

# G22.3520: Honors Analysis of Algorithms

## Problem Set 4+5

### Solutions

#### Problem 1

a) Perform binary search on each full array until the element is found. The worst-case run time is to do binary-search on each of the  $\log(n)$  arrays which takes at most  $\log^2(n)$  time.

b) The insertion is quite similar with the incrementing binary counter. When we insert an element, the binary representation of  $n$  changes exactly the way the binary counter does. For all  $t$  digits on the right that change from 1 to 0, we merge the  $t$  arrays (pairwise) to get an array of size  $2^{t+1}$  which we associate with the  $t + 1$ st digit which now becomes 1. The cost of this insertion is  $O(2^t)$  (recall: merging two sorted arrays of size  $m$  takes  $O(m)$  time). The worst case occurs when  $n = 2^{t+1}$  after the insertion (a 1 followed by all 0s in the binary representation). The cost in this case is  $O(n)$ .

For the amortized cost we define the following *negative* potential function:  $\Phi(D) = -\sum_{i=1}^k i|A_i|$  where  $k = \log(n)$  and  $|A_i|$  is the size of the array corresponding to bit  $i$  in the binary representation of  $n$  (i.e.  $|A_i| = 0$  or  $|A_i| = 2^{i-1}$ ). For any operation  $j$ , the true cost of the operation is  $c_j = O(2^t)$  and the change in potential is

$$\begin{aligned}\Phi(D_j) - \Phi(D_{j-1}) &= -\sum_{i=1}^k i|A'_i| - a(-\sum_{i=1}^k i|A_i|) = a\sum_{i=1}^k i(|A_i| - |A'_i|) \\ &= \sum_{i=1}^t i2^{i-1} - (t+1)2^t \\ &\leq t\left(\sum_{i=1}^t 2^{i-1}\right) - (t+1)2^t \\ &\leq -2^t\end{aligned}$$

where  $|A'_i|$  are the new array lengths and  $|A_i|$  are the old array lengths and  $t$  is the number of 1s to the right that change to 0. Therefore the accounted cost of an operation is  $\hat{c}_j = c_j + \Phi(D_j) - \Phi(D_{j-1}) = 0$  (we need to scale  $\Phi$  by some appropriate *constant* to make this work). Recall that the total cost of any  $n$  operations is:

$$\sum_{i=1}^n c_i \leq \sum_{i=1}^n \hat{c}_i + \Phi(D_0) - \Phi(D_n) \leq -\Phi(D_n) \leq \sum_{i=1}^k i|A_i| \leq n \log(n).$$

Therefore the amortized cost per operation is  $\log(n)$ .

Alternatively, we can use the aggregate method to show an  $O(\log n)$  amortized upper bound. As we described above, any insertion takes time  $2^{t+1}$  if we have  $t$  consecutive 1s in the least significant bits of the binary representation of  $n$ . This happens though at most  $\frac{n}{2^{t+1}}$  times for a sequence of  $n$  insertions. This means that the total cost of  $n$  operations is bounded by  $\sum_{t=0}^k \frac{n}{2^{t+1}} 2^{t+1}$  which is  $O(nk) = O(n \log n)$ . If we divide with the total number of  $n$  operations we get  $O(\log n)$  amortized cost per insert operation.

c) In order to implement the delete operation, we first search for the element to be deleted following the steps that we showed above. Then, if we find this element  $x$  in array  $A_i$  we remove it. Let  $t$  be the right-most bit which is a 1 in the binary representation of  $n$ . Then take an arbitrary element  $y$  from  $A_t$  and insert it into  $A_i$  in the appropriate position (so the number of elements in  $A_i$  will not change). Now take the array  $A_t$  and split it into  $t - 1$  sorted arrays of sizes  $1, 2, 4, 2^{t-1}$  respectively. These arrays correspond to the  $t - 1$  right-most bits in the representation of  $n - 1$ , which are now 1. This operation may take  $O(n)$  time - for example if  $n$  is a power of 2.

