

**Ex. 2.1:** Is it possible to find a general formula for  $p(C|A+B)$ , analogous to (2-48), from the product and sum rules? If so, derive it; if not, explain why this cannot be done.

**Answer:**

Sure. We've derived Bayes theorem without naming it, which is simply product rule and commutativity, so we can say

$$\begin{aligned} p(C|A+B) &= \frac{p(C)p(A+B|C)}{p(A+B)} \\ &= \frac{p(C)(p(A)+p(B)-p(AB|C))}{p(A)+p(B)-p(AB)}. \end{aligned}$$

**Ex. 2.2:** Now suppose we have a set of propositions  $\{A_1, \dots, A_n\}$  which on information  $X$  are mutually exclusive:  $p(A_i A_j|X) = p(A_i|X)\delta_{ij}$ . Show that  $p(C|(A_1+A_2+\dots+A_n)X)$  is a weighted average of the separate plausibilities  $p(C|A_i X)$ :

$$p(C|(A_1+\dots+A_n)X) = p(C|A_1 X + A_2 X + \dots + A_n X) = \frac{\sum_i p(A_i|X)p(C|A_i X)}{\sum_i p(A_i|X)}.$$

**Answer:** For  $i \neq j$ ,

$$\begin{aligned} p(C|A+B) &= \frac{p(A+B|C)p(C)}{p(A+B)} \\ &= \frac{(p(A)+p(B)-p(AB|C))p(C)}{p(A)+p(B)-p(AB|C)}. \end{aligned}$$

Iterating this over  $\sum_i A_i$  we get

$$p\left(\sum_i A_i \middle| X\right) = \sum_i p(A_i|X).$$

Thus, and in a similar manner to the above exercise,

$$\begin{aligned} p\left(C \middle| \sum_i A_i X\right) &= \frac{p(\sum_i A_i|CX)p(C|X)}{p(\sum_i A_i|X)} \\ &= \frac{\sum_i p(A_i|CX)p(C|X)}{\sum_i p(A_i|X)} \\ &= \frac{\sum_i p(A_i|X)p(C|A_i X)}{\sum_i p(A_i|X)}. \end{aligned}$$

**Ex.:** Let  $A_i$  mutually exclusive ( $P(A_i A_j) = P(A_i)\delta_{ij}$ ). Then,  $P(\sum_i A_i) = \sum_i P(A_i)$ .

**Answer:**

If  $i \in \{1, 2\}$ , then

$$P(\sum_i A_i) = P(A_1 + A_2) \tag{1}$$

$$= P(A_1) + P(A_2) - \cancel{P(A_1 A_2)} \xrightarrow{0} \tag{2}$$

$$= \sum_i P(A_i). \tag{3}$$

Then, by induction,

$$P\left(\sum_{i=1}^{N+1} A_i\right) = P\left(\sum_{i=1}^N A_i\right) + P(A_{N+1}) - P\left(\sum_{i=1}^N A_i A_{N+1}\right) \quad (4)$$

$$= \sum_{i=1}^N P(A_i) + P(A_{N+1}) - \sum_{i=1}^N \underbrace{(A_i A_{N+1})}_0 \quad (5)$$

$$= \sum_{i=1}^N P(A_i). \quad (6)$$

**Ex. 2.3: Limits on Probability Values:** As soon as we have the numerical values  $a = P(A|C)$  and  $b = P(B|C)$ , the product and sum rules place some limits on the possible numerical values for their conjunction and disjunction. Supposing that  $a \leq b$ , show that the probability of the conjunction cannot exceed that of the least possible proposition:  $0 \leq P(AB|C) \leq a$ , and the probability of the disjunction cannot be less than that of the most probable proposition:  $b \leq P(A + B|C) \leq 1$ . Then show that, if  $a + b > 1$ , there is a stronger inequality for the conjunction; and if  $a + b < 1$  there is a stronger one for the disjunction. These necessary general inequalities are helpful in detecting errors in calculations.

**Answer:**

For the conjunction,

$$P(AB|C) = P(A|C)P(B|AC) \quad (7)$$

$$= aP(B|AC) \leq a. \quad (8)$$

Since  $P(AB|C) = b + a - P(A + B|C)$ , if  $b + a > 1$  then  $P(AB|C) > 1 - P(A + B|C) = P(\overline{A + B}|C)$ .

