

# CS 251: Data Structures and Algorithms

Sum	Closed Form
$\sum_{k=0}^n ar^k \ (r \neq 0)$	$\frac{ar^{n+1} - a}{r - 1}, r \neq 1$
$\sum_{k=1}^n k$	$\frac{n(n + 1)}{2}$
$\sum_{k=1}^n k^2$	$\frac{n(n + 1)(2n + 1)}{6}$
$\sum_{k=1}^n k^3$	$\frac{n^2(n + 1)^2}{4}$
$\sum_{k=0}^{\infty} x^k,  x  < 1$	$\frac{1}{1 - x}$
$\sum_{k=1}^{\infty} kx^{k-1},  x  < 1$	$\frac{1}{(1 - x)^2}$

Figure 1: Formulas

## Big-O

Big- $\mathcal{O}$ :  $f(n) \in O(g(n)) \iff$

$\exists c > 0, \exists n_0 > 0 : \forall n \geq n_0, 0 \leq f(n) \leq c g(n).$

Big- $\Omega$ :  $f(n) \in \Omega(g(n)),$   
 $0 \leq c \cdot g(n) \leq f(n)$

Big- $\Theta$ :  $f(n) \in \Theta(g(n)),$   
 $0 \leq c_1 \cdot g(n) \leq f(n) \leq c_2 \cdot g(n)$

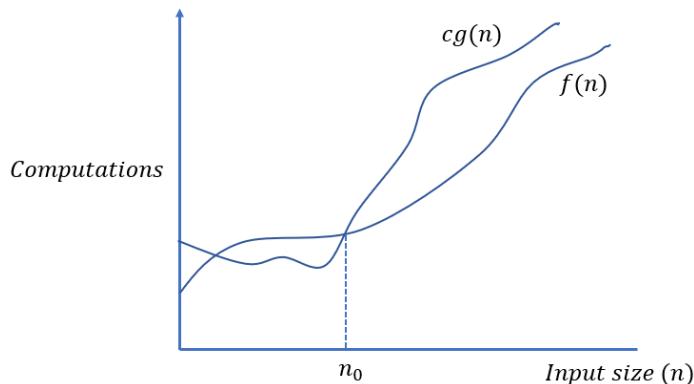


Figure 2: Big-O curve

## Arrays

Access:  $\Theta(1)$ , Insert:  $\Theta(1)$  (with tail)  $\Theta(n)$  (no tail).

Resize and space of resize:  $\Theta(n)$

## Binary trees

Max nodes in BT (total in complete BT) =  $2^{h+1} - 1$ . Leave nodes =  $2^h$

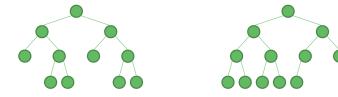


Figure 3: Full BT and Complete BT

→ Full BT: Each node has none or two children.

→ Complete BT: All levels full until last (left to right).

Balanced BT: At every node the height of the left and right subtree differs by at most 1:

$2^h \leq n \rightarrow h \leq \log_2 n \therefore h \in \mathcal{O}(\log(n))$ , where  $n$  is the number of nodes.

Balanced BT Insertions, Searches & Deletions:  $\mathcal{O}(\log(n))$

## Binary Heap

Complete binary tree  $\therefore h \in \Theta(\log(n))$ .

Insertions & Deletions:  $\mathcal{O}(\log(n))$  through single path.

→ Left:  $2i + 1$ , Right:  $2i + 2$ .

→ Parent:  $\lfloor \frac{i-1}{2} \rfloor$ .

Max heap (max on top). Min heap (min on top).

Insertions enter in complete BT order and then adjust (**sift up**).

Single insert is  $\mathcal{O} \log(n) \therefore n$  inserts  $\in \mathcal{O}(n \log(n))$

## Heapify

Sift down on all nodes in an array from  $[0, \lfloor \frac{n}{2} \rfloor - 1]$

Heap build:  $\mathcal{O}(n)$

## Heap sort

Call **Heapify** ( $\mathcal{O}(n)$ ) and then swap the root with the last item and fix the heap ("lock" the last element after swap).

