(R)^2 = \lambda + \lambda^2 = 24 + 24^2.$$

$$\text{Then, } E(Y) = 85(24 + 24^2) - 480$$

$$= \$50,520$$

\therefore The expected profit for a typical 8 hour workday is \$50,520.

- d) Suppose that you know that in an 8 hour workday, $E(R^4) = 416,472$. Determine the variance in profit. Can you determine what interval of profits they can expect to earn with at least 75% probability?

$$R \sim \text{Pois}(24).$$

$$E(Y) = 50,520$$

$$V(Y) = V(85R^2 - 480) = 85^2 \cdot V(R^2)$$

$$= 85^2 \cdot (E((R^2)^2) - E(R^2)^2)$$

$$= 85^2 \cdot (416,472 - (24 + 24^2)^2)$$

$$= \$^2 408,010,200$$

$$\sigma_Y = \sqrt{V(Y)} \approx \$20,199.26$$

By Chebyshev, $P(|Y - \mu_Y| < k\sigma_Y) \geq 1 - \frac{1}{k^2}$, and we want $1 - \frac{1}{k^2} = \frac{3}{4}$. That is, $k = 2$.

$$y \in (50520 - 2(20199.26), 50520 + 2(20199.26)) = (10121.48, 90918.52)$$

\therefore They are expected to earn within (10121.48, 90918.52) with at least 75% probability. ■

4.3.3 Geometric Distribution

The geometric distribution models the probability of observing some number of **consecutive failed trials before a “success”** is observed. The trials are independent and identical Bernoulli trials with fixed probability of success, p .

Definition 4.3.4 — Geometric Distribution. Let X be the random variable representing the *number of failures before the first success*. The probability, p , is fixed for each trial.

$$P(X = x) = (1 - p)^x \cdot p$$

A geometric distribution has expectation $E(X) = \frac{1-p}{p}$ and variance $V(X) = \frac{1-p}{p^2}$.

Note that $P(X \geq k) = q^k$, that is, if the number of trials is more than k , then the first k trials must have all been failures.

■ **Example 4.19** What is the probability you ask 3 people **before** you find someone born in December. You may assume that every month is equally likely to be a birth month.

$$\left(\frac{11}{12}\right)^3 \cdot \left(\frac{1}{12}\right)$$

How many people do you expect to have to ask before finding someone born in December? What is the expected total number of people surveyed?

Total people = number of failures + 1

$$E(X + 1) = E(X) + 1$$

$$\begin{aligned} &= \frac{1 - \frac{1}{12}}{\frac{1}{12}} + 1 \\ &= 12 \end{aligned}$$

Definition 4.3.5 — Memoryless Property. The geometric distribution has one additional property. It is a *memoryless distribution*. Suppose you have observed j consecutive failures. The probability that we will observe at least another k failures given that the first j must be failures is the probability of observing k failures. In other words:

$$P(X > j + k \mid X \geq j) = P(X > k)$$

■ **Example 4.20 — Memoryless Property.** The probability that a component will function for more than 5 years if it is already 2 years old is the same as the probability that it functions for more than 3 years (this is obviously unrealistic, but an example of memoryless property). ■

■ **Example 4.21** What is the probability that it takes at least 8 rolls of a fair die before you roll a 6, if you did not roll a 6 in the first 4 rolls?

Let D be the number of dice rolls before rolling a 6¹.

$$P(D \geq 8 \mid D \geq 4) = P(D \geq 8 - 4)$$

$$\begin{aligned} &= \left(\frac{5}{6}\right)^4 \\ &\approx 48.23\% \end{aligned}$$

4.3.4 Negative Binomial Distribution

Negative binomial random variables model the number of failed attempts before we observe a total of r successes.

¹If $D = 1$, then we roll a 6 on the second trial

The random variable can model either the number of failures or the total number of trials. Note that the two are linear shifts of each other:

$$n(\text{failure}) + ('r' \text{ successes}) = n(\text{trials})$$

Since the last success is guaranteed in the last trial, we can reduce the problem to finding the probability of $(r - 1)$ successes in the first $(x + r - 1)$ trials. This type of random variable X has a {term}negative binomial distribution.

