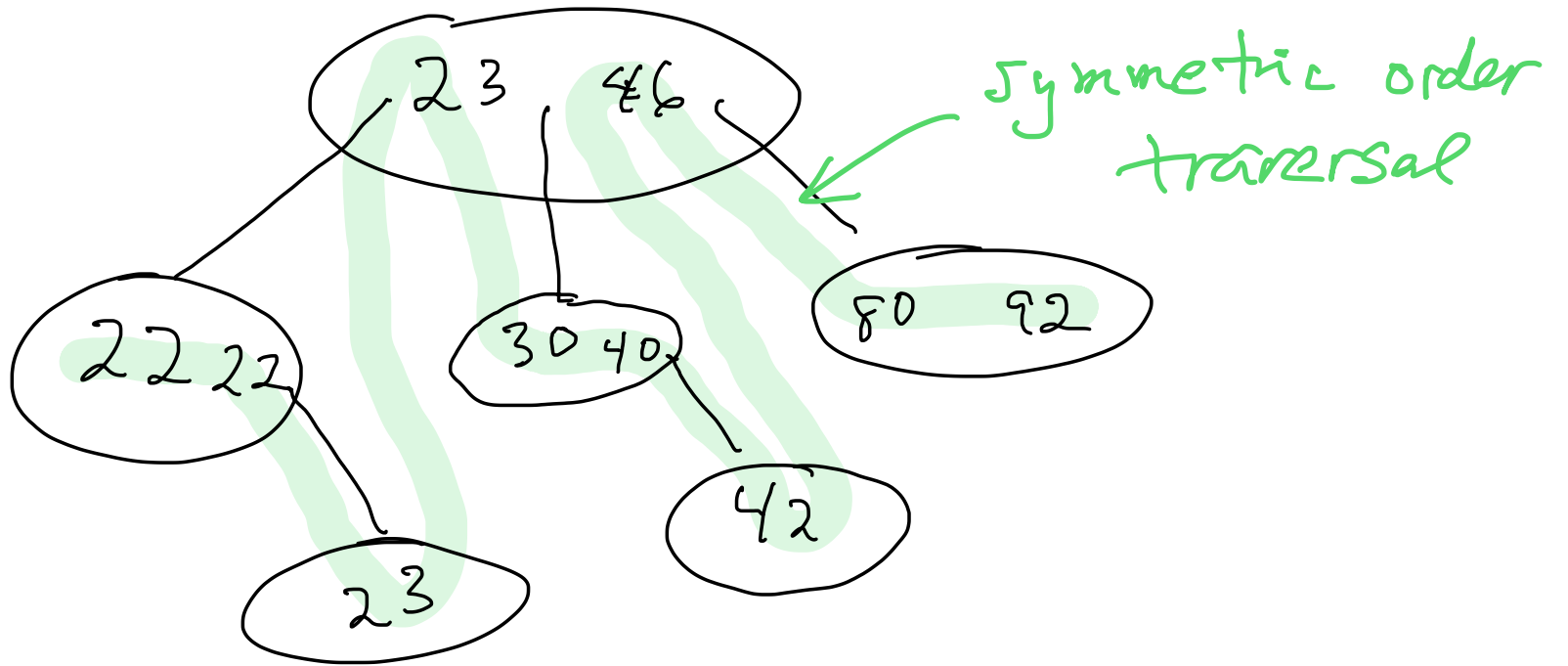


Data structures: package 2

Ternary trees.



$(((22 \ 22 (23)) \ 23 \ (30 \ 40 (42)) \ 46 \ (80 \ 92)))$

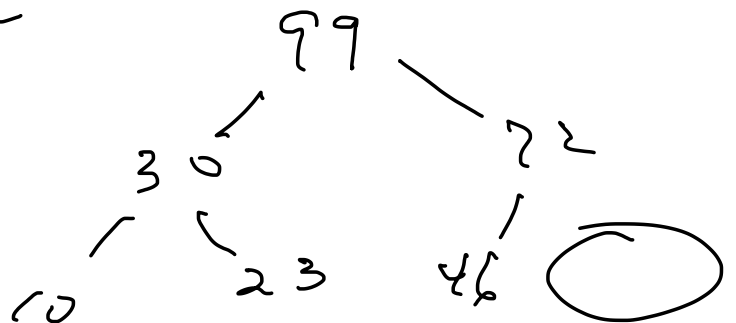
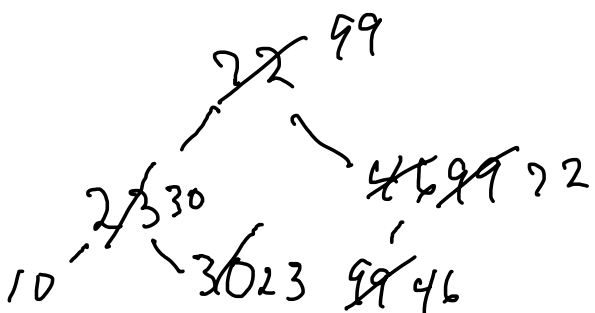
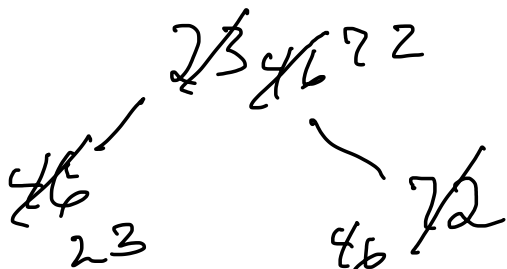
Heaps

(data structure,

as opposed to "the heap" = free storage allocation)

top-heavy: value at a parent \geq values at children.
 (top-light: \leq)

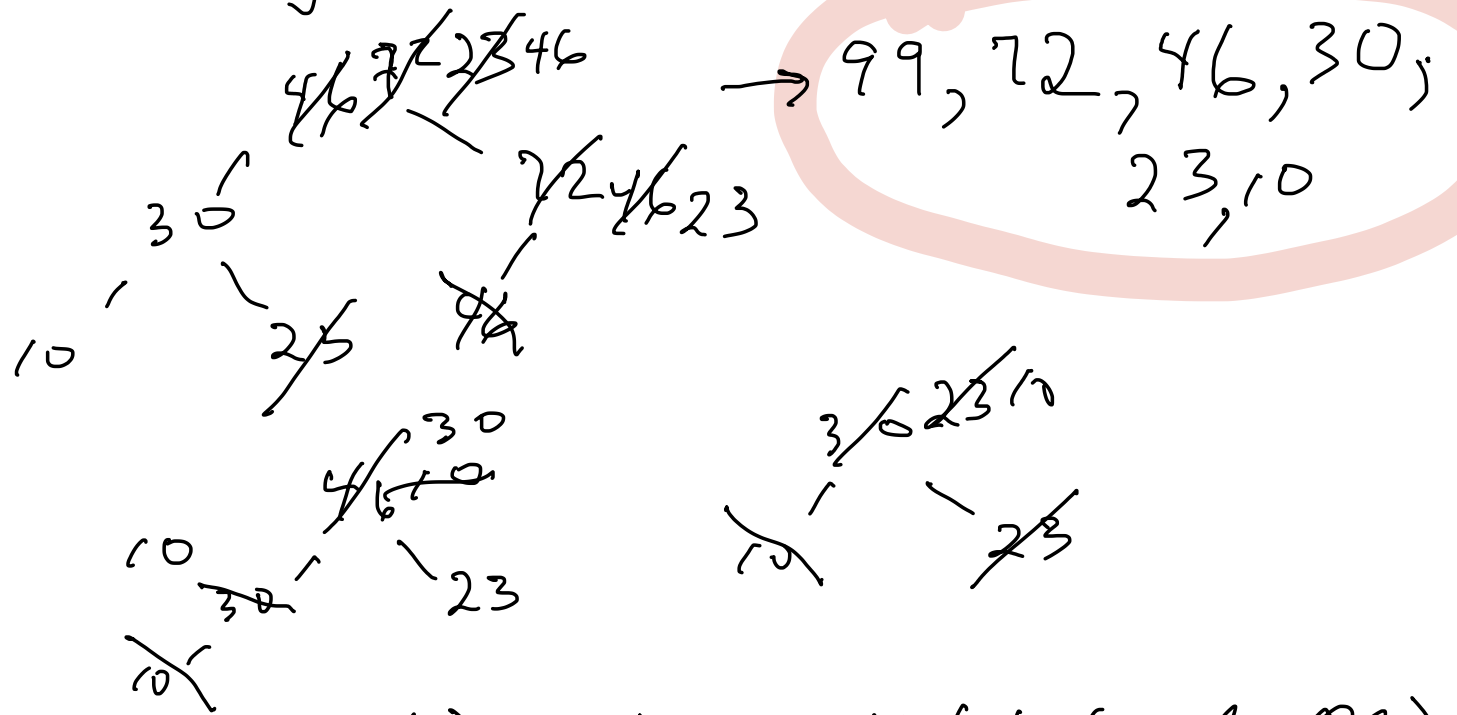
Sift up



insertion: placing $\Theta(1)$ sifting up: $\Theta(\log n)$
 $\Rightarrow \Theta(\log n)$

deletion: always returns top value.

sift down

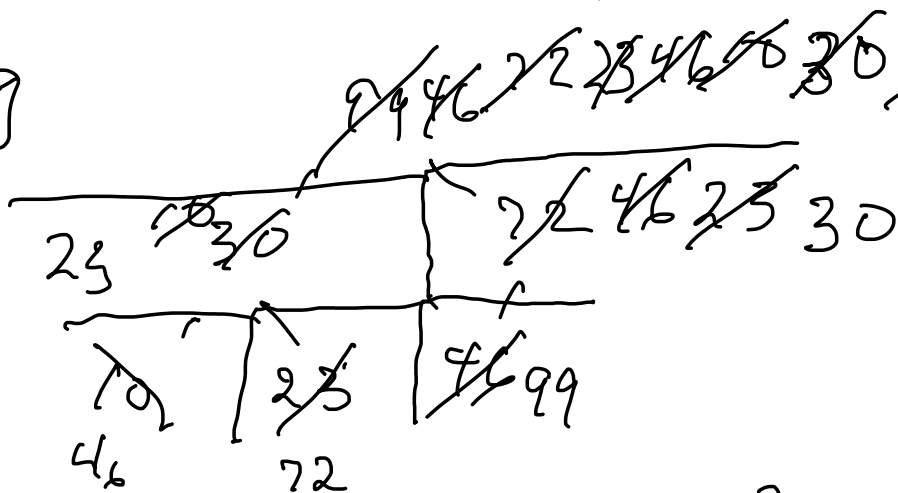


removal: $\Theta(1)$, replacing with last element: $\Theta(1)$,
sifting down: $\Theta(\log n)$
 $\Rightarrow \Theta(\log n)$

