

Wasserman: All of Statistics

Kenny Workman

April 2, 2025

1 Probability

1.1 Basics

Definition 1.1. The **sample space** Ω is the set of outcomes from an experiment. Each point is denoted ω and subsets, eg. $A \subset \Omega$ are called **events**.

Definition 1.2 (Axioms of Probability). A function $\mathbb{P} : \Omega \rightarrow \mathbb{R}$ that assigns a real number to each event $A \subset \Omega$ is called a **probability function** or **probability measure** if it satisfies these three axioms:

1. **Non-negativity.** $\mathbb{P}(A) \geq 0$ for every event A
2. **Normalization.** $\mathbb{P}(\Omega) = 1$.
3. **Additivity.** $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B)$ if $A \cap B = \emptyset$.

It is incredible, and not obvious, that much of probability is built up from these only these three axioms

Example 1.3. It's actually tricky to show $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$ with these three facts:

$$\begin{aligned}\mathbb{P}(A \cup B) &= \mathbb{P}(AB^c \cap AB \cap A^cB) \\ &= \mathbb{P}(AB^c) + \mathbb{P}(AB) + \mathbb{P}(A^cB) \\ &= \mathbb{P}(AB^c) + \mathbb{P}(AB) + \mathbb{P}(A^cB) + \mathbb{P}(AB) - \mathbb{P}(AB) \\ &= \mathbb{P}(AB^c \cup AB) + \mathbb{P}(A^cB \cup AB) - \mathbb{P}(AB) \\ &= \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(AB)\end{aligned}$$

Another simple idea is that events that are identical at the limit should have identical probabilities.

Theorem 1.4 (Continuity of Events). *If $A_n \rightarrow A$ then $\mathbb{P}(A_n) \rightarrow \mathbb{P}(A)$.*

Proof. Let A_n be monotone increasing: $A_1 \subset A_2 \subset \dots$. Let $A = \lim_{n \rightarrow \infty} A_n = \bigcup_{i=1}^{\infty} A_i$.

Construct disjoint sets B_i from each A_i where $B_1 = A_1$ and $B_n = \{\omega \in \Omega : \omega \in A_n, \omega \notin \bigcup_{i=1}^{n-1} A_i\}$. It will be shown that (1) each pair of B_i are disjoint, (2) $\bigcup_{i=1}^n A_i = \bigcup_{i=1}^n B_i$ and (3) $A = \bigcup_{i=1}^{\infty} A_i = \bigcup_{i=1}^{\infty} B_i$ (Exercise 1.1).

From Axiom 3: $\mathbb{P}(A_n) = \mathbb{P}\left(\bigcup_{i=1}^n A_i\right) = \mathbb{P}\left(\bigcup_{i=1}^n B_i\right) = \sum_{i=1}^n \mathbb{P}(B_i)$.

Then $\lim_{n \rightarrow \infty} \mathbb{P}(A_n) = \lim_{n \rightarrow \infty} \sum_{i=1}^n \mathbb{P}(B_i) = \sum_{i=1}^{\infty} \mathbb{P}(B_i) = \mathbb{P}\left(\bigcup_{i=1}^{\infty} B_i\right) = \mathbb{P}(A)$

□

Definition 1.5 (Conditional Probability). If $\mathbb{P}(B) > 0$, then the probability of A given B is

$$\mathbb{P}(A | B) = \frac{\mathbb{P}(AB)}{\mathbb{P}(B)}.$$

Theorem 1.6 (Total Probability). *If $A_1 \dots A_k$ partition Ω , $\mathbb{P}(B) = \sum_{i=1}^k \mathbb{P}(B | A_i) \mathbb{P}(A_i)$*

Note. It can be difficult to assign a probability to every subset of Ω . In practice, we only assign values to select subsets described by a **sigma algebra** denoted A . This is a subset algebra with three properties:

- Non empty. $\emptyset \in A$. "Measure of the impossible."
- Closed over unions. $A_1, A_2, \dots \in A \implies \bigcup_i A_i \in A$.
- Closed over complements. $A \in A \implies A^c \in A$.

Every set in A is considered **measurable** (by its membership in the sigma algebra). A along with Ω comprises a **measurable space**, denoted by the pair (A, Ω) . If the measure on A is a probability function, importantly $\mathbb{P}(\Omega) = 1$, this space is also a **probability space**, denoted by the triple (Ω, A, \mathbb{P}) .

When Ω is the real line, the measure is often the Lebesgue measure, assigning intuitive values of "set length", eg. $[a, b] \mapsto b - a$.

