

September 26, 2011

ASSIGNMENT 2

Problem 1. How to use random sampling ?

We shall be doing the following.

Let our original graph be $G = (V, E)$, we now need to find a subgraph $G_S = (V, E_S)$ which is a 3-approximate distance preserver, that is the distance between (u, v) in G_S is atmost 3 times the distance between (u, v) in the original graph G .

Algorithm

We assume that the original graph is maintained as an adjacency list, from this we will construct an augmented adjacency list where each entry u in the list of v also has a pointer to entry v in the list of u . This is done as follows.

For each entry i from 1 to n . If (i, j) is there in the original adjacency list, and $j > i$ we add (i, j) to list of i and (j, i) to the list of j . Since they are added at the same time we will have pointers from each to the other. Hence we get an augmented adjacency list, the reason we augmented the adjacency list is to support delete operation in $O(1)$ time.

- (a) Let $E_S = \phi$.
- (b) With a probability p , we shall select each vertex and add it to A . We shall fix on the value of p later. We shall also introduce the concept of clusters initially each of the vertex u will be a cluster by itself, and will be referred to as the center of the cluster.
- (c) For each vertex $u \notin A$ we shall do the following
 - i. If for u there is no edge (u, v) where $v \in A$ we shall add all the edges of the vertex to E_S .
 - ii. If u has an edge (u, v_i) where $v_i \in A$, then take the edge (u, v_i) of minimum weight add it to E_S and also add all edges (u, w_j) where $weight(u, w_j) < weight(u, v_i)$ where $w \notin A$, we will keep all other edges since in G since we might later need to add them. In short add the minimum weight edge of u that links it to A , and all edges of weight less than the minimum weight. We shall also add u to the cluster of v_i
- (d) For all edges (u, v) belonging to the same cluster, that has not been added. we remove the edge from the set. This completes our building of the cluster.
- (e) What now exists is the edges between the clusters, for each vertex $v \in C_i$ where C_i is a cluster, for all $j \neq i$ we add the minimum weight edge (u, v) where $u \in C_j$.
- (f) Remove all other edges.

Proof of Correctness

We shall now prove that the algorithm does indeed give a correct result.

To prove that the distance is atmost 3 times the actual distance all we need to prove that is that each edge (u, v) can be replaced by $(u, x), (x, y)(y, v)$ where $weight(u, x) < (u, v)$, $weight(x, y) < (u, v)$ and $weight(y, v) < (u, v)$.

- (a) Now if either u or v doesn't belong to any cluster then we have added all the edges and hence we will not have a problem.
- (b) Now if both u and v belong to the same cluster with cluster centre w ,
 - i. If $u == w$ or $v == w$ we are done every node in the cluster is directly linked to the centre.
 - ii. If $weight(u, v) < weight(u, w)$ or $weight(u, v) < weight(v, w)$ we have already added the edge when we added u or v to the cluster.
 - iii. Else we know that both $weight(u, w) < weight(u, v)$ and $weight(v, w) < weight(u, v)$ hence $weight(u, w) + weight(v, w) < 2 * weight(u, v)$. Hence the construction satisfies the condition necessary for a 3-approximate distance preserving subgraph.
- (c) Now we consider the case that u and v belong to different clusters, let there be an edge (u, w) where $w, v \in C_i$ and $weight(u, w) < weight(u, x) \forall x \in C_i$ where C_i is centred on x Now consider the following cases.
 - i. If $w == v$ we are done since we have added the edge (u, v) .
 - ii. Else consider the set of edges $(u, w), (w, x), (x, v)$. We already know that (u, w) is the least weight edge hence $weight(u, w) < weight(u, v)$, we also know that $weight(x, v) < weight(u, v)$ since otherwise we would have added the edge (u, v) when we added v to the cluster i , now we can also see that $weight(w, x) < weight(u, w) < weight(u, v)$ since (w, x) was added in the first phase and (u, w) in the second phase, hence by the proximity constraint. Hence we can replace (u, v) can be replaced by $(u, x), (x, y)(y, v)$ where $weight(u, x) < (u, v)$, $weight(x, y) < (u, v)$ and $weight(y, v) < (u, v)$.

Bounding expected number of edges

Before we go into the proof of that we find that $E[|A|]$ Since each element is added with probability p , $E[|A|] = np$ Now let $np = q$.