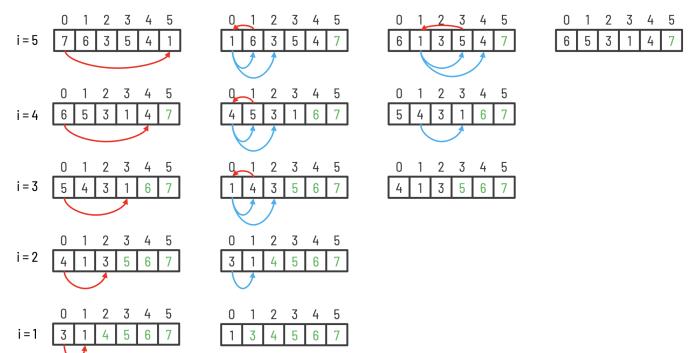


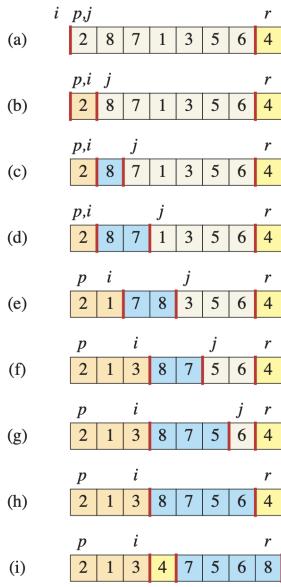
Figure 4: Heap sort example

## Quick sort

Not stable and pivot selection is arbitrary.

Best case ( $\mathcal{O}(n \log(n))$ ) → even partitions.

Worst case ( $\mathcal{O}(n^2)$ ) → one partition has  $n - 1$  elements.



**Figure 5:** Quick sort;  $i$  (return) is moved when comparison is lower and swaps at  $j$  (comparison index)

## Counting sort decision tree

Leaves are  $n!$  permutations of the array (all possibilities).

Internal nodes are the conditional checks.

As full BT will be  $2^h$  leaves  $\therefore 2^h \geq n!$

$h \geq \log_2 n! \Rightarrow h \geq n \log_2 n \therefore h \in \mathcal{O}(n \log(n))$

## Counting sort

Array → Freq. array → Cm. frq. array → Resulting array filled from cm. frq. array backwards.

Stable, bad with large max value of an array.

Runtime & Space:  $\mathcal{O}(n + k)$ , where  $n$ : array size;  $k$ : max value of array (freq. arr.).

## Bucket sort

Worst case  $\mathcal{O}(n^2) \rightarrow$  unsorted un-uniform distribution.

Best case  $\mathcal{O}(n)$  → one element per bucket or all sorted in one bucket.

Avg. case  $\mathcal{O}(n + \frac{n^2}{k} + k)$

If there is one key per bucket (i.e. there exists a bucket for any key) then,  $\mathcal{O}(n + k)$ , where  $n$ : elements in array;  $k$ : buckets.

Extra work comes from sorting within each bucket.

Better for uniform distributions: fits things equally.

## Radix sort

Sorting through LSD to MSD using a stable sorting algorithm (e.g. counting sort).

Runtime:  $\mathcal{O}(d(n+k))$ , where  $d$ : digits;  $n+k$ : counting sort.

Space:  $\mathcal{O}(n+k)$

Using bucket sort with binary digits, create 2 buckets for 2 keys (1 and 0).

If there are constant lengths of digits then  $\mathcal{O}(n)$  is possible.

E.g.  $\mathcal{O}(d(n+2)) \in \mathcal{O}(n)$  if there are constant lengths of binary numbers.

## Polynomial rolling hash

Where  $S$ : string of length  $m$ ;  $a$ : multiplier for polynomial.

$$H(S, a) = s_0 a^{m-1} + s_1 a^{m-2} + \dots + s_{m-1} a^0 = \sum_{i=0}^{m-1} s_i a^{m-i-1}$$

If using a fixed sized window  $m$  then do rolling updates on  $S$ :  $H_1 = (H_0 - s_{old} \cdot a^{m-old}) \times a + s_{new}$

## Hashing functions

Division method:  $h(k) = k \bmod m$  (avoid  $2^P$  and  $10^P = m$ ), where  $m$ : size of hash table.

Multiplication method:  $h(k) = \lfloor m(kA \bmod m) \rfloor$  ( $A \approx \frac{\sqrt{5}-1}{2}$ ;  $m$  is not critical)

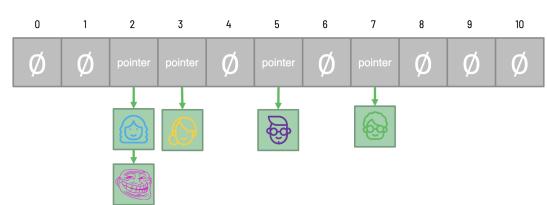
Aim to have keys be distributed uniformly (no clusters).

