

Probability

The setup

Random “experiment”: A well-defined process with an uncertain outcome (for example, toss three fair coins)

Sample space, \mathcal{S} : The set of possible outcomes (e.g., $\{ \text{HHH, HHT, HTH, THH, HTT, THT, TTH, TTT} \}$)

Event: A set of outcomes (a subset of \mathcal{S}) (e.g., $A = \{\text{exactly one head}\} = \{ \text{HTT, THT, TTH} \}$)

An event is said to have occurred if one of the outcome it contains occurs.

Basic rules of probability

1. $1 \geq \Pr(A) \geq 0$, for any event A.
2. $\Pr(\mathcal{S}) = 1$, for the sample space \mathcal{S} .
3. If A and B are *mutually exclusive*, $\Pr(A \text{ or } B) = \Pr(A) + \Pr(B)$.

More rules

$$\Pr(\text{not } A) = 1 - \Pr(A)$$

$$\Pr(A \text{ or } B) = \Pr(A) + \Pr(B) - \Pr(A \text{ and } B)$$

Conditional probability

$\Pr(A \mid B) = \text{“probability of A given B”} = \Pr(A \text{ and } B) / \Pr(B)$, provided $\Pr(B) > 0$.

Independence

Events A and B are *independent* if $\Pr(A \text{ and } B) = \Pr(A) \Pr(B)$; equivalently, $\Pr(A \mid B) = \Pr(A)$ and $\Pr(B \mid A) = \Pr(B)$.

Still more rules

$$\Pr(A \text{ and } B) = \Pr(B) \Pr(A \mid B) = \Pr(A) \Pr(B \mid A)$$

$$\begin{aligned} \Pr(A) &= \Pr(A \text{ and } B) + \Pr(A \text{ and not } B) \\ &= \Pr(B) \Pr(A \mid B) + \Pr(\text{not } B) \Pr(A \mid \text{not } B) \end{aligned}$$

Bayes rule

$$\begin{aligned} \Pr(A \mid B) &= \Pr(A) \Pr(B \mid A) / \Pr(B) \\ &= \Pr(A) \Pr(B \mid A) / [\Pr(A) \Pr(B \mid A) + \Pr(\text{not } A) \Pr(B \mid \text{not } A)] \end{aligned}$$

Examples

1. Mendel's peas.

Mendel's peas had either purple or white flowers. Flower color was due to a single gene with the purple allele (A) dominant to the white allele (a).

The F_1 hybrid (obtained by crossing two pure-breeding lines, one with purple and one with white flowers) has purple flowers and is heterozygous, Aa .

(a) Self the F_1 , pick an F_2 seed at random and grow it up.

i. $\Pr(F_2 \text{ plant has white flowers}) = (1/2) \times (1/2) = 1/4 = 25\%$.

ii. $\Pr(F_2 \text{ plant has purple flowers}) = 1 - (1/4) = 3/4 = 75\%$.

iii. $\Pr(F_2 \text{ plant has genotype } AA) = 1/4 = 25\%$.

iv. $\Pr(F_2 \text{ plant has genotype } AA \mid \text{it has purple flowers})$
 $= \Pr(F_2 \text{ has genotype } AA \text{ and purple flowers}) / \Pr(F_2 \text{ has purple flowers})$
 $= (1/4) / (3/4) = 1/3 \approx 33\%$.

(b) Self the F_2 . Grow up ten of the F_3 seeds at random. To save writing, define the following events:

$$P_2 = \{ F_2 \text{ has purple flowers} \}$$

$$H_o = \{ F_2 \text{ has genotype } AA \}$$

$$H_e = \{ F_2 \text{ has genotype } Aa \}$$

$$P_3 = \{ \text{all ten } F_3 \text{ plants have purple flowers} \}$$

Note that $P_2 = H_e$ or H_o . You may want to write out the following, using words. The use of this notation makes things compact but also somewhat harder to follow.

i. $\Pr(P_2) = 1 - (1/4) = 3/4 = 75\%$.

ii. $\Pr(P_3 \mid H_o) = 100\%$.

iii. $\Pr(P_3 \mid H_e) = (3/4)^{10} \approx 5.6\%$.

iv. $\Pr(P_3 \mid P_2) = \Pr(P_3 \text{ and } P_2) / \Pr(P_2)$
 $= \{ \Pr(P_3 \text{ and } H_o) + \Pr(P_3 \text{ and } H_e) \} / \Pr(P_2)$
 $= \{ \Pr(H_o) \times \Pr(P_3 \mid H_o) + \Pr(H_e) \times \Pr(P_3 \mid H_e) \} / \Pr(P_2)$
 $= \{ (1/4) \times 1 + (1/2) \times (3/4)^{10} \} / (3/4) \approx 37\%$

