

Algorithms and Data Structures Disjoint sets

Marco Pellegrini

Spring 2002, Lecture 9

4/2/03

1

Summary of lecture 9

- Summary of this lecture: (CLR 22.1-3)
Data Structures for Disjoint sets:
Union operation
Find operation

4/2/03

2

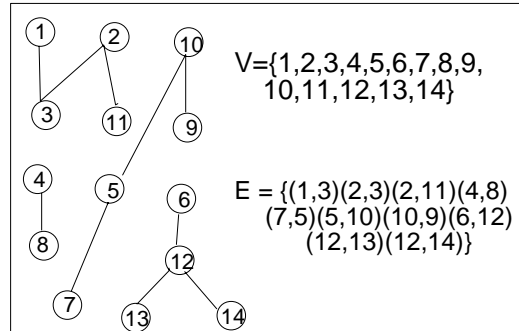
Disjoint set operations.

- We have a collection of elements, which we identify with numbers 1,2,3, etc.
- A Collection of sets S_1, S_2, \dots, S_k .
- The property we assume is that an element is at most in one set. That is the sets are disjoint. We shall see the following operations:
- Make-set(x): take a element x not in any set and generates the set containing x, that is: {x}.
- Union(x,y): makes the union of the set containing x and the set containing y.
- Find-set(x): returns a pointer to the representative of the set containing x. Each set has a unique representative.

4/2/03

3

Example: connected components in graphs



4/2/03

4

Connected-components.

- We want to find the components of the graph. A component is a group of nodes that we can reach by following edges starting from a node.
- ```
Connected-component(V,E)
 FOR each vertex v in V DO
 Make-set(v)
 END FOR
 FOR each edge (u,v) in E DO
 IF Find-set(u) ≠ Find-set(v)
 THEN Union(u,v)
 END IF
 END FOR
```
- ```
Same-component(u,v)
  RETURN [Find-set(u) = Find-set(v)]
```

4/2/03

5

Exercise

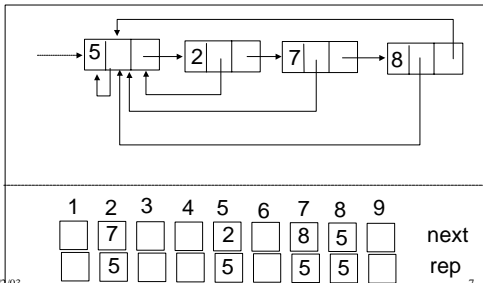
- Run the procedure Connected-component on the graph (that is V and E) given in the example.
- How many sets (that is components) do you have at the end.
- Hint: naturally I have not told you yet HOW make-set, union and find-set work but you know WHAT they are supposed to do.

4/2/03

6

Representing disjoint sets by lists.

- We can use circular singly linked lists and use the first element of the list as representative.



4/2/03

7

Operations

- Make-set and find-set can be done in $O(1)$ on a circular list with back-pointer to the head.
- Union(x,y) is trickier. Appending two circular lists can be done in $O(1)$, but to update the pointer to the representative we must scan one of the two lists.
- If n is the number of elements, m the number of operations, in the worst case we might use time $O(m^2)$ to perform the m operations.
- This happens when we grow a list by one element at a time and we always scan the longest list.

4/2/03

8

Size heuristic.

- Improvement: when we make the union scan the shortest of the two list.
- To know which is the shortest we carry a field SIZE linked to the representative which we update when we perform the union.
- This simple rule decreases the time complexity to $O(m + n \log n)$ for n elements, m operations.
- Proof. By a charging argument. Fix one element x , initially it is in a set of size 1. Then at every union the size of the set containing it at least doubles, that is $2, 4, 8, \dots, n$. This can happen only $\log n$ times. At every union we update $\text{repr}(x)$. In total $n \log n$ updates. All other costs are $O(m)$.

4/2/03

9

Homework

HW2.3 (CLR 22.2-1)

Write pseudocode for make-set, find-set, union using singly linked lists and the weighted union rule. Each object x , has:

- field $\text{repr}[x]$ pointing to the representative of the set containing x ,
- field $\text{last}[x]$ pointing to the last object in the list containing x ,
- field $\text{size}[x]$ giving the size of the list containing x .

$\text{Size}[x]$ and $\text{Last}[x]$ are correct only when x is a representative.

