

# **Analytical Epidemiology**

**Statistical and Causal Inference for Public Health**

Eben Kenah

January 7, 2025

# Table of contents

<b>Preface</b>	<b>4</b>
Who this book is for . . . . .	5
How to use this book . . . . .	5
Acknowledgements . . . . .	5
 <b>I   A. One-Sample Inference for Risks and Rates</b>	 <b>7</b>
 <b>1   Probability, Random Variables, and Disease Occurrence</b>	 <b>8</b>
1.1   Sets, experiments, and events . . . . .	8
1.1.1   Experiments and events . . . . .	9
1.1.2   Set operations and logic . . . . .	10
1.1.3   Venn diagrams . . . . .	11
1.1.4   Sequences of events* . . . . .	13
1.1.5   Algebra of sets* . . . . .	14
1.2   Probability . . . . .	14
1.2.1   Probability calculations . . . . .	15
1.3   Random variables . . . . .	16
1.3.1   Indicator variables . . . . .	17
1.4   R . . . . .	17
1.4.1   Probability distributions . . . . .	17
1.4.2   Mean . . . . .	18
1.5   R . . . . .	19
1.5.1   Variance . . . . .	19
1.5.2   Bernoulli distribution . . . . .	19
1.6   Joint and marginal distributions . . . . .	20
1.7   R . . . . .	21
1.7.1   Linear combinations* . . . . .	21
1.7.2   Variance and covariance* . . . . .	22
1.8   Probability and disease occurrence . . . . .	23
1.8.1   Prevalence . . . . .	24
1.9   R . . . . .	25
1.9.1   Risk (cumulative incidence) and the survival function . . . . .	25
1.10   R . . . . .	25
1.10.1   Prevalence and the duration of disease . . . . .	26

1.10.2 Descriptive and analytic epidemiology . . . . .	27
<b>2 Conditional Probability and Diagnostic Tests</b>	<b>35</b>
<b>3 Maximum Likelihood Estimation</b>	<b>36</b>
<b>4 Bayesian Estimation</b>	<b>37</b>
<b>5 Longitudinal Data and Rates</b>	<b>38</b>
<b>6 Survival Analysis</b>	<b>39</b>
 <b>II B. Two-Sample Inference and Study Design</b>	 <b>40</b>
 <b>III C. Principles of Causal Inference</b>	 <b>41</b>
 <b>IV D. Epidemiologic and Statistical Methods for Causal Inference</b>	 <b>42</b>
<b>References</b>	<b>43</b>

# Preface

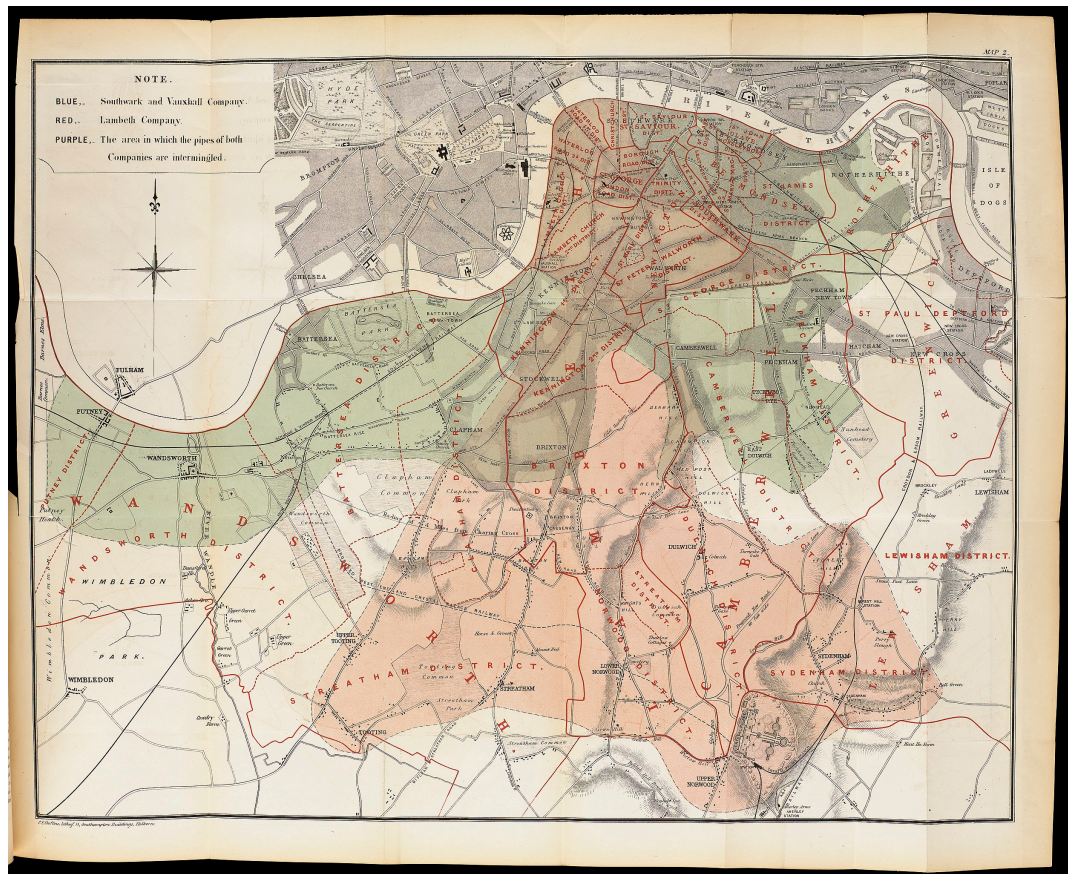


Figure 1: Areas of London supplied by the Southwark & Vauxhall (blue, now green) and Lambeth (red) water companies during the 1849-1854 cholera epidemic in London (Snow 1855). Source: Wellcome Collection via [Wikimedia Commons](#).

One day at lunch at the Harvard School of Public Health, I overheard Professor Murray Mittleman say: “I love epidemiology. It all fits together like a diamond.” As a second-year doctoral student in epidemiology, I was surprised to hear the subject described with such unstrained enthusiasm. It has taken years of study and experience for me to understand what he meant. On the way, I too have fallen in love.

## Who this book is for

This book is intended primarily for two audiences:

- Epidemiologists are often protected from the mathematical foundations of their field. The long-term price of this is “dogmatism, that is, a tendency to rigidly protect a partially understood theoretical heritage” (Morabia 2004). The mathematics needed for a deeper understanding of epidemiologic methods is within reach of anyone who has come far enough to need it. Whether you master this material or just learn to approach it with more patience than fear, you will be doing a service to epidemiology and to public health.
- Biostatisticians are familiar with probability and statistical inference, but applying statistics to solve scientific problems in public health requires skills different from those needed to prove mathematically that a method works under given assumptions. Epidemiology is a living example of the interplay between theory and applications in statistics, and epidemiologists have shown integrity, courage, and ingenuity in confronting causal questions with statistical tools.

Beyond these audiences, I hope to explain the logic of epidemiology to any interested reader. It is possible that epidemiologic research has already helped save your life.

## How to use this book

Difficult chapters, sections, subsections, and exercises are marked with an asterisk (\*). These can be skipped without harming the logical flow of the book, but none of them is beyond the reach of a determined reader. The starring is recursive: Starred sections can be skipped within a starred chapter, starred subsections can be skipped within a starred section, and so on. Footnotes offer context or hint at more advanced material. All of them can be ignored if they do not seem useful or interesting.

This is a work in progress. You may find that some parts are unfinished or just bad. Please report errors (including typos) or submit suggestions (especially good examples) at:

<https://github.com/ekenah/analyticalepi/issues>.

## Acknowledgements

This book is written in [LaTeX](#) and [Quarto](#) with calculations and figures generated in [R](#), [Python](#), and [Inkscape](#). These are free and open-source thanks to the work of many contributors.

Devesh Kapur, Paul Farmer, and James H. Maguire guided me to a career in public health when I was an undergraduate. James Robins, Miguel A. Hernán, Marc Lipsitch, and Stephen P. Luby helped me become an epidemiologist, biostatistician, and epidemic modeler in graduate

school. My career began under the mentorship of Ira M. Longini, Jr., and M. Elizabeth Halloran as a postdoctoral fellow at the University of Washington and an assistant professor at the University of Florida. My colleagues Yang Yang, Grzegorz Rempala, Forrest Crawford, and Patrick Schnell have all provided useful comments. For their patience with early versions of this material, I am grateful to the students of STA 6177/PHC 6937 (Applied Survival Analysis) at the University of Florida from 2013 to 2016 and PUBHEPI 8430 (Epidemiology 4) at The Ohio State University from 2019 to the present.

My parents, Chris and Kate Kenah, courageously allowed me to travel to places they had never been to and do things I had been told to avoid. These experiences in the United States, India, South Africa, and especially Bangladesh opened my eyes to the terrible importance of clear thinking in public health. My wife, Asma Aktar, and our sons Rafi, Rayhan, and Rabi remind me every day how important it is to destroy everything that stifles humanity. To that end, I hope this book is useful.

