

Lecture 11 February 24 2012

Lecturer: Abhishek Dang

Scribe: Anudhyan Boral and Arjun Arul

1 The Predecessor Problem

We are given a universe U of size 2^m and subset $S \subset U$, $|S| = n$. For $x \in U$, define the following functions:

Definition 1. $\text{pred}_S(x) = \max\{y \in S \mid y \leq x\}$.

Definition 2. $\text{rank}_S(x) = |\{y \in S \mid y \leq x\}|$.

Definition 3. $\oplus \text{rank}_S(x) = |\{y \in S \mid y \leq x\}| \pmod 2$.

Now given S , we wish to answer queries on S . The preprocessing algorithm should store information about S in an appropriate way so that given any $x \in U$, we can find $f_S(x)$ efficiently, where f could be one of the above functions.

Definition 4. A randomized $(s, w, t)_\epsilon$ storage scheme for $f_S(x)$ consists of

- a deterministic storage algorithm which takes as input $S \subset U$ and outputs a data structure T with s cells, each cell w bits long.
- a randomized query algorithm which, on input $x \in U$, probes at most t cells in T , and outputs $f_S(x)$ correctly with probability at least $(1 - \epsilon)$.

All these functions depend on S ; henceforth we drop the subscript S in pred_S , rank_S and $\oplus \text{rank}_S$ for convenience. Some departed from last time: we use m for the bit size of an element in the universe. We do not assume that the cell-word-size is m but allow an independent parameter w .

We have seen an $(\mathcal{O}(n), \mathcal{O}(m), \mathcal{O}(1))$ deterministic scheme for the dictionary problem using FKS hashing. We also saw an $(\mathcal{O}(mn), \mathcal{O}(m), \mathcal{O}(\log m))$ deterministic scheme for the predecessor problem, using X-tries and the dictionary solution. We also stated without proof that there is an $(\mathcal{O}(mn), \mathcal{O}(m), \min \left[\mathcal{O}\left(\frac{\log m}{\log \log m}\right), \mathcal{O}\left(\sqrt{\frac{\log n}{\log \log n}}\right) \right])$ deterministic scheme for the predecessor problem.

In this lecture we will show that the upper bound is almost tight.

Theorem 5. For $s \in \text{poly}(n)$, $w \in \text{poly}(m)$, if there is an $(s, w, t)_\epsilon$ randomized scheme for the predecessor problem, then $t \in \Omega \left[\frac{\log m}{\log \log m}, \sqrt{\frac{\log n}{\log \log n}} \right]$.

We will actually prove the theorem for the \oplus rank function. Then we will make use of the following observation:

Observation 6. *If there is an $(s, w, t)_\epsilon$ scheme for $\text{pred}(x)$, then there is an $(s + \mathcal{O}(n), w + \mathcal{O}(m), t + \mathcal{O}(1))_\epsilon$ scheme for $\text{rank}(x)$. This is because for each $y \in S$, if y is the predecessor of x , then $\text{rank}(x) = \text{rank}(y)$. So given x we first find $\text{pred}(x)$, and then query a dictionary to find $\text{rank}(\text{pred}(x))$. And for each $y \in S$, we can use the FKS scheme for the dictionary problem to store $\text{rank}(y)$. Similarly, under this hypothesis, $\bigoplus \text{rank}(x)$ also has an $(x + \mathcal{O}(n), w + \mathcal{O}(m), t + \mathcal{O}(1))_\epsilon$ scheme.*

Let (m, n) denote the size of the universe $|U| = 2^m$ and the size of the subset $|S| = n$. We carry these parameters as subscripts with the function.

To actually prove the theorem for the $\bigoplus \text{rank}_{m,n}$ function, we will consider the communication game associated with $\bigoplus \text{rank}_{m,n}$. Alice has an element $x \in U$, with $x = (x_1, \dots, x_m)$, each $x_i \in \{0, 1\}$. Bob has the subset $S = y_1, \dots, y_n \subset U$. They wish to determine $\bigoplus \text{rank}_{m,n}(x)$ with respect to S .

