

the following recursive procedure returns a random  $m$ -subset  $S$  of  $\{1, 2, 3, \dots, n\}$ , in which each  $m$ -subset is equally likely, while making only  $m$  calls to RANDOM:

```

RANDOM-SAMPLE( $m, n$ )
1  if  $m == 0$ 
2      return  $\emptyset$ 
3  else  $S = \text{RANDOM-SAMPLE}(m - 1, n - 1)$ 
4       $i = \text{RANDOM}(1, n)$ 
5      if  $i \in S$ 
6           $S = S \cup \{n\}$ 
7      else  $S = S \cup \{i\}$ 
8      return  $S$ 

```

---

## ★ 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, two of them share the same birthday. The second example examines what happens when we randomly toss balls into bins. The third investigates “streaks” of consecutive heads when we flip coins. 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

$$\begin{aligned}\Pr\{b_i = b_j\} &= \sum_{r=1}^n \Pr\{b_i = r \text{ and } b_j = r\} \\ &= \sum_{r=1}^n (1/n^2) \\ &= 1/n.\end{aligned}\tag{5.6}$$

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

$$\Pr\{B_k\} = \Pr\{B_{k-1}\} \Pr\{A_k \mid B_{k-1}\},\tag{5.7}$$

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,  $n - (k - 1)$  days are not taken. We iteratively apply the recurrence (5.7) to obtain

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

Inequality (3.12),  $1 + x \leq e^x$ , 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 (5.6), the probability that two people have matching birthdays is  $1/n$ , and thus by Lemma 5.1, we have

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

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

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

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 a process in which we randomly toss 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 we toss  $n$  balls, equation (C.37) tells us that the expected number of balls that fall in the given bin is  $n/b$ .

*How many balls must we 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.32), the expected number of tosses until success is  $1/(1/b) = b$ .

*How many balls must we 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.

Using the hits, we can 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,  $i - 1$  bins contain balls and  $b - i + 1$  bins are empty. 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 thus, by equation (C.32), we have

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

By linearity of expectation, we have

$$\begin{aligned} E[n] &= E\left[\sum_{i=1}^b n_i\right] \\ &= \sum_{i=1}^b E[n_i] \\ &= \sum_{i=1}^b \frac{b}{b - i + 1} \\ &= b \sum_{i=1}^b \frac{1}{i} \\ &= b(\ln b + O(1)) \quad (\text{by equation (A.7)}) . \end{aligned}$$

It therefore takes approximately  $b \ln b$  tosses before we can expect that every bin has a ball. This problem is also known as the ***coupon collector's problem***, which says that a person trying to collect each of  $b$  different coupons expects to acquire approximately  $b \ln b$  randomly obtained coupons in order to succeed.

### 5.4.3 Streaks

Suppose you flip a fair coin  $n$  times. What is the longest streak of consecutive heads that you expect to see? The answer is  $\Theta(\lg n)$ , as the following analysis shows.

We first prove that the expected length of the longest streak of heads is  $O(\lg n)$ . The probability that each coin flip is a head is  $1/2$ . Let  $A_{i,k}$  be the event that a streak of heads of length at least  $k$  begins with the  $i$ th coin flip or, more precisely, the event that the  $k$  consecutive coin flips  $i, i+1, \dots, i+k-1$  yield only heads, where  $1 \leq k \leq n$  and  $1 \leq i \leq n-k+1$ . Since coin flips are mutually independent, for any given event  $A_{i,k}$ , the probability that all  $k$  flips are heads is

$$\Pr\{A_{i,k}\} = 1/2^k. \quad (5.8)$$

For  $k = 2 \lceil \lg n \rceil$ ,

$$\begin{aligned} \Pr\{A_{i,2\lceil \lg n \rceil}\} &= 1/2^{2\lceil \lg n \rceil} \\ &\leq 1/2^{2\lg n} \\ &= 1/n^2, \end{aligned}$$

and thus the probability that a streak of heads of length at least  $2 \lceil \lg n \rceil$  begins in position  $i$  is quite small. There are at most  $n - 2 \lceil \lg n \rceil + 1$  positions where such a streak can begin. The probability that a streak of heads of length at least  $2 \lceil \lg n \rceil$  begins anywhere is therefore