Any mistakes are my own, and God knows best     ). (

## **Part I**

### **A. One-Sample Inference for Risks and Rates**

# 1 Probability, Random Variables, and Disease Occurrence

One sees, from this essay, that probability theory is basically common sense reduced to calculation; it makes us appreciate with exactitude that which fair minds sense with a sort of instinct, often without being able to account for it. (Laplace 1820)<sup>1</sup>

To begin at the beginning, we will start with probability. Morabia (2004) accurately observed that “Epidemiology came late in human history because it had to wait for the emergence of probability.” This is probably the most difficult chapter of the entire book, but it will make all subsequent chapters easier. You can use it as a reference and come back to the difficult parts when you need them. Learning to think clearly about probability will give you a compass to find your way through epidemiologic methods.

## 1.1 Sets, experiments, and events

To speak clearly about probabilities, we need some basic notation for sets. If  $A$  is a set that contains an **element**  $a$ , we write

$$a \in A. \tag{1.1}$$

If  $A$  and  $B$  are sets such that every element of  $A$  is also an element of  $B$ , we write

$$A \subseteq B. \tag{1.2}$$

to indicate that  $A$  is a **subset** of  $B$ . Sets  $A$  and  $B$  are equal if and only if  $A \subseteq B$  and  $B \subseteq A$ , which means they contain exactly the same elements. The *empty set* with no elements is denoted  $\emptyset$ . For any set  $A$ , it is true that  $A \subseteq A$  and  $\emptyset \subseteq A$ .

---

<sup>1</sup>[Pierre-Simone, marquis de Laplace](#) (1749-1827) is often called the Newton of France. He proved that the solar system is stable, developed theories of ocean tides and gravitational potential, proved one of the first general versions of the central limit theorem, and pioneered the Bayesian interpretation of probability. For just six weeks in 1799, he was Minister of the Interior for France under Napoleon. His is one of the 72 names on the Eiffel Tower.



We use  $\mathbb{R}$  to denote the real numbers. Intervals are subsets of  $\mathbb{R}$  that take one of the following forms:

$$\begin{aligned}[a, b] &= \{x \in \mathbb{R} : a \leq x \leq b\}, \\[a, b) &= \{x \in \mathbb{R} : a \leq x < b\}, \\(a, b] &= \{x \in \mathbb{R} : a < x \leq b\}, \\(a, b) &= \{x \in \mathbb{R} : a < x < b\}.\end{aligned}$$

An endpoint with a square bracket is included in the interval; an endpoint with a round bracket is not. We can have  $a = -\infty$  or  $b = \infty$  as long as we use a round bracket for the corresponding endpoint. For example, it is true that  $\mathbb{R} = (-\infty, \infty)$ . However,  $\mathbb{R} \neq [-\infty, \infty]$  because  $\pm\infty$  are not real numbers.

### 1.1.1 Experiments and events

In probability, an **experiment** is any process that will produce one outcome out of a set of possible outcomes. The set of possible outcomes is called the **sample space** and is traditionally denoted  $\Omega$ . An experiment produces a single outcome  $\omega \in \Omega$ . For example, the sample space for a single coin flip is

$$\Omega = \{H, T\}, \tag{1.3}$$

where  $\omega = H$  if we get heads and  $\omega = T$  if we get tails.

The outcomes in the sample space must determine everything about the random outcome of the experiment. If we flip a coin twice, the sample space cannot be  $\{H, T\}$  because each  $\omega \in \Omega$  must specify the outcome of both coin flips. Instead,

$$\Omega = \{HH, HT, TH, TT\} \tag{1.4}$$

where  $\omega = XY$  if we get  $X$  on the first flip and  $Y$  on the second. This helps us see, for example, that there are two ways to get one  $H$  and one  $T$  in two coin flips.

The purpose of probability is to summarize uncertainty about the outcomes of experiments. However, the outcomes themselves do not have probabilities. Probabilities are assigned to **events**, which are subsets of the sample space  $\Omega$ . If  $A$  is an event, then  $A$  occurs if and only if the outcome  $\omega$  produced by our experiment is an element of  $A$  (i.e., if and only if  $\omega \in A$ ). If we flip a coin twice, the event that we get two heads is  $\{HH\}$ , the event that we get one head is  $\{HT, TH\}$ , and the event that we get zero heads is  $\{TT\}$ . By definition, the event  $\Omega$  always occurs and the event  $\emptyset$  never occurs.

In experiments with a finite or countably infinite sample space,<sup>2</sup> the distinction between the

---

<sup>2</sup>The natural numbers  $\mathbb{N} = \{0, 1, 2, \dots\}$  are *countably infinite*, as are the integers  $\mathbb{Z}$  and the rational numbers  $\mathbb{Q}$ . The real numbers  $\mathbb{R}$  are *uncountably infinite*, as are the real numbers in any nonempty interval  $(a, b)$  and the irrational numbers. Uncountably infinite sets are infinitely larger than countably infinite sets. This distinction was discovered in the 1870s by the German mathematician [Georg Cantor](#) (1845–1918). It was considered shocking, but it has become a cornerstone of modern mathematics.

outcome  $\omega$  and the event  $\{\omega\}$  can be safely ignored. In more complex experiments (e.g., taking a random sample from a standard normal distribution), this distinction is important.<sup>3</sup> In all cases, experiments have outcomes and events have probabilities.

In epidemiology, it is often useful to think of the sample space  $\Omega$  as being a population and each  $\omega \in \Omega$  as an individual in this population. In this context, our experiment is to sample a person from  $\Omega$  and ask them questions, take measurements, or follow them over time to ascertain disease occurrence. Events would be subpopulations of  $\Omega$ , such as  $\{\omega \in \Omega : \omega \text{ lives in Ohio}\}$ . This event occurs if the sampled individual  $\omega$  lives in Ohio, and it does not occur if they live somewhere else.

### 1.1.2 Set operations and logic

There are three basic set operations that take one or more sets and define another set: complement, intersection, and union. Each operation has a simple interpretation in terms of logic.

- The **complement** of a set  $A$  is

$$A^c = \{\omega \in \Omega : \omega \notin A\}, \quad (1.5)$$

which can be interpreted logically as **not**  $A$ . If  $A$  is an event, then the event  $A^c$  occurs if  $\omega \notin A$ . For the same reason that “not not  $A$ ” means “ $A$ ”, we have  $(A^c)^c = A$ .

- The **intersection** of two sets  $A$  and  $B$  is

$$A \cap B = \{\omega \in \Omega : \omega \in A \text{ and } \omega \in B\}, \quad (1.6)$$

which can be interpreted logically as  $A$  **and**  $B$ . If  $A$  and  $B$  are events, then the event  $A \cap B$  occurs if  $\omega \in A$  and  $\omega \in B$ .

- The **union** of two sets  $A$  and  $B$  is

$$A \cup B = \{\omega \in \Omega : \omega \in A \text{ or } \omega \in B\}, \quad (1.7)$$

which can be interpreted logically as  $A$  **or**  $B$  as long as we use an *inclusive* “or” (i.e., and/or). If  $A$  and  $B$  are events, then the event  $A \cup B$  occurs if  $\omega \in A$  or  $\omega \in B$ .

---

<sup>3</sup>In experiments with uncountably infinite sample spaces, the probability of an event  $A$  cannot always be calculated by adding up the probabilities of  $\{\omega\}$  for all  $\omega \in A$ . For example: If we choose a number at uniformly at random in  $[0, 1]$ , the probability of getting any particular number  $\omega$  is zero. The sum of the probabilities of all  $\{\omega\} \subseteq A$  is zero (if  $A$  is countable) or undefined (if  $A$  is uncountable). By maintaining a distinction between outcomes and events and by limiting probability calculations to countable (i.e., finite or countably infinite) sums, we end up with something coherent and useful.

If  $A \subseteq B$ , then  $A \cap B = A$  and  $A \cup B = B$ . An important special case is that

$$A \cap A = A \cup A = A. \quad (1.8)$$

For the empty set  $\emptyset$ , we get  $A \cap \emptyset = \emptyset$  and  $A \cup \emptyset = A$ . For the sample space  $\Omega$ , we get  $A \cap \Omega = A$  and  $A \cup \Omega = \Omega$ .

Union and intersection are *commutative* operations like addition and multiplication, so the order of  $A$  and  $B$  does not matter:

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

and

$$A \cap B = B \cap A.$$

Events  $A$  and  $B$  are **disjoint** or **mutually exclusive** when  $A \cap B = \emptyset$ . If  $A$  and  $B$  are disjoint, then at most of one of them can occur in a single experiment. Any set and its complement are disjoint, and the empty set  $\emptyset$  is disjoint with itself and all other sets.