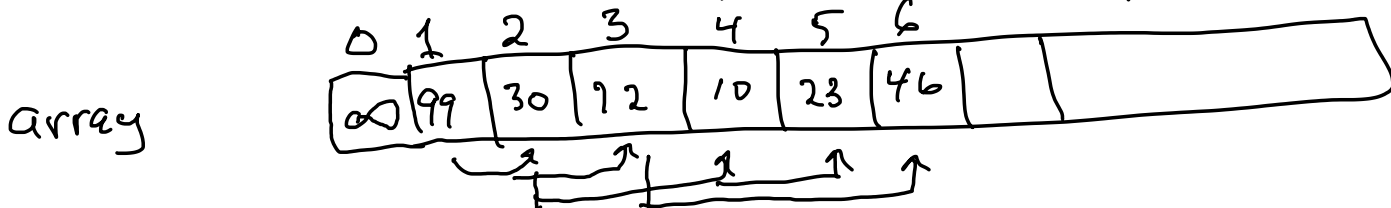
purposes:

priority queue: each element represents work to do
its value represents its priority

Sorting



How to store a heap on the computer?



index	children
1	2, 3
2	4, 5
3	6
n	2n, 2n+1

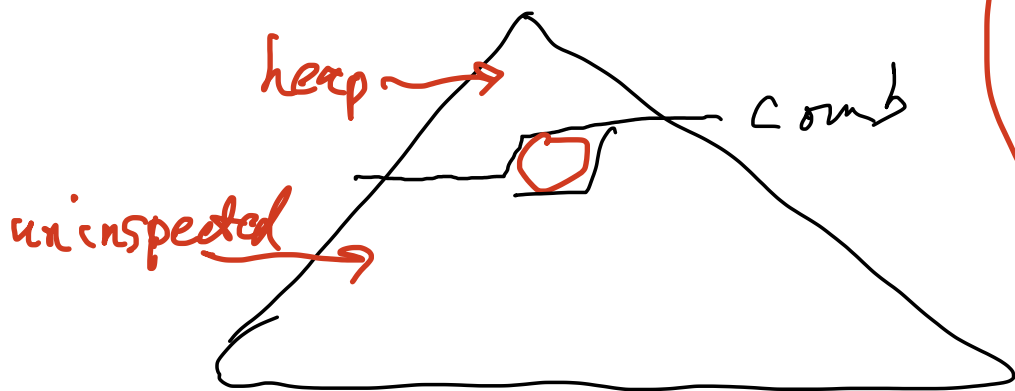
index	parent
1	
2	1
3	1
4	2
5	2
6	3
n	$\lfloor n/2 \rfloor$

to sort n numbers.

I: insert them 1 by 1 into an initially empty heap.
 delete them 1 by 1 into available space at end.

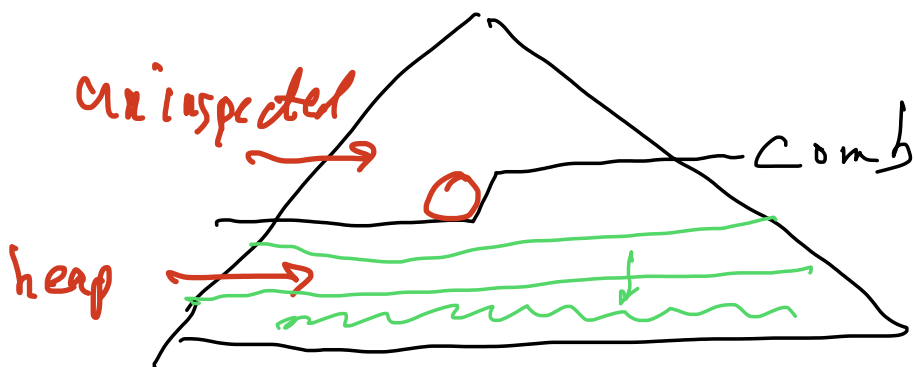
$$\Theta(n \log n) + \Theta(n \log n) = \Theta(n \log n)$$

II: insert them into an unsorted array $\Theta(n)$
 heapify the array: $\Theta(n \log n)$ or $\Theta(n)$
 delete as before: $\Theta(n \log n)$
 $= \Theta(n \log n)$



for each new element, sift up.

$$\Theta(n \log n)$$



for each new element, sift down.

trick: start with comb at half-way point

$\frac{1}{2}n$: no sifting needed

$\frac{1}{4}n$: sift 1 step

$\frac{1}{8}n$: sift 2 steps

$\frac{1}{16}n$: sift 3 steps
⋮

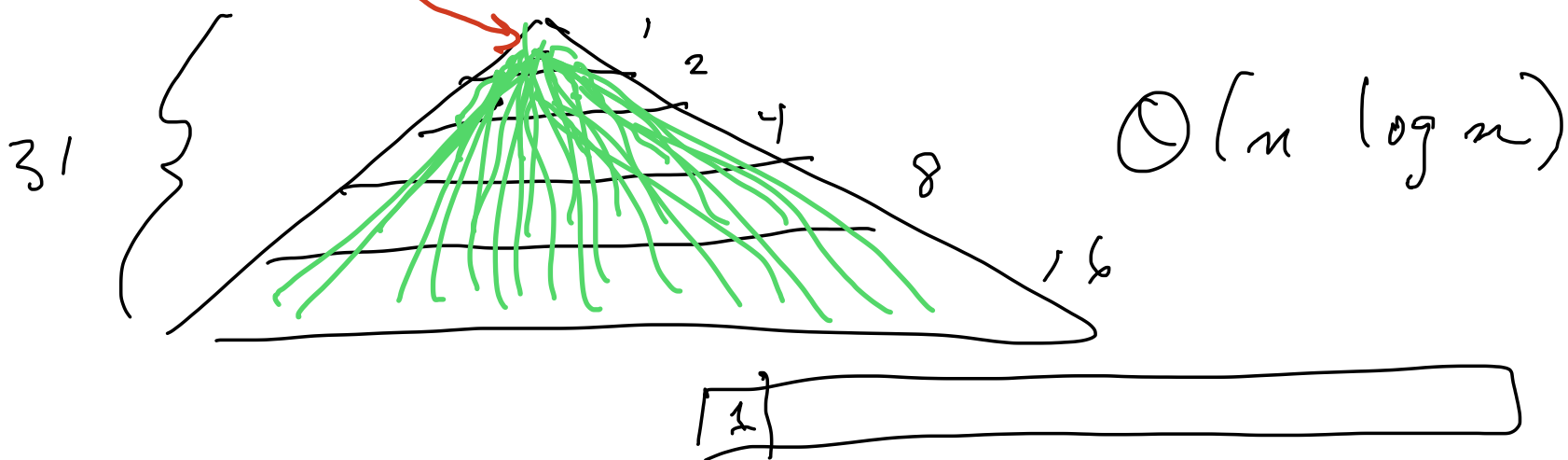
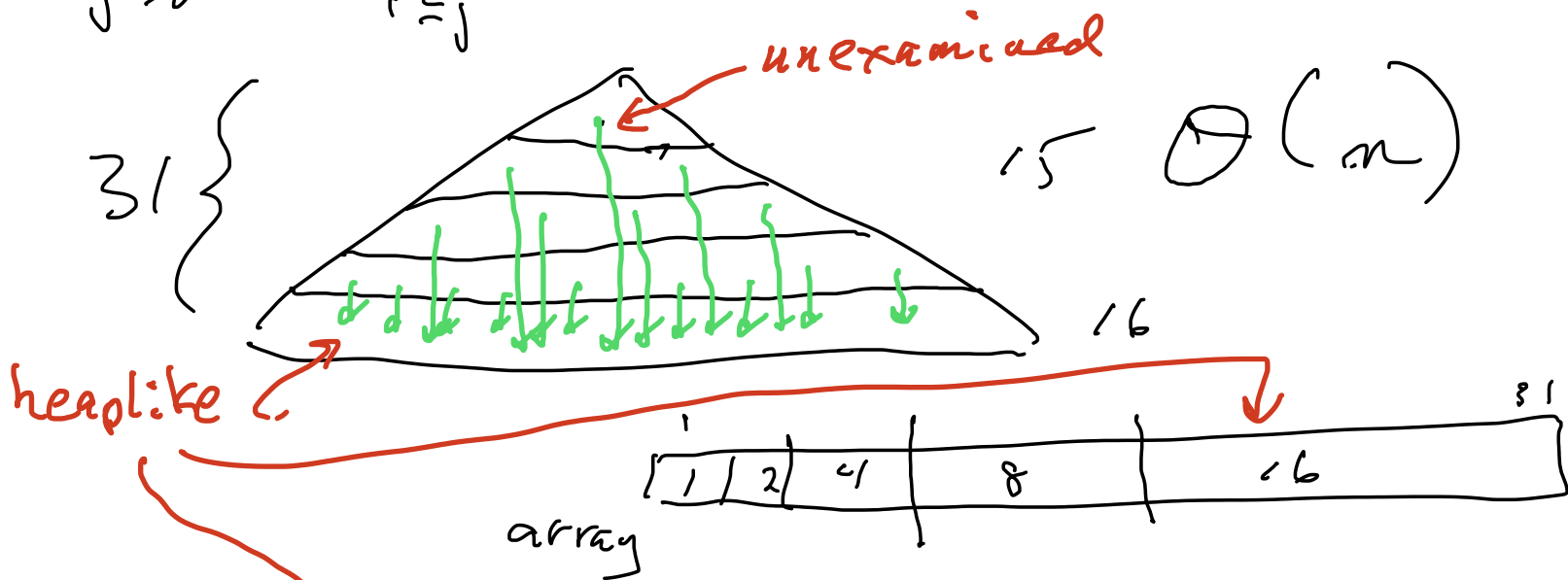
$$\frac{1}{2}n(0) + \frac{1}{4}n(1) + \frac{1}{8}n(2)$$

+ ... =

$$n\left(\frac{1}{4} + \frac{2}{8} + \frac{3}{16} + \frac{4}{32} + \dots\right)$$

$$\leq n$$

$$\lim_{j \rightarrow \infty} \frac{n}{2} \sum_{i=j}^{\infty} \frac{1}{2^i} \rightarrow n$$

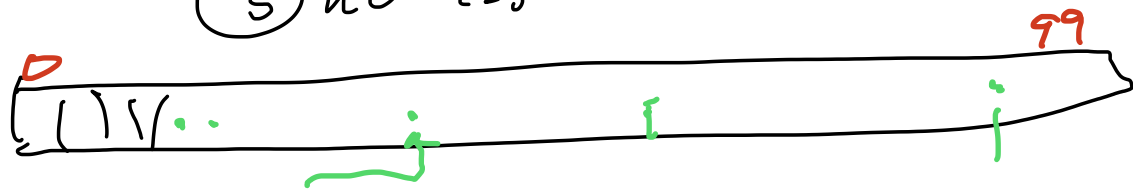


Bin sort

Assume: ① know small range of keys eg. 0-99

② no duplicates

③ no associated data



92, 42, 64, 10, 12, 14

Method: set the bit in array for each key
then read off which bits are set.

