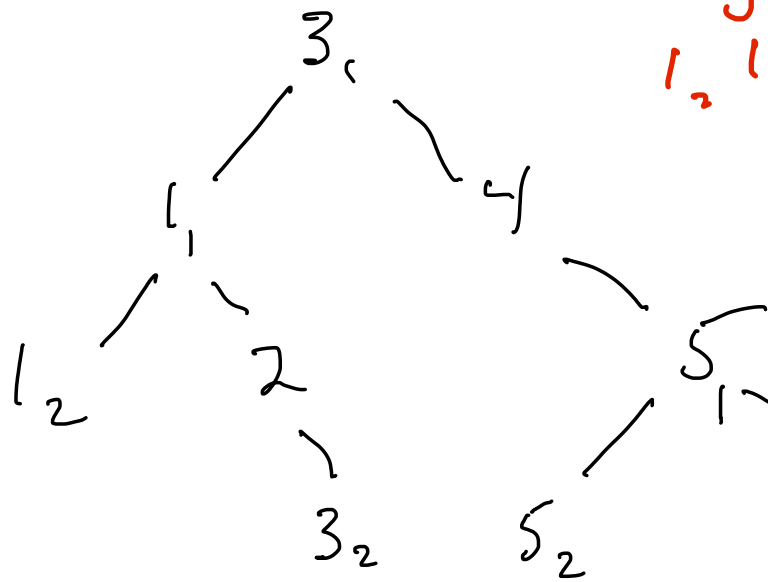


Data structures Packet 3

3₁, 1, 4 1₂, 5, 9 2 6 5₂ 3₂

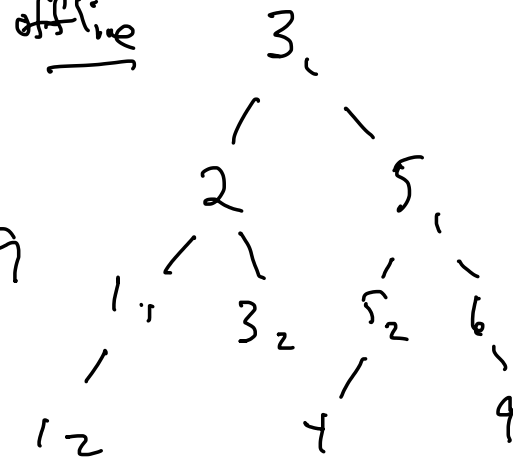
binary tree

online



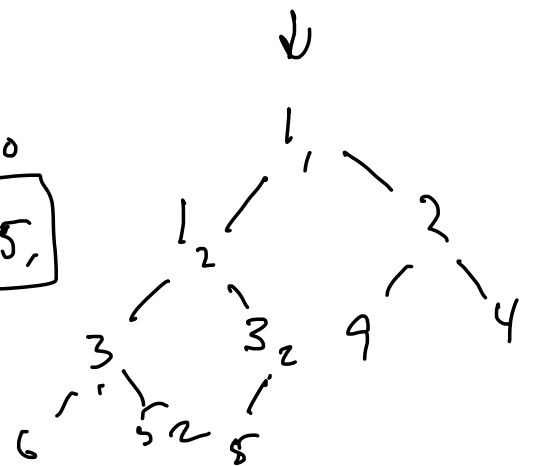
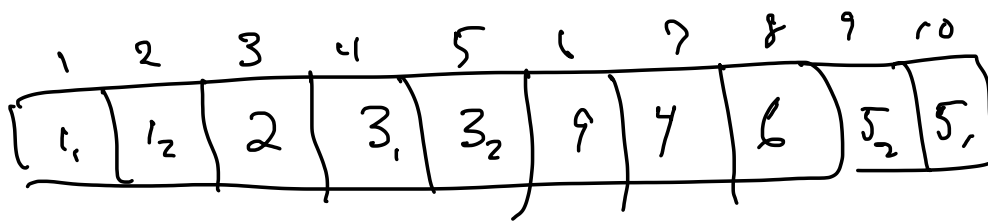
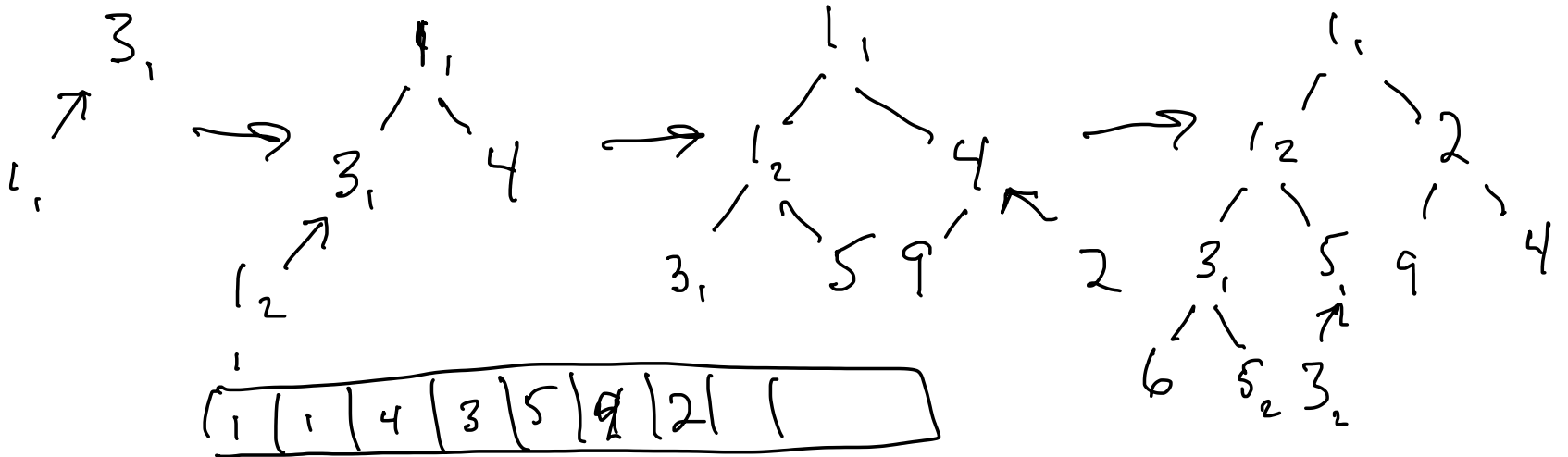
Symmetric (inorder)
1₂, 1, 2, 3₂, 3₁, 4, 5₂, 5, 6, 9

offline



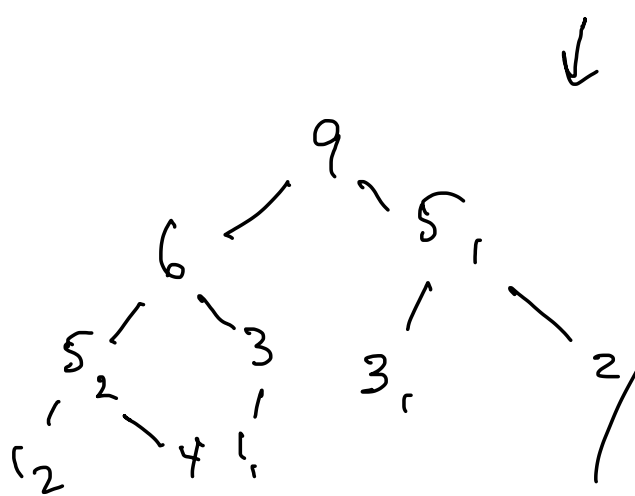
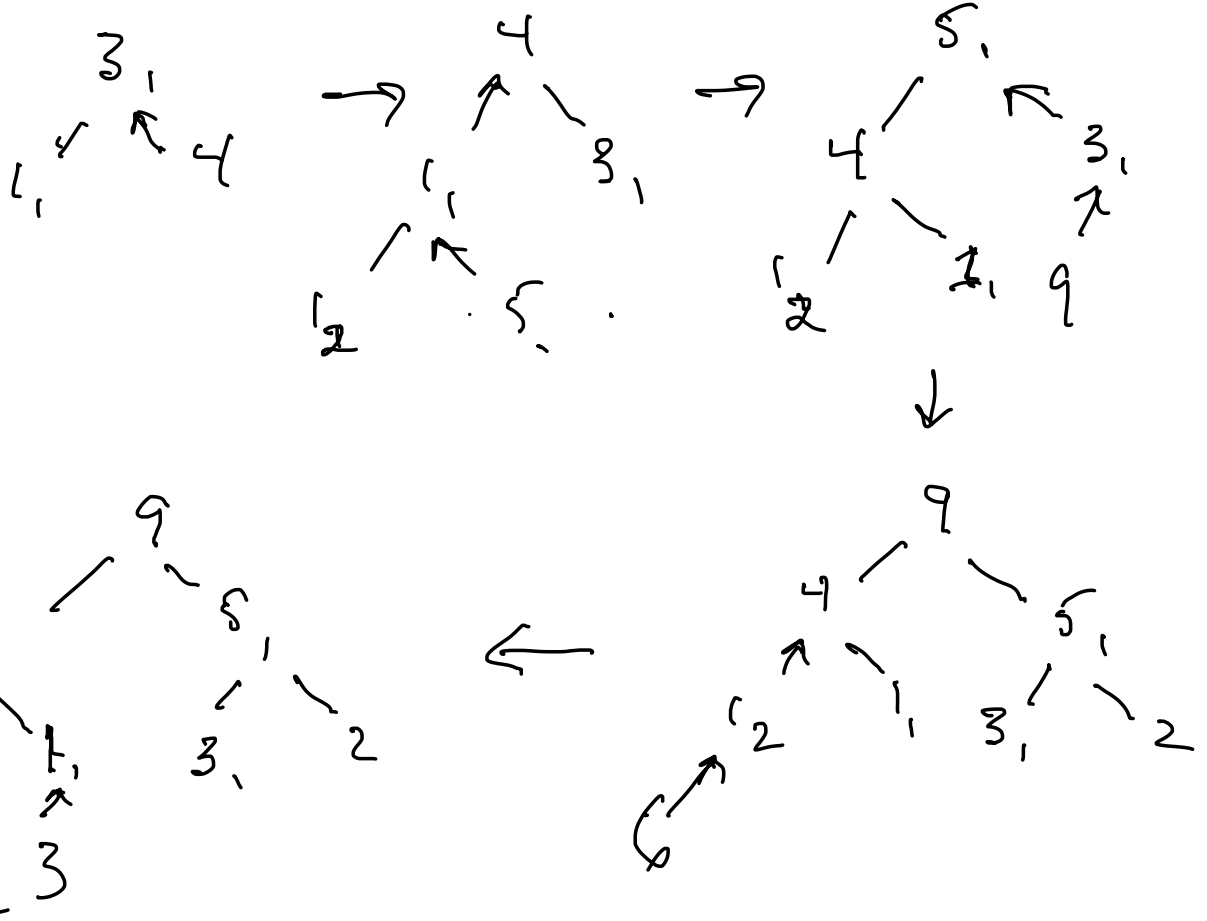
heap (top light)