For the disjunction,

$$P(A + B|C) = P(A|C) + P(B|C) - P(AB|C) \quad (9)$$

$$= b + a - P(AB|C) \quad (10)$$

$$\geq b. \quad (11)$$

And also,  $b + a < 1$  implies  $P(A + B|C) < 1 - P(AB|C) = P(\overline{AB}|C)$ .

**Ex. 3.1:** Why isn't the multiplicity factor  $(3 - 16)$  just  $n!$ ? After all, we started this discussion by stipulating that the balls, in addition to having colors, also carry labels  $(1 \cdots N)$ , so that different permutations of the red balls among themselves, which give the  $r!$  in the denominator of  $(3 - 16)$ , are distinguishable arrangements.

**Answer:**

We've decomposed the event "drawing  $r$  red balls in  $n$  draws" as the union of events of the form "draw  $r$  red balls in  $n$  draws with the particular permutation of orders  $p$ " for permutations  $p$ , which are mutually exclusive and identically plausible. At doing this, we're using the events  $R_i$  and  $W_i$ , with the probability given by  $(3 - 14)$ . If we were to use "drawing the following  $n$  balls, where  $r$  are red", we should be accounting for the probability  $B_{ij}$  of drawing the  $i$ th ball after  $j$  draws, which would incorporate the missing factors.

**Ex. 3.2: Probability of a Full Set:** Suppose an urn contains  $N = \sum_i N_i$  balls,  $N_1$  of color 1,  $N_2$  of color 2,  $\cdots$ ,  $N_k$  of color  $k$ . We draw  $m$  balls without replacement; what is the probability that we have at least one of each color? Supposing  $k = 5$ , all  $N_i = 10$ , how many do we need to draw in order to have at least 90% probability of getting a full set?

**Answer:**

We're asked  $P(\sum_{r \in I} \prod_i r_i)$ , where  $I = \{(r_i) : r_i \geq 1 \forall i\}$ . We can decompose this as

$$P\left(\sum_{r \in I} \prod_i r_i | B\right) = \sum_{r \in I} P(r_1 \cdots r_k | B) \quad (12)$$

$$= \frac{1}{\binom{N}{m}} \sum_{r \in I} \prod_i \binom{N_i}{r_i} \quad (13)$$

A calculation is probably better done by calculating the missing terms instead, which are those for which at least one of those are 0, of which there are less, since at least one of the summands is already 0, the rest varying between the same intervals.

The “3.2.py” script does the maths, and outputs  $m = 15$  as the solution, with a probability of about 0.91.

**Ex. 3.3: Reasoning Backwards:** Suppose that in the previous exercise  $k$  is initially unknown, but we know that the urn contains exactly 50 balls. Drawing out 20 of them, we find 3 different colors; now what do we know about  $k$ ? We know from deductive reasoning (i.e., with certainty) that  $3 \leq k \leq 33$ ; but can you set narrower limits  $k_1 \leq k \leq k_2$  within which it is highly likely to be?

**Answer:**

I suppose this is the first use case of Bayes theorem as it's usually used. By letting the hypotheses  $H_k =$  “There are  $k$  different colors in the urn” and the data  $D =$  “3 different colors have been seen after 20 draws”, we want  $P(H_k|D)$ , but what our model tells us with certainty is  $P(D|H_k)$ , so we'll use Bayes' theorem for this:

$$P(H_k|D) = \frac{P(D|H_k)P(H_k)}{\sum_k P(D|H_k)P(H_k)}. \quad (14)$$

I will assume each  $H_k$  is equally likely, since I've got no information on anything. Thus,

$$P(H_k|D) = \frac{P(D|H_k)}{\sum_k P(D|H_k)}. \quad (15)$$

We're now left with the task of finding the probability that 3 different colors are seen in a set of 20 taken from 50, given there are  $k$  different colors. To this end, we split  $D = \sum_N DH_N$ , where  $N \in \mathbb{N}^k$  is such that  $N_i$  is the number of balls with the  $i$ th color, and get

$$P(D|H_k) = P\left(\sum_N DH_N \middle| H_k\right) \quad (16)$$

$$= \sum_N P(DH_N|H_k) \quad (17)$$

$$= \sum_N P(D|H_N H_k) P(H_N|H_k). \quad (18)$$

We once again propose an homogeneous distribution over the  $N$ , which gives  $P(H_N|H_k) = \binom{50-k}{k-1}^{-1}$ .

