

★ 5.4 Probabilistic analysis and further uses of indicator random variables

This advanced section further illustrates probabilistic analysis by way of four examples. The first determines the probability that in a room of k people, some pair shares the same birthday. The second example examines the random tossing of balls into bins. The third investigates “streaks” of consecutive heads in coin flipping. The final example analyzes a variant of the hiring problem in which you have to make decisions without actually interviewing all the candidates.

5.4.1 The birthday paradox

Our first example is the *birthday paradox*. How many people must there be in a room before there is a 50% chance that two of them were born on the same day of the year? The answer is surprisingly few. The paradox is that it is in fact far fewer than the number of days in a year, or even half the number of days in a year, as we shall see.

To answer this question, we index the people in the room with the integers $1, 2, \dots, k$, where k is the number of people in the room. We ignore the issue of leap years and assume that all years have $n = 365$ days. For $i = 1, 2, \dots, k$, let b_i be the day of the year on which person i 's birthday falls, where $1 \leq b_i \leq n$. We also assume that birthdays are uniformly distributed across the n days of the year, so that $\Pr\{b_i = r\} = 1/n$ for $i = 1, 2, \dots, k$ and $r = 1, 2, \dots, n$.

The probability that two given people, say i and j , have matching birthdays depends on whether the random selection of birthdays is independent. We assume from now on that birthdays are independent, so that the probability that i 's birthday and j 's birthday both fall on day r is

$$\begin{aligned}\Pr\{b_i = r \text{ and } b_j = r\} &= \Pr\{b_i = r\} \Pr\{b_j = r\} \\ &= 1/n^2.\end{aligned}$$

Thus, the probability that they both fall on the same day is

$$\begin{aligned}\Pr\{b_i = b_j\} &= \sum_{r=1}^n \Pr\{b_i = r \text{ and } b_j = r\} \\ &= \sum_{r=1}^n (1/n^2) \\ &= 1/n.\end{aligned}\tag{5.7}$$

More intuitively, once b_i is chosen, the probability that b_j is chosen to be the same day is $1/n$. Thus, the probability that i and j have the same birthday is the same as the probability that the birthday of one of them falls on a given day. Notice,

however, that this coincidence depends on the assumption that the birthdays are independent.

We can analyze the probability of at least 2 out of k people having matching birthdays by looking at the complementary event. The probability that at least two of the birthdays match is 1 minus the probability that all the birthdays are different. The event that k people have distinct birthdays is

$$B_k = \bigcap_{i=1}^k A_i ,$$

where A_i is the event that person i 's birthday is different from person j 's for all $j < i$. Since we can write $B_k = A_k \cap B_{k-1}$, we obtain from equation (C.16) the recurrence

$$\Pr \{B_k\} = \Pr \{B_{k-1}\} \Pr \{A_k \mid B_{k-1}\} , \quad (5.8)$$

where we take $\Pr \{B_1\} = \Pr \{A_1\} = 1$ as an initial condition. In other words, the probability that b_1, b_2, \dots, b_k are distinct birthdays is the probability that b_1, b_2, \dots, b_{k-1} are distinct birthdays times the probability that $b_k \neq b_i$ for $i = 1, 2, \dots, k-1$, given that b_1, b_2, \dots, b_{k-1} are distinct.

If b_1, b_2, \dots, b_{k-1} are distinct, the conditional probability that $b_k \neq b_i$ for $i = 1, 2, \dots, k-1$ is $\Pr \{A_k \mid B_{k-1}\} = (n - k + 1)/n$, since out of the n days, there are $n - (k - 1)$ that are not taken. We iteratively apply the recurrence (5.8) to obtain

$$\begin{aligned} \Pr \{B_k\} &= \Pr \{B_{k-1}\} \Pr \{A_k \mid B_{k-1}\} \\ &= \Pr \{B_{k-2}\} \Pr \{A_{k-1} \mid B_{k-2}\} \Pr \{A_k \mid B_{k-1}\} \\ &\quad \vdots \\ &= \Pr \{B_1\} \Pr \{A_2 \mid B_1\} \Pr \{A_3 \mid B_2\} \cdots \Pr \{A_k \mid B_{k-1}\} \\ &= 1 \cdot \left(\frac{n-1}{n}\right) \left(\frac{n-2}{n}\right) \cdots \left(\frac{n-k+1}{n}\right) \\ &= 1 \cdot \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{k-1}{n}\right) . \end{aligned}$$

Inequality (3.11), $1 + x \leq e^x$, gives us

$$\begin{aligned} \Pr \{B_k\} &\leq e^{-1/n} e^{-2/n} \cdots e^{-(k-1)/n} \\ &= e^{-\sum_{i=1}^{k-1} i/n} \\ &= e^{-k(k-1)/2n} \\ &\leq 1/2 \end{aligned}$$

when $-k(k-1)/2n \leq \ln(1/2)$. The probability that all k birthdays are distinct is at most $1/2$ when $k(k-1) \geq 2n \ln 2$ or, solving the quadratic equation, when

$k \geq (1 + \sqrt{1 + (8 \ln 2)n})/2$. For $n = 365$, we must have $k \geq 23$. Thus, if at least 23 people are in a room, the probability is at least $1/2$ that at least two people have the same birthday. On Mars, a year is 669 Martian days long; it therefore takes 31 Martians to get the same effect.