- (a) In the second step we add the edge (u, v) which is the minimum for all $v \in A$, and all edges lesser than that, since each vertex is selected to the cluster with probability p , Assume that all the edges are ordered, then the edges incident on the vertices forming the cluster form a partition of the set of edges, since each vertex is taken uniformly and randomly the expected number of elements selected for cluster = δp where δ is the degree of the vertex. Now since we consider that they are ordered, the set of all edges with weight lower than edge (u, v) is nothing but the expected size of the first partition, which can be easily seen as $\frac{\delta}{\delta p} = \frac{1}{p}$. Now if we consider each vertex, then we can apply linearity of expectation and get the following result.

$$\text{Expected contribution of a single vertex} = \frac{1}{p} \quad (1-1)$$

$$\text{Total contribution by Linearity of Expectation} \leq \sum_{v \in V} \frac{1}{p} \quad (1-2)$$

$$\leq \frac{n}{p} \quad (1-3)$$

The inequality comes due to the fact that the same edge may be selected more than once. We know that this has to be $O(n^{3/2})$. Hence we set $p = \frac{1}{\sqrt{n}}$. We can then see that it will fit the other contributions also into the limit.

- (b) Now that we have fixed the probability p we can analyse the number of edges added by the vertices which does to belong to any cluster. Again we shall partition it into the number of edges added by each vertex.

$$\begin{aligned} \text{Expected number of edges added by single vertex} &= \delta Pr[\text{Vertex is selected}] \quad (1-4) \\ E[X_i] &= \delta Pr[\text{Vertex is selected}] \quad (1-5) \end{aligned}$$

Where δ is the degree of the vertex.

$$\begin{aligned} Pr[\text{Vertex is selected}] &= Pr[\text{None of it's neighbours belong to a cluster}] \quad (1-6) \\ &= \left(1 - \frac{1}{\sqrt{n}}\right)^\delta \quad (1-7) \end{aligned}$$

$$\text{Expected number of edges added by single vertex} = \delta \left(1 - \frac{1}{\sqrt{n}}\right)^\delta \quad (1-8)$$

We can have two cases $\delta \leq \sqrt{n}$ and $\delta > \sqrt{n}$. Now for $\delta \leq \sqrt{n}$

$$E[X_i] < \delta \quad (1-9)$$

$$< \sqrt{n} \quad (1-10)$$

This comes since the probability $\left(1 - \frac{1}{\sqrt{n}}\right)^\delta$ is atmost 1.

Now for $\delta > \sqrt{n}$. Let $\delta = k\sqrt{n}$ where $k > 1$.

$$E[X_i] = k\sqrt{n} \left(1 - \frac{1}{\sqrt{n}}\right)^{k\sqrt{n}}. \quad (1-11)$$

We know that asymptotic value of $\left(1 - \frac{1}{x}\right)^x = \frac{1}{e}$. Substituting that we get.

$$E[X_i] = k\sqrt{n} \left(\frac{1}{e}\right)^k \quad (1-12)$$

$$= \sqrt{n} \left(\frac{k}{e^k}\right) \quad (1-13)$$

$$< \sqrt{n} \cdot 1 \quad (1-14)$$

Hence for both cases we have $E[X_i] < \sqrt{n}$

So

$$E[X] \leq \sum_{\forall v \in V} \sqrt{n} \quad (1-15)$$

$$\leq n^{\frac{3}{2}} \quad (1-16)$$

The \leq sign comes due to the fact that few edges may be selected more than once.

- (c) Now we shall consider the edges added due to the inter cluster linkings. This is very much straight forward to see, each vertex adds a single edge for each cluster. If we consider a single vertex

$$E[X_i] = E[|A|] \quad (1-17)$$

We have already proved initially that, $E[|A|] = np$. Substituting $p = \sqrt{n}$ we get

$$E[X_i] = \frac{n}{\sqrt{n}} \quad (1-18)$$

$$= \sqrt{n}. \quad (1-19)$$

Again by linearity of expectation

$$E[X] \leq \sum_{\forall v \in V} \sqrt{n} \quad (1-20)$$

$$\leq n^{\frac{3}{2}} \quad (1-21)$$

The \leq sign comes due to the fact that few edges may be selected more than once.

Hence total number of edges selected on expectation $= 3n^{\frac{3}{2}}$.