3₁, 1, 4, 1₂, 5, 9, 2, 6, 5₂, 3₂

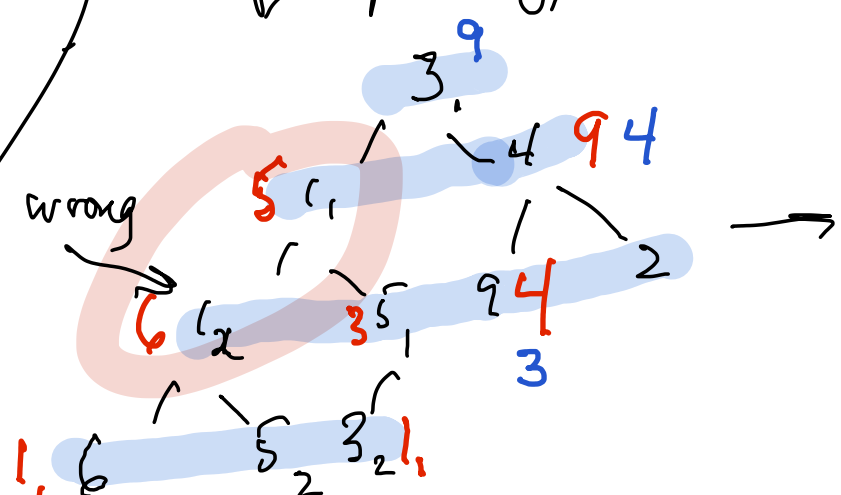


3, 1, 4, 1₂, 5, 9, 2, 6, 5₂, 3₂

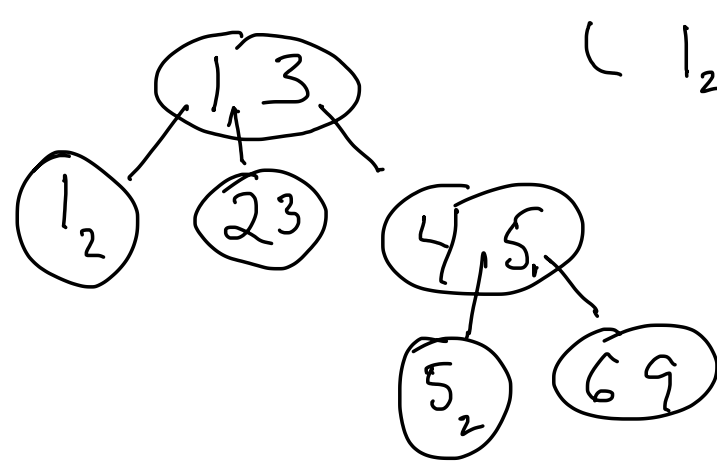
heap (top-heavy)
online



heap top heavy, offline

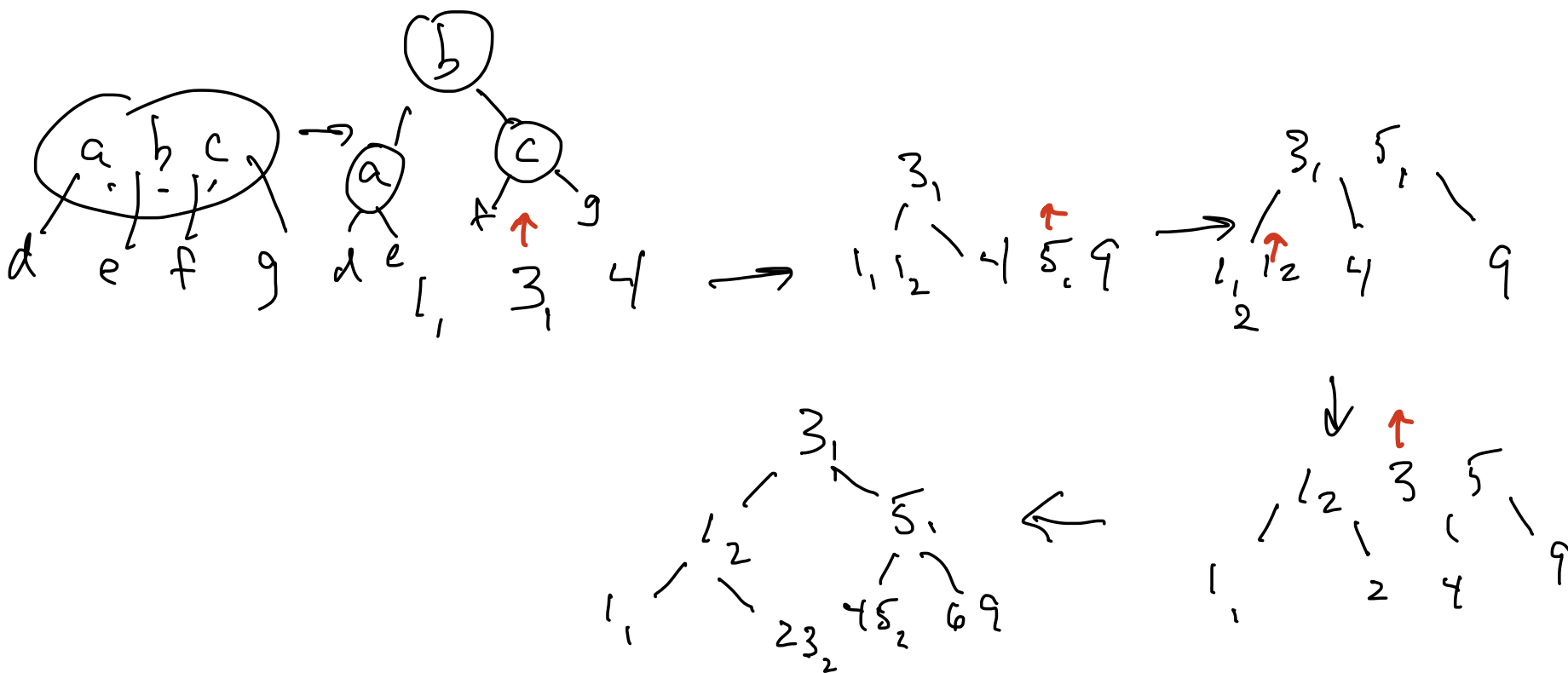


ternary tree (online) 3, 1, 4, 1₂, 5, 9, 2, 6, 5₂, 3₂



(1₂) 1, (2 3) 3 (4(5) 5(6,9))

2-3 tree (online) $3_1, 1_1, 4_1, 1_2, 5_1, 9_1, 2_2, 6_2, 5_2, 3_2$



pre-order traversal:

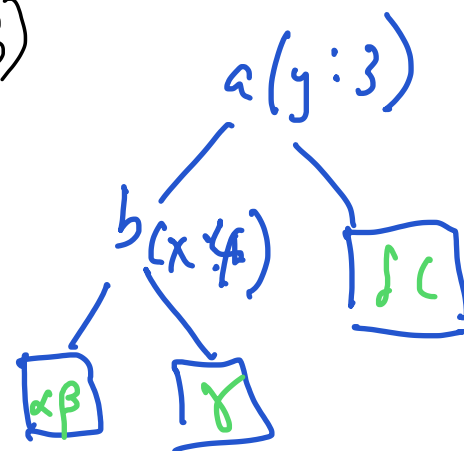
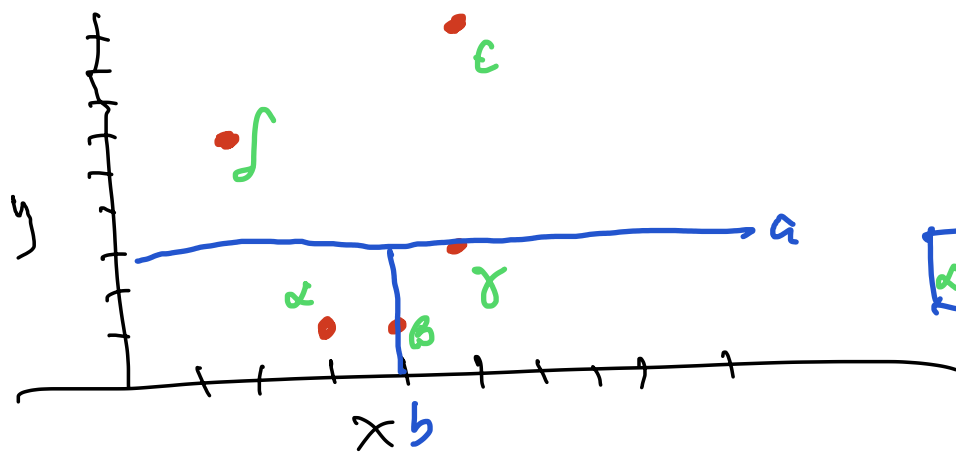
$3_1, 1_2, 1_1, 2, 3_2, 5_1, 4, 5_2, 6, 9$

post-order traversal:

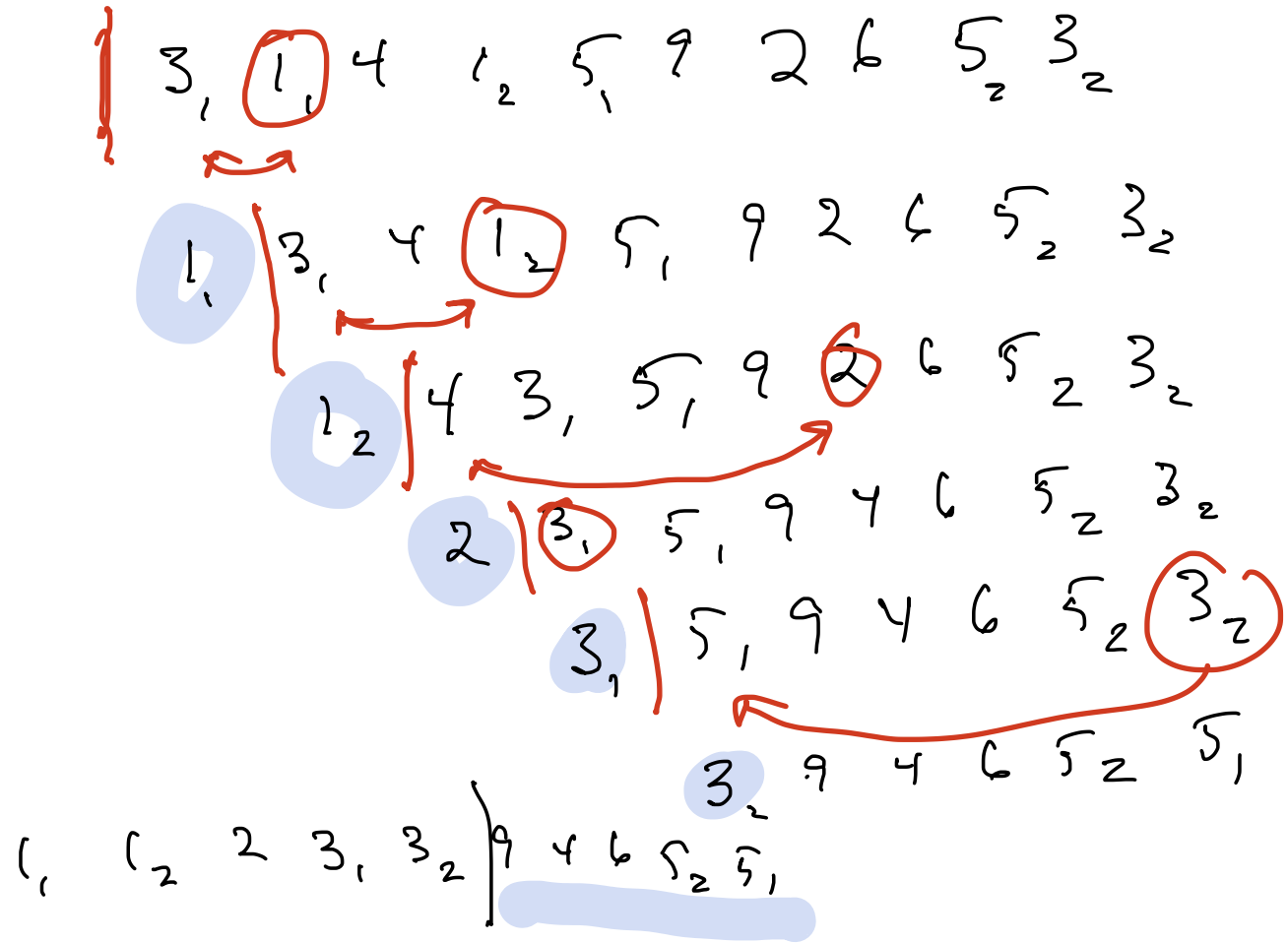
$1_1, 2, 3_2, 1_2, 4, 5_2, 6, 9, 5_1, 3$

k-d tree (2 dim, offline, bucket size 2)

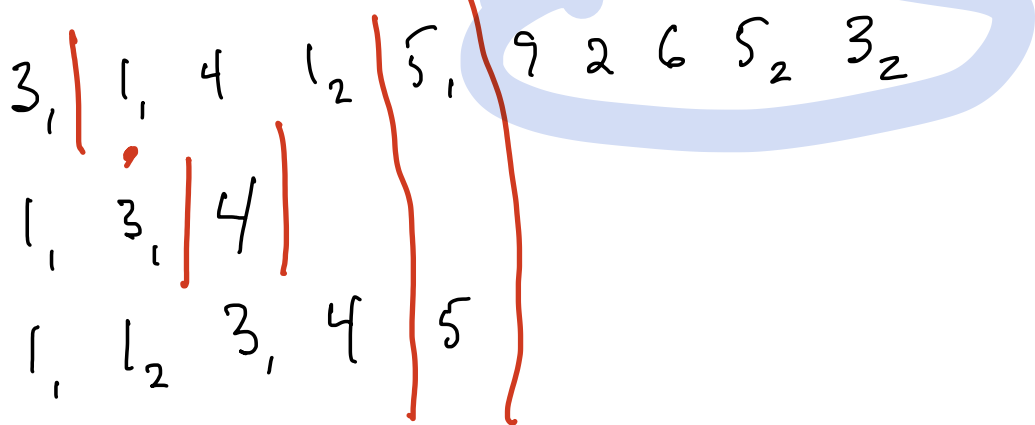
$(3, 1), (4, 1), (5, 9), (2, 6), (5, 3)$



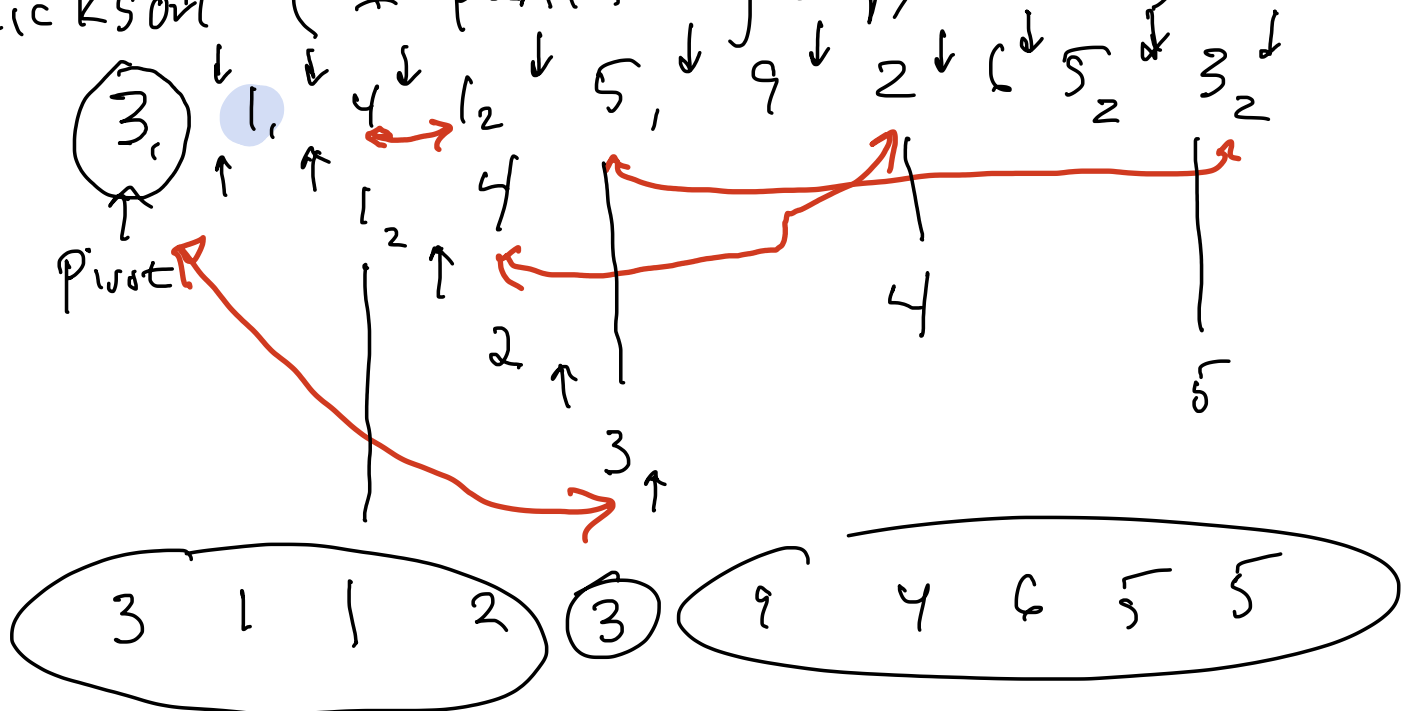
Selection sort : 5 steps.