Definition 4.3.6 — Negative Binomial Distribution. Let X be the number of failures before r^{th} success. Each trial is an independent and identical (i.i.d) Bernoulli trial with p probability of success.

$$P(X = x) = \binom{x + r - 1}{r - 1} \cdot p^r \cdot (1 - p)^x$$

A negative binomial distribution has expectation $E(X) = \frac{r(1 - p)}{p}$ and variance $V(X) = \frac{r(1 - p)}{p^2}$.

Note that if we model Y to be the number of trials to achieve r^{th} success, then we have a success in the Y^{th} trial, $Y - r$ failures and $r - 1$ successes among the first $Y - 1$ trials is

$$P(Y = y) = \binom{y - 1}{r - 1} \cdot p^r \cdot (1 - p)^{y - r}$$

■ **Example 4.22 — Negative Binomial Distribution.** Below are examples of Negative Binomial Distributions.

- The number of losing lotto numbers selected before you have 6 winning numbers
- The total numbers you have selected to achieve 6 winning numbers
- The number of tails before getting 4 heads in consecutive coin tosses
- The number of times you toss a coin to get 4 heads

■

■ **Example 4.23** Find the probability that you will have to survey $Y = D_1 + D_2 + D_3 = 10$ people NOT born in December before you have a total of 3 people born in December (i.e. that you survey a total of $10 + 3$ people to find a total of 3 December babies).

✓✓XXXXXXXXXX|✓

In order for 10 non-December individuals to be surveyed before this is fulfilled, this tells you that then 13th person is born in December, while 2 of the first 12 people are born in December.

$$\begin{aligned} P(\text{Survey 10 failures before 3 successes}) &= \binom{10 + 3 - 1}{3 - 1} \cdot \left(\frac{1}{2}\right)^3 \cdot \left(1 - \frac{1}{2}\right)^{10} \\ &= \binom{12}{2} \cdot \left(\frac{1}{2}\right)^3 \cdot \left(\frac{11}{12}\right)^{10} \\ &\approx 1.6\% \end{aligned}$$

■

4.3.5 Hypergeometric Distribution

A hypergeometric distribution can be used to calculate the probability of selecting k “successes” out of n selections **when the pool of selection N is small**. In other words,

- You have a population of size N made up of subpopulations: “successes” and “failures”.
- You are randomly selecting a group of n from this population (without replacement) and observing the number of “successes” in your group of n .

- Example: Tagging wildlife to be tracked over several years to observe population growth, territory expansion, etc.

If the pool of selection N is large relative to n , then the probability can be approximated using a binomial distribution.

- Sampling without replacement, the probability of subsequent selections is not greatly impacted
- We can model this with “ n ” trials, and $p = \frac{k}{N}$

■ **Example 4.24** Consider the following example, where N is big versus N is small.

- We have a lot of 5000 items, 10% of which are known to be defective. Find the probability of observing at least three defective items in a random selection of 10 items.

Let D be the number of defectives.

$$D \sim \text{Hypergeometric}(5000, n = 10, k = 500)$$

$$P(D \geq 3) = 1 - P(D \leq 2)$$

$$= 1 - \frac{\binom{4500}{10}}{\binom{5000}{10}} - \frac{\binom{500}{1} \cdot \binom{4500}{9}}{\binom{5000}{10}} - \frac{\binom{500}{2} \cdot \binom{4500}{8}}{\binom{5000}{10}}$$

$$= 7.00\%$$

- We have a lot of 50 items, 10% of which are known to be defective. Find the probability of observing at least three defective items in a random selection of 10 items.