Since we have the total number of edges count at the end of the algorithm, we shall see if this is greater than $6n^{\frac{3}{2}}$, if yes we shall drop all the result and run the algorithm again. We have the result by Markov Inequality

$$Pr \left[|E_S| > 6n^{\frac{3}{2}} \right] \leq \frac{3n^{\frac{3}{2}}}{6n^{\frac{3}{2}}} \quad (1-22)$$

$$Pr \left[|E_S| > 6n^{\frac{3}{2}} \right] \leq \frac{1}{2} \quad (1-23)$$

This now can be considered as a coin toss experiment with success probability $= \frac{1}{2}$. Then the expected number of trials before success $= 2 = O(1)$. Hence we get a Las Vegas algorithm by repeating the algorithm only $O(1)$ time. This completes our answer.

→ Answer

Problem 2. Random choices may work!

- (a) Deterministic algorithm will have a specific ordering for scanning the children. Without loss of generality let us assume that the order is from left to right. We can fix the deterministic algorithm to fix the values of the node.
- i. For the root set value as either **0** or **1**.
 - ii. Now for each of the children from left to right assign values as follows.
 - iii. leftChild.value = parent.value.
 - iv. middleChild.value = not (parent.value).
 - v. rightChild.value = parent.value

Proof. Now consider for the deterministic algorithm at all points of time when the left and middle children's values are not conclusive enough to decide the value of the parent hence we need to scan right child also. Hence we will need to scan all nodes of the tree. Now we need to prove the number of nodes in the tree. Thus tree with height 0 tree will have 1 node, tree with height 2 will have 3 leaves and tree with height 3 will have 9 leaves and so on. Thus tree with height h will have 3^h leaf nodes. Hence a deterministic algorithm must must read all $n = 3^h$ leaves when the left and middle children of the parent node have contradicting values. \square

(b) The algorithm can be stated as follows.

- i. Randomly choose any two leaf nodes which has same parent.
- ii. If they have same value then value of parent node is equal to that of value of these two children.
- iii. Else scan the third child and set the value of the parent is equal to the value of this child.

Analysis: For cases where leaf nodes are having values 0,0,0 and 1,1,1. we can immediately set the value after scanning any two nodes. We will now look at worst case scenario where we have one value is different.

The expected number query = $(2 * 1) + (1 * \text{Prob}(\text{first two values are different}))$

Now, $\text{Prob}(\text{first two values are different}) = \frac{2}{3}$.

The expected number query $\leq (2 * 1) + \frac{2}{3} = 2.67$

Now for the total tree of height h , we can see that total number of expected query = $(2.67)^h$

now we know that $3^h = n$

Height $h = \log_3 n$

The expected number query = $(2.67)^{\log_3 n}$

$\leq n^{\log_3 2.67}$

$\leq n^{.8939}$

→ Answer

Problem 3. Two-dimensional Pattern Matching We shall first show that the finger print function has the required properties.
Let the string $x = \{0, 1\}^*$

- (a) We shall now give a deterministic algorithm to generate the finger print function. Let $I = I_2$ the 2×2 Identity Matrix. Now for a string x we shall prove that $M(x)$ is defined by induction.
- For strings of length $0 \leq 1$ it is trivially true since $M(\epsilon), M(0), M(1)$ are well defined.
 - Now let us assume it is defined for all strings of length $\leq k$
 - Now let $x = \{0, 1\}^{k+1}$ be a string of length $k + 1$. Let $x = yz$ such that $y = \{0, 1\}^k$ and $z = 0, 1$. We know that both of them are well defined since they are strings of length $\leq k$. From the definition $M(x) = M(y) \times M(z)$, since $M(y)$ and $M(z)$ are well defined, $M(x)$ is also well defined.
- (b) To prove that

$$M(x) = M(y) \Rightarrow x = y \quad (3-1)$$

We shall define a deterministic algorithm to retrieve the string given the fingerprint function, proving that this can be done deterministically in effect proves that the function from string x to $M(x)$ is bijective and hence $M(x) = M(y) \Rightarrow x = y$. We shall once again use induction to prove our hypothesis.