Now, given a $(2^a, b, t)_\epsilon$ scheme for $\bigoplus \text{rank}_{m,n}$, there is a protocol for the communication game which satisfies the following,

- Messages from Alice to Bob are a bits long.
- Messages from Bob to Alice are b bits long.
- Alice begins, and there are $2t$ rounds.
- The protocol errs with probability at most ϵ .

The protocol is simple. Bob runs the preprocessing algorithm and constructs the data structure T . Alice runs the query algorithm. Whenever she needs to probe a cell, she sends the cell number to Bob, who responds with the contents of that cell in T . The randomness can be private or public, it is required only by Alice, while running the query algorithm.

We call any protocol with these properties a $(2^a, b, t)_{(\epsilon, m, n)}^A$ protocol for $\bigoplus \text{rank}_{m,n}$. A $(2^a, b, t)_{(\epsilon, m, n)}^B$ protocol for $\bigoplus \text{rank}_{m,n}$ is a similar $(2t - 1)$ -round protocol where Bob begins the communication. Note that a protocol for (m, n) is also a protocol for (m', n) for every $m' \leq m$.

The lower bound proof proceeds as follows. Suppose we have a $(2t, a, b)_\epsilon^A$ protocol for $\bigoplus \text{rank}_{m,n}$. Using round elimination we will then show that: $(2t, a, b)_\epsilon^A$ protocol for $\bigoplus \text{rank}_{m,n} \implies (2t - 1, a, b)_{\epsilon + \text{frac}112t}^B$ protocol for $\bigoplus \text{rank}_{\frac{m}{k}, n}$ [eliminate Alice's first message; still OK for slightly smaller universe] $\implies (2t - 2, a, b)_{\epsilon + \text{frac}16t}^A$ protocol for $\bigoplus \text{rank}_{\frac{m}{k} - \log l, \frac{n}{l}}$ [eliminate Bob's first message; still OK for slightly smaller set]

We will show that for $c_1 = 72 \ln 2$, $k = c_1 a t^2$, and $l = c_1 b t^2$, each round elimination adds no more than $1/6t$ to the error.

Consider the following parameters: m is any given value. Choose $n = 2^{\log^2 m / \log \log m}$. Set $c_1 = 72 \ln 2$, and let c_2, c_3 be any constants greater than 1. Choose $a = c_2 \log n$, $b = m^{c_3}$.

Let $t = \frac{\log m}{(c_1 + c_2 + c_3) \log \log m}$. Choose $k = c_1 a t^2$, $l = c_1 b t^2$. With these parameters, we can verify that:

- $\frac{m}{k} - \log l \geq \frac{m}{2k}$.
- $m' = \frac{m}{(2k)^t} \in m^{\Omega(1)}$.

- $n' = \frac{n}{l^t} \in n^{\Omega(1)}$.

Then, if we repeat round elimination t times, we obtain a $(0, a, b)_{\epsilon + \frac{1}{6}}$ protocol for $\bigoplus \text{rank}_{m', n'}$ for non-trivial m', n' . For $\epsilon < \frac{1}{3}$, we get a zero round protocol with error less than $\frac{1}{2}$. But this means that with no information whatsoever about the set S (since there is no communication between Alice and Bob), Alice can guess $\bigoplus \text{rank}(x)$ and be right with probability greater than $1/2$, which is a contradiction.

We now proceed to prove the round elimination theorem. Assume that the constants are chosen as above. Suppose P is a $(2t, a, b)_{\epsilon}^A$ protocol for $\bigoplus \text{rank}_{m, n}$. We will convert P into a $(2t - 2, a, b)_{\epsilon + \text{frac}16t}^A$ for $\bigoplus \text{rank}_{\frac{m}{k} - \log l, \frac{n}{l}}$.

