

WEEK 1

Why graphs? Why GRAPHS?

- ① Check if network is connected; example Kevin Bacon (~~number~~) movie network
- ② driving directions in Google Maps
- ③ formulate a plan where only certain changes can be applied to a system
- ④ Comprise the "pieces" (or components) of a graph.
 - clustering or structure of a certain graph.

GENERIC GRAPH SEARCH

- Goals:
- ① Find everything findable from a given start vertex: example
 - ② Don't explore anything twice Goal: $O(m+n)$

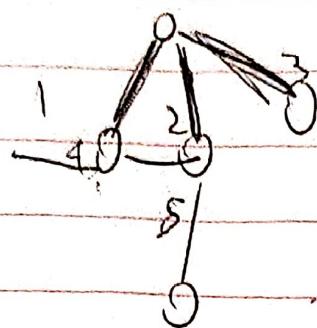
Generic algorithm (given graph G , vertex s)

- initially only s has been explored
- while possible:
 - choose edge (u,v) with u explored and v unexplored, (if none, halt)
 - mark v explored (append v)

Claim: At the end of the algorithm, v explored $\Rightarrow G$ has a path from s to v .

Note: how to choose between frontier edges?

BFS



etc, explore all nodes ~~at the same~~ of distance before exploring their children

DFS explores as deep as it can. Backtracks when necessary

BFS:
incomplete computing shortest paths
compute connected components of an undirected graph

DFS:

- Compute topological ordering of directed acyclic graph
- Compute "connected components" in directed graphs

Breadth First Search

THE CODE

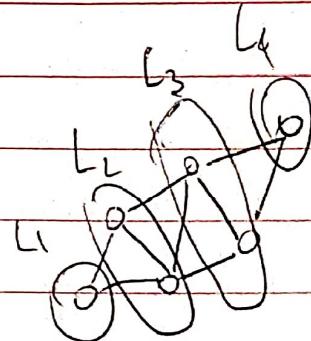
BFS (graph G, start vertex s)

explored = [] # all nodes initially unexplored

explored.append(s) # mark s as explored

Let Q = queue data structure (FIFO), initialized with s.

- while $Q \neq \emptyset$:
 - remove the first node of Q, call it v
 - for each edge (v, w)
 - if w unexplored
 - mark w as explored
 - add w to Q



APPLICATION OF BFS : SHORTEST PATHS

Goal: compute $d_{\text{list}}(v)$, the fewest # of edges on a path from s to v

to add to BFS list

Extra code) initialize $d_{\text{list}}(v) = \begin{cases} 0 & \text{if } v=s \\ \infty & \text{if } v \neq s \end{cases}$

- When considering edge (v, w) :

- if w unexplored, then set $d_{\text{list}}(w) = d_{\text{list}}(v) + 1$

WEEK 1 L4

APPLICATION OF BFS : UNDIRECTED CONNECTIVITY

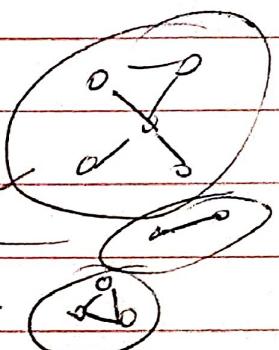
Directed / undirected connectedness are very different.

Let $G(V, E)$ be an undirected graph.

Connected components = the "pieces" of G

Formal definition

equivalence classes of the relation $v \sim v \Leftrightarrow \begin{cases} v=v \\ \text{path in } G \end{cases}$



Goal: compute all connected components

Why! :- Check if network is disconnected.

- Display graph?

- clustering

A

CONNECTED COMPONENTS V/P BFS

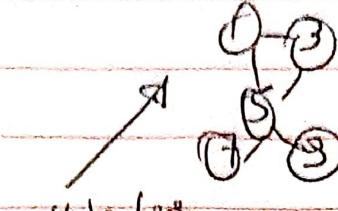
To compute all components (undirected case):

- all nodes are explored

- for $i = 1$ to n :

- if i not yet explored:

- $\text{BFS}(G, i)$ is a component



starts here

then here

then here

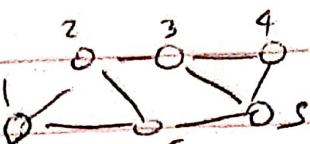
Running time: $O(n^2)$

$O(1)$ per edge
in each BFS

$O(1)$ per node

~~APPLICATION OF DFS, UNDIRECTED~~

DFS



Backtrack when necessary (we hit somewhere we've already been)

- computes & topological ordering of a directed acyclic graphs
- and strongly connected components of directed graphs.

Runtime is $O(n+m)$

THE CODE

DFS, (LIFO) instead of (FIFO)

DFS graph G, start vertex S:

LIFO = { }

- work set explored
- for every edge from $v \xrightarrow{\text{LIFO}} u$:
- if v unexplored, add to LIFO

Running Time $\Rightarrow O(n_s + m_s)$

$\#$ of edges reachable from s

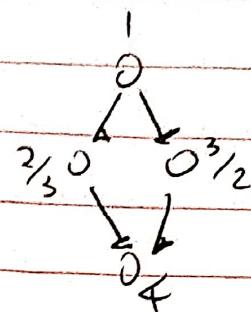
$\#$ of nodes reachable from s

6

FINALLY, WHAT THE HELL IS A TOPOLOGICAL ORDERING (of a directed acyclic graph)

Definition: A topological ordering of a directed graph G is a labeling F of G 's nodes s.t.:

- ① The $F(v)$'s are the set $\{1, 2, \dots, m\}$
- ② $\forall (v, u) \in G \Rightarrow F(v) < F(u)$



Motivation: Sequence tasks while respecting all precedence constraints

Q: You need the graph to be ~~cyclic~~ ...
(COND)

Theorem: You can compute Tord in linear time.
(if cond) Or you're screwed too!

STRAIGHTFORWARD SOLUTION

Note: every directed acyclic graph has at least one sink vertex
(but, maybe 0 → then solution)



To compute topological ordering:

The only candidates for 1st position are sink vertices

- set $f(v) = n$

- reverse or $G = \{v\}$

TOPOLOGICAL SORT via DFS (Slick).