Let D be the number of defectives.

$$D \sim \text{Hypergeometric}(50, n = 10, k = 5)$$

$$P(D \geq 3) = 1 - P(D \leq 2)$$

$$= 1 - \frac{\binom{45}{10}}{\binom{50}{10}} - \frac{\binom{5}{1} \cdot \binom{45}{9}}{\binom{50}{10}} - \frac{\binom{5}{2} \cdot \binom{45}{8}}{\binom{50}{10}}$$

$$\approx 4.83\%$$

- Approximate probability using binomial distribution for both, since $D \sim \text{Bin}(n = 10, p = 0.10)$, $P(D \geq 3) = 1 - P(D \leq 2) = 1 - \sum_{d=0}^2 \binom{10}{d} \cdot 0.10^d \cdot (1 - 0.10)^{10-d} \approx 7.02\%$

■

Definition 4.3.7 — Hypergeometric Distribution. Suppose you have a pool of N objects that can be partitioned into 2 (or more groups) by some characteristic. Suppose there are k objects of type A and $N - k$ objects of type B. In a random sample of size n (without replacement) from this pool of N objects, let X denote the random variable for the number of objects of type A that is selected.

$$P(X = x) = \frac{\binom{k}{x} \cdot \binom{N-k}{n-x}}{\binom{N}{n}}, \max(0, n - (N - k)) \leq x \leq \min(n, k)$$

X has an expected value and variance $E(X) = n \cdot \frac{k}{N}$ and $V(X) = n \cdot \frac{k}{N} \left(1 - \frac{k}{N}\right) \left(\frac{N-n}{N-1}\right)$

5. Continuous Random Variables

5.1 Probability Density Function

Recall that a random variable is **continuous** if it can take on values on an interval of real numbers. For example,

- The mass of a randomly selected salmon in kilograms
- The horizontal distance a figure skater travels on a specific type of jump in metres
- The random elapsed time between football plays in seconds

Definition 5.1.1 — Probability Density Function. The *probability density function* (PDF) of a continuous random variable X is a function $f(x)$ that has the following properties.

1. $f(x) \geq 0$ for all x in the support of X

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

3. $P(a \leq X \leq b) = \int_a^b f(x) dx$

Note that unlike in the discrete case, the density function $f(x) \neq P(X = x)$ for continuous RV X . Instead, $f(x)$ describes the probability *density* for various values of x and the **area under the probability density function** corresponds to the **probability over the interval**.

■ **Example 5.1** Verify that $f(y) = 3(y-1)^2$, $0 \leq y \leq 1$ ¹ is a valid probability density function.

1. $3 > 0$ and $(y-1)^2 \geq 0$. Thus, $f(y) = 3(y-1)^2 \geq 0$.

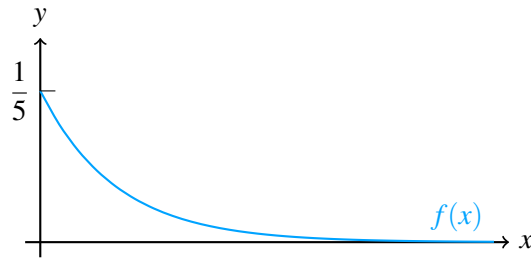
¹This is called the *support* of the random variable Y – the set of y that has **non-zero** probabilities. In other words,

$$f(y) = \begin{cases} 3(y-1)^2 & 0 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned}
2. \int_{-\infty}^{\infty} f(y) dy &= \int_{-\infty}^0 0 dy + \int_0^1 3(y-1)^2 dy + \int_1^{\infty} 0 dy \\
&= 0 + (y-1) \Big|_0^1 \\
&= (1-1)^3 - (0-1)^3 \\
&= 1
\end{aligned}$$

■ **Example 5.2 — MMSA Ex. 4.30.** The response time X in seconds at an online computer terminal, which is the elapsed time between the end of a user's inquiry and the beginning of the system's response, had a distribution given by the following probability density function

$$f(x) = \begin{cases} \frac{1}{5}e^{-x/5} & x > 0 \\ 0 & \text{otherwise} \end{cases}$$



- a) What is the probability the response time will be more than 10 seconds long if there was no response in the first 8 seconds?

$$\begin{aligned}
P(X > 10 | x \geq 8) &= P(X > 2) \\
&= \int_2^{\infty} \frac{1}{5}e^{-x/5} \\
&= e^{-2/5} \\
&\approx 67.03\%
\end{aligned}$$

Definition 5.1.2 — Expected Value. The *expected value* (sometimes called the *mean*) of a continuous random variable X with probability density function $f(x)$ is given by

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

For any real-valued function $g(X)$ of X ,

$$E(g(X)) = \int_{-\infty}^{\infty} g(x)f(x) dx$$

Definition 5.1.3 — Variance. The *variance* of a continuous random variable X with probability density function $f(x)$ is given by

$$V(X) = E((X - \mu)^2) = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$$

As was shown in the discrete case, we have an alternative, often shorter way to compute

variance:

$$V(X) = E(X^2) - E(X)^2$$

- The variance is the average *squared* deviation of X from its mean.
- The *standard deviation* of a random variable is the square root of the variance ($SD = \sqrt{V(X)}$).

