

Collaborators

Ben Nelson, Corbin Baldwin and I collaborated for this assignment.

Question 1: Balls and bins..... [12]

Consider the process we saw in class, of distributing n balls into n bins, independently and uniformly at random.

- (a) [4] What is the expected number of empty bins?
- (b) [2] Show that the probability that 80% of the bins are empty is $\leq 1/2$.
- (c) [2] Let X_i be the random variable that is 1 if bin i is empty and is 0 otherwise. Are X_1 and X_2 independent? Give an intuitive reasoning.
- (d) [4] Let Z be the random variable defined by $Z = \sum_i X_i$ (i.e., the number of empty bins). Suppose you are told that the variance of Z is $\leq n$. Use this to prove that the probability that 80% of the bins are empty is $< 8/n$. [Note that this is a much better bound than the one in part (b) when n is large.]

(a)

Let random variable $X_i = \begin{cases} 1, & \text{when slot } i \text{ is empty} \\ 0, & \text{otherwise} \end{cases}$

$$\text{Number of empty bins} = \sum_{i=1}^n X_i$$

$$\begin{aligned} \mathbb{E}[\text{Number of empty bins}] &= \mathbb{E}\left[\sum_{i=1}^n X_i\right] \\ &= \sum_{i=1}^n \mathbb{E}[X_i], \text{ due to linearity of expectation} \\ &= \sum_{i=1}^n (1 \times \Pr[X_i = 1] + 0 \times \Pr[X_i = 0]) \\ &= \sum_{i=1}^n \Pr[X_i = 1] \\ &= \sum_{i=1}^n \frac{(n-1)^n}{n^n} \end{aligned}$$

since n balls can be put into $n-1$ slots in $(n-1)^n$ ways

and n balls can be put into n slots in n^n different ways

$$\begin{aligned} &= n \frac{(n-1)^n}{n^n} \\ &= \frac{(n-1)^n}{n^{n-1}} \end{aligned}$$

(b)

$$\begin{aligned}
\Pr[80\% \text{ of the bins are empty}] &= \Pr[\text{Number of empty bins} \geq 0.8n] \\
&= \Pr[X_e \geq 0.8n], \text{ where } X_e = \text{number of empty bins} \\
&= \Pr\left[X_e \geq \frac{0.8n}{\mathbb{E}[X_e]} \mathbb{E}[X_e]\right] \\
&\leq \frac{\mathbb{E}[X_e]}{0.8n}, \text{ by Markov's Inequality} \\
&= \frac{(n-1)^n}{n^{n-1}} \frac{1}{0.8n} \\
&= \frac{1.25(n-1)^n}{n^n}
\end{aligned}$$

$$\begin{aligned}
\text{Let } y &= \frac{1.25(n-1)^n}{n^n} \\
&= 1.25(n-1)^n n^{-n}
\end{aligned}$$

$$\begin{aligned}
\frac{dy}{dn} &= 1.25n(n-1)^{n-1}n^{-n} - 1.25n(n-1)^n n^{-n-1} \\
&= 1.25n(n-1)^{n-1}n^{-n} (1 - (n-1)n^{-1}) \\
&= 1.25n(n-1)^{n-1}n^{-n} \left(1 - \left(1 - \frac{1}{n}\right)\right) \\
&= 1.25n(n-1)^{n-1}n^{-n} \frac{1}{n} \\
&\geq 0, \text{ for } n > 1
\end{aligned}$$

$$\begin{aligned}
\lim_{n \rightarrow \infty} \frac{1.25(n-1)^n}{n^n} &= \lim_{n \rightarrow \infty} 1.25 \left(\frac{n-1}{n}\right)^n \\
&= 1.25 \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n \\
&= 1.25 \frac{1}{e} \\
&= 0.4598 \\
&< \frac{1}{2}
\end{aligned}$$

Since $\frac{dy}{dn} > 0$ for $n > 1$ and $\lim_{n \rightarrow \infty} y < \frac{1}{2}$, $\Pr[80\% \text{ of the bins are empty}] < \frac{1}{2}$.

