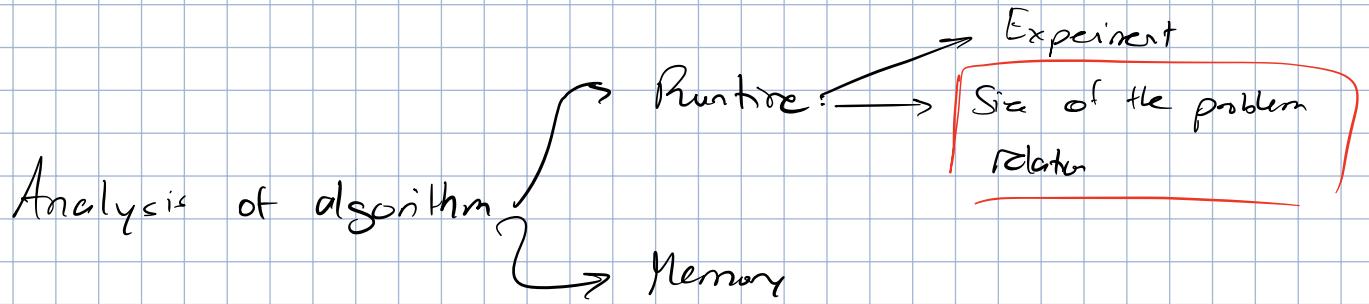


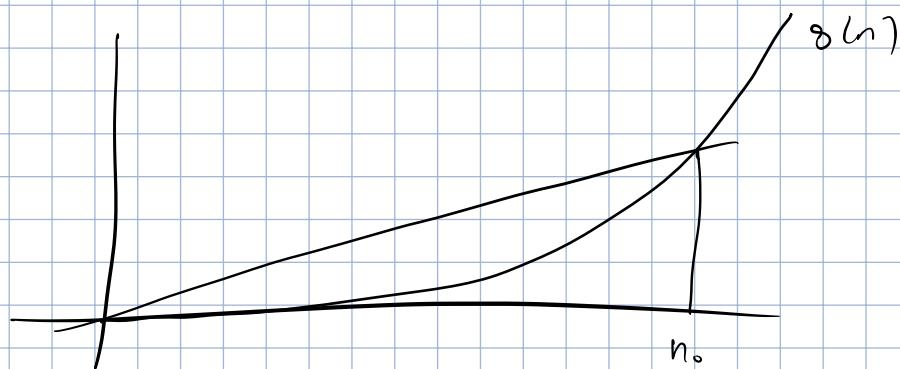
WEEK 1



Asymptotic: growth rate of algorithm.

① Big O: upper bound on g.r.

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



Prove:

1. Use inequality to derive c and n_0 .

$$2n^2 + 3n + 11 \in O(n^2)$$

$$\begin{aligned} 2n^2 + 3n + 11 &\leq 2n^2 + 3n^2 + 11n^2 \\ &\leq 16n^2 \quad \forall n \geq 1 \end{aligned}$$

Properties:

1. $a f(n) \in O(f(n))$ Also applies to Σ

2. $f(n) \in O(h(n)) \wedge h(n) \in O(g(n)) \Rightarrow f(n) \in O(g(n))$

3. $\max(f(n), g(n)) \in O(f(n) + g(n))$

4. $a_0 + a_1 n + \dots + a_k n^k \in O(n^k)$

5. $n^x \in O(a^n) \quad x > 0, a > 1$
Polynomial Exponential

$$6. (\log n)^2 \in O(n^2) \quad \text{as } \log n \leq n$$

$$7. O(f(n) + g(n)) \in O(\max\{f(n), g(n)\}) \rightarrow \text{Also appln to } \Omega$$

② Big Omega: lower bound on g.r.

$$\exists c > 0, n_0 > 0. \forall n > n_0 \Rightarrow |f(n)| \geq c |g(n)|, \text{ then } f(n) \in \Omega(g(n))$$

③ Big Theta: tight bound on g.r.

$$1. \text{ If } f(n) \in O(g(n)) \wedge f(n) \in \Omega(g(n)) \Rightarrow f(n) \in \Theta(g(n))$$

$$2. \exists c_1, c_2 > 0, n_0 > 0. \forall n > n_0 \Rightarrow c_1 |g(n)| \leq |f(n)| \leq c_2 |g(n)|$$

Tricks:

Key property: $\log n \leq n \quad \forall n \geq 2$

Fractional:

$$\frac{1}{2}n^2 \cdot S_n \in \Omega(n^2)$$

1. Smaller constant than leading term.

$$\frac{1}{2}n^2 \cdot S_n \geq \frac{1}{4}n^2$$

$$\frac{1}{4}n^2 + \frac{1}{4}n^2 - S_n \geq \frac{1}{4}n^2$$

$$\frac{1}{4}n^2 - S_n \geq 0$$

$$n^2 - 20n \leq 0$$

$$n(n-20) \leq 0$$

$$\therefore n \geq 20 \checkmark$$

$$\boxed{c = \frac{1}{4}, n_0 = 20}$$

④ Little o: $g(n)$ dominate $f(n)$

Smaller than Big O

$$\underline{\underline{\forall c > 0, \exists n_0 > 0. \forall n > n_0 \Rightarrow |f(n)| < c |g(n)|}}$$

Trick: set arbitrary $c \Rightarrow$ find n_0 to satisfy inequality
Use intermediate

$$2020n^2 + 1388n \in O(n^3)$$

Let $c > 0$. $h(n)$ \Rightarrow int. func.

Similar to more complicated func. + include c

$$2020n^2 + 1388n \leq h(n) \leq cn^3$$

$$2020n^2 + 1388n \leq \underbrace{5000n^2}_{\text{obv.}} \leq cn^3$$

obv.

Non-obv.

$$5000n^2 \leq n^3$$

$$5000n^2 \ll n^3$$

$$\frac{5000}{c} \leq n$$

$$n \geq \frac{5000}{c}$$

n_0 !

⑤ Little Omega: $f(n)$ dominate $g(n)$

Exact same defn., switch inequality as O

Limit technique:

Not sufficient for proof

$$L = \lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} \Rightarrow f(n) \in \begin{cases} o(n) & \Rightarrow L = 0 \\ \Theta(n) & \Rightarrow 0 < L < \infty \\ \omega(n) & \Rightarrow L = \infty \end{cases}$$

↳ Use L'Hôpital's

Doublings:

a) Constant:

$$T(n) = c \Rightarrow T(2n) = c$$

b) Logarithmic:

$$\begin{aligned} T(n) = c \log n &\Rightarrow T(2n) = c \log(2n) \\ &= c \log 2 + c \log n \\ &= T(n) + c \end{aligned}$$

c) Linear:

$$T(n) = cn \Rightarrow T(2n) = 2T(n)$$

d) Exponential:

$$\begin{aligned}
 T(n) = c \log n &\Rightarrow T(2n) = c(2n) \log(2n) \\
 &= 2cn \log 2 + 2cn \log n \xleftarrow{\text{Distributive}} 2T(n) \\
 &= 2T(n) + cn
 \end{aligned}$$

e) Quadratic:

$$T(n) = cn^2 \Rightarrow T(2n) = 4T(n)$$

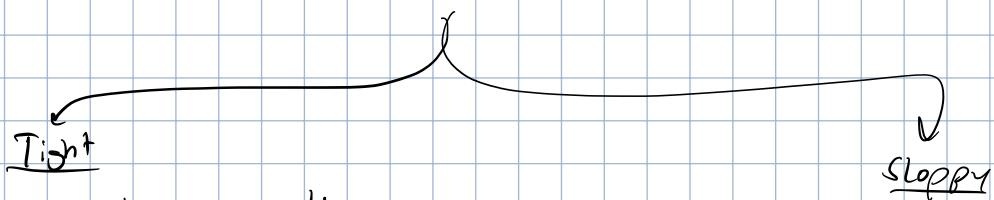
WEEK 2

How to create $T(n)$?

Strategy: count # of atomic ops $\Rightarrow \Theta(1)$

\hookrightarrow # of iterations = runtime.

Calculating # of iterations



Construct exact summation
+ evaluate.

Often much more difficult.

*Exclusive
+ sum func.*

Prove $O(1 \dots) \approx \sum (\dots)$

How? Mess around w/ sums +
limits of sums

O: sum our maximum value

$$\sum_{i=0}^n i \leq \sum_{i=0}^n n = n^2 \quad \hookrightarrow \text{Picks very few terms.}$$

Sum our as many iterations as possible.

S2: sum our minimum value

$$\sum_{i=0}^n i \geq \sum_{i=0}^n \frac{n}{2} = \frac{n^2}{2}$$

Sum our as little iterations as possible.

$$\sum_{i=1}^n \sum_{j=1}^i \sum_{k=i}^j 1 \geq \sum_{i=1}^{n/3} \sum_{j=1}^{n/3} \sum_{k=1}^{n/3} = \frac{n^3}{3}$$

While loop?



while cond:

... - - -

① What variable det. cond.

② start value + op. on value.

③ Let $k = \#$ of iterations until cond. breaks \Rightarrow solve for $k \Rightarrow$ nk in sums.

- sum $\leftarrow 0$
 for $i \leftarrow 1$ to n do
 $j \leftarrow i$
 while $j \geq 1$ do
 sum \leftarrow sum $\cdot \frac{1}{j}$
 $j \leftarrow \lfloor \frac{j}{2} \rfloor$

$$T(n) = \sum_{i=1}^n \# \text{ of while loops}$$

$$\textcircled{1} \quad j \leftarrow i \rightarrow 1$$

② Let $k = \# \text{ of while loops}$.

$$i \cdot \left(\frac{1}{2}\right)^k = 1 \quad \# \text{ of loops} \leq \lfloor k \rfloor + 1$$

$$\log i - k = 1$$

$$k = \log i - 1$$

$$T(n) = \sum_{i=1}^n k = \sum_{i=1}^n (\lfloor \log i \rfloor + 1)$$

0:

$$T(n) \leq \sum_{i=1}^n (\log n + 1) = n \log n + n$$

Σ:

$$\sum_{i=1}^{\lfloor \frac{n}{2} \rfloor} (\lfloor \log n \rfloor + 1) = \frac{n}{2} \log n + \frac{n}{2} \in \Sigma(n \log n + n)$$

* 2 instances of same size not guaranteed to have same $T(n)$ *

↳ Sorting.

↳ Consider: worst instance, increasing instance

Compare 2 diff algs? Use Θ , not O , use worst-case

Recurrences:

$$T(n) = 2T\left(\frac{n}{2}\right) + cn$$

①

$$= 2 \left(2T\left(\frac{n}{4}\right) + cn \right) + cn$$

$$= 4T\left(\frac{n}{4}\right) + cn + cn$$

$$= 4 \left(T\left(\frac{n}{8}\right) + cn \right) + cn + cn$$

$$n \times \left(\frac{1}{2}\right)^k = 1$$

$$\log n - k = 1$$

$$k = \log n - 1$$

$$= nT(1) + \left[2 \cdot \frac{n}{2} + 4 \cdot \frac{n}{4} + 8 \cdot \frac{n}{8} + \dots + \frac{n}{2} \cdot \frac{n}{2} \right]$$

$$= nT(1) + n \log n \in \Theta(n \log n)$$

Exped
+ find
pattern

Priority Queue ADT:

Insert item w/ priorities \Rightarrow access + delete highest pri item.

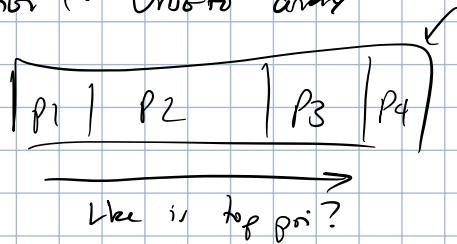
Operations:

- o ins +

- o delete Max

Note: minimum-oriented PQ

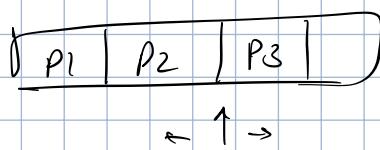
Realization 1: Unsorted array



Insert: $O(1)$

Deletion: $O(n)$

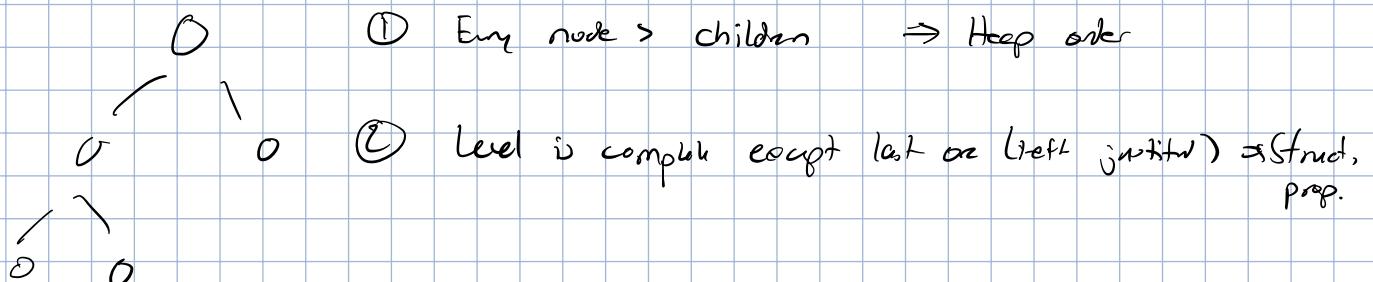
Realization 2: sorted array



Insert: $O(n)$

Deletion: $O(1)$

Realization 3: binary heap



Height: $O(\log n)$

Store priorities in binary heap?

Code: Array:

Left child of node $[i]$: node $[2i + 1]$

Right child: $[i]$: node $[2i + 2]$

Root: $A[1]$

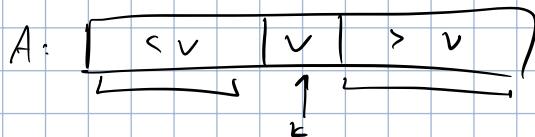
Last node: $A[n-1]$

Parent: $\left\lfloor \frac{i-1}{2} \right\rfloor$

MODULE 3: Sorting + Randomized Algos

Quick Select

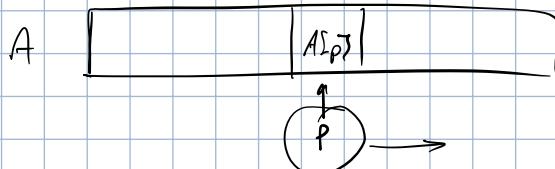
Given value v , find index of v if A sorts



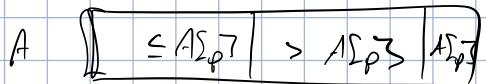
Given index k , what value must be in k if sorts?

Quick select:

1. Choosing pivot



2. Partition around pivot



How?

1. Choosing pivot: return the last index

2. Partition:

i) Simple

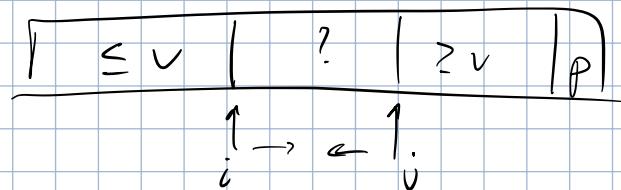
partition(A, p)
 small $\leftarrow \{ \}$
 equal $\leftarrow \{ \}$
 great $\leftarrow \{ \}$
 $v \leftarrow A[p]$
 loop through \rightarrow add elem in list
 append {small, equal, great}

Time: $\Theta(n)$

Space: $\Theta(n)$

\hookrightarrow Array with everything same!

ii) In-place: in-place

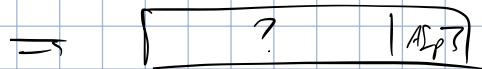


Swap values at $A[i]$ + $A[j]$ according to pivot $\Rightarrow i \geq j$

partition (A, p)

// Ensure pivot is last

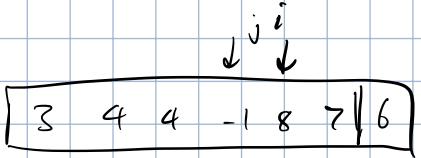
sweep ($A[n-1], A[p]$)



$i \leftarrow -1, j \leftarrow n-1$

$v \leftarrow A[n-1]$

loop



Efficient
for duplicates.

do $i \leftarrow i+1$ while $A[i] < v$

do $j \leftarrow j-1$ while $A[j] > v$ and $j \geq i$

// Make sure we didn't b/c of overlap

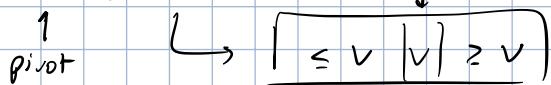
if $i \geq j$: break

sweep ($A[i], A[j]$)

end loop

// Everything below i is less than pivot \Rightarrow

sweep ($A[n-1], A[i]$)



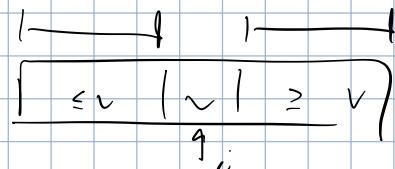
return i

Put all this together:

quick-select | (A, k)

$p \leftarrow \text{choose pivot}(A)$

$i \leftarrow \text{partition}(A, p) \rightarrow$



if $i = k$:

return $A[i]$

else if $k < i$:

return quick-select ($A[0 \dots i-1], k$)

else:

return quick-select ($A[i+1 \dots n-1], k-i+1$)

→ O. w. \Rightarrow overintend
length of \downarrow \downarrow \downarrow
cv subarr

Analysis: * Tricks *

1. Best-case analysis:

Think: what situation leads to fast outcome.

Don't recur: $i = k \Rightarrow \Theta(n)$

2. Worst-case analysis

Worst: situation goes poorly \Rightarrow maximize loops, recursions ...

$$T(n) = \begin{cases} cn + T(n-1) & \Rightarrow i = 1, \text{ or } 0 \quad (n \geq 1) \\ c & n=0 \end{cases}$$

① End out generalized expression

$$\begin{aligned} T(n) &= cn + T(n-1) \\ &= cn + c(n-1) + T(n-2) \\ &= cn + c(n-1) + c(n-2) + T(n-3) \\ &= \vdots \\ &= cn + c(n-1) + \dots + c(n-i) + T(n-i-1) \end{aligned}$$

② Figure out when generalized expression = best case

$$n-i-1 = 1$$

$$i = n-2$$

③ Substitute + write summation

$$T(n) = cn + \dots + 2c + T(1)$$

$$= c \sum_{i=2}^n i + c \in \Theta(n^2)$$

3. Average case analysis

$$T^{\text{avg}}(n) = \frac{1}{\# \text{ of instances of size } n} \sum_{I: \text{size}(I)=n} T(I)$$

① See if you can make assumptions

1. All elements are distinct
2. Sort from 0 - n-1

② Find out # of number of instanc.

$$\boxed{n \mid n-1 \mid n-2 \mid \dots \mid 1} \Rightarrow n! \# \text{ of instances of size } n$$

↳ NO DUPLICATES!

② Break out sum according to decision we've made.