4/2/03

10

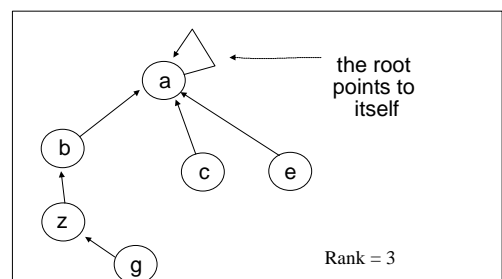
Forest-of-trees implementation.

- Can we do better than the linked lists DS?
- Yes, but it is hard to prove it!
- Idea: represent a set by a tree, we need only a pointer from a node to the parent.
- We could keep track of sizes, but for proof purposes we keep track of the *rank* of the trees.
- The rank is the length of the longest path.

4/2/03

11

Examples

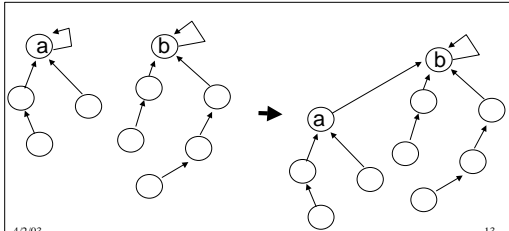


4/2/03

12

Union by rank

- The rank of the tree is the length of the longest path from root to leaf (without path compression)
- Union by rank is: attach the lowest rank tree under the root of the highest rank tree.



4/2/03

13

Union by rank

- If we merge two trees A,B of different rank the rank of the resulting tree is $\max(\text{rank}(A), \text{rank}(B))$.
- If they have the same rank, then the rank increases by 1.

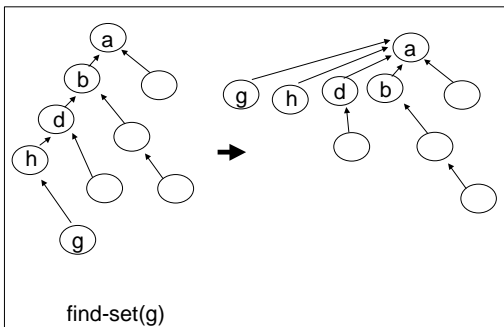
```
Make-set(x)
parent[x] := x
rank[x] := 0
```

```
Union(x,y)
a := find-set(x), b:= findset(y)
IF rank[a] > rank[b]
THEN parent[b] := a
ELSE parent[a] := b
IF rank[a] = rank[b]
THEN rank[b] := rank[b]+1
END IF
END IF
```

4/2/03

14

Find-set with path compression.



4/2/03

15

Find set + path compression.

- To implement $\text{find-set}(x)$ we follow pointers from x to the root and then we go back making every node we visit on the path a child of the root.

```
Find-set(x)
IF x ≠ p[x]
THEN parent[x] := Find-set(parent[x])
ENDIF
RETURN parent[x]
```

- Exercise: simulate this code on the example of the previous slide $\text{find-set}(g)$. Show for each call the node in input and the one in output and the change to the parent field.

4/2/03

16

Programming assignment

- Implement in C the Disjoints-sets DS using lists with weighted union rule. [DS1]
- Implement in C the Disjoint-set DS using trees with the rank rule and the path-compression. [DS2]
- Generate random Graph G with n nodes and n edges.
- Implement the connected components algorithm twice using DS1 the first time and DS2 the second time.[A1,A2]
- Run the A1 and A2 on the graph G (check that the output (that is the components) is the same.
- Plot the running times of A1 and A2 for $n=100,200,500,1000$, maximum number before the program crashes.

4/2/03

17

Performance of Union-Find DS

- n elements, m operations (union and find), f find operations.
- UF= Union find, PC= path compression, WH= weight heuristic, RH= rank heuristic.
- UF+WH is $O(m + n \log n)$.
- UF+RH is $O(m \log n)$
- UF+PC is $O(n + f \log n)$ if $f < n$, $O(f \log_{(1+f/n)} n)$ if $f > n$.
- UF+PC+RH is $O(m \log^* n)$

Where $\log^* n = \min\{i | \log^{(i)} n \leq 1\}$. And $\log^{(0)} n = \log n$,
 $\log^{(i)} n = \log(\log^{(i-1)} n)$.