## Collision strategies

Under the Simple Uniform Hashing Assumption

→ Average number of keys per bucket (Load Factor):  $\alpha = \frac{n}{m}$   
→  $n \in \mathcal{O}(m) \therefore \alpha \in \mathcal{O}(1)$ ,  $n$ : keys,  $m$ : buckets.

### Chaining



**Figure 6:** Chaining collision management

Worst case  $\mathcal{O}(n)$ : traversing all items in the table.

Runtime  $\Theta(1 + LengthOfChain) = \Theta(1 + \frac{n}{m})$ .

## Open addressing

Probing sequence for  $h(k, i) \Rightarrow h(k, 0), h(k, 1), \dots, h(k, m-1)$

→ Linear:  $h(k, i) = (h(k) + i) \bmod m$ .

→ Quadratic:  $h(k, i) = (h(k) + c_1 i + c_2 i^2) \bmod m$ , where  $c$ : coefficients chosen.

→ Double hashing:  $h(k, i) = (h(k) + ih_1(k)) \bmod m$ .

Element deletion will mess up the hash (offset from probing impossible to reach).

When load factor is too high there is a cluster overhead.

Insertions probe:  $\frac{1}{1-\alpha}$

Successful search probes:  $\frac{1}{\alpha} \ln(\frac{1}{1-\alpha})$

Unsuccessful search probes:  $\frac{1}{1-\alpha}$

## Rehashing

### Chaining:

→ Goal:  $\alpha$  constant.

→ Double  $m$  when  $\alpha \geq 8$ .

→ Halve  $m$  when  $\alpha \leq 2$ .

### Open addressing:

→ Goal:  $\alpha \leq \frac{1}{2}$ .

→ Double  $m$  when  $\alpha \geq \frac{1}{2}$ .

→ Halve  $m$  when  $\alpha \leq \frac{1}{8}$ .

## B-trees

The order (num. of children) is  $m$ , determined from  $m - 1$  keys in the parent.

Height:  $h \in \mathcal{O}(\log_m n)$ .

B-trees are always "complete"—all levels full.

Insertions do not increase height unless  $m$  order limit is reached in the root.



Figure 7: 2-3 B-tree insertion

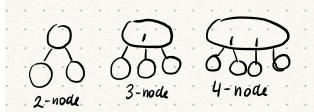


Figure 8: B-tree types

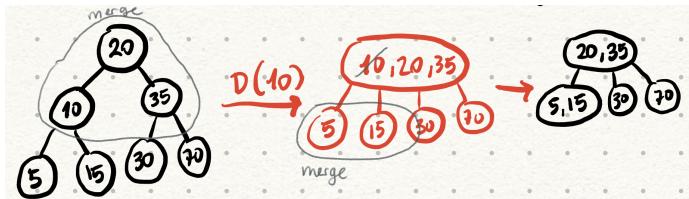


Figure 9: B-tree deletion with high redux

## Red-Black trees

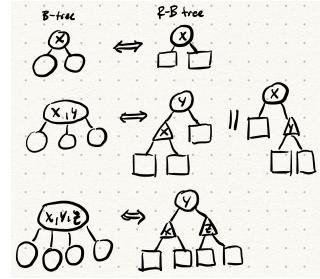


Figure 10: B-tree to R-B tree conversion

Root and null links are black.

Red parents have black children.

All paths, excluding null links, have the same number of black nodes  $h \in \mathcal{O}(\log(n))$  (black height).

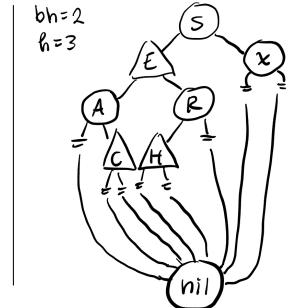
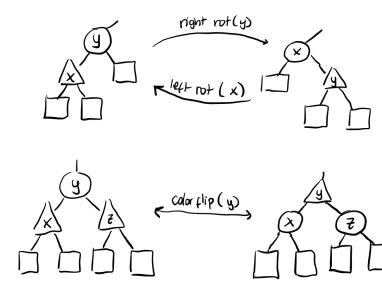


Figure 11: R-B Tree operations

Black height increase in an RB-tree, correlates to a height increase in a 4-node b-tree (2-3-4 tree).