Time complexity: $n = \#$ of elements to sort
 $r =$ range of possible values

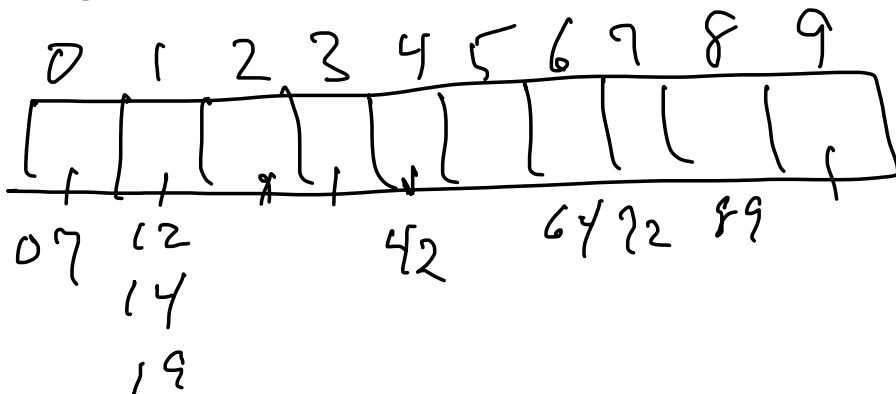
$$\Theta(n) + \Theta(r) = \Theta(n+r) = \Theta(r)$$

↑
because $n \leq r$

Relax assumptions?

- ① X
- ② instead of a bit, every cell of array is an integer count. Space is now $\Theta(r \cdot \log n)$
↑ cells ↑ width of a cell.
- ③ storing data in array instead of a bit
 Space is now $\Theta(r \cdot \text{sizeof}(\text{data})) = \Theta(r)$
 more precisely: $r \cdot \text{sizeof}(\text{data})$ bytes.

Radix sort



stable sort: duplicates in order

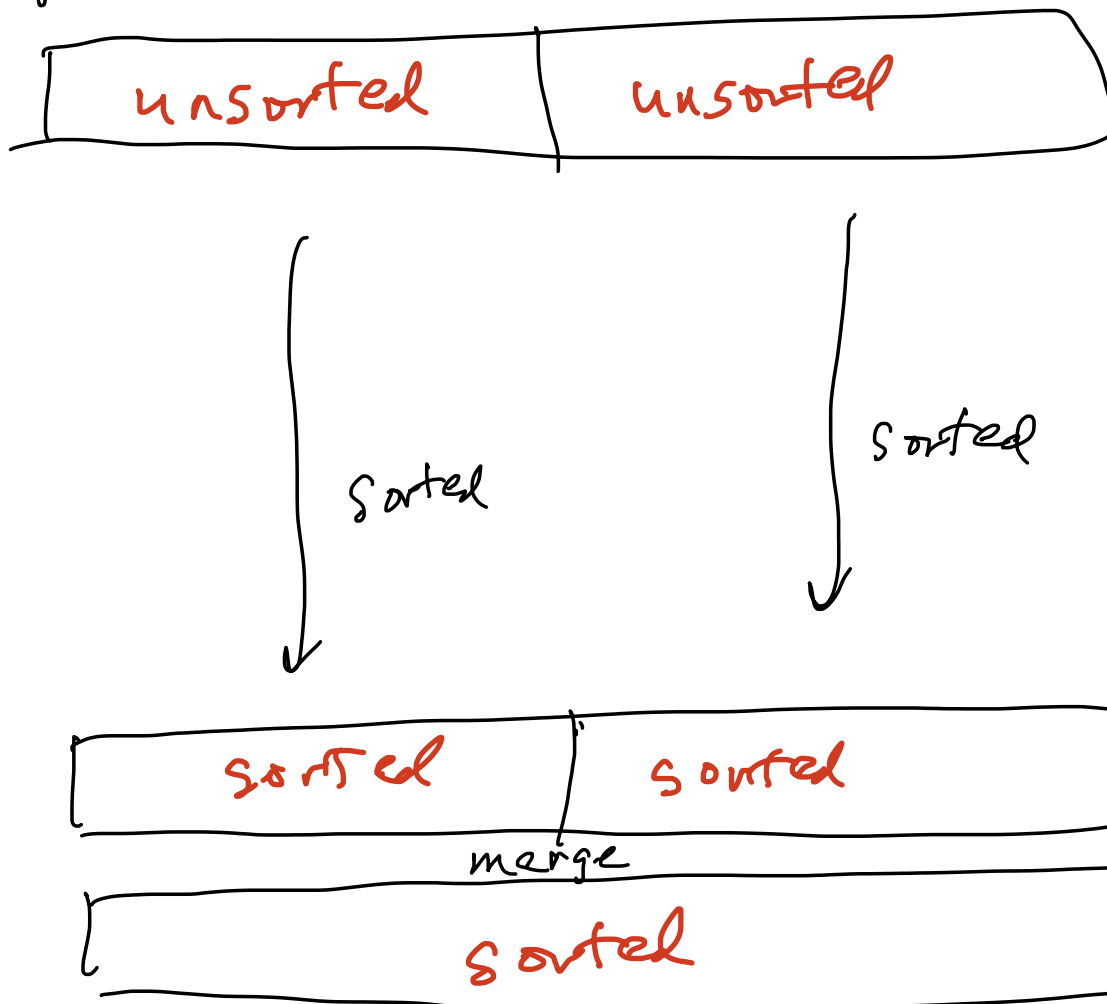


89	42	07
42	72	12
64	12	14
72	64	19
12	14	42
14	07	64
19	89	72
07	19	89
11	11	11
12	12	12

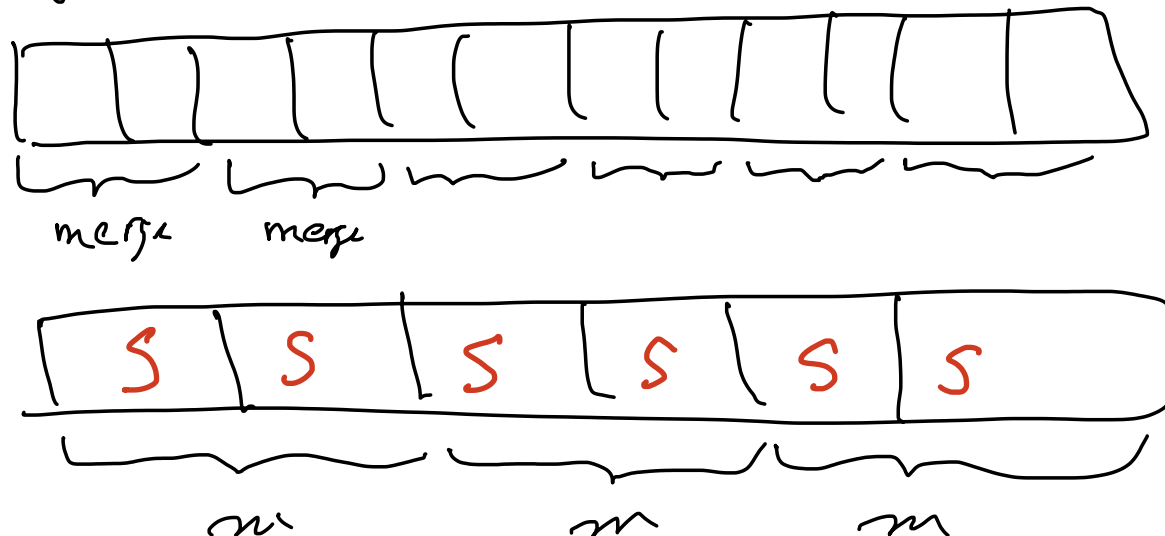
each pass uses a different place, adds at list ends.

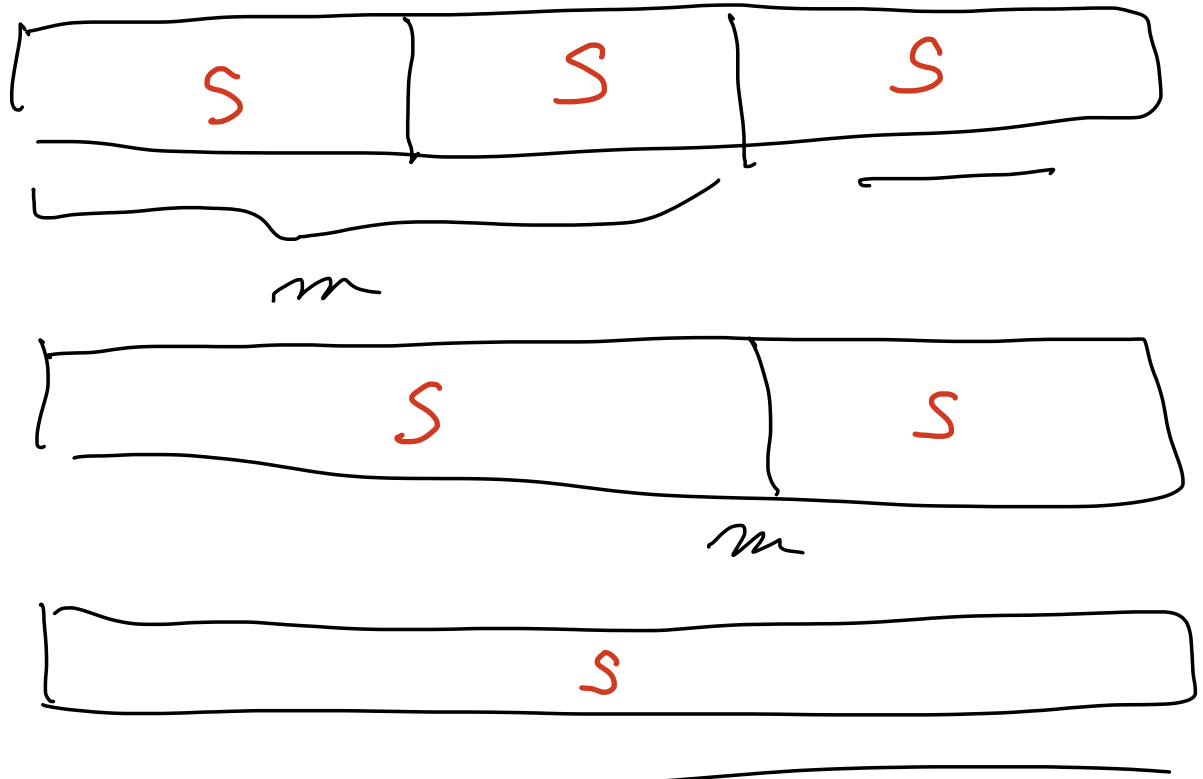
Complexity : each pass : n steps
of passes = # of digits = d
 $O(nd)$
 \uparrow
 $\log(n)$?