3 decisions:

1. Pivot Choose pivot to sum

2. Partition \Rightarrow arr, b/c easier to do

3. Recurr.

$$\sum_{i=0}^{n-1} \sum_{\substack{I: \text{size}(I)=n \\ \text{pivot} = i}} T(I)$$

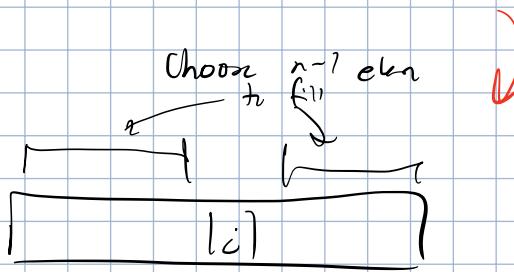
All possible partition

③ Create upper bound:

$$\sum_{\substack{I: \text{size}(I)=n \\ \text{pivot} = i}} T(I) \leq (n-1)! \cdot cn + (n-1)! \cdot \max\{T(i), T(n-i+1)\}$$

↑ Partition ↑ Recurr.

$(n-1)! : \# \text{ of possible arrays w/ pivot } i$



④ Simplify sum

$$T^{\text{ans}}(n) \leq \frac{1}{n!} \sum_{i=0}^{n-1} (n-1)! \cdot cn + (n-1)! \cdot \max\{ \dots \}$$

$$\leq \frac{(n-1)! \cdot cn}{n!} + \frac{1}{n!} \sum_{i=0}^{n-1} (n-1)! \cdot \max\{ \dots \}$$

$$\leq cn + \frac{1}{n} \sum_{i=0}^{n-1} \max\{ T(i), T(n-i+1) \}$$

$\therefore T^{\text{ans}}(n) \in O(n)$ Bound by 4 (A Induction)

Randomized Algorithms.

Involves RNG to make decision

Randomic thing \Rightarrow input doesn't matter \Rightarrow worst-case = avg. case

How to do randomized algo analysis (recursion)

$$T^{\text{exp}}(I) = E[T(I, R)] = \sum_{\text{all possible } R} T(I, R) \cdot P(R)$$

1. Figure out recurrence:

$$T(n) = \begin{cases} ① & n \geq \dots \\ ② & n \leq \dots \end{cases}$$

2. Figure out the probability of either of recurrences

$$P(①) = \dots \quad P(②) = \dots \quad \Rightarrow \text{With RNG: this is uniform probability.}$$

3. Expectation:

$$P(\text{1 elem}) = \frac{1}{n}$$

$$T^{\text{exp}}(n) = T(①) P(①) + T(②) P(②) + \dots$$

Solve recurrence like ab.

Aberrant worst case: unlucky w/ RNG \Rightarrow exact same thought as worst-case analysis

Quick select \Rightarrow randomized pivot selector

choose pivot():

return random $\in [0, n-1]$ \rightarrow Some random index in $\{0, \dots, n-1\}$

$$\begin{aligned} T^{\text{exp}}(n) &= \sum_{i=0}^{n-1} T^{\text{exp}}(n-i) \cdot P(i \text{ is pivot}) \quad \frac{1}{n} = \text{uniform} \\ &\leq \frac{1}{n} \sum_{i=0}^{n-1} cn + \max \{ T^{\text{exp}}(i), T^{\text{exp}}(n-i-1) \} \\ &\leq cn + \frac{1}{n} \sum_{i=0}^{n-1} \dots \end{aligned}$$

Quick Sort

quicksort(A):

if $n \leq 1$: return

$p \leftarrow \text{choose_pivot}(A)$
 $i \leftarrow \text{partition}(A, p)$

Quick-sort ($A[0:i-1]$)

Quick-sort ($A[i+1:n]$)

Time Analysis:

① Best case: little work in recursion \Rightarrow recursions must be eq.

$$T(n) = 2T\left(\frac{n}{2}\right) + cn \in \Theta(n \log n)$$

② Worst-case: Maximize recursion

$$T(n) = T(n-1) + cn \in \Theta(n^2)$$

③ Avg. case:

$$T(n) = \frac{1}{n} \sum_{i=0}^{n-1} (T(n) + T(i) + T(n-i+1))$$

Doesn't matter, only counting comp., not work

① Simplify

$$\begin{aligned} T(n) &= \frac{1}{n} \sum_{i=0}^{n-1} n + \frac{1}{n} \sum_{i=0}^{n-1} T(i) + \frac{1}{n} \sum_{i=0}^{n-1} T(n-i+1) \\ &= n + \frac{2}{n} \sum_{i=0}^{n-1} T(i) \end{aligned}$$

Symmetric

② Remove all fractions

$$nT(n) = n^2 + 2 \sum_{i=0}^{n-1} T(i) \rightarrow T(0) + T(1) + \dots + T(n-2) + T(n-1)$$

③ Remove the sum by considering next smallest / largest term.

$$\text{Smaller: } (n-1)T(n-1) = (n-1)^2 + 2 \sum_{i=0}^{n-2} T(i)$$

$$nT(n-1) - T(n-1) = n^2 - 2n + 1 + 2(T(0) + T(1) + \dots + T(n-2))$$

Subtract this from $T(n)$

$$nT(n) - nT(n-1) + T(n-1) = 2n - 1 + 2T(n-1)$$

$$nT(n) = 2n - 1 + (n+1)T(n-1)$$

Divide
+ simplify
into pseudo
recurrence

$$\frac{T(n)}{n+1} = \frac{T(n-1)}{n} + \frac{2n-1}{(n+1)n}$$

④ Change of variables:

$$A(n) = \frac{T(n)}{n+1}$$

$$\therefore A(n) = A(n-1) + \frac{2n-1}{n(n+1)}$$

⑤ Solve recursion:

$$A(n) = A(n-1) + \frac{2n-1}{n(n+1)}$$

$$= A(n-2) + \frac{2(n-1)-1}{(n-1)n} + \frac{2n-1}{n(n+1)}$$

\vdots

$$= \sum_{i=1}^n \frac{2i-1}{i(i+1)}$$

$$= \sum_{i=1}^n \frac{2}{i+1} - \sum_{i=1}^n \frac{1}{i(i+1)}$$

\uparrow
 $\Theta(\log n)$ $\mathcal{O}(1)$

$$\therefore A(n) \in \Theta(\log n)$$

⑥ Replace Co.v:

$$\frac{T(n)}{n+1} \in \Theta(\log n)$$

$$T(n) \in \Theta(n \log n)$$

Improvement:

1. Randomized partition
2. Tail cell elimination to reduce recursion
3. Almost-sorth: limit on depth of recursion
 - ↳ Insertion sort for even better perf.
4. Handle duplicates via 3-way partition.

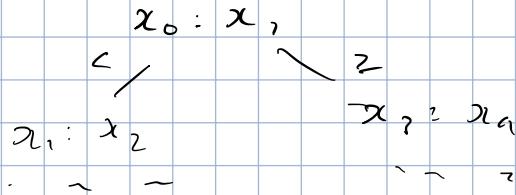
Lower Bound for Comparison-Based Sorting

Lower bound for comp-based: $\Omega(n \log n)$

↳ Sorting that $\forall i, j \ A[i] \leq A[j] \Rightarrow$

Proof: decision tree

Decision tree: shows all comparisons



Worst case of a sorting algo \Rightarrow max. number of comp.

↳ Equiv. to height of decision tree

Need to figure out height! Use # of leaves

↳ Leaf: represents unique output (e.g. unique sorting, unique selection)

1. Figure out # of unique outputs

2. $\text{decisions}^{\text{height}} \geq \text{leaves} \Rightarrow$ find lower bound on h .

Ex:// Assume n -distinct elements. Lower bound of sorting using comp-based algo

① Decisions / node

$\leq, \geq \Rightarrow 2$ decisions

② Leaves:

$[1, 1, 0] \Leftrightarrow [2, 2, 1]$ (same # of comp. \Rightarrow same leaf)

What matters \Rightarrow order of elems.

$\therefore n$ elems $\Rightarrow n!$ orderings.

③ Lower bound:

$\text{decisions}^{\text{height}} \geq \text{leaves}$

$2^h \geq n!$

$h \geq \log(n!)$

$h \in \Omega(n \log n)$

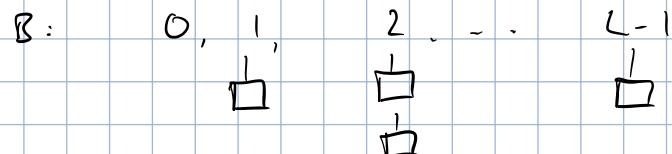
Non-comp. based algorithms

Almost always \Rightarrow assumptions made (e.g. range of elements, structural property, max size)

① Bucket sort

Assumption: array in range $[0, L-1]$

1. Create bucket array from $[0, \dots, L-1]$
2. Iterate through array + add to respective bucket (LL)



3. Iterate through bucket array + pop elements into array

Time: $\Theta(n)$, Space: $\Theta(L + n)$

\uparrow \uparrow
buckets nodes

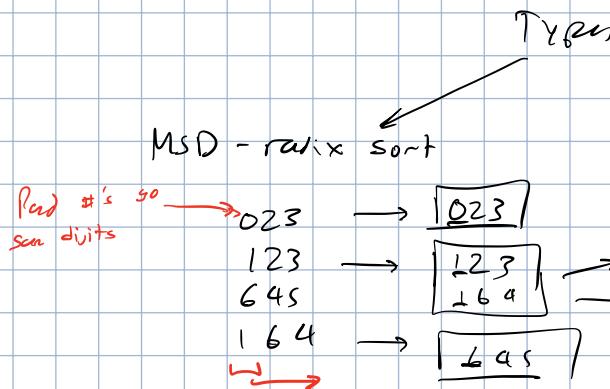
Great if $L \approx n$. If $L \gg n$, inefficient

② Digit-based sorting

Looks at digits one-by-one \Rightarrow bucket sort on digit \Rightarrow recursive

\hookrightarrow Stable sorting

\hookrightarrow Dependent on base representation. If base $R \Rightarrow R$ buckets



Precision-bound.

Runtime: $\Theta(n \log R)$

Space: $\Theta(n + m + R)$

LSD - radix sort

Bucket sort on multiple times \Rightarrow sort

$\Theta(m \log n + R)$

Why is this $\Theta(n)$?

If R is known + recr \Rightarrow m can be found.

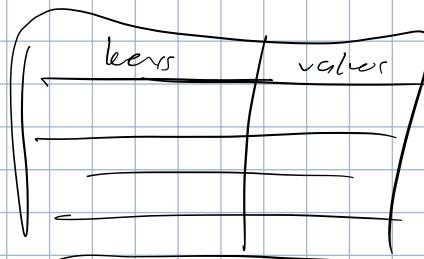
$\therefore [0, \dots, 10^4 - 1]$

$\therefore \boxed{m = 4} \boxed{R = 10}$

$\therefore \Theta(n)$

MODULE 5 - AVL Trees for Dicts

Dictionary ADT



search (k, v) \rightarrow value

insert (k, v)

delete (k)

Implementation:

Comparison + insert (k, v) is in constant time

1. Unordered array list:

• search: $\Theta(n)$, insert: $\Theta(1)$, deletion: $\Theta(n)$

2. Ordered array:

• Search: $\Theta(\log n)$, insert: $\Theta(n)$, deletion: $\Theta(n)$

3. AVL tree:

Binary search tree

Search: trivial

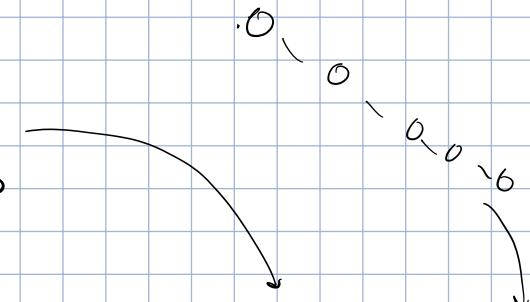
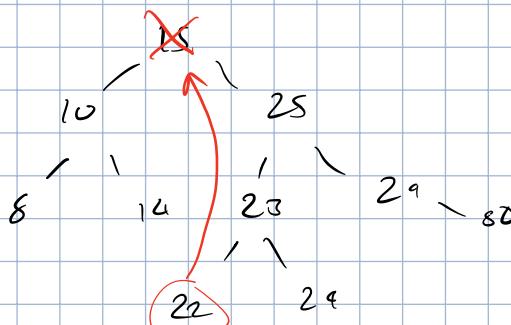
Insertion: search \rightarrow add into node

Deletion:

Case #1: deleting leaf \Rightarrow nothing but delete node

Case #2: deleting node w/ 1 child \Rightarrow move child to current node

Case #3: deleting w/ 2 children \Rightarrow swap w/ next node in in order trans.
(will be leaf) + delete



All operations in $\Theta(h)$. h depends on balance of tree. $\log n \leq h \leq n$

AVL Trees:

Self-balancing BST w/ height-balance prop:

Height-balance prop: \forall nodes in tree, $|\text{height}(\text{left}) - \text{height}(\text{right})| \leq 1$

$\hookrightarrow \text{balance}(v) := \text{height}(R) - \text{height}(L)$ ($\forall v \in T$, $\text{balance}(v) \in \{-1, 0, 1\}$)

Store height of each node

Prove that height of AVL tree $\in O(\log n)$.

1. Height of AVL tree $\in \mathcal{O}(\log n)$

Let h be height of tree. # of nodes $<$ # of nodes in complete binary tree $(2^{h+1} - 1)$

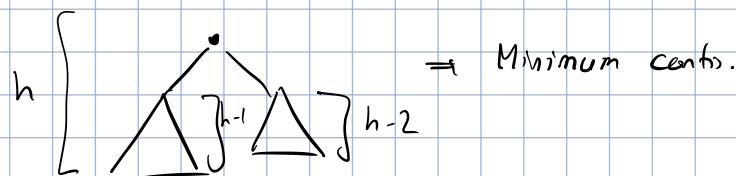
$$\therefore n \leq 2^{h+1} - 1$$

$$h \geq \log(n+1) - 1 \Rightarrow h \in \mathcal{O}(\log n)$$

2. Height of AVL tree $\in O(\log n)$

$N(h)$: nodes for AVL tree of height h (minimum)

$$N(h) = N(h-1) + N(h-2) + 1$$



Proof method!!
Minimum \Rightarrow balance = -1

Claim $N(h) \geq \sqrt{2}^h - 1$ (prove via induction)

$$\underline{n \geq N(h) \geq \sqrt{2}^h - 1}$$

$$n+1 \geq \sqrt{2}^h$$

$$\log(n+1) \geq h/2$$

$$h \leq 2\log(n+1) \in O(\log n)$$

Insertions:

insert $(k, v) \rightarrow z$ (node of structure)

1. Use regular BST root

2. Fix (unbalanced)

insert (k, v):

$z \leftarrow \text{bst_insert} (k, v)$

while (z is not null):

if $|z.\text{left} - z.\text{right}| > 1$

$y \leftarrow \text{taller child of } z$

$x \leftarrow \text{taller child of } y$

$z \leftarrow \text{restrict} (x, y, z)$

break

]} *No balancing fr root
of tree!!*

setHeightSubtrees (z)

$z \leftarrow z.\text{parent}$

} Precise up + fix heights

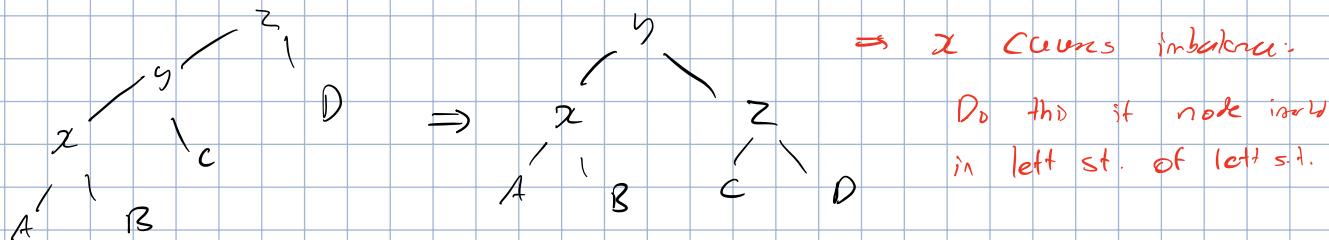
setHeightSubtrees (v):

Maximum subtree height

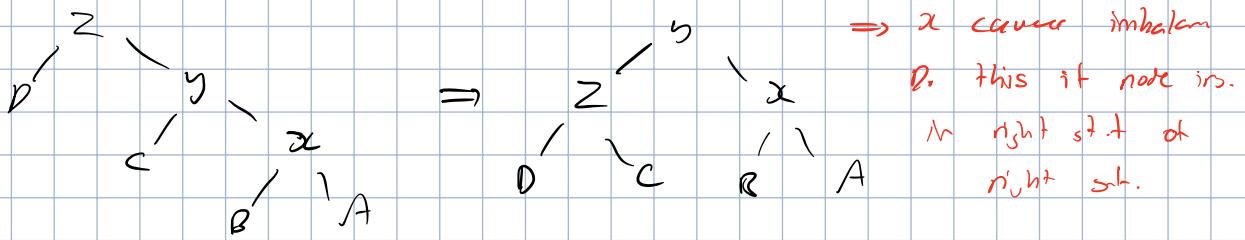
$v.\text{height} \leftarrow 1 + \max (v.\text{left}.\text{height}, v.\text{right}.\text{height})$

Restructuring: Rotation.

① Right rotation:

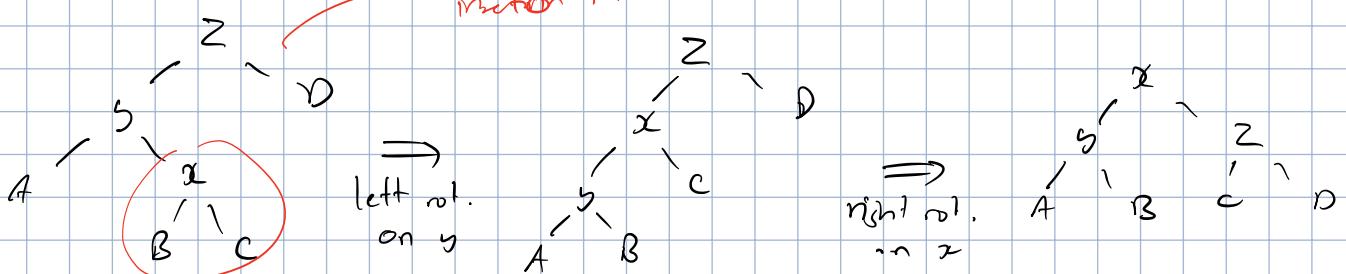


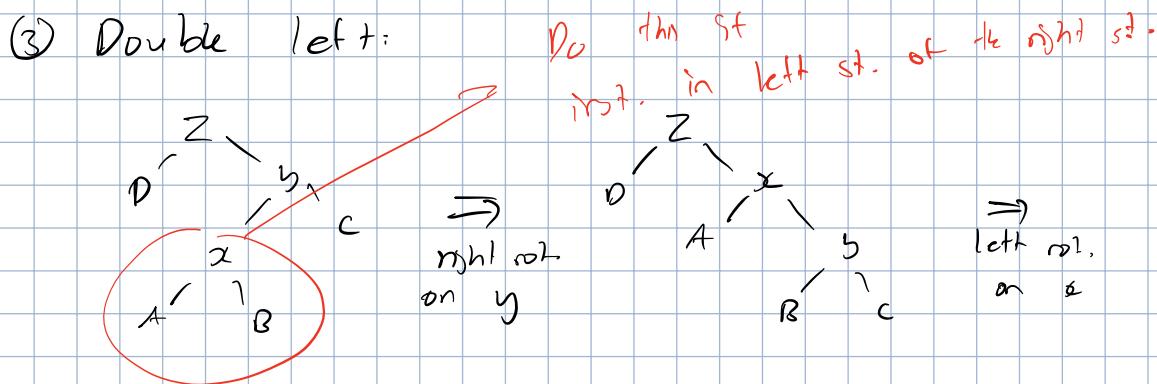
② Left rotation:



③ Double right:

Do this if
insert in right st. of left st.