- For strings of length $0 \leq 1$ it is true since $M(\epsilon), M(0), M(1)$ are well defined, and distinct.
- Now let us assume that for all strings of length $\leq k$ we will be able to find x given $M(x)$.
- Now let $x = \{0, 1\}^{k+1}$ be a string of length $k + 1$. Let $x = yz$ such that $y = \{0, 1\}^k$ and $z = 0, 1$.

$$\text{Let } M(y) = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad (3-2)$$

If $z = 0$.

$$M(x) = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \times [r]1011 \quad (3-3)$$

$$M(x) = \begin{bmatrix} a+b & b \\ c+d & d \end{bmatrix} \quad (3-4)$$

If $z = 1$.

$$M(x) = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \times [r]1101 \quad (3-5)$$

$$M(x) = \begin{bmatrix} a & a+b \\ c & c+d \end{bmatrix} \quad (3-6)$$

Since only for $b = 0$ or $c = 0$ we can see that if the last character was a 0 then the current matrix then atleast one entry in the left will be strictly greater than the value of the right, and vice versa if the last character was a 1 , hence given a matrix we can find out what the last character was deterministically, and since we have assumed that we can find out for the smaller matrix, inductively we can apply this until we reach the identity matrix. Hence we can define a bijective function from $M(x) \rightarrow x$. Hence

$$M(x) = M(y) \Rightarrow x = y \quad (3-7)$$

(c) Now to prove the bound of the fibonacci number we shall look again at the multiplication, as can be seen and was proved in the last example, each element in the matrix. We shall also prove this by induction.

- i. For strings of length 1 it is true since we can see that all entries of $M(1), M(0)$ are within 1 , the 1^{st} fibonacci number.
- ii. Now let us assume that for all strings of length $\leq k$ we will be able to find x given $M(x)$.
- iii. Now let us assume that

$$M(x) = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad (3-8)$$

Now given $z = 0, 1$ we can have

$$M(x) = \begin{bmatrix} a+b & b \\ c+d & d \end{bmatrix} \quad (3-9)$$

or

$$M(x) = \begin{bmatrix} a & a+b \\ c & c+d \end{bmatrix} \quad (3-10)$$

In both these cases the numbers are the sum of previous elements, and we also have proved that either $a > b$ and $c > d$ or vice versa, holds, so we can see that a, b are upper bounded by consecutive fibonacci terms, and hence their sum is bounded by the next fibonacci term.

One Dimensional Pattern Matching

Consider we have two string x of length n which is the original string and y of length m which is the pattern which we need to look for in x . We know that $n \geq m$.

The deterministic algorithm is as follows.

First we compute $M(z) \bmod p$ of the first m characters of x . Now for each new character we first divide by the $M(0)$ or $M(1)$ corresponding to the character m steps before, and then multiply with the next character, and check this with $M(y)$.

$$M(0)^{-1} = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix} \quad (3-11)$$

$$M(1)^{-1} = \begin{bmatrix} 1 & 0 \\ -1 & 1 \end{bmatrix} \quad (3-12)$$

$$M(z) = (M[z - m]^{-1} * M(z - 1) * M[z]) \bmod p \quad (3-13)$$

Here $M(z)$ is the fingerprint of the last m characters till z^{th} character. and $M[z]$ is the fingerprint of the z^{th} character.

If we do not have $\bmod p$ the algorithm will always give the correct solution, and also we can see that it has very low incremental cost. Now if we select p from range $0..i$ then we have the number of factors of $M(x)$ is restricted by m where m is the length of the pattern hence .

Probability of false positive = Probability that p has a factor in m hence

$$Pr[Error] = (m \log i) / i \quad (3-14)$$

If $i = cm \log m$ we have

$$Pr[Error] = 1/c \quad (3-15)$$

Two Dimensional Pattern Matching

Here for the first $m \times m$ entries we compute it deterministically, now for each of the other entries

$$A(i, j) = (A[z - m, j]^{-1} * M(i - 1, j) * M[i, j]) \bmod p \quad (3-16)$$

This problem is infact no different from matching a m^2 length string, from the point of view of the error probability. Again the probability holds except that here the length is $= m^2$

$$Pr[Error] = (m^2 \log i) / i \quad (3-17)$$

If $i = 2cm^2 \log m$ we have

$$Pr[Error] = 1/c \quad (3-18)$$

→ Answer

Problem 4. Approximate Ham Sandwich Cut Algorithm

- (a) Randomly and uniformly sample p red point and p blue points.
- (b) We shall now find the ham sandwich cut of these points using brute force, that is for all pairs of points we will check if the line passing through them is the ham sandwich cut or not, since if any other line which is a ham sandwich cut can be moved and rotated such that it passes through two of the points without disturbing any other points.
- (c) Hence for each pair of $\binom{2p}{2}$ points find the line, and check if it is the ham sandwich cut. Each point can be verified in $O(1)$ time, and hence for each line we have $O(p)$ time, hence total time taken to find ham sandwich cut of $2p$ points $= O(p^3)$.
- (d) Now we shall say that this is the approximate ham sandwich cut of the whole set of points, and based on the value of ϵ we shall fix the value of p

Since the complexity is $O(p^3)$ we shall select upto $O(n^{\frac{1}{3}})$. So we set $p = c \log n$

For the proof we shall draw parallels between this problem and that of finding the approximate median of a given sample, Now let the line found by the equation be $y = mx + c$ we shall now rotate all the point about the origin such that the line is now given by $x = C$, now we can easily see that C has to be the approximate median of the x co-ordinates of both red and blue points.