Alternatively, $\Pr(P_3 \mid P_2) = \Pr(P_3 \text{ and } H_o \mid P_2) + \Pr(P_3 \text{ and } H_e \mid P_2)$
 $= \Pr(H_o \mid P_2) \times \Pr(P_3 \mid H_o \text{ and } P_2) + \Pr(H_e \mid P_2) \times \Pr(P_3 \mid H_e \text{ and } P_2)$
 $= (1/3) \times 1 + (2/3) \times (3/4)^{10} \approx 37\%$

v. $\Pr(H_o \mid P_3, P_2) = \Pr(H_o \text{ and } P_3 \mid P_2) / \Pr(P_3 \mid P_2)$
 $= \Pr(P_3 \mid H_o \text{ and } P_2) \times \Pr(H_o \mid P_2) / \Pr(P_3 \mid P_2)$
 $= 1 \times (1/3) / [1 \times (1/3) + (3/4)^{10} \times (2/3)]$
 $\approx 90\%$.

2. Suppose a test for HIV correctly gives a positive result, if a person is infected, with probability 99.5%, and correctly gives a negative result, if a person is not infected, with probability 98%.

- (a) Suppose that 0.1% of a population are infected with HIV. Consider drawing a person at random and testing him or her for HIV infection. Calculate $\Pr(\text{infected} \mid \text{test is positive})$.

Let $I = \{ \text{the person is infected} \}$ and $P = \{ \text{the person tests positive} \}$.

From the information above, on the *sensitivity* and *specificity* of the test, we have $\Pr(P \mid I) = 0.995$ and $\Pr(\text{not } P \mid \text{not } I) = 0.98$. From this information, $\Pr(\text{not } P \mid I) = 1 - \Pr(P \mid I) = 0.005$ and $\Pr(P \mid \text{not } I) = 1 - \Pr(\text{not } P \mid \text{not } I) = 0.02$. Note further that $\Pr(I) = 0.001$, and so $\Pr(\text{not } I) = 0.999$.

We seek to calculate $\Pr(I \mid P)$. We use Bayes's rule. (Why use Bayes's rule here? because we want to "turn around the conditioning." We want to write $\Pr(I \mid P)$ in terms of things like $\Pr(P \mid I)$.)

$$\begin{aligned}\Pr(I \mid P) &= \Pr(I) \Pr(P \mid I) / [\Pr(I) \Pr(P \mid I) + \Pr(\text{not } I) \Pr(P \mid \text{not } I)] \\ &= 0.001 \times 0.995 / (0.001 \times 0.995 + 0.999 \times 0.02) \approx 4.7\%\end{aligned}$$

- (b) Consider a person drawn from a high-risk group, so that they have, *a priori*, probability 30% of being infected. Calculate $\Pr(\text{infected} \mid \text{test is positive})$.

In this case, we change $\Pr(I) = 0.3$ and $\Pr(\text{not } I) = 1 - \Pr(I) = 0.7$.

Thus, $\Pr(I \mid P) = 0.3 \times 0.995 / (0.3 \times 0.995 + 0.7 \times 0.02) \approx 96\%$

3. The “Monty Hall” problem.

You’re on a game show, and are presented with three doors. One hides a car; the other two hide goats. You’re allowed to choose a door. Monty Hall then opens one of the other two doors to reveal a goat. You are now allowed to either stick with the door you chose initially, or switch to the other closed door. What should you do (assuming that you are hoping to get the car, and not a goat): stick with your original choice, or switch?

The answer depends on Monty’s behavior.

- (a) **If you choose the door with the car, Monty opens one of the other two doors at random. If you choose a door hiding a goat, Monty opens the other door with a goat.**
- (b) **Monty Hall opens one of the other two doors at random. If he had revealed the car, you would have lost, but you happen to be in the situation where he revealed a goat.**

Let us call the doors A, B and C, and let us assume (“without loss of generality,” as mathematicians like to say) that you choose door A. Moreover, let us initially *not* condition on Monty revealing the goat.

There are six possible outcomes of the experiment, $\{ C\boxed{G}G, CG\boxed{G}, G\boxed{C}G, GC\boxed{G}, G\boxed{G}C, GG\boxed{C} \}$. (The three letters denote the objects behind the three doors, so that CGG means the car is behind door A. The boxes indicate which door Monty opens. Since you chose door A, he won’t be opening that door, but one of the other two.)