Why is this important? While overly pedantic at first glance, this is the structure that explains why continuous density functions (next section) have nonzero probabilities when integrated over intervals but assign 0 probability to single points. The continuous measure, eg. Lebesgue measure, defined on the underlying probability space assigns positive values to sets and 0 to single points.

Exercise 1.1. Fill in the details for Theorem 1.2 and extend to the case where A_n is monotone decreasing.

Proof. For any pair B_{n+1} and B_n , because $B_n \subset A_n$ and $B_{n+1} \cap A_n = \emptyset$, it follows that $B_{n+1} \cap B_n = \emptyset$.

Let $\bigcup_{i=1}^n B_i = \bigcup_{i=1}^{n+1} A_i$. Then $\bigcup_{i=1}^{n+1} B_i = (A_{n+1} \setminus \bigcup_{i=1}^n A_i) \cup (\bigcup_{i=1}^n A_i) = \bigcup_{i=1}^{n+1} A_i$.

For the monotone decreasing case, let A_n be a sequence where $A_1 \supset A_2 \supset A_3 \dots$

Observe $A_1^c \subset A_2^c \dots$ and $\lim_{n \rightarrow \infty} A_n = \Omega \setminus \bigcup^{\infty} A_i^c$. Construct disjoint B_n^c from A^c in the same way.

Then $\lim_{n \rightarrow \infty} \mathbb{P}(A_n) = 1 - \sum_{i=1}^{\infty} \mathbb{P}(B_i^c) = 1 - \mathbb{P}(A^c) = \mathbb{P}(A)$ □

Exercise 1.3. Let Ω be a sample space and A_1, A_2, \dots be events. Define $B_n = \bigcup_{i=n}^{\infty} A_i$ and $C_n = \bigcap_{i=n}^{\infty} A_i$.

(a) Show $B_1 \supset B_2 \supset B_3 \dots$ and $C_1 \subset C_2 \subset C_3 \dots$

(b) Show $\omega \in \bigcap_{n=1}^{\infty} B_n$ iff ω is in an infinite number of the events

(c) Show $\omega \in \bigcup_{n=1}^{\infty} C_n$ iff ω belongs to all of the events, except possibly a finite number of those events.

Proof. (a) Certainly $\bigcup_{i=1}^{\infty} A_i \supset \bigcup_{i=2}^{\infty} A_i \dots$ and $\bigcap_{i=1}^{\infty} A_i \subset \bigcap_{i=2}^{\infty} A_i \dots$

(b) Forward. Assume $\omega \in \bigcap_{n=1}^{\infty} B_n$. If ω does not belong to an infinite number of events A_i , there exists some index j past which $\omega \notin B_j$. Then certainly $\omega \notin \bigcap_{n=1}^{\infty} B_n$. Reverse. ω belonging to infinite events means there cannot exist such a j described previously so $\omega \in B_n$ for all n . Indeed $\omega \in \bigcap_{n=1}^{\infty} B_n$

(c) Forward. Assume $\omega \in \bigcup_{n=1}^{\infty} C_n$. Then $\omega \in C_j = \bigcap_{i=j}^{\infty} A_i$ for some j . This is another way of saying ω is in every single event except for perhaps a finite number in $A_{i < j}$. Reverse. Let j be the index of the largest event that ω is not in. Then $\omega \in C_{n > j}$ and certainly $\omega \in \bigcup^{\infty} C_n$. □

Note. The key idea above is this notion of "infinitely often" (i.o.) and "all but finitely often" (eventually) which are two distinct structures of infinite occurrence in sequences. Consider an ω that exists in every other event (eg. just the odd indices) for infinite events and revisit its inclusion in $\bigcap^{\infty} B_i$ and $\bigcup^{\infty} C_i$.

Note. $\lim \bigcap A_n$ is also referred to as the limit infimum of A_n . Similarly, $\lim \bigcup A_n$ is referred to as the limit supremum of A_n .

Exercise 1.7. Let $\mathbb{P}(\bigcup A_i) \leq \sum_{i=1}^n \mathbb{P}(A_i)$. Then $\mathbb{P}(A_{n+1} \cup (\bigcup_{i=1}^n A_i)) \leq \mathbb{P}(A_{n+1}) + (\sum_{i=1}^n \mathbb{P}(A_i)) - \mathbb{P}(A_{n+1} \cap (\bigcup_{i=1}^n A_i)) \leq \sum_{i=1}^{n+1} \mathbb{P}(A_i)$

Note. Expand a bit on the Boole inequality.

Exercise 1.9. For fixed B s.t. $\mathbb{P}(B) > 0$, show $\mathbb{P}(\cdot | B)$ satisfies the three axioms of probability.