M1

DFS - Loop (graph G)

marks all nodes very badly

and \oplus label = n (to keep track of ordering)

for each vertex $v \in G$

- if v not yet explored
in some previous DFS

call

- $\text{DFS}(G, v)$

M2

DFS (graph G, start vertex v)

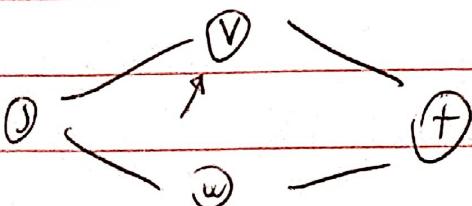
- mark, explored

- for every edge (u, v) :

- if v not yet explored
- $\text{DFS}(G, v)$

Crucial bit: remember what you've explored already

M1)



$V \rightarrow t$, sink so append t
 $\Rightarrow s(4)$

Back to v , sink so 3

Back to outer for loops

Skip v because already explored!

move to w

M2)

" "

" "

" "

Check v

:

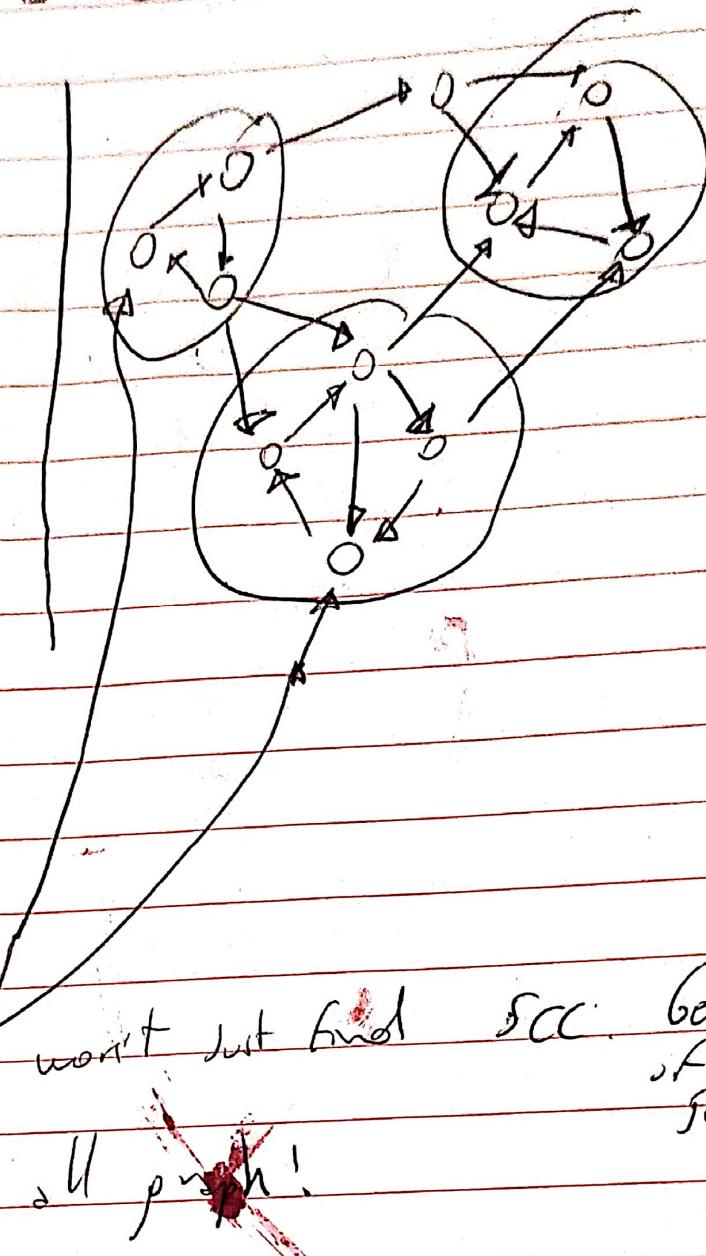
$O(m \cdot m)$, correct.

$y = y$
 $b = 1$
 $y = \text{whatever}, b^d = 2, \text{ case 1},$

STRONGLY CONNECTED COMPONENTS

in directed graphs

SCCs are cycles and nodes out
of cycles



Why DFS?

If we start from ~~any node~~, we won't just find SCC. Get a union of possible SCCs

~~Start from all nodes!~~

Need to choose where ~~we~~ we want to start.

KOSARAJU'S TWO PASS ALGORITHM

SCCs in $O(n+m)$ time

Algo:

- ① Reverse all arcs, let G^{rev} = G with all arcs reversed
- ② Run DFS on $G^{\text{rev}} \rightarrow \text{post}$: compute "lexical ordering" of nodes
- ③ Run DFS on $G \rightarrow \text{post}$: discover the SCCs one by one

DFS - Loop

DFS Loop (depth)

global variable $t = 0$
 (at start, procedure is)
~~not yet reached~~

global variable $S = \text{NULL}$
 (current source vertex) $\xrightarrow{\text{for ready}}$
 $\xleftarrow{\text{in loop}}$

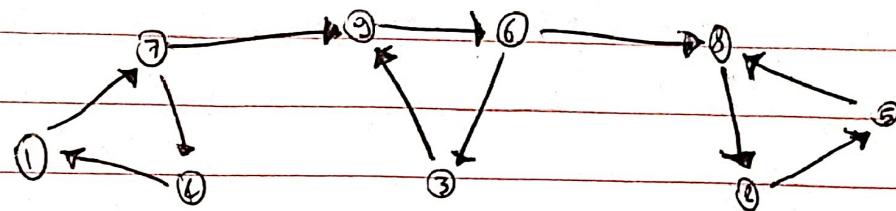
some nodes are labeled with
 finishing times t to n

for $i = n$ down to 1:

if i not yet explored.

DFS(G, i)

G^{rev}



Starting from 8; 9, 6, 3, finish! $F(3) = 1$

Back to 6: 6, 8, 2, 5, finish! $F(5) = 2$

Back to 2: 2, finish! $F(2) = 3$

Back to 8: 8, finish! $F(8) = 4$

Back to 6: 6, finish! $F(6) = 5$

Back to 9: 9, finish! $F(9) = 6$

All of these have leader 9!