## Undirected graphs

Handshake theorem:  $\sum_{v \in V} \deg(v) = 2 \cdot |E|$

Edge count:  $|E| \leq \frac{|V|(|V|-1)}{2}$ , considering  $|V| - 1$  edges on  $V$  edges.

Complete graph has vertexes with degree  $|V| - 1 \therefore |E| = \frac{|V|(|V|-1)}{2}$

## Paths and cycles

Simple path/cycle have no repeated edges/vertices.

Euler path visits all edges once.

Euler cycle is a Euler path that starts/ends on the same  $v$ .

Hamiltonian path visits all vertices once.

Hamiltonian cycle is a H. path that starts/ends on the same  $v$ .

Length is the number of edges.

## Edge representations

### Edge list

List:  $(u, v), (u, w), (v, x), \dots$

Adding edge:  $\mathcal{O}(1)$

Space, Adj. vertex check, & Adj. iteration:  $\mathcal{O}(|E|)$

### Adjacency matrix

Cells are 1 or 0, and  $A^k$  gives the num. of k-length paths between vertices.

Edges go  $i \rightarrow j$

	u	v	w	x	y	z
u	0	1	0	0	0	0
v	0	0	1	1	0	0
w	1	0	0	0	1	0
x	0	0	1	0	0	1
y	0	0	0	1	0	0
z	0	0	0	1	0	1

Figure 12: Adjacency matrix

Space:  $\mathcal{O}(|V|^2)$

Adj. iteration:  $\mathcal{O}(|V|)$

Adding edge & adj. vertex check:  $\mathcal{O}(1)$

### Adjacency list

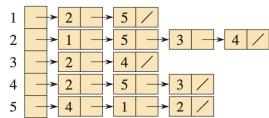


Figure 13: Adjacency list

Space:  $\mathcal{O}(|V| + |E|)$

Adding edge:  $\mathcal{O}(1)$

Adj. vertex check & Adj. iteration:  $\mathcal{O}(\deg(V))$

## Types of graphs

Connected graphs have paths for all pairs of vertices.

Spanning subgraph has all  $V$  in  $G$

Trees are undirected, connected, and acyclic graphs.

Forests are collections of trees (acyclic).

Spanning tree are spanning subgraphs of  $G$  that are trees.

## Graph characteristics

Sparse graphs  $\rightarrow |E| \ll |V|^2$  (use lists)

Dense graphs  $\rightarrow |E| \approx |V|^2$  (use matrices)

Strong connectivity is that every  $v$  is connected to every other  $v \in V$  (with reflexive, symmetric, and transitive properties).

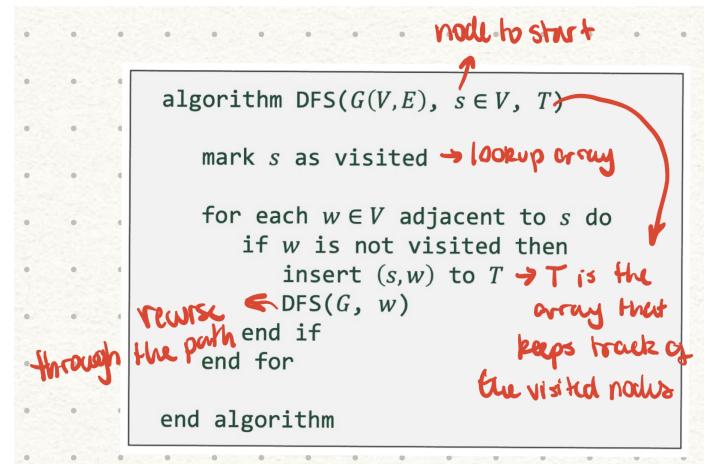
Test strong conn. by running DFS on  $G$  and  $G'$  (inverted pointing), if all nodes are visited then there is strong conn.

## Depth-first search (DFS)

Traverses  $|E|$  twice ( $2 \cdot |E|$ ) by backtracking.

If non-visited node remains, then  $G$  is not connected.