$$\begin{aligned} \Pr\left\{\bigcup_{i=1}^{n-2\lceil \lg n \rceil+1} A_{i,2\lceil \lg n \rceil}\right\} &\leq \sum_{i=1}^{n-2\lceil \lg n \rceil+1} 1/n^2 \\ &< \sum_{i=1}^n 1/n^2 \\ &= 1/n, \end{aligned} \quad (5.9)$$

since by Boole's inequality (C.19), the probability of a union of events is at most the sum of the probabilities of the individual events. (Note that Boole's inequality holds even for events such as these that are not independent.)

We now use inequality (5.9) to bound the length of the longest streak. For  $j = 0, 1, 2, \dots, n$ , let  $L_j$  be the event that the longest streak of heads has length exactly  $j$ , and let  $L$  be the length of the longest streak. By the definition of expected value, we have

$$E[L] = \sum_{j=0}^n j \Pr\{L_j\}. \quad (5.10)$$

We could try to evaluate this sum using upper bounds on each  $\Pr\{L_j\}$  similar to those computed in inequality (5.9). Unfortunately, this method would yield weak bounds. We can use some intuition gained by the above analysis to obtain a good bound, however. Informally, we observe that for no individual term in the summation in equation (5.10) are both the factors  $j$  and  $\Pr\{L_j\}$  large. Why? When  $j \geq 2 \lceil \lg n \rceil$ , then  $\Pr\{L_j\}$  is very small, and when  $j < 2 \lceil \lg n \rceil$ , then  $j$  is fairly small. More formally, we note that the events  $L_j$  for  $j = 0, 1, \dots, n$  are disjoint, and so the probability that a streak of heads of length at least  $2 \lceil \lg n \rceil$  begins anywhere is  $\sum_{j=2 \lceil \lg n \rceil}^n \Pr\{L_j\}$ . By inequality (5.9), we have  $\sum_{j=2 \lceil \lg n \rceil}^n \Pr\{L_j\} < 1/n$ . Also, noting that  $\sum_{j=0}^n \Pr\{L_j\} = 1$ , we have that  $\sum_{j=0}^{2 \lceil \lg n \rceil - 1} \Pr\{L_j\} \leq 1$ . Thus, we obtain

$$\begin{aligned}
 E[L] &= \sum_{j=0}^n j \Pr\{L_j\} \\
 &= \sum_{j=0}^{2 \lceil \lg n \rceil - 1} j \Pr\{L_j\} + \sum_{j=2 \lceil \lg n \rceil}^n j \Pr\{L_j\} \\
 &< \sum_{j=0}^{2 \lceil \lg n \rceil - 1} (2 \lceil \lg n \rceil) \Pr\{L_j\} + \sum_{j=2 \lceil \lg n \rceil}^n n \Pr\{L_j\} \\
 &= 2 \lceil \lg n \rceil \sum_{j=0}^{2 \lceil \lg n \rceil - 1} \Pr\{L_j\} + n \sum_{j=2 \lceil \lg n \rceil}^n \Pr\{L_j\} \\
 &< 2 \lceil \lg n \rceil \cdot 1 + n \cdot (1/n) \\
 &= O(\lg n) .
 \end{aligned}$$

The probability that a streak of heads exceeds  $r \lceil \lg n \rceil$  flips diminishes quickly with  $r$ . For  $r \geq 1$ , the probability that a streak of at least  $r \lceil \lg n \rceil$  heads starts in position  $i$  is

$$\begin{aligned}
 \Pr\{A_{i, r \lceil \lg n \rceil}\} &= 1/2^{r \lceil \lg n \rceil} \\
 &\leq 1/n^r .
 \end{aligned}$$

Thus, the probability is at most  $n/n^r = 1/n^{r-1}$  that the longest streak is at least  $r \lceil \lg n \rceil$ , or equivalently, the probability is at least  $1 - 1/n^{r-1}$  that the longest streak has length less than  $r \lceil \lg n \rceil$ .

