

1

Probability

2

Zvi M. Kedem

3

2014-11-09

4

Contents

5

1 Probability on finite sets 2

6

1.1 Introduction 2

7

1.2 Basic concepts 2

8

1.3 Conditional probability and Bayes' theorem 8

9

1.4 Practicing Bayesian thinking 10

10

1.5 Chain rule 14

11

1.6 Joint probability distribution and marginalization 15

12

2 Entropy 16

13

2.1 Example 16

14

2.2 A derivation 17

15

2.3 The case of general probabilities 18

16

2.4 Randomness and Entropy 19

17

2.5 Additional interpretations of entropy 19

18

2.6 Entropy obtained by randomly choosing from a set of symbols 20

19

2.7 Examples for C_0 , C_1 , and D_1 21

20

2.8 Getting random numbers 23

21

3 Bayesian inference 23

22

3.1 Hypothesis and evidence 23

23

3.2 Two coins example 23

24

4 Using Bayesian inference for incrementally increasing the confidence in the hypothesis that the coin tossed was Bad 25

25

4.1 The distribution of Good vs. Bad is known 25

26

4.2 The distribution of Good vs. Bad is not known 27

27

4.3 Recapitulation 29

28

29

5 Bayesian classification 30

30

5.1 Basic Bayesian classification 30

31

5.2 Naïve Bayesian classification 32

31

5.3 Laplacian correction 32

32

1 Probability on finite sets

1.1 Introduction

Discussion will be informal. Some of what we do may not be correct for infinite sets, which include the important case of the real line \mathbb{R} .

Question 1. We have a coin which with probability 50% comes up H (heads). What does this mean?

Question 2. The weather channel <http://www.weather.com/weather/today/10012> says that the probability of rain tonight is 50%. What does this mean?

As good examples, we will have several objects:

1. A fair coin C_0 , which with probability 0.5 each gives us H (heads) or T (tails).

Formally H and T are called *obverse* and *reverse* http://en.wikipedia.org/wiki/Obverse_and_reverse.

2. An unfair coin C_1 , which with probability 0.99 gives us H and with probability 0.01 gives us T

3. An extremely unfair coin C_2 , which with probability 1 gives us H and with probability 0 gives us T. How this is done is not important for us.

4. A fair die D_0 , which with probability $1/6$ gives us each of $1, \dots, 6$.

Observation 1. An event of probability 0 can still happen—it is just that “the probability of this is smaller than any positive number,” but if we assume informally that it never happens we will be OK, because we are dealing with finite sets.

But once you go to infinite sets you need to be much more careful. For example, if you pick a real number uniformly (each number with the same probability) from the interval $[0, 1]$ then of course each time you do it, you will get some real number.

But the probability of picking any specific number must be 0. Here is a simple argument, why it cannot be greater than 0. Let the probability of picking any specific number be p . As all numbers are picked with the same probability the sum of the probabilities over all the numbers, $\sum_{x \in [0,1]} p(x)$ must be infinity as there is an infinite number of real numbers in $[0, 1]$, which of course is impossible as this sum should be 1.

1.2 Basic concepts

We will assume that we have a (finite) set X and for each $x \in X$, we have a number p_x , such that $0 \leq p_x \leq 1$ and $\sum_{x \in X} p_x = 1$. p_x is the probability that if we pick “randomly” from X , we will get x . \mathcal{X} is X together with $\{p_x \mid x \in X\}$. This is called a *random variable*. We will also write $\Pr[x]$ or $\Pr[\mathcal{X} = x]$ for p_x . When there is no confusion, we may write X instead of \mathcal{X} . Formally

Definition 1. Let X be a finite set and let $p(x)$, for $x \in X$ satisfy

1. $p(x) \in [0, 1]$ for all $x \in X$

2. $\sum_{x \in X} p(x) = 1$

Then X together with p is a random variable \mathcal{X} . We may write p_x or $\Pr[x]$ instead of $p(x)$ and X instead of \mathcal{X} and similar obvious notational simplifications.

68 If we want to use pictures, then we could draw a square of area 1, and then for each x , p_x will be the area of
 69 an appropriate “part” of the square. See Fig. 1, Fig. 2, Fig. 3, Fig. 4.

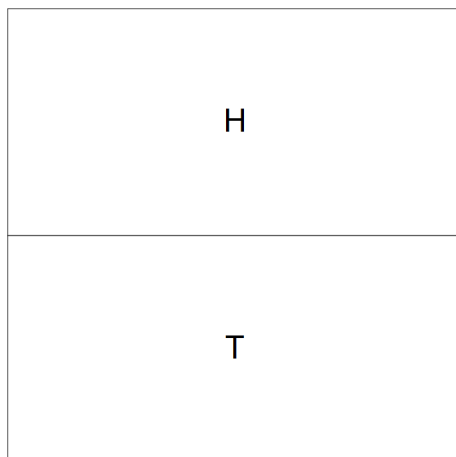


Figure 1: C_0 .

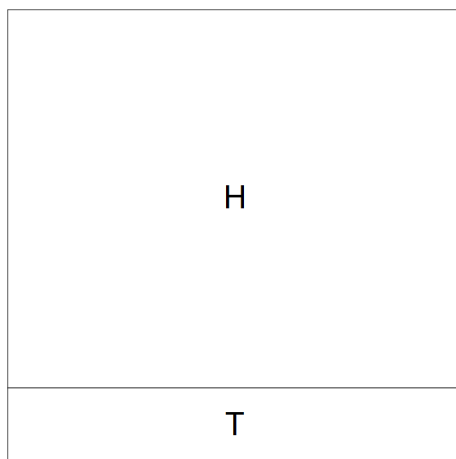


Figure 2: C_1 . Imagine that the lower stripe is only 1/100th the height of the total height.

70 Sometimes a real-valued function, say f , is defined on X . Then the *expectation* of f is $\text{Ex}[f] =$
 71 $\sum_{x \in X} \text{Pr}[x] \cdot f(x)$. Formally, just to repeat
 72 **Definition 2.** Given X and p as above and a real-valued function f on X , the expectation of f is

$$\text{Ex}[f] = \sum_{x \in X} \text{Pr}[x] \cdot f(x).$$

73

74 Expectation essentially means “average.” It does not mean that this is value that is expected to happen. For
 75 example the expectation of the number of children per woman in the US is about 2.2.

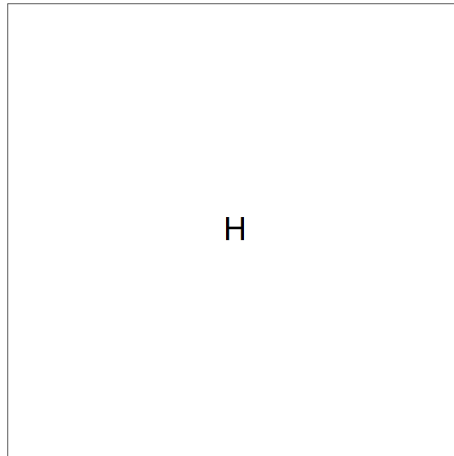


Figure 3: C_2 .

1	2	3
4	5	6

Figure 4: D_0 .

⁷⁶ **Example 1.** We toss C_1 and we get some payoff base on the result: if we toss x we get $f(x)$, which could be
⁷⁷ positive, zero, or negative. Say $f(H) = 2$ and $f(T) = -3$, then the expectation (how much money we make
⁷⁸ on the average per toss if we get paid f) is $0.99 \cdot 2 + 0.01 \cdot (-3) = 1.95$.

⁷⁹ **Note 1.** Sometimes, the value of expectation is quite “unexpected,” though we do not need to worry
⁸⁰ about this in our setting. The first example seems to have been provided by Nicolas Bernoulli [http:](http://en.wikipedia.org/wiki/Nicolas_Bernoulli)
⁸¹ [//en.wikipedia.org/wiki/Nicolas_Bernoulli](http://en.wikipedia.org/wiki/Nicolas_Bernoulli) and later called St. Petersburg paradox [http:](http://en.wikipedia.org/wiki/St._Petersburg_paradox)
⁸² [//en.wikipedia.org/wiki/St._Petersburg_paradox](http://en.wikipedia.org/wiki/St._Petersburg_paradox).

⁸³ Consider the following one-sided game in which Alice can only win but cannot lose. Alice tosses a fair coin
⁸⁴ until she gets the first T (tail). If the first T comes up in the first toss, the game ends and she gets \$1. If the
⁸⁵ first T comes up in the second toss, the game ends and she gets \$2. If the first T comes up in the third toss,
⁸⁶ the game ends and she gets \$4. In general, if the first T comes up in the i th toss, the game ends and she gets

87 $\$2^{i-1}$. What is the expected win?

88 In the general i th case, there was a sequence of $i - 1$ H's followed by a T. The probability of this sequence is
89 exactly 2^{-i} and she gets $\$2^{i-1}$. The expectation is $\sum_{i=1}^{\infty} 2^{-i} \cdot 2^{i-1} = \sum_{i=1}^{\infty} 2^{-1} = \sum_{i=1}^{\infty} 1/2 = \infty$.

90 There is nothing wrong with this derivation, the expectation does not have to be finite. The paradox arises
91 more from psychology/utility theory. Would you be willing to give all the money you have to play this game?
92 Would Bill Gates give all his money to play this game?

93 **Definition 3.** A subset of X is also called an event.

94 There is nothing interesting in referring to subsets by different names (events). But in this context we want
95 to talk about the probability that an event happened, which really means that some element in the event
96 happened.

97 **Example 2.** For C_0 , $E = \{H\}$ is an event, but in this case we may as well talk about H.

98 **Example 3.** An example of a more interesting event for D_0 is $E = \{2, 5\}$.