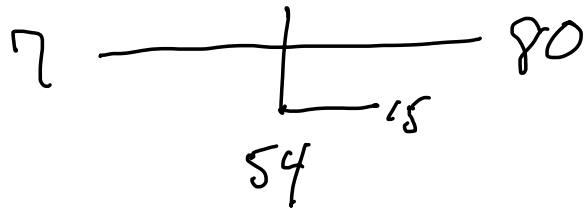
Insertion sort



Quicksort (1 partitioning step, Lomuto)



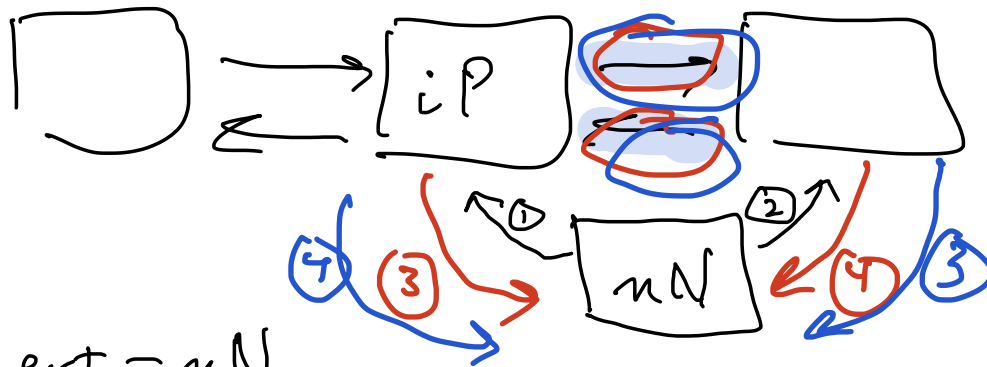
Score Range on Midexam



2b (423):

$\Theta(2) \Rightarrow \Theta(1)$

1a



- (3) $iP \rightarrow next = mN$
- (4) $mN \rightarrow next \rightarrow pred = mN$

- (3) $iP \rightarrow next \rightarrow pred = mN$
- (4) $iP \rightarrow next = mN$

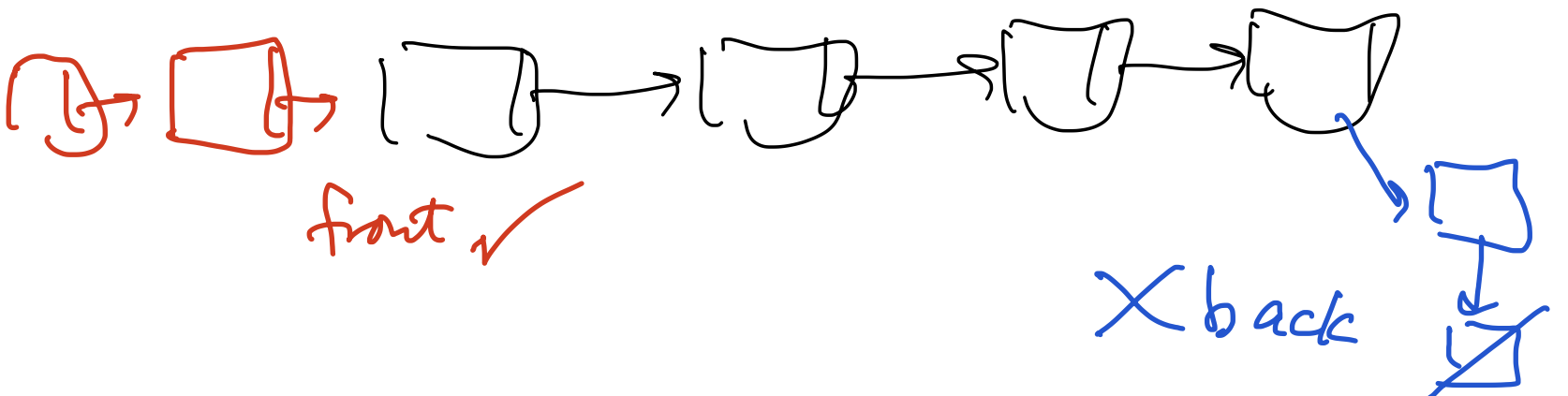
1b complexity: $\Theta(n)$ *average or worst-case*

insertion:

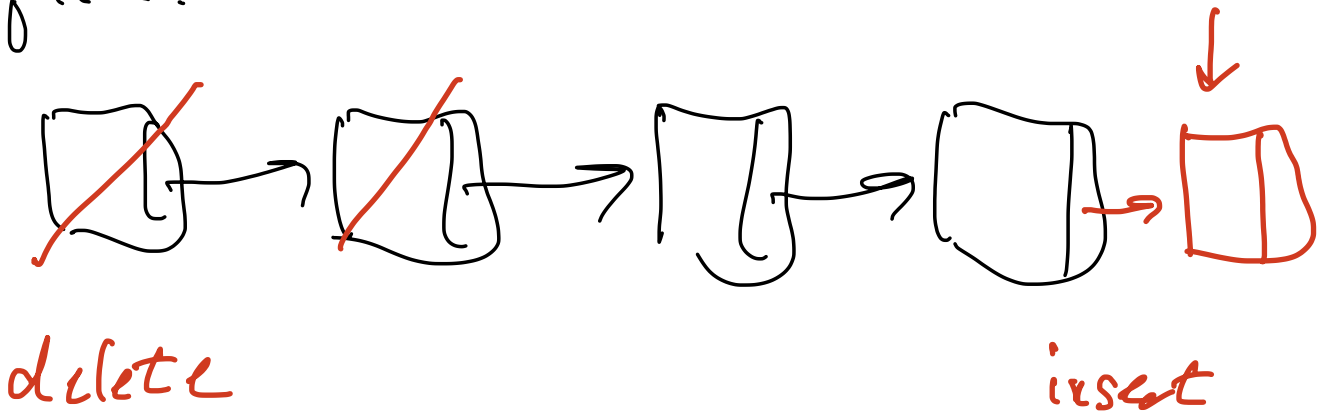
sorted list
singly-linked



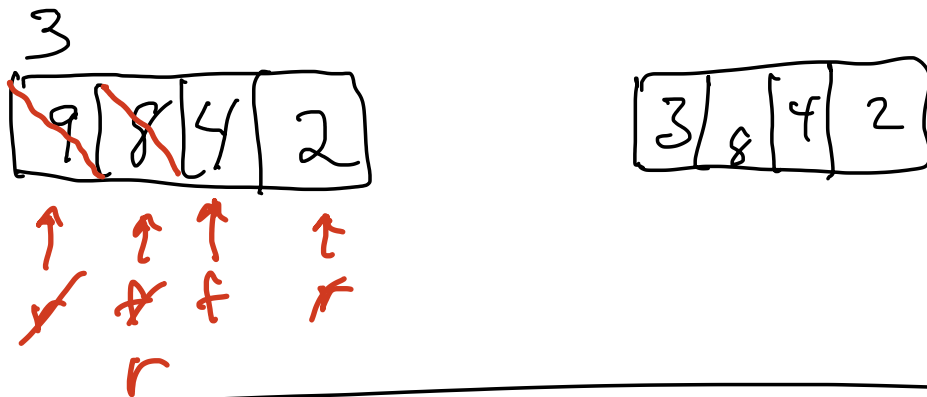
1c stack



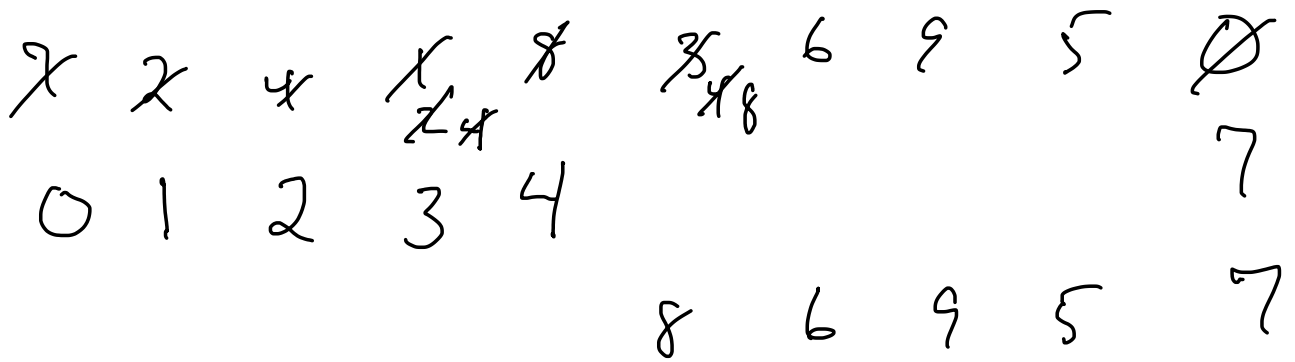
1d queue.