- (c) If bin 1 is known to be empty, then there is a higher chance that bin 2 is not empty, since all of the n balls would have been assigned to the $n-1$ remaining buckets. This means that X_1 and X_2 are not independent random variables.

(d)

$$\begin{aligned}
\mathbb{E}[Z] &= \frac{(n-1)^n}{n^{n-1}}, \text{ from part a above} \\
&= n \frac{(n-1)^n}{n^n} \\
&= n \left(1 - \frac{1}{n}\right)^n \\
&< \frac{n}{e}, \text{ since } \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n = \frac{1}{e}
\end{aligned}$$

$$\begin{aligned}
\Pr[80\% \text{ of the bins are empty}] &= \Pr[\text{Number of empty bins} \geq 0.8n] \\
&= \Pr[Z \geq 0.8n] \\
&= \Pr[Z - \mathbb{E}[Z] \geq 0.8n - \mathbb{E}[Z]] \\
&= \Pr\left[Z - \mathbb{E}[Z] \geq \frac{0.8n - \mathbb{E}[Z]}{\sigma} \sigma\right], \text{ where } \sigma^2 \text{ is the variance of } Z \\
&\leq \frac{1}{\left(\frac{0.8n - \mathbb{E}[Z]}{\sigma}\right)^2}, \text{ by Chebychev's inequality} \\
&= \frac{\sigma^2}{(0.8n - \mathbb{E}[Z])^2} \\
&\leq \frac{n}{(0.8n - \mathbb{E}[Z])^2} \\
&\leq \frac{n}{\left(0.8n - \frac{n}{e}\right)^2} \\
&= \frac{ne^2}{n^2 (0.8e - 1)^2} \\
&= \frac{1}{n} \frac{e^2}{(0.8e - 1)^2} \\
&= \frac{5.355}{n} \\
&< \frac{8}{n} \quad \square
\end{aligned}$$

Question 2: Median/order finding.....[10]

Consider the following simple algorithm for finding the k th smallest element in a given unsorted array $A[1], A[2], \dots, A[n]$ (assume that all the $A[i]$ are distinct):

1. pick a uniformly random index $i \in \{1, 2, \dots, n\}$
2. divide the array $A[]$ into two sub-arrays – $B[]$ and $C[]$, where the elements of $B[]$ are all $\leq A[i]$ and the elements of $C[]$ are $> A[i]$.
3. if $\text{length}(B) \geq k$, we recursively find the k th smallest element in $B[]$; else we find the $k - \text{length}(B)$ 'th smallest element in $C[]$.

(a) [2] Prove that the running time on the array A of length n can be bounded as

$$T(n) \leq \max\{T(\text{length}(B)), T(\text{length}(C))\} + O(n).$$

(b) [3] Give an input A, k , and an unlucky choice of indices i that leads to a running time larger than $n^2/4$.