99 **Definition 4.** Informally stating, the probability of an event E is the probability that the resulting “random”
100 chosen x is in E . More formally

101
$$\Pr[E] = \sum_{x \in E} \Pr[x].$$

102 So for our event: $\Pr[\{2, 5\}] = 1/6 + 1/6 = 1/3$. Compare with the areas in Fig. 4.

103 If we have two events A and B , then of course we can talk about events, \bar{A} , $A \cup B$ and $A \cap B$. For an example
104 of complementary events see Fig. 5, for an example of disjoint events, see Fig. 6; and for an example of not
105 disjoint events, see Fig. 7.

106 $A \cap B$ is also denoted by $\Pr[A, B]$. Of course, $\Pr[A, B] = \Pr[B, A]$.

107 $\bar{A} = X \setminus A$ (also written as $X - A$). Sometimes $\sim A$ or $\neg A$ are used to denote the complement of A .

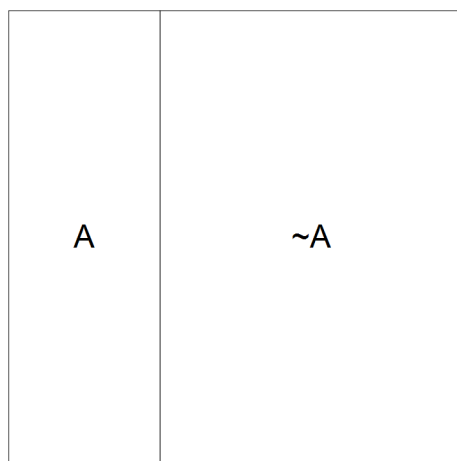


Figure 5: Complementary events.

108 Let us now toss first D_0 and then C_0 . Look at Fig. 8. Look at Fig. 9. It shows the probabilities explicitly as
109 areas.

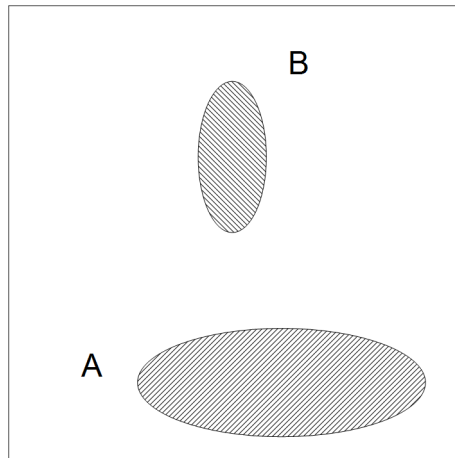


Figure 6: Disjoint events.

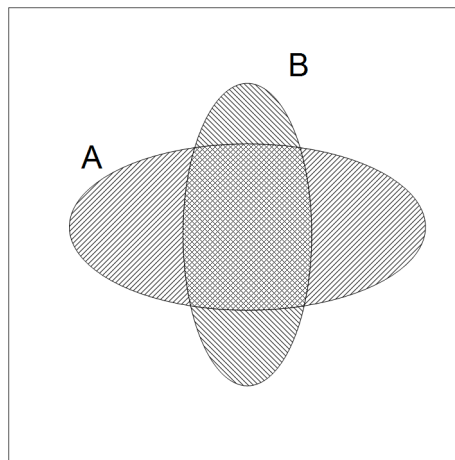


Figure 7: Not disjoint events.

110 Let us consider what is the probability of getting the pair (1, H). As we are here talking about a sequence of
 111 events *in order*, we will denote this by $\Pr[(1, H)]$, meaning we get first 1 and then H.

112 The result of tossing a die says nothing about the result of tossing the coin. Informally speaking for now, these
 113 two variables/events are *independent*. The probability of getting this result is just the product of probabilities
 114 of getting each of them separately, in this case $1/6 \cdot 1/2 = 1/12$.

115 Similarly, if we first toss D_0 and then C_1 , though the actual numbers are, of course, different.

116 We had a sequence of random variables. Let us look at a slightly more interesting one. We toss C_0 . If we get
 117 H we toss C_0 ; if we get T, we toss C_1 . What are the probabilities of getting HH, HT, etc. We just multiply
 118 probabilities, see Fig. 10. We get $\Pr[(T, H)] = 0.5 \cdot 0.99$, etc.

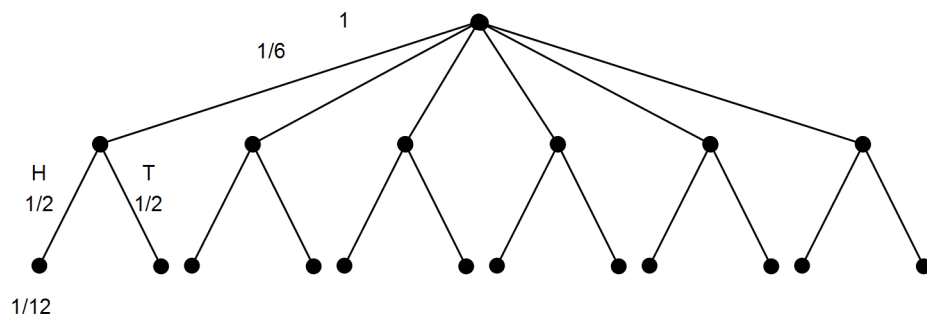


Figure 8: Tossing D_0 and then C_0 . Only some of the probabilities are written out. The probability is written below the result of the random variable.

1 H	2 H	3 H
4 H	5 H	6 H
1 T	2 T	3 T
4 T	5 T	6 T

Figure 9: Probabilities resulting from first tossing D_0 and then tossing C_0 .

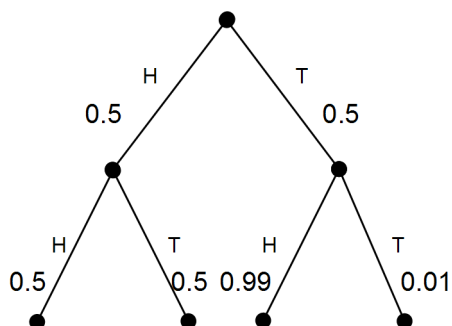


Figure 10: Tossing C_0 first and then C_0 or C_1 , depending on the result of the first toss.

1.3 Conditional probability and Bayes' theorem

Let us now talk about the very important concepts of independent variables and conditional probabilities.

Practically everything we need to know can be understood by carefully examining Fig. 7. For examples, it is good to look at the more specific figures dealing with the die and the various coins.

Assume that as the result of getting a random variable (value) (e.g., tossing a coin, or die, or both) we are in event B . What is the probability that we are also in event A ?

Definition 5. $\Pr[A | B]$ denotes the probability of event A given that event B took place.

Looking at Fig. 7 (and of course assuming that $\Pr[B] \neq 0$), this is

$$\Pr[A | B] = \frac{\Pr[A, B]}{\Pr[B]} \quad (1)$$

Note that using our pictorial representation, this is really

$$\frac{\text{area}(A \cap B)}{\text{area}(B)}.$$

Of course A and B do not have to be “not disjoint” for $\Pr[A | B]$ to make sense.

Example 4. What is $\Pr[A | B]$ for the situation depicted in Fig. 6?

$\Pr[A | B]$ denotes conditional probability: what is the probability of A given B . Recall that $\Pr[A, B]$ denotes the probability of both A and B , that is we find ourselves in $A \cap B$.

Now, let's look at some examples:

1. We toss D_0 . What is $\Pr[x > 1]$? Looking at Fig. 4, we see this is $5/6$ (5 rectangles out of 6).
2. We toss D_0 . What is $\Pr[x \text{ is even}]$? Looking at Fig. 4, we see this is $3/6$
3. We toss D_0 . What is $\Pr[x > 1 \text{ and is even}]$? Looking at Fig. 4, we see this is $3/6$
4. We toss D_0 . What is $\Pr[x \text{ is even} | x > 1]$? Looking at Fig. 4, we see this is $(3/6)/(5/6) = 3/5$
5. We toss D_0 . What is $\Pr[x > 1 | x \text{ is even}]$? Looking at Fig. 4, we see this is $(3/6)/(3/6) = 1$

Let us now return to the case of first tossing D_0 and then tossing C_0 , see Fig. 8 and Fig. 9.

The results of the toss seem independent (and they indeed are). Knowing the result of tossing the die does not tell us anything about the result of tossing the coin. Let us compute $\Pr[H | 1]$. It is $(1/12)/(2/12) = 1/2$.

In this case, $\Pr[H | 1] = \Pr[H]$. This is essentially (but not quite) the definition of the two events getting H and getting 1, being independent: the result of tossing the coin and the result of tossing the die are independent of each other.

Observation 2.

$$\Pr[A | B] + \Pr[\bar{A} | B] = 1. \quad (2)$$

Follows immediately, as given B either A or \bar{A} but not both.

Let us now prove a very important result:

149 **Theorem 1.** (Bayes) http://en.wikipedia.org/wiki/Bayes'_theorem.

150
$$\Pr[A | B] = \frac{\Pr[B | A] \cdot \Pr[A]}{\Pr[B]} \quad (3)$$

151 *Of course, we assume that $\Pr[B] \neq 0$, as we divide by it.*

152 For a nice example see “bowls and cookies” in http://en.wikipedia.org/wiki/Bayesian_inference#Probability_of_a_hypothesis.

154 *Proof.* We know that

155
$$\Pr[A | B] = \frac{\Pr[A, B]}{\Pr[B]}$$

156 and therefore

157
$$\Pr[A | B] \cdot \Pr[B] = \Pr[A, B] \quad (4)$$

158 Note that equation holds not only when $\Pr[B] \neq 0$, but also when $\Pr[B] = 0$ since then both sides of the equation are 0, though this is not of importance to us now.