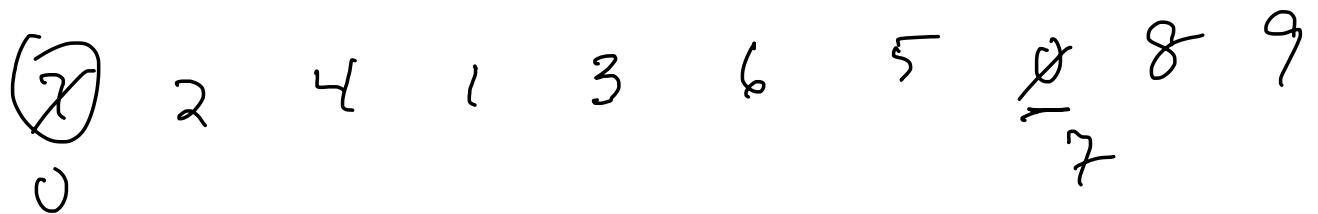
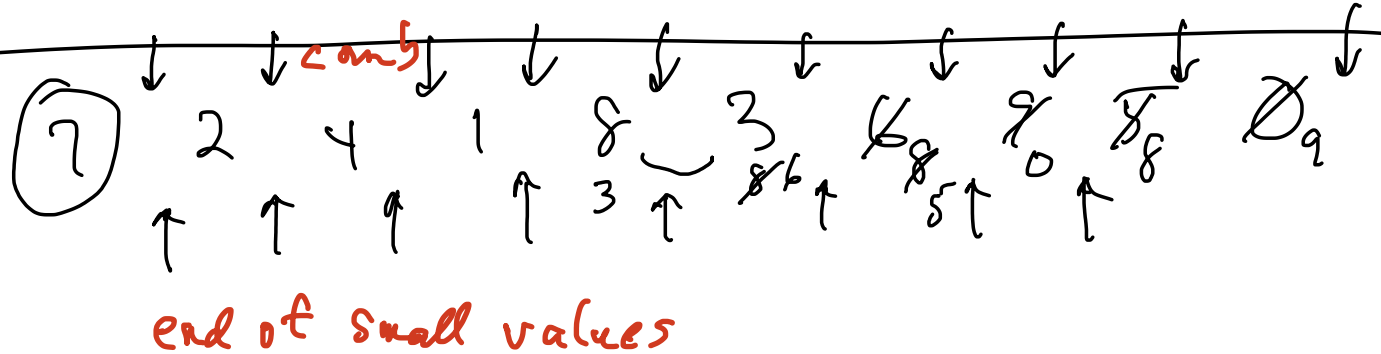
2a



2b

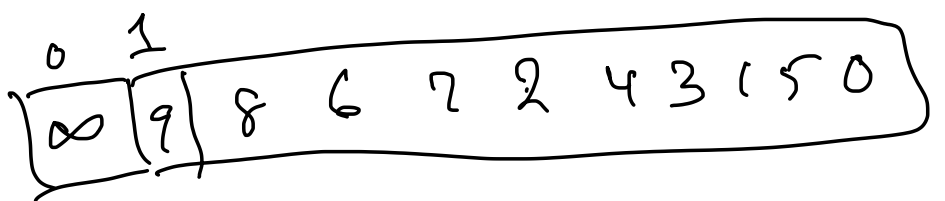
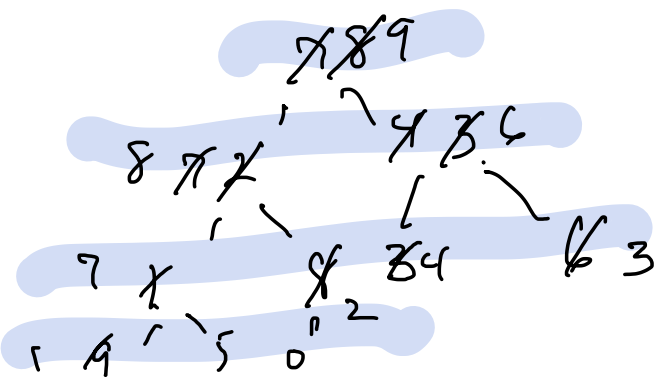


2c



2d

7 2 4 1 8 3 6 9 5 0



3 a

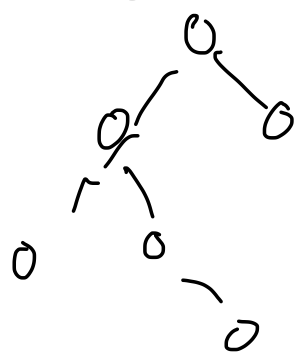
sorted binary tree

expected complexity

$\Theta(\log n)$

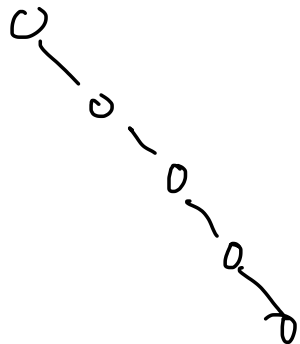
insertion

$\Theta(h)$



b

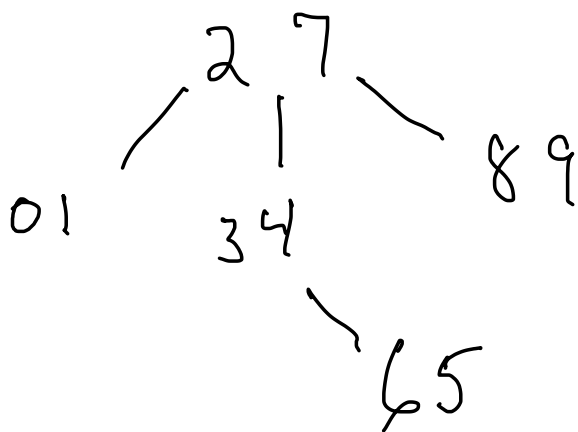
worst case



$\Theta(n)$

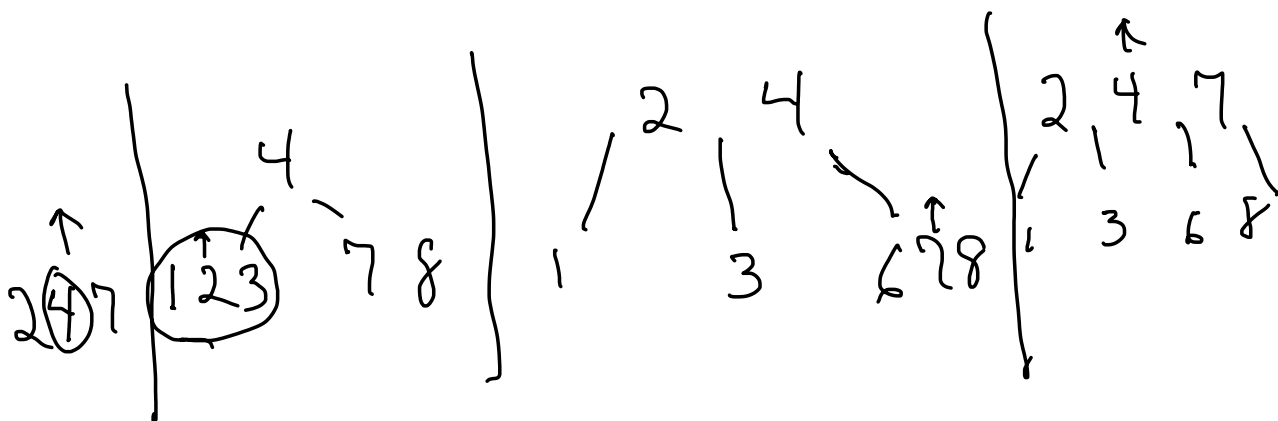
c.