Deletion:

Delete (k):

$$z \leftarrow \text{BST-delete}(k)$$

while (z not null)

if z imbalance:

$y \leftarrow$ taller chd of z

$x \leftarrow$ taller chd of y

$z \leftarrow \text{restruct}(z, y, x)$

Might have to do multiple restructuring.

set $lheight(z)$

$z \leftarrow z.parent$

Analysis:

Search: $\Theta(h) \Rightarrow \Theta(\log n)$

Restructure once

Insertion: $\Theta(h) \Rightarrow \Theta(\log n) \Rightarrow$ expensive!

Deletion: $\Theta(h)$, restructuring can be called $\Theta(h)$ tree $\Rightarrow \Theta(\log n)$

Restructuring: $\Theta(1)$ (scrapping, but exp.)

PROOF

Claim: no restriction \propto on insertn:

1. Tree after restructuring (T') is balanced

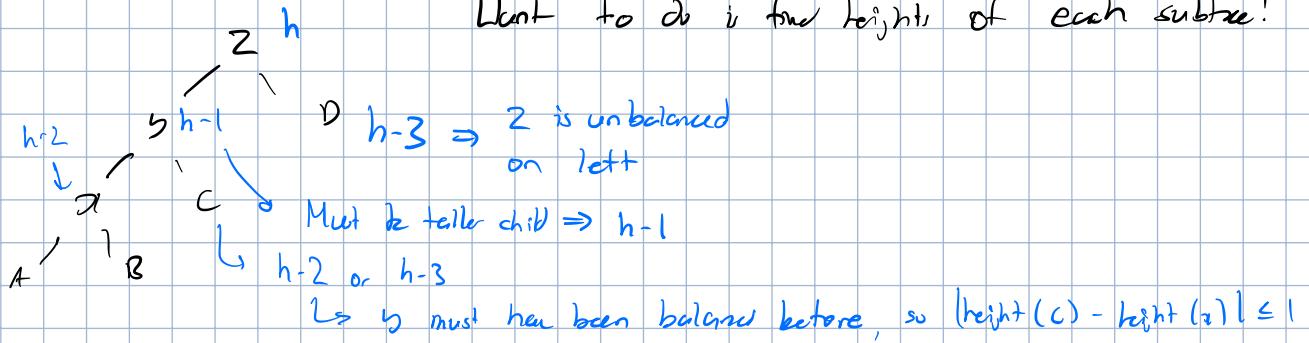
2. $height(T') = height$ of tree before insertn!

Proof method: case analysis on all restructuring methods

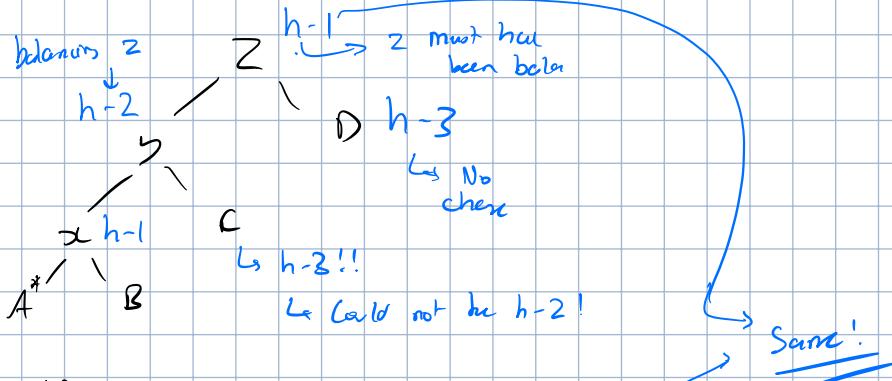
Case Analysis: before insertion, after insertion, after restructuring

Right rotation case:

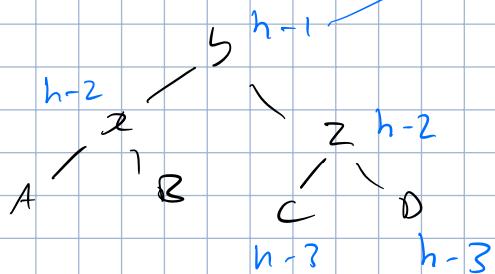
① After insertion:



② Before insertion:



③ After restructuring:



Continuation of Week 2:

Operations in Heap Tree

1. Insert

1. Insert at end

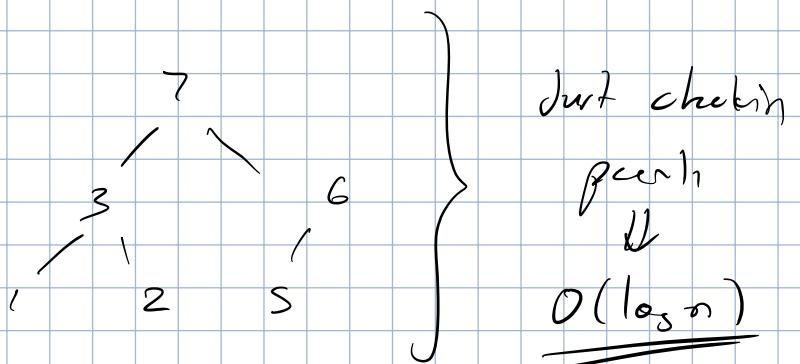
2. Fix up:

Is child < parent?

\hookrightarrow Swap

Else:

\hookrightarrow Done



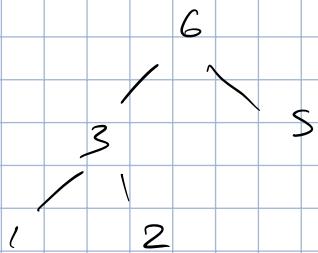
2. Deletion:

1. Replace top node w/ last node

2. Fix down

↳ Swaps parent w/ biggest child

↳ Do this until parent > children



PQ-Sort + Heap-Sort + Heapiify

PQ-Sort:

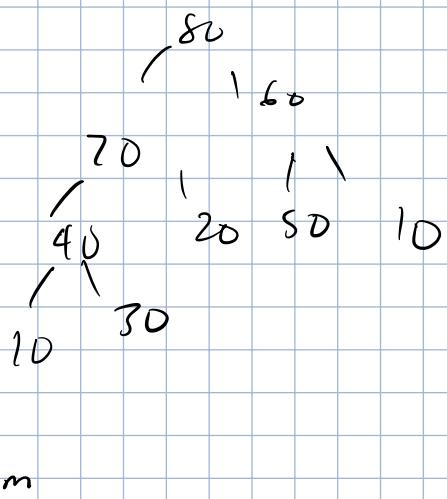
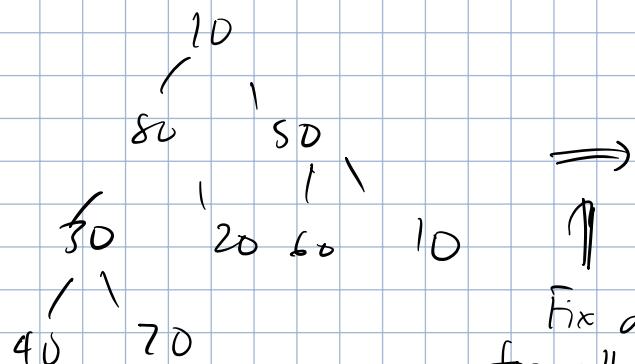
1. Put everything into heap ($\text{key} = \text{value} = A[\text{i}]$)

2. Delete max n times \Rightarrow back into Array starting at back

Time: $\Theta(n \log n)$. Space: $\Theta(n)$

Heapiify: takes array \rightarrow heap array $\Theta(n)$!

Initial



Why $\Theta(n)$?

$T(n) = \Theta(\text{worst case of steps})$

$$= \Theta(0 \cdot 2^h + 1 \cdot 2^{h-1} + 2 \cdot 2^{h-2} + \dots + h \cdot 2^{h-h})$$

$\# \text{ of levels}$ at bottom $\xrightarrow{2^h \text{ lat}} \leq 2$

$$= \Theta(2^h (0 + \frac{1}{2} + \frac{2}{2^2} + \frac{3}{2^3} + \dots + \frac{h}{2^h}))$$

$$= \Theta(2^h) \quad \xrightarrow{n = 2^{h+1} - 1}$$

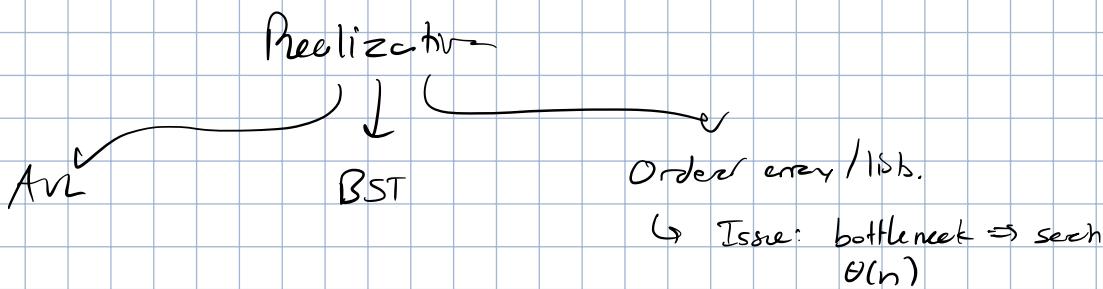
$$= \underline{\underline{\Theta(n)}}$$

Heap sort:

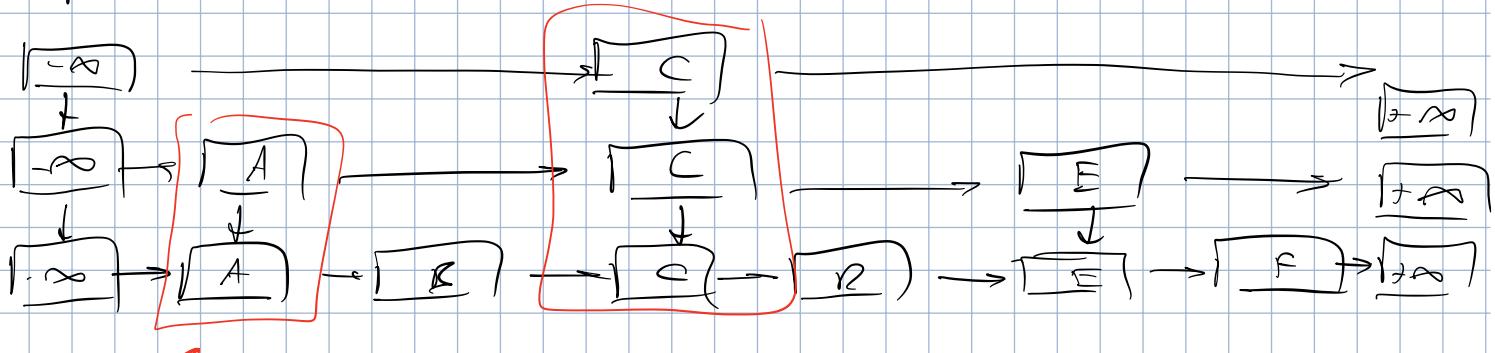
1. Heapsify

2. "Delete" but range keeps getting smaller so larger element at back of array

MODULE 5: SKIP LISTS



Skip: extension of LL \Rightarrow binary search.



Height 1 Tous. Height 2

Basic facts:

- Height: holding # of nodes at each height = $\sum_{i=1}^n \frac{n}{2^i} \in \Theta(h)$
- Search: visiting 2 nodes for each level \Rightarrow # of comp = $2 \log n \hookrightarrow \Theta(\log n)$

Insertion, deletion is diff \Rightarrow randomization

Randomization: randomly choose which node go to next level.

$\hookrightarrow S_0$: 0th level has all nodes

S_n : top level has only sentinels $(-\infty, +\infty)$

Height of tree = # of links in tree.

All nodes have following attriz prop:

1. After (next node)
2. Below (node below)

SEARCH

1. Init at top most $-\infty$

2. If $\text{after}(\text{node}) < \text{target}$

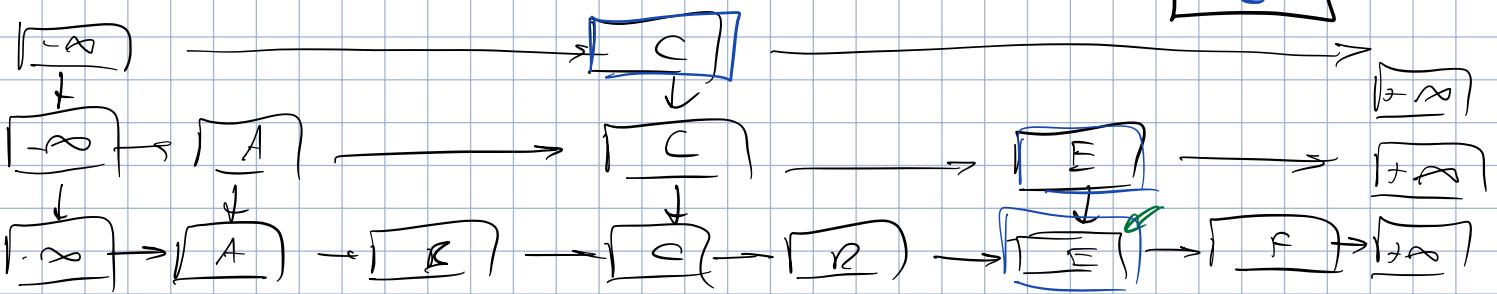
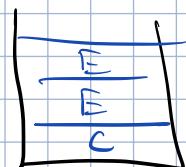
$\text{node} \leftarrow \text{node}. \text{after}$

3. Push node to stack P

4. $\text{node} \leftarrow \text{below}(\text{node})$

5. Do this until $\text{node} = \text{NULL}$

} set Predecessors (t_{pred})



Predecessor Stack: if t_{root} made into best tour stack
P holds which nodes would have son last

Search (t_{root}):

$P \leftarrow \text{set Predecessors} (t_{root})$

$\text{top} \leftarrow P. \text{top}()$

$\text{if } \boxed{\text{top. after}} \text{ key} == t_{root} : \text{return top. after.}$

$\text{return "not found"}$

Insertion:

Randomize w/ list of H, T

↳ Tour creation: # of times we get H before $T \Rightarrow$ # of times
node "promotes"

o $H, H, H, T \Rightarrow \text{insert in } S_0 \rightarrow S_1 \rightarrow S_2 \rightarrow S_3 \quad \square$
 $\boxed{\text{no match with}}$

insert (lc, v):

$i \leftarrow \# \text{ of heads in tails in randomization}$

$h \leftarrow \text{height of skip list so far. (incorrect for each time we do below()) on insertion)}$

If $h \leq i$ (tour height > existing)

Add additional sentinel layers.

$P \leftarrow \text{set Predecessors} (lc)$

// Insert into S_0

$P \leftarrow P. \text{pop}()$

$z \text{Below} \leftarrow \text{new Node with } b, v$

insertion (z_{Below}, p)

// Create tower:

while $i > 0$

$p \leftarrow P.\text{pop}()$

$z \leftarrow \text{new node with } b, v$

$z_{\text{below}} \leftarrow z_{\text{Below}}$

$z_{\text{Below}} \leftarrow z$

$i \leftarrow i - 1$

Deletion:

1. Get predecessor for key to delete

2. While $P.\text{top}.\text{after} == \text{key} \Rightarrow \text{delete } p.\text{after}$

↳ Checking entire tower

3. Remove any unnecessary sentinel.

delete (k):

$P \leftarrow \text{set Predecessor } (k)$

while P is not empty

$p \leftarrow P.\text{pop}()$

if $p.\text{after} = k$:

$p.\text{after} = p.\text{after}.\text{after}$

else:

break

$p \leftarrow \text{left sentinel of root}$

while $p.\text{below}.\text{after} = \infty$,

$p.\text{below} = p.\text{below}.\text{below}$

$p.\text{after} = p.\text{after}.\text{below}.\text{below}$

Remove useless
sentinel nodes.

Time analysis:

Like find $E[h] \Rightarrow$ height of skip list after all ops.

1. Define X_k : height of tower for key k

↳ P.R.V.! Randomization determine X_k

↳ Find $P(X_k \leq i)$

$$\circ P(X_{k^*} \geq 1) = P(1^{\text{st}} \text{ flip} > \text{head}) = \frac{1}{2} \quad (\text{above } S_0)$$

$$\circ P(X_{k^*} \geq 2) = P(1^{\text{st}} \text{ flip} \cap 2^{\text{nd}} \text{ flip} > \text{head}) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$$

$$\circ P(X_{k^*} \geq i) = \left(\frac{1}{2}\right)^i \quad (i \text{ flips needed})$$

2. Height is unbounded by $n \Rightarrow$ possible to have v. large height

↳ Expected height is \propto batch matrix

3. $|S_i| = \# \text{ of keys in level } i - \text{sentinels}$

↳ $|S_0| = n$

↳ $|S_i| = ?? \Rightarrow$ Use indicator variables!