Merge Sort

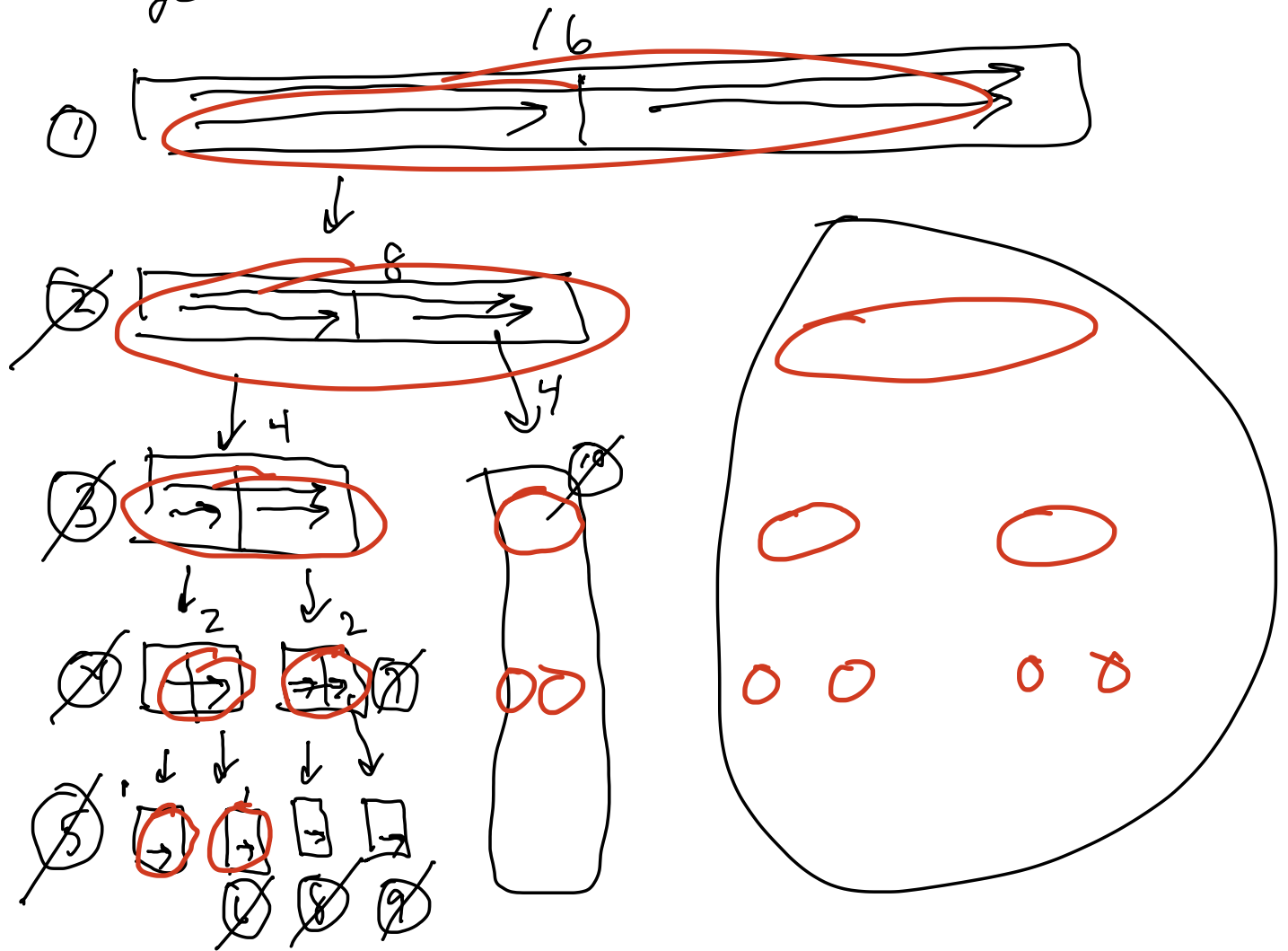


Iterative merge sort





Recursive merge sort



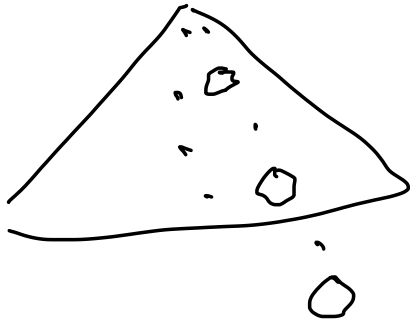
Red-black trees (Guibas + Sedgwick) (1978)

goal: binary trees that are mostly balanced.

"self-balancing" trees

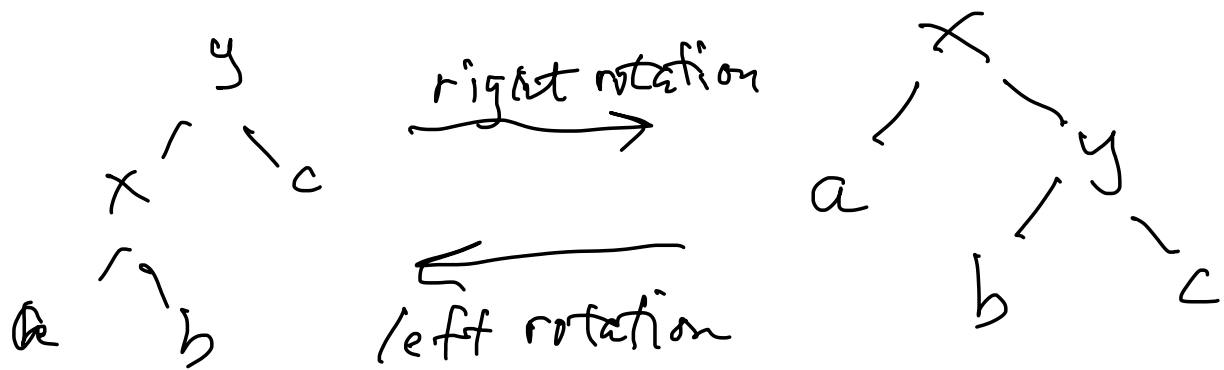
insertion may cause a balancing $\Theta(\log n)$

- Rules:
- 1) each node is red or black.
 - 2) each node is linked to left, right, parent.
 - 3) null nodes are black.
 - 4) root is black
 - 5) red nodes have only black children.
 - 6) all paths from root to a leaf have the same number of black nodes.



To insert:

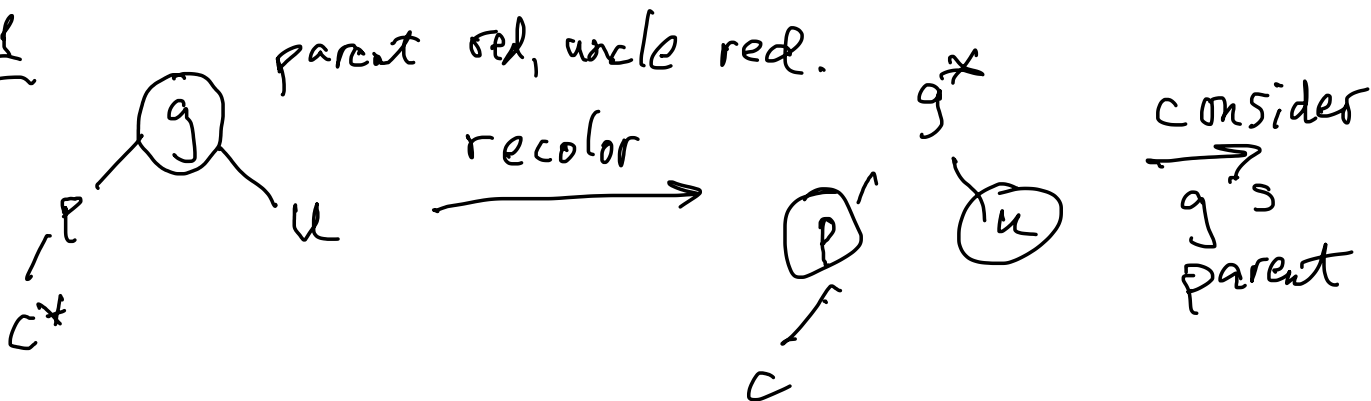
Sometimes we rotate at a node



put new node in the tree (standard insertion)
color it red.

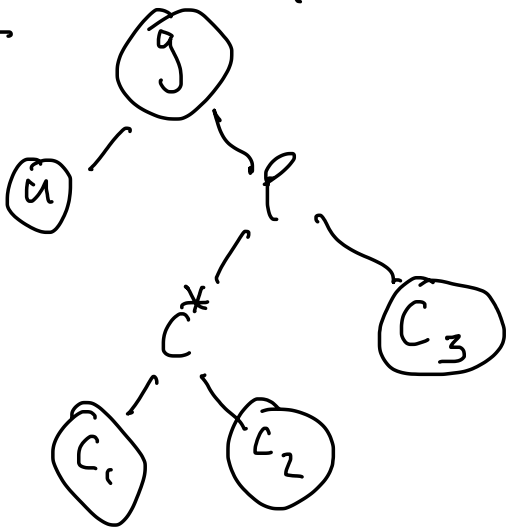
→ rotate as necessary
make root black (circled)