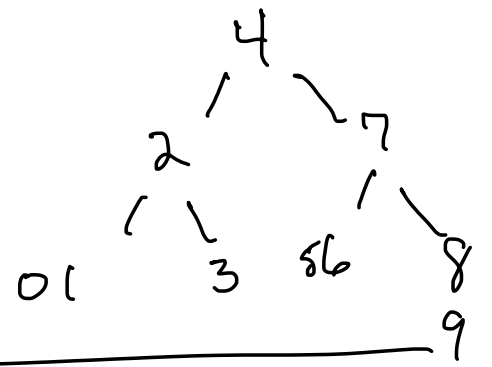
7 2 4 1 8 3 6 9 5 0



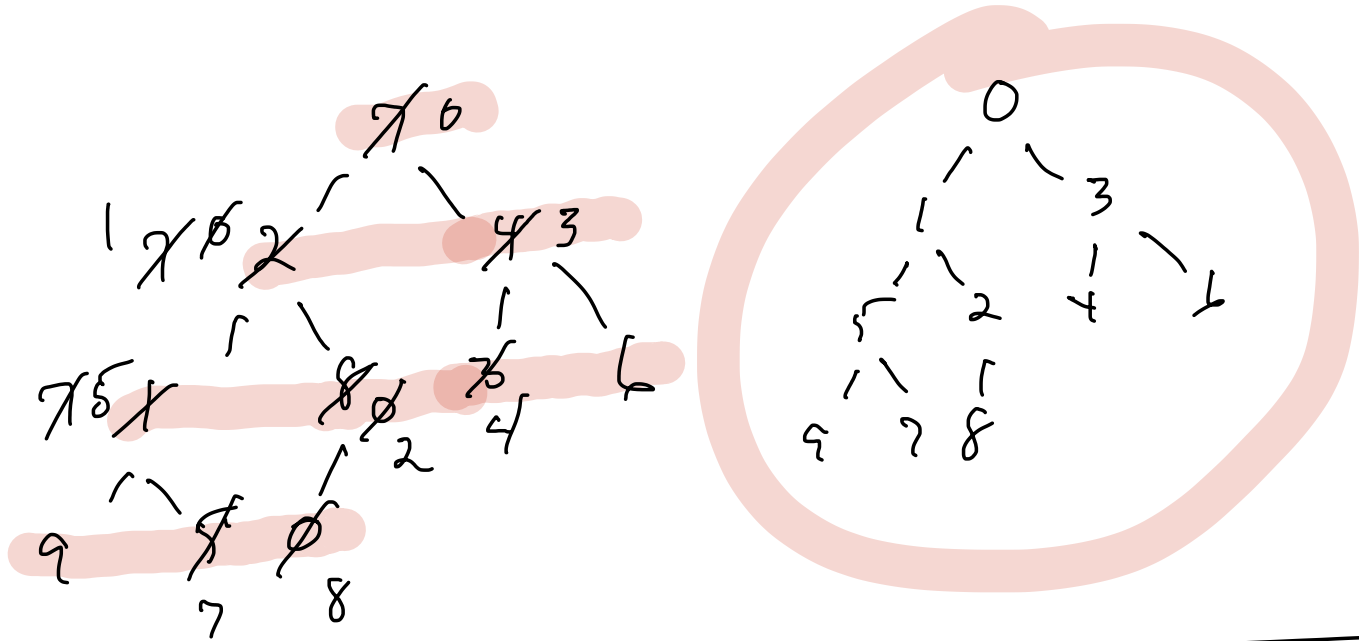
3 d

7 2 4 1 8 3 6 9 5 0





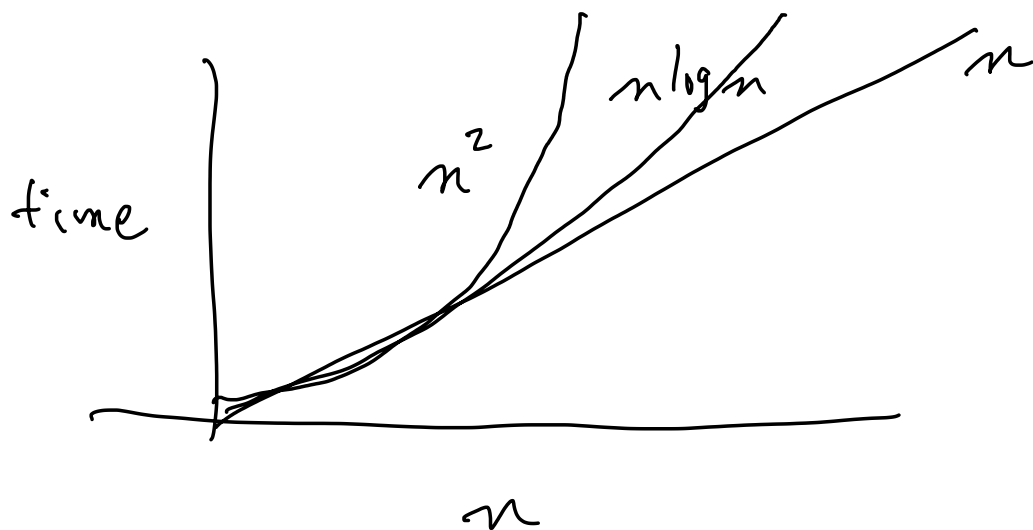
3 e 7 2 4 1 8 3 6 9 5 0



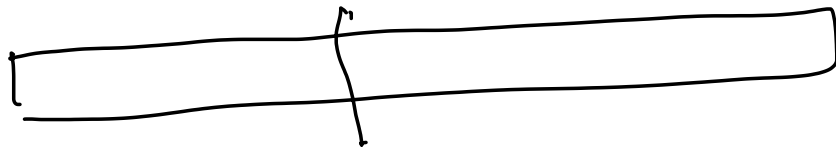
4 a) Why heap instead of bubble sort?

$\Theta(n \log n)$
guaranteed.

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

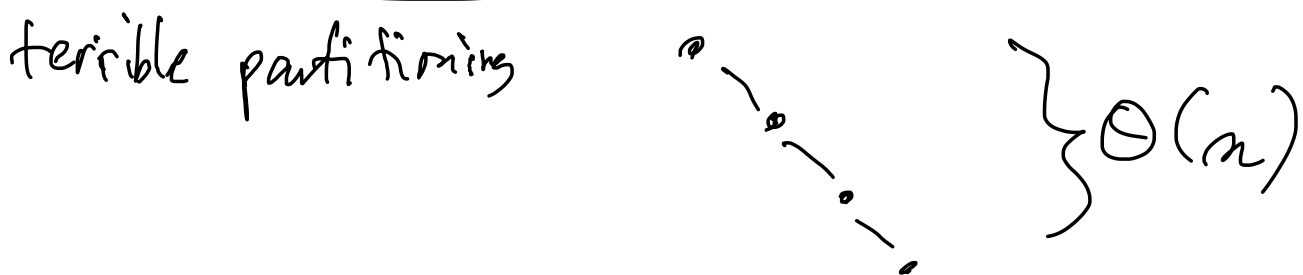
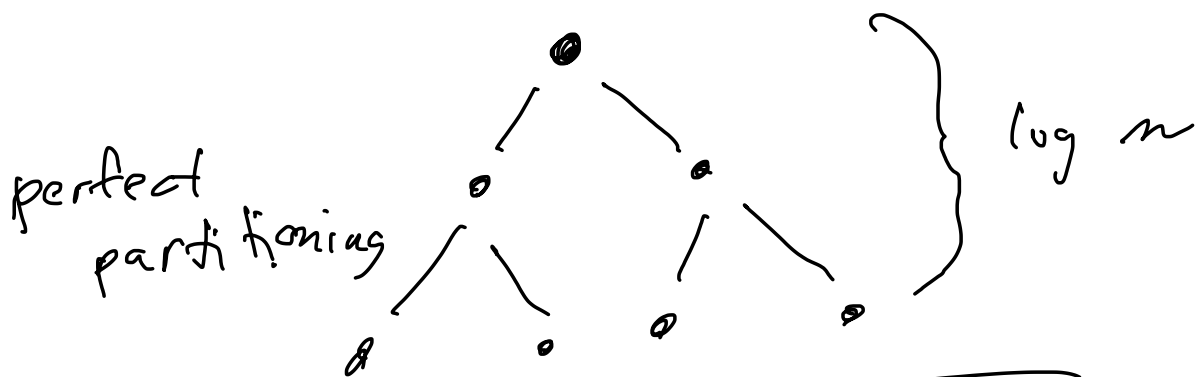


4b. why median of 3 (or 5)?



step 1: more likely to get a more balanced partition.

step 2: better balance \rightarrow shorter navigation tree



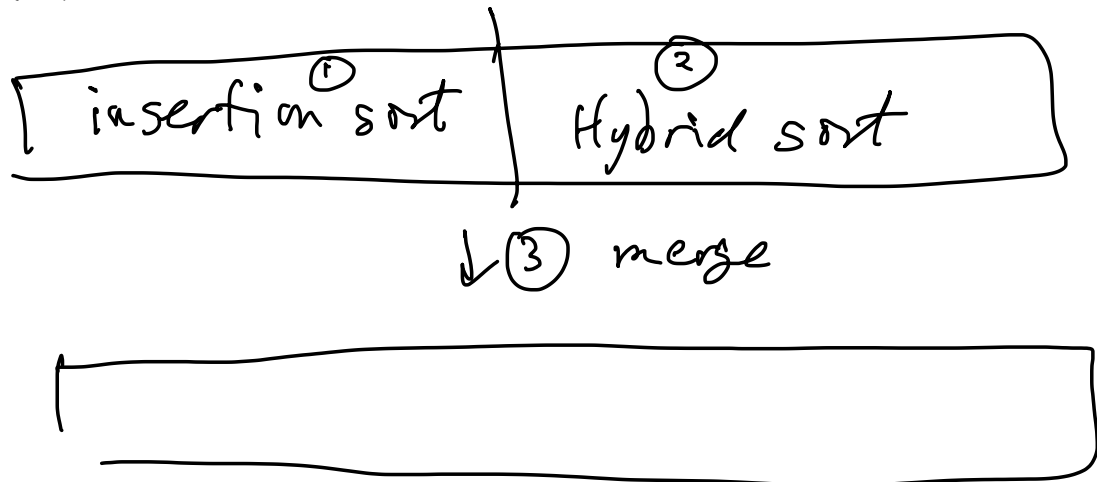
4c: 2nd largest element.

sort, index: $O(n \log n)$

quick select $O(n)$ expected.