4. Define indicator var to true expected level of keys in level.

$$I_{i,k} = \begin{cases} 0 & \text{if } k \text{ not in } S_i \\ 1 & \text{if } k \text{ in } S_i \end{cases} \Rightarrow P(k \text{ in } S_i) = \text{if height of } X_k \geq i \\ = P(X_k \geq i)$$

∴ Using indicator variable ths..

$$E[S_i] = \sum_{k=1}^n E[I_{i,k}] \Rightarrow \text{Go through all keys, calculate expected!}$$

$$= \sum_{k=1}^n P(X_k \geq i)$$

$$= \sum_{k=1}^n \left(\frac{1}{2}\right)^i$$

$$= \boxed{\frac{n}{2^i}}$$

S. Use $|S_i|$ to tree height.

Define another indicator variable.

$$I_i = \begin{cases} 0 & \text{if no keys in level } i \Rightarrow |S_i| = 0 \\ 1 & \text{if keys in level } i \Rightarrow |S_i| \geq 1 \end{cases}$$

$$\therefore h = 1 + \sum_{i \geq 1} I_i$$

$$E[h] = 1 + \sum_{i \geq 1} E[I_i]$$

Bounds on $E[I_i]$ \Rightarrow Done O from this:

$$1. I_i \leq |S_i| \Rightarrow E[I_i] \leq E[|S_i|] \leq \frac{n}{2^i}$$

$$2. I_i \leq 1 \Rightarrow E[I_i] \leq 1$$

Break the sum apart

$$\begin{aligned} E[h] &= 1 + \sum_{i=1}^{\log n} E[I_i] + \sum_{i=\log n+1}^{\infty} E[I_i] \\ &\leq 1 + \sum_{i=1}^{\log n} 1 + \sum_{i=\log n+1}^{\infty} \frac{n}{2^i} \\ &\leq 2 + \log n \end{aligned}$$

$$\therefore \underline{\underline{E[h] \in O(\log n)}}$$

Space:

Surveillance: $2h + 2$

$$\begin{aligned} \therefore E[\# \text{ of sentinel}] &= 4 + 2\log n + 2 \\ &\in \Theta(\log n) \end{aligned}$$

Space for node:

$|S_i| \Rightarrow \# \text{ of keys in level } i$

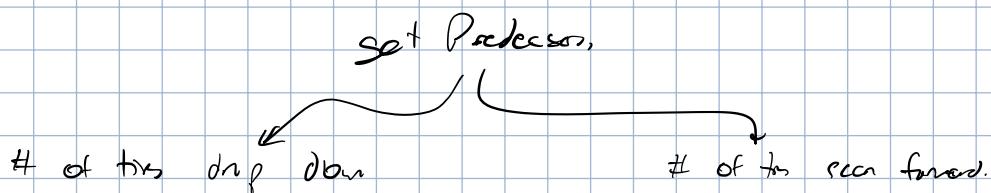
$$\hookrightarrow E[|S_i|] = \frac{n}{2^i}$$

$$\begin{aligned}
 \left\{ \left[\sum_{i \geq 0} |S_i| \right] \right\} &= \sum_{i \geq 0} E[|S_i|] = \sum_{i \geq 0} \frac{n}{2^i} \\
 &= n \sum_{i \geq 0} \frac{1}{2^i} \rightarrow \text{Converges to 2} \\
 &= 2n \in O(n)
 \end{aligned}$$

∴ Total space: $O(n)$

Search, insert, delete:

Dominates by set Predecessors.



Drop-down analysis:

Can't drop more than height

$$\# \text{ of drop downs} = O(h) = O(\log n)$$

Scan-forward analysis:

Let v be the leftmost node in S_i (what we're at)

Let w be the node after it

↳ $\text{height}(w) = i \Rightarrow \text{No way } \text{height}(w) > i \text{ b/c would have compared it already!}$

$$P(\text{tow } w=i \mid \text{tow of } w \geq i) = P(\text{set bits } T)$$

$$= \frac{1}{2}$$

$$\therefore P(\text{scanning for } w) \leq \frac{1}{2}$$

$$\therefore P(\text{scanning for } \ell \text{ bits}) \leq \left(\frac{1}{2}\right)^\ell$$

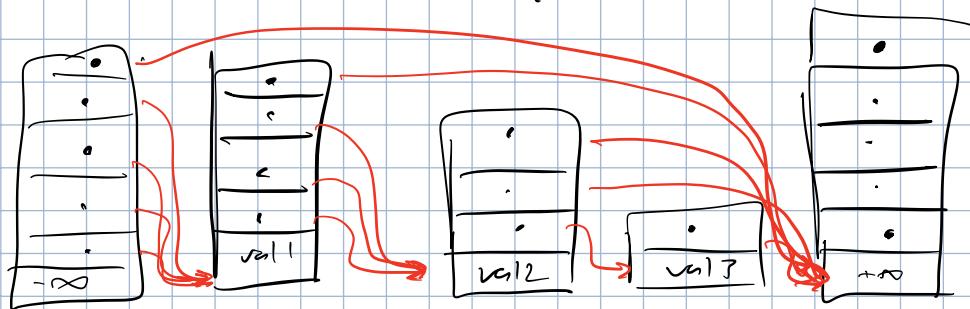
$$E \{ \# \text{ of scan towers} \} = \sum_{k \geq i} k \left(\frac{1}{2}\right)^k = 1$$

Scan towers $\in O(1)$

\therefore Get predecessors $= O(\log n)$
 \hookrightarrow Insert, search, deletion $\Rightarrow \underline{\log n}$

Optimization:

1. Make each tower an array:



Level = array index \Rightarrow find necessary links to next tower
 $\underline{\text{or re-decment in arr}}$

2. Biased coin flips

$$P(X_k \leq i) = \left(\frac{1}{2}\right)^i \Rightarrow \text{Fair}$$

$$\hookrightarrow P(X_k \leq i) = p^i \Rightarrow p < \frac{1}{2} \text{ (unfair)}$$

$$\text{Dom analysis again} = \log_{1/p} n < \log_2 n$$

Reordering Items:

Unsorted array / LL:

- \hookrightarrow Search: $O(n)$ Bottlenecked.
- \hookrightarrow Insert: $O(1)$
- \hookrightarrow Deletion: $O(n)$

One major optimization: put the most freq. accessed item at front.

Optimization

Access distr. known

Order keys from highest access prob. to lowest access prob.

Ex: //

keys	A	B	C	D	E
freq.	2	8	1	10	5
prob.	$\frac{2}{26}$	$\frac{8}{26}$	$\frac{1}{26}$	$\frac{10}{26}$	$\frac{5}{26}$

1 possible order:

C A B D E

Expect # of comp:

$$E[\sum C_i] = \frac{2}{26} \cdot 1 + 2 \cdot \frac{2}{26} + 3 \cdot \frac{8}{26} + 4 \cdot \frac{10}{26} + 5 \cdot \frac{5}{26}$$

$$\approx 3.261$$

Another: heuristics.

D B E A C

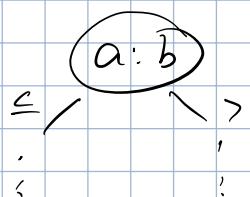
$$E[\sum C_i] = 2.54$$

MODULE 6: SPECIAL KEY DICT

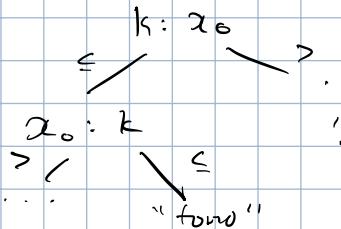
Lower Bound For Search

Comparison-based dict $\Rightarrow \Omega(\log n)$

Proof:



\Rightarrow How does this account for equality



Double check.

"tomb"

of leafs = # of distinct inputs/output,
 $= n + 1 \Rightarrow$ Any of n options to search in array
 or did not find any (+1).

$$\therefore 2^h \geq n + 1$$

$$h \geq \log(n+1)$$

$$\therefore h \in \Omega(\log(n+1))$$

Interpolation Search

Can we possibly speed up binary search?

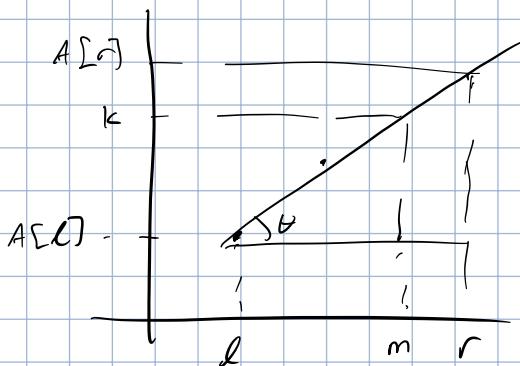
Recall binary search:

$$\text{Middle point. } \left\lfloor \frac{l+r}{2} \right\rfloor = l + \left\lfloor \frac{1}{2}(r-l) \right\rfloor$$

Interpolation search:

$$m = l + \left\lfloor \frac{k - A[l]}{A[r] - A[l]} (r-l) \right\rfloor$$

Idea: we don't need to check middle. If we know that our array grows linearly \Rightarrow height on side or middle.



$$\tan \theta = \tan \theta$$

$$\frac{A[r] - A[l]}{r-l} \approx \frac{k - A[l]}{m - l}$$

$$\therefore m = l + (r-l) \cdot \left(\frac{k - A[l]}{A[r] - A[l]} \right)$$

Analysis?

If keys are linear + uniform $\Rightarrow T^{\text{avg}}(n) = T^{\text{avg}}(\sqrt{n}) + \Theta(1) \in O(\log \log n)$

\hookrightarrow Worst-case: $\Theta(n)$.

Showing $T(n) = T(\sqrt{n}) + c \in O(\log \log n)$

Assume n is $\text{pow}(2) \Rightarrow 2 \rightarrow 2^2 \rightarrow 2^4 \rightarrow 2^8 \rightarrow 2^k \rightarrow 2^{2^i}$

Go backwards:

$$T(2^{2^i}) = T(2^{2^{i-1}}) + c$$

$$= T(2^{2^{i-2}}) + 2c$$

$$= T(2^{2^{i-3}}) + 3c$$

⋮

$$= \underbrace{T(2^{2^{i-i}})}_{O(1)} + ic$$

$$2^{2^i} = n \Rightarrow 2^i = \log n \Rightarrow i = \log \log n$$

Pseudocode:

interpolation-search (A, k, n):

$$l \leftarrow 0, r \leftarrow n-1$$

while ($l \leq r$):

 ↗ Bound check

 if ($A[k] < A[l]$ or $A[k] > A[r]$): "not found"

 if ($A[l] == A[r]$) return "found at l " \Rightarrow Avoid division by zero

$$m \leftarrow l + \frac{(k - A[l])}{A[r] - A[l]} (r - l)$$

 if $A[m] < A[k]$:

$$l \leftarrow m + 1$$

 else if $A[m] > A[k]$:

$$r \leftarrow m - 1$$

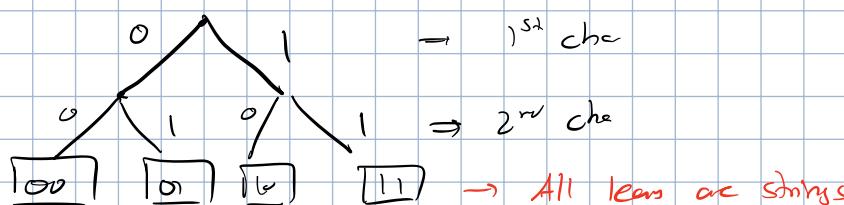
else:

 found at m

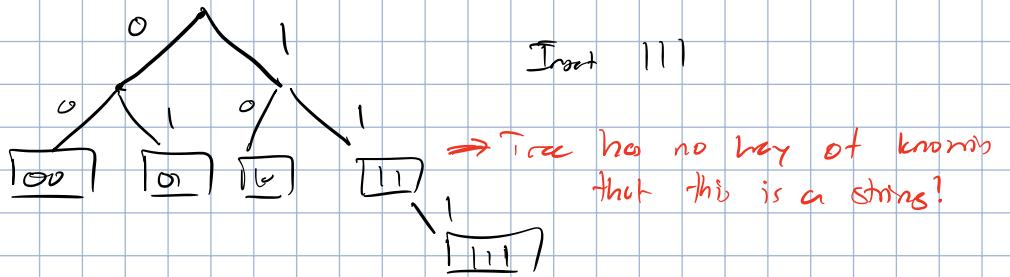
Tries

Dictionary for strings

Bitstring tries: bitstrings (0, 1, 001, 110, ...)

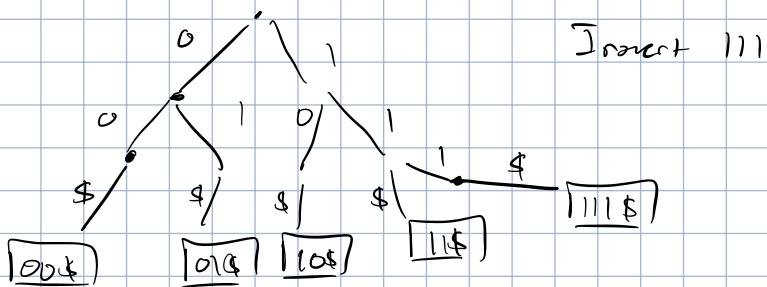


Problem: prefixes



Soln:

1. Enforce all strings of same length
2. More flexible: ensure no 2 strings can be a prefix \Rightarrow append end of word char (\$)



Operations:

① Search:

search (v, d, x)

the node \downarrow *digit* \leftarrow *string searching*

if v is a leaf: return v

else:

$v' \leftarrow$ child of v labeled $\sim x[d]$

if no v' : return None \Rightarrow Nr child \leftarrow digit

else:

return search $(v', d+1, x)$

② Insert:

Do search \rightarrow last successful node \rightarrow add more nodes until tree has

③ Deletion:

Do search \rightarrow remove node & ancestors until ancestor has
other child

Analysis:

• Time: $\Theta(\omega)$ \Rightarrow $\omega = 1$ string | (charachter comp.)

• Space: $\text{L} \times \text{W} \times \text{C}$ \Rightarrow each word has distinct path

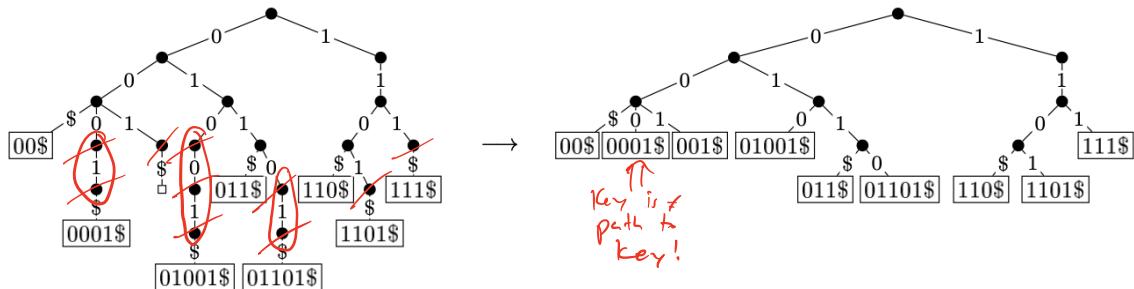
$O\left(\sum_{\text{word}} |\text{word}|\right)$ \Rightarrow Usually words share paths.

Variations:

① Don't store keys \Rightarrow path uniquely identifies string.

② Allow prefixes: have an intional flag if node is word or not
 \hookrightarrow No need for \$

③ Pruned tree: a node has a child \Leftrightarrow at least 2 descendants.

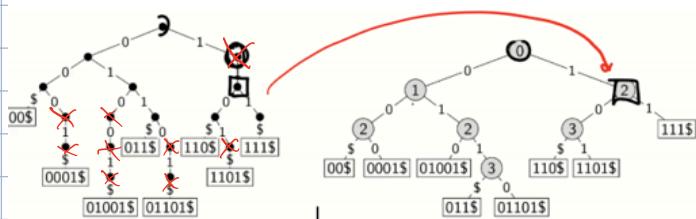


How to store keys

How to convert a tree \Rightarrow pruned tree

- ① For each node, check if it has ≥ 2 leaf descends
- ② If so, move on
- ③ Else, delete.

④ Compressed tree: only store node w 2 children



Additional thing: each node has index \Rightarrow depth of node.

Why? Skipping letters.

Compressed.

Trie \cup n keys \Rightarrow # of internal nodes $\leq n-1$ \Rightarrow Prove via induction.
 # of leaves $= n$
 # of nodes $\leq 2n-1$

Operations in compressed trie:

Search ($v \leftarrow \text{root}, x$)

if v is a leaf:
 return $x = v.\text{string}$

$d = v.\text{index}$

if $d > \text{len}(x)$:
 return "not found"

$v' \leftarrow \text{child of } v \text{ in path from } v \text{ to } x[d]$

if no child:
 return not found.

Search (v', x).

Paths don't uniquely determine str \Rightarrow how to compress key

Skipping nodes & length!
 We need to check if depth is too much.

$O(|x|)$

delete: search \rightarrow remove \rightarrow compress path

insert: search to maximal \rightarrow uncompress \rightarrow insert \rightarrow recompress.

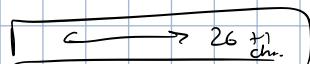
⑤ Multinary trie: over any alphabet Σ

Σ $\rightarrow |\Sigma| + 1$ end of string char.

Runtime: $O(|x| \cdot \text{time to find child})$

Soln 1: Array of size $|\Sigma| + 1 \Rightarrow$ search in char spot.

node: a



$O(1)$ for child search, but $O(|\Sigma| + 1)$ space / node.

Soln 2: children sorted by their letter.

$O(|\Sigma|)$ to search child, $O(1)$ space.

→ Soln 3: Dictionary of children

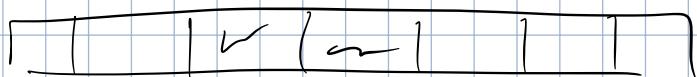
Hash, AVL, skip $\Rightarrow O(1) \rightarrow O(\log n)$ time, $O(|\Sigma|)$

MODULE 7: DICTIONARIES VIA HASHING.

Intro

$O(1)$ search:

- Direct addressing: know the range of keys



$O(1)$ time, $O(M)$ space $\Rightarrow M$ binds the range.

Disadvantages:

1. Only works for ints
2. Space wasted $\Rightarrow M \gg n$

- Hashing: keys $\rightarrow \{0, \dots, M-1\}$ for direct addressing

Need hash function

Disadvantage: collision (map to same key)

Dictionary that uses hashing \Rightarrow hash table.

Probability of collision:

n independent values uniformly

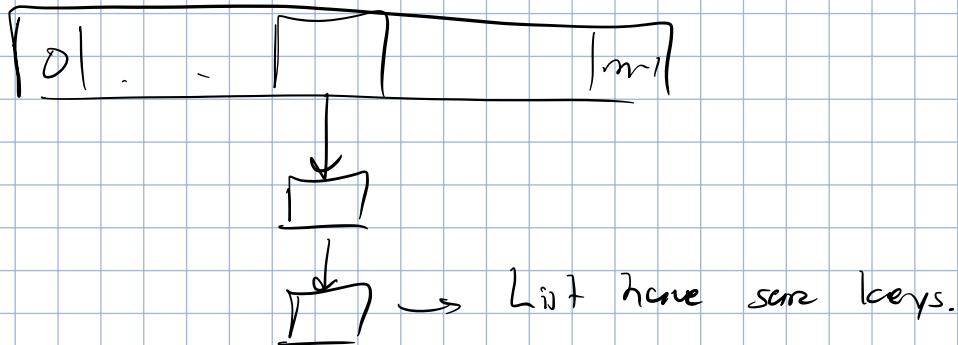
Prob. of no collision: $\frac{M(M-1)(M-2)\dots(M-(n-1))}{M^n}$

Prob. of collision = $1 - \frac{M(M-1)(M-2)\dots(M-(n-1))}{M^n}$

Can't compute:

What we can do \Rightarrow Fix probability α of collision & choose M accordingly

Separate Chaining:



Operations:

Search (k): Find k in $T[h(k)]$

Insert (k, v): Insert k, v in $T[h(k)]$

Deletion (k, v): search \rightarrow delete

Analysis:

Insertion: $O(1)$

- Hashing: $O(1)$

- Insertion in LL: front $\rightarrow O(1)$

Search:

- General: $O(1 + \text{size of bucket } T[h(k)])$

- Load factor:

$$\alpha = \frac{n}{M} \Rightarrow \text{Ave. bucket size.}$$

- Ave. case $\geq O(1 + \alpha)$

- Uniform Hashing Assumption:

\forall key k , $\forall j \in \{0 \dots M-1\} \Rightarrow h(k)=j$ has prob. $1/M$

- Analysis can be tight.