160 Also, by exchanging A and B , we get

161
$$\Pr[B | A] \cdot \Pr[A] = \Pr[B, A] \quad (5)$$

162 But of course,

163
$$\Pr[B, A] = \Pr[A, B]$$

164 and therefore from (4) and (5)

165
$$\Pr[A | B] \cdot \Pr[B] = \Pr[B | A] \cdot \Pr[A]$$

166 Therefore:

167
$$\Pr[A | B] = \frac{\Pr[B | A] \cdot \Pr[A]}{\Pr[B]}$$

168 □

169 **Definition 6.** (Formally,) A and B are independent iff

170
$$\Pr[A, B] = \Pr[A] \cdot \Pr[B] \quad (6)$$

171 Let us go immediately to the intuition.

172 **Theorem 2.** If $\Pr[B] \neq 0$, then A and B are independent iff

$$\Pr[A | B] = \Pr[A]$$

173

174 (This could have served as a more intuitive definition. What it tells us that knowing that we are “in” B does
175 not help us in “predicting” whether we are also in A . The reason we had definition in (6) was that this could
176 be written even for the case when $\Pr[B] = 0$, whereas we can talk about $\Pr[A | B]$ only when $\Pr[B] \neq 0$,
177 see (1).)

178 *Proof.* Let us start with (6) and continue

$$\Pr[A, B] = \Pr[A] \cdot \Pr[B]$$

179

180 iff (by rearranging)

$$\Pr[A] = \frac{\Pr[A, B]}{\Pr[B]}.$$

181

182 iff (by definition of $\Pr[A | B]$, (1))

$$\Pr[A] = \Pr[A | B].$$

183

184

□

185 **Observation 3.** Note that if A is independent of B then B is independent of A .

186 *Indeed,*

$$\Pr[A | B] = \Pr[A] \Leftrightarrow \Pr[A, B] = \Pr[A] \cdot \Pr[B] \Leftrightarrow \Pr[B, A] = \Pr[B] \cdot \Pr[A] \Leftrightarrow \Pr[B | A] = \Pr[B].$$

187

Observation 4.

188

$$\Pr[B] = \Pr[B | A] \cdot \Pr[A] + \Pr[B | \bar{A}] \cdot \Pr[\bar{A}]. \quad (7)$$

189 Look at Fig. 11. We will discuss only the case where B overlaps both A and \bar{A} , as in the figure; the other
190 cases are even simpler.

191 Using (1) and its complement,

$$\Pr[A | B] = \frac{\Pr[A, B]}{\Pr[B]} \quad \text{and} \quad \Pr[\bar{A} | B] = \frac{\Pr[\bar{A}, B]}{\Pr[B]}$$

193 the claim is equivalent to

194

$$\Pr[B] = \Pr[B, A] + \Pr[B, \bar{A}]$$

195 which of course is true.

196 **Remark 1.** There is a symbol sometimes used to denote independence. To write that A and B are independent,
197 we can write

$$A \perp\!\!\!\perp B$$

198

199 **Corollary 1.** From (3) and (7), we immediately get

$$\Pr[A | B] = \frac{\Pr[B | A] \cdot \Pr[A]}{\underbrace{\Pr[B | A] \cdot \Pr[A] + \Pr[B | \bar{A}] \cdot \Pr[\bar{A}]}_{\text{The denominator, ofcourse, is just } \Pr[B]}} \quad (8)$$

201 1.4 Practicing Bayesian thinking

202 **Example 5.** On a table there are two coins, one C_0 and the other C_2 . Without looking, you pick one and toss
203 it twice. You get the sequence HH. What is the probability that you picked C_2 ?

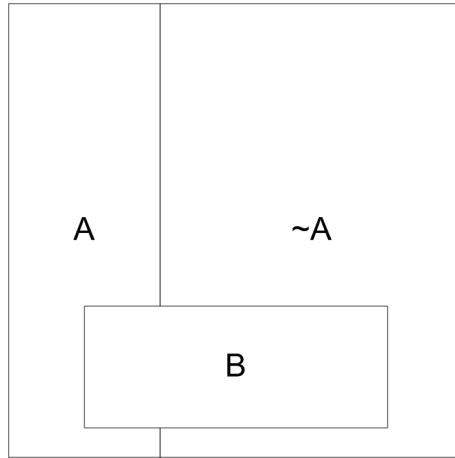


Figure 11: B and A and conditional probabilities.

204 Let C be the result of choosing the coin, so either C_0 or C_2 . Instead of writing $\Pr[C = C_0]$ we will write just
 205 $\Pr[C_0]$, and similarly elsewhere.

206 Using (3), we get

$$207 \quad \Pr[C_2 \mid HH] = \frac{\Pr[HH \mid C_2] \cdot \Pr[C_2]}{\Pr[HH]}$$

208 and

- 209 1. $\Pr[HH \mid C_2] = 1$, by the property of C_2
- 210 2. $\Pr[C = C_2] = 1/2$, by our random selection of the coin
- 211 3. $\Pr[HH] = 1/8 + 1/2$, from looking at Fig. 12 and checking what happens at the leaves

212 So, the final result is

$$213 \quad \frac{1 \cdot \frac{1}{2}}{\frac{1}{8} + \frac{1}{2}} = \frac{4}{5}.$$

214 .

215 **Example 6.** Let us redo Example 5.

- 216 1. as above
- 217 2. as above
- 218 3. We do not compute $\Pr[HH]$ and do not need to draw any figures. We use (2), which here is

$$219 \quad 1 = \Pr[C_2 \mid HH] + \Pr[C_1 \mid HH] = \frac{\Pr[HH \mid C_2] \cdot \Pr[C_2]}{\Pr[HH]} + \frac{\Pr[HH \mid C_0] \cdot \Pr[C_0]}{\Pr[HH]},$$

220 from which

$$221 \quad \Pr[HH] = \Pr[HH \mid C_2] \cdot \Pr[C_2] + \Pr[HH \mid C_0] \cdot \Pr[C_0],$$

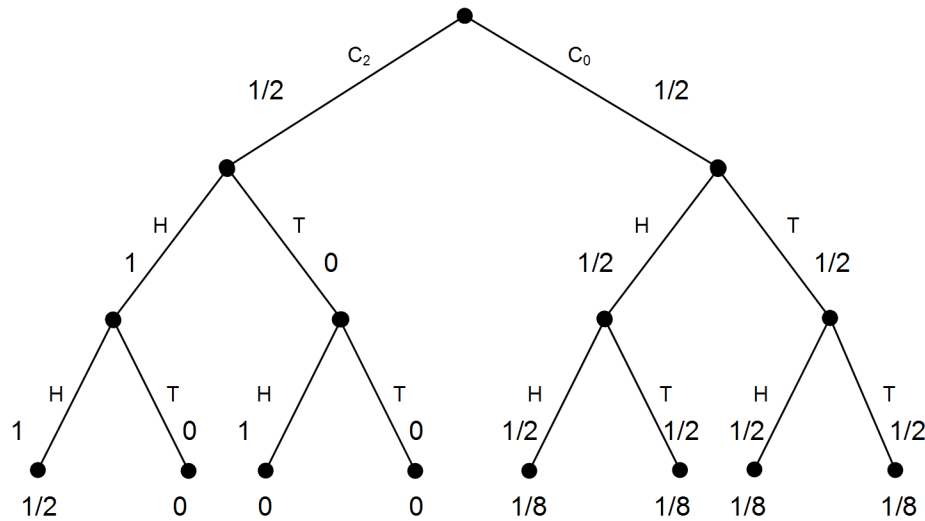


Figure 12: Tossing one of two coins

which is very easy to compute, without looking at Fig. 12. (Actually, we have rederived (8) for this case.)

Example 7. This is very important. Various businesses advertise tests such as total CT scan for symptomless people. After all, something could be lurking in the body that is dangerous but unknown to the person. But there could be false positive when the test declares a healthy person to be sick (or potentially sick, which may require additional, possibly life-threatening tests, like biopsies. Let us consider a scenario.

There exists a horrible disease and an imperfect test for it. There also exists an imperfect drug for the disease. Let us look at some numbers.

1. The probability of having the disease is 0.001 (very small)
2. If the disease is not treated, every person having it will die very quickly (very bad disease)
3. If the test is administrated to a sick person, it will correctly determine with probability 0.98 that the person is sick, and will declare with probability 0.02 that the person is healthy (very good test)
4. If the test is administrated to a healthy person, it will correctly determine with probability 0.99 that the person is healthy, and will declare with probability 0.01 that the person is sick (very good test)
5. If the drug is administered to a sick person, the person will be cured instantaneously (very good drug)
6. If the drug is administered to a healthy person, nothing will happen with probability 0.7, but the person will die instantenously with probability 0.3 (don't administer the drug to healthy people)

Should I take the test (and act on it, otherwise, why take it)?

