

CS 240 - S23

Data Structures and Data Management

Full Course Notes

Written from Lecture Slides

I wrote notes in class before midterms, and used lecture slides afterwards. Hope these help a little!

[Josiah Plett](#)

① Initial Definitions + Overview

Problems (input + output)

Sorting: given numbers, put them in order.

structured search: given data with keys, search by key.

unstructured search: given text search by string.

instance: one particular input to the problem.

algorithm: step-by-step process to mod input

recursive: algorithm that uses itself. not same as program, a specific implementation

solving: finishes in finite time, returning correct answer.

run time: the number of computing steps.

n: the size of the input (use $\text{size}(I)$)

worst case: even for worst possible

big O: ignore constant factors in the runtime bound

③ Algebra, Techniques

Transitivity: $f(n) \in O(g(n)) \wedge g(n) \in O(h(n)) \rightarrow f(n) \in O(h(n))$

Maximums: $f(n) + g(n) \in O(\max\{f(n), g(n)\})$ (if $f(n) > 0, g(n) > 0$)

Techniques

Limit

Suppose $f(n) > 0$ & $g(n) > 0, n \geq n_0$.

Assume $L = \lim_{n \rightarrow \infty} \frac{f(n)}{g(n)}$ exists.

Then:

$$f(n) \in \begin{cases} o(g(n)) & \text{if } L = 0 \\ \Theta(g(n)) & \text{if } 0 < L < \infty \\ \omega(g(n)) & \text{if } L = \infty \end{cases}$$

Polynomials

If $f(n)$ is polynomial of degree $d \geq 0$, then $f(n) \in \Theta(n^d)$.

Sine

Just use first principles I guess.

Growth Rates

⊖ If $f(n) \in \Theta(g(n))$, then growth rates of $f(n)$ & $g(n)$ are same!

- ⊖ If $f(n) \in o(g(n))$, then growth rate of $f(n)$ is less than $g(n)$.
- ⊖ If $f(n) \in \omega(g(n))$, then growth rate of $f(n)$ is greater than $g(n)$.

Relationships

$f(n) \in \Theta(g) \iff g(n) \in \Theta(f)$ etc.

- $\Theta(1)$ - Constant
- $\Theta(\log n)$ - Logarithmic
- $\Theta(n)$ - Linear
- $\Theta(n \log n)$ - Linearithmic
- $\Theta(n \log^k n)$ - Quasi-linear
- $\Theta(n^2)$ - Quadratic
- $\Theta(n^3)$ - Cubic
- $\Theta(2^n)$ - Exponential

④ Algorithm Complexity

(best)
 $T_A(n) = \max_{I: \text{size}(I)=n} \{T_A(I)\}$

(worst case)
 $T_A(n)^{\text{avg}} = \frac{1}{|\{I: \text{size}(I)=n\}|} \sum_{I: \text{size}(I)=n} T_A(I)$ (average case)

Recurrence Relations

Merge sort example:

$$T(n) = \begin{cases} T(\lfloor \frac{n}{2} \rfloor) + T(\lfloor \frac{n}{2} \rfloor) + \Theta(n) & \text{if } n > 1 \\ \Theta(1) & \text{if } n = 1 \end{cases}$$

exact recurrence: constants, not orders

sloppy recurrence: floor/ceiling removed

$$n = 2^k \text{ for analysis}$$

② Asymptotic Analysis

Algorithms will be presented using pseudocode and analyzed using order notation.

- For a problem Π , we can have several algorithms.
- For an algorithm A solving Π , several programs.

Solving Problems

- ① Design algorithm A that solves Π . (alg design)
- ② Assess correctness and efficiency of A . (alg analysis)
- ③ If acceptable, implement A . (write program)

- Efficiency
- ① Running time
 - ② Auxiliary space
 - ③ Size(I_{input})
- Time
- A # of primitive operations
 - B Growth rate (complexity)
 - C idealized computer
 - D RAM \rightarrow constant time

Order Notation

★ $f(n) \in O(g(n))$ if there exists $c > 0$ and $n_0 > 0$ s.t. $|f(n)| \leq c|g(n)|$ for all $n \geq n_0$.