Runtime:  $\mathcal{O}(|V| + |E|)$  (adj. list)

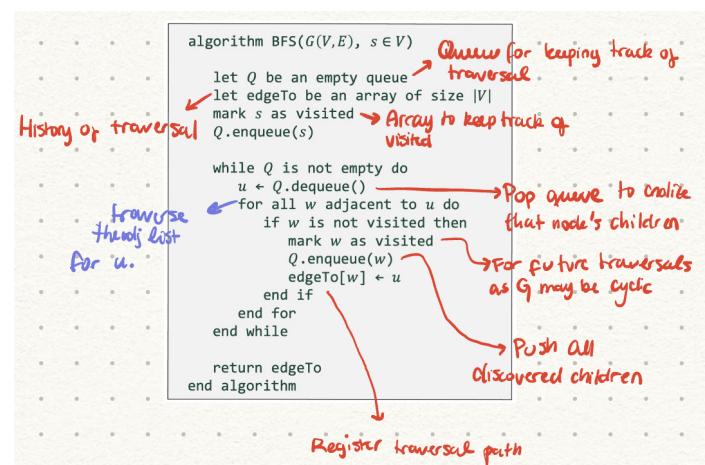


## Breadth-first search (BFS)

Traverses all  $V$  and  $E$ .

Runtime:  $\mathcal{O}(|V| + |E|)$

Terminates when the queue is empty or all vertices are discovered.



## Transitive closure

$G^*$  shows the paths where  $v \neq u$  has a directed path  $(u, v)$

Warshall uses 1s in diagonal if self loops or cycles.

Floyd-Warshall algorithm builds  $G^*$  matrix.

$R^k$  to denote progression from  $k = -1$ .

```

algorithm Floyd-Warshall( $M$ : adjacency matrix representing  $G(V,E)$ )
 $R^{(1)} \leftarrow M$ 
 $n \leftarrow |V|$ 
for  $k$  from 0 to  $n-1$  do
    for  $i$  from 0 to  $n-1$  do
        for  $j$  from 0 to  $n-1$  do
             $R^{(k)}[i,j] \leftarrow R^{(k-1)}[i,j]$  or  $(R^{(k-1)}[i,k] \text{ and } R^{(k-1)}[k,j])$ 
        end for
    end for
end for
return  $R^{(n-1)}$ 
end algorithm

```

$\boxed{O(n^3)}$

where  
m is the  
adj matrix  
and g is the  
graph

$m[i][j] = g[i][j] \text{ || } (g[i][k] \&& g[k][j]);$

final = init || (indirect path  $i \rightarrow k \rightarrow j$ );

Space:  $O(|V|^2)$  (matrix)

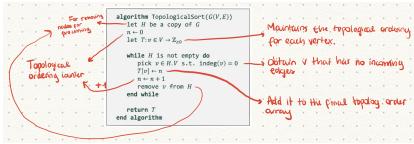
Runtime:  $O(|V|^3)$  (i, j, k)

## Topological ordering

Directed acyclic graph (DAG): digraph with no cycles.

Numbering  $v_1, \dots, v_n$  such that every  $(v_i, v_j)$  is  $i < j$ .

By inspection, remove nodes with `indeg(0)` in order.



Runtime:  $O(V + E)$  using DFS,  $O(1)$  for insertion lists.

## Dijkstra's algorithm

Greedy shortest path, that uses min-heap for inspected nodes and their weights.

→ Removes when discovered, updated when finding path with lesser cost.

Min-heap removal:  $O(|E| \log(|V|))$

Update:  $O(\log(|V|))$

Runtime:  $O((|V| + |E|) \log(|V|))$

```

algorithm DijkstrasShortestPath( $G(V,E)$ ,  $s \in V$ )
    let dist:  $V \rightarrow \mathbb{Z}$ 
    let prev:  $V \rightarrow V$  → To keep track of path traversed
    let  $Q$  be an empty priority queue → Min-heap
    dist[s] ← 0
    for each  $v \in V$  do
        if  $v \neq s$  then
            dist[v] ← ∞
        end if
        prev[v] ← -1
        Q.add(dist[v], v)
    end for

    while  $Q$  is not empty do
         $u \leftarrow Q.\text{getMin}()$ 
        for each  $w \in V$  adjacent to  $u$  still in  $Q$  do
             $d \leftarrow dist[u] + weight(u, w)$ 
            if  $d < dist[w]$  then
                dist[w] ← d
                prev[w] ← u
                Q.set(d, w)
            end if
        end for
    end while

    return dist, prev
end algorithm

```

for every  $v$  its shortest dist from src

src (init vertex)

To keep track of path traversed

Min-heap

initializing all min-heap to  $\infty$  except for src (0)

Remove min value from heap

until min-heap is empty

Change the value in min heap

Find → which set does each vertex belong to

Union → Merges sets if an edge connected is added

Extract the next edge with smallest weight

## Bellman-Ford algorithm

Shortest path with negative weights, that iterates  $V - 1$  times +1 for negative cycle checking.

