



## Exercise Set VII, Advanced Algorithms 2022

These exercises are for your own benefit. Feel free to collaborate and share your answers with other students. Solve as many problems as you can and ask for help if you get stuck for too long. Problems marked \* are more difficult but also more fun :).

These problems are taken from various sources at EPFL and on the Internet, too numerous to cite individually.

- 1 Design a one-pass (streaming) algorithm that, for a stream that possesses a majority element (appearing more than  $m/2$  times), terminates with this element. Prove the correctness of your algorithm.

**Solution:** The algorithm is as follows: We store the first item and a counter, initialized to 1. For each subsequent item, if it is the same as the currently stored item, increment the counter. If it differs, and the counter is zero, then store the new item and set the counter to 1; else, decrement the counter.

Say that  $s$  is the unknown value that occurs more than  $m/2$  times. The idea of the algorithm is that if you could pair up elements of the stream so that distinct values are paired up, and then you “kill” these pairs, then  $s$  will always survive. The way this algorithm pairs up the values is by holding onto the most recent value that has no pair (implicitly, by keeping a count how many copies of that value you saw). Then when you come across a new element, you decrement the counter and implicitly account for one new pair.

- 2 (*MinHashing*) Suppose we have a universe  $U$  of elements. For  $A, B \subseteq U$ , the Jaccard distance of  $A, B$  is defined as

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}.$$

This definition is used in practice to calculate a notion of similarity of documents, webpages, etc. For example, suppose  $U$  is the set of English words, and any set  $A$  represents a document considered as a bag of words. Note that for any two  $A, B \subseteq U$ ,  $0 \leq J(A, B) \leq 1$ . If  $J(A, B)$  is close to 1, then we can say  $A \approx B$ .

Let  $h : U \rightarrow [0, 1]$  where for each  $i \in U$ ,  $h(i)$  is chosen uniformly and independently at random. For a set  $S \subseteq U$ , let  $h_S := \min_{i \in S} h(i)$ . **Show that**

$$\Pr[h_A = h_B] = J(A, B).$$

Now, if we have sets  $A_1, A_2, \dots, A_n$ , we can use the above idea to figure out which pair of sets are “close” in time essentially  $O(n|U|)$ . We can also obtain a good approximation of  $J(A, B)$  with high probability by using several independently chosen hash functions. Note that the naive algorithm would take  $O(n^2|U|)$  to calculate all pairwise similarities.

**Solution:** First, let us simplify the situation a little by noticing that with probability 1, all elements  $h(i)$  for  $i \in U$  are different. This is because  $\Pr[h(i) = h(j)] = 0$  for  $i \neq j$  (recall that each  $h(i)$  is uniform on the interval  $[0, 1]$ ).

Given this, let us see where  $\min_{i \in A \cup B} h(i)$  is attained:

- if it is attained in  $A \cap B$ , then  $h_A = h_B = h_{A \cup B} = h_{A \cap B}$ ,
- otherwise, say it is attained in  $A \setminus B$ : then  $h_A < h_B$ .

Therefore the event  $h_A = h_B$  is (almost everywhere) equal to  $h_{A \cup B} = h_{A \cap B}$ . Furthermore, notice that for any set  $S \subseteq U$  and any  $i \in S$  we have  $\Pr[h(i) = h_S] = 1/|S|$  due to symmetry. Therefore

$$\Pr[h_A = h_B] = \Pr[h_{A \cap B} = h_{A \cup B}] = \sum_{i \in A \cap B} \Pr[h(i) = h_{A \cup B}] = |A \cap B| \cdot \frac{1}{|A \cup B|} = J(A, B).$$

- 3 In this problem we are going to formally analyze the important median trick. Suppose that we have a streaming algorithm for distinct elements that outputs an estimate  $\hat{d}$  of the number  $d$  of distinct elements such that

$$\Pr[\hat{d} > 3d] \leq 47\% \quad \text{and} \quad \Pr[\hat{d} < d/3] \leq 47\%,$$

where the probabilities are over the randomness of the streaming algorithm (the selection of hash functions). In other words, our algorithm overestimates the true value by a factor of 3 with a quite large probability 47% (and also underestimates with large probability). We want to do better!

An important and useful technique for doing better is the median trick: run  $t$  independent copies in parallel and output the median of the  $t$  estimates (it is important that it is the median and *not* the mean as a single horrible estimate can badly affect the mean). Prove that if we select  $t = C \ln(1/\delta)$  for some large (but reasonable) constant  $C$ , then the estimate  $\hat{d}$  given by the median trick satisfies

$$d/3 \leq \hat{d} \leq 3d \quad \text{with probability at least } 1 - \delta.$$

*Hint: an important tool in this exercise are the Chernoff Bounds, which basically say that sums of independent variables are highly concentrated.* Two such bounds can be stated as follows. Suppose  $X_1, X_2, \dots, X_n$  are independent random variables taking values in  $\{0, 1\}$ . Let  $X$  denote their sum and let  $\mu = \mathbb{E}[X]$  denote the sum's expected value. Then for any  $\delta \in (0, 1)$ ,

$$\Pr[X \leq (1 - \delta)\mu] \leq e^{-\frac{\delta^2 \mu}{2}} \quad \text{and} \quad \Pr[X \geq (1 + \delta)\mu] \leq e^{-\frac{\delta^2 \mu}{3}}.$$

**Solution:** Let  $d_i$  be the estimate of the  $i$ -th copy of the algorithm and  $\bar{d}$  be the median of  $d_i$ . We also define  $X_i = 1$  if  $d_i \geq 3d$  and zero otherwise. Moreover,  $X = \sum_{i=1}^t X_i$ . Since  $\Pr[\hat{d} > 3d] \leq 0.47$ , the expected number of answers that exceed  $3d$  is  $0.47t$ . If the median is larger than  $3d$ , this means at least half of  $t$  individual answers exceed  $3d$ . Then, we have:

$$\Pr(X > t/2) = \Pr(X > (1 + 3/47)0.47t) \leq e^{-\frac{3/47^2 0.47t}{3}} \leq e^{-0.00063t} = e^{-0.00063C \ln(1/\delta)} \leq \delta/2$$

with  $C$  being a large enough constant.