If  $\Omega$  is a population, these set operations allow us to define subpopulations in terms of multiple traits. If the event  $A = \{\omega \in \Omega : \omega \text{ lives in Ohio}\}$ , then its complement  $A^c$  contains all individuals in  $\Omega$  who live outside Ohio. If the event  $B = \{\omega \in \Omega : \omega \text{ is 42 years old}\}$ , then the intersection  $A \cap B$  contains everyone in  $\Omega$  who is 42 years old and lives in Ohio. If  $\Omega$  does not contain any 42-year-old Ohio residents, then  $A$  and  $B$  are disjoint. The union  $A \cup B$  contains everyone in  $\Omega$  who lives in Ohio or is 42 years old. This could include both a 24-year-old who lives Ohio and a 42-year-old who lives Michigan.

### 1.1.3 Venn diagrams

A useful tool for understanding events and set operations is the **Venn diagram**.<sup>4</sup> An example is shown in Figure 1.1. The rectangle represents  $\Omega$ , and the circles  $A$  and  $B$  represent events.  $A^c$  is everything in  $\Omega$  outside the circle  $A$ , and  $B^c$  is everything outside the circle  $B$ . Their intersection  $A \cap B$  is the area where the two circles overlap. Their union  $A \cup B$  is everything contained in at least one of  $A$  or  $B$ .

---

<sup>4</sup>Named after [John Venn](#) (1834-1923), an English logician and philosopher who was one of the pioneers of the frequentist interpretation of probability. He was ordained as an Anglican priest in 1859 but resigned from the church in 1883. He was a prize-winning gardener of roses and white carrots and a prominent supporter of women's right to vote. From 1903 until his death, he was President of Fellows in Gonville and Caius College at the University of Cambridge, where he is commemorated with a Venn diagram in a stained glass window.

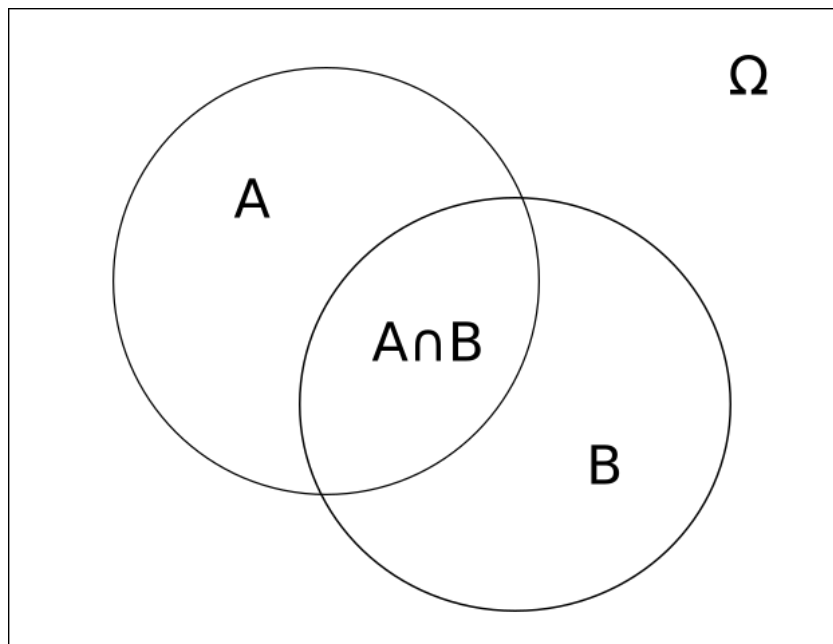


Figure 1.1: Venn diagram showing events  $A$  and  $B$ . The area contained in both events is their intersection  $A \cap B$ . The union  $A \cup B$  is all area contained in at least one of  $A$  and  $B$ , including  $A \cap B$ .

### 1.1.4 Sequences of events\*

Intersections can be written for more than two events. The intersection of  $A_1, A_2, \dots, A_n$  is

$$I_n = \bigcap_{i=1}^n A_i. \quad (1.9)$$

Because set intersection is commutative and associative, any ordering of  $A_1, \dots, A_n$  produces the same intersection. The event  $I_n$  occurs if and only if all of the events  $A_1, \dots, A_n$  occur. Each new event makes the intersection smaller (i.e., never larger) in the sense that

$$\bigcap_{i=1}^{n+1} A_i \subseteq I_n.$$

whenever  $A_{n+1}$  is another event.

Similarly, unions can be written for more than two events. If  $A_1, A_2, \dots, A_n$  is a set of events, then their union is

$$U_n = \bigcup_{i=1}^n A_i. \quad (1.10)$$

Because set union is commutative and associative, any ordering of  $A_1, \dots, A_n$  produces the same union. The event  $U_n$  occurs if and only if at least one of the events  $A_i$  occurs. Each new event makes the union bigger (i.e., never smaller) in the sense that

$$U_n \subseteq \bigcup_{i=1}^{n+1} A_i$$

whenever  $A_{n+1}$  is another event.

Both unions and intersections can be defined for infinite sequences of events.<sup>5</sup> To describe this, we let  $n = \infty$  in the notation from Equation 1.9 or Equation 1.10. The union of any finite sequence of events can be turned into the union of an infinite sequence of events by adding an endless sequence of empty sets to the finite sequence. The new sequence is still a sequence of disjoint events, and each empty set  $\emptyset$  leaves the union unchanged. If  $(A_1, A_2, \dots)$  is an infinite sequence of events such that  $A_i = \emptyset$  for all  $i > n$ , then

$$\bigcup_{i=1}^{\infty} A_i = \bigcup_{i=1}^n A_i.$$

This turns out to be useful when we try to give a mathematically rigorous definition of probability.

---

<sup>5</sup>In probability, we only consider unions and intersections of finite or countably infinite sets of events. Although unions and intersections can be defined for uncountably infinite sets of events, it can be impossible to assign probabilities to the resulting sets (see the [Banach-Tarski paradox](#)). As an epidemiologist, this should not keep you up at night.

### 1.1.5 Algebra of sets\*

Unions, intersections, and complements can be combined in complex ways. Fortunately, there are a few basic principles that can be used to simplify these calculations. We have already seen that unions and intersections are commutative. Unions and intersections are also *associative*, so

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

and

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

for any sets  $A$ ,  $B$ , and  $C$ .

*De Morgan's laws* describe how complements affect unions and intersections. If  $A$  and  $B$  are sets, then

$$(A \cap B)^c = A^c \cup B^c \tag{1.11}$$

because you are outside  $A \cap B$  if and only if you are outside  $A$  or outside  $B$ . Similarly,

$$(A \cup B)^c = A^c \cap B^c. \tag{1.12}$$

because you are outside  $A \cup B$  if and only if you are outside  $A$  and outside  $B$ . Note that each of these equations implies the other if we replace  $A = (A^c)^c$  with  $A^c$  and replace  $B = (B^c)^c$  with  $B^c$ . They are two sides of the same coin, but it is helpful to remember them both.

The *distributive properties* describe how unions and intersections interact with each other. Recall that multiplication distributes over addition, so  $a(b + c) = ab + ac$ . For any sets  $A$ ,  $B$ , and  $C$ , we have the following distributive properties:

- Intersections distribute over unions, so

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

- Unions distribute over intersections, so

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

Intersections and unions also distribute over themselves. However, this is a consequence of commutativity, associativity, and Equation 1.8, not a separate property like the distributive rules above.

## 1.2 Probability

A *probability measure* is a function that takes an event  $A \subseteq \Omega$  and returns a number  $\Pr(A) \in [0, 1]$  in any way that conforms to the following rules:

- $\Pr(\Omega) = 1$ .
- $\Pr(A) \in [0, 1]$  for any event  $A \subseteq \Omega$ .<sup>6</sup>
- The **addition rule**: If  $(A_1, A_2, \dots)$  is any sequence of disjoint events, then

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

The addition rule is stated in terms of an infinite sequence of disjoint events because this implies the addition rule for any finite sequence of disjoint events (see Section 1.1.4).

It is useful to think of probability as a generalization of our intuitions about area or volume. When there is no overlap in a set of two-dimensional shapes, we can get the total area they cover by adding up the areas of the individual shapes. Similarly, we can get the total volume taken up by a set of bowling balls by adding up their individual volumes.

There is a lot of debate about the meaning of probability, but its definition does not assume any particular interpretation. Probability calculations are based on the rules above no matter what we think it all means, and any interpretation consistent with these rules is valid.

### 1.2.1 Probability calculations

Several useful properties of probability follow immediately from the definition above. A short proof follows each result. To follow the proofs, it helps to draw Venn diagrams.