Proof. • **Non-negativity.** If $\mathbb{P}(B) > 0$ and $\mathbb{P}(AB) > 0$ for any $A \subset \Omega$, certainly $\frac{\mathbb{P}(AB)}{\mathbb{P}(B)} > 0$.

- **Normalization.** $\frac{\mathbb{P}(\Omega \cap B)}{\mathbb{P}(B)} = \frac{\mathbb{P}(B)}{\mathbb{P}(B)} = 1$

- **Additivity.** Let $AB \cap CB = \emptyset$, then $\mathbb{P}(AB \cap CB) = \mathbb{P}(AB) + \mathbb{P}(CB)$. Indeed $\frac{\mathbb{P}(AB \cap CB)}{B} = \frac{\mathbb{P}(AB)}{B} + \frac{\mathbb{P}(CB)}{B}$

□

Exercise 1.11. Suppose A and B are independent events. Show that A^c and B^c are also independent.

Proof. We are given $\mathbb{P}(AB) = \mathbb{P}(A)\mathbb{P}(B)$. Then $\mathbb{P}(A^c)\mathbb{P}(B^c) = (1 - \mathbb{P}(A))(1 - \mathbb{P}(B)) = 1 - \mathbb{P}(A) - \mathbb{P}(B) + \mathbb{P}(A)\mathbb{P}(B) = 1 - \mathbb{P}(A \cup B) = \mathbb{P}(A^c B^c)$. The second to last equality uses independence of $P(AB)$. The last equality uses the property of set complements $P(A \cup B) = P(A^c \cap B^c)$. □

Exercise 1.13. Suppose a fair coin is tossed repeatedly until heads and tails is each encountered exactly once. Describe Ω and compute the probability exactly three tosses are needed.

Proof. • The sample space is the set of binary strings with exactly one 0 and 1. For strings of length greater than 2, these are repeated strings of 0 or 1 capped with a 1 or 0 respectively.

- By independence, each n -string has an identical probability $\frac{1}{2^n}$. There are two such 3-strings: 001 and 110. Using additivity, $\mathbb{P}(3 \text{ tosses}) = \frac{1}{8} + \frac{1}{8} = \frac{1}{4}$

□

Exercise 1.15. The probability a child has blue eyes is $\frac{1}{4}$. Assume independence between children. Consider a family with 3 children.

- If it is known that at least one of the children have blue eyes, what is the probability that at least two of the children have blue eyes?
- If it is known that the youngest child has blue eyes, what is the probability that at least two of the children have blue eyes?

Proof. • Straightforward conditional probability. Let A be the event where at least one child has blue eyes and B be the event where at least two children have blue eyes. Consider first, $\mathbb{P}(A) = 1 - \mathbb{P}(\text{no child has blue eyes}) = 1 - \frac{27}{64} = \frac{37}{64}$. Compute $\mathbb{P}(A \cap B)$ by enumerating events 101, 111, 110 and using additivity: $2 \cdot \frac{1}{4}^2 \cdot \frac{3}{4} + \frac{1}{4}^3 = \frac{10}{64}$. Then $\mathbb{P}(B | A) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(A)} = \frac{10}{64} \cdot \frac{64}{37} = \frac{10}{37}$

- Similar procedure. Let A be the event where the youngest child has blue eyes and B be as before. Using independence, $\mathbb{P}(A) = \frac{1}{4}$. (To see this rigorously, enumerate the sample space and see $\mathbb{P}(\Omega | \text{first child blue}) = 1$). Now $\mathbb{P}(B \cap A)$ describe events 110, 101, 111 only. $\frac{7}{64} \cdot \frac{4}{1} = \frac{7}{16}$.

□

Exercise 1.17. Show $\mathbb{P}(ABC) = \mathbb{P}(A | BC)\mathbb{P}(B | C)\mathbb{P}(C)$

Proof. By straightforward application of the definition of conditional probability: $\frac{\mathbb{P}(ABC)}{\mathbb{P}(BC)} \frac{\mathbb{P}(BC)}{\mathbb{P}(C)} \mathbb{P}(C) = \mathbb{P}(ABC)$ □

Exercise 1.19. Suppose 50% of computer users are Windows. 30% are Mac. 20% are Linux. Suppose 65% of Mac users, 82% of Windows users and 50% of Linux users get a virus. We select a person at random and learn they have the virus. What is the probability they are a Windows user?

Proof. Let each $\omega \in \Omega$ be a distinct user. Then $W, M, L \subset \Omega$ are the users with Windows, Mac + Linux machines. $V, N \subset \Omega$ are the users with and without viruses.