Claim: Each key will collide w/ $\frac{n-1}{M}$ other keys \Rightarrow avg. cost $\in O(1 + \alpha)$.

\downarrow
 $\frac{n-1}{M}$ bucket size!

Proof:

$\forall i, i' \in \{1, \dots, n\}, i \neq i'$

$$\begin{aligned} P(h(k_i) = h(k_{i'})) &= \sum_{j=0}^{M-1} P(h(k_i) = h(k_{i'}) = j) \\ &= \sum_{j=0}^{M-1} \frac{1}{M} \cdot \frac{1}{M} \\ &= \frac{1}{M} \end{aligned}$$

\therefore Prob. of collision = $\frac{1}{M}$

$$\text{Let } X_{ii'} = \begin{cases} 0 & \text{if } h(k_i) \neq h(k_{i'}) \\ 1 & \text{if } h(k_i) = h(k_{i'}) \end{cases}$$

Expected # of collisions: $\sum_{i \neq i'} E[X_{ii'}]$

$$= \sum_{i \neq i'} P(\text{collision})$$

$$= \sum_{i \neq i'} \frac{1}{M}$$

$$= \frac{n-1}{M}$$

o Key: keep α small

o Say α too big $\Rightarrow M$ is bigger \Rightarrow rehash into bigger table.

↳ Really, ignore

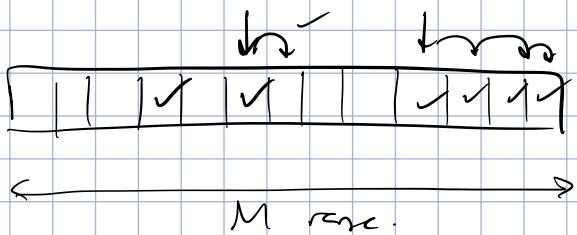
o Say α too small $\Rightarrow M$ is smaller \Rightarrow rehash into smaller table.

o If $\frac{1}{4} < \alpha < 2 \Rightarrow M \in O(n)$

$\therefore \alpha \in O(1) \Rightarrow$ Cost of chaining is $O(1)$

+ Space: $O(n)$

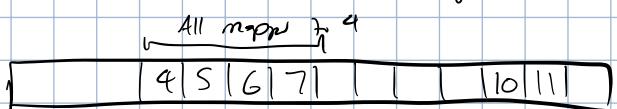
Probe Sequences



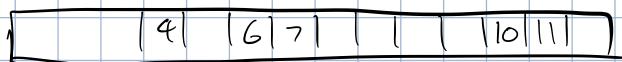
Operations:

Search/insert: Map to $T[h(k)] \rightarrow$ move toward until key found.

Deletion: do search \rightarrow mark spot del ✓ flag.



1. Delete \rightarrow (empty)



2. Search for 6

\hookrightarrow Map to 4

\hookrightarrow Stop @ empty spot!

How to continually move to next spot

$$h(k, i) = (h(k) + i) \bmod M$$

Problem: clustering issues (all keys map to similar hash \rightarrow cluster of keys to search through).

1. Quadratic hashing:

$$h(k, i) = (h(k) + c_1 i + c_2 i^2) \bmod M$$

Still have clustering issue, but diff keys have diff patterns, so not as bad.

Make sure c_1, c_2 are chosen well

2. Double hashing:

$$h(k, i) = (h_0(k) + i h_1(k)) \bmod M$$

Key: $h_1(k)$ is non-zero (o.w. collision will fire)

$$\& h_0(k) \perp h_1(k)$$

How to make $h_0(k) \perp h_1(k)$?

↪ Multiplicative method.

$$h_1(k) = 1 + \lfloor (M-1) (kA - \lfloor kA \rfloor) \rfloor \xrightarrow{0 \rightarrow M}$$

\uparrow Non-zero. \uparrow Fractional: $[0, 1]$

Typically $A = \varphi$

Make M prime reach all spots in hash table & h_1 is rel. prime to h_0

Ex/ll Insert 194. $M = 11$. Double hashing.

① Make indep. hash functions via multiplicative multi.

$$h_0(k) = k \bmod 11$$

$$h_1(k) = 1 + \lfloor 10 (k\varphi - \lfloor k\varphi \rfloor) \rfloor$$

$$\therefore h(k, i) = (h_0(k) + i h_1(k)) \bmod 11$$

② Probe sequence:

$$h(k, 0) = 7$$

$$h(k, 1) = 5$$

$$h(k, 2) = 3$$

$$h(k, 3) = 1$$

$$7 + h_1(194)$$

$$11$$

$$1 + \lfloor 10 \times (194 \cdot \varphi - \lfloor 194 \cdot \varphi \rfloor) \rfloor$$

$$= 9$$

$$7 + 18 = 25$$

$$7 + 27 = 34$$

③ Go down probe sequence &

find first spot available.

Cuckoo Hashing

2 tables: $T_0(h_0)$, $T_1(h_1)$

A key is in either T_0 or T_1 , cannot be in both.

Operations:

Search: $T_0[h_0(k)]$ or $T_1[h_1(k)]$

Delete: Search \rightarrow delete.

Insertion:

① Map $h_0(k)$. If not occupied \rightarrow done

② If occupied, replace k' w/ k and hash $T_1[h_1(k')]$

③ If not occupied, done

④ If occupied, replace k'' w/ k' and hash $T_0[h_0(k'')]$

◦ Operation: $O(1)$ time if α is small.

◦ Pseudo code:

insert (k, v) :

$i \leftarrow 0$

do at most $2n$ times:

if $T_i[h_1(k)]$ is empty:

$T_i[h_1(k)] = (k, v)$

return "success"

swap $((k, v), T_i[h_1(k)])$

$i \leftarrow i + 1$

return "failure"

If inserted 2n+1 times,
it will go forever

Ex:// $M = 11$, $h_0(k) = k \bmod 11$, $h_1(k) = \lfloor 11(\psi_k - \lfloor \psi_k \rfloor) \rfloor$

	T_0		T_1
0	44	0	
1		1	
2		2	
3		3	
4	26	0	
5		5	51
6		6	
7	51	7	95
8		8	
9	92	9	
10		10	

① insert (95)

1. T_0 hashes

$$h_0(95) = 7$$

2. kick out \rightarrow hash to T_1 ,

$$h_1(51) = 5$$

② insert (26)

1. T_0

$$h_0(26) = 4$$

$$2. T_1 \\ h_1(s_1) = S$$

$$3. T_0 \\ h_0(s_1) = T$$

$$4. T_1 \\ h_1(s_2) = T$$

③ search (s₂)

$$h_0(s_2) = 4 \rightarrow ? \rightarrow x \\ h_1(s_2) = 8 \rightarrow ? \rightarrow \checkmark$$

Notes:

1. $|T_1| \neq |T_0| \Rightarrow$ can be diff., can be same

2. Load factor:

$$\alpha = \frac{n}{|T_1| + |T_0|}$$

If $\alpha < 1/2 \Rightarrow O(1)$ insertion.

3. Order of T_0 / T_1 doesn't matter.

Programmer does load
chores.

All operations w/ hashing take avg. time of $O(1)$. Worst case: $O(n)$

Hash Function Strategies

Prime objective: we want uniform hashing assumption

↳ We want 2 things:

① Hash functions \perp data

Want to hash pairs (x, y)

② Depend on all key parts

Combo (x, y) is unique, not just x or just y

$$\therefore h_0((x, y)) = x + y \bmod M$$

Int. keys

Modular:

$$h(k) = k \bmod M$$

M is prime

Multiplication

$$h(k) = \lfloor M (kA - \lfloor kA \rfloor) \rfloor \text{ for some } 0 \leq A < 1$$

No hash func. is perfect \Rightarrow exists some seq. of keys that do poorly

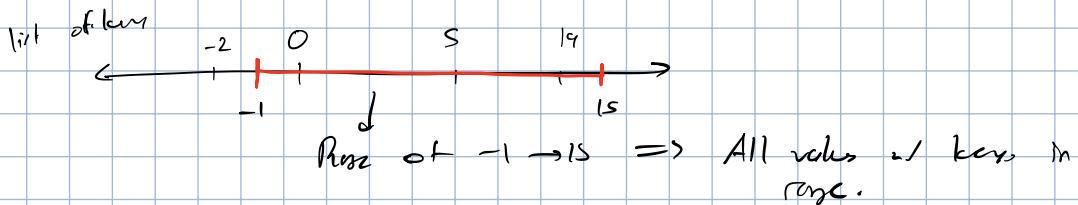
Multidimensional keys: flatten to int.

- Ex:// $APPLE \Rightarrow 65R^4 + 80R^3 + 80R^2 + 76R + 68$ (ASCII)
 \hookrightarrow Mod M

MODULE 8: RANGE SEARCHING IN DICTS

Range Searching:

- 1D case:



Note: We have multiple items to search, insert, delete

\hookrightarrow Analysis based on output size: s

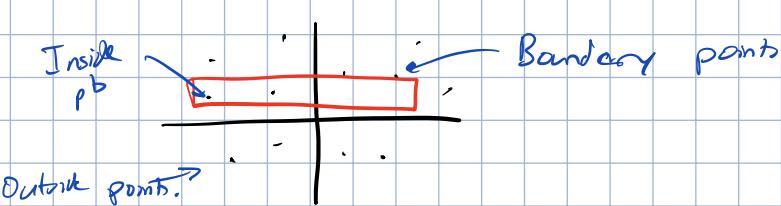
1D implementations:

- Unsorted array/hash table/list: check each item in D.S.
 - $\Omega(n)$
- Sorted array/list: binary search for endpoints \Rightarrow return output
 - $O(\log n + s)$
- BST: almost ~~sort~~ $\Rightarrow O(h + s)$

- Multidimensional search:

- Aspects: dimensions
 - ≥ 2 dimension keys. Ex:// geographical coordinates.
- Multiple aspects \Rightarrow search across all aspects.

\hookrightarrow Define query polygon (2D case, do a rectangle).

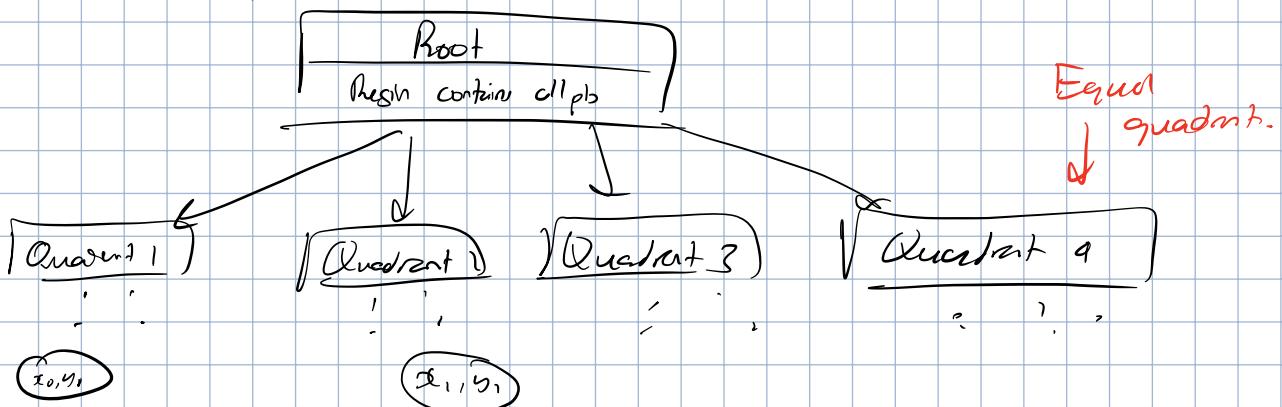


Implementation:

1. Try to map aspect \rightarrow 1 dimension: very complicated
2. Use dictionary for aspect \rightarrow get points \rightarrow union: inefficient
3. Other data structures.

Quadtrees

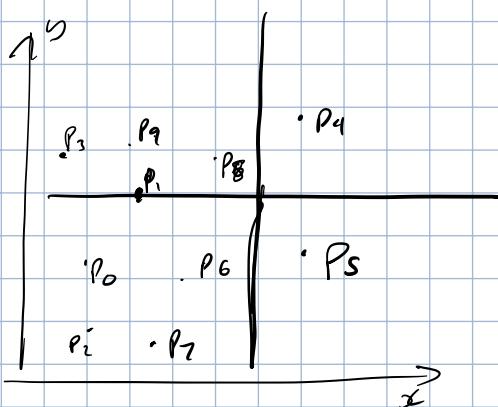
Defn: quadtree of n points $\{(x_0, y_0), (x_1, y_1), \dots\}$



o Recursively define struct:

1. If quadrant has 1/0 pts \Rightarrow point as leaf
2. Split into equal size quads if not
 - ↳ Split point set to match split in quads.
 - ↳ Convention: points on line \Rightarrow go to top/right quadrant

Ex://



Draw quadtree:

- ① Root \Rightarrow all points

Find out max/min of x, y

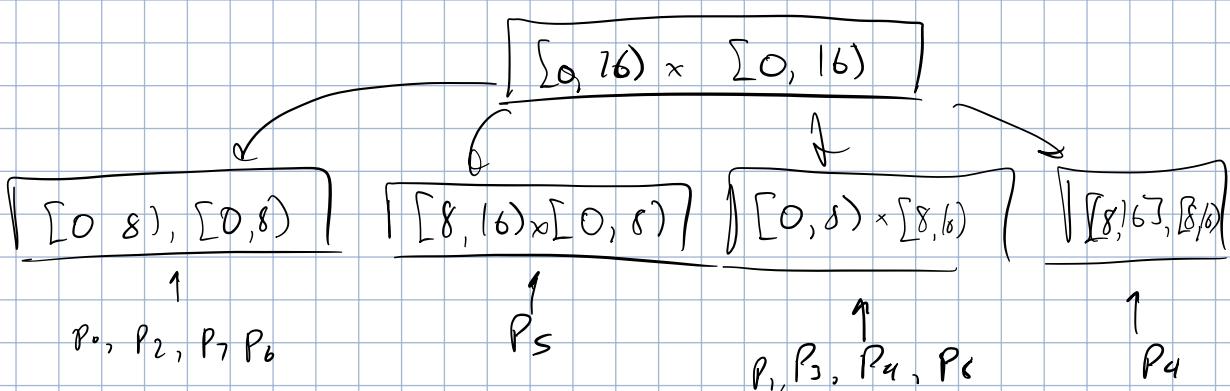
$$[0, 16) \times [0, 16)$$

- ② Subdivide

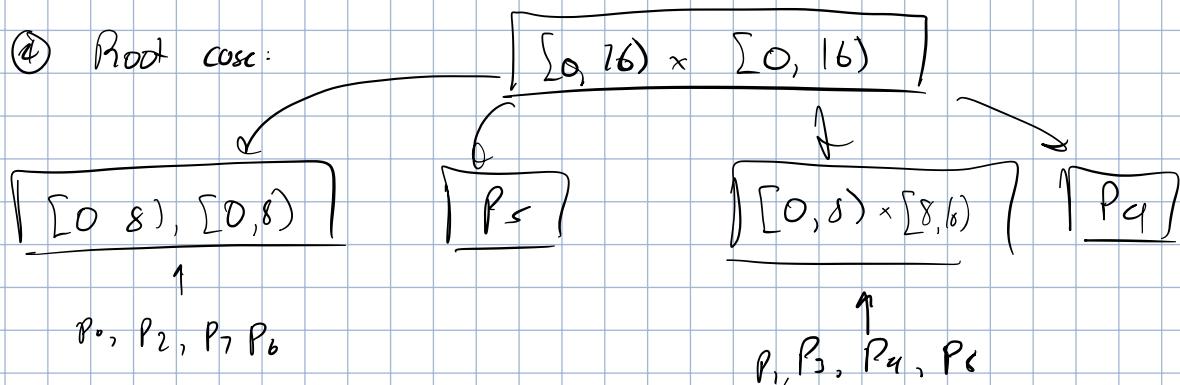
$$[0, 16) \times [0, 16)$$

$$[0, 8) \times [0, 8) \quad [8, 16) \times [0, 8) \quad [0, 8) \times [8, 16) \quad [8, 16) \times [8, 16)$$

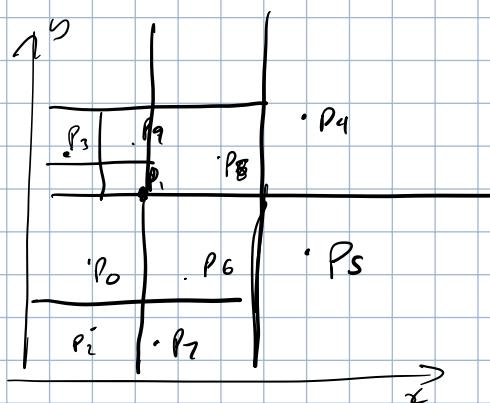
- ③ Subdivide points.



④ Root case:

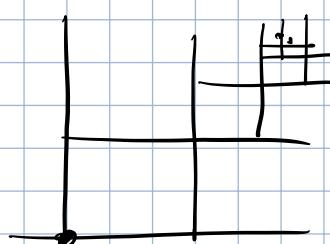


⑤ Recur until pts are leaves.

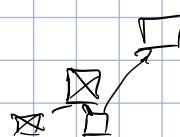


Operations on Quadtree:

1. Single point search: BST search w/ x & y comparisons.
2. Insert 1 point: search \rightarrow leaf \rightarrow split leaf in 4 quadrant until have 1 point.



3. Delete: search \rightarrow delete \rightarrow recursively delete parent if only 1 child.



4. Range search:

range search (r, A):

$R \leftarrow$ region of root

if $R \subseteq A$: R completely within A : inside node
return everything below r

if $R \cap A = \emptyset$: R is not overlapping w/ A : outside node
return

if r is a leaf:

$p \leftarrow$ point at r

if $p \in A$:

return p

else:

return range search (all children, A)

} Boundary node.

Analysis?

Note: Height of quadtree \propto # of points?