Next to be explored node: 7

DFS($G, p, h, 6, \text{node } i$):

- mark: as explored

- set $\text{leader}(i) = \text{node } S$

~~do-ex~~

- For each arc $(i, j) \in G$:

- if j not yet explored:

- DFS(G, j)

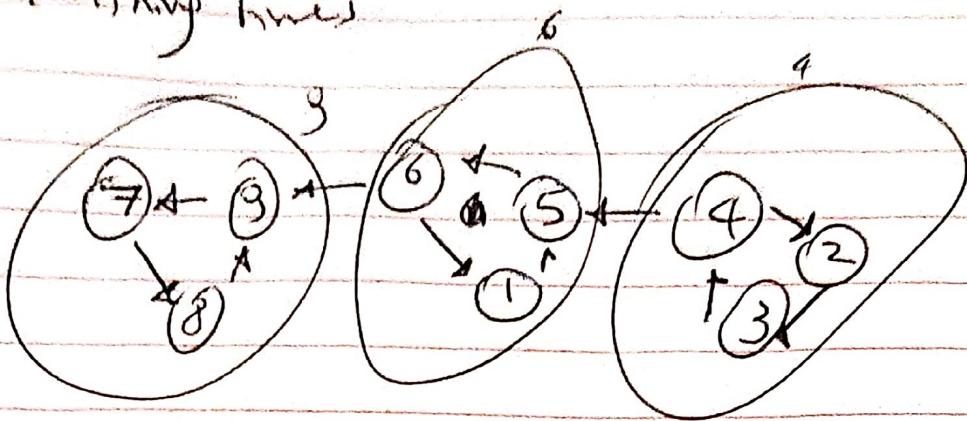
$t++$

- set $f(i, j) = t$

$\begin{cases} \text{---} \\ \text{---} \end{cases}$
 finishing time

$$T(n) = [7, 3, 1, 8, 2, 5, 9, 4, 6]$$

All the edges are reversed, & nodes names change to finishing times



Run DFS again ^{for} from both borders. top → bottom

from 3: 3, 7, 8, done!, (order 3)

from 8, 7, 6: 6, 1, 5, done!, (order 6)

from 4, 1: 4, 2, 3, done!, (order 4)

Running time: $2 \times \text{DFS} = O(m+n)$

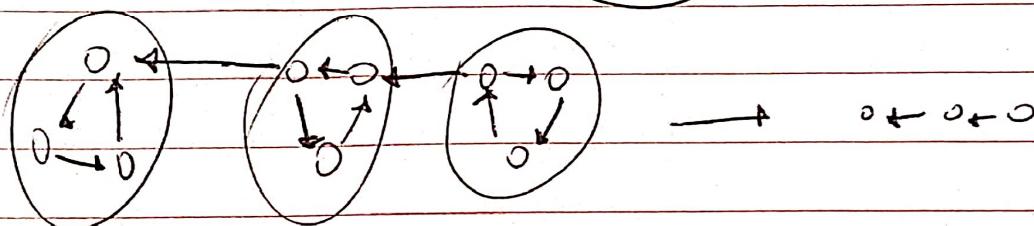
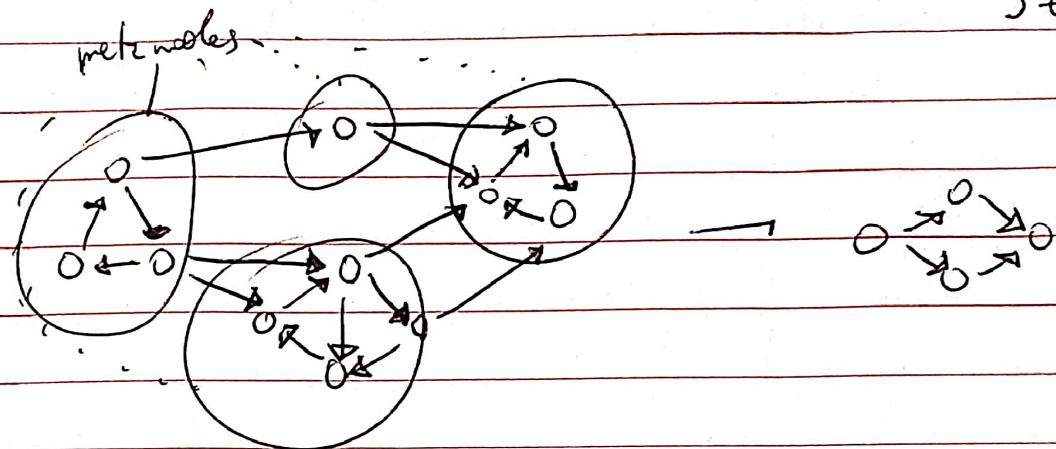
GENERAL ARGUMENT FOR LINEAR TIME

OBSERVATION

Claim: The SCCs of a directed graph induce an acyclic "meta-graph"

- meta ^{nodes}
~~graph~~: the SCCs
- $\exists \text{ arc } i \rightarrow j \Leftrightarrow \exists \text{ arc } (i, j) \in G$, with $i \in G$
 $j \in G$

Example:



The meta-graph is always acyclic, any cycle of SCCs is just one SCC.

THE SCCs ARE THE same IN G and G^{rev}

(and in the second pass of DFS in Karp's algo too!)

Key lemma:

We'll prove this later

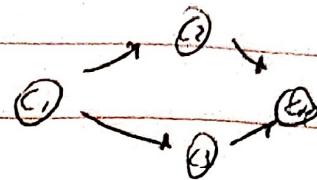
Consider the "bottom" SCCs S_1, S_2, \dots, S_k

(out)

end part of DFS (sp begins somewhere in a sink) see c

First cell discovery C^* and nothing more.

After that goes back up "to C_2, C_3 ,
etc etc.



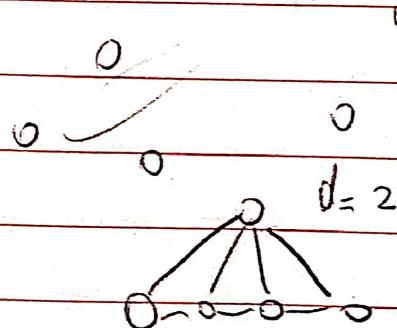
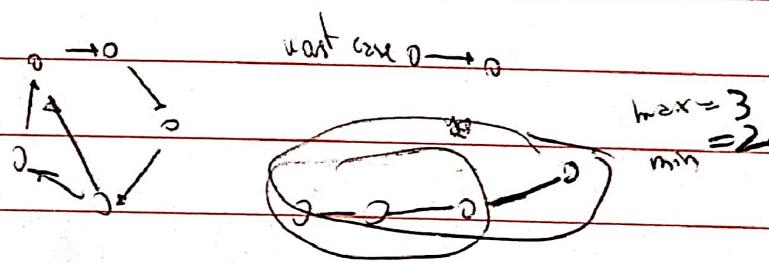
Got it? Now let's prove the lemma.

