

# Disjoint Sets

## Union Find / Quick Find: (Based on id's)

+  $\{$   $lst = [0, 1, 2, 3, 4]$  - implicit form

To start  $\{$   $(0) (1) (2) (3) (4)$  - explicit form

$\Omega(n)$   $\{$   $union(0, 2) \rightarrow (0) (1) (0) (3) (4)$   $\rightarrow$  for  $i$  in range( $len(lst)$ ):  
 $O(n)$   $\{$   $union(0, 4) \rightarrow (0) (1) (0) (3) (0)$  if  $lst[i] == lst[b]$ :  
 $lst[i] = lst[a]$

$\Omega(1)$   $\{$   $connected(2, 4) \rightarrow [st[2] == st[4]] = True$

## Quick Union: (Basic Tree, Based on root's)

+  $\{$   $lst = [0, 1, 2, 3, 4, 5, 6]$  - implicit form

To start  $\{$   $(0) (1) (2) (3) (4) (5) (6)$  - explicit form

$\Omega(1)$   $\{$   $union(0, 3) \rightarrow (5) (0) (1) (2) (4) \rightarrow [st[root(3)] = root(0)]$

$O(n)$   $\{$   $union(3, 1) \rightarrow (6) (5) (2) (4) \rightarrow [st[root(1)] = root(3)]$

$\Omega(1)$   $\{$   $root(6) \rightarrow (6) (5) (0) \rightarrow$  while ( $st[a] != a$ ):  
 $O(n)$   $\{$   $a = st[a]$  return  $a$

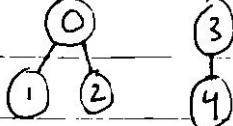
$\Omega(1)$   $\{$   $connected(5, 4) \rightarrow (6) (5) (3) (1) (2) (4) \rightarrow$  def  $connected(st, a, b)$ :  
 $O(n)$   $\{$   $root_a = root(st, a)$   
 $root_b = root(st, b)$   
 $return root_a == root_b$

\* Can do the same thing with union by height but has issues with integrating with path compression.

## Weighted Quick Union / Union by Size (Attach smaller to larger, Tree Based)

$\Omega(1)$

$O(n)$  { union(0, 3) }  
 $\Theta(n \log n)$



{ 3 }

{ 1 }

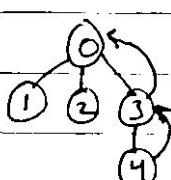
{ 2 }

{ 4 }

def union(lst, a, b)  
 root\_a = root[lst, a]  
 root\_b = root[lst, b]

$\Omega(1)$

$O(n)$  { root(4) }



def root(lst, a):

while lst[a] != a:  
 a = lst[a]

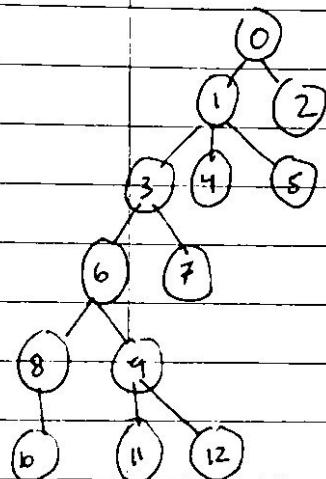
return a

if s2[root\_a] < s2[root\_b]:  
 lst[root\_b] = root\_a  
 s2[root\_a] += s2[root\_b]  
 else:  
 lst[root\_a] = root\_b  
 s2[root\_b] += s2[root\_a]

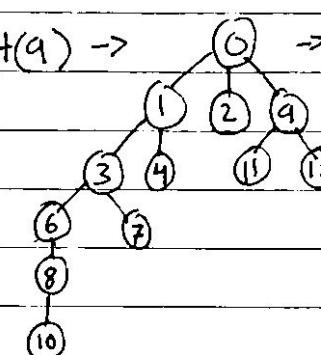
$\Omega(n)$

~~connected(0, 4)  $\rightarrow$  root(a) == root(b)~~

## Path Compression:



root(a)  $\rightarrow$



def root(lst, a):

while lst[a] != a:

lst[a] = lst[lst[a]]

a = lst[a]