- ① Build g out of $n \geq ? \rightarrow \# \leq ? \rightarrow \# \leq g(n)$
- ② Build f out of $\# \leq ? \rightarrow f \leq \#$
- ③ Conclude $f \leq \# \leq g(n)$ for $n_0 = ?$

Ω -Notation

(asymptotic lower bound)

★ $f(n) \in \Omega(g(n))$ if there exists $c > 0$ and $n_0 \geq 0$ s.t. $c|g(n)| \leq |f(n)|$ for all $n \geq n_0$.

Θ -Notation

(asymptotic tight bound)

★ $c_1|g(n)| \leq |f(n)| \leq c_2|g(n)|$ for all $n \geq n_0$.

o -Notation

(asymptotically strictly smaller)

★ if for all $c > 0$, there exists $n_0 \geq 0$ s.t. $|f(n)| \leq c|g(n)|$ for all $n \geq n_0$.

- RELATIONS
- $T(n) = T(n/2) + \Theta(1) \rightarrow T(n) \in \Theta(\log n)$
 - $T(n) = 2T(n/2) + \Theta(n) \rightarrow T(n) \in \Theta(n \log n)$
 - $T(n) = 2T(n/2) + \Theta(\log n) \rightarrow T(n) \in \Theta(n)$
 - $T(n) = T(cn) + \Theta(n) \rightarrow \text{OCCUR} \rightarrow T(n) \in \Theta(n)$
 - $T(n) = 2T(n/4) + \Theta(1) \rightarrow T(n) \in \Theta(\sqrt{n})$
 - $T(n) = T(\sqrt{n}) + \Theta(\sqrt{n}) \rightarrow T(n) \in \Theta(\sqrt{n})$
 - $T(n) = T(\sqrt{n}) + \Theta(1) \rightarrow T(n) \in \Theta(\log \log n)$

4) Cont'd: Priority Queues

ADT: describes information and operations (interface)
 realization: specifies data structure and algorithms
 data structure: how information is stored
 algorithms: how the operations are performed

STACK

- push()
- pop()
- size()
- isEmpty()
- top()
- LIFO

QUEUE

- enqueue()
- dequeue()
- size()
- isEmpty()
- front()
- FIFO

PRIORITY QUEUE

- insert() ← tagged w/ priority
- deleteMax()
- "priority" AKA "key"

Empty tree height: -1

5) Priority Queue Realization

Unsorted Arrays: insert: $O(1)$
 PQ-sort → selection(n^2) deleteMax: $O(n)$

Sorted Arrays: insert: $O(n)$
 PQ-sort → insertion(n^2) deleteMax: $O(1)$

Binary Heap

Heap: ① All levels of heap are full, except last is left-justified.

② For any node i , the key of the parent j is $j \geq i$.

("max-oriented binary heap")

Heap height: $O(\log n)$

heap.insert: $O(\log n)$

heap.deleteMax: $O(\log n)$

PQ-sort → $(n \log n)$

Heaps in Arrays

• Store root in $A[0]$, and read left → right then top → down

Left Child of i is $2i+1$

Right Child of i is $2i+2$

Parent of i is $\lfloor \frac{i-1}{2} \rfloor$

6) Heapify, Heapsort

Heapify

$O(n)$

Heapsort

- swap with last unused.
- fix down.
- To sort in literal order.

① Fix parent relative to children

bottom-to-top, left-to-right

② call fixdown

7 Runtime Analysis (Avg) & Randomized

Complexity

PQ-sort w/ binary heaps $O(n \log n)$

eg: Given n items (A), build heap.

Using fixups, we get $O(n \log n)$

Using fixdowns, we get $O(n)$

8 Avg Runtime

(Randomized Algs)

EG

$$T^{avg}(n) = \frac{1}{n!} \sum_{\pi \in \Pi_n} T(\pi) = \frac{1}{n!} \left(\sum_{\pi \in \Pi_n: \pi \text{ good}} T(\pi) + \sum_{\pi \in \Pi_n: \pi \text{ bad}} T(\pi) \right)$$

$A = [12, 8, 0, 10, 4]$
 $\pi = \langle 2, 4, 1, 3, 0 \rangle$
 $\pi^i = \langle 4, 2, 0, 3, 1 \rangle$

Only Permutations Matter

Average Runtime Strats

- 1 Apply the def
- 2 Break up the summation
- 3 Insert (recursive) e.v. formulae
- 4 Simplify the sums

Avg Runtime

$$T^{avg}(n) = \frac{\sum_{I: |I|=n} T(I)}{\# \text{ of size of } n} = \frac{\sum_{I \in I_n} T(I)}{|I_n|}$$

$T(I, R)$ is the runtime of a randomized algorithm A for an instance I and the sequence of random choices R .