**Theorem 1.1.** *If  $A$  is an event,  $\Pr(A^c) = 1 - \Pr(A)$ .*

*Proof.* Because  $\Omega = A \cup A^c$  and  $A$  and  $A^c$  are disjoint, we have

$$\Pr(A) + \Pr(A^c) = \Pr(\Omega) = 1$$

by the addition rule. The result follows when we subtract  $\Pr(A)$  from both sides.  $\square$

**Theorem 1.2.** *If  $A$  and  $B$  are events such that  $A \subseteq B$ , then  $\Pr(A) = \Pr(B) - \Pr(B \cap A^c)$ . This implies that  $\Pr(A) \leq \Pr(B)$ .*

---

<sup>6</sup>Technically, we assign probabilities only to events in a set  $\mathcal{F}$  of subsets of  $\Omega$  that is required to contain  $\Omega$  and to be closed under complements and countable unions. “Closed under complements” means that  $A^c \in \mathcal{F}$  whenever  $A \in \mathcal{F}$ . For example,  $\emptyset = \Omega^c$  must be in  $\mathcal{F}$  because  $\Omega \in \mathcal{F}$ . “Closed under countable unions” means that  $\bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$  whenever  $(A_1, A_2, \dots)$  is a sequence of events in  $\mathcal{F}$ . The set  $\mathcal{F}$  is called a  $\sigma$ -algebra, and this restriction on the domain of probability helps avoid internal contradictions like the mind-blowing [Banach-Tarski paradox](#).

*Proof.* Each element of  $B$  either is or is not in  $A$ , so

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

where the second equality follows from the fact that  $B \cap A = A$  because  $A \subseteq B$ . The two sets on the right-hand side are disjoint, so we have

$$\Pr(B) = \Pr(A) + \Pr(B \cap A^c)$$

by the addition rule. The result follows if we subtract  $\Pr(B \cap A^c)$  from both sides. This implies that  $\Pr(A) \leq \Pr(B)$  because  $\Pr(B \cap A^c) \geq 0$ .  $\square$

**Theorem 1.3.** *If  $A$  and  $B$  are events,  $\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B)$ .*

*Proof.* We can break  $A \cup B$  into three disjoint sets: elements of  $A$  and not  $B$ , elements of  $B$  and not  $A$ , and elements of both  $A$  and  $B$ . In set notation, this is

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

By the addition rule,

$$\Pr(A \cup B) = \Pr(A \cap B^c) + \Pr(B \cap A^c) + \Pr(A \cap B). \quad (1.13)$$

By Theorem 1.2, we have

$$\Pr(A \cap B^c) = \Pr(A) - \Pr(A \cap B),$$

because  $A \cap B \subseteq A$  and

$$\Pr(B \cap A^c) = \Pr(B) - \Pr(A \cap B).$$

because  $A \cap B \subseteq B$ . The result follows from substituting these back into Equation 1.13 and collecting terms involving  $\Pr(A \cap B)$ . Intuitively,  $\Pr(A) + \Pr(B)$  includes the overlap  $\Pr(A \cap B)$  twice, so we have to subtract out one of them.  $\square$

## 1.3 Random variables

The outcomes of an experiment are not necessarily numbers. A **random variable** is a real-valued function defined on a sample space  $\Omega$ . In other words, a random variable  $X$  is a function that takes an *argument*  $\omega \in \Omega$  as input and returns a *value*  $X(\omega) \in \mathbb{R}$ . Traditionally, random variables are written as capital letters and possible values are written as lower-case letters, so  $\Pr(X = x)$  denotes the probability of the event

$$\{\omega \in \Omega : X(\omega) = x\}.$$

For simplicity, random variables are usually written without the argument  $\omega$ .

The distinction between outcomes and random variables is useful because we can define multiple random variables on the same sample space. For example, the height, weight, and age of an individual  $\omega$  sampled from a population  $\Omega$  are different random variables defined on the same sample space.



### 1.3.1 Indicator variables

The simplest random variables are **indicator variables**. For an event  $A$ , the indicator variable

$$\mathbb{1}_A(\omega) = \begin{cases} 1 & \text{if } \omega \in A, \\ 0 & \text{if } \omega \notin A. \end{cases}$$

Indicator variables are **binary** random variables, which take exactly two values. In practice, these values should be zero and one unless there is a specific reason to do otherwise. When sampling from a population, we can define indicator variables for membership in different subpopulations.

All of the basic set operations above can be expressed in terms of indicator variables for sets.

- The indicator function for the complement of  $A$  is

$$\mathbb{1}_{A^c} = 1 - \mathbb{1}_A. \quad (1.14)$$

- If  $B$  is another event and  $\mathbb{1}_B$  is its indicator variable, then the indicator variable for the intersection  $A$  and  $B$  is the product of their indicator variables:

$$\mathbb{1}_{A \cap B} = \mathbb{1}_A \mathbb{1}_B. \quad (1.15)$$

- The indicator variable for the union  $A \cup B$  is

$$\mathbb{1}_{A \cup B} = 1 - (1 - \mathbb{1}_A)(1 - \mathbb{1}_B) = \mathbb{1}_A + \mathbb{1}_B - \mathbb{1}_{A \cap B}. \quad (1.16)$$

This follows from Equation 1.12 because  $A \cup B = (A^c \cap B^c)^c$ .

## 1.4 R

### 1.4.1 Probability distributions

The set of possible values of a random variable  $X$  is called the *support* of  $X$  and denoted  $\text{supp}(X)$ .<sup>7</sup> For example, the support of an indicator variable is  $\{0, 1\}$ . In this section, we will focus on **discrete** random variables, which have a support on a finite or countably infinite set. There are two standard ways to describe the distribution of a discrete random variable:

---

<sup>7</sup>Technically, the support of  $X$  is the smallest closed set  $S_X$  such that  $\Pr(X \in S_X) = 1$ . For a discrete random variable with support on a finite set, it is just the set of possible values. For a discrete random variable with support on a countably infinite set, it can include points whose probability mass is zero—a pathological case that we can safely ignore. For a continuous random variable, it can include values whose probability density is zero—a case that is not unusual or pathological.

- The **probability mass function** (PMF) of a discrete random variable  $X$  is

$$f(x) = \begin{cases} \Pr(X = x) > 0 & \text{if } x \in \text{supp}(X), \\ 0 & \text{if } x \notin \text{supp}(X). \end{cases}$$

Because  $\Pr(\Omega) = 1$ , we always have

$$\sum_{x \in \text{supp}(X)} f(x) = 1.$$

- The **cumulative distribution function** (CDF) of  $X$  is

$$F(x) = \Pr(X \leq x).$$

$F(x)$  is monotonically increasing in  $x$ , which means that  $F(a) \leq F(b)$  whenever  $a < b$ . It has a jump upward of size  $f(x)$  at each  $x \in \text{supp}(X)$ , and its value at each such  $x$  is the value that it jumps to—not the value that it jumps up from. For sufficiently small  $x$ ,  $F(x)$  can be made arbitrarily close to zero. For sufficiently large  $x$ ,  $F(x)$  can be made arbitrarily close to one. More formally, we say that  $\lim_{x \downarrow -\infty} F(x) = 0$  and  $\lim_{x \uparrow \infty} F(x) = 1$ .

The PMF and CDF provide equivalent descriptions of the distribution of  $X$  in the sense that either of these functions can be used to calculate the other. Given the PMF  $f$ , the CDF is defined by

$$F(x) = \sum_{\substack{v \in \text{supp}(X): \\ v \leq x}} f(v).$$

where the sum is taken over all  $u \in \text{supp}(X)$  such that  $u \leq x$ . Given the CDF  $F$ , the PMF is defined by

$$f(x) = F(x) - \max_{v \leq x} F(v)$$

where the maximum is  $F(v)$  for the largest  $v \in \text{supp}(X)$  such that  $v < x$ .

## 1.4.2 Mean

The **mean** or *expected value* of a random variable  $X$  is

$$\mathbb{E}(X) = \sum_{x \in \text{supp}(X)} x \Pr(X = x) = \sum_{x \in \text{supp}(X)} x f(x),$$

where  $f$  is the PMF of  $X$ . The mean is often written  $\mu$ , and it is often described as a measure of the “location” or “central tendency” of  $X$ .

Indicators are an extremely useful for calculating probabilities using means. For any event  $A$ , its probability is the mean of the indicator variable  $\mathbb{1}_A$ :

$$\Pr(A) = 0 \Pr(\mathbb{1}_A = 0) + 1 \Pr(\mathbb{1}_A = 1) = \mathbb{E}(\mathbb{1}_A).$$

This is a common way to calculate probabilities in data analyses.

## 1.5 R

### 1.5.1 Variance

If  $X$  has  $\mathbb{E}(X) = \mu$ , then  $(X - \mu)^2$  is another random variable. The **variance** of  $X$  is the expected value of  $(X - \mu)^2$ :

$$\text{Var}(X) = \mathbb{E}[(X - \mu)^2] = \sum_{x \in \text{supp}(X)} (x - \mu)^2 f(x).$$

Because  $(x - \mu)^2 \geq 0$  with equality if and only if  $x = \mu$ , we always have  $\text{Var}(X) \geq 0$ . We have  $\text{Var}(X) = 0$  if and only if  $X = \mu$  with probability one. The variance is often written  $\sigma^2$ , and it is often described as a measure of the dispersion of  $X$  around the mean.