↳ Distribution of points determines height

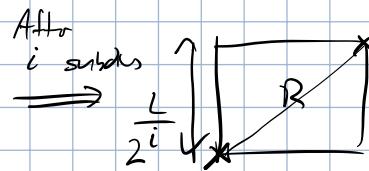
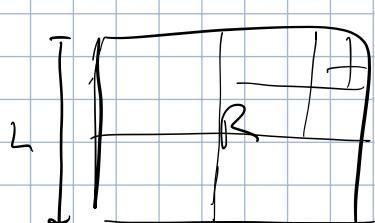
Distribution measure: spread factor.

$$\beta(s) = \frac{\text{sidelength of } R}{\text{min. distance of pts in } s}$$

$$h \in \Theta(\log(\beta(s)))$$

↳ Proof that $h \in \Theta(\log(\beta(s)))$:

Region R :



Thus:

Max distance between pts: $(\sqrt{2} \cdot \frac{L}{2^i})$

$$d_{\min} \leq \sqrt{2} \cdot \frac{L}{2^i}$$

$$\begin{aligned}
 2^i &\leq \sqrt{2} \frac{L}{d_{\min}} \\
 2^i &\leq \sqrt{2} \beta(s) \\
 \# \text{ of subdts} &= \text{height of tree} \Rightarrow i = \log(\sqrt{2} \cdot \beta(s)) \\
 \therefore h &\in O(\log(\beta(s)))
 \end{aligned}
 \quad \boxed{\text{Root technique:}}$$

↳ Try to subdivide
 R into side lengths \Rightarrow
 spread factor

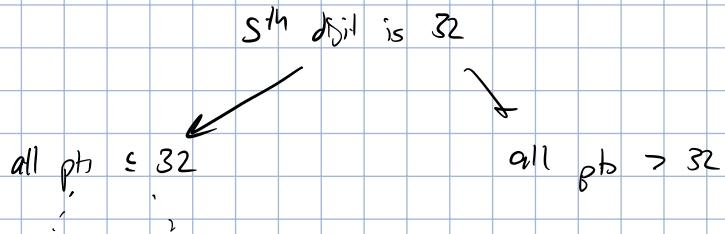
Operation analysis:

Time to build quadtree: $\Theta(nh)$

Range search: $\Theta(\ln h)$ \Rightarrow search through all nodes

Multidimensional Quadtrees:

1D: put pts in box 2 \Rightarrow use tree logic to subdivide



3D: octree

Variations:

1. Multiple pts in 1 leaf
2. Stop splitting prematurely

kd Trees

Idea: split points into 2 \Rightarrow equal # of points / not equal size

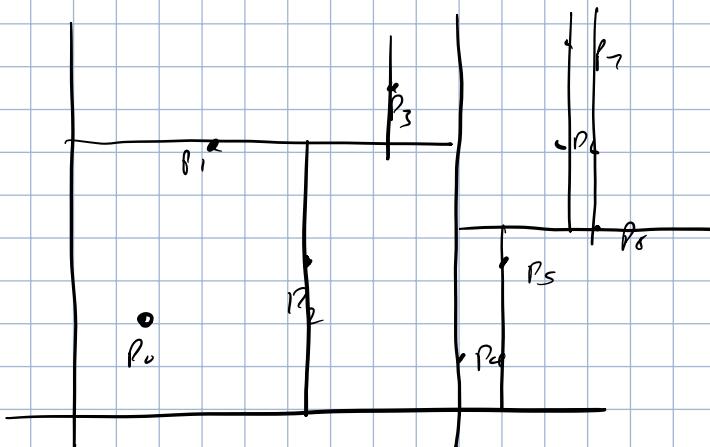
Construction:

- (1) If $|S| = 1 \Rightarrow$ leaf
- (2) Else: Do quick select to find median x-coord

↳ Partition into $S_{x < \bar{x}}, S_{x \geq \bar{x}}$

- (3) Recursively create kd tree on left/right subtrees but alternate y coordinates.

Ex: //



$$x < p_q \cdot x$$

Y

$$\boxed{y < p_1 \ldots}$$

$$d < p_2 \cdot x$$

$$\downarrow x < p_3 \cdot x$$

A hand-drawn diagram on lined paper showing a vertical cylinder with a horizontal top cap. A curved line extends from the top right corner of the cap, representing a gauge connection. Below the cylinder, the label P_0 is written.

TR

P₁)

Running Time Analysis:

Binary search $\Rightarrow \Theta(\log n)$ (quick select)

Precison:

$$T^{cof} = 2T\left(\frac{n}{2}\right) + O(n) \xrightarrow{\text{Merge}} T(n) \in O(n \log n)$$

↓ ↓
 Recuse constructu

Height: $O(\log n)$

Operations:

1. Search for point : BST $\Rightarrow O(\log n)$

2. Insertion: BST insertion will impact root \Rightarrow unbalanced, so $n \notin O(\log n)$

3. Deletion: Best deletion (Solution: 

- ① Prebuild tree
 - ② Point $cur \Rightarrow$ reverse searching tree layer.

Bottom line: not good for insertion/deletion.

4. Preje search: exact score also as good tree, use 2 children.

↳ Time complexity = # of points returned * # of recursion.

b
D(s)

$$Q(n) \leq 2Q\left(\frac{n}{q}\right) + O(1)$$

$\hookrightarrow O(\sqrt{n})$

$$\therefore O(\leq \sqrt{n})$$

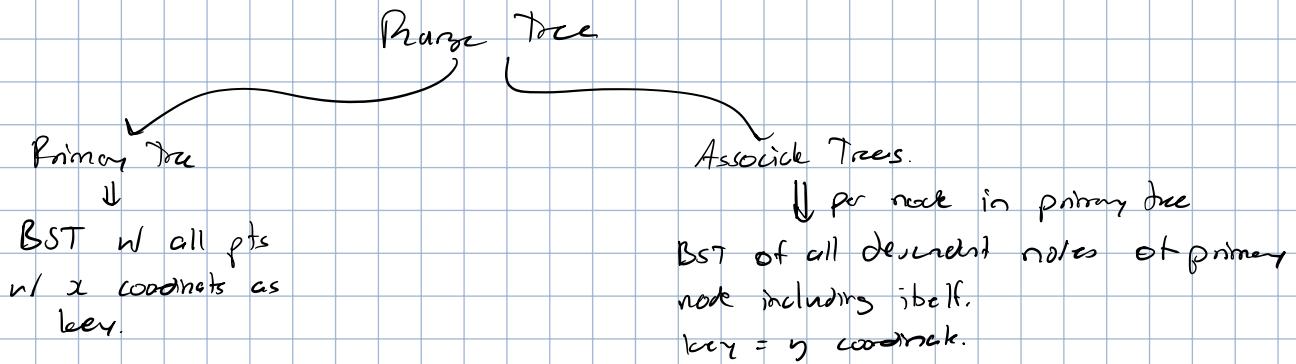
Multidimensional:

Attack all d dimensions when doing splits

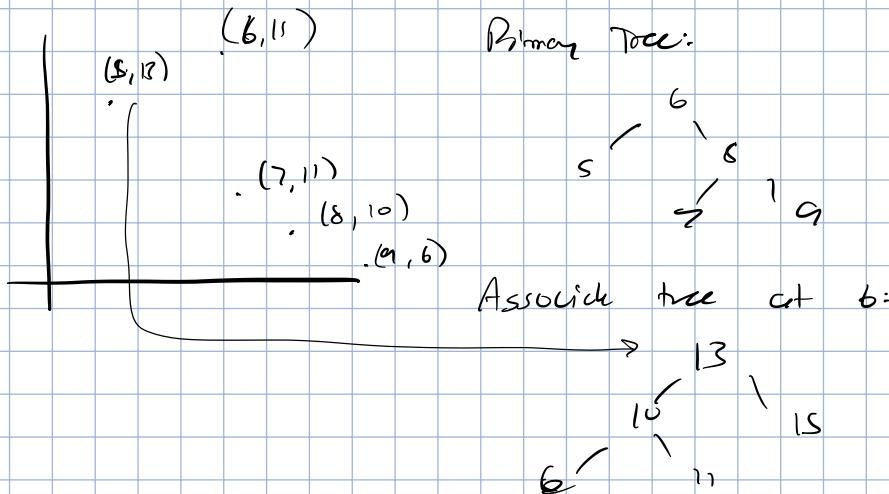
Space complexity: $O(n)$, Construction time: $O(n \log n)$

Range search: $O(s + n^{1-\frac{1}{d}})$

Range Trees



Ex://



Space complexity:

Primary tree: $O(n)$

Associate tree: $T(v)$ will have $O(|P(v)|)$

Size of all associate trees: $\sum_v |P(v)| = \sum_{v, w} \delta_{v, w} \Rightarrow \delta_{v, w} = \begin{cases} 1 & \text{if } w \in P(v) \\ 0 & \text{if } w \notin P(v) \end{cases}$

\uparrow $v \rightarrow$ center of w

\uparrow $v \rightarrow$ not center of w

$$= \sum_v \left[\sum_w \delta_{v, w} \right] \Rightarrow \# \text{ of ancestors for node } w \text{ (fixing } w \Rightarrow \text{ going through all nodes!!)}$$

Think: # of ancestors of node? Path to root: $O(\log n)$

of points = n

Size of all associate trees: $O(n \log n)$

Operations:

1. Search for single coordinate: BST using x
 - if x is in tree, do
 - if not, do a search on associate tree via B.S.
2. Insert point: do search for x -coordinate \Rightarrow insert \Rightarrow walk to root \rightarrow insert in each associate tree.
3. Delete point: ??
4. Range search:

1D Range Search via BST:

a) Code-first way:

RangeSearch ($r \leftarrow \text{root}, x_1, x_2$):

if r not a node:

return {} \Rightarrow No descendants

if $x_1 \leq r.\text{val} \leq x_2$:

$L \leftarrow \text{RangeSearch}(r.\text{left}, x_1, x_2)$

$R \leftarrow \text{RangeSearch}(r.\text{right}, x_1, x_2)$

return $L \cup \{r.\text{val}\} \cup R \Rightarrow$ Inorder traversal \Rightarrow in a BST, ordered set of points!

Current node is
small \Rightarrow go bigger

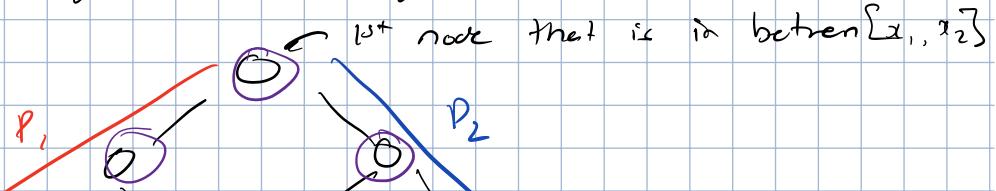
return RangeSearch($r.\text{right}, x_1, x_2$)

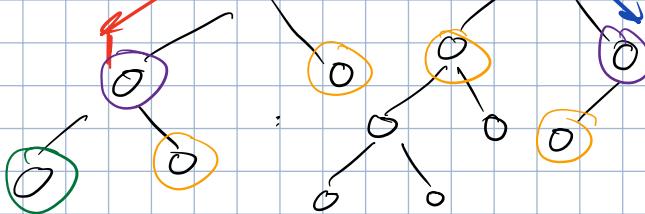
Opposite:
elit $r.\text{val} > x_2$:

return RangeSearch($r.\text{left}, x_1, x_2$)

b) Analysis-friendly way:

Try to create paths P_1 (goes to smallest value in range) and P_2 (largest value in range).





Goal:

- 2 o Boundary nodes: nodes on boundary
- 3 o Outside nodes: nodes that are left of P_1 / right of P_2
- 4 o Inside nodes: nodes that are right of P_1 & left of P_2
 - ↳ Just return all descendants!

Analysis:

Time to find paths: assuming balanced = $O(\log n)$

↳ # of boundary nodes = $O(\log n)$

↳ At very least, have to check all boundary nodes $\Rightarrow \log n$ factor.

↳ # of inside nodes: children of $\log n$ boundary nodes,

so there are $\log n$ inside nodes that we check in $O(1)$

↳ returning descendants = $O(s)$

$T(n) =$ time to get boundary nodes \Rightarrow time to get inside nodes
+ time to get descendants

$$= c_1 \log n + c_2 \log n + s$$

$$= O(\log n + s)$$

2D Range Search:

Also:

1. Do a 1D range search on primary tree $([x_1, x_2])$
 - o Check all boundary nodes
 - o Get inside nodes \Rightarrow range search $([y_1, y_2])$
 - ↳ Optimize: use top most node bc it already has descendants.

Analysis: $O(\log^2 n + s)$

Summary of analysis: let d be dimensions.

1. Space: $O(n(\log n)^{d-1})$

2. Construct tree: $O(n(\log n)^d)$

3. Prez search: $O(s + (\log n)^d)$

} Not good for higher dimensions.

Prez Searching

Quadtrees

- + Simple
- Don't work if skewed
- Space waste

Kd trees

- + Linear space
- + Great as dimensions \uparrow
- + Search time: $O(\sqrt{n} + s)$
- Insertion / deletion

Prez Trees

- + $O(\log^2 n + s)$ rep search
- Waste space.
- Insertion & deletion
- Insertion / deletion

MODULE 9: STRING MATCHING

Introduction:

Goal: search for some pattern $P[0 \dots m-1]$ in text $T[0 \dots n-1]$

↳ Find 1st occurrence of pattern in T (index where pattern starts).

Terminology:

- Substring: $T[i \dots j]$
 - length: $j - i + 1$
- Prefix: substrings with $T[0 \dots j]$
- Suffix: substrings with $T[i \dots n-1]$

Also Terminology:

- Guess: posn i where P might start at $T[i]$
 - ↳ Valid guesser \Leftrightarrow in range $0 \leq i \leq n-m$
↳ Must fit m characters!

$T[n-m+1 \dots n-1]$

$\hookrightarrow n-1 - (n-m+1) + 1$

$\hookrightarrow n-1 - n + m - 1 + 1$

$\hookrightarrow m-1$ (cannot fit P)

- Check: see if $T[i+j] = P[j]$ for guess i , $0 \leq j \leq m$

↳ Need to perform m checks.

Diagrams:

	T	a	b	b	a	.
0		PS03	PS13	PS27		
1			PS07	PS13		
2				PS03		
3					PS13	..
.						

Brute force: check every substring of m in text T

- Time: $\Theta(nm)$

- Worst input: $T = a \dots a a$, $P = a \dots a b$.

↳ Has to guess all $n-m+1$ characters, up to m checks



Karp-Rabin Algorithm

Can we use hashes to make this quicker?

Naive version of hashing K-R:

Ex: // $P = 59265$ $T = 31415926535$

Diagram:

3 1 4 1 5 9 2 6 5 3 5

$h(314159) =$ []

[]

[]

Can hash values equal the hash value of pattern?

Brute with $h(P) = h(\text{substring})$

↳ Collisions! \Rightarrow still check if $P = \text{substring}$

Analysis:

1. Construct pattern hash $\Rightarrow O(1)$

2. Construct hash value for each possible substring:

$O(mn)$

Build string
↓ hash

Repeat for all
possible substrings.

Use then back face:

We have to check all characters when doing hash, no early exit if mismatch

Improvement: compute hash value from prev. in $O(1)$ time.

Ex: // $S_1 = 41592$ $S_2 = 15926$

$41592 \xrightarrow{-4 \cdot 10^5} 1592 \xrightarrow{\times 10} 15920 \xrightarrow{+6} 15926$

$h(15926) = h(h(41592) - h(4 \cdot 10^5)) \cdot 10 + 6$

Rolling hash

Formula: $h_{S_2} = h(h_{S_1} - T[i] \cdot s) \cdot 10 + T[i+m]$

First character of S_1 $h(10^{m-1})$ Next character

Pseudocode:

1. $h(T, P):$

$M \leftarrow \text{prime} \#$

$h_P \leftarrow h(P[0 \dots m-1])$

$h_T \leftarrow h(T[0 \dots m-1]) \Rightarrow O(m)$ time

$s \leftarrow 10^{m-1} \bmod M$

for $i: 0 \rightarrow n-1:$

if $h_T = h_P$: $\xrightarrow{O(m)}$

if $T[i \dots i+m-1] = P:$

return i

if $i < n-m$

$h_T = \text{formula to precompute.}$

Analysis:

Decent: practically fast

Worst case: $\Theta(mn)$

↳ Lots of collisions.

Knuth - Morris - Pratt Also

Terminology:

- i : person. in T that is α companion
 - j : α P β

Basic idea:

if $T[i] \neq P[j]$ and $j = 0$: } Special case of j at beginning. Continue checking
 $i += 1$ } P , no need to increment i

edit $T[j] = P[j]$:
 $i += 1$
 $j += 1$ } Move to next char.

else ($T[i] \neq P[j]$ and $j > 0$):

Use P to find shift to move to

$$\text{Ex: // } P = a b a b a c a$$

Ex: 11 $P = a b a b a | c a$

Mismatch

↓

c a b a b a [a] b a b
 a a b a b a [c]
 |
 a b a b
 a b
 a b

Longest suffix of pattern is prefix

Skipped out on shifts:

Idea: Use the longest suffix of $P[1..j-1]$ s.t. $P[1..j-1]$ is a prefix to shift \Rightarrow update j accordingly ($j = \text{len}(P[1..j-1])$)

Failure array: longest valid suffix of $P[1..j]$ at $F[j]$

↳ If failure at j , check $F\{j-1\}$

Ex:// $P = ababaca$

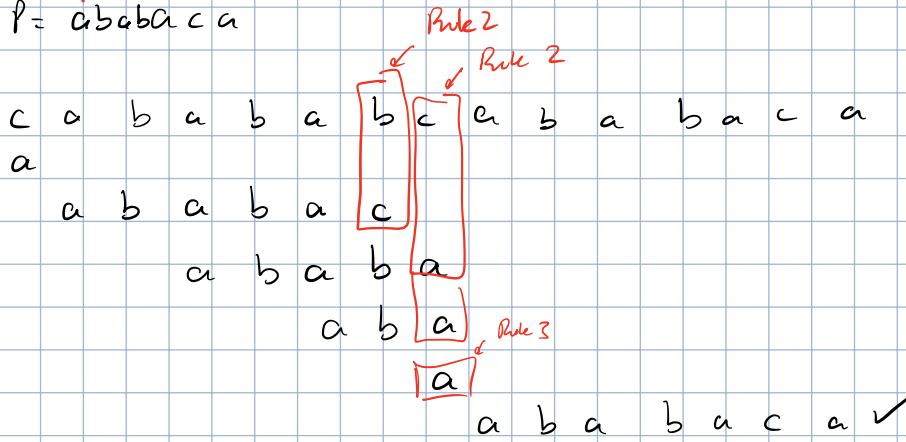
$$F = \begin{bmatrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ 0 & 0 & 1 & 2 & 3 & 0 & 1 \end{bmatrix}$$

$\beta \{1..1\} = b$ $\beta \{1..2\} = b\alpha$ a is also a prefix
 Note α is prefix α is prefix

Naive: $O(m^2)$ to construct F

Find KMP also:

Ex:// $P = ababaca$



$$F = [0, 0, 1, 2, 3, 0, 1]$$

Rule 1: match $\underbrace{j}_{j=1} \xrightarrow{j+1}$

Rule 2: mismatch at $j > 0$

$$\hookrightarrow j = F[j-1]$$

Rule 3: mismatch at $j = 0$

$$j+1$$

Analysis:

$$Time = \# \text{ of horizontal shifts } (i \uparrow) + \text{ vertical shifts } (j \uparrow)$$

$\underbrace{\hspace{10em}}$

$\mathcal{O}(n)$: n characters
| limit

$\underbrace{\hspace{10em}}$

$\mathcal{O}(n)$: n iterations max.

$$\in \mathcal{O}(n)$$

Failure array computation:

Idea: Run KMP on pattern against itself

Building up failure array \Rightarrow use failure array during building (DP).

Notation: $\ell =$ position in pattern that we are checking.

Ex:// Compute failure array of $ababaca$

① Initialize F , $F[0] = 0$

② Initialize $j = 1 \Rightarrow$ no need to process $P[0]$.

$\ell = 0 \Rightarrow$ matching w/ 1st char in pattern.

③ Cases:

a) $P[j] = P[\ell]$: \Rightarrow From $0 \rightarrow \ell$, we have a suffix that matches.