PROOF OF KEY LEMMA

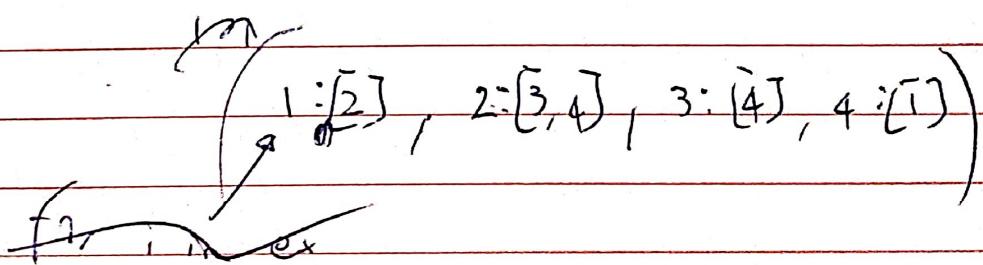
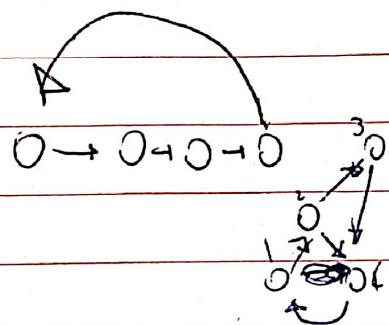
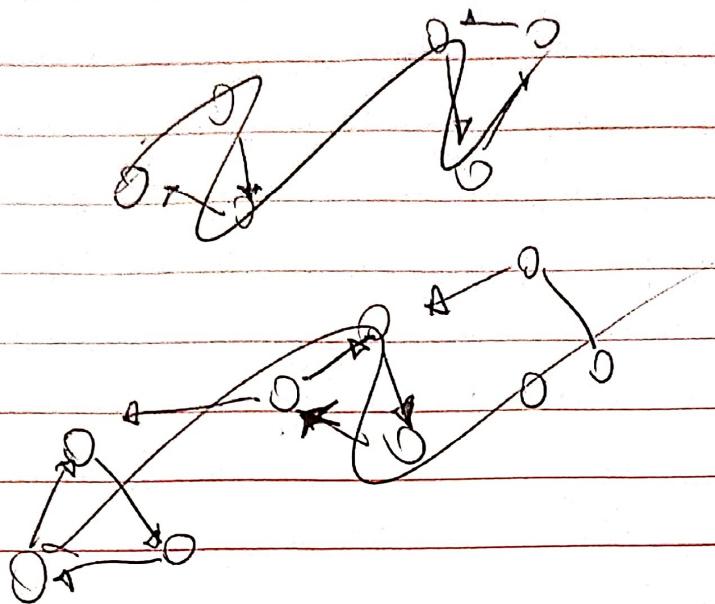
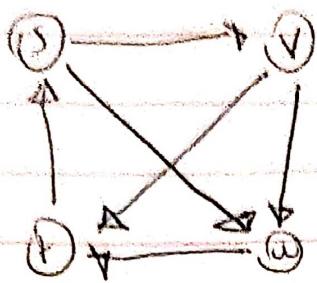
Case $[v \in C_1]$: we'll see all of C_1 before C_2

Case $[v \notin C_1]$: DFS (G_{new}, v) won't finish until $C_1 \cup C_2$
completely explored

PROB PERIBBLINGS



n times rows of length $n \times n^2$



2

$f_2(2:)$

WEEK 1

→ Dijkstra's shortest path algorithm.

Input: directed graph $G_p = (V, E)$

- Each edge has non negative length
- Source vertex s

Output: For each $v \in V$, compute

$$L(v) = \text{length of shortest } s-v \text{ path}$$

Assumptions:

- ① [for conciseness], consider strongly connected; $\exists s-v$ path for all edge length (for concnca)
- ② $l_e \geq 0, \forall e \in E$

WHY DO WE NEED ANOTHER SHORTEST PATH ALGO?

ALGO? CAN'T WE ALREADY USE BFS IN $O(n^2)$?

Yes. IF $l_e = 1$ for every edge.

Question: why not replace each edge by multiple edges of length 1.

Stupid question

DJIKSTRA'S ALGORITHM

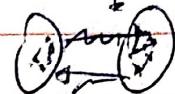
Initialize:

- $X = \{s\}$ \Rightarrow vertices processed so far

- $A[s] = 0$ \Rightarrow computed shortest path distance

- $B[s] = \text{empty path}$ (computed so far)

This is only for show not included in real implementation
we can omit this



- while $X \neq V$:

- sweep all edges $(v, w) \in E$ with $v \in X$, $w \notin X$,
pick the edge (v, w) that minimizes the following criterion:
 $A[v] + l_{vw}$ (Dijkstra's greedy algorithm criterion)

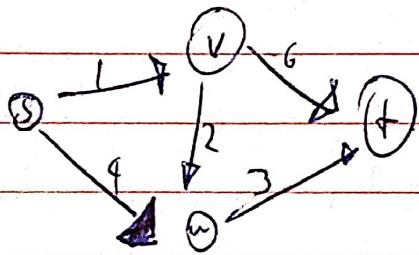
- add w^* to X

- set $A[w^*] = A[v^*] + l_{v^*w^*}$

$\Rightarrow B[w^*] = B[v^*] \cup [v^*, w^*]$

WEITERE

EXAMPLE



First iteration: all edges outside X

Dijkstra's criterion: (s, v) has score of $0 + 1$, $(s, u) = 0 + 4$

So we go to s, v

Second iteration

All edges leading to outside of X : (S, u) , (v, w) , (v, t)
Distances increase for each: $(0+4)$, $(1+2)$, $(1+6)$

We then add w to X , etc etc ($A(w) = 1+2 = 3$)

Third iteration

?