Let us define some events:

1. A : I am healthy
2. B : the test came positive

243 What I am interested in are false positives, because if the test came positive, and I am healthy, and I take
 244 the drug, I have some probability of dying. What I want to know is what is the probability that I am actually
 245 healthy even though the test comes positive, so I want to know: $\Pr[A | B]$. Very easy, let's substitute into (8),
 246 which we repeat here:

$$247 \quad \Pr[A | B] = \frac{\Pr[B | A] \cdot \Pr[A]}{\Pr[B | A] \cdot \Pr[A] + \Pr[B | \bar{A}] \cdot \Pr[\bar{A}]}$$

248 We need to compute the terms appearing there:

- 249 1. $\Pr[B | A] = 0.01$
- 250 2. $\Pr[A] = 0.999$
- 251 3. $\Pr[B | \bar{A}] = 0.98$
- 252 4. $\Pr[\bar{A}] = 0.001$

253 So, plugging it in:

$$254 \quad \Pr[\text{healthy} | \text{test positive}] = \frac{0.01 \cdot 0.999}{0.01 \cdot 0.999 + 0.98 \cdot 0.001} = 0.91$$

255 We have many false positives. Let us now see what happens when every person who gets a positive test result
 256 is treated with the drug. It is best to look at Fig. 13

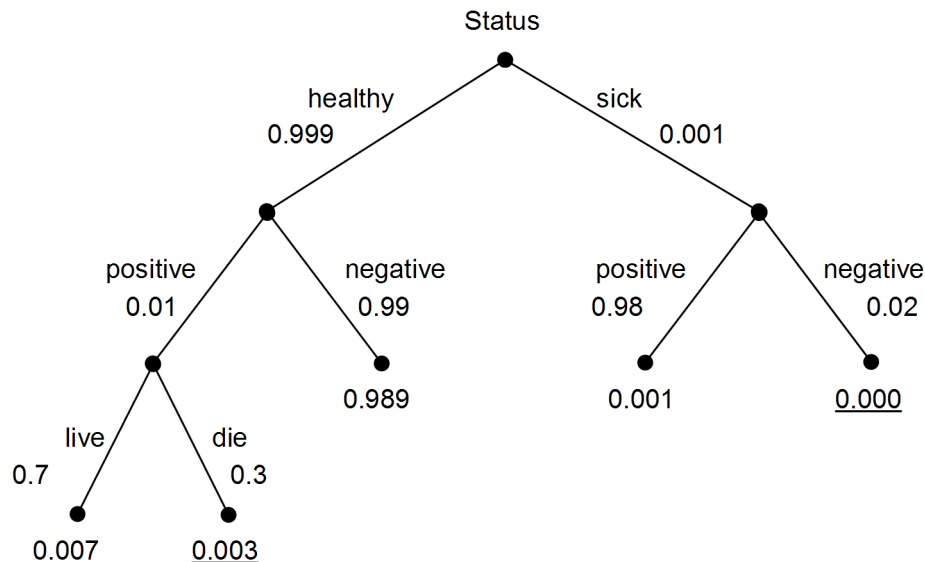


Figure 13: The result of testing and treating. The probability at a leaf is the product of probabilities along the path from the root, computed up to three decimal places. The cases resulting in death are underlined.

257 We see that if people are tested and treated, three times more will die than in the case when nobody is tested,
 258 as we have the probability of death of $0.003 + 0.000$, as opposed to 0.001 . In other words, if we do not test,

259 1 person out of 1000 will die, but if we test and treat based on the result of the test, 3 persons out of 1000 will
 260 die.

261 **Conclusion 1.** Do not test!

262 (For another version of this example, see [http://en.wikipedia.org/wiki/Bayes'_theorem#](http://en.wikipedia.org/wiki/Bayes'_theorem#Example_.232:_.2D_Drug_testing.)
 263 [Example_.232:_.2D_Drug_testing.](http://en.wikipedia.org/wiki/Bayes'_theorem#Example_.232:_.2D_Drug_testing.))

264 This was an extremely important example, of fundamental importance to physicians, but unfortunately not
 265 understood by many of them as they do not understand Bayes' theorem. They think that if the test detects a
 266 large fraction of bad cases it is a very reliable test. There are documented cases of people who were told
 267 they surely had AIDS and committed suicide, when they did not have AIDS but the result of the test was not
 268 understood.

269 **Example 8.** For an actual example, see [http://www.nytimes.com/2014/11/06/health/study-](http://www.nytimes.com/2014/11/06/health/study-warns-against-overdiagnosis-of-thyroid-cancer.html?module=Search&mabReward=relbias%3As)
 270 [warns-against-overdiagnosis-of-thyroid-cancer.html?module=Search&mabReward=](http://www.nytimes.com/2014/11/06/health/study-warns-against-overdiagnosis-of-thyroid-cancer.html?module=Search&mabReward=relbias%3As)
 271 [relbias%3As](http://www.nytimes.com/2014/11/06/health/study-warns-against-overdiagnosis-of-thyroid-cancer.html?module=Search&mabReward=relbias%3As).

272 1.5 Chain rule

273 See also [http://en.wikipedia.org/wiki/Chain_rule_\(probability\)](http://en.wikipedia.org/wiki/Chain_rule_(probability)), on which this writeup is
 274 based.

275 Let us look again at (1).and write it as

$$276 \Pr[A, B] = \Pr[A | B] \cdot \Pr[B]. \quad (9)$$

277 Assume now that event B is the intersection (conjunction) of events C and D . That is,

$$278 B = C \cap D \quad \text{or} \quad B = C \wedge D,$$

279 which are two ways of writing the same condition. We have, therefore,

$$280 \Pr[A, B] = \Pr[A, C, D]$$

281 (as we interested in the probability/area of $A \cap C \cap D$) and

$$282 \Pr[B] = \Pr[C, D]$$

$$283 \Pr[B] = \Pr[C, D] = \Pr[C | D] \cdot \Pr[D].$$

284 from (9). Putting these together

$$285 \Pr[A, C, D] = \Pr[A | C, D] \cdot \Pr[C | D] \cdot \Pr[D].$$

286 And if we have events A_1, A_2, \dots, A_n

$$\begin{aligned} \Pr[A_n, A_{n-1}, \dots, A_1] &= \Pr[A_n | A_{n-1}, A_{n-2}, \dots, A_1] \cdot \Pr[A_{n-1}, A_{n-2}, \dots, A_1] \\ &= \Pr[A_n | A_{n-1}, A_{n-2}, \dots, A_1] \cdot \Pr[A_{n-1} | A_{n-2}, \dots, A_1] \cdot \Pr[A_{n-2}, \dots, A_1] \\ &= \dots, \end{aligned}$$

287 or writing concisely

$$\Pr[\cap_1^n A_k] = \prod_1^n \Pr[A_k | \cap_1^{k-1} A_k]$$

288

289 Quoting “[The chain rule] permits the calculation of any member of the joint distribution of a set of random
290 variables using only conditional probabilities. The rule is useful in the study of Bayesian networks, which
291 describe a probability distribution in terms of conditional probabilities.”

292 1.6 Joint probability distribution and marginalization

293 This concept is meaningful if we have at least two random variables. So let us talk about this simple case: two
294 random variables. See also http://en.wikipedia.org/wiki/Joint_probability_distribution.

295 Assume we have a coin and a die and they are “somewhat connected”. When we toss them, one of the twelve
296 possibilities takes place and the probabilities of each of the twelve are given in Fig. fig:joint01. So, e.g.,
297 $\Pr[1, H] = 0.05$. This is *joint probability distribution*.

die \ coin	coin	
	H	T
1	0.05	0.10
2	0.05	0.15
3	0.00	0.05
4	0.20	0.05
5	0.10	0.15
6	0.05	0.10

Figure 14: Joint probability distribution for a coin and a die.

Given the information in Fig. 14, we can compute all interesting probabilities, such as

$$\Pr[H] = \sum_{i=1}^6 \Pr[i, H] \quad (10)$$

298 and

$$\Pr[1 | H] = \frac{\Pr[1, H]}{\Pr[H]}.$$

299

300 Equation (10) is an example of *marginalization*. We compute the probability of H by summing the proba-
301 bilities of the elementary events in which H occurs. This procedure is called so because customarily, such
302 probabilities were written on the margin. So in our case we get Fig. 15.

303 The probabilities $\Pr[H]$ and $\Pr[1]$ were not independent (it is clear that the first probability refers to the coin
304 and the second to the die) because

$$\Pr[1, H] \neq \Pr[1] \cdot \Pr[H].$$

305

<div>die \ coin</div>		coin		H + T
		H	T	
	1	0.05	0.10	0.15
	2	0.05	0.15	0.20
	3	0.00	0.05	0.05
	4	0.20	0.05	0.25
	5	0.10	0.15	0.20
	6	0.05	0.10	0.15
1 + 2 + 3 + 4 + 5 + 6		0.45	0.55	1.00

Figure 15: Joint probability for a coin and a die with marginalization.

2 Entropy

Entropy is a fundamental concept both in physics and information theory. We are only interested in the latter. In information theory, very roughly speaking, entropy quantifies how many bits it takes to describe a system.

To capture the value of information (actually he was interested in transmission of phone calls), Shannon http://en.wikipedia.org/wiki/Claude_Shannon introduced the concept of entropy http://en.wikipedia.org/wiki/Information_entropy#Entropy_as_information_content; he actually invented information theory. (As an aside, before that work, he wrote quite a remarkable MS thesis.)

We will start with an example, which will explain what the core issue is.

2.1 Example

Assume that we have a board of size 8×8 ; that is it has 64 equal squares. It is partitioned into 4 stripes of sizes: 8×4 colored red (R), 8×2 colored blue (B), 8×1 colored green (G), and 8×1 colored yellow (Y).

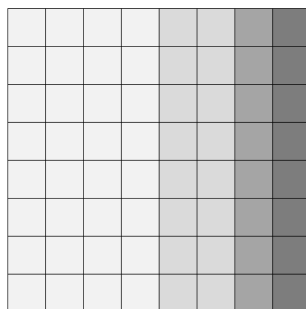


Figure 16: The square being covered by a toss. Different shades correspond to different colors.

Alice tosses a 1×1 cube on the board which ends up completely covering one square, with the probability of each square being covered of $1/64$. Alice tells Bob what was the color of the square. How much information did Bob get? We will start by asking how much information is needed to specify the square and how much out of this information the knowledge of its color provides.

321 It is necessary to have $\log_2 64 = 6$ bits to specify any square. There are four cases:

- 322 1. If the color was red, then Bob needs additional $\log_2 32 = 5$ bits to find out which of the 32 squares the
323 cube landed on, so he only has 1 bit.
- 324 2. If the color was blue, then Bob needs additional $\log_2 16 = 4$ bits, to find out which of the 16 squares
325 the cube landed on so he only has 2 bits.
- 326 3. If the color was green to find out which of the 8 squares the cube landed on, then Bob needs additional
327 $\log_2 8 = 3$ bits, so he only has 3 bits.
- 328 4. If the color was yellow, then Bob needs additional $\log_2 8 = 3$ bits to find out which of the 8 squares
329 the cube landed on, so he only has 3 bits.