■ **Example 5.3** A gas station operates a pump that can pump up to 20,000 gallons of gas a month. The total amount of gas that is pumped at a station in a month can be modelled with a random variable G (measured in 10,000 gallons) with PDF

$$f(g) = \begin{cases} g & 0 < g < 1 \\ 2 - g & 1 \leq g < 2 \\ 0 & \text{otherwise} \end{cases}$$

- a) Find the expected number of gallons of gas pumped per month.

$$\begin{aligned} E(F) &= \int_{-\infty}^{\infty} g \cdot f(g) dg \\ &= \int_0^1 g \cdot g dg + \int_1^2 g(2 - g) dg \\ &= \frac{1}{3}g^3 \Big|_0^1 + \left(g^2 - \frac{1}{3}g^3 \right) \Big|_1^2 \\ &= \frac{1}{3} + \left(\left(4 - \frac{8}{3} \right) - \left(1 - \frac{1}{3} \right) \right) \\ &= 3 - \frac{6}{3} \\ &= 1 \end{aligned}$$

- b) Find the variance in the number of gallons of gas pumped per month.

$$\begin{aligned} V(G) &= E((G - \mu)^2) \\ &= E(G^2) - E(G)^2 \\ \text{Note that } E(G^2) &= \int_{-\infty}^{\infty} g^2 \cdot f(g) dg \\ &= \int_0^1 g^2 \cdot g dg + \int_1^2 g^2(2 - g) dg \end{aligned}$$

■

5.2 Common Continuous Distributions

5.2.1 Uniform Distribution

The previous example uses a common continuous distribution called the ‘*Uniform Distribution*’. A distribution with constant density tells us that any intervals of the same width in the support have equal probability of occurrence. Such distributions are called **uniform**.

■ **Definition 5.2.1 — Uniform Distribution.** A continuous random variable X follows a *uniform distribution* on the interval $a \leq X \leq b$ if it has probability density function

$$f(x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$

A uniform distribution has mean and variance defined by its endpoints

$$E(X) = \frac{b+a}{2} \quad V(X) = \frac{(b-a)^2}{12}$$

5.2.2 Cumulative Distribution Function

Definition 5.2.2 — Cumulative Distribution Function. The *cumulative distribution function* (CDF) of a continuous random variable X is the function $F(x)$ that returns the cumulative probability of a random variable X for any value $x \in \mathbb{R}$

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

The function $F(x)$ can be explicitly found by finding an expression for $\int_{-\infty}^x f(w)$

The derivative of the CDF $F'(x)$ is the PDF of X : $F'(x) = f(x)$

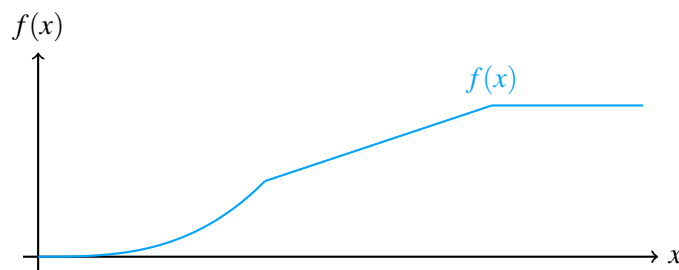
Proposition 5.2.1 — Properties of CDF. Below are some properties of CDF.