4/2/03

18

$$F(i) = 2^{2^{2^{\dots^2}}}$$

i times

$\text{Log}^*(F(x)) = x.$

$F(1) = 1$

$F(2) = 4$

$F(3) = 16$

$F(4) = 65,536$

$F(5) = 2^{65536} > 10^{80}$

4/2/03
19

Analysis of UF+PC+RH

The function $G(n) = \log^* n$ is the pseudo-inverse of $F(n)$.

For all practical inputs $\log^* n < 5$.

The main idea is that to estimate the number of nodes *visited* by find operations, we can instead count the number of find operations *visiting* any node.

4/2/03
20

Lemma 1. For a node x that is not root, $\text{rank}[x] < \text{rank}[\text{parent}[x]]$.

Proof: by construction.

Lemma 2. For a node x , $\text{rank}[x]$ is an increasing function of time (non-decreasing). Constant when x ceases to be a root.

Proof: by construction.

4/2/03
21

For a root x , $\text{size}(x)$ is the number of nodes in the tree rooted at x .

Lemma 3. $\text{Size}(x) \geq 2^{\text{rank}(x)}$.

Proof: by induction. Initially each node is a root and each node has rank 0, so the lemma is true.

Inductive step: suppose $\text{rank}(x) < \text{rank}(y)$. After the union we have size at the root $\text{size}(x) + \text{size}(y) > \text{size}(y) \geq 2^{\text{rank}(y)}$. The rank of the root after the union is $\text{rank}(y)$.

Suppose $\text{rank}(x) = \text{rank}(y)$. After the union we have size at the root $\text{size}(x) + \text{size}(y) \geq 2^{\text{rank}(y)} + 2^{\text{rank}(y)} \geq 2^{\text{rank}(y)+1}$. The rank of the root after the union is $\text{rank}(y)+1$.

4/2/03
22

Lemma 4. For each integer r , there are *at most* $n/2^r$ nodes of rank r at any given time.

Proof. Since the rank is an increasing function over each path from leaf to the root, no two nodes of same rank r are one the ancestor of the other. So for x and y both of rank r the subtrees rooted at x and y are *disjoint*.

Each subtree has at least 2^r nodes, so there are at most $n/2^r$ nodes of rank r .

Lemma 5. No vertex has rank greater than $\log n$.

4/2/03
23

Charging scheme.

We split the possible ranks $r = [0, \dots, \lceil \log n \rceil]$ into groups.

Rank r goes into group $G(r)$.

Equivalently group j contains ranks $[F(j-1)+1, \dots, F(j)]$.

Case 1. If find visits x and $\text{group}(x) \neq \text{group}(\text{parent}(x))$, or x is the root, charge 1 to the find operation.

Case 2. If find visits x and $\text{group}(x) = \text{group}(\text{parent}(x))$ charge 1 to the node x . (x is not the root).

4/2/03
24

Costs

Case 1. A single find operation incur in case (1) when crossing the boundaries between two groups, so at most $G(\log n) = G(n) - 1$ times. So n find operations are charged $O(nG(n))$.

Case 2. While x and $\text{parent}(x)$ are in the same group g , since the rank is increasing, we can charge x at most $F(g) - F(g-1)$ times $< F(g)$.

4/2/03

25

How many nodes do we have in group g ?

$$\sum_{r=F(g-1)+1}^{F(g)} n/2^r \leq (n/2^{F(g-1)+1}) \sum_{i=0}^{F(g)-F(g-1)-1} 2^{-i} \leq (n/2^{F(g-1)}) = n/F(g)$$

So the total charge to group g is n . The number of groups is $G(n)$. The total charge by case 2 is $O(nG(n))$.

Theorem 6. Union Find with Path Compression and Rank Heuristic, n elements and operations takes time $O(n \log^* n)$.

4/2/03

26

Conclusions

- We have seen two DS for Disjoint-sets, sometimes called (Union-find ADS).
- The first based on lists, the second based on trees with union-by-rank and path-compression.
- The analysis tells us the second is asymptotically faster. Your programming assignment should test whether this happens also in experimental tests.
- Actually an even more complex analysis can prove that the right bound is $O(m \alpha(n,m))$ where $\alpha(n,m)$ is a function that is even slower than $\log^*(n)$.

4/2/03

27