1.1 Round Elimination: Eliminating Alice's message

We will first convert P into a $(2t - 1, a, b)_{\epsilon + \text{frac}112t}^B$ protocol Q for $\bigoplus \text{rank}_{\frac{m}{k}, n}$. To do so we will use the randomized version of Yao's lemma which states, $R_{\epsilon}(f) = \max_{\mu} D_{\epsilon}^{\mu}(f)$ where the protocols D_{ϵ}^{μ} are randomized. We will show that for any distribution μ over (x, S) , there is a $(2t - 1, a, b)_{\epsilon + \frac{1}{12t}}^B$ protocol Q that solves $\bigoplus \text{rank}_{\frac{m}{k} - \log l, \frac{n}{l}}$ well when the inputs are distributed according to μ . Recall that P works well for all distributions; in particular, it works well for (m, n) distributions that somehow extend μ .

Choose any distribution μ over (x, S) where $|U| = 2^{\frac{m}{k}}$ and $|S| = n$. We first design a protocol $(2t, a, b)_{\epsilon}^A$ protocol Q' for $\bigoplus \text{rank}_{\frac{m}{k}, n}$ with respect to μ . Then we adapt Q' to obtain Q .

1.1.1 The Protocol Q

Consider a run of the protocol P . Let Alice's input be $x' = x_1, \dots, x_k$ where x' is broken up into blocks of length m/k , and each block x_i is drawn according to μ . Let M be the first message that Alice sends in the protocol P while using randomness R .

$$\begin{aligned}
I(x' : MR) &= I(x' : R) + I(x' : M|R) \\
&\leq 0 + H(M|R) \\
&\leq H(M) \\
&\leq |M| = a
\end{aligned}$$

Therefore,

$$\begin{aligned}
a &\geq I(x_1, \dots, x_k : MR) \\
&= I(x_1 : MR) + I(x_2 : MR|x_1) + \dots + I(x_k : MR|x_1, \dots, x_{k-1})
\end{aligned}$$

Therefore, there is a block numbered $i \in [k]$ such that

$$I(x_i : MR | x_1, \dots, x_{i-1}) \leq \frac{a}{k}$$

That is, the first message from Alice and the public randomness together give Bob very little information about the i th block, even if Bob knows the strings in all the preceding blocks. Fix such an i . By definition,

$$E_{x_1=u_1, \dots, x_{i-1}=u_{i-1}}[I(x_i : MR | x_1 = u_1, \dots, x_{i-1} = u_{i-1})] \leq \frac{a}{k}$$

So $\exists u_1, \dots, u_{i-1} [I(x_i : MR | x_1 = u_1, \dots, x_{i-1} = u_{i-1})] \leq \frac{a}{k}$

Fix these u_1, \dots, u_{i-1} .

Now we start designing Q' . Alice gets $x \in U = 2^{\frac{m}{k}}$ and Bob gets a set $S \subset U$ of size n , where (x, S) are drawn according to μ . To run P , they must *extend* their inputs to look like inputs to P . The idea is to embed x and S into the i th block of suitable chosen longer strings, so as to make the first message almost irrelevant.

Bob extends his set by prefixing each element of S with $u_1 \dots u_{i-1}$ with suffixing it with zeroes. That is, he constructs the set $S' = u_1 \dots u_{i-1} y 0^{(k-i)\frac{m}{k}} | y \in S$.

Alice constructs the element x' by prefixing x with $u_1 \dots u_{i-1}$ and suffixing it with $k-i$ blocks each chosen according to μ using private randomness. Thus $x' = u_1 \dots u_{i-1} x x_{i+1} \dots x_k$, where x_{i+1}, \dots, x_k are drawn according to μ .