The square root of the variance is called the **standard deviation**, which is often written  $\sigma$ . If a random variable  $X$  has units (e.g., length, weight, or time), the mean and the standard deviation have the same units as  $X$ . For example, the mean and standard deviation of a length in meters both have units of meters but the variance has units of meters<sup>2</sup>.

### 1.5.2 Bernoulli distribution

The distribution of an indicator variable is called the **Bernoulli distribution**.<sup>8</sup> A random variable with the Bernoulli( $p$ ) distribution has the PMF

$$f(x) = \begin{cases} 1 - p & \text{if } x = 0 \\ p & \text{if } x = 1. \end{cases}$$

Equivalently, it has the CDF

$$F(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 - p & \text{if } x \in [0, 1) \\ 1 & \text{if } x \geq 1. \end{cases}$$

If a random variable  $X$  has a Bernoulli( $p$ ) distribution, we write  $X \sim \text{Bernoulli}(p)$ . The indicator variable for an event  $A$  has a Bernoulli distribution with  $p = \Pr(A)$ .

If  $X \sim \text{Bernoulli}(p)$ , then it has mean

$$\mathbb{E}(X) = 0 \times (1 - p) + 1 \times p = p$$

---

<sup>8</sup>Named after [Jacob Bernoulli](#) (1655-1705), a Swiss mathematician who derived the first version of the law of large numbers and discovered the constant  $e \approx 2.718281828$ , which is the base for natural logarithms. He and his younger brother Johann Bernoulli (1667-1748) were some of the first mathematicians to try to understand and apply calculus, but their relationship eventually curdled into a jealous rivalry. A lunar impact crater called Bernoulli is named jointly after them.

and variance

$$\text{Var}(X) = (0 - p)^2(1 - p) + (1 - p)^2p = p(1 - p).$$

Its standard deviation is  $\sqrt{p(1 - p)}$ , which is greater than zero unless  $p = 0$  or  $p = 1$ . If  $p = 0$ , then  $X = 0$  with probability one. If  $p = 1$ , then  $X = 1$  with probability one.

## 1.6 Joint and marginal distributions

If  $X$  and  $Y$  are random variables defined on the same probability space, then their **joint** probability mass function is

$$f(x, y) = \Pr(X = x \text{ and } Y = y) = \Pr(\{\omega : X(\omega) = x \text{ and } Y(\omega) = y\}).$$

The **marginal** probability mass functions are the PMFs of  $X$  or  $Y$  individually, which can be calculated from the joint PMF. The marginal PMF of  $X$  is

$$f_X(x) = \sum_{y \in \text{supp}(Y)} f(x, y),$$

and the marginal PMF of  $Y$  is

$$f_Y(y) = \sum_{x \in \text{supp}(X)} f(x, y).$$

These are called marginal distributions by analogy to the margins of a table. The distinction between joint and marginal distributions is extremely important in epidemiology and other applications of probability.

For example, Table 1.1 shows the joint and marginal PMFs for two binary random variables  $X$  and  $Y$ . By definition,

$$f(0, 0) + f(0, 1) + f(1, 0) + f(1, 1) = 1.$$

In the table, it is clear that the joint distribution determines the marginal distributions. However, there are many different joint distributions that are consistent with the same marginal distributions. Thus, the marginal distributions do not determine the joint distribution.<sup>9</sup>

---

<sup>9</sup>This becomes a fundamental insight when we discuss hypothesis tests for independence as well as confounding and selection bias.

Table 1.1: Joint and marginal PMFs for binary random variables  $X$  and  $Y$ .

	$Y = 0$	$Y = 1$	$X$ margin
$X = 0$	$f(0, 0)$	$f(0, 1)$	$f_X(0) =$ $f(0, 0) + f(0, 1)$
$X = 1$	$f(1, 0)$	$f(1, 1)$	$f_X(1) =$ $f(1, 0) + f(1, 1)$
$Y$ margin	$f_Y(0) =$ $f(0, 0) + f(1, 0)$	$f_Y(1) =$ $f(0, 1) + f(1, 1)$	1

## 1.7 R

Joint distributions can be defined for more than two random variables. If  $X_1, X_2, \dots, X_n$  are random variables defined on the same sample space, then their joint PMF is

$$f(x_1, x_2, \dots, x_n) = \Pr(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n).$$

The marginal distribution of each  $X_i$  can be found by adding up the PMF over the support of all the other random variables. For example,

$$f_{X_2}(x_2) = \sum_{x_1 \in \text{supp}(X_1)} \sum_{x_3 \in \text{supp}(X_3)} f(x_1, x_2, x_3).$$

when  $n = 3$ . In this same case, we can talk about the joint distribution of any two variables marginalized over the third. For example,

$$f_{X_2, X_3}(x_2, x_3) = \sum_{x_1 \in \text{supp}(X_1)} f(x_1, x_2, x_3).$$

For larger  $n$ , the formulas gets uglier but the ideas are the same.

### 1.7.1 Linear combinations\*

If  $a$  and  $b$  are constants, then  $aX + bY$  is another random variable on  $\Omega$ . It is called a *linear combination* of  $X$  and  $Y$ . Linear combinations can be defined for more than two random variables. If  $X_1, \dots, X_n$  are random variables defined on a sample space and  $a_1, \dots, a_n$  are constants, then

$$\sum_{i=1}^n a_i X_i = a_1 X_1 + a_2 X_2 + \dots + a_n X_n$$

is a linear combination of  $X_1, \dots, X_n$ . The constants can be any real numbers, including one and zero.

Section 1.3.1 contains both examples and non-examples of linear combinations of random variables.

- The indicator function for  $A^C$  in Equation 1.14 is a linear combination of  $\mathbb{1}_A$  and the random variable  $\mathbb{1}_\Omega$ , which equals one for all  $\omega \in \Omega$ .
- The indicator function for  $A \cup B$  in Equation 1.16 is linear combination of the indicator variables  $\mathbb{1}_A$ ,  $\mathbb{1}_B$ , and  $\mathbb{1}_{A \cap B}$ .
- The indicator function for  $A \cap B$  in Equation 1.15 is not a linear combination of  $\mathbb{1}_A$  and  $\mathbb{1}_B$  because we have to multiply these two variables.

If  $X$  and  $Y$  are random variables defined on the same sample space and  $a$  and  $b$  are constants, the mean of the linear combination  $aX + bY$  is

$$\mathbb{E}(aX + bY) = a \mathbb{E}(X) + b \mathbb{E}(Y).$$

This is a direct consequence of the definition of expected value:

$$\begin{aligned} \mathbb{E}(aX + bY) &= \sum_{x \in \text{supp}(X)} \sum_{y \in \text{supp}(Y)} (ax + by) f(x, y) \\ &= a \sum_{x \in \text{supp}(X)} \left( x \sum_{y \in \text{supp}(Y)} f(x, y) \right) + b \sum_{y \in \text{supp}(Y)} \left( y \sum_{x \in \text{supp}(X)} f(x, y) \right) \\ &= a \sum_{x \in \text{supp}(X)} x f_X(x) + b \sum_{y \in \text{supp}(Y)} y f_Y(y). \end{aligned}$$

The algebra is not pretty, but the logic is straightforward. We split up the sum into parts depending only on  $x$  and only on  $y$  outside the joint PMF. In each part, we factor out a constant and find the marginal PMF. This same logic extends to a linear combination of any number of random variables.

### 1.7.2 Variance and covariance\*

The variance of  $aX + bY$  is

$$\text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y)$$

where

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

is called the **covariance** of  $X$  and  $Y$ . Note that  $\text{Cov}(X, Y) = \text{Cov}(Y, X)$ . Because  $\text{Var}(X) = \text{Cov}(X, X)$ , variance is a special case of covariance.

The joint distribution of  $X$  and  $Y$  has a **covariance matrix** which is

$$\begin{bmatrix} \text{Var}(X) & \text{Cov}(X, Y) \\ \text{Cov}(X, Y) & \text{Var}(Y) \end{bmatrix}$$

The variances are along the diagonal of the matrix, and the covariances appear off the diagonal. Because  $\text{Cov}(X, Y) = \text{Cov}(Y, X)$ , covariance matrices are always symmetric (i.e., symmetric across the diagonal). Covariance matrices are an extremely useful tool for calculating the variances of linear combinations of random variables. For example:

$$\text{Var}(aX + bY) = \begin{pmatrix} a & b \end{pmatrix} \begin{bmatrix} \text{Var}(X) & \text{Cov}(X, Y) \\ \text{Cov}(X, Y) & \text{Var}(Y) \end{bmatrix} \begin{pmatrix} a \\ b \end{pmatrix}$$

in matrix and vector notation from [linear algebra](#). This logic extends to linear combinations of any number of random variables.

The covariance is the numerator of the *Pearson correlation coefficient*,<sup>10</sup> which is

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

Because of the [Cauchy-Schwarz inequality](#), it turns out that  $\rho_{XY} \in [-1, 1]$ .