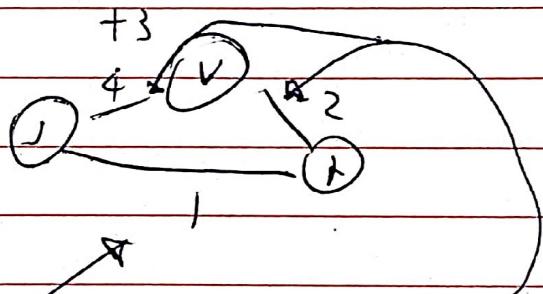
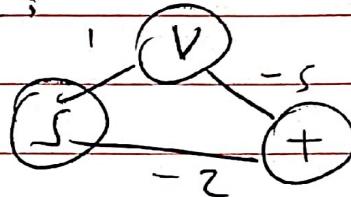
All edges: ~~t, z~~ (v, t) , (w, t) greedy $(1+6)$ $(3+3)$

Bring in t ,

$$A(t) = 3+3 = 6$$

~~Breaks~~ The while loop now breaks due to $X = 1$

(Can't reduce negative length to non negative by adding
a large number to all lengths because:



reduces this, but it's actually not

WEEK 2 C3

WHY IS DIJKSTRA'S ALGO CORRECT

Proof by induction:

Base case: $A(S) = L(S) = \emptyset$ correct ✓

Then on and on. Seems pretty simple.

Inductive hypothesis: all previous iterations correct.

→ In current iteration, we set $B(u^*) = B(v) \cup \{v\}$
↓ bounded length

DIJKSTRA AND RUNNING TIME

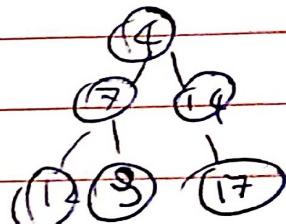
Naively appears to be $O(n^2)$ (roughly)

HEAP OPERATIONS

Reason of efficiency of heap = perform insert, Extract-min in $O(\log n)$ time

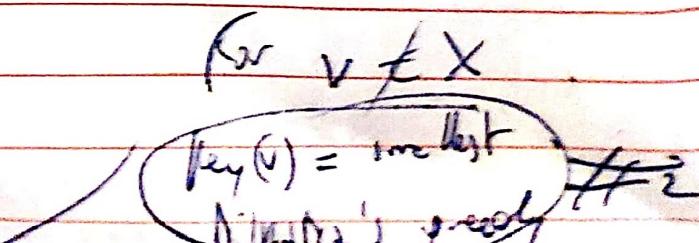
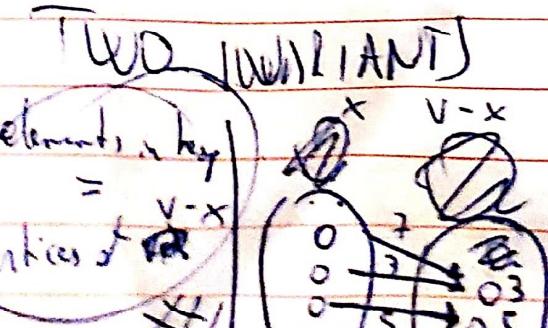
[rest of video assumes familiarity with heaps]

- Conceptually, perfectly balanced binary tree
- heap property: at every node, key \leq children's keys
- Extract-min by plucking off the root, +~~up~~ bubble down
- Insert via bubbling down



Also: will need ability to delete from middle of heap.

Because you want to keep it as a pbbt, height $\approx \log_2 n$
All operations run in $O(\log n)$ time.



Knockout tournament.

Two rounds: At first each voter runs a local round to find smallest edge from X pointing to it, put them in the heap.

Then min (heap) finds best out of these local wins. We'll be tracking the min on a silver platter, root of the heap.

To maintain #2

that $\forall v \notin X$, $\text{key}[v] = \min_{u \in X} \text{greedy score of all edges } (u, v)$,

- When w is extracted from heap (related to x) .

- for each edge $(w, v) \in E$:

- Insert if $v \notin V - x$

- delete from heap

- recompute $\text{key}[v] = \min[\text{key}[v], \text{key}[w] + d_{vw}]$

- reinsert v into heap

(maybe it slides down)

Running Time is $O(\log n)$ for each heap op question

- $(n-1)$ extract, min

- each edge (v, w) results in at most one Delete/Insert combination of vertices to X first.

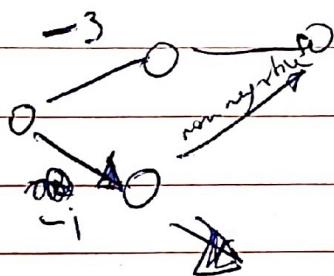
So # of heap operations is $O(n+m) = O(n)$

we have much if work \Rightarrow

Problem Sunday



- ✓
- ✓
- ✗
- ✗



Can't go back, so yeah. Can add positive constant to all other edges

Or

Insert is mostly open, but can't be n
Time $O(n)$

Kill always
remove

Opposite for nearly open

with negative shift mixed result

Largest elem $O(n)$
Median same
1st, nth best
✓ $\log(n)$

Open tree.

WEEK 3 L1

DATA STRUCTURES

Point: organize data so it can be accessed quickly and safely

Examples: lists, stacks, queues, heaps, search trees, hash tables, bloom filters, union find etc.

Why so many?

They will support different operations with different times.

So they are suitable for different kinds of tasks.

Rule of thumb: choose the "minimal" data structure that supports all the operations that you need.

TAKING IT TO THE NEXT LEVEL

- (LEVEL 0) - What's global structure?
- (LEVEL 1) - "Gotta know knowledge, kind of stuff, can't use in tech whenever"
- (LEVEL 2) - "This problem calls for a heap"
- (LEVEL 3) - "I only use data structures I wrote myself"