#### Problem 2

a) Perform in-order traversal of the tree rooted at  $x$  to get a sorted array of elements. Then we build a balanced tree out of this array as follows: choose the median element and make it a root. Then recursively

build a left subtree from the portion of the array that's lower than the median, and a right subtree from the portion of the array which is greater than the median.

b) If  $T(n)$  is the amount of time that it needs in order to search a  $n$ -node  $\alpha$ -balanced binary search tree, then we have  $T(n) = 1 + T(\alpha n)$ . If  $\alpha$  is a constant then  $T(n) = O(\log n)$  (for example, by the master theorem).

c) Since  $c$  must be a positive constant and given that  $\Delta(x)$  is never negative (absolute value), we conclude that the potential can't take negative values. By the definition of  $1/2$ -balanced trees we know that for every vertex  $x$  we have  $size[right[x]] \leq \frac{1}{2}size[x]$  and  $size[left[x]] \leq \frac{1}{2}size[x]$ . If there existed a vertex  $x$  for which  $\Delta(x) \geq 2$ , then one of the two children, say  $right[x]$  would have at least two elements more than  $left[x]$ . This would mean that  $2size[right[x]] \geq size[x] + 1$  which in turn means that  $size[right[x]] \geq \frac{1}{2}size[x] + \frac{1}{2}$  which is a contradiction. Therefore  $\Delta(x) < 2$  and  $\Phi(T) = 0$ .

d) The amortized cost is  $\hat{c}_i = c_i + \Phi(D_i) - \Phi(D_{i-1})$  and we want it to be  $O(1)$ . Since after each rebuild we have  $\Phi(D_i) = 0$  and the rebuilding costs  $m$ , we want  $O(1) = m - \Phi(D_{i-1}) \Rightarrow \Phi(D_{i-1}) \geq m$ . We want to find the minimum value of the potential that could cause the need for rebuilding. This would mean that one of the subtrees, say the right one, would violate the  $\alpha$ -balanced constraint and therefore

$size[right[x]] > \alpha m$ . We also know that  $size[right[x]] + size[left[x]] = m - 1$  so

$size[left[x]] < m - 1 - \alpha m$ . In this case, we have

$\Delta(x) = size[right[x]] - size[left[x]] > \alpha m - (m - 1 - \alpha m)$  and therefore  $\Delta(x) > 2\alpha m - m + 1$ .

Summing up all the above we want  $m \leq c(\alpha m - m + 1) \Rightarrow c \geq \frac{m}{\alpha m - m + 1} \geq \frac{1}{2\alpha}$ .

e) Now that we have shown how to rebuild a subtree in  $O(1)$  amortized time, we only need to worry about the actual cost of an insertion or deletion and about the increase of the potential function. It is easy to see that the actual cost of these operations is  $O(\log n)$  just like we showed the same upper bound for searching above. What remains to be shown is an upper bound for the increase of the potential function after some insertion or deletion. Such an operation only affects the potential function of the vertices that belong to the traversed path from the root. To be more specific, the worst case is when the new element is added to the biggest subtree or deleted from the smallest subtree and therefore leads to an increase of  $\Delta(x)$  by one. Since there are only  $O(\log n)$  vertices that get affected though, this means that the potential function cannot increase by more than  $O(\log n)$  and this proves that the amortized time complexity of both operations is  $O(\log n)$ .

### **Problem 3**

1) We start from some arbitrary node and perform depth first search<sup>1</sup>. At each step of this process, we mark the node that we have just visited and then we check all the "unused" edges from its adjacency list. If one of these edges  $(u, v)$  leads to a node  $v$  which is already marked, this can only mean one thing. Since this edge was not already used although node  $v$  was marked, then node  $v$  is a predecessor of  $u$  in the DFS tree. This makes it very easy to output the cycle by just traversing the path of the predecessors of  $u$  until we reach  $v$ . It is obvious that if a cycle exists within the graph, then a DFS search would inevitably visit the same vertex twice and since we check every vertex and every edge once, the running time is  $O(m + n)$ .

2) They can just assign each butterfly to a vertex in a graph and connect the corresponding vertices for each judgment either with a blue edge if they are believed to be the same or with a red edge if they are believed to be different. We start by marking some arbitrary vertex with the letter A, implying that the corresponding butterfly belongs to this species. We then perform breadth first search, marking all adjacent vertices with different letters from  $\{A, B\}$  if they are connected with red edges and with the same letters otherwise. If at some point during the execution of the algorithm we try to label some vertex with the opposite letter from the one that it's already labelled with, the algorithm rejects the input. If this algorithm terminates successfully, it returns a consistent classification of the specimens based on the consistent judgements.