As an example, for  $n = 1000$  coin flips, the probability of having a streak of at least  $2 \lceil \lg n \rceil = 20$  heads is at most  $1/n = 1/1000$ . The chance of having a streak longer than  $3 \lceil \lg n \rceil = 30$  heads is at most  $1/n^2 = 1/1,000,000$ .

We now prove a complementary lower bound: the expected length of the longest streak of heads in  $n$  coin flips is  $\Omega(\lg n)$ . To prove this bound, we look for streaks

of length  $s$  by partitioning the  $n$  flips into approximately  $n/s$  groups of  $s$  flips each. If we choose  $s = \lfloor (\lg n)/2 \rfloor$ , we can show that it is likely that at least one of these groups comes up all heads, and hence it is likely that the longest streak has length at least  $s = \Omega(\lg n)$ . We then show that the longest streak has expected length  $\Omega(\lg n)$ .

We partition the  $n$  coin flips into at least  $\lfloor n / \lfloor (\lg n)/2 \rfloor \rfloor$  groups of  $\lfloor (\lg n)/2 \rfloor$  consecutive flips, and we bound the probability that no group comes up all heads. By equation (5.8), the probability that the group starting in position  $i$  comes up all heads is

$$\begin{aligned} \Pr \{A_{i, \lfloor (\lg n)/2 \rfloor}\} &= 1/2^{\lfloor (\lg n)/2 \rfloor} \\ &\geq 1/\sqrt{n} . \end{aligned}$$

The probability that a streak of heads of length at least  $\lfloor (\lg n)/2 \rfloor$  does not begin in position  $i$  is therefore at most  $1 - 1/\sqrt{n}$ . Since the  $\lfloor n / \lfloor (\lg n)/2 \rfloor \rfloor$  groups are formed from mutually exclusive, independent coin flips, the probability that every one of these groups *fails* to be a streak of length  $\lfloor (\lg n)/2 \rfloor$  is at most

$$\begin{aligned} (1 - 1/\sqrt{n})^{\lfloor n / \lfloor (\lg n)/2 \rfloor \rfloor} &\leq (1 - 1/\sqrt{n})^{n / \lfloor (\lg n)/2 \rfloor - 1} \\ &\leq (1 - 1/\sqrt{n})^{2n / \lg n - 1} \\ &\leq e^{-(2n / \lg n - 1) / \sqrt{n}} \\ &= O(e^{-\lg n}) \\ &= O(1/n) . \end{aligned}$$

For this argument, we used inequality (3.12),  $1 + x \leq e^x$ , and the fact, which you might want to verify, that  $(2n / \lg n - 1) / \sqrt{n} \geq \lg n$  for sufficiently large  $n$ .

Thus, the probability that the longest streak exceeds  $\lfloor (\lg n)/2 \rfloor$  is

$$\sum_{j=\lfloor (\lg n)/2 \rfloor + 1}^n \Pr \{L_j\} \geq 1 - O(1/n) . \quad (5.11)$$

We can now calculate a lower bound on the expected length of the longest streak, beginning with equation (5.10) and proceeding in a manner similar to our analysis of the upper bound:



$$\begin{aligned}
E[L] &= \sum_{j=0}^n j \Pr\{L_j\} \\
&= \sum_{j=0}^{\lfloor (\lg n)/2 \rfloor} j \Pr\{L_j\} + \sum_{j=\lfloor (\lg n)/2 \rfloor + 1}^n j \Pr\{L_j\} \\
&\geq \sum_{j=0}^{\lfloor (\lg n)/2 \rfloor} 0 \cdot \Pr\{L_j\} + \sum_{j=\lfloor (\lg n)/2 \rfloor + 1}^n \lfloor (\lg n)/2 \rfloor \Pr\{L_j\} \\
&= 0 \cdot \sum_{j=0}^{\lfloor (\lg n)/2 \rfloor} \Pr\{L_j\} + \lfloor (\lg n)/2 \rfloor \sum_{j=\lfloor (\lg n)/2 \rfloor + 1}^n \Pr\{L_j\} \\
&\geq 0 + \lfloor (\lg n)/2 \rfloor (1 - O(1/n)) \quad (\text{by inequality (5.11)}) \\
&= \Omega(\lg n).
\end{aligned}$$