- (c) [5] Let $f(n)$ denote the *expected running time* on an input array of length n . Derive a “probabilistic recurrence” analogous to the one we saw in class for quicksort, and show that $f(n) = O(n)$. [Hint: use part (a).]
- (a) To divide the array into two sub-arrays, we would need to compare each element with the one corresponding to the position of the randomly chosen index. This will take $O(n)$ time. In the worst case, the element corresponding to the k^{th} smallest element in the original array, would lie in the larger of the two sub-arrays. This means that the running time on the array A of length n can be bounded as $T(n) \leq \max\{T(\text{length}(B)), T(\text{length}(C))\} + O(n)$.
- (b) Let $k = 7$ and $A = \{6, 2, 5, 4, 3, 1, 9\}$
 $B = \{1\}$ and $C = \{6, 2, 5, 4, 3, 9\}$ when $i = 6$
 Since $\text{length}(B) = 1 < k$, we need to find the $7 - 1 = 6^{th}$ smallest element in C
 $\{2\}$ and $\{6, 5, 4, 3, 9\}$ are the two sub-arrays when $i = 2$
 Since $\text{length}(\text{first sub-array}) = 1 < 6$, we need to find the $6 - 1 = 5^{th}$ smallest element in the second sub-array
 $\{3\}$ and $\{6, 5, 4, 9\}$ are the two sub-arrays when $i = 4$
 Since $\text{length}(\text{first sub-array}) = 1 < 5$, we need to find the $5 - 1 = 4^{th}$ smallest element in the second sub-array
 $\{4\}$ and $\{6, 5, 9\}$ are the two sub-arrays when $i = 3$
 Since $\text{length}(\text{first sub-array}) = 1 < 4$, we need to find the $4 - 1 = 3^{rd}$ smallest element in the second sub-array
 $\{5\}$ and $\{6, 9\}$ are the two sub-arrays when $i = 2$
 Since $\text{length}(\text{first sub-array}) = 1 < 3$, we need to find the $3 - 1 = 2^{nd}$ smallest element in the second sub-array
 $\{6\}$ and $\{9\}$ are the two sub-arrays when $i = 1$
 Since $\text{length}(\text{first sub-array}) = 1 < 2$, we need to find the $2 - 1 = 1^{st}$ smallest element in the second sub-array
 We get 9 as the answer, which is the $k = 7^{th}$ smallest element in the original array $A = \{6, 2, 5, 4, 3, 1, 9\}$
- If we add up the lengths of the arrays worked on, we get $7 + 6 + 5 + 4 + 3 + 2 = 27 > \frac{n^2}{4} = \frac{7^2}{4} = \frac{49}{4} = 12.25$.
- (c) Consider the base case of an array of length 1. $T(1) = O(1)$ since it takes constant time to find the 1^{st} smallest element in the array. Since the recurrence is bounded as

$$T(n) \leq \max\{T(\text{length}(B)), T(\text{length}(C))\} + O(n).$$

Consider the case where array B has a length of $n - r$ and array C has a length on r . The recurrence will be

$$T(n) \leq \max\{T(n - r), T(r)\} + O(n).$$

$$\begin{aligned}
\mathbb{E}[T(n)] &\leq \sum_{r=0}^n \Pr[r] (\max\{T(n-r), T(r)\} + O(n)), (r \text{ starts at } 0 \text{ when largest element is chosen}) \\
&= \sum_{r=0}^n \frac{1}{n} (\max\{T(n-r), T(r)\} + O(n)), \text{ since the probability of any } r \text{ is } \frac{1}{n} \\
&= \sum_{r=0}^n \frac{1}{n} (\max\{O(n-r), O(r)\} + O(n)), \text{ assuming } T(n-r) = O(n-r) \text{ and } T(r) = O(r) \\
&= \sum_{r=0}^n \frac{1}{n} (\max\{(n-r), r\} + n), \text{ taking the constant to be } 1 \\
&= \frac{1}{n} \left(2 \left(\frac{n(n+1)}{2} - \frac{\frac{n}{2}(\frac{n}{2}+1)}{2} \right) + n^2 \right), \text{ as we take the larger sub-array for worst case} \\
&= \frac{1}{n} \left(n^2 + n - \frac{n^2}{4} - \frac{n}{2} + n^2 \right) \\
&= \frac{1}{n} \left(\frac{7}{4}n^2 + \frac{n}{2} \right) \\
&= \frac{7}{4}n + \frac{1}{2} \\
&= O(n) \quad \square
\end{aligned}$$

Question 3: Quicksort revisited [6]

Recall the quicksort procedure we saw in class, where we pick a random element of the input array as the *pivot*, recursively sort the elements less than and larger than the pivot, and concatenate the solutions. We saw that the expected running time is $O(n \log n)$.