Observe that $\bigoplus \text{rank}_{\frac{m}{k}, n}(x, S) = \bigoplus \text{rank}_{m, n}(x', S')$. So Alice and Bob can now run the protocol P to determine $\bigoplus \text{rank}_{\frac{m}{k}, n}(x, S)$. This is the $(2t, a, b)_\epsilon^A$ protocol Q' for $\bigoplus \text{rank}_{\frac{m}{k}, n}$.

1.2 The Protocol Q

Observe that because of the way we constructed the protocol Q' , the first message M sent by Alice to Bob contains very little information about x , i.e. $I(x : MR) \leq \frac{a}{k}$. Since M contains so little information about x , Bob might as well replace it with an "average" message. This will introduce some additional error, but we can keep this within bounds using the following:

Theorem 7. (Average Encoding Theorem) *Let X, Y be correlated random variables with joint distribution $r_{x,y}$. Let F be the marginal distribution of Y . For any x , let F^x denote the distribution of Y conditioned of the event $X = x$. Then,*

$$\sum_x Pr[X = x] \|F^x - F\|_1 \leq \sqrt{2 \ln(2) I(X : Y)}$$

Proof. Consider the definitions of these quantities:

$$F(y) = \sum_{x'} r_{x',y}; \quad F^x(y) = \frac{r_{x,y}}{\sum_{y'} r_{x,y'}}; \quad Pr(X = x) = \sum_{y'} r_{x',y}$$

Define the following distributions on XY :

$$P(x, y) = Pr[X = x]F^x(y) \qquad Q(x, y) = Pr[X = x]F(y)$$

The first distribution P is exactly the joint distribution $r_{x,y}$. The second distribution Q is a product distribution: imagine independent random variables X', Y' distributed according to the marginals, and consider their joint distribution. Therefore,

$$\text{LHS in Theorem} = \|P - Q\|_1 \leq \sqrt{(2 \ln(2))D(P\|Q)} = \sqrt{2 \ln(2)I(X : Y)}$$

Here, $D(P\|Q)$ is the relative entropy or Kullback-Leibler distance between P and Q . This gives the inequality above.

Now we define the $(2t - 1, a, b)$ protocol Q for $\bigoplus \text{rank}_{\frac{m}{k}, n}$, where (x, S) are drawn according to distribution μ .

Alice gets a string x of $\frac{m}{k}$ bits. Bob gets a set S of size n . Bob constructs $S' = \{u_1 \dots u_{i-1}y0^{(k-1)\frac{m}{k}}|y \in S\}$. Bob then uses public randomness R to construct the "average" message. That is, using public randomness he samples U_i, \dots, U_k according to μ , and then simulates the protocol P to generate the first message Alice would have sent if her input were $u_1 \dots u_{i-1}U_i \dots U_k$. We call this the "average" message M' .

Observe that Alice also knows M' , because Bob uses public randomness R . Now Alice does a "reverse engineering" of M' . Using private randomness, she samples V_{i+1}, \dots, V_k according to μ , conditioned on the message being M' and V_i being x . She then constructs $x' = u_1 \dots u_{i-1}xV_{i+1} \dots V_k$. This ensures that Alice and Bob now have "consistent" states with input x, S and first message M' , and Bob still has very little information about x .

Now Alice and Bob proceed using the protocol Q' (which itself used P) from the second message onwards.

1.3 Calculating the Error

Assume Alice's input is x . Consider the following distributions on the set of first messages that can be sent by Alice. Let F^x be the distribution in protocol Q' , when Alice has input x and F be the distribution in protocol Q where Bob samples an average first message. By the Average Encoding Theorem,

$$\begin{aligned} \sum_x Pr[X = x] \|F^x - F\|_1 &\leq \sqrt{2 \ln(2)I(X : MR)} \\ &\leq \sqrt{(2 \ln(2))\frac{a}{k}} \end{aligned}$$

$$\begin{aligned}
Pr[Q \text{ errors}] &= Pr[Q \text{ errors} | M = M'] * Pr[M = M'] + Pr[Q \text{ errors} | M \neq M'] * Pr[M \neq M'] \\
&\leq Pr[Q \text{ errors}] + Pr[M \neq M'] \\
&\leq \epsilon + \sum_x Pr[X = x] Pr[M \neq M' | X = x] \\
&\leq \epsilon + \sum_x Pr[X = x] \frac{1}{2} \|F^x - F\|_1 \\
&\leq \epsilon + \frac{1}{2} \sqrt{2 \ln(2) \frac{a}{k}}
\end{aligned}$$