As with the birthday paradox, we can obtain a simpler but approximate analysis using indicator random variables. We let  $X_{ik} = I\{A_{ik}\}$  be the indicator random variable associated with a streak of heads of length at least  $k$  beginning with the  $i$ th coin flip. To count the total number of such streaks, we define

$$X = \sum_{i=1}^{n-k+1} X_{ik}.$$

Taking expectations and using linearity of expectation, we have

$$\begin{aligned}
E[X] &= E\left[\sum_{i=1}^{n-k+1} X_{ik}\right] \\
&= \sum_{i=1}^{n-k+1} E[X_{ik}] \\
&= \sum_{i=1}^{n-k+1} \Pr\{A_{ik}\} \\
&= \sum_{i=1}^{n-k+1} 1/2^k \\
&= \frac{n-k+1}{2^k}.
\end{aligned}$$

By plugging in various values for  $k$ , we can calculate the expected number of streaks of length  $k$ . If this number is large (much greater than 1), then we expect many streaks of length  $k$  to occur and the probability that one occurs is high. If

this number is small (much less than 1), then we expect few streaks of length  $k$  to occur and the probability that one occurs is low. If  $k = c \lg n$ , for some positive constant  $c$ , we obtain

$$\begin{aligned}
 E[X] &= \frac{n - c \lg n + 1}{2^{c \lg n}} \\
 &= \frac{n - c \lg n + 1}{n^c} \\
 &= \frac{1}{n^{c-1}} - \frac{(c \lg n - 1)/n}{n^{c-1}} \\
 &= \Theta(1/n^{c-1}).
 \end{aligned}$$

If  $c$  is large, the expected number of streaks of length  $c \lg n$  is small, and we conclude that they are unlikely to occur. On the other hand, if  $c = 1/2$ , then we obtain  $E[X] = \Theta(1/n^{1/2-1}) = \Theta(n^{1/2})$ , and we expect that there are a large number of streaks of length  $(1/2) \lg n$ . Therefore, one streak of such a length is likely to occur. From these rough estimates alone, we can conclude that the expected length of the longest streak is  $\Theta(\lg n)$ .

#### 5.4.4 The on-line hiring problem

As a final example, we consider a variant of the hiring problem. Suppose now that we do not wish to interview all the candidates in order to find the best one. We also do not wish to hire and fire as we find better and better applicants. Instead, we are willing to settle for a candidate who is close to the best, in exchange for hiring exactly once. We must obey one company requirement: after each interview we must either immediately offer the position to the applicant or immediately reject the applicant. What is the trade-off between minimizing the amount of interviewing and maximizing the quality of the candidate hired?

We can model this problem in the following way. After meeting an applicant, we are able to give each one a score; let  $score(i)$  denote the score we give to the  $i$ th applicant, and assume that no two applicants receive the same score. After we have seen  $j$  applicants, we know which of the  $j$  has the highest score, but we do not know whether any of the remaining  $n - j$  applicants will receive a higher score. We decide to adopt the strategy of selecting a positive integer  $k < n$ , interviewing and then rejecting the first  $k$  applicants, and hiring the first applicant thereafter who has a higher score than all preceding applicants. If it turns out that the best-qualified applicant was among the first  $k$  interviewed, then we hire the  $n$ th applicant. We formalize this strategy in the procedure `ON-LINE-MAXIMUM( $k, n$ )`, which returns the index of the candidate we wish to hire.

ON-LINE-MAXIMUM( $k, n$ )

```

1  bestscore =  $-\infty$ 
2  for  $i = 1$  to  $k$ 
3      if  $\text{score}(i) > \text{bestscore}$ 
4           $\text{bestscore} = \text{score}(i)$ 
5  for  $i = k + 1$  to  $n$ 
6      if  $\text{score}(i) > \text{bestscore}$ 
7          return  $i$ 
8  return  $n$ 