- [2] Using a basic implementation (base case being a singleton), find the *constant* in the $O()$ notation for the algorithm above. (You may do this by picking any array, repeatedly running the procedure above, and averaging the values of running time divided by $n \log n$.)
- [2] Now, consider the following procedure: for $k = 1, 2, 3$, first pick $(2k + 1)$ random indices, and choose their *median* as the pivot. Now report the constant in the $O()$ notation.
- [2] Explain your observations intuitively.
- The constant that was computed for various random integer array sizes for runtimes averaged over 20 iterations is shown below

Array Size	Constant
10^5	5.21×10^{-6}
10^6	3.47×10^{-6}
10^7	2.74×10^{-6}
10^8	2.69×10^{-6}

- The value of the constant for various array sizes is shown below

Array Size	$k = 1$	$k = 2$	$k = 3$
10^5	6.08×10^{-6}	6.08×10^{-6}	6.95×10^{-6}
10^6	4.85×10^{-6}	5.14×10^{-6}	5.57×10^{-6}
10^7	4.20×10^{-6}	4.52×10^{-6}	4.98×10^{-6}
10^8	3.63×10^{-6}	4.17×10^{-6}	4.43×10^{-6}

- The value of the constant is least for the case of the single random pivot and higher for the cases where a median of random elements is used as a pivot. It looks like the work involved in computing the median of elements negated the expected boost in performance from choosing a supposedly

better pivot. The more the elements that a pivot is chosen from, the higher the value of the constant and the run time. The best performance was from choosing a single random pivot as it involved minimum extra work and performance was good as it was improbable that many pivots would be bad and adversely impact run time.

Question 4: Checking matrix multiplication..... [12]

Matrix multiplication is one of the classic algorithmic problems. Consider the problem of multiplying two $n \times n$ matrices over the field \mathbf{F}_2 (i.e., we have matrices with entries 0/1, and we perform all computations modulo 2; e.g., $0*0 + 1*1 + 1*1 + 1*0 = 0$).

The best known algorithms here are messy and take time $O(n^{2.36\dots})$. However, the point of this exercise is to prove a simpler statement. Suppose someone gives a matrix C and claims that $C = AB$, can we *quickly verify* if the claim is true?

- (a) [4] First prove a warm-up statement: suppose a and b are any two 0/1 vectors of length n , and suppose that $a \neq b$. Then, for a random binary vector $x \in \{0, 1\}^n$ (one in which each coordinate is chosen uniformly at random), prove that $\Pr[\langle a, x \rangle \neq \langle b, x \rangle \pmod{2}] = 1/2$.
- (b) [6] Now, design an $O(n^2)$ time algorithm that tests if $C = AB$ and has a success probability $\geq 1/2$. (You need to prove the probability bound.)
- (c) [2] Show how to improve the success probability to $7/8$ while still having running time $O(n^2)$.
- (a) If a and b can be any vectors $\in \{0, 1\}^n$ and x is also $\in \{0, 1\}^n$ with each coordinate is chosen uniformly at random, then the probability of any term $a_i \times x_i$ or $b_i \times x_i$ where $0 \leq i \leq n$ in the inner product being 0 is $\frac{3}{4}$ corresponding to the pairs $\{00, 01, 10\}$. The probability of any term being 1 is $1 - \frac{3}{4} = \frac{1}{4}$.

I got the hint for the following recurrence from math.stackexchange.com. But I can't find the page now to list as a reference.