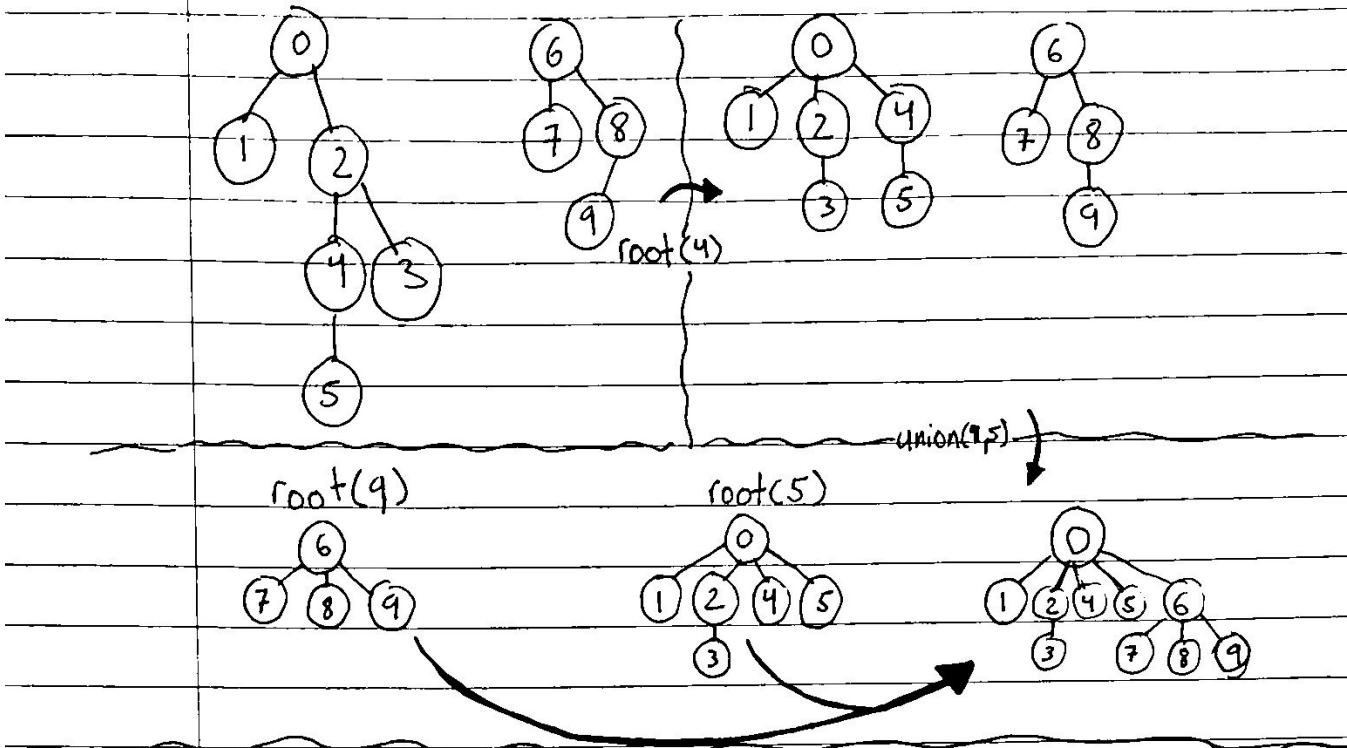
union  $\rightarrow$  def union(lst, a, b):

root\_a = root(a)

root\_b = root(b)

lst[root\_b] = root\_a

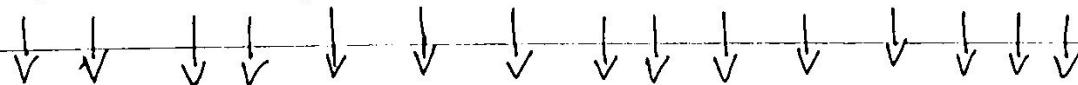
## Visualize Path Compression:



\* When we find the root of a node, we drag it up & attach it to the root directly to flatten the tree. This makes future look ups faster.