We want $\mathbb{P}(W | V) = \frac{\mathbb{P}(V | W)\mathbb{P}(W)}{\mathbb{P}(V)}$. Compute $\mathbb{P}(V) = \sum_{X \in \{W, M, L\}} \mathbb{P}(V | X)\mathbb{P}(X) = 0.705$. Then $\mathbb{P}(W | V) = \frac{0.82 \cdot 0.50}{0.705} = 0.581$. □

Exercise 1.20. A box contains 5 coins, each with a different probability of heads: 0, 0.25, 0.5, 0.75, 1. Let C_i be the event with coin i and H_i be the event that heads is recovered on toss i . Suppose you select a coin at random and flip it.

- What is the posterior probability $\mathbb{P}(C_i | H_1)$ for each coin?
- What is $\mathbb{P}(H_2 | H_1)$?
- Let B_i be the event that the first heads is recovered on flip i . What is $\mathbb{P}(C_i | B_i)$ for each coin?

Proof. • $\mathbb{P}(H_1) = \frac{1}{2}$. For each coin, $\mathbb{P}(C_i | H) = \frac{\mathbb{P}(H | C_i)\mathbb{P}(C_i)}{\mathbb{P}(H)}$. $\mathbb{P}(H)$ can be worked out using total probability: $\sum_i \mathbb{P}(H|C_i)\mathbb{P}(C_i) = \frac{1}{2}$. Then eg. the posterior $\mathbb{P}(C_4 | H) = \frac{3}{4} \cdot \frac{1}{5} \cdot \frac{2}{1} = \frac{3}{10}$.

- Note that both tosses are conditionally independent: $\mathbb{P}(H_2 H_1 | C_i) = \mathbb{P}(H_2 | C_i)\mathbb{P}(H_1 | C_i)$.
 $\mathbb{P}(H_2 | H_1) = \frac{\mathbb{P}(H_2 H_1)}{\mathbb{P}(H_1)} = \frac{\sum_i \mathbb{P}(H_2 H_1 | C_i)\mathbb{P}(C_i)}{\sum_i \mathbb{P}(H_1 | C_i)\mathbb{P}(C_i)}$. Because $\mathbb{P}(C_i)$ is uniform, we can simply to
 $\frac{\sum_i \mathbb{P}(H_2 H_1 | C_i)}{\sum_i \mathbb{P}(H_1 | C_i)}$. The result is $\frac{\sum_i p_i^2}{\sum_i p_i}$.
- Similar idea to (a).

□

Note. Important to see that independent events are not conditionally independent in general. Try to construct an example.

1.2 Random Variables

– (kenny) TODO fix counters

Definition 1.7 (Random Variable). A random variable X is a function mapping the sample space to real numbers: $X : \Omega \rightarrow \mathbb{R}$.

It is important to think of the relationship between the random variable and its underlying sample space when computing probabilities: eg. $\mathbb{P}(X = x) = \mathbb{P}(X^{-1}(x))$ and $\mathbb{P}(X \in A) = \mathbb{P}(X^{-1}(A))$.

Definition 1.8 (Cumulative Distribution Function). The CDF is the function $F_X : \mathbb{R} \rightarrow [0, 1]$ where $F_X(x) = \mathbb{P}(X \leq x)$. Equivalently $F_X(x) = \mathbb{P}(X^{-1}((-\infty, x]))$.

The CDF contains "all the information" in a random variable. This is articulated by the following theorem:

Theorem 1.9. For random variables X and Y with CDFs F and G , if $F(x) = G(X) \forall x \in [0, 1]$, then $X = Y$ ($\mathbb{P}(X \in A) = \mathbb{P}(Y \in A)$ for each $A \subset \mathbb{R}$).

And the behavior of the CDF, including "all of its information" is uniquely determined by just three properties:

Theorem 1.10. A function $F : \mathbb{R} \rightarrow [0, 1]$ is a CDF iff it satisfies three properties:

- Non-decreasing. $x_2 > x_1 \implies F(x_2) \geq F(x_1)$
- Normalization. $\lim_{y \rightarrow -\infty} F(y) = 0$ and $\lim_{y \rightarrow \infty} F(y) = 1$
- Right-continuous. For any $x \in \mathbb{R}$, $F(x) = F^+(x)$ where $F^+(x) = \lim_{y \rightarrow x, y > x} F(y)$

Proof. Starting with (iii) from the text, let $A = (-\infty, x]$ and y_1, y_2, \dots be a sequence where $y_1 < y_2 \dots$ and $\lim_i y_i = x$. By the definition of the CDF, $F(y_i) = \mathbb{P}(A_i)$ and $F(x) = \mathbb{P}(A)$, where $\lim_i F(y_i)$ is equivalent to $\lim_{y \rightarrow x, y > x} F(y)$. Observe $\cap_i A_i = A$ so $\mathbb{P}(A) = \mathbb{P}(\cap_i A_i) = \lim_i \mathbb{P}(A_i) = \lim_i F(y_i) = F(x)$ as desired.