Let a_N be the probability that a product of two random binary vectors has an even number of 1's. This happens when for the first term to be summed in the inner product is 0 followed by an even number of 1's, or the other way round. We can see that $a_0 = 1$ and for $N \geq 1$

$$a_N = p_0 a_{N-1} + p_1 (1 - a_{N-1})$$

where p_0 is the probability of a term in the dot product being 0 and p_1 is the probability of a term in the dot product being 1. The solution to this recurrence can be shown to be equal to the following by induction

$$a_N = \frac{1}{2} + \frac{(p_0 - p_1)^N}{2}$$

For the base case where $N = 0$, we get $a_0 = \frac{1}{2} + \frac{(p_0 - p_1)^0}{2} = \frac{1}{2} + \frac{1}{2} = 1$. Let us assume that this is true for $N - 1$. Then

$$\begin{aligned}
a_N &= p_0 \left(\frac{1}{2} + \frac{(p_0 - p_1)^{N-1}}{2} \right) + p_1 \left(1 - \left(\frac{1}{2} + \frac{(p_0 - p_1)^{N-1}}{2} \right) \right) \\
&= \frac{p_0 + p_1}{2} + (p_0 - p_1) \frac{(p_0 - p_1)^{N-1}}{2} \\
&= \frac{1}{2} + \frac{(p_0 - p_1)^N}{2}, \text{ which proves the recurrence} \\
&= \frac{1}{2} + \frac{(\frac{3}{4} - \frac{1}{4})^N}{2}, \text{ by plugging in } p_0 = \frac{3}{4} \text{ and } p_1 = \frac{1}{4} \\
&= \frac{1}{2} + \frac{1^N}{2}
\end{aligned}$$

$$\begin{aligned}
\Pr[\langle a, x \rangle = \langle b, x \rangle \pmod{2}] &= \Pr[\text{both inner products terms have odd 1's or both have even 1's}] \\
&= a_N a_N + (1 - a_N)(1 - a_N) \\
&= a_N^2 + 1 - 2a_N + a_N^2 \\
&= 2a_N^2 - 2a_N + 1 = \frac{1}{2}, \text{ make this assumption for the time being}
\end{aligned}$$

$$\Rightarrow 2a_N^2 - 2a_N + \frac{1}{2} = 0$$

$$4a_N^2 - 4a_N + 1 = 0$$

$$(2a_N - 1)^2 = 0$$

$$a_N = \frac{1}{2}, \text{ which cannot be true by } a_N \text{ given above}$$

$$\Rightarrow \Pr[\langle a, x \rangle = \langle b, x \rangle \pmod{2}] \neq \frac{1}{2}$$

- (b) Although I was not able to prove that $\Pr[\langle a, x \rangle = \langle b, x \rangle \pmod{2}] = 1/2$. I will use this to be a known fact for the purposes of this part.

Using Freivalds' algorithm [1], create a random binary vector $x \in \{0, 1\}^n$ (one in which each coordinate is chosen uniformly at random). A and B are the two $n \times n$ matrices over the field \mathbf{F}_2 . Evaluate $P = A \times (Bx) - Cx$. If $P = \{0\}^n$, we output *true* signifying that $AB = C$, else we output *false*.

When $AB = C$, $\Pr[P = \{0\}^n] = 1$ since any value of x will result in $P = A \times (Bx) - Cx$.

When $AB \neq C$, for P to be $\{0\}^n$, we would need every element in $A \times (Bx)$ to be the same as the corresponding element in $Cx \pmod{2}$. Since $A \times (Bx) = (A \times B)x$, $\Pr[P = \{0\}^n] = (\frac{1}{2})^n$, since every element in vector $(A \times B)x$ would need to be the same as every element in vector $Cx \pmod{2}$. Each of these will be equal with probability $\frac{1}{2}$ as per part a above. And for all of them to be equal, it will be $(\frac{1}{2})^n$. So for the case when $n = 2$, $\Pr[P = \{0\}^n] = (\frac{1}{2})^2 = \frac{1}{4}$ or in other words the success probability is $1 - \frac{1}{4} = \frac{3}{4}$. In the general case the success probability is therefore $\geq \frac{1}{2}$.