Now, let's think of the probability of drawing from only the colors 1, 2, and 3 in all draws. For the  $m$ th draw, the probability is

$$\frac{51 - m - \sum_{l=3}^{k-2} N_l}{51 - m} = \frac{51 - m - N_{\text{rem}}}{51 - m}, \quad (19)$$

so for the 20 draws, we get

$$\frac{(50 - N_{\text{rem}})! 30!}{(30 - N_{\text{rem}})! 50!}. \quad (20)$$

So, the chances of drawing only the first 3 colors are the sum over all of the  $N$  of this term. For each choice of  $N_{\text{rem}}$ , there are  $\binom{47 - N_{\text{rem}}}{2} = (47 - N_{\text{rem}})(46 - N_{\text{rem}})$  such terms. And, by symmetry of label permutation, for each choice of 3 colors the result is the same, so finally, by splitting  $D$  into the sum of “drawing only for these 3 particular colors”, we sum over these  $k!/(k-3)!$  choices, and get our final result:

$$P(D|H_k) = \frac{k!}{(k-3)!} \sum_{N_{\text{rem}}=k-3}^{47} (47 - N_{\text{rem}})(46 - N_{\text{rem}}) \frac{(50 - N_{\text{rem}})!}{(30 - N_{\text{rem}})!} \frac{30!}{50!} \binom{50 - k}{k-1}^{-1} \quad (21)$$

The “3.3.py” script implements this solution, giving the probability distribution shown in Figure 1

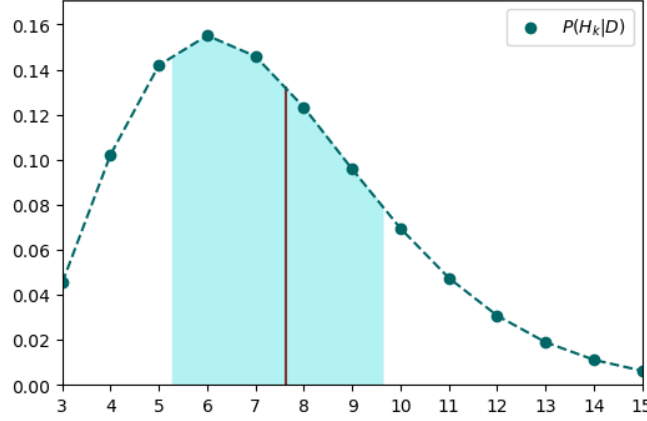


Figure 1: Probability distribution  $P(H_k|D)$  for different values of  $k$ . The shaded region goes between first and third quartile, and the vertical line denotes the (linearly extrapolated) median.

**Ex. 3.4: Matching:** The  $M$  urns are now numbered 1 to  $M$ , and  $M$  balls, also numbered 1 to  $M$ , are thrown into them, one in each urn. If the numbers of a ball and its urn are the same, we have a match. Show that the probability of at least one match is

$$P = \sum_{k=1}^M (-1)^{k+1} / k!.$$

As  $M \rightarrow \infty$ , this converges to  $1 - 1/e = 0.632$ . The result is surprising to many, because however large  $M$  is, there remains an appreciable probability of no match at all.

**Answer:**

I give up. The problem turns out to be that of counting derangements, which Wikipedia explains:

Let the urn 1 receive ball  $i$ . If urn  $i$  receives ball 1, then we’re left with the problem of counting derangements for  $M - 2$ . If urn  $i$  receives a ball other than 1, then it’s the problem of counting derangements for  $M - 1$ , since for each of the other urns there is 1 ball they may not receive (urn  $i$  shall not receive ball 1). Thus,

$$!n = (n-1)(!(n-1)) + (n-2)(!(n-2)), \quad (22)$$

which may now be solved by induction.

**Ex. 3.5: Occupancy:**  $N$  balls are tossed into  $M$  urns; there are evidently  $M^N$  ways this can be done. If the robot considers them all equally likely, what is its probability that each urn receives at least one ball?