```

We wish to determine, for each possible value of  $k$ , the probability that we hire the most qualified applicant. We then choose the best possible  $k$ , and implement the strategy with that value. For the moment, assume that  $k$  is fixed. Let  $M(j) = \max_{1 \leq i \leq j} \{\text{score}(i)\}$  denote the maximum score among applicants 1 through  $j$ . Let  $S$  be the event that we succeed in choosing the best-qualified applicant, and let  $S_i$  be the event that we succeed when the best-qualified applicant is the  $i$ th one interviewed. Since the various  $S_i$  are disjoint, we have that  $\Pr\{S\} = \sum_{i=1}^n \Pr\{S_i\}$ . Noting that we never succeed when the best-qualified applicant is one of the first  $k$ , we have that  $\Pr\{S_i\} = 0$  for  $i = 1, 2, \dots, k$ . Thus, we obtain

$$\Pr\{S\} = \sum_{i=k+1}^n \Pr\{S_i\} . \quad (5.12)$$

We now compute  $\Pr\{S_i\}$ . In order to succeed when the best-qualified applicant is the  $i$ th one, two things must happen. First, the best-qualified applicant must be in position  $i$ , an event which we denote by  $B_i$ . Second, the algorithm must not select any of the applicants in positions  $k + 1$  through  $i - 1$ , which happens only if, for each  $j$  such that  $k + 1 \leq j \leq i - 1$ , we find that  $\text{score}(j) < \text{bestscore}$  in line 6. (Because scores are unique, we can ignore the possibility of  $\text{score}(j) = \text{bestscore}$ .) In other words, all of the values  $\text{score}(k + 1)$  through  $\text{score}(i - 1)$  must be less than  $M(k)$ ; if any are greater than  $M(k)$ , we instead return the index of the first one that is greater. We use  $O_i$  to denote the event that none of the applicants in position  $k + 1$  through  $i - 1$  are chosen. Fortunately, the two events  $B_i$  and  $O_i$  are independent. The event  $O_i$  depends only on the relative ordering of the values in positions 1 through  $i - 1$ , whereas  $B_i$  depends only on whether the value in position  $i$  is greater than the values in all other positions. The ordering of the values in positions 1 through  $i - 1$  does not affect whether the value in position  $i$  is greater than all of them, and the value in position  $i$  does not affect the ordering of the values in positions 1 through  $i - 1$ . Thus we can apply equation (C.15) to obtain

$$\Pr\{S_i\} = \Pr\{B_i \cap O_i\} = \Pr\{B_i\} \Pr\{O_i\} .$$

The probability  $\Pr\{B_i\}$  is clearly  $1/n$ , since the maximum is equally likely to be in any one of the  $n$  positions. For event  $O_i$  to occur, the maximum value in positions 1 through  $i-1$ , which is equally likely to be in any of these  $i-1$  positions, must be in one of the first  $k$  positions. Consequently,  $\Pr\{O_i\} = k/(i-1)$  and  $\Pr\{S_i\} = k/(n(i-1))$ . Using equation (5.12), we have

$$\begin{aligned} \Pr\{S\} &= \sum_{i=k+1}^n \Pr\{S_i\} \\ &= \sum_{i=k+1}^n \frac{k}{n(i-1)} \\ &= \frac{k}{n} \sum_{i=k+1}^n \frac{1}{i-1} \\ &= \frac{k}{n} \sum_{i=k}^{n-1} \frac{1}{i} . \end{aligned}$$

We approximate by integrals to bound this summation from above and below. By the inequalities (A.12), we have

$$\int_k^n \frac{1}{x} dx \leq \sum_{i=k}^{n-1} \frac{1}{i} \leq \int_{k-1}^{n-1} \frac{1}{x} dx .$$

Evaluating these definite integrals gives us the bounds

$$\frac{k}{n}(\ln n - \ln k) \leq \Pr\{S\} \leq \frac{k}{n}(\ln(n-1) - \ln(k-1)) ,$$

which provide a rather tight bound for  $\Pr\{S\}$ . Because we wish to maximize our probability of success, let us focus on choosing the value of  $k$  that maximizes the lower bound on  $\Pr\{S\}$ . (Besides, the lower-bound expression is easier to maximize than the upper-bound expression.) Differentiating the expression  $(k/n)(\ln n - \ln k)$  with respect to  $k$ , we obtain

$$\frac{1}{n}(\ln n - \ln k - 1) .$$

Setting this derivative equal to 0, we see that we maximize the lower bound on the probability when  $\ln k = \ln n - 1 = \ln(n/e)$  or, equivalently, when  $k = n/e$ . Thus, if we implement our strategy with  $k = n/e$ , we succeed in hiring our best-qualified applicant with probability at least  $1/e$ .