An analysis using indicator random variables

We can use indicator random variables to provide a simpler but approximate analysis of the birthday paradox. For each pair (i, j) of the k people in the room, we define the indicator random variable X_{ij} , for $1 \leq i < j \leq k$, by

$$\begin{aligned} X_{ij} &= \mathbf{I}\{\text{person } i \text{ and person } j \text{ have the same birthday}\} \\ &= \begin{cases} 1 & \text{if person } i \text{ and person } j \text{ have the same birthday,} \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

By equation (5.7), the probability that two people have matching birthdays is $1/n$, and thus by Lemma 5.1, we have

$$\begin{aligned} \mathbf{E}[X_{ij}] &= \Pr\{\text{person } i \text{ and person } j \text{ have the same birthday}\} \\ &= 1/n. \end{aligned}$$

Letting X be the random variable that counts the number of pairs of individuals having the same birthday, we have

$$X = \sum_{i=1}^k \sum_{j=i+1}^k X_{ij}.$$

Taking expectations of both sides and applying linearity of expectation, we obtain

$$\begin{aligned} \mathbf{E}[X] &= \mathbf{E}\left[\sum_{i=1}^k \sum_{j=i+1}^k X_{ij}\right] \\ &= \sum_{i=1}^k \sum_{j=i+1}^k \mathbf{E}[X_{ij}] \\ &= \binom{k}{2} \frac{1}{n} \\ &= \frac{k(k-1)}{2n}. \end{aligned}$$

When $k(k-1) \geq 2n$, therefore, the expected number of pairs of people with the same birthday is at least 1. Thus, if we have at least $\sqrt{2n} + 1$ individuals in a room, we can expect at least two to have the same birthday. For $n = 365$, if $k = 28$, the expected number of pairs with the same birthday is $(28 \cdot 27)/(2 \cdot 365) \approx 1.0356$.

Thus, with at least 28 people, we expect to find at least one matching pair of birthdays. On Mars, where a year is 669 Martian days long, we need at least 38 Martians.

The first analysis, which used only probabilities, determined the number of people required for the probability to exceed $1/2$ that a matching pair of birthdays exists, and the second analysis, which used indicator random variables, determined the number such that the expected number of matching birthdays is 1. Although the exact numbers of people differ for the two situations, they are the same asymptotically: $\Theta(\sqrt{n})$.

5.4.2 Balls and bins

Consider the process of randomly tossing identical balls into b bins, numbered $1, 2, \dots, b$. The tosses are independent, and on each toss the ball is equally likely to end up in any bin. The probability that a tossed ball lands in any given bin is $1/b$. Thus, the ball-tossing process is a sequence of Bernoulli trials (see Appendix C.4) with a probability $1/b$ of success, where success means that the ball falls in the given bin. This model is particularly useful for analyzing hashing (see Chapter 11), and we can answer a variety of interesting questions about the ball-tossing process. (Problem C-1 asks additional questions about balls and bins.)

How many balls fall in a given bin? The number of balls that fall in a given bin follows the binomial distribution $b(k; n, 1/b)$. If n balls are tossed, equation (C.36) tells us that the expected number of balls that fall in the given bin is n/b .

How many balls must one toss, on the average, until a given bin contains a ball? The number of tosses until the given bin receives a ball follows the geometric distribution with probability $1/b$ and, by equation (C.31), the expected number of tosses until success is $1/(1/b) = b$.

How many balls must one toss until every bin contains at least one ball? Let us call a toss in which a ball falls into an empty bin a “hit.” We want to know the expected number n of tosses required to get b hits.

The hits can be used to partition the n tosses into stages. The i th stage consists of the tosses after the $(i - 1)$ st hit until the i th hit. The first stage consists of the first toss, since we are guaranteed to have a hit when all bins are empty. For each toss during the i th stage, there are $i - 1$ bins that contain balls and $b - i + 1$ empty bins. Thus, for each toss in the i th stage, the probability of obtaining a hit is $(b - i + 1)/b$.

Let n_i denote the number of tosses in the i th stage. Thus, the number of tosses required to get b hits is $n = \sum_{i=1}^b n_i$. Each random variable n_i has a geometric distribution with probability of success $(b - i + 1)/b$ and, by equation (C.31),

$$E[n_i] = \frac{b}{b - i + 1}.$$