We'll get to here

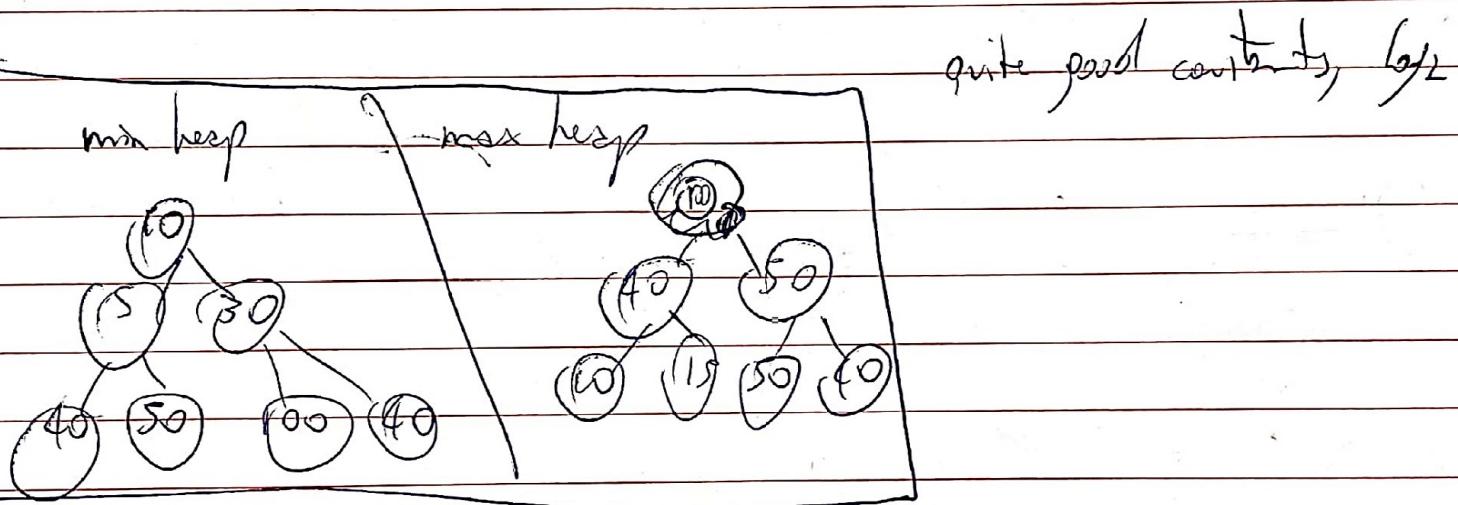
WIBA 3 L2

Heap HEAP : SUPPORTED OPERATIONS

- > container for objects (not here key)
- employe networks, network edges, events etc

INSERT: add object to a heap Runtime $O(\log n)$

EXTRACT + MIN (or MAX): remove an object with a minimum key value (tree broken arbitrarily) Runtime $O(\log n)$



IN-HEAPIFY : (n batched inserts)
 $\in O(n)$ time).

You don't need to create an empty heap and put n objects in. n times $\log n$ time needed for insertion $O(n \log n)$,
We can do it in $O(n)$ only

DELETE an arbitrary element from middle of a heap in $O(\log n)$

~~Smartest uses of heaps~~

CANONICAL USES OF HEAP

Canonical use: fast way to do repeated minimum computations.

Example: Selection Sort $\approx \Theta(n)$ linear scans, $\Theta(n)$ time on array of length n

Heap sort:

- ① Insert all n array elements into a heap
- ② Extract-Min to pick out elements in sorted order

Running time = ~~on~~ $2n$ heap operations = $\Theta(n \log n)$ time

Can we do better? No. $n \log n$ is lower bound for any sort

ANOTHER APPLICATION: EVENT MANAGER.

"Priority queue" - synonymous for a heap

Imagine video game simulation:

- Objects = event records, event-triggering events added with time-stamps key
- Extract-min \rightarrow next scheduled event.

APPLICATION: MEDIAN MAINTENANCE

I give you an increasing amount of numbers, one by one.

I want the median each time I step;

Need to do it in $\log n$ (?) time

WEEK 3 L4

BALANCED SEARCH TREES

~~Sorted Array: supported operations~~

3	6	10	11	17	23	30	36
---	---	----	----	----	----	----	----

OPERATIONS	RUNNING TIME
SEARCH (bin search)	$O(\log n)$
SELECT i^{th} statistic	$O(1)$
MIN / MAX	$O(1)$
PREDCESSOR / SUCCESSOR	$O(1)$
RANK (rank of 23 is 6)	$O(\log n)$
OUTPUT IN SORTED ORDER	$O(n)$ or $O(1)$ n times, $O(n)$
INSERTION	$O(n)$
DELETION	$O(n)$

That's a lot for insertion and deletion unless you barely ever do them

That's why we make:

BALANCED SEARCH TREES

BST's

Reason of tree: sorted array + fast (logarithmic) inserts and deletes

OPERATIONS	RUNNING TIME
SEARCH	$O(\log n)$
SELECT	$O(\log n)$
MIN / MAX	$O(\log n)$
PRED / SUCC	$O(\log n)$
RANK	$O(\log n)$
OUTPUT IN SORTED	$O(n)$

Up from $O(1)$

SAME FOR HEAP
or
MIN/MAX $O(n)$ searches

WEEK 3 LS

Reason of tree: like sorted array \rightarrow fast logarithmic

BSTs

- exactly one node per key

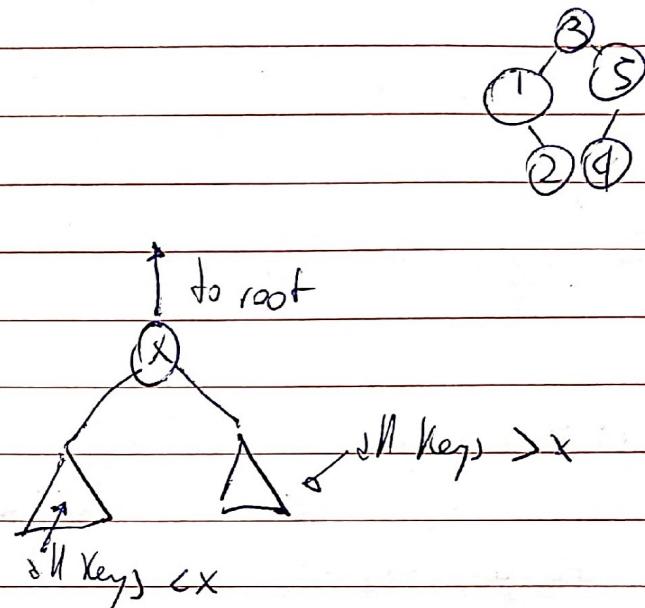
- most basic version:

each node has:

- left child pointer

- right child pointer

- parent pointer

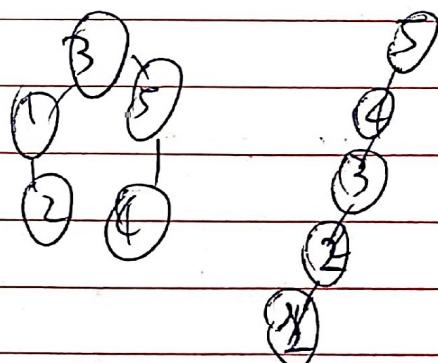


HEIGHT OF BST

best case, perfectly balanced worst case: chain

Roughly $\Theta(\log_2(n))$ usually, up to $\Theta(n)$

Many possible trees (for a set of keys):



SEARCH FOR KEY IN TREE T

- start at the root

- traverse left/right child pointers as needed

- return node with key k or null as appropriate

TO INSERT

- search for k (unsuccessfully)

- put new node with key k at the end of pointer

WEEK 3

Search or Insert operations in a binary tree take at most $\Theta(\text{height})$

MIN, MAX, PRED AND SUC

To compute min: start at root, if you follow left pointer until the end
 $\Theta(\text{height})$ again

PREDCESSOR OF KEY K

- easy case: if ~~is~~ left subtree to K, return K's max in left subtree
- otherwise: if left subtree empty, father pointer points until you get to a key less than K

PRINT KEYS IN INCREASING ORDER

- Let r = root of search tree with subtrees T_L and T_R
- recurse on T_L
- print r
- recurse on T_R

$O(n)$ time, $\frac{n}{\log n}$ recursive calls. $O(n)$

DELETE A KEY FROM A SEARCH TREE

- SEARCH FOR KEY

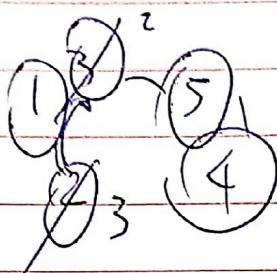
Then: three cases

(easy) no children. Just delete

(medium) splice out 's K' node. Its one child takes its place

BIT1 CULT (Note (it's node has 2 children))

- Compute K's predecessor & until you get to top, then jump up



Then delete 3.

The running time $\Theta(\text{height})$

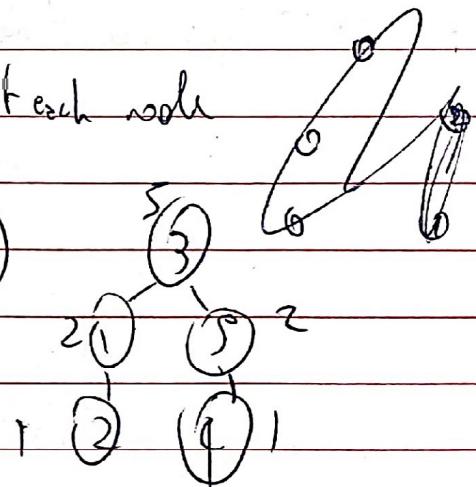
SELECT AND RANK

Ideas: store a little bit of extra info at each tree node.
about the tree itself

Example: store # of nodes in subtree rooted at each node

Note: if x has children y and z , $(\text{size}(y) + \text{size}(z) + 1)$

Also: Easy to keep sizes up to date in $\Theta(\text{height})$



TO SELECT:

- Start at root x , with children y and z
- let $s = \text{size}(y)$ [$s=0$ if x has no left child]
- if $s \geq i-1$ return x 's key
- if $s \geq i$ recursively compute on root y
- if $s \leq i-1$ (i.e. $i-1$)th on z

WEEK 3 (7)

Isasi Enve height slumpy layer

Example: Red-Black Trees

RED BLACK INVARIANTS

- ① Each node red or black
- ② Root is black
- ③ No 2 reds in a row
- ④ Every root-null path has the same number of black nodes

~~EXCEPT FOR LEAVES~~

NO CHAIN CAN BE AN RB~~BT~~ TREE

HEIGHT GUARANTEE:

$$\text{height} \leq 2 \log_2(n+1)$$

WEEK 4 L1

HASH TABLES

Purpose: immediate random access, through keys

Insert: add new record

Delete: delete existing record

Lookups: check for a particular record

Using key

ONLY THESE OPERATIONS, BUT GUARANTEED O(1) TIME (if implemented well)

- (1) Properly implemented
- (2) Non pathological

APPLICATION: DE-DUPLICATION

Input: "stream" of objects

Goal: ignore duplicates

hashtable Hash lookups.

If not found, insert x into H

APPLICATION: 2-SVM PROBLEM

Input: Unsorted array of n integers. Target sum t

Goal: determine whether two numbers in Array $x + y = t$

Name solution: $O(n^2)$ ($n(n-1)$)

SUM PROG CND

- Even better:
- ① insert elements of A in H . $O(n)$
 - ② for each x , look in H for $t-x$ $O(n)$

FURTHER OBVIOUS APPLICATIONS

- historical application: symbol tables in compilers
- blocking network traffic
- search algorithms (game tree exploration), can use to avoid exploring more than once

WEEK 4 L2

HASH TABLE IMPLEMENTATION

High Level rules:

Setup: Universe U (generally, really big)

Goal: Maintain evolving set $S \subseteq U$ (of reasonable size)

Naive Solution ①: Array based solution indexed by n ($O(1)$ ops, $O(n)$ space)

Naive Solution ②: List based solution ($O(1)$ space, $O(s)$ lookups)

Solution:

- ① Pick $n = \# \text{ of buckets}$ to store data in with $n \geq |S|$
splits out position in say
- ② Choose a hash function $h: U \rightarrow \{0, 1, 2, \dots, n-1\}$
for simplicity we'll assume S doesn't vary too much
- ③ Use array A of length n , store x in $A[h(x)]$

What happens when there are collisions? $h(\text{"Bob"}) = h(\text{"Alice"})$

Think of birthdays! 23 people!

RESOLVING COLLISIONS:

Collision: $h(x) = h(y)$ with $x, y \in U$.

Solution #1: Separate chaining. Keep linked list in bucket, perform operations on lists.

Solution #2: Open Addressing (one object per bucket), h now specifies probe sequence

~~Double hashing~~: use buckets as hash lists themselves, takes lots of time!

(Linear probing, quadratic probing)

WEEK 4 C3

HASH TABLES IMPLEMENTATION DETAILS

WHAT IS A GOOD HASH FUNCTION?

Note: in hash table with chaining, insert $O(1)$ most of the time

Performance depends on the choice of hash function!