Similarly, by defining the probability that the median is below  $\frac{d}{3}$  we get:

$$\Pr(X < t/2) = \Pr(X < (1 - 3/53)0.53t) \leq e^{-\frac{3/53^2 0.53t}{2}} \leq \delta/2$$

with  $C$  being a large enough constant.

This means the probability that  $d/3 \leq \bar{d} < 3d$  is at least  $1 - \delta$ .

- 4 (\*, Pairwise independent random variables) Let  $y_1, y_2, \dots, y_n$  be uniform random bits. For each non-empty subset  $S \subseteq \{1, 2, \dots, n\}$ , define  $X_S = \oplus_{i \in S} y_i$ . Show that the bits  $\{X_S : \emptyset \neq S \subseteq \{1, 2, \dots, n\}\}$  are pairwise independent.

This shows how to stretch  $n$  truly random bits to  $2^n - 1$  pairwise independent bits.

*Hint: Observe that it is sufficient to prove  $\mathbb{E}[X_S] = 1/2$  and  $\mathbb{E}[X_S X_T] = 1/4$  to show that they are pairwise independent. Also use the identity  $\oplus_{i \in A} y_i = \frac{1}{2} (1 - \prod_{i \in A} (-1)^{y_i})$ .*

**Solution:** Recall the definition of pairwise independence: for any non-empty  $S$  and  $T$  such that  $S \neq T$  and two bits  $b_S$  and  $b_T$ , we have

$$\Pr[X_S = b_S \wedge X_T = b_T] = 1/4.$$

We now first argue that  $\mathbb{E}[X_S] = 1/2$ ,  $\mathbb{E}[X_T] = 1/2$  and  $\mathbb{E}[X_S X_T] = 1/4$  implies that they are pairwise independent. We have

$$\begin{aligned} \Pr[X_S = 1 \wedge X_T = 1] &= \mathbb{E}[X_S X_T] = 1/4, \\ \Pr[X_S = 1 \wedge X_T = 0] &= \mathbb{E}[X_S] - \mathbb{E}[X_S X_T] = 1/4, \\ \Pr[X_S = 0 \wedge X_T = 1] &= \mathbb{E}[X_T] - \mathbb{E}[X_S X_T] = 1/4, \\ \Pr[X_S = 0 \wedge X_T = 0] &= \text{“remaining probability”} = 1 - 3 \cdot 1/4 = 1/4. \end{aligned}$$

We thus complete the proof by showing that  $\mathbb{E}[X_S] = \mathbb{E}[X_T] = 1/2$  and  $\mathbb{E}[X_S X_T] = 1/4$ . In both calculations we use the identity  $\oplus_{i \in A} y_i = \frac{1}{2} (1 - \prod_{i \in A} (-1)^{y_i})$ . For the former,

$$\mathbb{E}[X_S] = \mathbb{E}[\oplus_{i \in S} y_i] = \mathbb{E}\left[\frac{1}{2} \left(1 - \prod_{i \in S} (-1)^{y_i}\right)\right] = \frac{1}{2} \left(1 - \prod_{i \in S} \mathbb{E}[(-1)^{y_i}]\right) = \frac{1}{2}.$$

The second to last equality is due to the independence of the random bits  $y_i$  and the last equality follows because  $y_i$  is an uniform random bit. The same calculation also shows that  $\mathbb{E}[X_T] = 1/2$ .

For the latter,

$$\begin{aligned} \mathbb{E}[X_S X_T] &= \mathbb{E}[\oplus_{i \in S} y_i \cdot \oplus_{i \in T} y_i] \\ &= \mathbb{E}\left[\frac{1}{2} \left(1 - \prod_{i \in S} (-1)^{y_i}\right) \cdot \frac{1}{2} \left(1 - \prod_{i \in T} (-1)^{y_i}\right)\right] \\ &= \frac{1}{4} \left(1 - \mathbb{E}\left[\prod_{i \in S} (-1)^{y_i}\right] - \mathbb{E}\left[\prod_{i \in T} (-1)^{y_i}\right] + \mathbb{E}\left[\prod_{i \in S} (-1)^{y_i} \prod_{i \in T} (-1)^{y_i}\right]\right) \\ &= \frac{1}{4} \left(1 + \mathbb{E}\left[\prod_{i \in S} (-1)^{y_i} \prod_{i \in T} (-1)^{y_i}\right]\right) \quad (\text{by independence of } y_i\text{'s}) \\ &= \frac{1}{4} \left(1 + \mathbb{E}\left[\prod_{i \in S \Delta T} (-1)^{y_i}\right]\right) \quad (\text{recall } S \Delta T = S \setminus T \cup T \setminus S) \\ &= \frac{1}{4} \quad (S \Delta T \neq \emptyset \text{ and again using independence of } y_i\text{'s}). \end{aligned}$$