In scenario (a), Monty always opens a door with a goat. Thus, if the door you chose hid the car (that is, if the car was behind door A), he’ll choose at random between the doors B and C. On the other hand, if the car is behind door B, he’ll always open door C, and vice versa. Thus, under scenario A, the events $\{ C\boxed{G}G \}$ and $\{ CG\boxed{G} \}$ each have probability 1/6; events $\{ GC\boxed{G} \}$ and $\{ G\boxed{G}C \}$ each have probability 1/3, and events $\{ G\boxed{C}G \}$ and $\{ GG\boxed{C} \}$ each have probability 0.

In scenario (b), in which Monty chooses at random between the two remaining doors, the six outcomes listed above each have probability 1/6.

Now, define the events $A = \{ \text{the car is behind door A} \} = \{ C\boxed{G}G, CG\boxed{G} \}$, $B = \{ \text{the car is behind door B} \} = \{ G\boxed{C}G, GC\boxed{G} \}$, $C = \{ \text{the car is behind door C} \} = \{ G\boxed{G}C, GG\boxed{C} \}$, and $D = \{ \text{Monty opens a door with a goat} \} = \{ C\boxed{G}G, CG\boxed{G}, GC\boxed{G}, G\boxed{G}C \}$.

Our objective: we want to know the chance that the door we originally chose hides the car, *given that Monty has opened a door revealing a goat*. This is $\Pr(A | D)$. If this probability is $< 1/2$, we should switch doors; if it is $> 1/2$, we should stick with our first choice; if it is $= 1/2$, it doesn’t matter whether we switch or stick.

The event “A and D” (the intersection between A and D) is simply event A, and $\Pr(A) = 1/3$ in either case (a) or (b), so $\Pr(A \text{ and } D) = 1/3$.

In case (a), $\Pr(D) = 1$, and so $\Pr(A | D) = 1/3$. We should switch doors.

In case(b), $\Pr(D) = 2/3$, and so $\Pr(A | D) = 1/2$. It doesn’t matter what we do.

The appropriate action to take depends on what Monty is doing!

4. Mendel, revisited.

Mendel's peas had either purple or white flowers; flower color is due to a single gene, for which the purple allele (A) is dominant to the white allele (a).

We cross two pure-breeding lines (one purple and one white) to produce the F_1 hybrid. We self the F_1 and choose an F_2 seed at random. We grow and self the F_2 and choose two F_3 seeds at random.

Consider the following events.

$$\begin{aligned} P &= \{F_2 \text{ has purple flowers}\} \\ O &= \{F_2 \text{ is homozygous } AA\} \\ E &= \{F_2 \text{ is heterozygous } Aa\} \\ W &= \{F_2 \text{ has white flowers}\} \\ A_1 &= \{F_3 \text{ number 1 has purple flowers}\} \\ A_2 &= \{F_3 \text{ number 2 has purple flowers}\} \end{aligned}$$

Are A_1 and A_2 independent?

A_1 and A_2 are *not* independent! A_1 and A_2 are *conditionally independent*, given the genotype of the F_2 plant. But, as we will see, in the absence of information about the F_2 parent's genotype, the flower color of the F_3 plants are not independent.

First, note that the events E , O and W are *mutually exclusive*, and that $[\Pr(E) + \Pr(O) + \Pr(W)] = 1$.

$$\begin{aligned} \Pr(A_1) &= \Pr(A_2) = \Pr(A_1 \mid E) \Pr(E) + \Pr(A_1 \mid O) \Pr(O) + \Pr(A_1 \mid W) \Pr(W) \\ &= (3/4) \times (1/2) + (1) \times (1/4) + 0 \times (1/4) = 5/8 \approx 63\%. \end{aligned}$$

$$\begin{aligned} \Pr(A_1 \text{ and } A_2) &= \Pr(A_1 \text{ and } A_2 \mid E) \Pr(E) + \Pr(A_1 \text{ and } A_2 \mid O) \Pr(O) \\ &\quad + \Pr(A_1 \text{ and } A_2 \mid W) \Pr(W) \\ &= (3/4)^2 \times (1/2) + 1 \times (1/4) + 0 \times (1/2) = 17/32 \approx 53\%. \end{aligned}$$

Thus $\Pr(A_2 \mid A_1) = \Pr(A_1 \text{ and } A_2) / \Pr(A_1) = (17/32) / (5/8) = 17/20 = 85\%$, which is considerably greater than $\Pr(A_2)$. And so, A_2 and A_1 are not independent.