Properties of good hash function:

- ① Should spread data out well (gold standard: completely random)
- ② Should be simple enough, not take too much time (must be $O(1)$)

EXAMPLES OF BAD HASH FUNCTIONS

Ex: Keys = phone numbers $h(x) = 10^{10}$

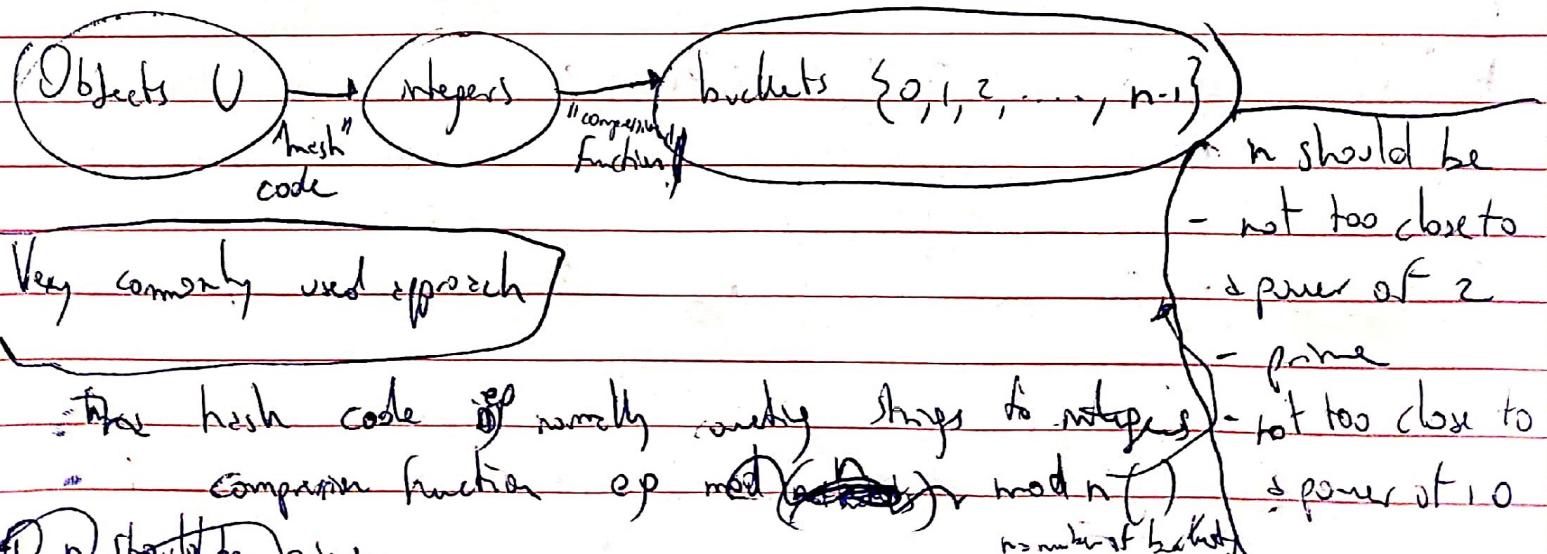
terrible: $h(x) = \text{1st 3 digits of } x \quad (n = 10^3)$

mediocre: $h(x) = \text{last 3 digits of } x \quad (n = 10^3)$

Example: keys = memory locations (power of 2)

bad hash func: $h(x) = \text{mod 1000 } (x)$ ('all slot buckets are empty')

QUICK AND DIRTY HASH FUNCTIONS



WEEK 4 L4

THE LOAD OF A HASH TABLE

Definition: The load factor of a hash table is:

$$\alpha = \frac{\text{# of objects in buckets}}{\text{# of buckets}}$$

Conditions for good performance:

- ① $\alpha = O(1)$, else won't run in constant time. (NEED TO BE ACC
TO COMPUTE α in O(1))
- ② for open addressing, $\alpha < 1$ middle

For good performance, we need to control load

Also, THE PERFECT HASH FUNCTION DOES NOT EXIST. FOR EVERY $H(x)$ THERE IS A PATHOLOGICAL DATA SET

PATHOLOGICAL DATA IN THE REAL WORLD

~~Remember~~ Mem Point : Can analyze several real-world systems by exploiting badly designed hash functions.

If the hash function ~~is~~ ^{is} open source and ~~too~~ simple, ~~it's~~ easy anyone can read it and reverse engineer it.

SOLUTIONS

- ① Use a cryptographic hash function which can't be reverse engineered
- ② Use randomization: design a family H of hash functions and pick one at random

BLOOM FILTERS

A variant on hash tables, save space, have errors.

PROS

- Space

CONS

- Can't store associated object
- No deletions
- Small false positive probability

Digital APPLICATIONS

Digital: early spell checkers

(anonymize): List of forbidden passwords

Modern: Network routers

- limited memory, need to be super fast still

WHAT IS A BLOOM FILTER?

Ingredients:

- An array of n bits. $\frac{n}{|S|} = \# \text{ of bits per object in the set}$
- K hash functions h_1, h_2, \dots, h_K . $K = \text{small constant}$

for $i = 1, 2, \dots, K$, $A[h_i(x)] = 1$, whether or not it was 1 already

Lookup (x) return TRUE if $A[h_i(x)] = 1$ for every $i = 1, 2, \dots, K$

There are false positives, but no false negatives and it saves space

WEEK 4 L 9

HEURISTIC ANALYSIS OF BLOOM FILTERS

There is a trade off between % of false positives and space saved.

After inserting object s , $1 - \left(1 - \frac{1}{b}\right)^{K(S)}$ probability for a given bit to have been set to 1

$$\text{prob} \in 1 - e^{-\frac{K}{b}} \quad b = \frac{n}{|S|}, \text{ number of bits per object}$$

This means the false positive prob. is bounded by $(1 - e^{-\frac{K}{b}})^K$

how to find K ? for fixed b , the false positive probability is minimised by setting $K \approx (\ln 2) \cdot b$

Example : $b=8$, $K=5$ or 6 , error prob. $\approx 2\%$

Problem Solving

$$h(n \text{ keys}) \rightarrow \sqrt{\log_2 m}$$



like with the n tasks for m series

$\Omega(n)$

$\Omega(1)$



$$\frac{n(n-1)}{2} \cdot \frac{1}{m} = \frac{n(n-1)}{2m}$$