### **Problem 4**

We begin by creating a directed graph which has two vertices for each person  $P_i$ , one that corresponds to his birth: " $P_i$  born" and one to his death: " $P_i$  dead". We add a directed edge from " $P_i$  born" to " $P_i$  dead"

<sup>1</sup>We follow the exact same procedure for every connected component

implying chronological order. For each fact of type “ $P_i$  died before  $P_j$  was born”, we add a directed edge from “ $P_i$  dead” to “ $P_j$  born”. For each fact “the lives of  $P_i$  and  $P_j$  overlap” we add one directed edge from “ $P_i$  born” to “ $P_j$  dead” and another edge from “ $P_j$  born” to “ $P_i$  dead”.

We then do a topological sort and quit if we find a cycle. It is clear that, if a cycle exists, then there is an inconsistency (all edges should point to future times). Otherwise, if no cycle exists then a topological sort orders the sequence of births and deaths in a way that is consistent with all facts.

### **Problem 5**

The initial residual graph is the same as the flow network. The Ford-Fulkerson algorithm could choose the  $(s, a, b, t)$  augmenting path and push a flow of 1 through it. Then it could do the same thing for the path  $(s, b, a, t)$ , therefore, after two steps, the residual graph would be similar to the initial one with 999 instead of 1000. This sequence could continue for another 1998 steps until the residual capacities of the four peripheral edges become zero, so the running time would be 2000 steps.

### **Problem 6**

See <http://www.cs.nyu.edu/courses/fall06/G22.3520-001/lec20.pdf> slides 8-11.

### **Problem 7**

1) Every edge of the matching  $M$  must have at least one of its vertices in  $S$ . Therefore, every edge in  $M$ , can be mapped to some *unique* vertex in  $S$ .

2) A fully connected graph of three vertices has matchings of size 1 but its vertex cover must be of size at least two.

3) In order to verify that  $(U \setminus A) \cup B \cup C$  is a vertex cover, we notice that all possible edges have at least one of their endpoints in this set. All possible edges must either start from  $(U \setminus A)$  and end in  $W$  or start from  $U$  and end in  $B$  or start from  $A$  and end in  $C$ .

It is easy to see that the size of this set is equal to the number of edges that cross the minimum cut. We have exactly  $|U \setminus A|$  edges from  $s$  to  $(U \setminus A)$ , exactly  $|B|$  edges from  $B$  to  $t$  and exactly  $C$  edges from  $A$  to  $C$ . Therefore, this is the capacity of the cut and it is equal to the value of the maximum flow. Since we already know that the value of the maximum flow equals the size of the maximum matching, we conclude that it suffices to just run the Ford-Fulkerson algorithm for the new graph.

### **Problem 8**

We will use a randomized algorithm that colors each vertex uniformly at random. For each edge of the graph, the probability that this edge is satisfied is  $\frac{2}{3}$  since the only case when it is not satisfied is if both its vertices are colored with the same color. If  $X_{i,j}$  is the indicator variable that is 1 if edge  $(i, j)$  is satisfied and 0 otherwise, then  $E[X_{i,j}] = \frac{2}{3}$ . The expected number of edges that our algorithm satisfies is therefore equal to  $E[X] = E[\sum_{(i,j) \in E} X_{i,j}]$  which by linearity of expectation is equal to

$$\sum_{(i,j) \in E} E[X_{i,j}] = \sum_{(i,j) \in E} \frac{2}{3} = \frac{2}{3}|E| \text{ which is obviously at least } \frac{2}{3}c^*.$$

### **Problem 9**

Let us order the bids  $b_1, \dots, b_n$  in *decreasing* order. If some bid  $b_k$  modifies the value of  $b^*$ , it must be chosen before any of the  $k - 1$  bids  $b_1, \dots, b_{k-1}$ . There are  $k!$  ways to order the first  $k$  bids and  $(k - 1)!$  of these orderings have her bid  $b_k$  in the first place. This means that the probability that the bid  $b_k$  results in an update of  $b^*$  is equal to  $\frac{1}{k}$ . If  $X_k$  is the indicator variable which is 1 if the bid of  $b_k$  affects  $b^*$  and 0 otherwise, we have that  $E[X_k] = \frac{1}{k}$ . The expected number of changes of  $b^*$  is therefore equal to  $E[X] = E[\sum_{i=1}^n X_k]$  which by linearity of expectation is equal to  $\sum_{i=1}^n \frac{1}{k} = \Theta(\log n)$ .