330 Taking into account the probabilities of ending in different colors, the expected (average) amount of informa-
331 tion Bob has after being told the color is

$$\sum_{i \in \text{Red, Blue, Green, Yellow}} (\text{the probability of landing on } i) \cdot (\text{the number of bits acquired by landing on } i)$$

332 which is

$$\frac{1}{2}1 + \frac{1}{4}2 + \frac{1}{8}3 + \frac{1}{8}3 = \frac{1}{2}\log_2 2 + \frac{1}{4}\log_2 4 + \frac{1}{8}\log_2 8 + \frac{1}{8}\log_2 8 = 1.75.$$

333 So the actual formula is

$$\sum_{i \in \text{Red, Blue, Green, Yellow}} (\text{the probability of landing on } i) \cdot \log_2 \left(\frac{1}{\text{the probability of landing on } i} \right)$$

334 This is the *entropy* \mathcal{H} of the toss.

335 Frequently in the literature, the Greek capital letter eta <http://en.wikipedia.org/wiki/Eta> (should be
336 pronounced “heta” <http://en.wikipedia.org/wiki/Heta>), H, is used, which looks exactly like capital
337 Latin H <http://en.wikipedia.org/wiki/H>. To avoid confusion, I will use calligraphic H, that is \mathcal{H} .

338 Note that when you make an observation/experiment, the higher the probability of the outcome, the less you
339 learn.
340

341 2.2 A derivation

342 Consider a rectangle a of some n squares, b_1, \dots, b_n , each of size 1×1 , with the actual shape of the rectangle
343 immaterial. The squares are partitioned into k subsets c_j , $j = 1, \dots, k$ with c_j containing n_j squares. Of
344 course, $\sum_j n_j = n$. Both Alice and Bob know this. We assume that, n, n_1, \dots, n_k are all powers of 2.

345 A cube of size $1 \times 1 \times 1$ can be thrown at a and it will end up in a random square exactly covering it. The
346 probability of each square being covered by any such single throw is $1/n$.

347 Alice throws a cube on a and it lands up on square b_i . Alice sees which square it is, so she knows the complete
348 state of the universe and records i by using exactly the optimal number of bits, $\log_2 n$. Let $c_i \in b_j$. She tells
349 Bob the value of j . How much information did Bob gain from this?

In order for Bob to know the complete state of the universe he needs to know which one of its n_j squares of b_j was covered. For this he needs $\log_2 n_j$ bits, about which he knows nothing. Therefore to completely know the state of the universe, he misses $\log_2 n_j$ bits. So he knows

$$\log_2 n - \log_2 n_j = \log_2 \frac{n}{n_j}$$

bits.

The $\Pr[a_i \in b_j] = n_j/n$ and we will write p_j for $\Pr[a_i \in b_j]$. Therefore the information that Bob got from Alice is just $\log_2 1/p_j = -\log_2 p_j$. Averaging over all the outcomes of the throw of the die, the expectation of information that Bob got is exactly:

$$\mathcal{H} = - \sum_j p_j \log_2 p_j. \quad (11)$$

Note that everything was *not* the function of n, n_1, \dots, n_k , but *only* of the ratios $n_1/n, \dots, n_k/n$, that is of the probabilities p_1, \dots, p_k .

2.3 The case of general probabilities

Even when the probabilities do not satisfy “the inverse of power of 2” condition, (11) stating the expected amount of information obtained still holds, but to show that requires more work. And it holds for any “experiment” in which there are some n outcomes with outcome i occurring with probability p_i , that is

$$\mathcal{H} = \sum_{i=1}^n p_i \log_2 \frac{1}{p_i} = - \sum_{i=1}^n p_i \log_2 p_i.$$

The latter formula, with “—” is customarily used as it is easier to typeset without increasing the height of the line.

Note that in some other fields a base different from 2 is used for the logarithm. This is just a choice of a different constant and may be better when the information is not necessarily specified in bits. Typically, in physics, natural log (to the base of “e” is used).

There is a simple calculator for entropy on the web at <http://planetcalc.com/2476/>.

Note 2. To compute $\log_2 z$, if you can only get $\log_{10} z$ from your calculator, use

$$\log_2 z = \frac{\log_{10} z}{\log_{10} 2} \quad (12)$$

This follows from $z = 2^{\log_2 z}$, taking \log_{10} from both sides. There is also a free calculator with a “log₂” key <http://www.bestsoftware4download.com/software/t-free-esbcalc-freeware-calculator-download-hisqvad.html>.

Note 3. Sometimes (not as natural for us) log to a different base is used, which just changes multiplicative constants. The latter is analogous to (12), as the basis 10 could be replaced by a different basis.

379 **Note 4.** *If some event is of probability 0, we get a term*

$$0 \cdot \log_2 1/0 = 0$$

380

381 *It is correct to set the value to 0, even though $1/0$ is infinity and therefore $\log_2 1/0$ is infinite also and therefore*
382 *formally undefined.*

383 *We can say that because*

$$\lim_{x \rightarrow +0} x \log_2 x = 0.$$

384

385 **2.4 Randomness and Entropy**

386 We are given one bit. It is, of course 0 or 1. Can we determine whether it was obtained from a random process
387 such as tossing a relatively fair coin (we write “relatively fair”, as very unfair coin is not very random, loosely
388 speaking)? We cannot.

389 We are given a long sequence of bits, $\mathbf{b} = b_1, b_2, \dots, b_n$, with n big. Can we determine whether it was
390 obtained from a random process such as tossing a fair coin?

391 Let us think about a program that could print out \mathbf{b} . It is easy to do, we just put \mathbf{b} as a constant and the
392 program prints it out. But this program’s length is about n . So here is a possible definition

393 Sequence \mathbf{b} is random if and only if any program that prints it has the length of at least about n .

394 So, there is no short process to produce this sequence

395 Another definition could state that the entropy of \mathbf{b} is about n or perhaps

396 Sequence \mathbf{b} is random if and only if any lossless compressed version of it has the length of at
397 least about n .

398 Entropy is also related to efficient coding: short strings coding longer strings: compression. This is the topic
399 of lossless compression and not lossy compression, which is used, e.g., in jpg format.

400 **2.5 Additional interpretations of entropy**

401 We return to our example in Sec. 2.1.

402 Alice again tosses the cube on the square and Bob wants to find out what was the color on which it landed by
403 asking questions to which there are Yes/No answers. He may want to use the binary search tree shown in
404 Fig. 17. What is the expected number of questions Bob has to ask? He has to ask

$$\text{Number of questions Bob has to ask} = \begin{cases} 1 & \text{color is Red} \\ 2 & \text{color is Blue} \\ 3 & \text{color is Green} \\ 3 & \text{color is Yellow.} \end{cases}$$

405

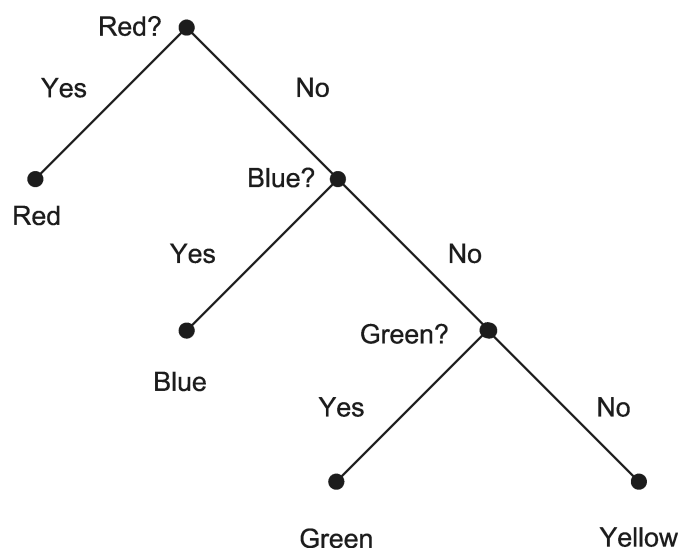


Figure 17: The binary search tree used to determine the color.

So the expected number of question that he has to ask is

$$\frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 = 1.75,$$

which is our \mathcal{H} .

2.6 Entropy obtained by randomly choosing from a set of symbols

These tables are taken from http://en.wikipedia.org/wiki/Password_strength.

i	Symbols in set	Set size	Entropy per symbol
1	Arabic numerals (0–9) (e.g. PIN)	10	3.322 bits
2	hexadecimal numerals (0–9, A–F) (e.g. WEP keys)	16	4.000 bits
3	Case insensitive Latin alphabet (a–z or A–Z)	26	4.700 bits
4	Case insensitive alphanumeric (a–z or A–Z, 0–9)	36	5.170 bits
5	Case sensitive Latin alphabet (a–z, A–Z)	52	5.700 bits
6	Case sensitive alphanumeric (a–z, A–Z, 0–9)	62	5.954 bits
7	All ASCII printable characters	95	6.570 bits
8	All extended ASCII printable characters	218	7.768 bits
9	Diceware word list	7 776	12.925 bits

Figure 18: Entropy per symbol for different symbol sets.

Desired password entropy	1	2	3	4	5	6	7	8	9
32	10	8	7	7	6	6	5	5	3
40	13	10	9	8	8	7	7	6	4
64	20	16	14	13	12	11	10	9	5
80	25	20	18	16	15	14	13	11	7
96	29	24	21	19	17	17	15	13	8
128	39	32	28	25	23	22	20	17	10
160	49	40	35	31	29	27	25	21	13
192	58	48	41	38	34	33	30	25	15
224	68	56	48	44	40	38	35	29	18
256	78	64	55	50	45	43	39	33	20
384	116	96	82	75	68	65	59	50	30
512	155	128	109	100	90	86	78	66	40
1 024	309	256	218	199	180	172	156	132	80

Figure 19: Length of random passwords to achieve desired entropy. Integer labels in the top row refer to i in Fig. 18.