Case 1

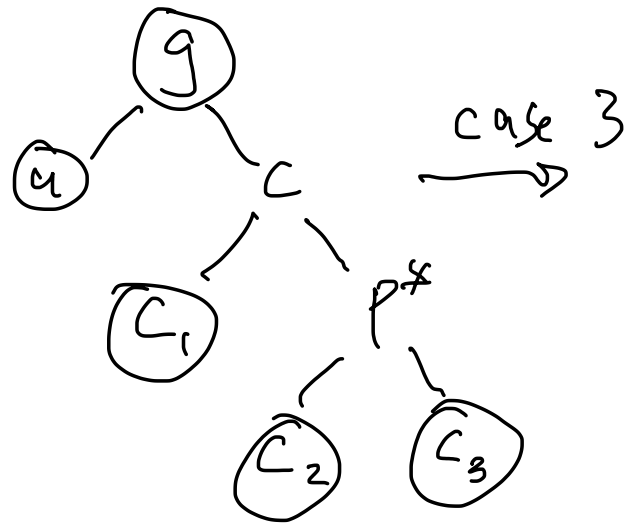


Case 2

parent red, uncle black, c inside



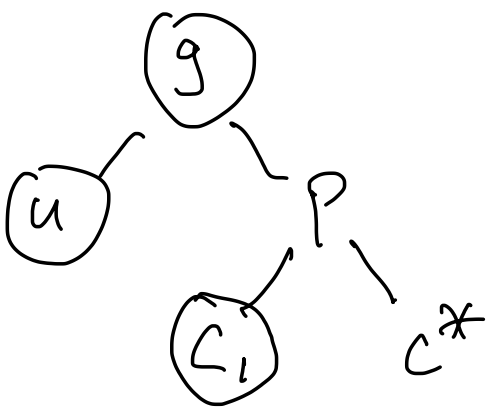
rotate c
u p, p
down



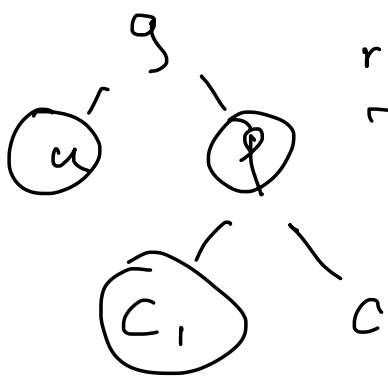
case 3

Case 3

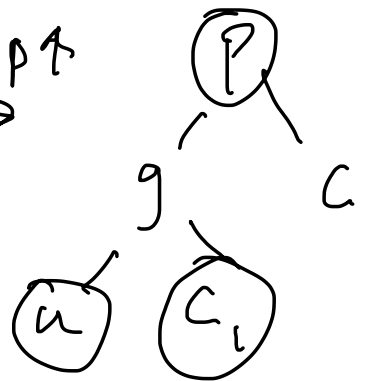
parent red, uncle black, c outside



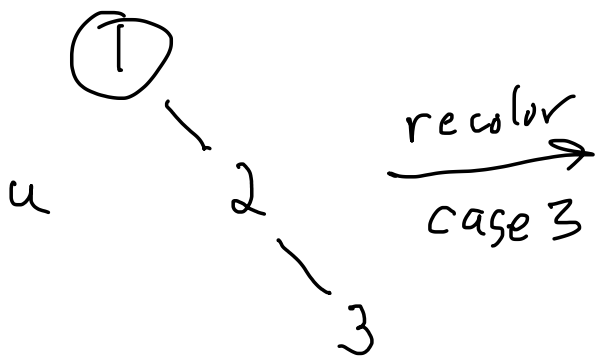
recolor



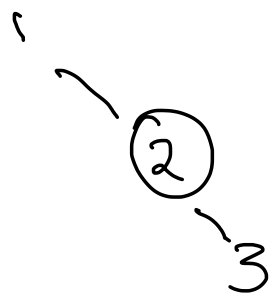
rotate p ↑
g ↓



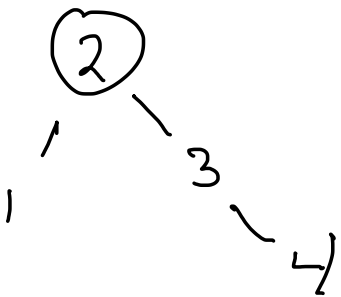
Example:
1, 2, 3, 4, 5, 6



recolor
case 3



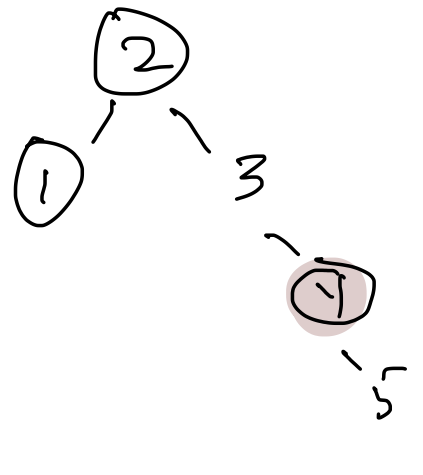
rotate



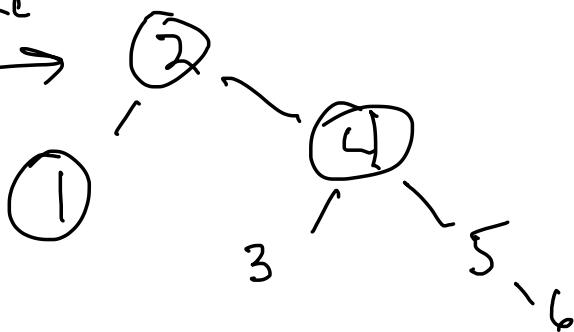
case 1
recolor



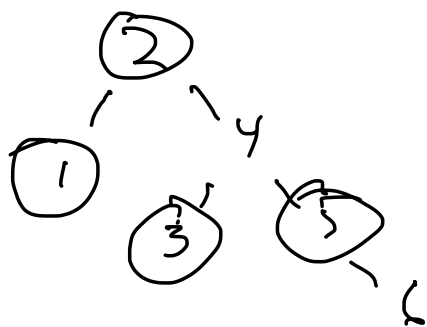
recolor
case 3



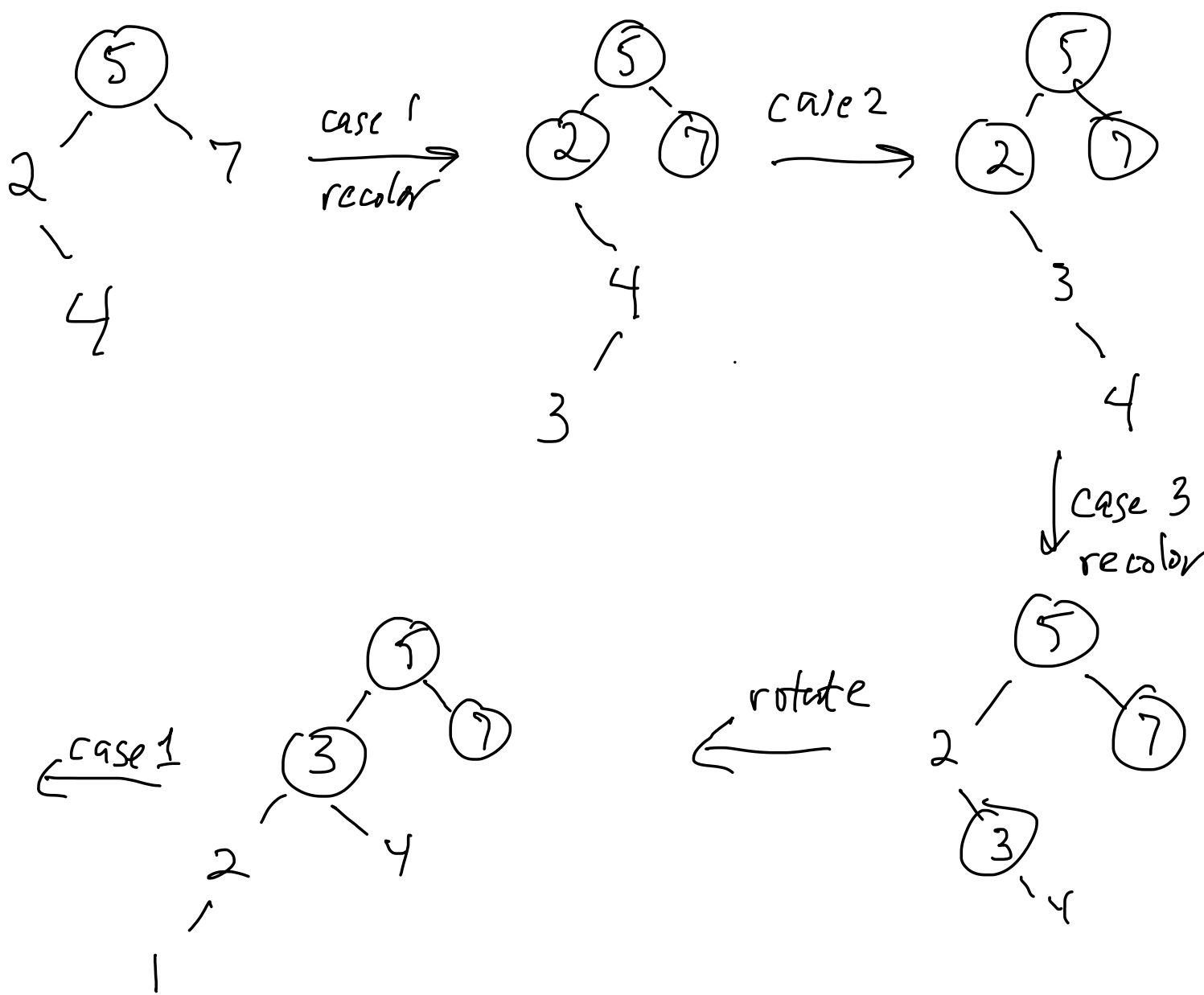
rotate



case 1
recolor



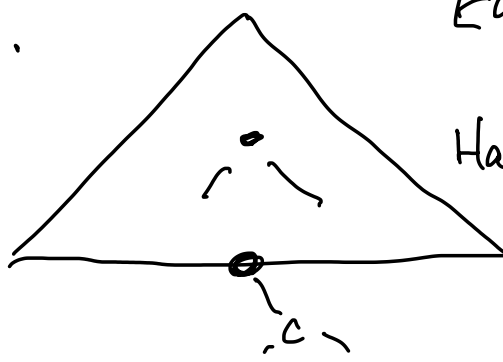
Example: 5 2 7 4 3 1



Properties of binary trees

- 1) Expected worst path is $O(\log n)$ deep.
If random data.
- 2) Can have $O(n)$ worst path length.
- 3) Can use a self-balancing tree (red black, AVL)
- 4) Insertion: $O(\log n)$
- 5) Deletion: ugly.

delete node d



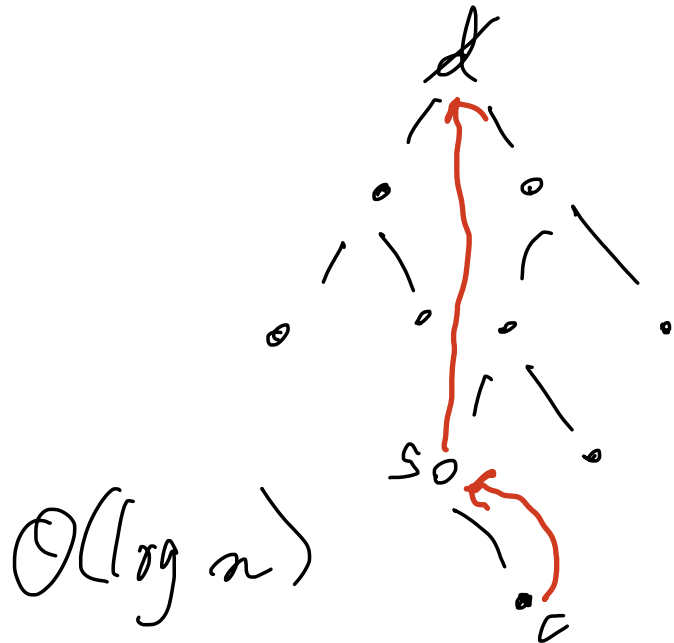
Easy: d is a leaf
remove it.

Harder: d has 1 child c:
move c into d's
place

Hardest: d has 2 children.

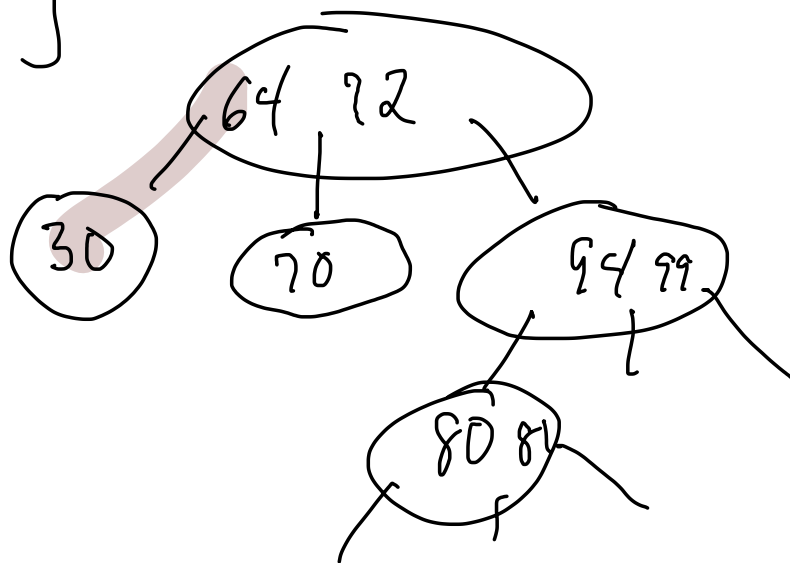
find d 's successor s : $s = \textcircled{RL^*}$

replace d with s ,
reparent d 's right child
(c) in place of s .