Expected Running Time for I "probability"

$$T^{exp}(I) = E[T(I, R)] = \sum T(I, R) \cdot P(R)$$

Expected Running Time for A

$$T^{exp}(n) := \max_{I \in I_n} T^{exp}(I)$$

Module 3 (PDF) Slides 26-35

QuickSelect

choose-pivot(A): return index p in A (semi-hard-coded)

partition(A, p): rearrange A and return pivot index i s.t.

- ▶ the pivot-value v is in $A[i]$
- ▶ all items left of i are $\leq v$
- ▶ all items right of i are $\geq v$

$A \leq v \quad v \quad \geq v$

9 Randomized

QuickSelect Random

choose-pivot(A): random pivot

$T^{exp}(A) \in O(n)$

QuickSort (randomized)

$T^{worst}(n) \in \Omega(n^2)$

$T^{best}(n) \in O(n \log n)$

$T^{avg}(n) \in O(n \log n)$

$T^{exp}_{(rand)}(n) \in O(n \log n)$

Best Sort in practice

- Randomize the pivot
- While loop, not recursion
- Stop recursing when $n \leq 10$. Run InsertionSort then!
- If duplicates: $\leq v = v = \geq v$
- Pass range of indices, not A .
- Keep explicit stack (no recursion!)

CS 240 (2)

(10)(11) Radix Sort

radix-R: $\{0, \dots, R-1\}$ are all the possible values in input.

Stable: equal keys return in the same order.

Bucket Sort

Everyone's Mom Sort
Time: $O(n+R)$
Space: $O(n+R)$

Count Sort

Bucket sort but counting # of items in each bucket instead of moving.

LSD-Radix-Sort

1. for $d \leftarrow m$ down to 1
do bucket-sort(A, d)

Time: $O(m(n+R))$

Space: $O(n+R)$

MSD-Radix Sort

Group by first digit, then next, recursively.

Time: $O(Rn)$ ^{stack}
Space: $O(n+R) + O(m)$
of digits \leftarrow

MSDRS($A, n, l=r, r=n-1, d=1$)

- l, r , range $A[l, \dots, r]$ to sort
- d is digit we are sorting
- if $(l < r)$ {
 - bucket-sort(A, n, l, r, d)
 - if $(d < m)$ {
 - // find subarrays w/ same d^{th} digit and recurse
 - int $L = l$
 - while $(L < r)$ {
 - int $q = L$
 - while $(q < r \ \& \ A[q+1] = A[L])$ {
 - $q++$
 - MSDRS($A, n, L, q, d+1$)
 - $L = q+1$

Biased Search Requests

key	A	B	C	D	E
access-frequency	2	8	1	10	5
access-probability	$\frac{2}{26}$	$\frac{8}{26}$	$\frac{1}{26}$	$\frac{10}{26}$	$\frac{5}{26}$

Dynamic: move forward if recently accessed?

Expected Access-Cost:

Expected # of comparisons needed to find random key.

"competitive analysis" is online (vs) offline.

Self-Adjusting Lists/Arrays

Move-to-front heuristic (lists): exactly what you think.

Move-to-front heuristic (arrays): same thing, but make the "front" of your unsorted array the end (where you insert).

Transpose heuristic: if accessed, move up by 1.

Special-Key Dictionaries

Lemma 6.2: If the keys are uniformly distributed, interpolation search is $O(\log \log n)$ in average case

more on next page

(12)(13)(14)(15) Stuff

BSTs.