- $P(X = c) = \int_c^c f(x) dx = 0 \implies P(X \leq c) = P(X < c)$
- $P(a \leq X \leq b) = \int_a^b f(x) dx = F(b) - F(a)$
- $\lim_{x \rightarrow \infty} F(x) = 1$ and $\lim_{x \rightarrow -\infty} F(x) = 0$

■ **Example 5.4** The distribution function of a random variable X is as follows:

$$F(x) = \begin{cases} 0 & x < 0 \\ \frac{x^3}{2} & 0 \leq x < 1 \\ \frac{x}{2} & 1 \leq x \leq 2 \\ 1 & x > 2 \end{cases}$$

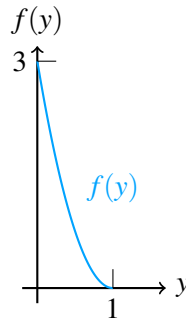
Find $P(0.5 \leq X \leq 1.5)$.



$$\begin{aligned} P(0.5 \leq X \leq 1.5) &= P(X \leq 1.5) - P(X < 0.5) \\ &= F(1.5) - F(0.5) \\ &= \frac{1.5}{2} - \frac{0.5^3}{2} \\ &\approx 68.75\% \end{aligned}$$

■ **Example 5.5** For a random variable Y with probability density function $f(y) = 3(y-1)^2$, $0 \leq y \leq 1$

a) Sketch the PDF.



b) Find the cumulative distribution function of Y .

CDF of Y : $P(Y \leq y)$

$$\text{For } y \in [0, 1), P(Y \leq y) = \int_0^y 3(u-1)^2 du$$

$$= (u-1)^3 \Big|_0^y$$

$$= (y-1)^3 - (0-1)^3$$

$$= 1 + (y-1)^3$$

$$\text{Thus, we have } F_Y(y) = \begin{cases} 0 & y < 0 \\ 1 + (y-1)^3 & 0 \leq y < 1 \\ 1 & y \geq 1 \end{cases}$$

c) Find $P(0.1 \leq Y \leq 0.4)$

• PDF: Area under PDF from 0.1 to 0.4, $\int_{0.1}^{0.4} f(y) dy$

• CDF: $P(0.1 \leq Y \leq 0.4) = P(Y \leq 0.4) - P(Y < 0.1)$
 $= F(0.4) - F(0.1)$

d) Find the probability that $Y = 0.6$.

0

■

5.2.3 Exponential Distribution

Definition 5.2.3 — Exponential Distribution. A continuous random variable X is *exponentially distributed* with mean parameter $\theta > 0$ (OR rate of $\lambda > 0$) if it has the probability density function

$$f(x) = \begin{cases} \frac{1}{\theta} e^{-\frac{x}{\theta}} \text{ OR } \lambda e^{-\lambda x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

X has mean and variance

$$E(X) = \theta \text{ OR } \frac{1}{\lambda} \quad B(X) = \theta^2 \text{ OR } \frac{1}{\lambda^2}$$

We say $X \sim \text{Exp}(\theta)$ or $X \sim \text{Exp}(\lambda)$ to define the distribution.

Definition 5.2.4 Given an exponential random variable X with parameter θ , derive the cumulative distribution function (CDF) of X . Are exponentially distributed RV memoryless?

1. CDF: $P(X \leq x), X \sim \text{Exp}(\theta)$

$$\begin{aligned}
 F_X(x) &= \int_0^x f(u) du \\
 &= \int_0^x \frac{1}{\theta} e^{-\frac{u}{\theta}} du \\
 &= -e^{-\frac{u}{\theta}} \Big|_0^x \\
 &= 1 - e^{-\frac{x}{\theta}}
 \end{aligned}$$

$$\text{That is, } F_X(x) = \begin{cases} 0 & x \leq 0 \\ 1 - e^{-x/\theta} & x > 0 \end{cases}$$

2.

■ **Example 5.6** Previously, we modeled the number of shiba arrivals at the dog park over a day by a Poisson distribution. Suppose on average there are 5 Shiba arrivals per day.

- a) A Shiba just arrived at the dog park. What is the probability that another one will arrive in the next day (full 24 hours)?