2.7 Examples for C_0 , C_1 , and D_1

Alice and Bob play a (one-sided) game. Alice tosses a coin. Bob guesses the result. If he guesses correctly, Alice gives him \$100. If he guesses incorrectly, he gets nothing. Bob knows which coin is used.

Eve sees the result before Bob guesses and offers to sell him the result, so he can “guess” the answer correctly. How much should Bob pay Eve to increase the amount of money he will actually make.

If $C = C_0$, then no matter what Bob says, he will win (on the average) in half the times, so he will win, on the average \$50 after each toss. So if he pays Eve any amount below \$50, he makes more money net. If Bob pays Eve \$49, Alice pays Bob \$100 and he will win exactly \$51 on each toss.

If $C = C_1$, then if Bob always says H, he will win in 99 out of 100 tosses (on the average). So he benefits by paying Eve only if it is less than \$1. Therefore, the value of knowing the result of the toss is very small (relatively speaking).

Let us consider one toss of C_0 . Computing, we get

$$\mathcal{H} = \overbrace{\frac{1}{2} \log_2 \frac{2}{1}}^{\text{heads}} + \overbrace{\frac{1}{2} \log_2 \frac{2}{1}}^{\text{tails}} = \frac{1}{2} + \frac{1}{2} = 1$$

Intuitively, knowing the toss result is knowing the difference between two outcomes, 1 bit of information.

Let us consider one toss of C_1 . Computing, we get

$$\mathcal{H} = \overbrace{\frac{99}{100} \log_2 \frac{100}{99}}^{\text{heads}} + \overbrace{\frac{1}{100} \log_2 \frac{100}{1}}^{\text{tails}} \approx 0.08 \quad (13)$$

427 **Example 9.** Let D_1 be a die in which 1 comes up with probability 0.95, and every other number with
 428 probability 0.01.

429 Alice tosses D_1 . Bob asks her yes/no questions to find out what the result of the toss was.

430 Bob uses the algorithm described in Fig. (20). It is easy to see that the expected number of questions to ask is
 431 1.1 and the entropy is 0.4 (one decimal point of accuracy).

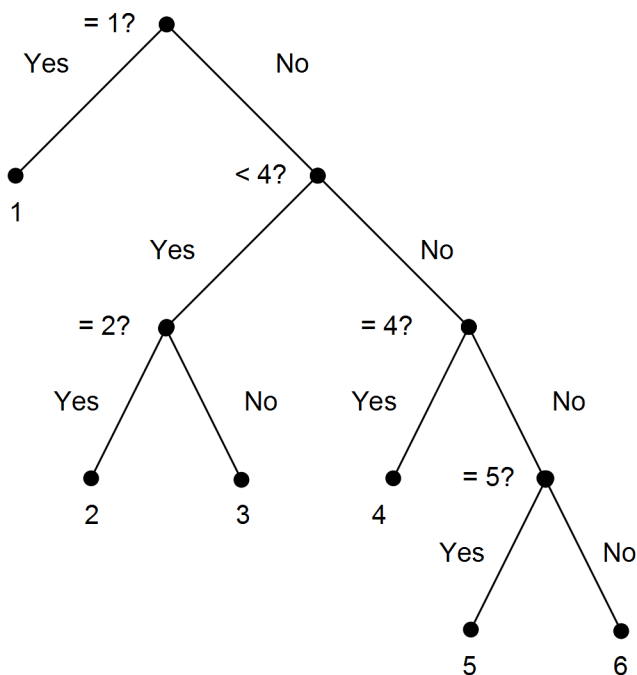


Figure 20: The algorithm Bob uses for determining the result of tossing D_1 .

432 There is an important theorem (we state a little informally), which we will not prove, which highlights the
 433 importance of entropy and its meaning as it tells how what are possible optimal binary search trees.

434 **Theorem 3.** Let \mathcal{X} be a random variable (that is, we get x_1 with probability p_1 , we get x_2 with probability
 435 p_2, \dots , etc.).

436 Then if you ask yes/no questions in “an optimal way,” the expected number of questions lies between $H(\mathcal{X})$
 437 and $H(\mathcal{X}) + 1$.

438 Looking at (13), we note that the expected number of question to find what the result of the toss was was
 439 between 0.08 and 1.08. And of course, we have to ask at least 1 question when there are at least two outcomes
 440 possible, so in fact the expected number of questions was between 1.00 and 1.08.

441 Consider coin C_1 . The entropy of the toss is low as event T comes rarely, as we get H 99% of the time.
 442 Nevertheless, we have to ask at least one question to find out the result of the toss, and it could simply be
 443 “Was the toss H?”

444 There is a popular game, called “Twenty Questions” http://en.wikipedia.org/wiki/Twenty_

445 **questions.** Let us say the “guesser” always guesses in exactly 20 questions. What is the entropy of this
446 game?

447 Interestingly, once the correct answer was provided after 0 questions were posed.

448 An interesting read: <http://danielwilkerson.com/entropy.html>.

449 2.8 Getting random numbers

450 It important in many applications to get truly random number or very good pseudo-random number.
451 To read more about getting random number and connection with choosing passwords and entropy, see
452 http://en.wikipedia.org/wiki/Password_strength.

453 Some sources you may want to explore if you ever need random numbers: http://en.wikipedia.org/wiki/Random_number_generator, <http://www.random.org/>, http://www.elmwoodmagic.com/full/Magic-Tricks-Magic-Books-Magic-DVDs-Red-Casino-Dice-Pack-of-4__3222.htm,
456 <http://www.casinochips3000.com/dice2.php>.

457 If you want to do some work, you can get essentially perfect uniformly distributed random bits from a
458 biased coin with unknown bias using the Von Neumann extractor http://en.wikipedia.org/wiki/Randomness_extractor#Von_Neumann_extractor.
459

460 3 Bayesian inference

461 3.1 Hypothesis and evidence

462 We have already seen, and proved, Bayes’ theorem. You can also look at http://en.wikipedia.org/wiki/Bayes'_theorem. We repeat it here with a different notation of the various probabilities:

$$464 \Pr[H | E] = \frac{\Pr[E | H] \cdot \Pr[H]}{\Pr[E]} \quad (14)$$

465 “ H ” refers to “hypothesis” and “ E ” to “evidence.” We will return to this later.

466 3.2 Two coins example

467 We will discuss an example of two coins, good (denoted by Good), which is our old C_0 ; and bad (denoted by
468 Bad), which is essentially our old C_2 . Good is a fair coin and Bad is a coin with heads on both sides. So the
469 probabilities of the result of the toss are:

- 470 • If Good is tossed then the result is H with probability 1/2 and T with probability 1/2
- 471 • If Bad is tossed then the result is H with probability 1 and T with probability 0

472 Alice picks up one of the coins randomly and tosses it twice, getting H both time. What is the probability that
473 the coin she has picked is Bad?

Using (14), and substituting Bad for H and HH for E we can write

$$\Pr[\text{Bad} | \text{HH}] = \frac{\Pr[\text{HH} | \text{Bad}] \cdot \Pr[\text{Bad}]}{\Pr[\text{HH}]}$$

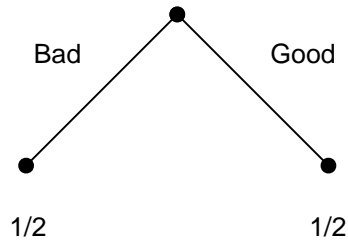


Figure 21: The probabilities after choosing a coin. ($\Pr[\text{Good}] = 1/2$ and $\Pr[\text{Bad}] = 1/2$)

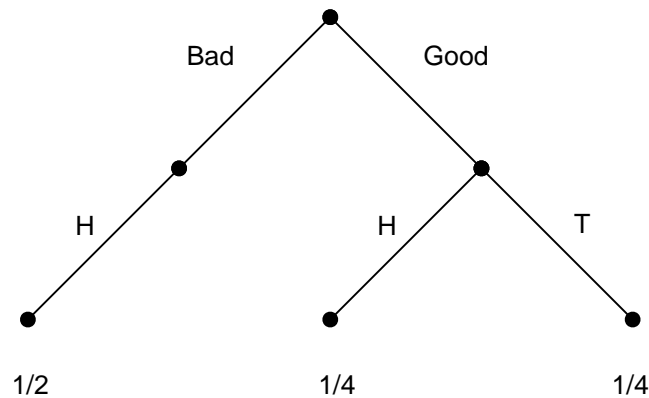


Figure 22: The probabilities after one toss of the coin. ($\Pr[\text{Good}] = 1/2$ and $\Pr[\text{Bad}] = 1/2$)

474 We now compute various probabilities on the right hand side. We look at Fig. 23, we see that:

- 475 • $\Pr[\text{HH} \mid \text{Bad}] = 1$
- 476 • $\Pr[\text{Bad}] = 1/2$
- 477 • $\Pr[\text{HH}] = 5/8$

We conclude that

$$\Pr[\text{Bad} \mid \text{HH}] = \frac{1 \cdot 1/2}{5/8} = \frac{4}{5}.$$

478 We have seen this before.

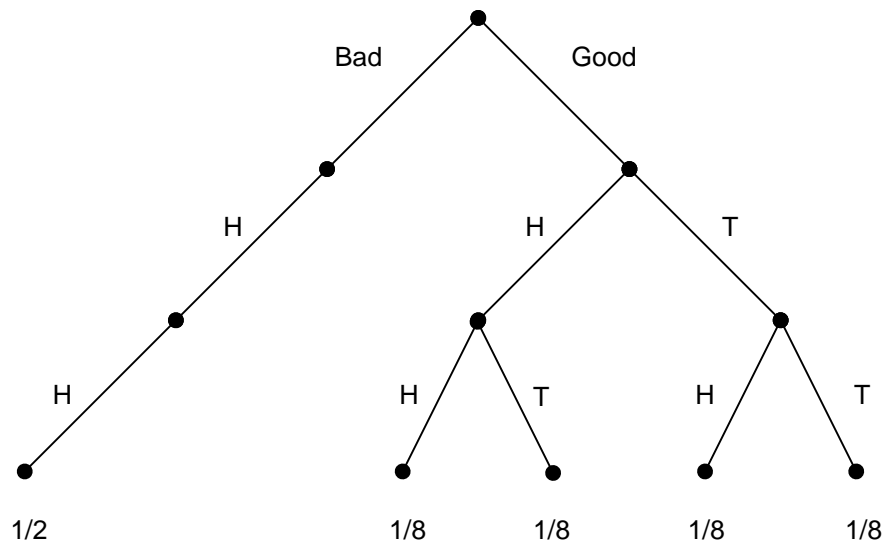


Figure 23: The probabilities after two tosses of the coin. ($\Pr[\text{Good}] = 1/2$ and $\Pr[\text{Bad}] = 1/2$)