1. Increment ℓ

$$\ell + 1$$

\Rightarrow length of the suffix $\xrightarrow{\ell=2} abaca$

2. Update failure array

$$\ell = 2$$

3

$$F[j] = l$$

3. Continue to next text:

$$j += 1$$

b) If not match & $l > 0$ (not at beginning.)

$$l = F[l-1] \Rightarrow \text{Precomputed already}$$

c) If not match & $l = 0$

$$F[j] = 0 \Rightarrow \text{No suffix to match}$$

$$j += 1$$

This runs in $\Theta(m)$!

Overall: KMP runs in $\Theta(m+n)$ worst case.

Boyer-Moore Algorithm

Fastest algo for string matching.

Idea: reverse search of pattern

Algo:

1. Look @ $m+1$ character in pattern & $T[i+m-1]$
2. If match \Rightarrow look @ $i-1$ char... until 0
3. If no match \Rightarrow bad character heuristic.

If text char not in pattern \rightarrow move $m-1$ steps ahead (check new substrings).

If text char in pattern \Rightarrow shift to last occurrence of shift character in pattern. Check at end of pattern.

Ex:// P=aaron T=acranapple

	a	c	r	a	n	a	p	p	l	e
0				<u>0</u>	n					
1.	a	a	r	o	<u>n</u>				a	a

Derivation:

1. Construct last character occurrence array:

$LC(c) = \text{last index of } 'c' \text{ in } P \quad (c \notin P \Rightarrow LC(c) = -1)$

Iterate in string \Rightarrow put char index in array when iterating.

$\Theta(m)$

2. If mismatch at i in text and j in pattern.

Case #1: $LC(c) > 0$ (not -1) $\& L[c] < j^{\text{old}}$

$j = m - 1 \Rightarrow$ starting @ back again

$i = i + m - 1 - LC(c) \Rightarrow$ text pawn to start comp.

Case #2: $LC(c) > 0 \quad \& \quad L[c] > j^{\text{old}}$

$j = m - 1$

↳ Already passed last char.

$i = i + m - 1 - (j^{\text{old}} - 1)$

Case #3: $L[c] = -1$

Use Case #1

United: $i = i + m - 1 - \min(L(c), j^{\text{old}} - 1)$

Worst case: $O(nm) \Rightarrow$ How string like pattern almost there, but first char incorrect.

↳ Good suffix heuristic: $O(n + m + 1\Sigma 1)$

Suffix Trees

Create trie of suffixes on text

Search pattern \Rightarrow it will occur as prefix in trie!

Space: $O(n \cdot 1\Sigma 1)$

Construction: $O(n \cdot 1\Sigma 1)$ ↳ Searching for character insert

Searching: $O(m)$

Conclusion:

- Good
- Slow for constr. & takes up space.

Suffix Arrays

Also:

1. List all suffixes & put in order of appear.

$T = aaron a \$$

- 0 aaron a \\$
- 1 aaron a \\$
- 2 ron a \\$
- 3 on a \\$
- 4 na \\$
- 5 a \\$
- 6 \\$

2. Sort alphabetically

- 0 \\$
- 1 a \\$
- 2 aaron a \\$
- 3 aron a \\$
- 4 na \\$
- 5 on a \\$
- 6 ron a \\$

3. Suffix array: map sorted posn \rightarrow appearance ord in 1

$$A^S : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline 6 & 5 & 0 & 1 & 4 & 3 & 2 \\ \hline \end{array}$$

4. Matching of pattern:

Binary search on suffix array!

search (A^S, P) :

$$l = 0, r = n-1$$

while $l < r$

$$v \leftarrow \lfloor \frac{l+r}{2} \rfloor$$

$$i \leftarrow A^S[v] \quad | \quad \text{Same pattern size}$$

$$s \leftarrow \text{strcmp}(T(i \dots i+m-1), P)$$

if $s < 0$:

$$l \leftarrow v + 1$$

else if $s > 0$:

$$r \leftarrow v - 1$$

etc:
return four of λ 's

Analysis:

Construction:

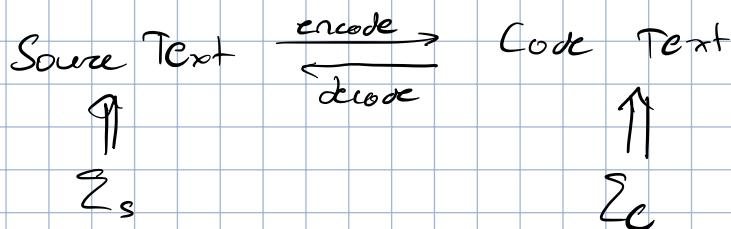
MSD radix sort $\Rightarrow \Theta(n^2)$ worst case, but can improve to $\Theta(n \log n)$

Search:

$\Theta(\log n)$ comp, $\Theta(m)$ string comp $\Rightarrow \Theta(m \log n)$

MODULE 10: COMPRESSION

Encoding & Decoding



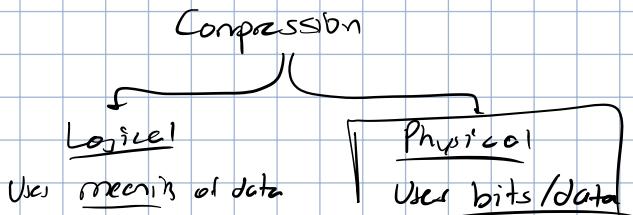
Texts are stored as streams \Rightarrow encode one char at a time.

Analysis: compression ratio:

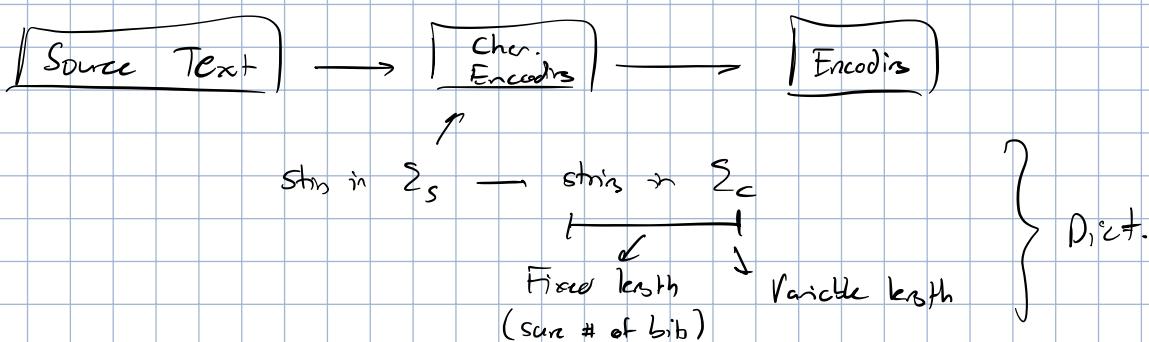
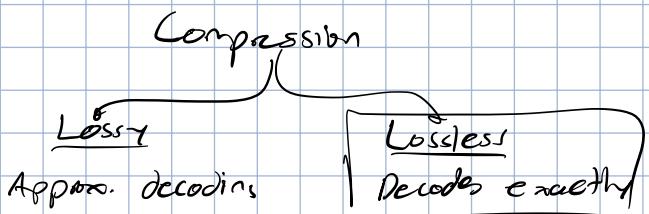
$$\frac{|C| \cdot \log |\Sigma_c|}{|S| \cdot \log |\Sigma_s|} = \text{Good compression} \Rightarrow \text{ratio} < 1$$

Types:

① Logical vs. physical



② Lossy vs. lossless



Encoding \rightarrow Char. Encoding \rightarrow Source

Could be diff. rep (diff D.S) of encoders
Prefix free
trie

Analysis:

• Encode: $O(|T|)$

• Decode: $O(|C| + |\Sigma_S|)$ \Rightarrow if imp? as trie

Huffman Encoding

Guarantees to get shortest strings if using binary

Also:

① Calculate character frequencies

↳ Need access to entire stream

② Put character & frequency as trie nodes

③ Join 2 least frequent nodes into tries. Tie \Rightarrow alphabetical ordering.

↳ 0 for left child, 1 for right child

↳ New node has freq. = sum of left + right freq.

④ Repeat until all joined.

Ex://

Text: G: 2 E: 4 Y: 1
R: 2 N: 2

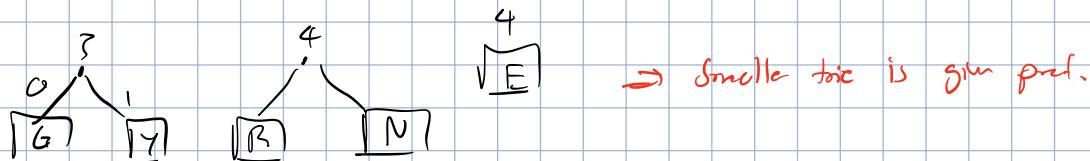
①



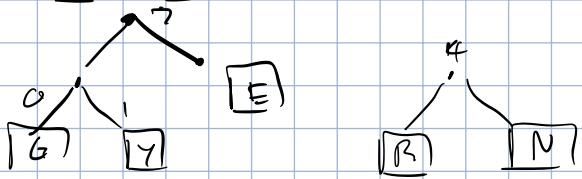
②



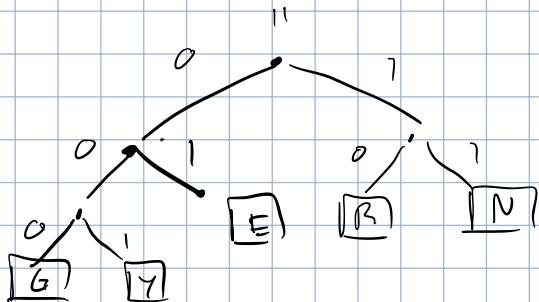
③



④



(S)



Use min-heap to find min. freq nodes. $\Rightarrow O(|\Sigma_s| \log(|\Sigma_s|))$

The actual tree encoding $\Rightarrow O(|\Sigma_s| + c)$

Total encoding time: $O(|\Sigma_s| \log(|\Sigma_s|) + c)$

|| decoders ||: $O(|c|)$

Notes:

- ① Trc is not unique
- ② Trc must be encoded for decoding
- ③ Doesn't use sturm \Rightarrow needs entire text (reading order $\times 2$)

Run Length Encoding

Uses multi-character encoding: multiple source chars \rightarrow 1 codeword

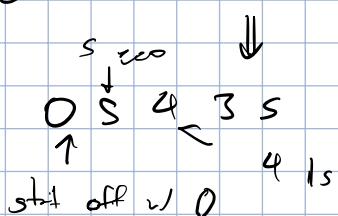
Assumption: $\Sigma_s = \{0, 1\}$

Strength: multiple same chars in string $(0\ldots 0, 1\ldots 1\ldots)$

Name:

1. First bit indicates bit that we start off with
2. Following #'s tell us sequence of bit

Ex: // 0 0 0 0 0 1 1 1 0 0 0 1 1 1



Elias Gamma

To encode run length $k \Rightarrow \lfloor \log k \rfloor$ copies of 0s + binary rep. of k

To decode: take # of 0s + \Rightarrow # of bits to read next for repeat of bit.

Ex:// 1111110010000001. . . 0 1. — 1

① Start by putting 1st bit in encoding.

$$C = 1$$

② For $k = 7$

③ Encode:

$$\lfloor \log 7 \rfloor = 2$$

$$C = 1 + 00111$$

④ $k = 2$

$$\hookrightarrow C = 1 + 00111 + 010$$

⑤ $k = 1$

$$C = 1 + 00111 + 010 + 1$$

.

.

.

Ex:// Decode RLE 0 0001101001001010

① Read initial bit

0

② Read as many 0s

3 0s \Rightarrow next 3+1 bits for decode

③ Read 3+1 bits.

$$1101 = 13$$

Run length of 13 for 0 \Rightarrow switch to 1 for next decode.

$$S = 0 \dots 0$$

$\overbrace{1 \dots 1}^{13}$

④ Read as many 0s

2 0s \Rightarrow next 2+1 bits for runlen+2

⑤ Read 2+1 bits

$100 = 4 \Rightarrow$ switch to 0 for next decod.

$$S = 0 \dots 0 \quad 1 \dots 1$$

$\overbrace{1 \dots 0}^1 \quad \overbrace{1 \dots 1}^0$

Analysis:

◦ Encode: $(|S| + |C|)$

◦ Decode: $(|S| + |C|)$

Notes:

- It can fail \Rightarrow we don't have enough bits when looking for $k+1$ bits in decoding.
- Depends on input: if no repeats, doesn't work well
 - ↳ If $k < 6 \Rightarrow$ no compression
 - ↳ $k=2, k=4 \Rightarrow$ expansion.

LZW

Major idea: encode substrings \Rightarrow if repeated, we already have encoding!

Really good if you have repeating substrings

This allows dictionary to be created during reading/decoding

↳ No need to send dictionary!

Encoding Also

1. Initialize D_0 for Σ_s

$D_0: \{ \text{char: } 0, \text{char: } 1, \dots, \text{char: } |\Sigma_s| - 1 \}$

2. At each iteration i

a) Pop off char i

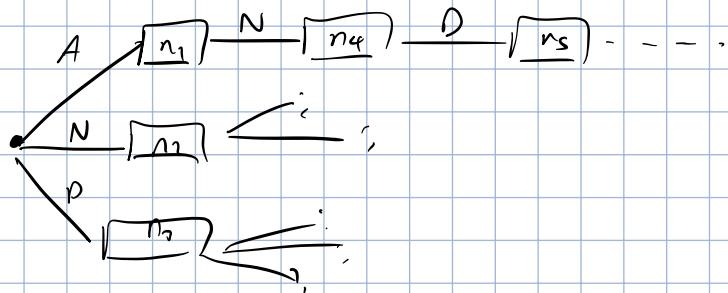
b) Find longest substring starting at i & already in dict

c) $D_{i+1} = D_i - \text{append}(\{ T_i \dots T_{i+k}, \text{next avail code} \})$

$\overbrace{\text{Substrings that}}^{\text{just beyond } D_i}$

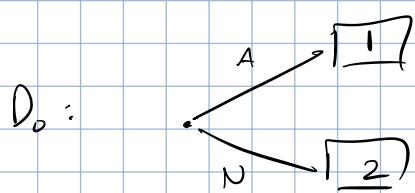
Dict implementation for encode? trie

- Ex: $\Sigma_s = \{A, N, D\}$



Ex: $\Sigma_s = \{A, N\}$. Encode: ANANANA NNA

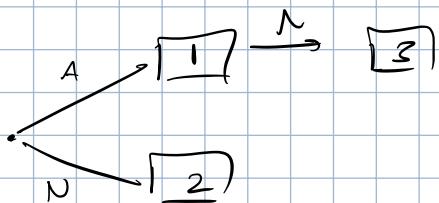
- ① Initialize:



- ② Consume as much chars as we can.

$A \rightarrow \text{code: 1}$

Add AN $\rightarrow \text{code: 2}$

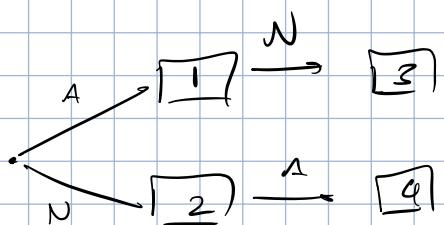


String to encode
 $O(|s|)$

② $N \rightarrow A^x$

$\hookrightarrow \text{code: 1. append}(2)$

$\hookrightarrow \text{Dict: } D_1. \text{append}('A', 2)$



③ $A \rightarrow N \rightarrow A^x$

$\hookrightarrow \text{code. append}(12). \text{append}(3)$

$\hookrightarrow \text{Dict: } D_2. \text{append}('ANA', 8)$

Decoding Alg

Create dict while encoding, but will be behind by 1 char.

Ex:// $C = 65, 78, 128, 180, 78, 129$

① Initialize D_0 :

We can use a reg. Dict.

$$D_0 = \{ \text{code word} \rightarrow \text{string} \}$$
$$\begin{matrix} 1 & & 1 \\ 0 - 65 & & \hookrightarrow \dots - 2 \end{matrix}$$

② Consur:

a) 65

$$D_0(65) = A$$

$$S = A$$

$$S_prev = A$$

b) 78

$$D_0(78) = N$$

$$S = AN$$

$D_0 \cdot \text{add}(\text{next_char_code}, S_prev + S[0])$

$\hookrightarrow D_0 \cdot \text{add}(128, AN)$

$$S_prev = N$$

c) 128

$$D_0(128) = AN$$

$$S = ANAN$$

$D_0 \cdot \text{add}(129, ANA)$

$$S_prev = AN$$

d) 180

180 not in dict

$S_append(S_prev + S_prev[0])$

$\hookrightarrow (ANA)$

:

:

Time: $O(|s|)$

bzip 2

Uses transforms \Rightarrow makes text easier to decode



MTF Transform

Store dict in array: $\{c, b, c, d, \dots\}$

\leftarrow code = index of char

Each time it's used \rightarrow move character to front

\hookrightarrow Repeated characters \Rightarrow whole bunch of 0s

Encoding

1. Get text in array
2. Put index in string
3. Move to front

We get skewed distribution



Really nice for Huffman.

Decoding

1. start w/ same dict
2. Do exact same ops.

Ex. 11 MISSISSIPPI

$D_0: \{a: 0, b: 1, \dots\} \quad 2: 2\bar{2}$

① Consider M:

$$M = 13 \Rightarrow D = 13, \dots$$

Move M to front: $\{M, A, B, \dots\}$

② I:

$$I = 9 \Rightarrow D = 13, 9$$

$\{I, M, A, \dots\}$

③ S:

$$S = 18 \Rightarrow D = 13, 9, 18$$

$\{S, M, A, \dots\}$

④ S:

$$S = 0 \Rightarrow D = 13, 9, 18, 0$$

$\{S, M, A, B, \dots\}$

Burrows Wheeler Transform.

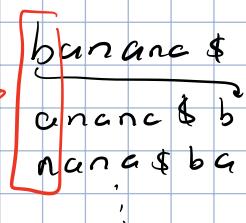
Obj: source text \rightarrow coded text (contain exact same letters, in diff. order)

Assumption: EOS shr \$

Use cyclic shifts: banana \$
↓
a \$ banana

Encoding

① Create an array of all cyclic shifts.

Same as source text!
→ 
!:

② Sort array alphabetically

↳ This is what makes this a useful for MTF function \Rightarrow same chars. will line up.

↳ Efficient: suffix sorting (MSD sort in $O(n \log n)$)

③ Use last column for encoding.

$$O(n \log n + n) \in O(n \log n)$$

$\xrightarrow{\text{Sort}}$ $\xrightarrow{\text{Cyclic, last char}}$

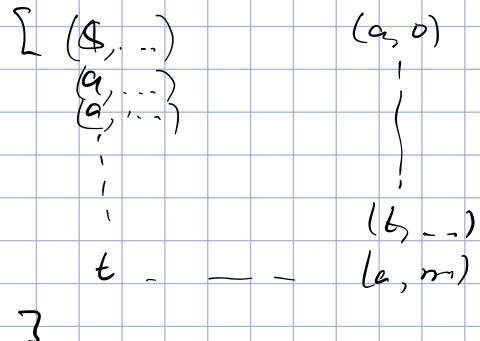
Decoding

① Create array w/ last character = coded character.