- We get  $\rho_{XY} = -1$  if and only if  $Y = cX$  for some negative constant  $c$ .
- We get  $\rho_{XY} = 1$  if and only if  $Y = cX$  for some positive constant  $c$ . For example,  $\rho_{XX} = 1$  for any random variable  $X$ .
- We get  $\rho_{XY} = 0$  if  $X$  and  $Y$  are *independent* in the sense that the value of one tells us nothing about the value of the other.<sup>11</sup> However, it is possible to have  $\rho_{XY} = 0$  when  $X$  and  $Y$  are not independent.

If we divide each entry  $\text{Cov}(X, Y)$  in a covariance matrix by  $\sqrt{\text{Var}(X) \text{Var}(Y)}$ , when we get a *correlation matrix*. Any correlation matrix is symmetric, and the entries along its diagonals are all ones.

## 1.8 Probability and disease occurrence

In epidemiology, there are two fundamental measures of disease occurrence that are probabilities: **prevalence** and **risk**. In both cases, our experiment is to sample an individual  $\omega$  from a population  $\Omega$ . The *disease outcome* is a binary random variable

$$D(\omega) = \begin{cases} 1 & \text{if } \omega \text{ has the disease outcome,} \\ 0 & \text{otherwise.} \end{cases}$$

<sup>10</sup>Named after [Karl Pearson](#) (1857-1936), an English mathematician who founded the modern discipline of mathematical statistics. In 1911, he started the world's first university department of statistics at University College London. He was an outspoken socialist and supporter of women's rights, but he was also a vocal proponent of social Darwinism and eugenics who opposed Jewish immigration into Britain.

<sup>11</sup>We will define independence of random variables more rigorously when we discuss conditional probabilities in Chapter 2.

The set of individuals in  $\Omega$  who have  $D(\omega) = 1$  is an event in  $\Omega$ , and our measure of disease occurrence is

$$\Pr(\{\omega \in \Omega : D(\omega) = 1\}).$$

The most important difference between prevalence and risk is the role of time in the definition of  $D$ .

There is an important technical detail to remember when we talk about disease onset and recovery. When a person has disease onset at time  $t^{\text{ons}}$  and recovers at time  $t^{\text{rec}}$ , they have disease for each  $t \in [t^{\text{ons}}, t^{\text{rec}})$ . We assume that  $t^{\text{rec}} > t^{\text{ons}}$  so this interval is nonempty. We let the onset and recovery times for person  $i$  be  $t_i^{\text{ons}}$  and  $t_i^{\text{rec}}$ , respectively. If a person has multiple episodes of the disease, each episode has its own  $t^{\text{ons}}$  and  $t^{\text{rec}}$ . For example, the  $j^{\text{th}}$  episode in person  $i$  would have onset time  $t_{ij}^{\text{ons}}$  and recovery time  $t_{ij}^{\text{rec}}$ .

The time scale used to define disease onset is flexible, and this flexibility is useful. The most obvious time scale is *calendar time* or *absolute time*. Another common time scale is age, which is an important determinant of the risk of many diseases. In some cases, time since an event is a useful time scale. The event that defines time scale could be a single event (e.g., exposure to contaminated food at a party) or an event that occurs at different times for different individuals (e.g., time since menopause). In general, it is wise to choose the time scale that corresponds to the most important time-varying determinant of disease onset. The chosen time scale is often called the *analysis time scale*.

### 1.8.1 Prevalence

For prevalence, the disease outcome is defined by choosing a time  $t$  and letting

$$D(\omega) = \begin{cases} 1 & \text{if } \omega \text{ has disease at time } t, \\ 0 & \text{otherwise.} \end{cases}$$

In other words, it is the proportion of the population  $\Omega$  that disease at time  $t$ . This includes individuals who have disease onset at time  $t^{\text{ons}} = t$  but not individuals who recover from disease at time  $t^{\text{rec}} = t$ . This is often called the **point prevalence** at time  $t$ .

Another version of prevalence is **period prevalence**. For period prevalence, we choose a nonempty time interval  $(t_a, t_b]$  and define

$$D(\omega) = \begin{cases} 1 & \text{if } \omega \text{ has disease at any time } t \in (t_a, t_b], \\ 0 & \text{otherwise.} \end{cases}$$

In other words, it is the proportion of the population that has disease at any time in the interval  $(t_a, t_b]$ . This includes prevalent cases at time  $t_a$  and cases with disease onset in  $(t_a, t_b]$ . The period prevalence in  $(t_a, t_b]$  is the point prevalence at  $t_a$  plus the risk of disease onset in  $(t_a, t_b]$ , to which we now turn.



## 1.9 R

### 1.9.1 Risk (cumulative incidence) and the survival function

To define **risk** or **cumulative incidence**, we first choose a nonempty time interval  $(t_a, t_b]$ . The disease outcome is defined as

$$D(\omega) = \begin{cases} 1 & \text{if } \omega \text{ has } t^{\text{ons}} \in (t_a, t_b], \\ 0 & \text{otherwise.} \end{cases}$$

In the population that is disease-free and at risk of disease at time  $t_a$ , it is the proportion who have disease onset at  $t^{\text{ons}} \leq t_b$ . The risk is sometimes called the *incidence proportion*.

The risk depends on a specified interval  $(t_a, t_b]$ . We can always define our time scale so that  $t_a = 0$ , so the risk in  $(t_a, t_b]$  on the original time scale is the same as the risk in the interval  $(0, t_b - t_a]$  on the analysis time scale. On the analysis time scale, the **cumulative incidence function**  $F(t)$  is the risk of disease in  $(0, t]$  for any possible  $t$ . The corresponding **survival function** is

$$S(t) = 1 - F(t),$$

which is the probability of no disease onset in  $(0, t]$ . In practice, it is often easier to calculate the survival function than to calculate the cumulative incidence function directly. There is only one way to survive disease-free through the interval  $(0, t]$ , but you can have disease onset at any time.

## 1.10 R

The survival function has several important properties:

- $S(0) = 1$  because  $(0, 0]$  is an empty interval where no one can have disease onset.
- Because  $S(t)$  is a probability,  $S(t) \in [0, 1]$  for all  $t$ .
- $S(t)$  monotonically decreases (i.e., never increases) with increasing  $t$ . If  $t_a < t_b$ , then the time interval  $(0, t_a]$  is contained  $(0, t_b]$ . Everyone who survives disease-free through  $(0, t_b]$  must have survived disease-free through  $(0, t_a]$ , but some people who survived through  $(0, t_a]$  might not make it all the way through  $(0, t_b]$ . Thus,  $S(t_a) \geq S(t_b)$  whenever  $t_a < t_b$ .
- If the disease or event occurs eventually for all individuals in our population  $\Omega$  (e.g., death), then  $S(t) \rightarrow 0$  as  $t \rightarrow \infty$ .

Each of these probabilities follows directly from the definition of  $S(t)$ . Similarly, the cumulative incidence function  $F$  has  $F(0) = 0$  and  $F(t) \in [0, 1]$ , and it is monotonically increasing (i.e., never decreasing) with increasing  $t$ . If the disease or event occurs eventually in all individuals, then  $F(t) \rightarrow 1$  as  $t \rightarrow \infty$ . Figure 1.2 shows the survival and cumulative hazard curves for the data generated in the prevalence example above.

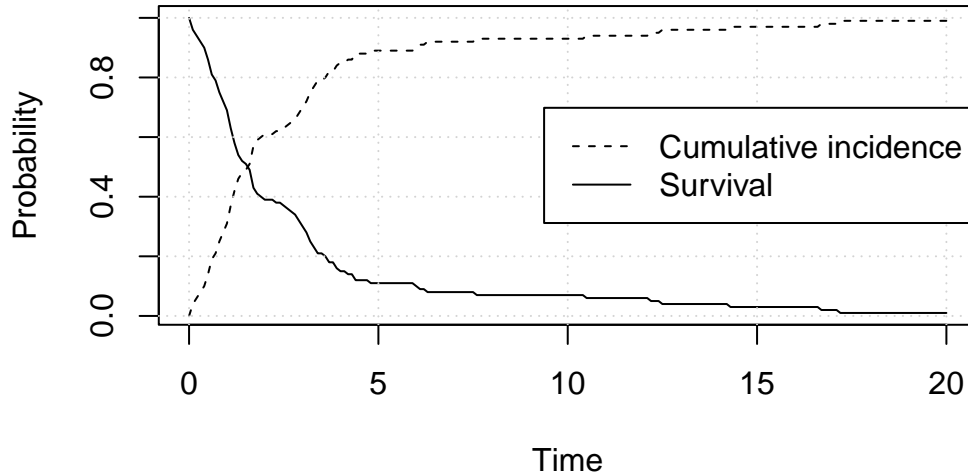


Figure 1.2: Survival and cumulative incidence curves for the data from the prevalence example.