For a suitable choice of k such as $72 \ln(2) a t^2$, we will get the error to be less than $\epsilon + \frac{1}{12t}$.

1.4 Eliminating Bob's Message

We have a $(2t-1, a, b)_{\delta}^B$ protocol for $\bigoplus \text{rank}_{M,N}$, where $M = \frac{m}{k}$ and $N = n$. Following a similar strategy as above, we will convert P into a $(2t-2, a, b)_{\delta + \frac{1}{12t}}^B$ protocol Q for $\bigoplus \text{rank}_{M - \log(l), \frac{N}{t}}$.

Choose a distribution μ on (x, S) where $x \in 2^{M - \log(l)}$ and $|S| = \frac{N}{t}$. Now let Bobs input in protocol P be $S = [1].S_1 \cup \dots \cup [l].S_l$, where the S_i are chosen according to μ , $[i]$ is the representation of i using $\log(l)$ bits and $[i].S_i = \{[i], y | y \in S_i\}$. Let M be the first message that Bob sends in protocol P while using randomness R . Then,

$$b \geq I(S : MR) = \sum_x I(S_i : MR | S_1, \dots, S_{i-1})$$

So $\exists i$ such that $I(S_i : MR | S_1 \dots S_{i-1}) \leq \frac{b}{k}$. Fix such an i . By definition,

$$\frac{b}{k} \geq E_{S_1=s_1 \dots S_{i-1}=s_{i-1}} I[S_i : MR | S_1 = s_1, \dots, S_{i-1} = s_{i-1}]$$

So, $\exists s_1, \dots, s_{i-1}$ such that $I(S_i : MR | S_1 = s_1, \dots, S_{i-1} = s_{i-1}) \leq \frac{b}{k}$. Fix these sets s_1, \dots, s_{i-1} .

Now the $(2t-1, a, b)_{\delta}^B$ protocol Q' for $\bigoplus \text{rank}_{M - \log(l), \frac{N}{t}}$ is as follows. Bob and Alice embed their inputs into inputs suitable for protocol P .

Bob gets a set S of size $\frac{N}{t}$. Bob draws sets S_{i+1}, \dots, S_l according to μ using public randomness, and constructs $S' = [1].s_1 \cup \dots \cup [i-1].s_{i-1} \cup [i].S \cup [i+1].S_{i+1} \cup \dots \cup [l].S_l$.

Alice gets a string x of length $M - \log(l)$. Alice constructs the string $x' = [i]x$.

Now observe that $\bigoplus \text{rank}_{S'}(x') = \bigoplus \text{rank}_S(x)$. Therefore Alice and Bob run the protocol P on (x', S') . This is the protocol Q' for $\bigoplus \text{rank}_{M - \log(l), \frac{N}{t}}$.

In this protocol, Alice knows the sets s_1, \dots, s_{i-1} since they are fixed. By choice of the index i and these sets, knowing this and after getting the first message from Bob, she still has very little (at most $\frac{b}{k}$) information about S . So if the first message is dispensed with and replaced with an average message, the error wont increase much. This gives the protocol Q : As before, Alice will sample the average first message M' with public randomness, and Bob will reverse engineer the process to sample S_{i+1}, \dots, S_l conditioned on M' and S .

To bound the error, as before, use the Average Encoding Theorem. For a suitable choice of l (at least $72(\ln(2))bt^2$), we will get the error to be less than $\delta + \frac{1}{12t}$.