

6.046 Lecture #1

Design and Analysis of Algorithms

- Sign up for a recitation section on stellar!

Modules:

- Divide and Conquer
- Optimization: greedy, dynamic programming
- Network Flow
- Intractability: approximation
- Advanced topics

Very similar problems can have very different complexity

Recall: P: class of problems solved in polynomial time $O(n^k)$ for some constant k

NP: " " verifiable "

NP-Complete: in NP and as hard as any problem in NP

Interval Scheduling

single resource

Resources and requests 1, ..., n

$s(i)$ start time

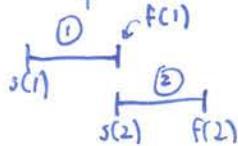
$f(i)$ finish time

$s(i) < f(i)$ can satisfy with single
resource simultaneously

Two requests i and j are compatible if they don't overlap.

$f(i) \leq s(j)$

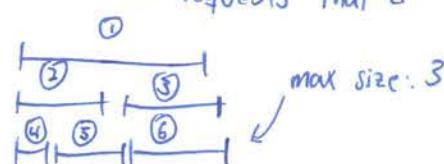
or
 $f(j) \leq s(i)$



Goal: select compatible subset of requests that is

of maximum size.

Ex



Claim: We can solve this problem using a greedy algorithm.

maximize now, don't look ahead.

A greedy algorithm is a myopic algorithm that processes input one piece at a time with no apparent look ahead

Greedy Interval Scheduling

1. Use a simple rule to select a request i
2. Reject all requests incompatible with i
3. Repeat until all requests are processed

Heuristics:

- each request find # incompatible requests, select the one with the min # incompatibles.



heuristic would select this one ... not optimal ∇

- select $R(i)$ with earliest finish time. Scan through $F(i)$, pick $\min^{\text{compatible}}$.
- \hookrightarrow Always works \rightarrow but need to prove with mathematical rigor.

Claim: Given a list of intervals L , greedy algorithm with earliest finish time produces k^* intervals, where k^* is maximum.

Proof: By induction on K^*

Base Case: $k^* = 1$ Any interval works

Inductive Step:

Suppose the claim holds for k^* and we are given a list of intervals whose optimal schedule whose optimal schedule has $k^* + 1$ intervals

$$S^* [1, 2, \dots, k^* + 1] = \langle (s(j_1), f(j_1)), (s(j_{k^*+1}), f(j_{k^*+1})) \rangle$$

$$S [1, \dots, k] = \langle (s(i_1), f(i_1)), (s(i_k), f(i_k)) \rangle$$

i and k^* not comparable $f(i_1) \leq f(j_1)$

$$S^{**} = \langle (s(i_1), f(i_1)), (s(j_2), f(j_2)), \dots, (s(j_{k^*+1}), f(j_{k^*+1})) \rangle$$

S^{**} is also optimal.

Define $L' = \text{set of intervals } s(i_j) \leq f(i_j)$

Since S^{**} is optimal for L' , $S^{**} [2, \dots, k^* + 1]$ is optimal for L

↳ subsets of optimal solution must be optimal
 \therefore optimal schedule for L' has k^* size

By inductive hypothesis, run the greedy algorithm on L' should produce a schedule of size k^*

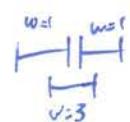
By construction, the greedy alg on L' gives $S[2..k]$ of size $k-1 = k^*$
 so $S[1..k]$ of size $k^* + 1$
 optimal because $k = k^* + 1$

□

Weighted Interval Scheduling

each request has weight $w(i)$

schedule subset of requests with maximum weights



↳ earliest time won't work now!

Dynamic Programming

$$R^{x^*}_{\text{subproblems}} = \left\{ \text{request } j \in R \mid s(j) \geq x \right\}$$

* $x = f(i)$ $R^{f(i)}$ requests later than $f(i)$

n : # of requests in original problem

subproblems = n

solve each subproblem once and memorize

subproblems * time/subproblem
 (assume $O(1)$ for lookups)

DP Guessing

Try each request i as a possible first request

$$\text{opt}(R) = \max_{1 \leq i \leq n} (w_i + \text{opt}(R^{f(i)}))$$

complexity $O(N^2)$ ✓

One small change can change the whole problem

non-identical machines

$$T = \{T_1, \dots, T_m\}$$

weight = 1 for all requests

BUT

$Q(i) \subseteq T$ is a set of machines that
i runs on.

↳ so some requests can only run on some machines.

This problem is NP-complete.

↳ can $k \leq n$ requests be scheduled.

Divide and Conquer

Paradigm
Convex Hull
Median Finding

Paradigm : Given a problem of size n . Divide it into "a" subproblems of size n/b , $b > 1$.

Solve each subproblem recursively.