Here, I will generally use the word “risk” to refer to the probability of disease onset in a specified interval. When there is possible confusion about the meaning of “risk”, I will use “cumulative incidence” instead. The terms “cumulative incidence function” and “survival function” are standard in survival analysis, which is the branch of statistics that studies times to events. The creative use of “risk” in public health and medicine should not make you shy away from using the word correctly.

### 1.10.1 Prevalence and the duration of disease

Point and period prevalence are both affected by the duration of disease. Both measures will increase if the duration of disease increases. A simple illustration of this is given in Figure 1.3. For a fixed set of onset times, the point prevalence of disease at any time  $t$  either stays the same or increases when the duration of disease increases. The prevalence at time  $t = 5$  is  $\frac{2}{5} = 0.4$  under the shorter duration of disease but  $\frac{3}{5} = 0.6$  under the longer duration of disease. Period prevalence over any interval  $(t_a, t_b]$  is affected by the duration of disease because it is the point prevalence at  $t_a$  (which is affected by disease duration) plus the risk of disease onset over  $(t_a, t_b]$ . In a given population, the relationship between prevalence, frequency of disease onset (incidence), and the duration of disease can be complex (Freeman and Hutchison 1980; Preston 1987; Keiding 1991; Alho 1992). The risk of disease in any given interval is not affected by the duration of disease.

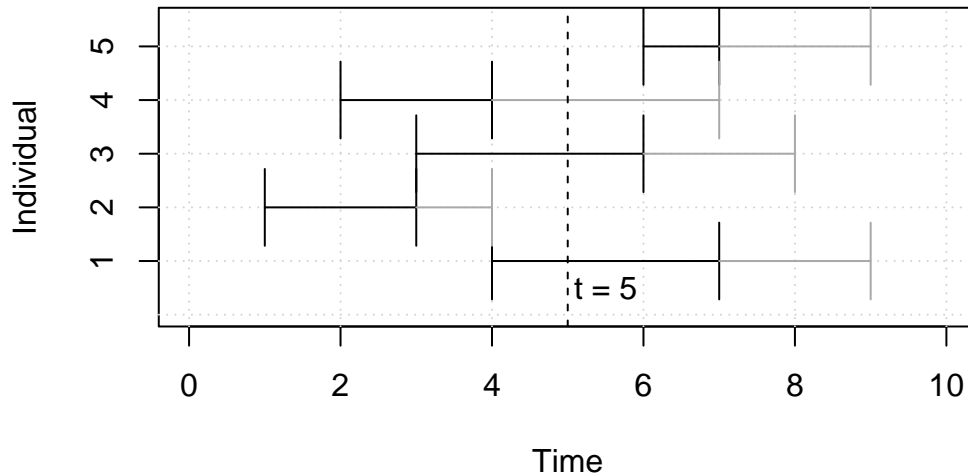


Figure 1.3: Each black horizontal line shows the onset of disease and recovery from disease in a single individual. The gray lines show recoveries from disease if the disease duration increases.

### 1.10.2 Descriptive and analytic epidemiology

Prevalence is often a useful measure for **descriptive epidemiology**, which measures the distribution of disease over person, place, and time. Because prevalence depends on both incidence and duration of disease, a change in the prevalence of disease can generally be explained several different ways (MacMahon and Terry 1958; Dunn Jr 1962). For example, an increase in prevalence of human immunodeficiency virus (HIV) infection could be caused by an increase in the incidence of HIV infection (which is bad) or an increase in the life expectancy of HIV-infected people (which is good).

Risk (cumulative incidence) is generally more useful than prevalence for **analytic epidemiology**, which attempts to identify the causes of a disease. Another advantage of risk is that it can be used for outcomes that begin and end very quickly (e.g., traffic accidents or being hit by lightning) and for outcomes that remove individuals from the population (e.g., emigration or death). Prevalence is not a useful measure of the public health impact of these events.

---

**Listing 1.1** indicators.R

---

```
## Indicator variables for events A and B, etc.

# Setting the seed ensures that everyone gets the same random samples.
# Functions are called using parentheses (round brackets).
# The function rbinom() is a random sample from a binomial distribution.
set.seed(42)
n <- 100
dat <- data.frame(A = rbinom(n, 1, 0.3))
dat$B <- rbinom(n, 1, 0.6)

# inspecting a data frame
names(dat) # variables in the data frame
nrow(dat)  # number of rows (individuals)
ncol(dat)  # number of columns (variables)
dim(dat)   # rows and columns in the data frame
str(dat)   # summary of the data frame structure (variables and types)

# inspecting columns of a data frame (or vectors)
# Our sample space or population consists of 100 individuals.
# Square brackets are used for indices, which can be numbers or TRUE/FALSE.
dat$A      # indicator for A for all 100 individuals
dat$A[10]  # indicator for A in individual 10
dat$A[2:6] # indicator variables for individuals 2 to 6
dat$A[c(10, 20, 30)] # A indicators for individuals 10, 20, and 30
which(dat$A == 1)   # which individuals are in event A
which(dat$A == 0)   # which individuals are not in event A

# indicator variable for A complement
# In R (and many other languages), "!" means "not".
# The function as.integer() changes TRUE/FALSE to 1/0.
dat$Acomp <- as.integer(!dat$A)

# indicator variable for A intersection B
# In R (and many other languages), "&" means "and".
dat$ABintersect <- as.integer(dat$A & dat$B)

# indicator variable for A union B
# In R (and many other languages), "|" means "or".
dat$ABunion <- as.integer(dat$A | dat$B)

# save the data frame as a CSV file
# The file argument can be a path (e.g., "./data/indicators.csv" in Linux).
write.csv(dat, file = "indicators.csv", row.names = FALSE)
```

---

**Listing 1.2** probabilities.R

---

```
## Indicator variables and probability calculations

# read in CSV file with indicator variables using the function read.csv()
# The argument can be a path (e.g., "./data/indicators.csv" in Linux).
dat <- read.csv("indicators.csv")

# calculate probabilities from indicator variables using the function mean()
# This will also work with TRUE/FALSE (i.e., logical) variables, which are
# converted to TRUE = 1 and FALSE = 0 in calculations.
prob_A <- mean(dat$A)
prob_B <- mean(dat$B)
prob_Acomp <- mean(dat$Acomp)
prob_ABintersect <- mean(dat$ABintersect)
prob_ABunion <- mean(dat$ABunion)

# Pr(A complement) = 1 - Pr(A)
prob_Acomp
1 - prob_A

# Pr(A union B) = Pr(A) + Pr(B) - Pr(A intersect B)
prob_ABunion
prob_A + prob_B - prob_ABintersect

# Beware of numerical error when comparing floating-point numbers!
# This example is from The R Inferno by Patrick Burns.
# https://www.burns-stat.com/pages/Tutor/R_inferno.pdf
0.1 == 0.3 / 3
sprintf("%.20f", 0.1)
sprintf("%.20f", 0.3 / 3)

# math can be more accurate than computers (which is not their fault)
prob_ABunion == prob_A + prob_B - prob_ABintersect
sprintf("%.20f", prob_ABunion)
sprintf("%.20f", prob_A + prob_B - prob_ABintersect)
```

---

---

**Listing 1.3** jointdist.R

---

```
## Joint and marginal distributions of indicators for events A and B

# read indicator variable data from the CSV file
dat <- read.csv("indicators.csv")
n <- nrow(dat)

# tables of counts
# Putting "<name> = " before the vector creates a label.
table(A = dat$A)
table(B = dat$B)

# joint table of counts
# In table(), the first argument defines rows and the second defines columns.
# The addmargins() functions adds the row, column, and overall sums.
table(A = dat$A, B = dat$B)
addmargins(table(A = dat$A, B = dat$B))

# tables of probabilities
# Table margins match the distributions of A (rows) and B (columns).
table(Adist = dat$A) / n      # marginal distribution of A indicator
table(Bdist = dat$B) / n      # marginal distribution of B indicator
addmargins(table(A = dat$A, B = dat$B)) / n  # joint distribution
```

---

---

**Listing 1.4** prevalence.R

---

```
## Point and period prevalence

# generate onset and recovery data for 100 individuals
# Setting the seed ensures that everyone gets the same random numbers,
# but it is strictly optional.
# The function rexp() randomly samples from an exponential distribution.
set.seed(42)
cohort <- data.frame(onset = rexp(100, rate = 0.4))
cohort$duration <- rexp(100, rate = 2)
cohort$recovery <- cohort$onset + cohort$duration

# statistical summaries (mean, quartiles, range)
summary(cohort$onset)
summary(cohort$duration)
summary(cohort$recovery)

# highest and lowest recovery times
# The function sort() sorts the vector from lowest to highest.
# head() returns the first 6 values of a vector; tails() returns the last 6.
min(cohort$onset)
head(sort(cohort$onset))      # lowest 6 values (first 6 in the sorted vector)
tail(sort(cohort$onset))      # highest 6 values (last 6 in the sorted vector)
max(cohort$onset)

# With a long vector, sorting repeatedly can be slow.
# You can also control the number of elements returned by head() or tail().
onset_ordered <- sort(cohort$onset)
head(onset_ordered, n = 10)
tail(onset_ordered, n = 10)

# seeing rows and columns of the data frame
cohort[1:10, c("onset", "duration", "recovery")]
cohort[c(10, 20, 50), c("onset", "recovery")]
cohort[which(cohort$recovery < 1), c("onset", "recovery")]
cohort[, c("onset", "recovery")]      # all rows
cohort[c(2, 3, 5, 7, 11), ]          # all columns

# point prevalence
prev <- function(t) {
  # vector of TRUE/FALSE for prevalent cases at time t
  prevalent <- cohort$onset <= t & cohort$recovery > t
  mean(prevalent)
}

prev(0)
prev(1)
prev(2)
prev(6)

# period prevalence
# The parentheses around the logical tests are just for readability.
```

