

August 29, 2011

ASSIGNMENT 1

Problem 1. Analysis of number of comparisons in Random Quick Select We shall first prove that the expected number of comparisons is less than $4n$.

Lemma 1. *Given a x randomly selected from S , the expected size of the larger set among $S_{<x}$ and $S_{>x}$ is $\frac{3|S|}{4}$*

Proof. Now let us look at the smaller set, the size of the smaller set can vary anywhere from 0 to $\frac{|S|}{2}$, since x is taken uniformly at random

$$E[\text{—Smaller Set—}] = \frac{|S|}{4} \quad (1-1)$$

This gives us that the expected size of the larger set is

$$E[\text{—Larger Set—}] = \frac{3|S|}{4} \quad (1-2)$$

□

Now consider the first x chosen it will have to make $|S|$ comparisons, and finally it will be splitting the set S into two and we can see that we will be discarding at least one of them. To ensure that we get the worst case expected probability we assume that the smaller set is discarded and hence the size of the instance now reduces to $\frac{3|S|}{4}$. Now taking this as the set the same step is repeated hence.

$$E[\text{Number of comparisons}] \leq |S| + \frac{3|S|}{4} + \left(\frac{3|S|}{4}\right)^2 + \dots \quad (1-3)$$

This is an infinite geometric series with $r = \frac{3}{4}$, and hence taking $|S| = n$ we get.

$$\text{Sum} = \frac{n}{1 - \frac{3}{4}} = 4n \quad (1-4)$$

Hence

$$E[\text{Number of comparisons}] \leq 4n \quad (1-5)$$

→ Answer

Problem 2. Candidate Selection Problem

We shall do this problem by careful application of partition theorem. Now we need to find the probability of Professor Dixon hiring the best qualified candidate.

$$P[\varepsilon] \quad (2-1)$$

Where ε is the event that the most eligible candidate is hired

Since computing this value will be very hard we shall apply partitioning on it. We consider the event X_i where the best candidate is hired provided that he is at position i .

It is easy to see that since the best candidate can occur in any position from 1 to n with equal probability (since we have assumed that the applicants appear in a uniformly random order)

If we denote

$$p_i = Pr[X_i] \quad (2-2)$$

then we can see that

$$P[\varepsilon] = \frac{1}{n} \sum_{i=1}^n p_i \quad (2-3)$$

We shall now construct a qualification array A where A_i = number of people behind him + 1. It is easy to see that for the best qualified candidate $A[i] = n$.

Since Professor Dixon rejects all candidates appearing in positions from 1 to k . We can see that

$$p_i = 0 \forall i = [1 \dots k] \quad (2-4)$$

Now we shall assume the case of $i > k$.

Lemma 1. *The best candidate at i^{th} position is hired iff the maximum element of $A[1 \dots i-1]$ is at a position $\leq k$.*

Proof. If the maximum element of $A[1 \dots i-1]$ is at a position $\leq k$, then no person from $k+1$ to $i-1$ will be better qualified than the most qualified person from 1 to k .

Now if the maximum element of $A[1 \dots i-1]$ is at a position $> k$, then either that person will be hired, or some other person from the $(k+1)^{th}$ to this person will be hired, and hence the person at i^{th} position won't be hired. \square

We can now define p_i as the probability that the best qualified person in $A[1 \dots (i-1)]$ appears from 1 to k . Since all permutations of $n-1$ persons are possible in the $i-1$ positions.

$$Pr[\text{Maximum of } A[1 \dots (i-1)] \text{ is from } 1 \text{ to } k] = \frac{k}{i-1} \quad (2-5)$$

Hence we can see that

$$p_i = \frac{k}{i-1} \quad (2-6)$$

Now we can express $P[\varepsilon]$

$$P[\varepsilon] = \frac{1}{n} \sum_{i=1}^n p_i = \frac{1}{n} \left(\sum_{i=1}^k 0 + \sum_{i=k+1}^n \frac{k}{i-1} \right) \quad (2-7)$$

Now we can take k outside, and we get

$$P[\varepsilon] = \frac{k}{n} \sum_{i=k+1}^n \frac{1}{i-1} \quad (2-8)$$

Expressing this in terms of the Harmonic Number H_n .

$$P[\varepsilon] = \frac{k}{n} \left(\sum_{i=1}^{n-1} \frac{1}{i} - \sum_{i=1}^{k-1} \frac{1}{i} \right) \quad (2-9)$$

$$P[\varepsilon] = \frac{k}{n} (H_{n-1} - H_{k-1}) \quad (2-10)$$

We know that $H_x \approx \ln x + \gamma$ where γ is the Euler- Mascheroni Constant.

$$P[\varepsilon] = \frac{k}{n} (\ln(n-1) - \ln(k-1)) \quad (2-11)$$

$$P[\varepsilon] = \frac{k}{n} \left(\ln\left(\frac{n-1}{k-1}\right) \right) \quad (2-12)$$

Which is the required answer for part (i).