$\Theta(\log n)$

Ternary trees



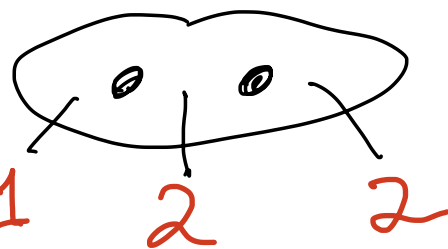
Complexity for searching

lucky: balanced. Depth = $\Theta(\log n)$

$$\log_3 n$$

work at each level

comparisons



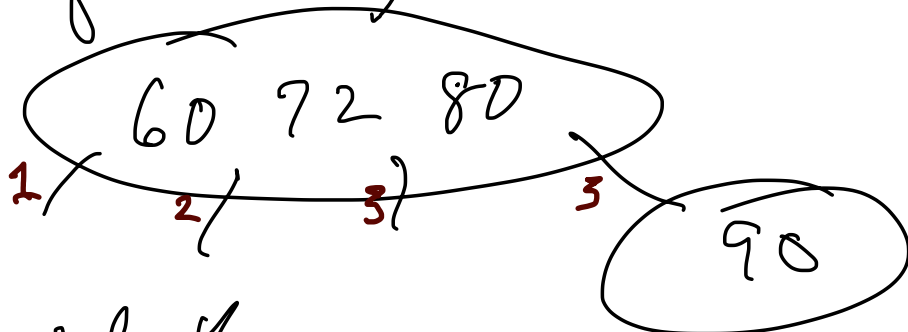
avg # comparisons = $\frac{5}{3}$

total work for searching $\frac{5}{3} \log_3(n)$

compare to binary tree: $1 \log_2(n)$

ternary tree is 1.05 cost of binary tree
5% degradation.

what about quaternary tree?



$\log_4 n$: depth

$\frac{9}{4}$ work at each level : $\frac{9}{4} \log_4 n$

12.5% degradation in comparison with
binary.

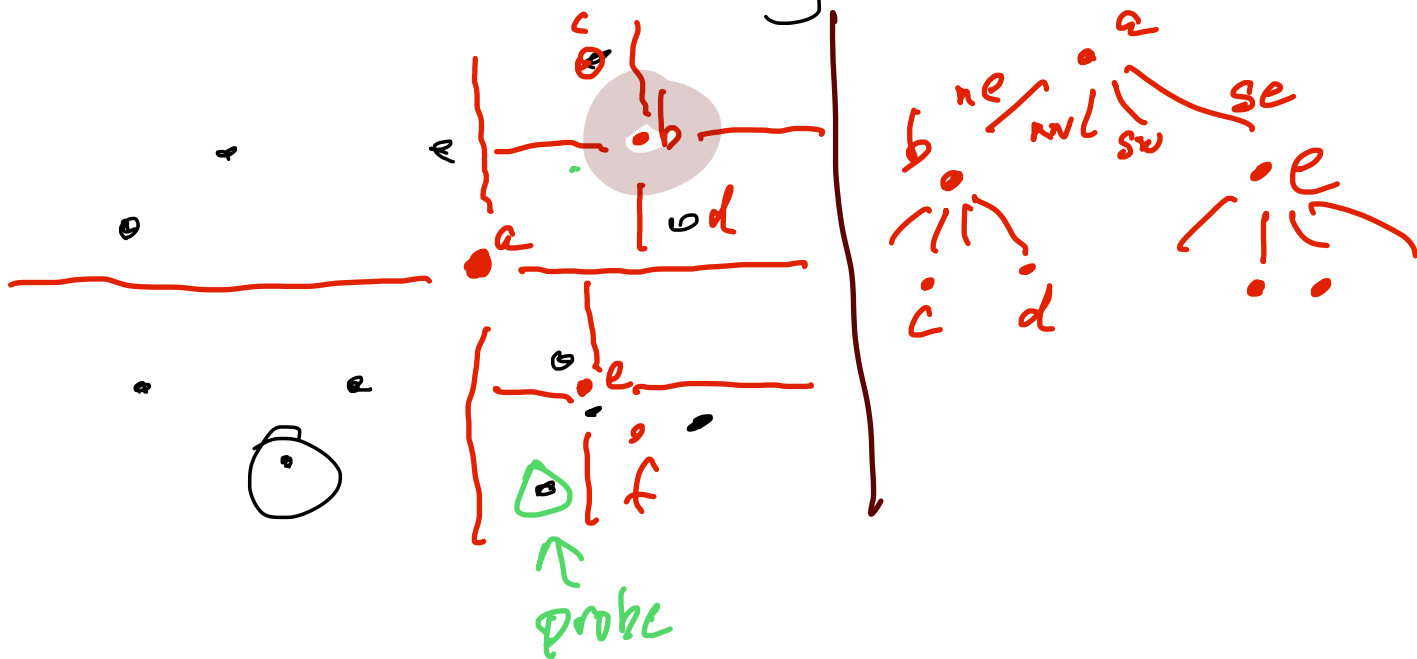
Other generalization?

Quad trees (Finkel, 1973)

2-dimensional data.

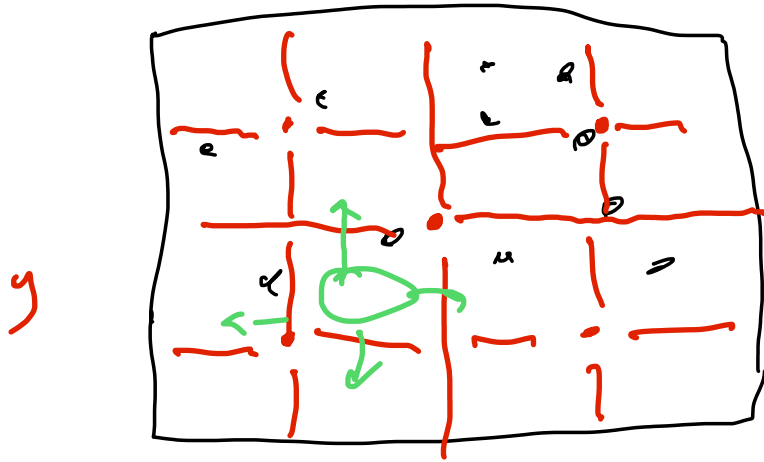
online : data appear 1 by 1.

vs offline : all data are already available.



options:

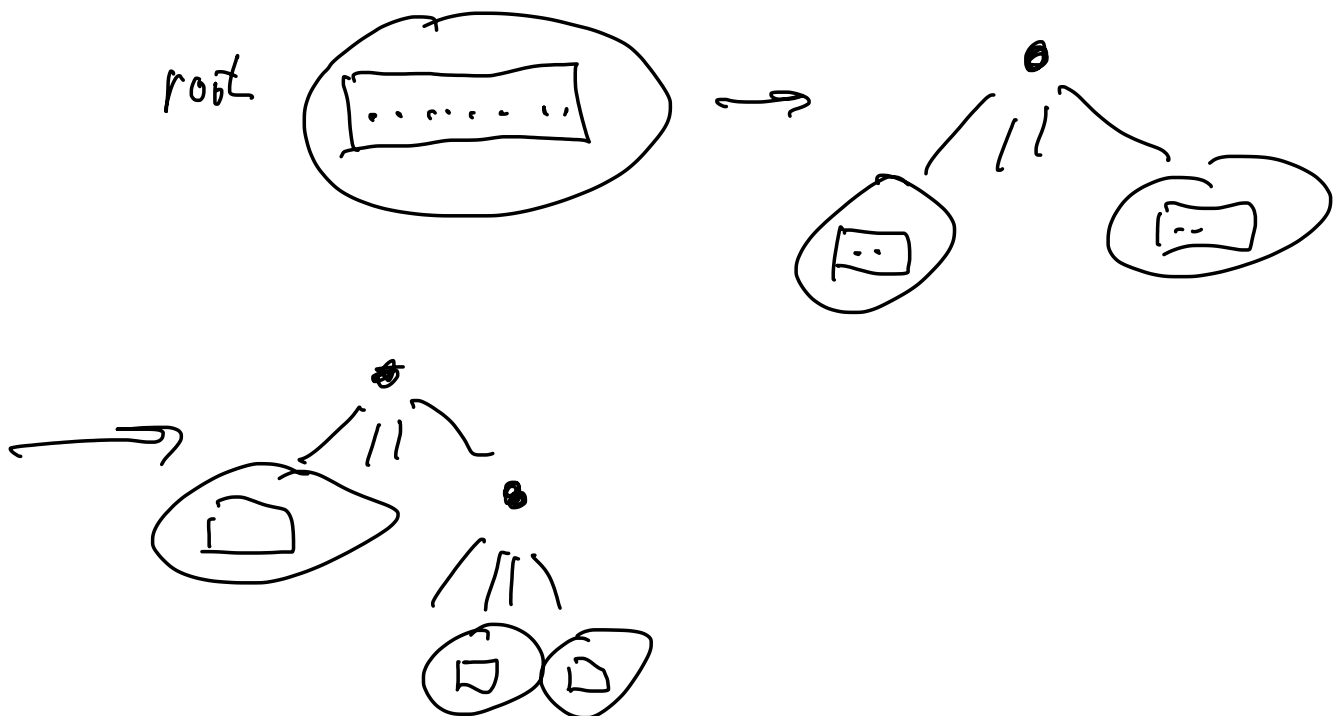
- ① only use points as tree nodes
- ② stop subdividing when region holds only a few points
about 10 is a good "bucket" size.
- ③ place tree nodes in center of their region.



- ④ place tree node s at "median" location
trying to balance the # of points in
each quadrant.

Online: put all first points in to a bucket.

when it fills, split it and build a root of tree.



good for

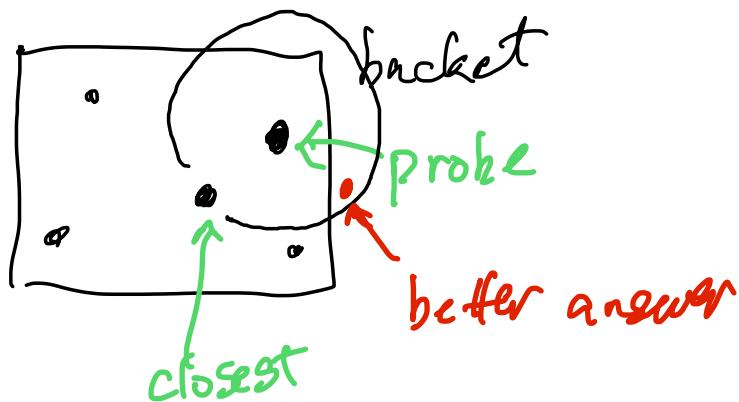
- ① searching
- ② (Jon Bentley) nearest-neighbor search.

probe: traverse down tree to reach bucket that the probe would be in.

Those bucket points are close to probe.