We can also note that after the rotation we need not bother about the y co-ordinates since they are partitioned only by the x co-ordinate.

We shall now consider the red points alone, we know that \mathbf{C} is the median of the \mathbf{x} co-ordinates of \mathbf{p} red points sampled. Now we shall look at the probability that this is a $(1 + \epsilon)$ approximation of the median of \mathbf{x} co-ordinates of the points. Let \mathbf{R} be the set of Red points and \mathbf{B} be the set of blue points. Let \mathbf{R}_S be the set of sampled points of \mathbf{R} and similarly \mathbf{B}_S . We shall now define the random variable

$$\mathbf{X} = \text{The no. of elements in } \mathbf{R}_S \text{ which have rank } \leq (1 - \epsilon) \frac{n}{2} \quad (4-1)$$

Since we can see that each point is independently chosen without considering the other choices we have

$$\mathbf{X}_i = \begin{cases} 1 & \text{if the } i^{th} \text{ sampled element has rank } \leq (1 - \epsilon) \frac{n}{2} \\ 0 & \text{otherwise} \end{cases} \quad (4-2)$$

(4-3)

$$\mathbf{X} = \sum_{i=1}^p \mathbf{X}_i \quad (4-4)$$

$$\mathbf{E}[\mathbf{X}] = \sum_{i=1}^p \mathbf{E}[\mathbf{X}_i] \quad (4-5)$$

$$= \sum_{i=1}^p \frac{(1 - \epsilon) \frac{p}{2}}{p} \quad (4-6)$$

$$= (1 - \epsilon) \frac{p}{2} \quad (4-7)$$

Since the median is a very symmetric condition we can see that

$$\mathbf{P}[\text{Algorithm is wrong}] = 2\mathbf{P}[\mathbf{X} > \frac{p}{2}] \quad (4-8)$$

We shall later bound this to a suitable value to get a very high probability.

$$\mathbf{P}[\mathbf{X} > \frac{p}{2}] = \sum_{i=\frac{p}{2}}^p \binom{p}{i} \left(\frac{1 - \epsilon}{2}\right)^i \left(1 - \frac{1 - \epsilon}{2}\right)^{p-i} \quad (4-9)$$

The following conversion is possible since $(\frac{1}{2})^p$ is the maximum value for the product.

$$\mathbf{P}[\mathbf{X} > \frac{p}{2}] \leq \sum_{i=\frac{p}{2}}^p \binom{p}{i} \left(\frac{1}{2}\right)^p \quad (4-10)$$

Applying Sterlings approximation and the fact that $\frac{p}{i} \leq 2$ we get.

$$\mathbf{P}[\mathbf{X} > \frac{p}{2}] \leq \sum_{i=\frac{p}{2}}^p \frac{e^p}{4} \quad (4-11)$$

$$\mathbf{P}[\mathbf{X} > \frac{p}{2}] \leq \frac{e^p}{4} \quad (4-12)$$

Setting $p = c \log n$

$$P[X > \frac{p}{2}] \leq C' n^{-c \log 4/e} \quad (4-13)$$

$$P[\text{ErrorRed}] = 2 * P[X > \frac{p}{2}] = 2 \times C' n^{-c \log 4/e} \quad (4-14)$$

Asymptotically this value approaches zero hence we have very high probability bound, Similarly we can get the same bound for the set of Blue points hence total error probability

$$P[\text{Error}] \leq P[\text{ErrorRed}] + P[\text{ErrorBlue}] \quad (4-15)$$

By Union Theorem

$$P[\text{Error}] = 4 \times C' n^{-c \log 4/e} \quad (4-16)$$

Hence the running time of our algorithm is $O(\log^3 n)$

Since we will anyway use $O(n)$ time to read the input, if we set $p = n^{\frac{1}{3}}$ then we shall actually get much better bounds without losing out on the time complexity

→ Answer

Submitted by Abhimanyu M A (11111002) , Sumesh T A (11111065) on September 26,2011.