AVLs.

insert delete

restructure

Randomized BST

Expected height: $O(\log n)$
as long as no delete();

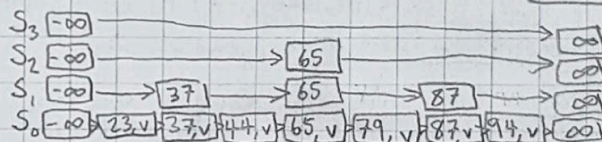
At this point I totally stopped going to lectures :)

Skip Lists

(aka "levels" "layers")

DEF: A hierarchy S of ordered lists S_i s.t.:

- S_i contains keys $-\infty$ & ∞ , "sentinels"
- S_0 contains key-value pairs of S in non-decreasing order.
- Each list is subsequence of prev: $S_0 \supseteq S_1 \supseteq \dots \supseteq S_k$
- S_k contains only the sentinels.



getPredecessors(k)

eg to get 79, we do: result: $\{s_0, s_1, s_2\}$
 $s_3[0] \rightarrow s_2[0] \rightarrow s_1[1] \rightarrow s_0[2] \rightarrow s_0[3]$
 $79 = s_0[5] \leftarrow s_0[4]$

search(k)

basically getPredecessors(k), then check \exists .

insert(k, v)

insert, then randomly choose ^{stack} height.

delete(k)

remove it from all S_i (after finding it)

Interpolation Search: faster than binary search when values are numerical.
Guess the approximate location! Instead of (for $k=100$) $\frac{110}{40} > 75$, do $\frac{110}{40} \sim 100$

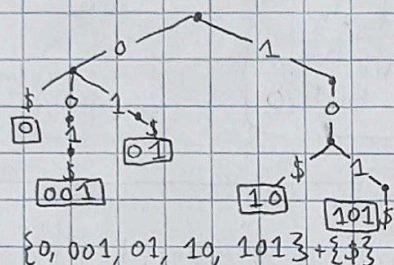
Dictionary for Words

\$ → end-of-word character.
 $\text{strcmp}(w_1, w_2) \rightarrow \{-1, 0, 1\}$

Binary Tries

We assume:
 • elements are $\Sigma^* = \{0, 1\}^*$
 • no word is a prefix of another.
 Second is true if we attach \$, or all words are the same length.

Binary Trie → Ternary Tree



$\{0, 001, 01, 10, 101\} \cup \{11\}$
 All operations have runtime $O(|w|)$.

Variations

① No leaf Labels

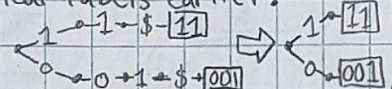
Don't store words at end.

② Allow Proper Prefixes

Tree becomes binary, holding value of whether in dictionary.

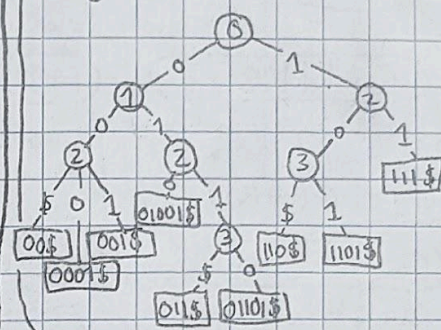
③ Pruned tries

If more traversals won't be needed, just store the leaf labels earlier!



Compressed Tries

Like pruned trees on steroids. Each node stores the lower limit on the length of its "sub-words."



Hashing Assumptions

- We know the universe of keys
- Have hash function $U \rightarrow \{0, 1, \dots, M-1\}$
- Dictionary is in array T , $|T| = M$.

Uniform Hashing Assumption:

Any hash function is equally likely.

Then: $P(h(k) = i) = \frac{1}{M}$
 for any k and i .

HASHING

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

Direct Addressing: The key k goes directly to the map.

(If k is known bounded $0 \leq k \leq M$, everything is $\Theta(1)$.)

→ Can't be used if: • keys aren't ints • ? smth else

Hash Collision: When a new key maps to a location already taken.

Solving Hash Collisions

① Chaining: each slot is a bucket. Use unsorted linked list w/ MTF.

→ average bucket size is $\frac{n}{M} = \alpha$ (load factor). Insert is still $\Theta(1)$ but search and delete are $\Theta(1 + \alpha)$.

→ Re-hashing: when α gets too big, double the hash table size.

② Probe Sequencing: Take an extra parameter for probing.

Search & insert are normal but delete needs to not leave holes,

→ so: delete: → move later items back
 → mark slot as deleted, not NIL.

→ Linear Probing: $h(k, i) = (h(k) + i) \bmod M$

③ Double Hashing: open addressing (meaning look for new

address when collide) with a second hash function. $h_0(k)$ & $h_1(k)$.

→ We need $h_1(k) \neq 0$, and $h_1(k)$ relative prime w/ M for all k .