1 pass: $O(n)$

4d: Hybrid sort



$$C_n = f(n) + a C_{n/b}$$

$$n^2 + n + 1 C_{n/2}$$

$$C_n = n^2 + C_{n/2}$$

$$\left. \begin{array}{l} k \geq 2 \\ a = 1 \\ b = 2 \end{array} \right\}$$

$$\begin{array}{cc} a & b^k \\ 1 & 4 \\ \hline \Theta(n^k) \end{array}$$

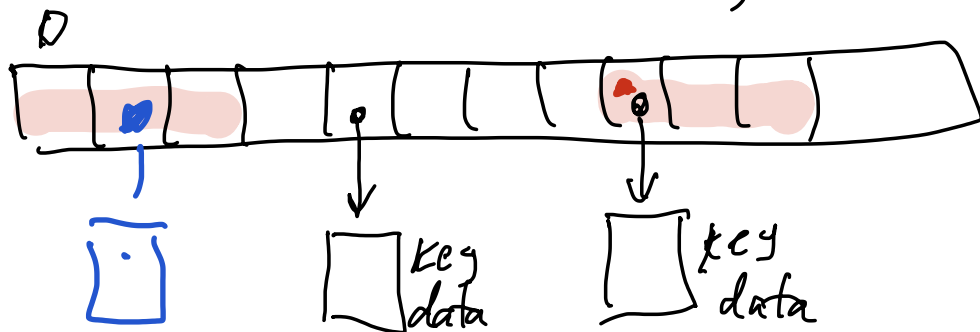
$$\Theta(n^2)$$

Hashing

Key (typically a string) k

apply a hash function $h(k)$

returns number, used as an array index



open addressing: if there is a collision, place the new key in some other cell.

method: linear probing

$$p_i = (h(k) + i) \bmod s \quad i = 0 \dots$$

very bad: clusters

advice: $s \gg 3 \cdot (\text{max \# of entries})$

to search:

for each $p_i = (h(k) + i) \bmod s$ {

if $A[p_i]$ is empty, return failure.

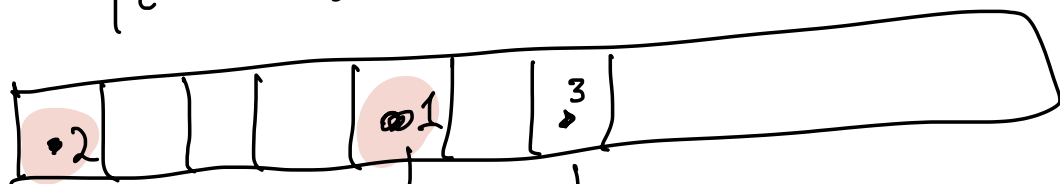
if $A[p_i] \rightarrow \text{key} == k$, return success ($A[p_i] \rightarrow \text{data}$)

}

// very bad case: table is full
return failure.

method: family of hash functions $h_0(), h_1(), \dots$

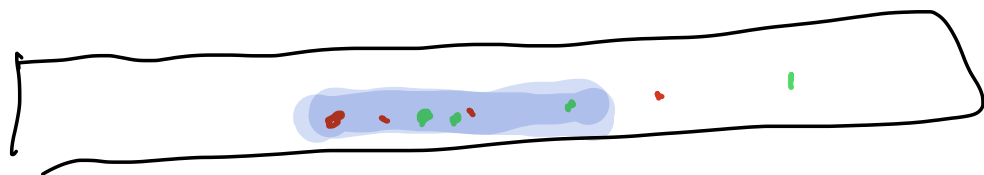
$$p_i = h_i(k)$$



$$h_0(k_1) = h_1(k_2) =$$
$$h_1 \neq \neq$$

method: quadratic probing

$$p_i = (h(k) + i^2) \bmod s$$



still clustering, secondary

method: add-the-hash rehash

$$p_i = (h(k) + (i+1)) \bmod S$$

don't use 0 index of array.

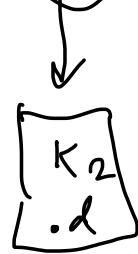
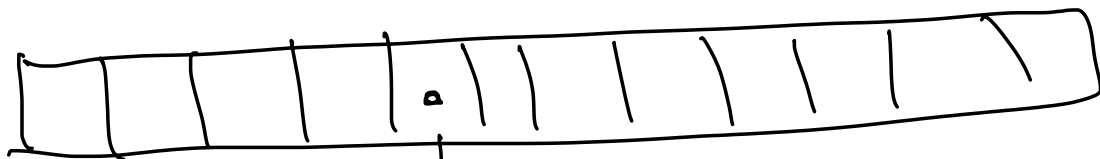
S should be prime so that p_i eventually try all cells.

method: double hashing

$$p_i = (h_1(k) + i h_2(k)) \bmod S$$

↓
never 0

Alternative to open addressing: external chaining



} could be very long.
in practice, short.

$S \approx \#$ of elements to insert

on average, each chain has length 1.

cost of insertion: $\Theta(1)$ (average)

to insert k : (without duplicating)

let $i = h(k)$

for each node d in list headed by $A[h(k)]$

verify $d.key \neq k$.

Build new node $(a) = \{k, data\}$
 Insert (a) in list headed by $A[h(k)]$
 at beginning, especially if search/insert
 show "locality of reference"
 time line



advice: insert at start of chain
 when searching, move success node
 to start of its chain

to search for k :

Let $i = h(k)$

for each node d in list headed by $A[i]$ {
 if $d.key == k$, return success

}

return failure.

Alternative to representing chain as a list

binary tree, 2-d tree, ...

better: use larger s .

What is a good $h()$?

- fast.

- examine all of K . (up to a reasonable limit)

- uniform: $0 \dots s-1$ equally likely.

- ? spreading (not important for external chaining)

similar key should hash to very different values.

In practice: key is a sequence of bytes, b_i

$$\text{method: } \left(\sum_i b_i \right) \underline{\text{mod } s}$$

$$j = \lceil \log_2 n \rceil$$

↑
expected capacity

↓
fast if $s = 2^j$ for some j
mask the lower j bits
eg: $j = 4$ mask (binary)
with 0000 1111

method:
for b_i do:
value \oplus value *
237 + b_i ;
answer \geq value
mod s ;

$$\text{method: } \left(\sum_i b_i \ll i \right) \underline{\text{mod } s}$$

$$\text{method: } \left(\text{XOR } b_i \ll i \right) \underline{\text{mod } s}$$

wisdom: not important what $h()$ you use,
as long as it's not silly.

How big should s (array size) be?

prime if using quadratic probing, ...
(open addressing)

2^j so that mod is fast

if expect n elements, set $s \geq 2n$.
(external chaining)

What if s turns out to be too small?

- 1) Live with it.
- 2) Rehash all elements into a bigger table
(pause in computation)
- 3) extendible hashing: split chains
binary tree, trie (split on last bit)

Modern programming languages have hash tables built in.
(associative arrays)

Perl: $\$a\{k\} = \text{data}$

JavaScript: $a[1]$ $a[\text{'string'}]$
 ↑ ↑
 index key

Java: library HashMap

Python: $\text{Foo} = \text{dict}()$

$\text{Foo}[\text{'this'}] = \text{'that'}$
 ↑ ↑
 key data

you
don't specify
S.

Cryptographic hashes (digests)

$h(\text{text}) = \text{number (represented as a string of hex digits)}$

purpose: uniquely identify text.

goals:

fast computation

uninvertable

given $h(k)$, ~~impossible~~ to derive k .
infeasible

collision-proof:

infeasible to generate a collision

Examples: MD5 $h(k)$ is 128 bits.

Practical attack in 2008

SHA-1 160 bits

Generating a collision in
 2^{69} operations. (2005)

SHA-2

variants: SHA 256 (256 bits)

SHA 512 (512 bits)

Uses:

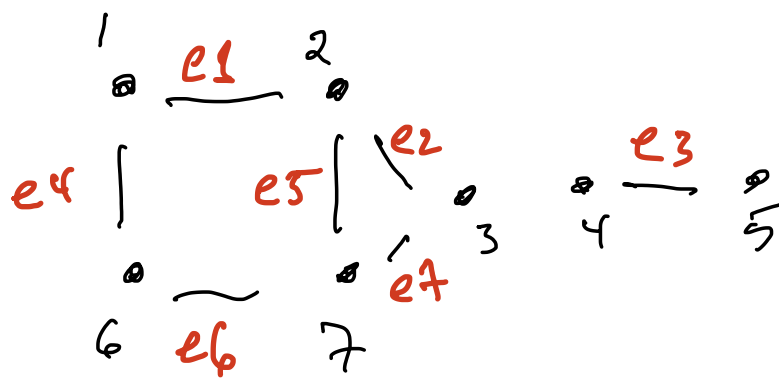
① Storing passwords: store user name, digest (password)
/etc/passwd

② catching copying: verify that excerpts of
submissions are unique.