Let T be the number of days until the next Shiba arrival.

$$R(T \leq 1) = F_T(1)$$

$$\begin{aligned}
 &= \int_0^1 \frac{1}{\theta} e^{-\frac{t}{\theta}} dt \\
 &= 1 - e^{-\frac{1}{\theta}} \\
 &= 1 - e^{-5} \\
 &\approx 0.9933
 \end{aligned}$$

- b) It's been 3 days since you've last seen your a Shiba at the dog park. What's the probability you'll have to wait at least five days before another Shiba arrives at the park?

$$\begin{aligned}
 P(T \geq 5 | T \geq 3) &= \frac{P(T \geq 5 \wedge T \geq 3)}{P(T \geq 3)} \\
 &= \frac{P(T \geq 5)}{P(T \geq 3)}
 \end{aligned}$$

$$\begin{aligned}
 \text{OR } P(T \geq 2) &= 1 - P(T < 2) \\
 &= 1 - F_T(2) \\
 &= 1 - (1 - e^{-10}) \\
 &\approx 0.00454\%
 \end{aligned}$$

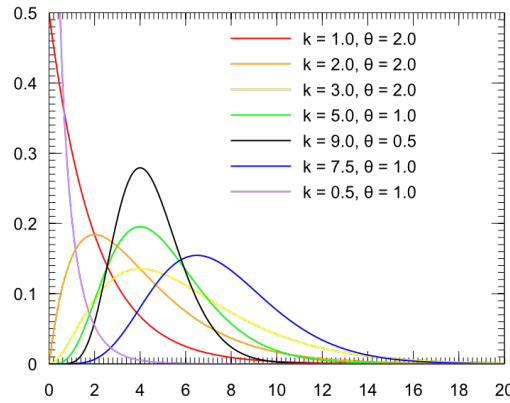
■

5.2.4 Gamma Distribution

We have learned that the exponential distribution models the random interval until the **next** Poisson arrival (or in other words, the random interval between two consecutive arrivals). Can it also be used to model the total interval between any two Poisson arrivals?

What would the distribution of total time between the second and sixth shiba arrival look like? Will it also be exponential? What do you think?

The exponential distribution that you learned, which can sometimes be used to model wait time, is quite limiting in its shape. The related *gamma distribution* that you're about to learn about is **defined by two parameters** which offers more options and flexibility in modeling quantities that tend to be right-skewed.



Definition 5.2.5 — Gamma Distribution. A *gamma distribution* with **shape parameter** $\alpha > 0$ and **scale parameter** $\beta > 0$ for a continuous random variable X has probability density function

$$f(x) = \begin{cases} \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-\frac{x}{\beta}} & x > 0 \\ 0 & \text{otherwise} \end{cases}$$

X has mean and variance

$$E(X) = \alpha\beta \quad V(X) = \alpha\beta^2$$

Note:

1. $\Gamma(\alpha)$ is the *gamma function* defined by

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

2. For $\alpha > 1$, $\Gamma(\alpha) = (\alpha - 1) \cdot \Gamma(\alpha - 1)$
3. For $\alpha \in \mathbb{Z}_+$, $\Gamma(\alpha) = (\alpha - 1)!$
4. $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$
5. When $\alpha = 1$, $\text{Gamma}(1, \beta) = \text{Exp}(\beta)$. That is, the exponential distribution is a special case of the gamma distribution².
6. The sum of n **independent** $\text{Exp}(\theta)$ random variables $T = X_1 + X_2 + \cdots + X_n$, where $X_i \sim \text{Exp}(\theta)$, results in a gamma random variable with $T \sim \Gamma(n, \theta)$.
7. The CDF of $\Gamma(\alpha, \beta)$ has no closed form when $\alpha, \beta \notin \mathbb{Z}_+$, making it difficult to compute probabilities without R.

Let's derive the expected value of Gamma.

Let X be a random variable, $X \sim \Gamma(\alpha, \beta)$.

$$\text{Then, } f(x) = \begin{cases} \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-x/\beta} & x > 0 \\ 0 & \text{otherwise} \end{cases}$$