4 Using Bayesian inference for incrementally increasing the confidence in the hypothesis that the coin tossed was Bad

4.1 The distribution of Good vs. Bad is known

Please pay careful attention to the difference between H and H . The first stands for “heads” and the second stands for “hypothesis”. This is an instance of the standard mathematics convention for typesetting:

- H is a constant and therefore in roman font
- H is a single-letter variable and therefore in italics

Alice does the same thing as before but she thinks a little differently. She picks a coin out of the two on the table, Good and Bad, without knowing which one is picked. There are two hypotheses:

1. The coin is Good
2. The coin is Bad

Now Alice assigns initial probabilities for these hypotheses, and she assigns

1. $\Pr[\text{Good}] = 1/2$
2. $\Pr[\text{Bad}] = 1/2$

This is not surprising, of course, as she knows what the two coins were. This is the situation depicted in Fig. 21.

495 She will now get additional evidence that may make her change the probabilities. So, we now interpret the
 496 terms in (14) as follows (see also http://en.wikipedia.org/wiki/Bayesian_inference.)

- 497 • H is hypothesis
- 498 • E is new evidence
- 499 • $\Pr[H]$ is *prior probability* of the hypothesis (or maybe our prejudice) H , before the new evidence
- 500 • $\Pr[E | H]$ is the *conditional probability* for the evidence to occur if the hypothesis H is true
- 501 • $\Pr[E]$ is the *a priori probability* of the evidence under all hypotheses (we will clarify soon and will
 502 have it in (15))
- 503 • $\Pr[H | E]$ is the *posterior probability* of H after the new evidence

Alice tosses the coin and sees H. This is her new evidence. She computes:

$$\Pr[\text{Bad} | H] = \frac{\Pr[H | \text{Bad}] \cdot \Pr[\text{Bad}]}{\Pr[H]} = \frac{1 \cdot (1/2)}{\Pr[H]}$$

$$\Pr[\text{Good} | H] = \frac{\Pr[H | \text{Good}] \cdot \Pr[\text{Good}]}{\Pr[H]} = \frac{(1/2) \cdot (1/2)}{\Pr[H]}$$

So,

$$\Pr[\text{Bad}] = \frac{1/2}{\Pr[H]} \quad \text{and} \quad \Pr[\text{Good}] = \frac{1/4}{\Pr[H]}$$

But what is $\Pr[H]$, which is $\Pr[E]$ as E is the new evidence? We do not want to look at Fig. 22. But we know that, of course:

$$\Pr[\text{Bad}] + \Pr[\text{Good}] = 1$$

that is

$$\frac{1/2}{\Pr[H]} + \frac{1/4}{\Pr[H]} = 1$$

and therefore:

$$\Pr[H] = 1/2 + 1/4 = 3/4$$

So,

$$\Pr[\text{Bad}] = \frac{1/2}{3/4} = 2/3 \quad \text{and} \quad \Pr[\text{Good}] = \frac{1/4}{3/4} = 1/3$$

Let us go back to the formulas

$$\Pr[\text{Bad} | H] = \frac{\Pr[H | \text{Bad}] \cdot \Pr[\text{Bad}]}{\Pr[H]} = \frac{\Pr[H | \text{Bad}] \cdot \Pr[\text{Bad}]}{\Pr[H | \text{Bad}] \cdot \Pr[\text{Bad}] + \Pr[H | \text{Good}] \cdot \Pr[\text{Good}]}$$

$$\Pr[\text{Good} | H] = \frac{\Pr[H | \text{Good}] \cdot \Pr[\text{Good}]}{\Pr[H]} = \frac{\Pr[H | \text{Good}] \cdot \Pr[\text{Good}]}{\Pr[H | \text{Bad}] \cdot \Pr[\text{Bad}] + \Pr[H | \text{Good}] \cdot \Pr[\text{Good}]}$$

Alice continues tossing the coin as long as she likes, incrementally modifying the probabilities based on the new evidence she gets.

What to do in the general case: We have some hypotheses: H_1, H_2, \dots, H_n . Then we can write Bayes' formula as:

$$\Pr[H_i | E] = \frac{\Pr[E | H_i] \cdot \Pr[H_i]}{\sum_{j=1}^n \Pr[E | H_j] \cdot \Pr[H_j]} \quad \text{for } i = 1, \dots, n \quad (15)$$

which is just a generalization of (8) when there are more than two cases, and we do not need to know $\Pr[E]$, if we know the conditional probabilities.

We will discuss this again in section 4.3.

4.2 The distribution of Good vs. Bad is not known

Bob comes to Alice hands her a coin and tells her the following story.

There is a barrel full of coins, some of them are Good and some of them are Bad. I know that more of them are Bad than there are Good, but I do not know the fraction. I picked out a coin randomly and did not look at it, and here it is. Please toss it a number of times and based on the results give me your estimate of the probabilities that the coin is either Bad or Good.

Before tossing the coin, Alice thinks to herself as follows:

I have to start with some initial probabilities. The fraction of the Bad is more than 1/2 and less than 1. So let me assume as prior probabilities that $\Pr[\text{Bad}] = 3/4$ and $\Pr[\text{Good}] = 1/4$, essentially in the middle of the possibilities, as I, Alice, cannot think of anything better. (Note this is subjective, somebody else may have reasons to think of different initial probabilities.)

She tosses the coin and it comes up H. She computes:

$$\begin{aligned} \Pr[\text{Bad} | H] &= \frac{\Pr[H | \text{Bad}] \cdot \Pr[\text{Bad}]}{\Pr[H]} = \frac{1 \cdot (3/4)}{\Pr[H]} \\ \Pr[\text{Good} | H] &= \frac{\Pr[H | \text{Good}] \cdot \Pr[\text{Good}]}{\Pr[H]} = \frac{(1/2) \cdot (1/4)}{\Pr[H]} \end{aligned}$$

So,

$$\Pr[\text{Bad}] = \frac{3/4}{\Pr[H]} \quad \text{and} \quad \Pr[\text{Good}] = \frac{1/8}{\Pr[H]}$$

But what is $\Pr[H]$? We do not want to look at the equivalent of Fig. 22 with correct probabilities because we cannot create it. But we know that

$$\Pr[\text{Bad}] + \Pr[\text{Good}] = 1 \quad (16)$$

that is

$$\frac{3/4}{\Pr[H]} + \frac{1/8}{\Pr[H]} = 1$$

and therefore we can compute a value for $\Pr[H]$ just to make sure that (16): and in our case

527

$$\frac{1 \cdot (3/4)}{(1/2) \cdot (1/4)}$$

$$\Pr[H] = 3/4 + 1/8 = 7/8$$

So,

$$\Pr[\text{Bad}] = \frac{3/4}{7/8} = 6/7 \quad \text{and} \quad \Pr[\text{Good}] = \frac{1/8}{7/8} = 1/7$$

528 **Note 5.** One can also discuss a related concept odds ratio, that is $\Pr[\text{Bad} | H]$ vs. $\Pr[\text{Good} | H]$. This
529 describes the current odds for Bad vs. Good based on current beliefs/hypothesis.

530 For this it is enough to look at

$$\Pr[H | \text{Bad}] \cdot \Pr[\text{Bad}] \quad \text{vs.} \quad \Pr[H | \text{Good}] \cdot \Pr[\text{Good}],$$

531

532 and we do not need to concern ourselves with the denominator.

533 Notice that this combines her *priors* (her prior beliefs) with the new evidence to get new beliefs.

As Alice continues tossing the coin, and assuming she always gets H, $\Pr[\text{Bad}]$ goes up and $\Pr[\text{Good}]$ goes down. Let us see what happens if she ever gets T as her new evidence.

$$\begin{aligned} \Pr[\text{Bad} | T] &= \frac{\Pr[T | \text{Bad}] \cdot \Pr[\text{Bad}]}{\Pr[T]} = \frac{0 \cdot \Pr[\text{Bad}]}{\Pr[T]} \\ \Pr[\text{Good} | T] &= \frac{\Pr[T | \text{Good}] \cdot \Pr[\text{Good}]}{\Pr[T]} = \frac{(1/2) \cdot \Pr[\text{Good}]}{\Pr[T]} \end{aligned}$$

so

$$\Pr[\text{Bad}] = 0 \quad \text{and} \quad \Pr[\text{Good}] = 1$$

534 and if Alice continues tossing the coin and using (15), the probabilities (not surprisingly) do not change. In
535 this case, once a coin is proved to be Good there is no way that it is the case that it is going to actually turn
536 out to be Bad.

537

$$\Pr[\text{Bad} | E] = \frac{\Pr[E | \text{Bad}] \cdot \Pr[\text{Bad}]}{\Pr[E]} = \frac{\Pr[E | \text{Bad}] \cdot 0}{\Pr[E]} = 0$$

538 Let's generalize, if at any point one of the hypotheses has probability 0, it will continue being 0 and no
539 amount of subsequent contrary evidence can change it to > 0 .