To see (ii), $\lim_{y \rightarrow -\infty} F(y) = 0$, define a sequence y_1, y_2, \dots where $y_1 > y_2 \dots$ as before and $y_1 = y$. Let $A_i = (\infty, y_i]$. Then $\cap_i A_i = \emptyset$ and $\mathbb{P}(\cap_i A_i) = \mathbb{P}(\emptyset) = 0$. Indeed $\lim_{y \rightarrow -\infty} F(y) = \lim_i \mathbb{P}(A_i) = \mathbb{P}(\cap_i A_i) = 0$. A similar argument shows the limit to the other direction.

For (iii), if $x_2 > x_1$ then $P((-\infty, x_2]) \geq P((-\infty, x_1])$ and $F(x_2) \geq F(x_1)$. □

The interesting direction is the reverse: a function satisfying these properties uniquely determines a probability function. It is difficult to show in general. A concrete example is the Cantor function ([Devil's staircase](#)) which satisfies non-decreasing, normality and right-continuous properties but from which is difficult to derive a measure that satisfies eg. countable additivity.

Note. A deeper measure theory course will approach this problem by defining the probability function on an algebra of subsets rather than on each subset directly. Refer to tools like [Caratheodory's extension theorem](#).

It is from these random variables that we build "distributions", essentially functions $\mathbb{R} \rightarrow [0, 1]$ that obey the three probability axioms.

Definition 1.11. If X "takes" countably many values (eg. has a countable range) it is **discrete**. $f_X(x) = \mathbb{P}(X = x)$ is its **probability mass function** or PMF.

Definition 1.12. X is **continuous** if it has some f_X that obeys three properties:

- $\int_{-\infty}^{\infty} f_X(x) dx = 1$
- $\forall x \in \mathbb{R} : f_X(x) \geq 0$
- $\mathbb{P}(a < X < b) = \int_a^b f_X(x) dx$

f_X is called the **probability density function** or PDF. Additionally, $F_X(x) = \int_{-\infty}^x f_X(t) dt$ and $f_X(x) = F'_X(x)$ for all points x where F_X is differentiable.

The formal relation between the density function and the sample space is a bit tricky, especially when X is continuous. In practice, we often just produce a function and deal with it directly while assuming the underlying sample space with a well defined measure is lurking around.

Note. We learned the probability function is defined on a well-defined sample space by measuring events / sets.

We will proceed with some important mass functions.

Definition 1.13 (The Point Mass Distribution). If $X \sim \sigma_a$ (reads " X has a point mass distribution at a "), $f_X(a) = 1$ while $f_X(x) = 0$ for all $x \neq a$.

$$F_X(x) = \begin{cases} 0, & x < a, \\ 1, & x \geq a \end{cases}$$

Definition 1.14 (The Uniform Distribution).

Definition 1.15 (The Bernoulli Distribution).

Definition 1.16 (The Binomial Distribution).

Definition 1.17 (The Geometric Distribution).

Definition 1.18 (The Poisson Distribution).

Exercise 2.2. Show $\mathbb{P}(X = x) = F(x^+) - F(x^-)$

Proof. The key here is to see $\lim_{z \leftarrow x, z \rightarrow x} F(z) = \mathbb{P}(X \in \cup_i (\infty, z_i]) = \mathbb{P}(X < x)$ for some sequence z_1, z_2, \dots where $\lim_i z_i = x$. While $\lim_{y \rightarrow x, y \rightarrow x} F(y) = \mathbb{P}(X \in \cap_i (\infty, y_i]) = \mathbb{P}(X \leq x)$.

Pay attention to the behavior of converging sets and the boundary. In the right-continuous case, the sequence is approaching the boundary x from above and each sequence is closed on x . Therefore in the limit, they include x .

In the left-continuous case, the sequence is approaching the boundary x from below and each sequence excludes x . Therefore in the limit, they exclude x .

To conclude $F(x^+) - F(x^-) = \mathbb{P}(X \leq x) - \mathbb{P}(X < x) = \mathbb{P}(X = x)$. Of course, if X is continuous, $F(x^+) = F(x^-)$ and $\mathbb{P}(X = x) = 0$, showing once again that every real value has no probability mass. \square