Now to find the value for which the probability is maximum we differentiate, we first express

$$P[\varepsilon] \leq \frac{k}{n} \ln\left(\frac{n}{k}\right) \quad (2-13)$$

Differentiating this *w.r.t* k we get

$$\frac{\partial}{\partial k} \left(\frac{k}{n} \ln\left(\frac{n}{k}\right) \right) = \frac{1}{n} \left(\ln\left(\frac{n}{k}\right) - 1 \right) = 0 \quad (2-14)$$

Since we know that n is finite then maximizing this

$$\ln\left(\frac{n}{k}\right) = 1 \quad (2-15)$$

$$\frac{n}{k} = e \quad (2-16)$$

$$k = \frac{n}{e} \quad (2-17)$$

Now substituting for this value of k we get

$$P[\varepsilon] \leq \frac{1}{e} \quad (2-18)$$

Problem 3.

Randomization in Error Rectification

The algorithm we define is as follows,

- (a) Pick a number uniformly and randomly from 0 to $n - 1$ let this be x .
- (b) Now we compute $y = (z - x) \bmod n$. This gives us two random variables x and y such that $(x + y) \bmod n = z$.
- (c) Substituting this value for z , we get $F((x + y) \bmod n) = (F(x) + F(y)) \bmod m$. So now we do not have to find $F(z)$, but instead we find $F(z) = (F(x) + F(y)) \bmod m$. This completes our algorithm.

We shall now analyze the probability of $F(z)$ being corrupt

Since x is a uniformly and randomly picked number

$$P[F(x) \text{ has been corrupted}] = \frac{1}{5} \quad (3-1)$$

$$P[F(x) \text{ has not been corrupted}] = \frac{4}{5} \quad (3-2)$$

We see that since x is a uniformly random variable, y can also be any of the numbers uniformly (but not randomly). We can see that here we do not require randomness, but the possibility that all values of y are equally possible, this means that adversary will NOT be able to change certain numbers alone, and give a failure probability greater than the fraction of table entries corrupted, This enables us to say that the distribution of $F(y)$ will not change with the change according to the value z chosen and hence

$$P[F(y) \text{ has been corrupted}] = \frac{1}{5} \quad (3-3)$$

$$P[F(y) \text{ has not been corrupted}] = \frac{4}{5} \quad (3-4)$$

Since we find that we cannot argue the independence of the random variables. We shall use Union Theorem. $F(z)$ is corrupted if either $F(x)$ or $F(y)$ is corrupted.

$$P[F(z) \text{ has been corrupted}] \leq P[F(x) \text{ has been corrupted}] + \quad (3-5)$$

$$P[F(y) \text{ has been corrupted}] \quad (3-6)$$

$$P[F(z) \text{ corrupted}] \leq \frac{2}{5} \quad (3-7)$$

$$P[F(z) \text{ is correct}] \geq \frac{3}{5} \quad (3-8)$$

This completes our answer to part (i) of the problem.

For the second part, we run the same algorithm three times, and take the value than has been repeated most number of times, if no value has been repeated, then we output the first value.

We check the following probabilities.

Probability that all three values of $\mathbf{F}(\mathbf{z})$ are correct =

$$P[\text{All 3 correct}] \geq \left(\frac{3}{5}\right)^3 = \frac{27}{125} \quad (3-9)$$

Since all three runs of the algorithm are independent of each other we have been able to apply the product rule.

Probability that any two values of $\mathbf{F}(\mathbf{z})$ are correct =

$$P[\text{Some 2 correct}] \geq \binom{3}{2} \left(\frac{3}{5}\right)^2 \left(\frac{2}{5}\right) = \frac{54}{125} \quad (3-10)$$

Probability that the first value is correct and the rest two are wrong.

$$P[\text{Some 2 correct}] \geq \left(\frac{3}{5}\right) \left(\frac{2}{5}\right)^2 = \frac{12}{125} \quad (3-11)$$

We can see that the following is a partition of the probable cases and hence if we are allowed to repeat our algorithm 3 times with independent values of \mathbf{x} each time, then the total success probability.

$$P[\mathbf{F}^{(3)}(\mathbf{z}) \text{ is correct}] \geq \frac{27}{125} + \frac{54}{125} + \frac{12}{125} = \frac{93}{125} \quad (3-12)$$

$$\geq 0.744 \quad (3-13)$$

→ Answer

Problem 4.