Checks every node, and updates neighbors weights until no nodes are updated.

Runtime:  $O(|V| \cdot |E|)$

```

algorithm BellmanFord( $G(V,E)$ ,  $s \in V$ )
    let dist:  $V \rightarrow \mathbb{Z}$  → Array of all distances for  $V \in V$ 
    let prev:  $V \rightarrow V$ 

    for each  $v \in V$  do
        dist[v] ← ∞
        prev[v] ← -1
    end for
    dist[s] ← 0

    for  $i$  from 1 to  $|V| - 1$  do
        for each  $e = (u,v) \in E$  do
             $d \leftarrow dist[u] + weight(e)$ 
            if  $d < dist[v]$  then
                dist[v] ← d
                prev[v] ← u
            end if
        end for
    end for

    for each  $e = (u,v) \in E$  do
        if  $dist[u] + weight(e) < dist[v]$  then
            error "Negative Weight Cycle"
        end if
    end for

    return dist, prev
end algorithm

```

History of Parent that provided the weight

Initialization

$\boxed{V-1 \text{ iterations}}$

For all edges

Update if edge produces a lower dist, then current

One additional iteration for negative cycle checking

## MST & Kruskal's MST algorithm

Minimum spanning trees (MSTs) are spanning trees with minimal weights.

→ Splitting vertices into two groups, the min weight edge is part of the MST of  $G$ .

→ On cycles the biggest weight is never in the MST of  $G$ . Kruskal's algorithm finds the MST by using Union-Find on nodes at increasing weights.

Uses min-heap to obtain the weights in increasing order.

Runtime:  $O(|E| \log(|V|))$

```

algorithm KruskalMST( $G(V,E)$ )
    let  $Q$  be an empty min-heap
    let  $UF$  be a Union-Find with  $|V|$  components

    for each  $e \in E$  do
        Q.insert(weight(e), e)
    end for

     $T \leftarrow \emptyset$  → store edges that will form the MST

    while  $|T| < |V| - 1$  do
         $(u,v) \leftarrow Q.\text{getMin}()$  → Take out the min edge
        if  $u$  and  $v$  are not connected in  $UF$  then
             $T.\text{insert}((u,v))$ 
             $UF.\text{union}(u,v)$  → joins or merges nodes into a set
        end if
    end while

    return  $T$ 
end algorithm

```

Find → which set does each vertex belong to

Union → Merges sets if an edge connected is added

Extract the next edge with smallest weight

Initialize min heap

For all vertices

Take out the min edge

Joins or merges nodes into a set

## Union-find

Quick-Find merges by grouping nodes under a "parent" node label, and finds by an array lookup.

```

UF(n)
count <= n
for i from 0 to n-1 do
    id[i] <- i
end for

function union(p:item, q:item) joining sets
    exit if id[p] = id[q]
    idP + id[p] > idQ
    idQ + id[q]
    for i from 0 to n-1 do
        if id[i] = idP then
            id[i] <- idQ
        end if
    end for
    count + count - 1
end function

function find(p:item)
    return id[p]
end function

function connected(p:item, q:item)
    return id[p] = id[q]
end function

function count()
    return count
End function

```

Number of sets

Initialize the identifier array (id) with each elem being its own set

Hold the sets before union

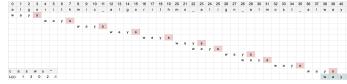
check if they are in the same set

Element at pos. i is in the same set as p thus needs to be converted to new set q.

Returns the set of element at index p.

If the two items share the set.

Number of sets

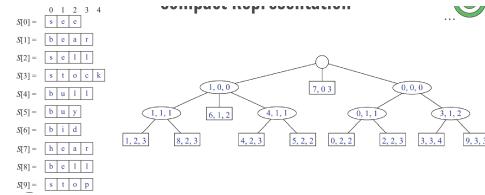


## Tries

Prefix/suffix tries contain all possible combinations.