2) Consider them all, find closest among them.

3) go to neighboring regions.



Extend quad trees to higher dimensions?

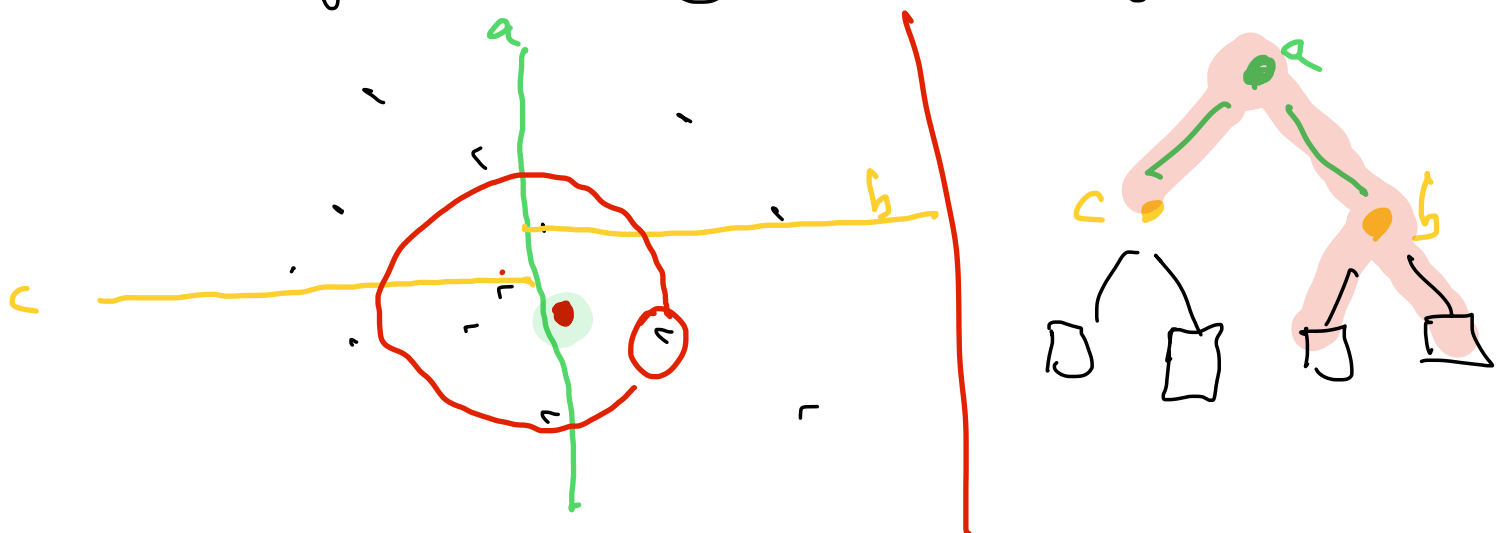
yes: 3d - Oct Trees



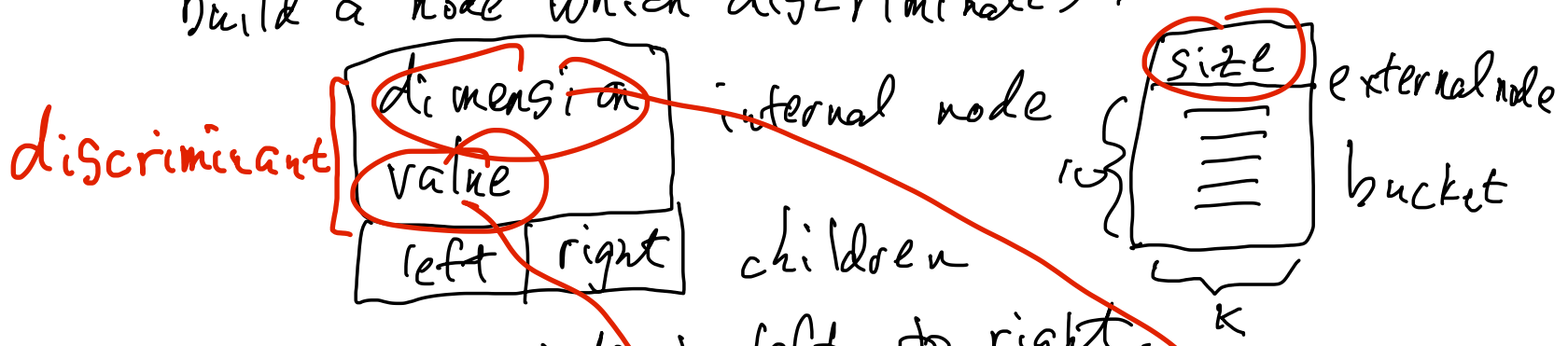
each internal node has 8 children.

people use oct trees to represent 3-d images and objects.

k-d trees (Bentley): especially good for high dimensional



Rule: given a cloud of points:
 find the dimension with greatest range.
 in that dimension, find median value.
 build a node which discriminates:



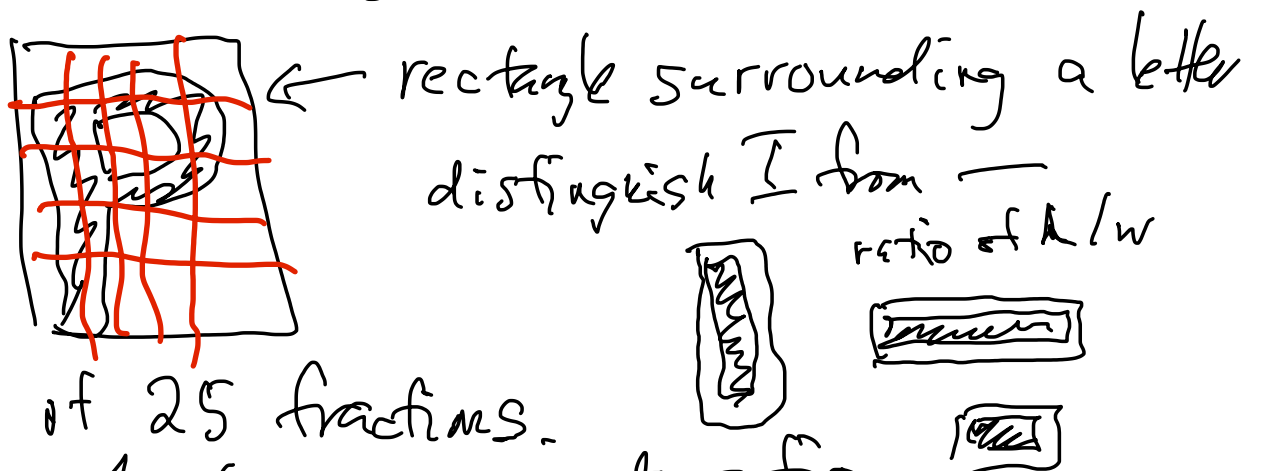
recurse on points to left, to right.

those whose value in the dimension are less than or equal to the discriminating value.

Result: balanced tree,
 leaf nodes are buckets (with $b \approx 10$ at most points)

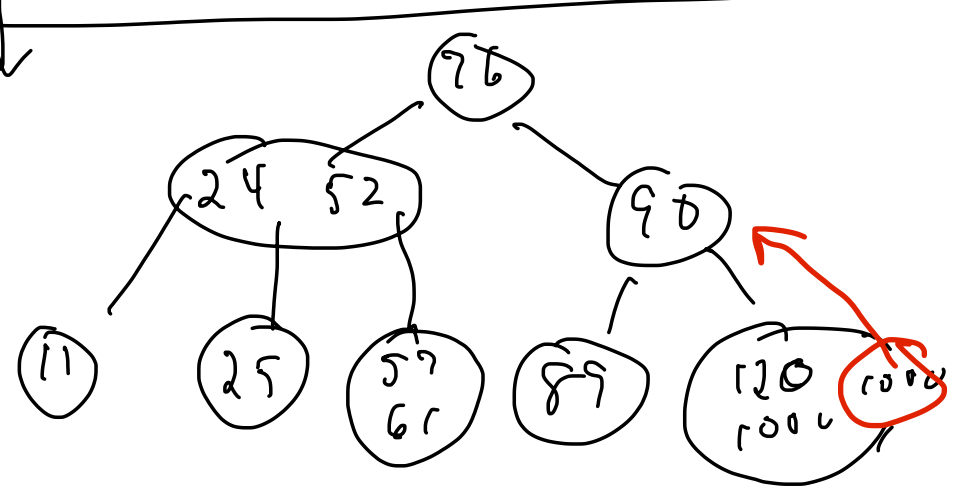
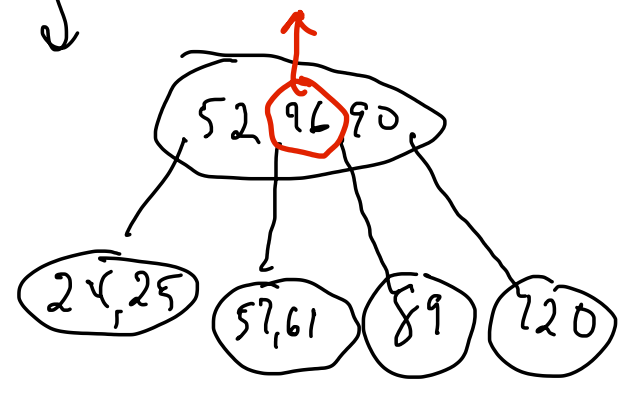
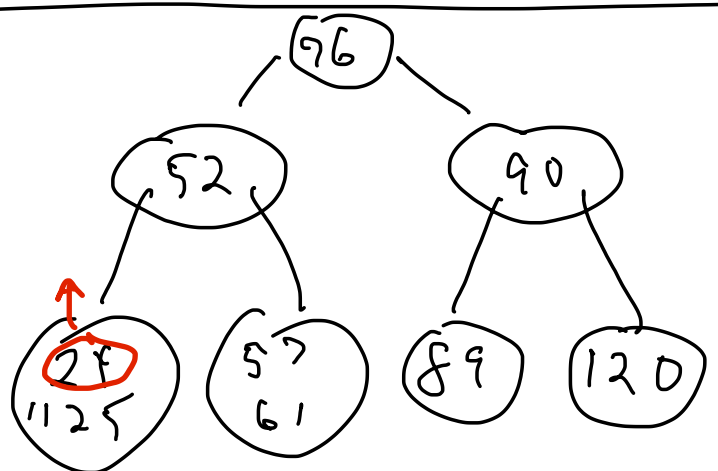
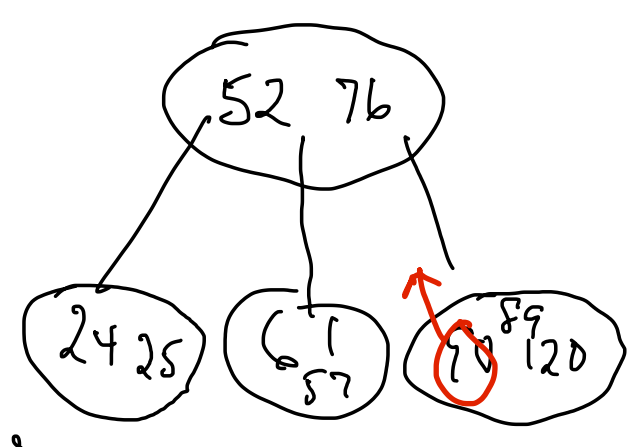
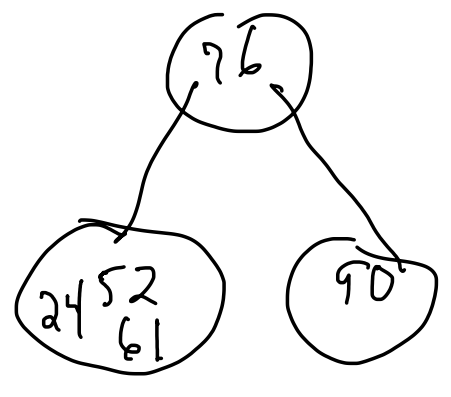
Good for:
 representing clouds of high-dimensional data
 finding nearest neighbors.