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

④ Cuckoo Hashing: take h_0, h_1 , use T_0, T_1 tables. Item with

key k can only be at $T_0[h_0(k)]$ or $T_1[h_1(k)]$, so we do booting if there's a collision! Rehash everything if too many boots.

→ load factor $\alpha = n / (|T_0| + |T_1|)$.

→ for insertion expected runtime to be constant we need $\alpha < \frac{1}{2}$.

Independent hash functions:

$h_1(k)$ and $h_2(k)$ should be independent otherwise, StUpId.

Multiplication Method:

$$h(k) = \lfloor M \cdot (kA - \lfloor kA \rfloor) \rfloor, 0 < A < 1$$

Carter-Wegman's Universal Hashing

All keys $\{0, 1, \dots, p-1\}$ for a (big) prime p . Choose $M < p$. Choose random a, b , $a \neq b$, in key range.

$$h(k) = ((ak + b) \bmod p) \bmod M$$

keys collide with probability at most $\frac{1}{M}$.

CS 240 ③

Range Searches

S = output size
 I = interval (x, x')

- If unsorted, runtime $\Omega(n)$ since we have to search each individually.

RangeSearch implementation:

- ↳ binary search to find the two limits, i, i'
- ↳ report $A[i+1, \dots, i'-1]$ and $A[i], A[i']$ if in range.

Multidimensional (d) Range Search

Given a query-rectangle, return all points in it!

Assumption: No two x- or y-coors are the same.

Kd-Trees

Suppose n points... use \sqrt{d} splitting lines at halfway:

- ① assume first split is vertical.
- ② height is $O(\log(n))$
- ③ maintain height by occasionally rebuilding.
- ④ Range search complexity: $O(s + Q(n))$

output size # of "boundary" nodes visited

Range Tree

Literally just a balanced binary search tree for x-coordinates that stores a balanced binary search tree for y-coordinates.

Each node in the main tree is a point and also another tree, that holds the same children.

Space: $O(n \log n)$

Algorithms (for pattern matching)

Karp-Rabin: Compute fingerprint (hashing) for each guess.

- ↳ if fingerprint different from P's, no need to check.
 - ↳ key insight: update fingerprints in constant time.
- This allows runtime to just be $O(m+n)$, worst: $O(mn)$.

Boyer-Moore: Better pattern matching on english.

- ① Reverse order (where)
- ② Bad char jumps
- ③ Good suffix jumps (② & ③ are both "move as far as can").

Implement w/a last occurrence array. Example:

0 1 2 3 4 → char | p | a | e | r | others
L[i] | 2 | 1 | 3 | 4 | -1

$O(m \cdot \Sigma)$ construction

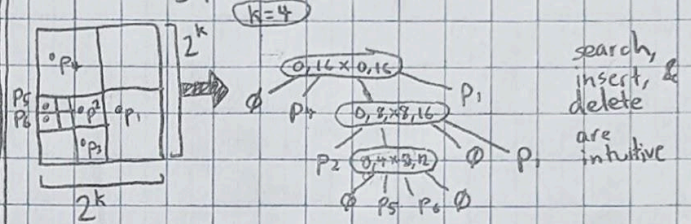
② ★ To shift after guess at i_{old} , $i_{new} = i_{old} + (m-1) - \min\{L[char], j_{old}-1\}$

③ Good suffix is like "good" character but in groups (no details).

Quad Trees

Bounding Box: $R = [0, 2^k) \times [0, 2^k)$.

↳ make by picking smallest k that works w/data.



search, insert, & delete are intuitive

Quadtree Range Search

Get given a rectangle, do expected stuff!

spread factor: $\beta(S) = \frac{\text{sidelength of } R(\text{region})}{\text{min distance between points}}$

height: $h = O(\log(\beta(S)))$

building complexity: $O(hn)$

Module 9 Pattern Matching

Definitions

Pattern Matching: find a string in a lot of text.

$T[0 \dots n-1]$: The text (haystack)

$P[0 \dots m-1]$: The pattern (needle)

returns first occurrence

If P doesn't occur in T , return FAIL.

Substring: Exactly what you think. eg. $T[4 \dots 7]$

Prefix/Suffix: substring at start/end.

Guess or Shift: position i s.t P might start at $T[i]$.

Check: position j , $0 \leq j < m$, to check

$T[i+j] = P[j]$

Use DFA to read the letters linearly



some real animals lol

these what you gonna look like if you keep not showering during exam season
take care of yourself