For a single string:  $\mathcal{O}(|w|)$

For PATRICIA tries, a node is redundant and can be grouped if it's not the root and has one child (include \$).



## Huffman

Greedy, min-heap, that encodes high freq. with short code-words.

Runtime:  $\mathcal{O}(n \log(n))$ .

```

algorithm HuffmanCoding(S:string)
    C <- distinctCharacters(S) → Σ
    F <- computeFrequencies(S, C) → Greedy heuristic
    let Q be an empty min-heap → Quick access to smallest
        for each c ∈ C do
            let T be a new tree node
            T.char <- c
            T.freq <- F[c]
            Q.insert(F[c], T)
        end for
        while Q.size() > 1 do
            let T be a new tree node → Node that holds the freq. of its children
            T.left <- Q.getMin() → Extract least freq. two nodes from min heap
            T.right <- Q.getMin() → Insert it back into min-heap for iterative building
            T.freq <- T.left.freq + T.right.freq
            Q.insert(T.freq, T)
        end while
        return Q.getMin() → Last node in min-heap is the root of the encoding tree
    end algorithm

```

Adding all chars and their freq into min heap

Last node in min-heap is the root of the encoding tree

## K-d trees & quadtrees

K-dimensions where non-leaf nodes split into half-spaces. Each level alternates dimensions (e.g. 0:x, 1:y, 2:x).

Quadtrees split into 4 spaces (NW, NE, SW, SE).

Insertions use level discriminant to locate parent child nil (e.g. Insert(70, 50) uses 70 for x levels, and 50 for y levels).

Deletions swap node with minimum value of discriminant at the level in the right subtree (if the min node has children then the same process is done for that node). Leaf is deleted.

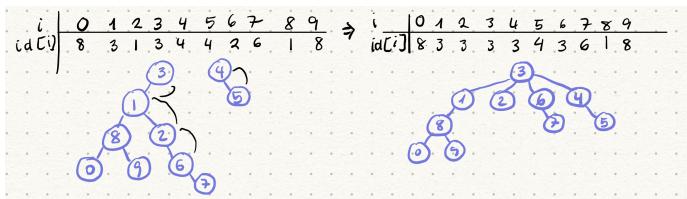
The expected depth of the quadtree at each level:  $4^D \approx n \therefore D \approx \frac{1}{2} \log_2(n)$ , where  $n$ : num. points.

Range query result is determined by square query  $s$  intersection with square region:  $\frac{\text{UpperBound}}{2^D} = \frac{\text{UpperBound}}{n^{1/2}}$ .

$\frac{s}{1/n^{1/2}} = s \cdot n^{1/2} \therefore s^2 \cdot n$  are the points intersecting (leaf nodes).

Quick-Union, using path compression, joins nodes into their corresponding tree.

Union(6, 5) will look up to the root of 6 in Find(6): 6 → 2 → 1 → 3 and will attach all intermediaries directly to the root (3). Does the same with Find(5) 5 → 4 and joins the smaller tree as a child of the bigger tree (4 child of 3).



```

UF(n)
count <= n
for i from 0 to n-1 do
    id[i] <- i
end for

function union(p:item, q:item)
    idP + find(p)
    idQ + find(q)
    exit if idP = idQ
    idP <- idQ
    count + count - 1
end function

function find(p:item)
    while p < id[p]
        p <- id[p]
    end while
    return p
end function

function connected(p:item, q:item)
    return find(p) = find(q)
end function

function count()
    return count
end function

```

```

function find(p:item)
    if p = id[p] then
        return p
    end if
    id[p] <- find(id[p])
    return id[p]
end function

```

Quick-Find runtime:  $\mathcal{O}(\text{tree\_depth})$  or worst degen  $\mathcal{O}(n)$

Quick-Union runtime:  $\mathcal{O}(1)$  and path comp:  $\mathcal{O}(\alpha(n))$

## Brute-force string matching

Takes a substring and compares it contiguously.

Runtime:  $\mathcal{O}(m) \cdot \mathcal{O}(n) \in \mathcal{O}(nm)$

## Boyer-Moore algorithm

$T$  (text) and backwards  $P$  (pattern) comparison.

Creates a numeric rep. for the occurrence of each char, and then compares and moves at offset:  $m - \min(j, 1 + L[c])$ .

Number of comparisons are the total times the last char is compared + offsets (example has 15).