**Answer:**

If they're equally likely, we can just count. We'll count the negative, summing over the number  $k$  of boxes with no matches. There are  $\binom{M}{k}$  such choices, and the rest of boxes allow for  $(M - k)^N$  choices, so there are

$$\sum_{k=1}^M \binom{M}{k} (M - k)^N \quad (23)$$

instances where no box receives balls, and thus

$$P(\text{"Each urn receives at least one ball"}) = 1 - \sum_{k=1}^M \binom{M}{k} \left(1 - \frac{k}{M}\right)^N \quad (24)$$

**Ex. 4.1:** Show that there is no such nontrivial extension of the binary case. More specifically, prove that if (4.28) and (4.29) hold with  $n > 2$ , then at most one of the factors

$$\frac{P(D_1|H_i X)}{P(D_1|\overline{H_i} X)} \cdots \frac{P(D_m|H_i X)}{P(D_m|\overline{H_i} X)}$$

is different from unity, therefore at most one of the data sets  $D_j$  can produce any updating of the probability for  $H_i$ .

**Answer:**

The equations mentioned are that of independence in  $D_1$  to  $D_m$  conditioned to both the hypothesis and its negation:

$$P(D_1, \dots, D_m | A X) = \prod_j P(D_j | A X) \quad \forall A \in \{H_i, \overline{H_i}\} \forall i \in [1, n]_{\mathbb{N}},$$

under the hypothesis that the  $H_i$  are exhaustive.

This means

$$\begin{aligned} P(H_i | D_j X) + P(\overline{H_i} | D_j X) &= 1, \\ (P(D_j | H_i X) - P(D_j | \overline{H_i} X)) P(H_i | X) &= P(D_j | X) - P(D_j | \overline{H_i} X). \end{aligned}$$

I believe the proof has to do with the fact that, for a proposition  $D_j$  such that  $P(D_j | H_i X) \neq P(D_j | \overline{H_i} X)$ ,

$$P(H_i | X) = \frac{P(D_j | X) - P(D_j | \overline{H_i} X)}{P(D_j | H_i X) - P(D_j | \overline{H_i} X)},$$

but so far it's leading me nowhere.

Maybe use that

$$\begin{aligned} P(D_k | \overline{H_i}) - P(D_k | H_i) &= \sum_{j \neq i} \frac{P(D_k | H_j) P(H_j)}{P(\overline{H_i})} - P(D_k | H_i) \\ &= \frac{1}{P(\overline{H_i})} \sum_j [P(D_k | H_j) - P(D_k | H_i)] P(H_j), \\ [P(D_k | \overline{H_i}) - P(D_k | H_i)]_k &= \frac{1}{P(\overline{H_i})} [P(D_k | H_j) - P(D_k | H_i)]_{kj} [P(H_j)]_j \end{aligned}$$

**Ex. 4.2:** Calculate the exact threshold of skepticism  $f_t(x, y)$ , supposing that proposition  $C$  has instead of  $10^{-6}$  an arbitrary prior probability  $P(C|X) = x$ , and specifies instead of 99/100 an arbitrary fraction  $y$  of bad widgets. Then discuss how the dependence on  $x$  and  $y$  corresponds - or fails to correspond - to human common sense.

**Answer:**

Okay, let's do this.

Evidence is defined by

$$\begin{aligned}
e(H|DX) &= 10 \log_{10} O(H|DX) \\
&= 10 \log_{10} \left[ \frac{P(H|DX)}{P(\bar{H}|DX)} \right] \\
&= 10 \log_{10} \left[ \frac{P(D|HX)P(H|X)}{P(D|\bar{H}X)P(\bar{H}|X)} \right] \\
&= e(H|X) + 10 \log_{10} \left[ \frac{P(D|HX)}{P(D|\bar{H}X)} \right].
\end{aligned}$$

So, for multiple mutually exclusive and exhaustive hypotheses,

$$e(H_i|DX) = e(H_i|X) + 10 \log_{10} \left[ \frac{P(D|H_iX)}{P(D|\sum_{j \neq i} H_jX)} \right],$$

and we can replace

$$P\left(D \middle| \sum_{j \neq i} H_jX\right) = P\left(D \sum_{k \neq i} H_k \middle| \sum_{j \neq i} H_jX\right) \quad (25)$$

$$= \sum_{k \neq i} \frac{P(DH_k \sum_{j \neq i} H_j|X)}{P(\sum_{j \neq i} H_j|X)} \quad (26)$$

$$= \frac{\sum_{k \neq i} P(DH_k|X)}{\sum_{j \neq i} P(H_j|X)} \quad (27)$$

$$= \frac{\sum_{k \neq i} P(D|H_kX)P(H_k|X)}{\sum_{k \neq i} P(H_k|X)}. \quad (28)$$