③ authentication: want to prove that I
sent message m .
shared secret s between sender (A) receiver (B)
 $h(m + s)$

④ intrusion detection: store hash of every essential
program in a protected file.
Compare $h(\text{file content})$ with stored value.
(tripwire) (the digest of a program is called
its signature)

Graphs



vertex (vertices)

edge

graph properties: connected? (no)

directed? (no)

edges have a direction \longrightarrow

weighted? (no)

edges have a numeric value (weight)

vertex properties: sparse? (somewhat) / dense (not very)

fanout / fanin degree

directed

undirected

vertex 2

has degree 3.

graphs represent

1) streets in a city (vertices are corners, edges are streets)

want to compute paths. (sequence of edges)

2) airline routes.

weight of an edge: cost of ticket

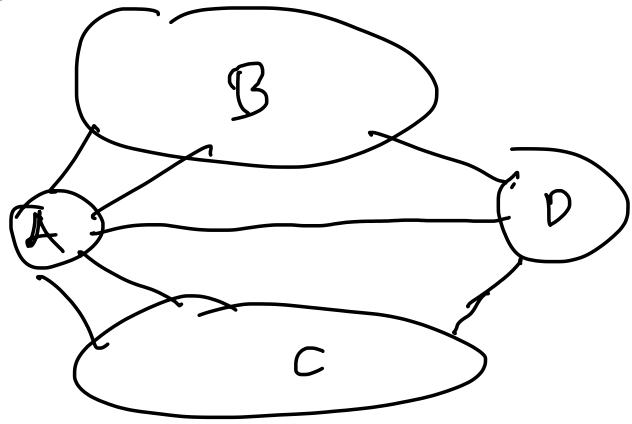
minimal-cost cycles

Hamiltonian cycle: visit every vertex once.

Eulerian cycle: visit every edge once.

1707-1783

3) Bridges of Königsburg (later Kaliningrad)

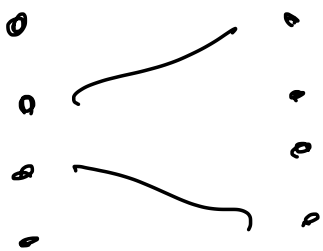


Eulerian cycle exists
iff # vertices
with odd degree
0 or 2, (?)

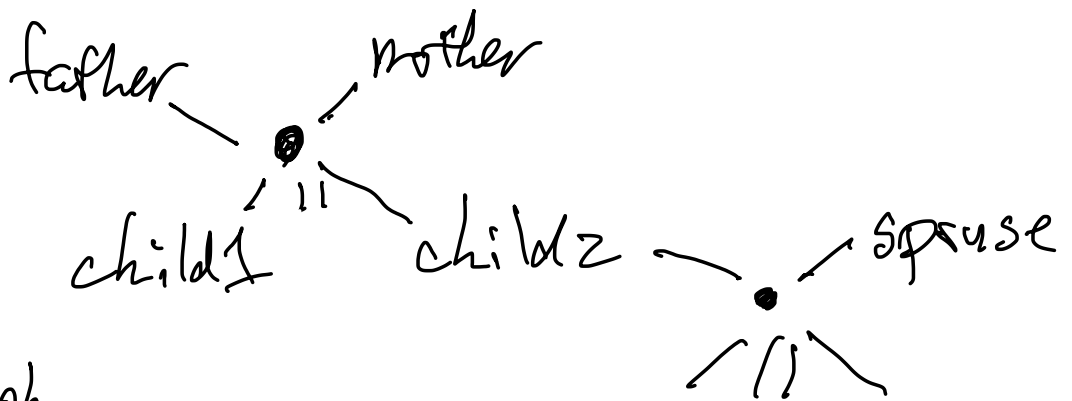
4) Family trees.

Bipartite graph: 2 kinds of vertices

Edges only connect vertices of different kinds.

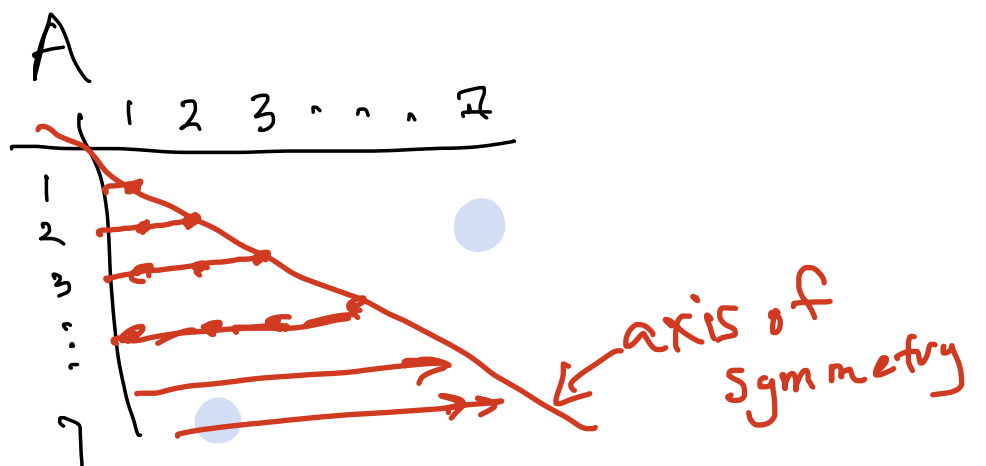


Two kinds: people
families



To represent a graph.

Adjacency matrix



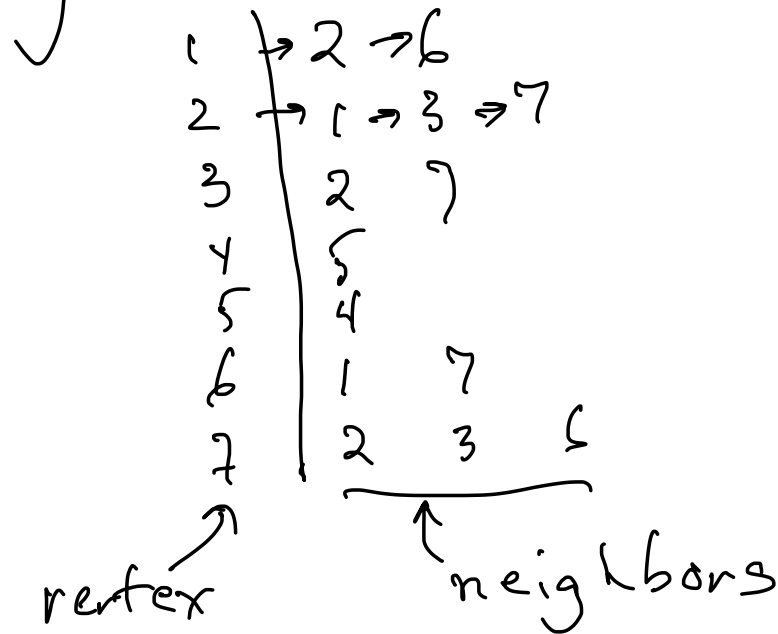
$A[i, j] = \text{true}$ iff \exists edge connecting vertex i to vertex j .
use a 2-dimensional representation!

$A[i, j]$ stored at $A(\frac{i(i-1)}{2} + j)$

$A[3, 7]$

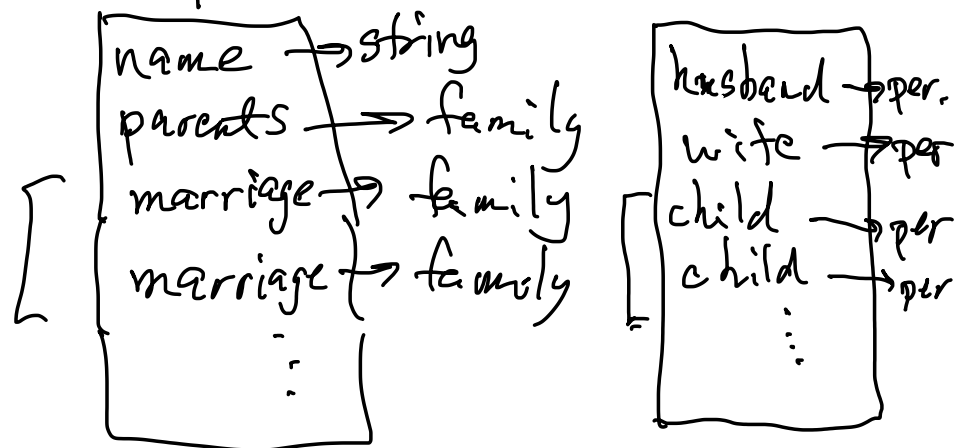
1 dimensional.
 $A[\frac{3 \cdot 2}{2} + 7] = A[10]$

• Adjacency list



• Specialty representations

family tree: structs for people, families



Compute the degree of all vertices
 matrix number of vertices

for each vertex $(0 \dots v-1)$ {
 degree[vertex] = 0;
 for each neighbor $(0 \dots v-1)$
 if $A[\text{vertex}, \text{neighbor}]$
 degree[vertex] += 1;
 }

$O(v^2)$

Adjacency list

$O(v + e)$

```
foreach vertex (0 .. v-1) {  
    degree[vertex] = 0;  
    for (neighbor ∈ L[vertex];  
        neighbor ≠ null;  
        neighbor = neighbor → next) {  
        degree[vertex] += 1  
    }  
}
```

Connected component of vertex i $\vdots \vdots \dots$

Algorithm: Depth-First Search (DFS)

```
void DFS(vertex i)  
    // assume visited[*] = false at start  
    foreach neighbor (i) {  
        visited[neighbor] = true  
        it , visited[neighbor]  
        DFS(neighbor);  
    }  
}
```