Set Balancing

We shall first look at a single entry. Let \mathbf{A}_{ij} be the element at the i^{th} row and j^{th} column of the matrix \mathbf{A} . We also define \mathbf{Ab}_i as the i^{th} value of the column vector \mathbf{Ab} .

Proving that $\|\mathbf{Ab}\|_\infty$ is of order $O(\sqrt{n \log n})$ is equivalent to proving that

$$\|\mathbf{Ab}\|_\infty \leq C\sqrt{n \log n} \quad (4-1)$$

for some finite C with very high probability.

Lemma 1. $|Ab_i| \leq c\sqrt{n \log n}$ with very high probability where Ab_i is the i^{th} value of the Column Vector Ab .

Proof. We can write Ab_i as

$$Ab_i = \sum_{j=1}^n A_{ij}b_j \quad (4-2)$$

Where A_{ij} is from the set $\{0, 1\}$ and b_j is randomly picked from $\{-1, +1\}$

It is very easy now to see the connection with the random walk problem.

In short we can see that if $A_{ij} = 0$ then it has no effect to the sum, this can be seen as equivalent to "not moving" in a step in the random walk. We shall now reorder the sum in such a way that all the 0's occur at the end, since b_i is uniformly random this will not change the problem. Let us assume,

$$m = \sum_{j=1}^n A_{ij} \quad (4-3)$$

This will count the total number of non-zero entries in the row, now we can see that this problem is equivalent to the random walk problem.

We shall now define a new row A'_i of size m with all zero entries removed, and B with all the entries in the same index as the zero entries of A_i removed.

Now let us define a random variable

$$Y_i = \{1 \text{ if } B_i = +1, 0 \text{ if } B_i = -1\} \quad (4-4)$$

Expected value of $Y = \sum Y_i$

$$E[Y] = \frac{m}{2} \quad (4-5)$$

We can see that Y is the result of a series of bernouli trials and hence we can apply Chernoff bound on this.

We need to prove that

$$P[Y > (1 + \delta)\frac{m}{2}] \leq \frac{1}{m^c} \quad (4-6)$$

Applying Chernoff Bound, we get

$$e^{\frac{\mu\delta^2}{2}} \leq \frac{1}{m^c} \quad (4-7)$$

$$e^{\frac{m\delta^2}{4}} \leq e^{c \ln m} \quad (4-8)$$

$$\frac{m\delta^2}{4} \leq c \ln m \quad (4-9)$$

$$(4-10)$$

Now we find the value of δ that will make this valid.

$$\delta = \sqrt{\frac{4c \ln m}{m}} \quad (4-11)$$

Applying this value for δ we also set $c = 2$ get

$$P[Y > (1 + \sqrt{\frac{8 \ln m}{m}}) \frac{m}{2}] \leq \frac{1}{m^2} \quad (4-12)$$

Rewriting

$$P[(Y - \frac{m}{2}) > \sqrt{8m \ln m}] \leq \frac{1}{m^2} \quad (4-13)$$

From this we can get

$$P[Ab_i > \sqrt{8m \ln m}] \leq \frac{1}{m^2} \quad (4-14)$$

Now m was the number of non-zero entries in the row and setting $E[m] = \frac{n}{2}$

$$P[Ab_i > 2\sqrt{n \ln n}] \leq \frac{4}{n^2} \quad (4-15)$$

And hence we say $|Ab_i| \leq cn \log n$ with very high probability. \square

Now we apply union theorem on the result, and hence

$$P[\text{Max}(Ab_i) > 2\sqrt{n \ln n}] \leq n \frac{4}{n^2} \quad (4-16)$$

$$P[\text{Max}(Ab_i) > 2\sqrt{n \ln n}] \leq \frac{4}{n} \quad (4-17)$$

Hence we can say that $\|Ab\|_\infty = O(\sqrt{n \log n})$ with very high probability

→ Answer

Submitted by Abhimanyu M A (1111002) , Sumesh T A (1111065) on August 29, 2011.