### **Problem 10**

a) The probability that some particular machine will still remain idle is equal to  $(\frac{k-1}{k})^k$ . Therefore, if we let a random variable  $X_i$  be an indicator variable for machine  $i$  being idle, then  $E[X_i] = (\frac{k-1}{k})^k$  and  $N(k) = E[\sum X_i] = k(\frac{k-1}{k})^k$  and therefore  $\frac{N(k)}{k} = (\frac{k-1}{k})^k$ . The limit of this fraction as  $k$  goes to infinity is  $\frac{1}{e}$ .

b) Since the number of jobs that are assigned is equal to the number of machines, the number of jobs that get rejected must be equal to the number of machines that remain idle. Based on the previous subproblem we conclude that the limit is again equal to  $\frac{1}{e}$ .

c) We first calculate the expected number of rejected jobs for each machine and then add them to find the value of  $R_2(k)$ . Since the size of the buffer is now 2, in order to have  $i$  rejected jobs from some machine, there must have been exactly  $i + 2$  jobs assigned to it. If we count all possible ways that this can happen we get that the expected number of rejections per machine is equal to

$$\sum_{i=3}^k (i-2) \binom{k}{i} \left(\frac{1}{k}\right)^i \left(1 - \frac{1}{k}\right)^{k-i} = \sum_{i=3}^k i \binom{k}{i} \left(\frac{1}{k}\right)^i \left(1 - \frac{1}{k}\right)^{k-i} - 2 \sum_{i=3}^k \binom{k}{i} \left(\frac{1}{k}\right)^i \left(1 - \frac{1}{k}\right)^{k-i}$$

We notice that the first part of the sum resembles the expected value of the binomial distribution and the second part resembles the probability of the binomial distribution. For the specific case that we are working on, both these would be equal to 1 and we only need to deal with the missing values for  $i = 0, 1, 2$ . After some computations we conclude that the expected number of rejections per machine is equal to  $-1 + 2\left(1 - \frac{1}{k}\right)^k + \left(1 - \frac{1}{k}\right)^{k-1}$ . By once again using the same arguments as above we conclude that  $\lim_{k \rightarrow \infty} \frac{R_2(k)}{k} = \frac{3}{e} - 1$ .

### **Problem 11**

See <http://www.cs.nyu.edu/courses/fall06/G22.3520-001/lec15.pdf> for an  $O(|V| + |E|)$  algorithm. Here we give a simpler poly-time algorithm.

To begin with, we perform a depth first search from each vertex and we aggregate all the child vertices that we reached to the root after backtracking. This way, in polynomial time, for each vertex we have a list of all the vertices that can be reached from it through directed paths. We pick some arbitrary vertex  $u \in V$ . For every vertex  $v$  in the reachability list of  $u$ , we check if  $u$  also exists in the reachability list of  $v$ . We put  $u$  along with all the vertices of its list which prove to be equivalent with it in a set  $P_1$ . We delete all the vertices of  $P_1$  from  $V$  and continue by picking an arbitrary next vertex and repeating the above process, defining one new set  $P_i$  each time.

Now that we have a decomposition of the graph into strongly directed components, we assume that there exists some directed cycle in the induced graph  $G'$ . It is easy to see that if we have an edge in  $G'$  between two sets  $P_i$  and  $P_j$ , this means that there exists a directed path from every vertex of the first set to every vertex of the second set. This obviously means that the existence of a directed cycle in  $G'$  implies the existence of a directed cycle in  $G$  which passes through vertices that are not equivalent. This can't be true though, since all vertices that belong to a directed cycle must be equivalent with each other and therefore  $G'$  must be acyclic.