## Exercises

### 5.4-1

How many people must there be in a room before the probability that someone has the same birthday as you do is at least  $1/2$ ? How many people must there be before the probability that at least two people have a birthday on July 4 is greater than  $1/2$ ?

### 5.4-2

Suppose that we toss balls into  $b$  bins until some bin contains two balls. Each toss is independent, and each ball is equally likely to end up in any bin. What is the expected number of ball tosses?

### 5.4-3 ★

For the analysis of the birthday paradox, is it important that the birthdays be mutually independent, or is pairwise independence sufficient? Justify your answer.

### 5.4-4 ★

How many people should be invited to a party in order to make it likely that there are *three* people with the same birthday?

### 5.4-5 ★

What is the probability that a  $k$ -string over a set of size  $n$  forms a  $k$ -permutation? How does this question relate to the birthday paradox?

### 5.4-6 ★

Suppose that  $n$  balls are tossed into  $n$  bins, where each toss is independent and the ball is equally likely to end up in any bin. What is the expected number of empty bins? What is the expected number of bins with exactly one ball?

### 5.4-7 ★

Sharpen the lower bound on streak length by showing that in  $n$  flips of a fair coin, the probability is less than  $1/n$  that no streak longer than  $\lg n - 2 \lg \lg n$  consecutive heads occurs.

---

## Problems

### 5-1 Probabilistic counting

With a  $b$ -bit counter, we can ordinarily only count up to  $2^b - 1$ . With R. Morris's *probabilistic counting*, we can count up to a much larger value at the expense of some loss of precision.

We let a counter value of  $i$  represent a count of  $n_i$  for  $i = 0, 1, \dots, 2^b - 1$ , where the  $n_i$  form an increasing sequence of nonnegative values. We assume that the initial value of the counter is 0, representing a count of  $n_0 = 0$ . The INCREMENT operation works on a counter containing the value  $i$  in a probabilistic manner. If  $i = 2^b - 1$ , then the operation reports an overflow error. Otherwise, the INCREMENT operation increases the counter by 1 with probability  $1/(n_{i+1} - n_i)$ , and it leaves the counter unchanged with probability  $1 - 1/(n_{i+1} - n_i)$ .

If we select  $n_i = i$  for all  $i \geq 0$ , then the counter is an ordinary one. More interesting situations arise if we select, say,  $n_i = 2^{i-1}$  for  $i > 0$  or  $n_i = F_i$  (the  $i$ th Fibonacci number—see Section 3.2).

For this problem, assume that  $n_{2^b-1}$  is large enough that the probability of an overflow error is negligible.

- a. Show that the expected value represented by the counter after  $n$  INCREMENT operations have been performed is exactly  $n$ .
- b. The analysis of the variance of the count represented by the counter depends on the sequence of the  $n_i$ . Let us consider a simple case:  $n_i = 100i$  for all  $i \geq 0$ . Estimate the variance in the value represented by the register after  $n$  INCREMENT operations have been performed.

### 5-2 Searching an unsorted array

This problem examines three algorithms for searching for a value  $x$  in an unsorted array  $A$  consisting of  $n$  elements.

Consider the following randomized strategy: pick a random index  $i$  into  $A$ . If  $A[i] = x$ , then we terminate; otherwise, we continue the search by picking a new random index into  $A$ . We continue picking random indices into  $A$  until we find an index  $j$  such that  $A[j] = x$  or until we have checked every element of  $A$ . Note that we pick from the whole set of indices each time, so that we may examine a given element more than once.

- a. Write pseudocode for a procedure RANDOM-SEARCH to implement the strategy above. Be sure that your algorithm terminates when all indices into  $A$  have been picked.