² β is the scale parameter, not rate – recall scale θ and rate λ from exponential distribution

$$\begin{aligned}
E(X) &= \int_0^{\infty} x \cdot f(x) dx \\
&= \int_0^{\infty} \frac{1}{\Gamma(\alpha)\beta^\alpha} x^\alpha e^{-x/\beta} dx \\
&= \frac{1}{\Gamma(\alpha)\beta^\alpha} \int_0^{\infty} x^\alpha e^{-x/\beta} dx \quad \text{Let } u = \frac{x}{\beta}, dx = \beta du \\
&= \frac{1}{\Gamma(\alpha)\beta^\alpha} \int_0^{\infty} (u\beta)^\alpha e^{-u} \beta du \\
&= \frac{1}{\Gamma(\alpha)\beta^\alpha} \int_0^{\infty} \beta^{\alpha+1} u^\alpha e^{-u} du \\
&= \frac{\beta}{\Gamma(\alpha)} \int_0^{\infty} u^\alpha e^{-u} du
\end{aligned}$$

R Code Computation – Gamma

Since the CDF of a gamma random variable has no closed form solution, we can use R statistical software to compute probabilities when $\alpha, \beta \notin \mathbb{Z}$.

Given $X \sim \Gamma(\alpha, \beta) \dots$

- $P(X \leq k)$ is computed with `pgamma(k, shape= α , scale= β)`. You MUST indicate whether you are using scale or rate parameters when using R since it allows for either: $\Gamma(\text{shape}, \text{scale})$ or $\Gamma(\text{shape}, \text{rate} = \frac{1}{\text{scale}} = \frac{1}{\beta})$
- To simulate n data values from X , we use `rgamma(n, shape= α , scale= β)`

■ **Example 5.7** The number of service requests fulfilled by a technician at a support centre is Poisson distributed with an average of 10 per workday (8 hours). How would you model the random time it takes to complete two unrelated and independent service requests by this technician? Find the probability that two service requests takes more than 3 hours to fulfill.

Let R be the number of service requests per workday (8 hours).

$R \sim \text{Pois}(\lambda = 10)$.

Let T_i be the time in hours to fulfill the i^{th} service request.

$T_i \sim \text{Exp}\left(\theta = \frac{8}{10}\right)$.

Want to model $Y = T_1 + T_2 \sim \Gamma(\alpha = 2, \beta = 0.8)$.

$$f(y) = \begin{cases} \frac{1}{\Gamma(2)0.8^2} y^{2-1} e^{-y/0.8} & y > 0 \\ 0 & \text{otherwise} \end{cases} = \begin{cases} \frac{1}{0.64} y e^{-y/0.8} & y > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned}
P(Y > 3) &= \int_3^{\infty} \frac{1}{0.64} y e^{-y/0.8} dy \\
&= \frac{1}{0.64} \int_3^{\infty} y e^{-y/0.8} dy \\
&= \frac{1}{0.64} \left[-0.8 y e^{-y/0.8} - 0.64 e^{-y/0.8} \right] \Big|_3^{\infty} \\
&= \frac{1}{0.64} \left[0 - 0 + 0.8 y e^{-3/0.8} + 0.64 e^{-3/0.8} \right] \Big|_3^{\infty} \\
&\approx 11.17\%
\end{aligned}$$

\therefore There is 11.17% chance that 2 service requests take more than 3 hours.

Using R to compute $P(T > 3) = 1 - P(T \leq 3)$, we use `1-pgamma(3, shape=2, scale=0.8)`.

■

5.3 Percentile

Definition 5.3.1 — Percentile. The k^{th} percentile is the value of a random variable for which $k\%$ of values are less than or equal to it. For a random variable X with PDF $f(x)$, the k^{th} percentile $x_{k/100}$ is the value of X at which $k\%$ of outcomes lie below it. That is,

$$P(X \leq x_{k/100}) = \frac{k}{100} \quad \text{OR} \quad x_{k/100} = F^{-1}\left(\frac{k}{100}\right)$$

Special Percentiles

- The *median* is also called the 50th percentile
- The *first quartile* and *third quartile* are the 25th and 75th percentiles, respectively.
- Measure of spread called Inter Quartile Range (IQR) = $Q_3 - Q_1$.

■ **Example 5.8** Find the median of a random variable M with PDF $f(m) = 3m^2$, $0 \leq m < 1$.