C = \$aaaaaaefffellst

[- - - - (a, 0)
(f, 1)
(f, 2)
:
(e, n-1)
(a, n-1)

② Generate first char for each string by taking last col \rightarrow sort strings.



③ Find row w/ last character as \$ (row \rightarrow strings).

④ Take 1st char \rightarrow decode it.

⑤ Take 1st char row number \rightarrow find row # of sub string & get first char

⑥ Repeat until we get \$.

Ex:// C = and \$ r c a a a a b b

① Array at end:

end = A: 0 1 2 3 4 5 6 7 8 9 10 11
a, 0 | r, 1 | d, 2 | \$, 3 | r, 4 | c, 5 | a, 6 | a, 7 | a, 8 | a, 9 | b, 10 | b, 11

② Sort:

A = 0 1 2 3 4 5 6 7 8 9 10 11
\$, 3 | a, 0 | a, 6 | a, 7 | a, 8 | a, 9 | b, 10 | b, 11 | c, 5 | d, 2 | r, 1 | r, 4

③ Look @ end to find \$ \Rightarrow index 3 in end.

④ $A[3] = (a, b) \Rightarrow S = ab$
index = 6

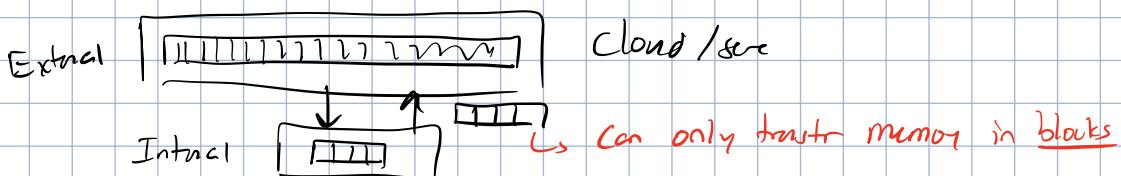
⑤ $A[\text{index}] = (b, 10) \Rightarrow S = ab$
index = 10

⑥ $A[\text{index}] = (r, 1) \Rightarrow S = abr$
index = 1

Time complexity: $O(n + 1 \sum_s 1)$

MODULE 11: EXTERNAL MEMORY

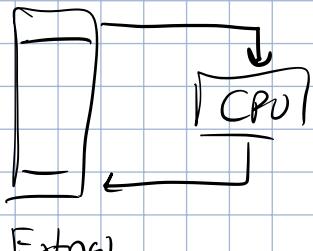
Motivation



Block transfers are \$\$ \Rightarrow \text{minimize block transfers.}

Block transfers dominate running time \Leftarrow all analysis on block transfers.

Stream Based Algorithms



External memory: stream blocks

Efficient

\hookrightarrow Input: $O\left(\frac{n_{\text{input}}}{B}\right)$

\hookrightarrow block size

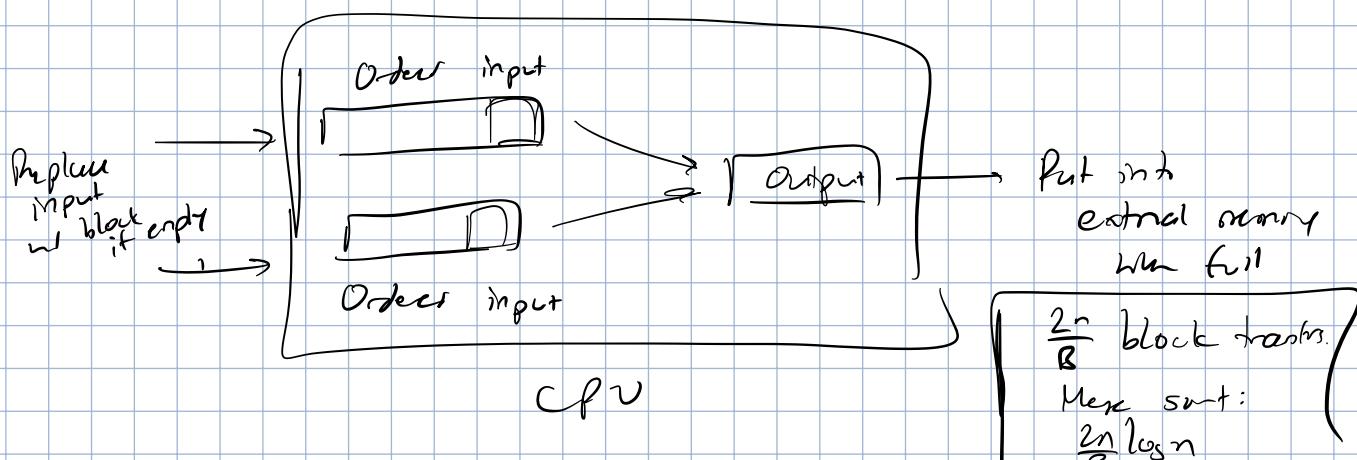
Output: $O\left(\frac{s_{\text{output}}}{B}\right)$

Any stream-based algo cannot do better

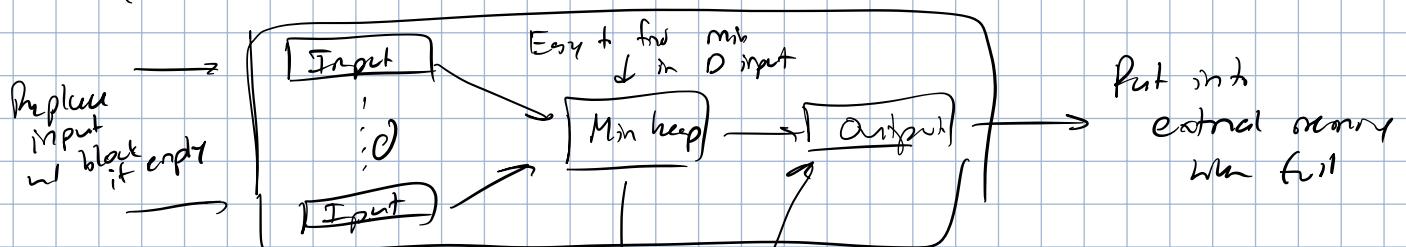
\hookrightarrow Text match & compression.

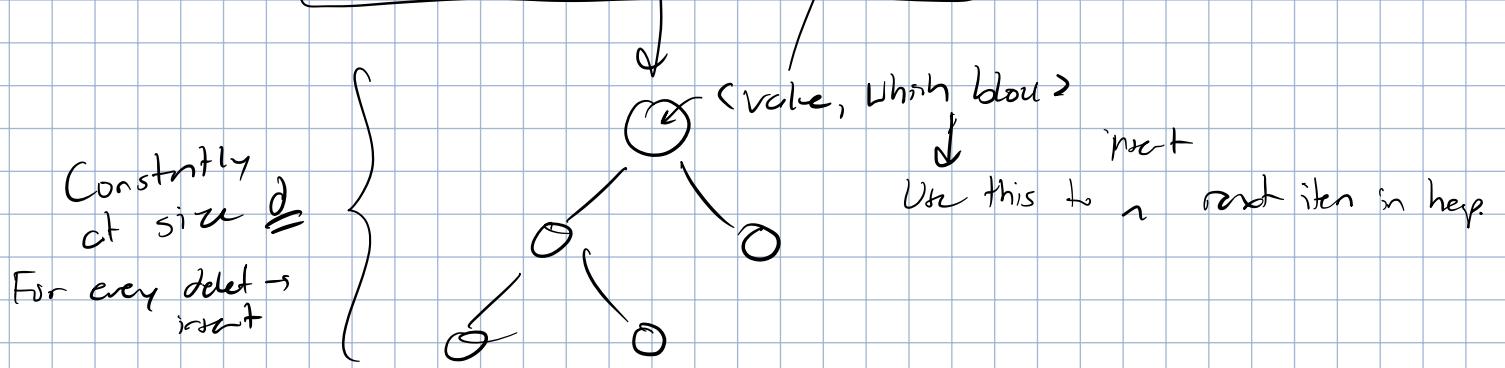
External Sorting

Merging is a pretty good candidate:



D-way merging:





Principle:

$\Theta(n \log d)$ Just looks at RAM
Taking n items \uparrow \downarrow $\log d$ deletion/insertion

of block transfers:

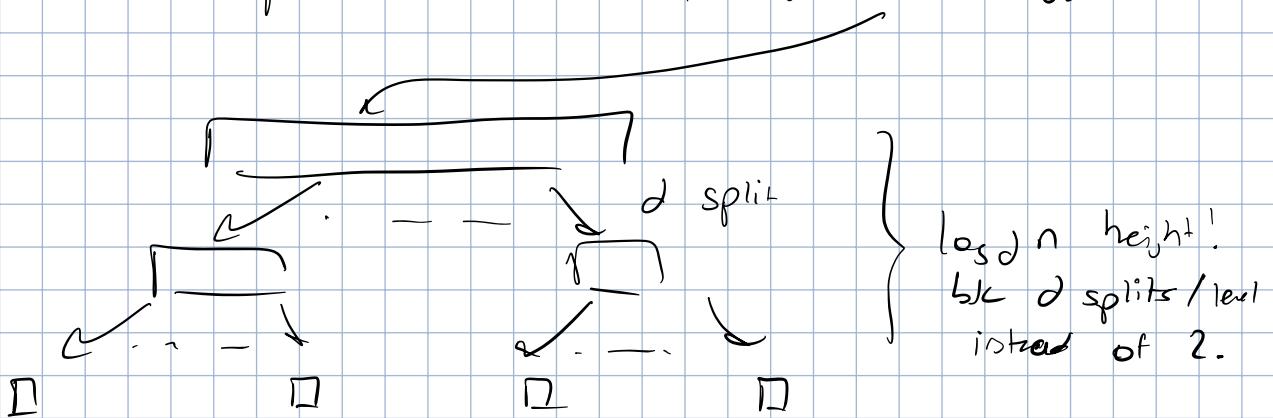
We can only put $d+1$ blocks into internal memory
input \uparrow \downarrow Output

1 block transfer $\Rightarrow \Theta(n/B)$
↳ stayed same!

d -way merge sort

Alg.:

1. Split input into d arrays
 2. Merge each d puts together \Rightarrow d -way merge
- $\log_d n$ height!
 $O(\log_d n)$



Principle:

$$\log_d n \cdot n \log d = n \frac{\log n}{\log d} \in \Theta(\log n)$$

depth merge

Has same runtime as merge sort, but far better for external.

of block transfers = recursion depth \times # of block transfers/level

$$\leq \log_B n \times \frac{2n}{B}$$

We can improve:

Stop at block level $M \Rightarrow M$ bits can fit into array. No need to swap.

\therefore # of block transfers = recursion depth \times # of block trans. / level

$$= [\log_B n - \log_B M] \times \frac{2n}{B}$$

$$= \log_B \frac{n}{M} \times \frac{2n}{B}$$

$$\vdots$$

$$\vdots$$

$$\in \Theta\left(\frac{\log(n/M)}{\log(M/B)} \times \frac{n}{B}\right)$$

\Rightarrow Really fast
Minimizing transfers.

This is optimal. Cannot be better.

2-4 Trees

Motivation: AVL tree has poor memory locality \Rightarrow # of block transfers \uparrow

2-4 tree is good for internal memory

Idea:



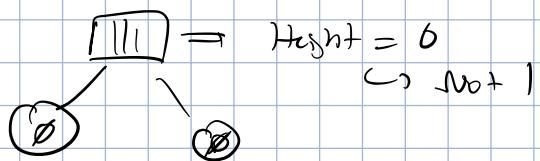
Properties:

- Each node is either 1-node, 2-node, 3-node
 - 1-node: 1 key, 2 subtrees
 - 2-node: 2 keys, 3 subtrees
 - 3-node: 3 keys, 4 subtrees
- All empty trees are in same level
- Order property: BST property w/ max tree.

$$\langle T_0, n_1, T_1, n_2, \dots, T_3 \rangle$$

$\xrightarrow{\text{in order of keys}}$

Note: not counting empty trees for height



Operations :

① Search (k): compare k among node keys \Rightarrow BST

o If k not in tree \Rightarrow return parent of stopping empty node.

② Insert (k):

1. Search \rightarrow get leaf node (non-empty)

2. Insert KVP into tree & create another subtree.

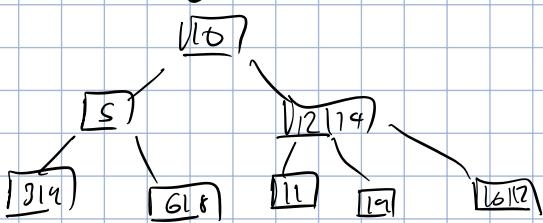
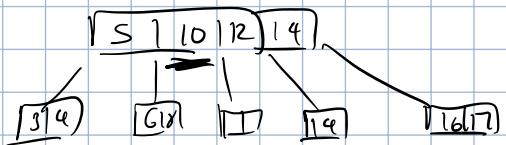
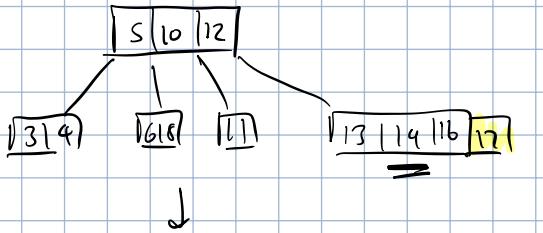
3. Split up to restore 1-3 node property.

a) Find middle node \rightarrow push to parent

b) Split up node into 2 subtrees.

4. Repeat until root is balanced. (new root if necessary)

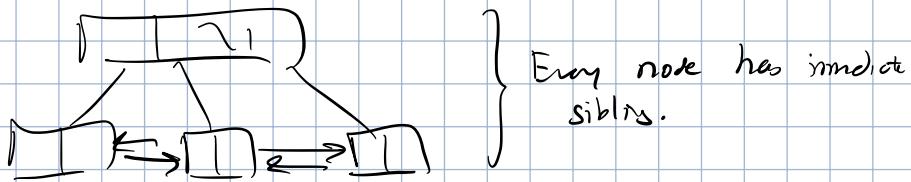
Ex. //



③ Delete (k):

Terminologies.

1. Immediate sibling: closest left/right node



2. Inorder successor: smallest key in right child. (always a leaf).

Also:

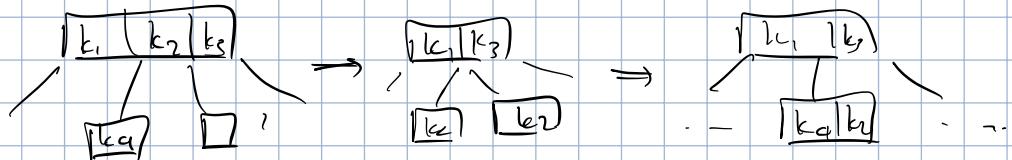
1. Search for key
2. If node has non-empty children \rightarrow swap w/ inorder successor
 \Downarrow Node is a leaf
3. Delete key & empty subtree
4. If undel. (< 1 -node in tree):

Case #1: If immediate sibling is "rich" (2/3-node):

- a) Transfer key to parent (closest leaf)
- b) Parent will transfer appropriate key to leaf \Rightarrow no longer empty
- c) Done \Rightarrow no more repair needed.

Case #2: If both siblings poor

- a) Take in-between key in parent \rightarrow move in leaf w/ immediate sibling.
- b) Repeat until parent not empty
 \hookrightarrow If root becomes empty \Rightarrow delete it



Analysis:

height $\in O(\log n)$

\hookrightarrow Ops: $O(\log n)$

Not better than AVL tree, can generalize.

a, b - Tree

Properties:

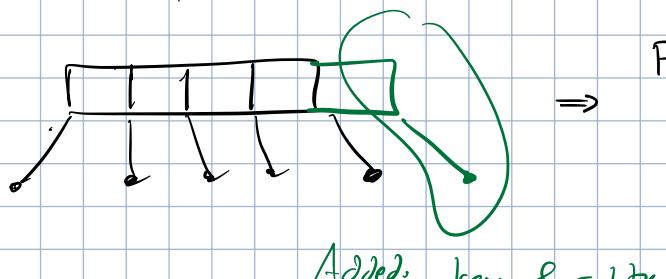
- * o a: minimum # of subtrees / node
 - * o b: maximum
 - o Empty leaves @ same level
 - o Order property
 - * o $a \leq \lceil \frac{b}{2} \rceil = \lfloor \frac{b+1}{2} \rfloor$
- } Unless $\text{root} = \text{cen}$ then min. of 2 min.
} a doesn't apply to root
↓
- Used if storing very few keys, (i.e. storing 1 key, but $a=3$)

Operations: Exact Search

↳ Clarity insertion:

Too many keys in a node? Whenever # of subtrees $> b$

Ex: // (3, 5) tree



For this node:

6 subtrees! $\Rightarrow > b$

Split \Rightarrow Splitting into $\lfloor \frac{b+1}{2} \rfloor$ subtrees / split (restores property).

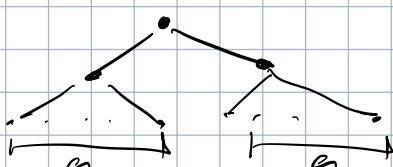
Why do we even need $a \leq \lceil \frac{b}{2} \rceil$?

If we had (4, 5)-tree \Rightarrow overflow \Rightarrow split will create (2, 5) subtree which violates (4, 5)-tree!

Analysis:

Q: What is height of (a, b) tree?

Find smallest # of KVPs \Rightarrow height (will give us O.l.)



$$\min \# \text{ of KVPs} = 1 + 2 + 2a + 2a^2 + \dots$$

Slippery root \Rightarrow Root

$$= 1 + \sum_{i=1}^{h-1} 2a^i (a-1)$$

$$= 2a^h - 1$$

∴ If n KVPs stored:

$$n \geq 2a^h - 1$$

$$h \in O(\log_a n) \in O(\log n)$$

INTERESTING NOTE:

$$\# \text{ of KVPs in } (a, b) \text{ tree} = \# \text{ of empty subtrees} - 1$$

Super useful in proofs.

Q: What is the cost of ops?

$$\text{Height} \in O(\log_a n)$$

↳ Maximize $a \rightarrow$ at min. height. $a = \lceil \frac{b}{2} \rceil$, $a = \underbrace{b/2}_{\text{Max.}}$

To actually find keys in node \Rightarrow AVL tree

↳ # of KVPs in node $\in O(b)$

$$\therefore \text{Search in AVL tree} = O(\log b)$$

In total:

$$O(\text{find the node} \cdot \text{find a leaf}) = O(\log_a n \cdot \log b) \\ \in O(\log n)$$

Not faster than AVL, but useful for memory.

B Tree

Idea: specialize (a, b) tree s.t. it can do well in external mem.

Want to minimize height

↳ $a = \frac{b}{2} \Rightarrow b$ dictates the entire tree (root)

What should b be? $b \in \Theta(\text{block size})$

Choose b s.t. node w/ $b-1$ leaves can fit in 1 block

$\therefore 1$ block $\rightarrow 1$ node to search in CPU

Why is this nice:

of block trans = height of search

$$= \Theta(\log_b n) \Rightarrow \text{SuperEfficient.}$$