540 **Note 6.** That is why your initial priors should never include 0, as you will never get out of it. So (I do not
541 recall who came up with the following) you should never assign initial probability of 0 to the claim that the
542 moon is made of green (meaning young, such as cottage cheese) cheese. If you do that, no amount of evidence
543 such as thousands of astronauts you send to the moon coming back with green cheese as evidence will make
544 this probability positive.

545 **Note 7.** *If at the beginning Bob also tells her that he had previously conducted this experiment with many*
 546 *barrels and in a large majority of time of the times the coin turned out to be Bad because presumably, the*
 547 *fractions of Bad coins is generally very high, then Alice may assign $\Pr[\text{Bad}] = 9/10$ and $\Pr[\text{Good}] = 1/10$*
 548 *as her initial probabilities.*

549 4.3 Recapitulation

550 We review what we already know, but phrase it somewhat differently and more emphatically.

We have some number n of *prior* hypotheses, *priors*, each with some initial prior probability, that is we start with

H_1 with probability $\Pr[H_1]$
 H_2 with probability $\Pr[H_2]$
 \dots
 H_n with probability $\Pr[H_n]$

551 and, of course,

$$\sum_i \Pr[H_i] = 1.$$

552

553 where the priors were assigned based on our belief about how the universe works and what we know about
 554 the situation before obtaining new *evidence*, which could be perhaps *experimental* or *observational*.

We get new evidence E . It has impact on the probabilities and we now have probabilities *conditioned on* the new evidence we have, which presumably we know how to compute

H_1 with probability $\Pr[H_1 | E]$
 H_2 with probability $\Pr[H_2 | E]$
 \dots
 H_n with probability $\Pr[H_n | E]$

555 and of course

556

$$\sum_i \Pr[H_i | E] = 1, \tag{17}$$

557 as the probabilities deal with disjoint “alternatives,” and these alternatives cover “everything”: exactly one of
 558 the hypotheses is true.

Our new probabilities (which will be the new priors: for the next “iteration”) are

$\Pr[H_1] \leftarrow \Pr[H_1 | E]$
 $\Pr[H_2] \leftarrow \Pr[H_2 | E]$
 \dots
 $\Pr[H_n] \leftarrow \Pr[H_n | E]$

We know that these are

$$\begin{aligned}\Pr[H_1 | E] &= \frac{\Pr[E | H_1] \cdot \Pr[H_1]}{\Pr[E]} \\ \Pr[H_2 | E] &= \frac{\Pr[E | H_2] \cdot \Pr[H_2]}{\Pr[E]} \\ &\dots \\ \Pr[H_n | E] &= \frac{\Pr[E | H_n] \cdot \Pr[H_n]}{\Pr[E]}.\end{aligned}$$

Unfortunately, we do not know what $\Pr[E]$ is. Fortunately, we do not need to know it (directly) as we can compute it using (17).

$$\frac{\Pr[E | H_1] \cdot \Pr[H_1]}{\Pr[E]} + \frac{\Pr[E | H_2] \cdot \Pr[H_2]}{\Pr[E]} + \dots + \frac{\Pr[E | H_n] \cdot \Pr[H_n]}{\Pr[E]} = 1$$

from which

$$\Pr[E] = \Pr[E | H_1] \cdot \Pr[H_1] + \Pr[E | H_2] \cdot \Pr[H_2] + \dots + \Pr[E | H_n] \cdot \Pr[H_n]$$

Stating concisely, the new priors are

$$\frac{\Pr[E | H_i] \cdot \Pr[H_i]}{\sum_{i=1}^n \Pr[E | H_i] \cdot \Pr[H_i]} \quad \text{for } i = 1, 2, \dots, n.$$

5 Bayesian classification

See also http://en.wikipedia.org/wiki/Naive_Bayes_classifier.

We will have an extended example from which the general case will be clear. We will invent all the data.

5.1 Basic Bayesian classification

We know that 52% of people are female and 48% are male. We also have some data about some, presumably randomly selected people. The data are listed in Figs. 24 and 25.

Shoe Size	Number of Males	Number of Females
8	7	12
9	6	5
10	1	0

Figure 24: Number of males and females who had a certain shoe size.

A person comes. We measure and we get that the shoe size is 8 and the height is 67. What are the probabilities for the sex of the person? We have (we omit obvious labels: we do not write “shoe size = 8” and just write 8,

Height	Number of Males	Number of Females
65	8	10
66	7	7
67	20	16

Figure 25: Number of males of females who had a certain height (in inches).

etc.).

$$\Pr[\text{male} \mid 8, 67] = \frac{\Pr[8, 67 \mid \text{male}] \cdot \Pr[\text{male}]}{\Pr[8, 67]}$$

$$\Pr[\text{female} \mid 8, 67] = \frac{\Pr[8, 67 \mid \text{female}] \cdot \Pr[\text{female}]}{\Pr[8, 67]}$$

572 So what do we do with the right-hand sides?

573 1. $\Pr[8, 67 \mid \text{male}]$, we consider shortly

574 2. $\Pr[\text{male}]$ is our (current) prior and it is 0.48

575 3. $\Pr[8, 67]$ is our evidence and as we know we do not need to think about it and just derive it by normal-
576 ization (there is no need to consider the semantics of the situation: what is the meaning of this), which
577 means $\Pr[\text{male} \mid 8, 67] + \Pr[\text{female} \mid 8, 67] = 1$

578 4. $\Pr[8, 67 \mid \text{female}]$, we consider shortly

579 5. $\Pr[\text{female}]$ is our (current) prior and it is 0.52

580 Let us now consider people in general. We believe (and let's assume it is actually true) that tall people in
581 general have a large shoe size. Therefore, for people in general very likely we believe that

$$\Pr[8, 67] \neq \Pr[8] \cdot \Pr[67],$$

582

583 that is shoe size and height are not independent random variables.

This should be true (we believe) if we look just at males and just at females (conditioning on the person being a male or conditioning on the person being a female), that is

$$\Pr[8, 67 \mid \text{male}] \neq \Pr[8 \mid \text{male}] \cdot \Pr[67 \mid \text{male}]$$

$$\Pr[8, 67 \mid \text{female}] \neq \Pr[8 \mid \text{female}] \cdot \Pr[67 \mid \text{female}],$$

584 but we do not know how to compute $\Pr[8, 67 \mid \text{male}]$ and $\Pr[8, 67 \mid \text{female}]$, which we would like to know.

5.2 Naïve Bayesian classification

We pretend to be naïve (having or showing a lack of experience, judgment, or information; credulous, see <http://dictionary.reference.com/browse/naive>), we *pretend* that

$$\begin{aligned}\Pr[8, 67 \mid \text{male}] &= \Pr[8 \mid \text{male}] \cdot \Pr[67 \mid \text{male}] \\ \Pr[8, 67 \mid \text{female}] &= \Pr[8 \mid \text{female}] \cdot \Pr[67 \mid \text{female}],\end{aligned}$$

from which

$$\begin{aligned}\Pr[\text{male} \mid 8, 67] &= \frac{\Pr[8 \mid \text{male}] \cdot \Pr[67 \mid \text{male}] \cdot \Pr[\text{male}]}{\Pr[8, 67]} \\ \Pr[\text{female} \mid 8, 67] &= \frac{\Pr[8 \mid \text{female}] \cdot \Pr[67 \mid \text{female}] \cdot \Pr[\text{female}]}{\Pr[8, 67]}\end{aligned}$$

and looking at Figs. 24 and 25

1. $\Pr[8 \mid \text{male}] = 7/(7 + 6 + 1)$
2. $\Pr[67 \mid \text{male}] = 20/(8 + 7 + 20)$
3. $\Pr[8 \mid \text{female}] = 12/(12 + 5 + 0)$
4. $\Pr[67 \mid \text{female}] = 16/(10 + 7 + 16)$

and from these we can compute the needed probabilities for the new person being male or female.

5.3 Laplacian correction

All is nice so far but a new person comes and we measure that the shoe size is 10 and the height is 65. What are the probabilities for the person to be male or female? Using our procedure we look at

$$\begin{aligned}\Pr[\text{male} \mid 10, 65] &= \frac{\Pr[10 \mid \text{male}] \cdot \Pr[65 \mid \text{male}] \cdot \Pr[\text{male}]}{\Pr[10, 65]} \\ \Pr[\text{female} \mid 10, 65] &= \frac{\Pr[10 \mid \text{female}] \cdot \Pr[65 \mid \text{female}] \cdot \Pr[\text{female}]}{\Pr[10, 65]}.\end{aligned}$$

By looking at Fig. 24 we see that $\Pr[10 \mid \text{female}] = 0$ and therefore the prior for female for the next step is 0 and it will remain 0 *no matter what the new evidence is*.

For example, assume that we now obtain evidence about the weight of a sample of males and females and based on this have another table with weights of males and females. No matter what we learn about a weight of a new person, a person whose shoe size was 10 will never be classified as female no matter what the new evidence about the weight is, which could indicate that the person very likely is a female.

It is unacceptable, as we have seen, for the probability to be fixed to 0 no matter what the new evidence might say. So the prior *should never be* 0.

We will therefore apply the *Laplacian correction* to the priors so that no prior has probability of 0. One simple way of doing that is to artificially add 1 to each entry in Figs. 24 and 25 getting Figs. 26 and 27 and then proceeding as before using Figs. 26 and 27. We do not want to “discriminate” against any segment of the

604 population so we apply the correction to all, though maybe proportional as opposed to absolute correction
605 could have been better, by applying a reasonably equal percentage correction to all the entries as opposed to
606 just adding 1.

Shoe Size	Number of Males	Number of Females
8	8	13
9	7	6
10	2	1

Figure 26: Number of males and females who had a certain shoe size after Laplacian correction.

Height	Number of Males	Number of Females
65	9	11
66	8	8
67	21	17

Figure 27: Number of males of females who had a certain height (in inches) after Laplacian correction.