Note: root happens as well when we union two nodes, those calls to root() also help to flatten the tree.

\* Path compression can be used in collaboration with other methods to enhance it. A typical one is Weighted Path Compression as it works well together. However, union by height does not work well with path compression in its basic form since path compression changes the tree heights so you may have to recompute them. We can instead use an estimate of the height instead to solve this problem. When we do this estimation it becomes just as efficient as union by size in theory.



\* Using path compression the amortized cost of nearly constant. This is because path compression greatly limits the height of a tree especially after many root() calls. Even say that the tree is 4 deep then it takes 3 lookups to find the root which is essentially  $O(1)$ . So, finally  $O(4) = O(1)$ .

### Summary:

- The purpose of disjoint set problem is to efficiently support the root and union operations.

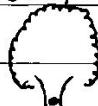
- Used to track elements partitioned into non-overlapping sub-sets.

- Implicit = the array representation

- Explicit = the node drawing representation

- 4 Main techniques: (Some may be combined)

1. Union-Find



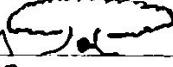
2. Quick-Union



3. Weighted Quick Union (union by size)



4. Path Compression



- With path compression & other optimization amortized cost is essentially constant as the tree is very flat.  $O(1)$

- Practical Applications

- Can be used to solve mazes

- Network connectivity

- Kruskal minimum spanning tree

## Algorithm Analysis

$$\cdot 1 < \log n < \sqrt{n} < n < n \log n < n^2 < n^3 < \dots < 2^n < 3^n < \dots < n^n$$

### Big Oh (Upper bound)

def:  $f(n) = O(g(n))$  iff  $\exists$  pos constants  $c \& n_0$ , such that  $f(n) \leq c * g(n) \ \forall n \geq n_0$

example:  $f(n) = 2n + 3$

$\sqrt{2n+3} \leq 10n \quad n \geq 1$

$f(n) \leq c * g(n)$        $2n+3 \leq 2n+3n$   
 $f(n) = O(n)$        $2n+3 \leq 5n$

\* Since big-oh is an upper bound if a function  $f(n) = O(n)$  then anything greater than that are true such as  $f(n) = O(n^2)$  or  $O(2^n)$  or  $O(n^n)$

$$\cdot 1 < \log n < \sqrt{n} < n < n \log n < n^2 < n^3 < \dots < 2^n < 3^n < \dots < n^n$$

### Big Omega (Lower bound)

def:  $f(n) = \Omega(g(n))$ , iff  $\exists$  pos constants  $c \& n_0$ , such that  $f(n) \geq c * g(n) \ \forall n \geq n_0$

example:  $f(n) = 2n + 3$

$$2n+3 \geq ln \quad \forall n \geq 1$$

$$f(n) \geq c * g(n)$$

$$f(n) = \Omega(n)$$

\* Since big omega is a lower bound if a function  $f(n) = \Omega(n)$  then anything smaller than that are true such as  $f(n) = \Omega(\sqrt{n})$  or  $\Omega(1)$ .

$$\log n < \sqrt{n} < n < n \log n < n^2 < n^3 < \dots < 2^n < 3^n < \dots < n^n$$

## Big Theta (average)

def:  $f(n) = \Theta(g(n))$ , iff  $\exists$  positive constants  $c_1, c_2$ , and  $n_0$ , such that  $c_1 * g(n) \leq f(n) \leq c_2 * g(n)$

example :  $f(n) = 2n + 3$

$$1 \times n \leq 2n+3 \leq 5 \times n$$

$$C_1 \times g(n) \leq f(n) \leq C_2 \times g(n)$$

$$f(n) = \Theta(n)$$

Average

$$1 < \log n < \sqrt{n} < n < n \log n < n^2 < n^3 < \dots < 2^n < 3^n < \dots < n^n$$