Optical character recognition:



vector of 25 fractions.
 + 1 element: aspect ratio
 \Rightarrow point in 26-dimensional space

2-3 trees :



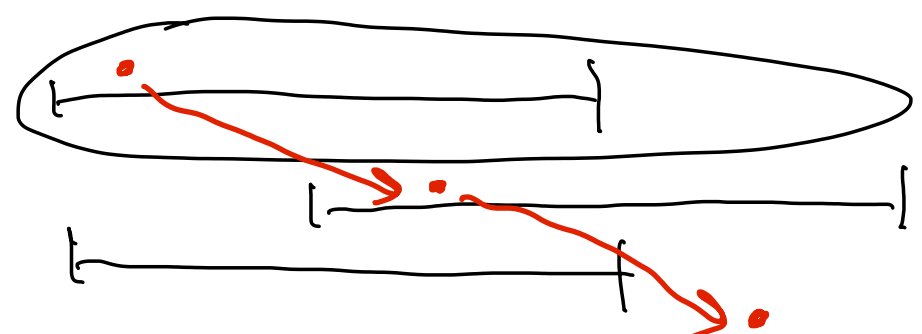
duplicates are tricky:
can we always break ties
to the left?

⇒ no duplicates
allowed.

⇒ or recursive search
must look 2 ways after tie

Insertion, search: $\Theta(\log n)$

Stooge sort



Analysis:

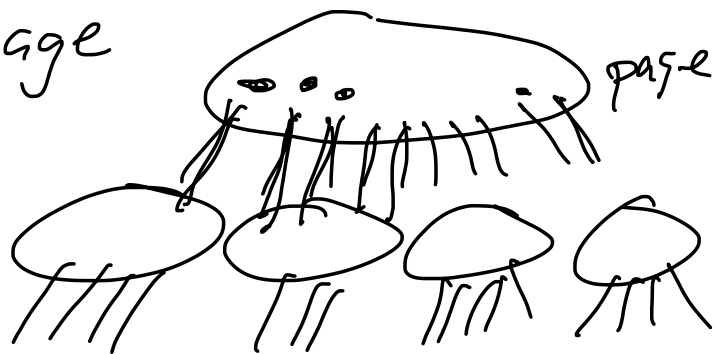
$$C_n = 1 + 3k + \frac{2}{3}n$$

$a = 3$
 $b = 3/2$
 $k = 0$ (because $1 = n^0$)

$$\begin{aligned}
 & \left. \begin{array}{l} a \\ 3 > 1 \end{array} \right\} \Theta(n^{\log_b a}) \\
 & \left. \begin{array}{l} b \\ \Theta(n^{\log_{3/2} 3}) \end{array} \right\} \\
 & \approx \Theta(n^{2.71})
 \end{aligned}$$

B-trees: extension of 2-3 trees.

Bucket: fits in a disk page



block: 4KB

element $\left\{ \begin{array}{l} \text{value: 4B} \\ \text{index to a disk block: 4B} \end{array} \right.$

$$4\text{KB} / 8\text{B} = 512 \text{ elements / block}$$

$$m = 512 \text{ (m "capacity" of a node)}$$

$$m = 3 \Rightarrow 2\text{-}3 \text{ tree}$$

Use $g = \lceil m/2 \rceil$ as the "fill" amount for new nodes.

Capacity of a node:

$m-1$ values, m pointers

$$g = \lceil \frac{m}{2} \rceil \text{ (half size)}$$

internal nodes: $g \dots m$ children.

exception: root: $1 \dots m$ children.

insertion: 1) find the right leaf node

2) insert value in node

3) if node is overfull (m values)

split it: Hoist middle value to parent,

remaining values: g of them as a new node (left) others as a new

node (right). Adjust pointers in parent

to point to the resulting nodes.

4) If parent is overfull, split it,

and if necessary, continue up the tree to root (which itself can

split, generating a new root).

Height of tree (full): $\log_m(m)$

$\Theta(\log m)$.

Deletion of a value v . Difficult.

1) Internal node: Replace v in that node with $\text{successor}(v)$. Then continue with deletion from leaf.

2) Leaf:

common \rightarrow good case: leaf is still g or more values.

rare \rightarrow bad case: leaf node is under-full ($< g$ values)

steal a value from a neighbor.

(B \neq tree: nodes also point to neighbors)

if all neighbors are almost under-full:
merge with one neighbor, taking a value
from parent, ...

Summary of trees

in
use →

Binary (sorted, insertion $O(\log n)$, deletion hard)

traversals (inorder, postorder, ...)

Higher arity (ternary, ...)

not practical: cost of dealing with nodes
overwhelms benefit of shallower tree

in
use →

Self-balancing (Red-black, AVL)

guarantee of $O(\log n)$

insertion is more complicated.

in
use →

Higher dimension
quad, oct
k-d

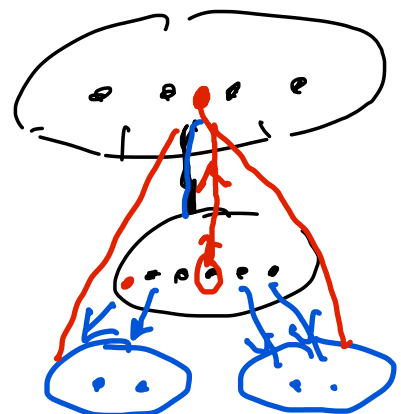
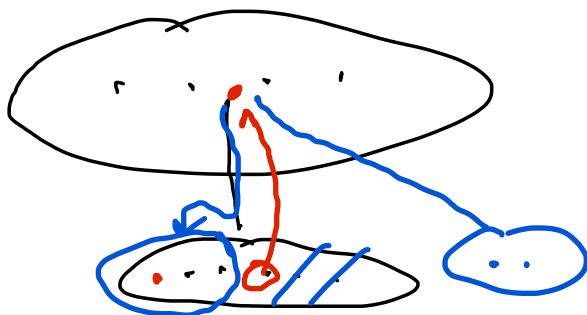
Other organizations to guarantee balance

2-3 tree

B-tree

Moral:

trees are essential
recursion is essential.

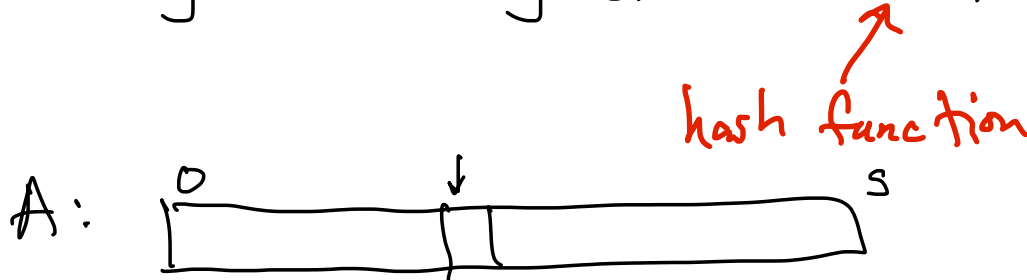


Searching: $\Theta(n)$ array
 $\Theta(\log n)$ sorted array
 tree
 $\Theta(1)$

Hashing $\left\{ \begin{array}{l} \text{expected, not guaranteed} \\ n \text{ elements} \Rightarrow \text{key is at least } (\log_2 n) \text{ bits.} \end{array} \right.$

Basic idea: use an array $A[]$

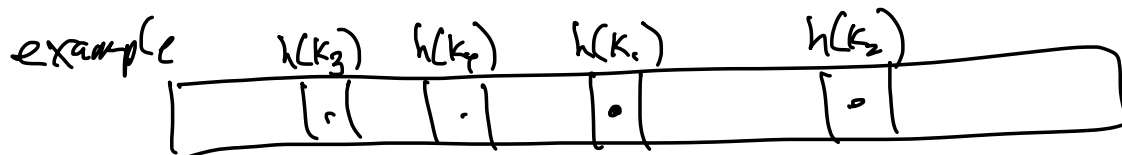
place key k in array at index $h(k)$



Keys are usually strings

$h(\text{string}) = \text{small integer } (0 \dots s-1)$

h should be fast.



$k_1 = \text{once}$

$k_3 = a$

$k_2 = \text{upon}$

$k_4 = \text{midnight}$

problems: collisions. Two keys that hash to the same index.

Collisions are unavoidable.

Birthday paradox

$\text{hash}(\text{person}) = \text{birthdate}$



Prob (no collisions with j students)

$$365 \cdot 364 \cdot 363 \cdot \dots (365 - j + 1)$$

$$\frac{\quad}{365^j}$$

$$= \frac{365!}{(365-j)!} \cdot \frac{1}{365^j}$$

$$j = 23 \Rightarrow \text{Prob} < 1/2$$

$$j = 50 \Rightarrow \text{Prob} = 0.029$$

Dealing with collisions.

open addressing: if there is a collision, use another location in the array.

disadvantage: clustering

disadvantage: deletion

(perfect hashing: if know all keys in advance, design $h()$ to avoid collisions.)

1) Linear probing, Multiple tries for key k .

$$p_0 = h(k)$$

$$p_1 = (h(k) + 1) \bmod S.$$

$$p_j = (h(k) + j) \bmod S.$$

insert: at first empty cell starting at $h(k)$.

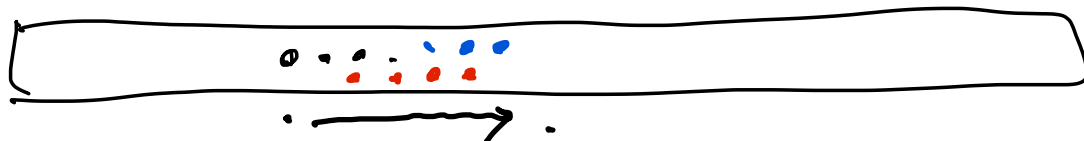
search: p_0 (check!), p_1 (check!) ...

stop when

1) find key! success

2) find empty slot. failure

3) back to p_0 . failure (full)



2) Family of hash functions.

$$p_0 = h_0(k) \quad \dots \quad p_j = h_j(k)$$

$$p_i = h_i(k)$$

hope: avoid clusters merging.

worse and worse as table fills.