So, letting  $w_k^{(i)} = P(H_k|X)/\sum_{j \neq i} P(H_j|X)$ ,

$$e(H_i|DX) = e(H_i|X) + 10 \log_{10} \left[ \frac{P(D|H_iX)}{\sum_{k \neq i} w_k^{(i)} P(D|H_kX)} \right].$$

This is, we're comparing the likelihood of  $H_i$  with the sum of the likelihoods  $H_k$  for  $k \neq i$  weighted by their prior likelihood. So yeah, I guess we can keep the  $j$ th term of the denominator as long as it's 10 times bigger than the rest of them.

In the problem at hand,

$$\begin{aligned}
P(A|X) &= \frac{1}{11}(1-x), \\
P(B|X) &= \frac{10}{11}(1-x), \\
P(C|X) &= x.
\end{aligned}$$

Furthermore, if after  $m$  measurements we find  $fm$  of them are bad,

$$\begin{aligned}
P(D|AX) &= \binom{m}{m_b} \left(\frac{1}{3}\right)^{fm} \left(\frac{2}{3}\right)^{(1-f)m}, \\
P(D|BX) &= \binom{m}{m_b} \left(\frac{1}{6}\right)^{fm} \left(\frac{5}{6}\right)^{(1-f)m}, \\
P(D|CX) &= \binom{m}{m_b} y^{fm} (1-y)^{(1-f)m}.
\end{aligned}$$

This all results in

$$\begin{aligned}
e(C|DX) &= e(C|X) + 10 \log_{10} \left[ \frac{P(D|CX)}{w_A^{(C)} P(D|AX) + w_B^{(C)} P(D|AX)} \right] \\
&= 10 \log_{10} \left[ \frac{x}{1-x} \right] + 10 \log_{10} \left[ \frac{y^{fm} (1-y)^{(1-f)m}}{\frac{1}{11} \left(\frac{1}{3}\right)^{fm} \left(\frac{2}{3}\right)^{(1-f)m} + \frac{10}{11} \left(\frac{1}{6}\right)^{fm} \left(\frac{5}{6}\right)^{(1-f)m}} \right] \\
&= 10 \log_{10} \left[ \frac{x}{1-x} \right] + 10 \log_{10} \left[ \frac{\left( \left( \frac{y}{1-y} \right)^f (1-y) \right)^m}{\frac{1}{11} \left( \left( \frac{1}{2} \right)^f \frac{2}{3} \right)^m + \frac{10}{11} \left( \left( \frac{1}{5} \right)^f \frac{5}{6} \right)^m} \right].
\end{aligned}$$

I find the last form more enlightening for the problem at hand. All of those terms are of the form something<sup>m</sup>. For  $m \rightarrow \infty$ , the term with the higher base will govern the denominator. To see which, we note the quotient of them is

$$\left( \frac{5}{2} \right)^f \frac{4}{5},$$

so setting this to unity yields a critical value  $f = \log \left( \frac{5}{4} \right) / \log \left( \frac{5}{2} \right) \approx 0.24$ . Given that the term in parentheses is greater than one, this means hypothesis  $A$  governs when  $f$  goes above this value, and  $B$  governs otherwise. If we want the threshold for  $y > \frac{1}{3}$ , it is clear that we need to work in the range where hypothesis  $A$  is more likely. So the threshold of skepticism is given by the value  $f$  such that the quotient

$$\left( 2 \frac{y}{1-y} \right)^f \frac{3}{2} (1-y)$$

is equal to 1. This yields the exact value

$$f_t(x, y) = \frac{\log \left( \frac{2}{3} \frac{1}{1-y} \right)}{\log \left( 2 \frac{y}{1-y} \right)}.$$

For  $y = 99/100$ , this yields approximately 0.794155, in conflict with the value of 0.793951 given in the book. But the solution book I'm contrasting with in github gives this same value, 0.7941etc.

The general case for estimating this threshold is now obvious:

$$f_t = \frac{\log \left[ \frac{P(\text{good}|CX)}{P(\text{good}|AX)} \right]}{\log \left[ \frac{P(\text{bad}|CX)P(\text{good}|AX)}{P(\text{good}|CX)P(\text{bad}|AX)} \right]}.$$