lower bayno

## Better Examples

$$f(n) = 2n^2 + 3n + 4$$

$$2n^2 + 3n + 4 \leq 9n^2 \rightarrow f(n) = O(n^2)$$

$$2n^2 + 3n + 4 \geq 1 \times n^2 \rightarrow f(n) = \Omega(n^2)$$

$$(1n^2 \leq 2n^2 + 3n + 4 \leq 9n^2 \Rightarrow f(n) = \Theta(n^2))$$

$$f(n) = n^2 \log n + n$$

$$\bullet \quad n^2 \log n \leq n^2 \log n + n \leq 10n^2 \log n$$

$f(n) = \Theta(n^2 \log n)$

$$f(n) = \Theta(n^2 \log n)$$

\* If we do not get same term like  $n^2$  or  $n^2 \log n$  on both sides we can not find average or  $\Theta$  of  $f(n)$ .

## Asymptotic properties

### General

- if  $f(n)$  is  $O(g(n))$  then  $\alpha * f(n)$  is  $O(g(n))$

example:  $f(n) = 2n^2 + 5 \rightarrow f(n) = O(n^2)$

$$\text{then } 7 * f(n) = 7 \times (2n^2 + 5) = 14n^2 + 35 = O(n^2)$$

- The same rule applies for omega as big-oh above.

- The same rule applies for theta.

### Reflexive

- if  $f(n)$  is given then  $f(n)$  is  $O(f(n))$

example:  $f(n) = n^2 \rightarrow O(n^2)$

### Transitive

- if  $f(n) = O(g(n))$  and  $g(n) = O(h(n))$

then  $f(n) = O(h(n))$

example:  $f(n) = n$     $g(n) = n^2$     $h(n) = n^3$

$n$  is  $O(n^2)$  and  $n^2$  is  $O(n^3)$

then  $n$  is  $O(n^3)$

### Symmetric

- if  $f(n)$  is  $\Theta(g(n))$  then  $g(n)$  is  $\Theta(f(n))$

example:  $f(n) = n^2$     $g(n) = n^2$

$f(n) = \Theta(n^2)$

$g(n) = \Theta(n^2)$

### Transpose Symmetric

- if  $f(n)$  is  $\Omega(g(n))$  then  $g(n)$  is  $\Omega(f(n))$

example:  $f(n) = n$     $g(n) = n^2$

then  $n$  is  $\Omega(n^2)$  and

$n^2$  is  $\Omega(n)$

## K-SUM K=2

"Obvious" K=2  $O(n^2)$

• check every index for a pair with every other index.

• nested for loops  $O(n^2)$ , returns duplicates without "smart iteration" method

def dumbtwosum(lst):

```

    res = []
    for i in range(len(lst)):
        for j in range(len(lst)):
            if i != j and lst[i] + lst[j] == 0:
                res.append([lst[i], lst[j]])

```

Smart Iteration K=2  $O(n^2)$

• simply change inner loop indexing to  $i+1$  instead of  $i$

```

    for i ... blah blah
        for j in range(i+1, len(lst))
            ...

```

• Still  $O(n^2)$  but will perform better, avoids duplicates

Sorting Pointers K=2  $O(n \log n)$

• sort first then use pointers at each end

def pointer\_ksum(lst):

{res, fp, ep, lst = [], 0, len(lst)-1, sorted(lst)}

while fp < ep:

currSum = lst[fp] + lst[ep]

if currSum == 0: look almost constant only because sorted hides cost well.

res.append([lst[fp], lst[ep]])

fp += 1

ep -= 1

elif currSum < 0:

fp += 1

else: ... - 1

Caching k=2  $O(n)$

- cache seen nums
- if we've seen inverse of curr number then we found a pair.
- takes a bit more space  
cache = set()

check if  $-lst(i)$  in cache

k-sum k=3

k=3 using k=2  $O(n^2)$

- duplicates included
- adjust k=2 sum to accept a target value not just 0.
- slice the array you give two sum so it doesn't include the current value

k=3 with 2sum cache, k=3 with pointer 2sum

- Both are  $O(n^2)$
- nothing to much to note here just combining the two functions