Let $m_{0.5}$ be the median.

$$\int_0^{m_{0.5}} 3m^2 dm = 0.5$$

$$m^3 \Big|_0^{m_{0.5}} = 0.5$$

$$m_{0.5}^3 = 0.5$$

$$m_{0.5} = \sqrt[3]{0.5}$$

$$\approx 0.7937$$

Thus, 50% of the time, M will take on a value larger than 0.7937. ■

6. Bivariate Distributions

6.1 Multivariate and Joint Distributions

Thus far in the course, we have been studying **univariate distributions** – i.e. how we can model the behaviour or outcomes of a single variable. This is not always realistic or practical as we may sometimes be interested in how two (or usually more) variables occur together.

Consider the studio of colourblindness among males and females:

	Men (M)	Women (M')	Total
Colourblind (C)	0.04	0.002	0.042
Not Colourblind (C')	0.47	0.488	0.958
Total	0.51	0.49	1.00

The table is an example of a *joint probability distribution*.

Definition 6.1.1 — Joint Probability Distributions. Given two **discrete** random variables X and Y , the joint PMF is defined as

$$f(x,y) = P(X = x, Y = y) = P(X = x \cap Y = y)$$

where x, y are used to denote particular values in the support of X and Y .

Given two **continuous** random variables X and Y , we have the joint probability function $f(x,y)$ to describe the probability density over the region of (X,Y) outcomes. Recall that *probability density \neq probability mass*!

Any function satisfying the below two properties is a valid probability mass / density function.

Discrete

1. $0 \leq f(x, y) \leq 1$
2. $\sum_{x \in X} \sum_{y \in Y} P(X = x, Y = y) = 1$

Continuous

1. $f(x, y) \geq 0$ ^a
2. $\int_{x \in X} \int_{y \in Y} f(x, y) dy dx = 1$

^aThe density can exceed 1, as long as the area ≤ 1

Definition 6.1.2 — Marginal Distributions. The *marginal probability mass/density functions* can be extracted from the joint distribution, which returns the probability mass/density of one variable only

- $f_X(x) = P(X = x) = \sum_{y \in Y} P(X = x, Y = y)$
- $f_X(x) = P(X = x) = \int_{y \in Y} P(X = x, Y = y) dy$
- $f_Y(y) = P(Y = y) = \sum_{x \in X} P(X = x, Y = y)$
- $f_Y(y) = P(Y = y) = \int_{x \in X} P(X = x, Y = y) dx$

If $\forall (x, y) \in (X, Y)$, $f(x, y) = f_X(x) \cdot f_Y(y)$, the X and Y are *independent*. otherwise, X and Y are said to be *dependent*.



Bibliography

Textbooks

- [DBC21] Jay L. Devore, Kenneth N. Berk, and Matthew A. Carlton. *Modern Mathematical Statistics with Applications*. Springer, Apr. 2021.
- [WD21] Amy S. Wagaman and Robert P. Dobrow. *Probability with Applications and R*. 2nd edition. WILEY, July 2021.

Index

A

Associative Law	13
Axiom 1	15
Axiom 2	15
Axiom 3	15

B

Bayes' Rule	29
Bernoulli Trials	39
Binomial Distribution	40

C

Chebyshev's Inequality	38
Combinations	22
Commutative Law	13
Complement Probability	15
Conditional Probability	26
Continuous Random Variable	31
Cumulative Distribution Function	36, 52

D

DeMorgan's Laws	14
Discrete Random Variable	31
Disjoint	13
Distributive Law	13

E

Event	12
Complement Event	12
Complex Event	12
Simple Event	12
Expected Value	33
Expected Value (Continuous RV)	50
Exponential Distribution	53

G

Gamma Distribution	55
Gamma Function	55

I

Independent Events	13, 27
Intersection	12

J

Joint Probability Distributions	59
---------------------------------------	----

L

Law of Total Probability	28
--------------------------------	----

M

Marginal Distributions	60
------------------------------	----

Markov's Inequality	38
Memoryless Property	45
Mutually Exclusive	13, 27

P

Percentile	57
Permutations	20
Poisson Distribution	42
Probability	14
Probability Function	15
Probability Mass Function	31

R

Random Experiment	11
Random Variable	31

S

Sample Space	11
Standard Derivation	34
Support (of RV)	49

U

Uniform Distribution	51
Union	12

V

Variance	34
Variance (Continuous RV)	50