---

**Listing 1.5 risk.R**

---

```
## Risk, survival function, and cumulative incidence function

# read data from CSV file
# Change or remove ".R/" in the path as needed to locate the cohort.csv file.
# You can also re-generate the data as in prevalence.R using the same seed.
cohort <- read.csv("./R/cohort.csv")

# risk (cumpulative incidence)
risk <- function(t) {
  # vector of TRUE/FALSE for incident cases in (0, t]
  incident <- cohort$onset <= t
  mean(incident)
}

risk(0)
risk(1)
risk(2)
risk(6)

# cumulative incidence function
# Vectorize() takes a function like risk() that takes a single number as input
# and creates a function that can take a number or vector as input.
cuminc <- Vectorize(risk)
cuminc(c(0, 1, 2, 6))

# survival function
# A simple function can be put on one line.
# It takes the same input as cuminc(), so it can take a vector
surv <- function(t) 1 - cuminc(t)
surv(c(0, 1, 2, 6))

# plot the survival and cumulative incidence functions
t <- seq(0, 20, by = 0.1)
plot(t, surv(t), type = "l",
      xlab = "Time", ylab = "Probability")
lines(t, cuminc(t), lty = "dashed")
grid()
legend("right", bg = "white", lty = c("dashed", "solid"),
      legend = c("Cumulative incidence", "Survival"))
```

---



---

**Listing 1.6** surv-fig.R

---

```
## Plot of survival and cumulative incidence functions

# read data from CSV file
# Change or remove ".R/" in the path as needed to locate the cohort.csv file.
# You can also re-generate the data as in prevalence.R using the same seed.
cohort <- read.csv("./R/cohort.csv")

# risk (cumpulative incidence)
risk <- function(t) {
  # vector of TRUE/FALSE for incident cases in (0, t]
  incident <- cohort$onset <= t
  mean(incident)
}

# cumulative incidence function
cuminc <- Vectorize(risk)

# survival function
surv <- function(t) 1 - cuminc(t)

# plot the survival and cumulative incidence functions
t <- seq(0, 20, by = 0.1)
plot(t, surv(t), type = "l",
      xlab = "Time", ylab = "Probability")
lines(t, cuminc(t), lty = "dashed")
grid()
legend("right", bg = "white", lty = c("dashed", "solid"),
      legend = c("Cumulative incidence", "Survival"))
```

---

---

**Listing 1.7** prevdur-fig.R

---

```
## R code for prevalence and duration plot
plot(0, 0, type = "n", xlim = c(0, 10), ylim = c(0, 5.5),
     xlab = "Time", ylab = "Individual", yaxt = "n")
Axis(side = 2, at = 1:5, labels = 1:5)
grid()
start <- c(4, 1, 3, 2, 6)
stop1 <- c(7, 3, 6, 4, 7)
stop2 <- c(9, 4, 8, 7, 9)
arrows(x0 = start, y0 = 1:5, x1 = stop1, code = 3, length = 0.2, angle = 90)
arrows(x0 = stop1, y0 = 1:5, x1 = stop2, code = 2, length = 0.2, angle = 90,
      col = "darkgray")
abline(v = 5, lty = "dashed")
text(5.5, 0.5, label = "t = 5")
```

---

## 2 Conditional Probability and Diagnostic Tests

The probability that two subsequent events will happen is a ratio compounded of the probability of the 1st and the probability of the 2d on supposition the 1st happens. (Bayes 1763)<sup>1</sup>

---

<sup>1</sup>Thomas Bayes (1701-1761) was an English Presbyterian minister from a family of Nonconformists (i.e., Protestants who did not observe the rules of the Church of England). He studied logic and theology at the University of Edinburgh and served as a minister in Tunbridge Wells near Kent, England. He was elected a Fellow of the Royal Society in 1742 for his defense of Newton's calculus against a 1734 book called *The Analyst: A Discourse Addressed to an Infidel Mathematician* by Bishop George Berkeley (1685-1753). Late in life, Bayes became interested in probability and "inverse probability" (statistics). This essay was published posthumously. See [https://en.wikipedia.org/wiki/Thomas\\_Bayes](https://en.wikipedia.org/wiki/Thomas_Bayes).

### 3 Maximum Likelihood Estimation

Far better an approximate answer to the *right* question, which is often vague, than an *exact* answer to the wrong question, which can always be made precise. (Tukey 1962)<sup>1</sup>

---

<sup>1</sup>John Tukey (1915-2000) was an American mathematician and statistician who worked at Bell Labs and Princeton University. He developed the box plot, Tukey's range test for multiple comparisons, and the fast Fourier transform. In 1947, he coined the term "bit" as shorthand for "binary digit". See [https://en.wikipedia.org/wiki/John\\_Tukey](https://en.wikipedia.org/wiki/John_Tukey).

## 4 Bayesian Estimation

In the null hypothesis schema we are trying only to nullify something: “The null hypothesis is never proved or established but is possibly disproved in the course of experimentation.” But ordinarily evidence does not take this form. With the *corpus delicti* in front of you, you do not say, “Here is evidence against the hypothesis that no one is dead.” You say, “Evidently someone has been murdered.” (Berkson 1942)<sup>1</sup>

---

<sup>1</sup>Joseph Berkson (1899–1982) was an American physician and statistician at the Mayo Clinic in Rochester, Minnesota. He helped develop and popularize the use of logistic regression for binary outcomes, coining the term “logit” for the log odds in 1944. He also pioneered the study of selection bias, a special case of which is called “Berkson’s bias”. Later, he became a prominent opponent of the idea that smoking causes lung cancer. See [https://en.wikipedia.org/wiki/Joseph\\_Berkson](https://en.wikipedia.org/wiki/Joseph_Berkson).

## 5 Longitudinal Data and Rates

## 6 Survival Analysis

## **Part II**

# **B. Two-Sample Inference and Study Design**



## **Part III**

### **C. Principles of Causal Inference**

## **Part IV**

### **D. Epidemiologic and Statistical Methods for Causal Inference**

# References

- Alho, Juha M. 1992. "On Prevalence, Incidence, and Duration in General Stable Populations." *Biometrics* 48 (2): 587–92.
- Bayes, Thomas. 1763. "LII. An Essay Towards Solving a Problem in the Doctrine of Chances. By the Late Rev. Mr. Bayes, FRS Communicated by Mr. Price, in a Letter to John Canton, AMFRS." *Philosophical Transactions of the Royal Society of London* 53: 370–418.
- Berkson, Joseph. 1942. "Tests of Significance Considered as Evidence." *Journal of the American Statistical Association* 37 (219): 325–35.
- Dunn Jr, John E. 1962. "The Use of Incidence and Prevalence in the Study of Disease Development in a Population." *American Journal of Public Health* 52 (7): 1107–18.
- Freeman, Jonathan, and George B Hutchison. 1980. "Prevalence, Incidence and Duration." *American Journal of Epidemiology* 112 (5): 707–23.
- Keiding, Niels. 1991. "Age-Specific Incidence and Prevalence: A Statistical Perspective." *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 154 (3): 371–96.
- Laplace, Pierre Simon. 1820. *Théorie Analytique Des Probabilités*. Vol. 7. Courcier.
- MacMahon, Brian, and William D Terry. 1958. "Application of Cohort Analysis to the Study of Time Trends in Neoplastic Disease." *Journal of Chronic Diseases* 7 (1): 24–35.
- Morabia, Alfredo. 2004. "Epidemiology: An Epistemological Perspective." In *A History of Epidemiologic Methods and Concepts*, edited by Alfredo Morabia, 3–125. Springer.
- Preston, Samuel H. 1987. "Relations Among Standard Epidemiologic Measures in a Population." *American Journal of Epidemiology* 126 (2): 336–45.
- Snow, John. 1855. *On the Mode of Communication of Cholera*. Second edition. John Churchill. <https://wellcomecollection.org/works/uqa27qrt>.
- Tukey, John W. 1962. "The Future of Data Analysis." *The Annals of Mathematical Statistics* 33 (1): 1–67.