To compute vector Bx , we would need to do n inner products and each inner product will take n work. So time complexity will be $O(n^2)$. Then to compute $A \times (Bx)$, where we multiply matrix A with vector Bx , the time complexity will again be $O(n^2)$. To compute $P = A \times (Bx) - Cx$ after this, we would need to do $O(n^2)$ work to compute Cx and then n subtractions. So total time complexity is $O(n^2)$.

- (c) To improve the success probability to $7/8$ while still having running time $O(n^2)$, repeat the algorithm described in part b three times and output *true* signifying that $AB = C$ only if $P = \{0\}^n$ each of these three times. So for the algorithm to output true when $AB \neq C$, the probability would be $\leq (\frac{1}{2})^3 = \frac{1}{8}$. The success probability will be $\geq 1 - \frac{1}{8} = \frac{7}{8}$. The time complexity will still be $O(n^2)$ even with repeating an $O(n^2)$ algorithm 3 times.

Question 5: Coloring graphs [10]

One of the popular problems in graph theory is that of graph coloring: we are given an undirected graph $G = (V, E)$, and the aim is to assign each vertex u a *color* $c(u)$, with the guarantee that for any edge u, v , $c(u) \neq c(v)$. (i.e., end-points of edges must receive different colors.) The goal is to color the graph using the fewest total number of colors, subject to the above condition.

Coloring is known to be a very hard problem. Suppose we are OK with something weaker: suppose that we consider an assignment of colors *acceptable* if $c(u) \neq c(v)$ holds for $\geq 90\%$ of the edges.

- (a) [7] Suppose that we *randomly* assign a color in the range $\{1, 2, \dots, 20\}$ to the vertices. Prove that we obtain an acceptable assignment with probability $\geq 1/2$. (Note that the algorithm is quite remarkable – it doesn't even look at the edges of the graph!)
- (b) [3] What happens above when we randomly assign colors in the range $\{1, 2, \dots, 11\}$? Can we still obtain an algorithm that succeeds with probability $\geq 1/2$?

(a)

Let random variable $E_i = \begin{cases} 1, & \text{when vertices at the two ends of edge } E_i \text{ are of different colors} \\ 0, & \text{otherwise} \end{cases}$

Let E_T be the number of edges where vertices at the two ends are of different colors

$$\begin{aligned} \mathbb{E}[E_T] &= \mathbb{E}\left[\sum_{i=1}^n E_i\right], \text{ where } n \text{ is the number of vertices} \\ &= \sum_{i=1}^n \mathbb{E}[E_i], \text{ by linearity of Expectation} \\ &= \sum_{i=1}^n (1 \times \Pr[E_i = 1] + 0 \times \Pr[E_i = 0]) \\ &= \sum_{i=1}^n \Pr[E_i = 1] \\ &= \sum_{i=1}^n \frac{19}{20}, \text{ as once a color is given to a vertex, one of 19 can be given to the other} \\ &= \frac{19}{20}n \\ &= 0.95n \end{aligned}$$

Since the expected value or mean of the number of edges where vertices at the two ends are of different colors is $0.95n$, the probability that the number of edges where vertices at the two ends are of different colors is $\geq 0.95n$ is $\frac{1}{2}$. This means that the probability that the number of edges where vertices at the two ends are of different colors is the $\geq 0.90n$ is definitely $\geq \frac{1}{2}$ since this covers more than half the probability mass.

- (b) In the case where we randomly assign colors in the range $\{1, 2, \dots, 11\}$, $\mathbb{E}[E_T] = \frac{10}{11}n = 0.91n$ based on a computation similar to the one shown above. So the the probability that the number of edges where vertices at the two ends are of different colors is the $\geq 0.90n$ is definitely $\geq \frac{1}{2}$ since this covers more than half the probability mass.

References

- [1] "Freivalds' Algorithm." *Wikipedia*, Wikimedia Foundation, 30 Sept. 2018, https://en.wikipedia.org/wiki/Freivalds\%27_algorithm.