(combine solutions of subproblems into overall solution.)

$$T(n) = aT(\frac{n}{b}) + [\text{work for merge}]$$

$a \geq 1$

"merge"

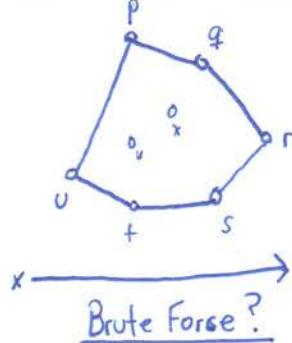
Convex Hull (2D)

Given n points in plane

$$S = \{(x_i, y_i) \mid i=1, 2, \dots, n\}$$

assume no two have the same x or y coord, no three in same line

(Convex Hull - smallest convex polygon containing all points in S) $CH(S)$



Sequence of points on the boundary in clockwise order as doubly linked list.

$$\textcirclearrowleft p \leftrightarrow q \leftrightarrow r \leftrightarrow s \leftrightarrow t \leftrightarrow u \leftrightarrow v \textcirclearrowright$$

- Draw a line between any two points

$O(n^2)$ segments

- check if all points are in one of the halves

$O(n)$ test complexity

- if yes \rightarrow segment is part of convex hull

$\Rightarrow O(n^3)$

Divide and Conquer

- Sort the points by x -coords

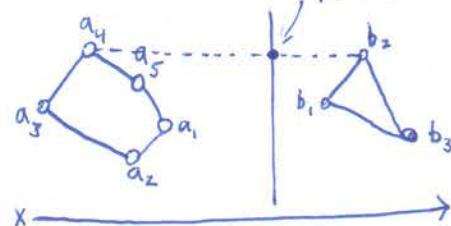
- For input set S

Divide into left half A and right side B by x -coords

compute $CH(A)$ and $CH(B)$

combine.

How to merge?



Obvious merge algorithm: Look at all pairs of points $\Rightarrow \Theta(n^2)$

Two-Finger and string algorithm:

$i=1$
 $j=1$
while ($y(i, j+1) > y(i, j)$ or $y(i-1, j) > y(i, j)$):

if $y(i, j+1) > y(i, j)$: # more right finger

$j=j+1 \text{ (mod } p)$

else:

$i=i-1 \text{ (mod } p)$

return (a_i, b_j) as upper tangent

Complexity:

On every iteration one of the counters will update
 $\Rightarrow \Theta(n)$ merge.

Overall Complexity: $T(n) = aT\left(\frac{n}{2}\right) + \Theta(n)$
 $= \Theta(n \log n)$ ✓

cut and paste? - see notes

Median Finding

Median finding in better than $\Theta(n \log n)$ time.

Given a set of n numbers, define rank of (x) as #'s in the set that are $\leq x$

Find element of rank $\lfloor \frac{n+1}{2} \rfloor$ lower median $\lceil \frac{n+1}{2} \rceil$ upper median

select(S, i)

Pick $x \in S$

Divide $x \in S$

compute $k = \text{rank}(x)$

$B = \{y \in S \mid y < x\}$

$C = \{y \in S \mid y > x\}$

B		x		C
$k-1$ elm				$n-k$ elm

$\Theta(n^2)$ if bad selection

in which inbalanced B and C

```

if k=i: return x
else if k > i: return select(B,i)
else if k < i: return select(C,i-k)

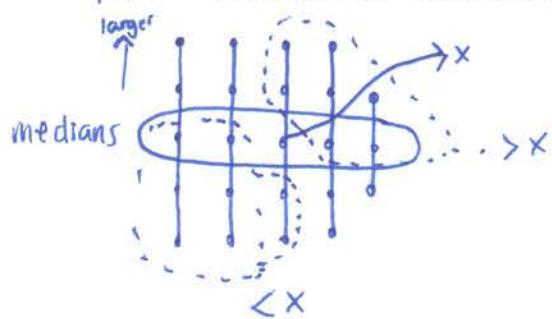
```

Pick x cleverly

Arrange S into columns of size 5 ($\lceil \frac{n}{5} \rceil$ cols)

Sort each col (big elms in top) linear time (since only 5)

Find "median of medians" of x



How many elements are guaranteed to be $> x$?

Half of the $\lceil \frac{n}{5} \rceil$ groups contribute at least

3 elements $> x$ except for 1 group with less than 5 elements and 1 group that contains x

At least $3(\lceil \frac{n}{10} \rceil - 2)$ elements are $> x$

" " " < x

Recurrence :

$$T(n) = \begin{cases} O(1) & \text{for } n \leq 140 \\ T\left(\lceil \frac{n}{5} \rceil\right) + \overbrace{\left(\frac{7n}{10} + 6\right)}^{\substack{\text{finding median of} \\ \text{medians}}} + \underbrace{\Theta(n)}_{\substack{\text{discarding} \\ \text{elements}}} & \text{otherwise} \end{cases}$$

\nearrow sorting of all columns

Deepak Narayanan deepakn@mit.edu

OH: Wed 7-9 pm

- Weighted Interval Scheduling
- Divide and conquer (Matrix multiplication)
- Master's Theorem

Weighted Interval Scheduling

n requests, each associated with start time, end time and a weight
Want to maximize the total weight of intervals picked

$O(n^2)$ Dynamic Programming Solution

$$R^x = \{ \text{request } j \in R \mid s(j) \geq x \}$$

If $x = f(i)$, then $R^{f(i)}$ represents the set of requests that follow i :

$$R^{f(1)}, R^{f(2)}, R^{f(3)}, \dots, R^{f(n)}$$

$$\text{opt}(R) = \max_{1 \leq i \leq n} \{ w_i + \text{opt}(R^{f(i)}) \}$$

subproblems: n

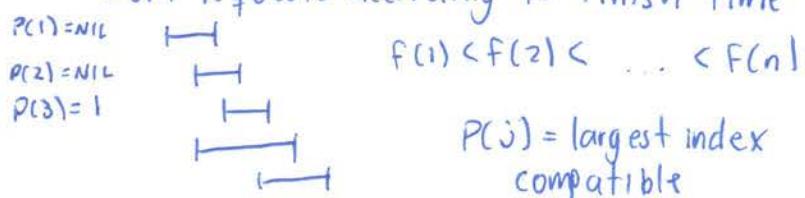
time/problem: $O(n)$

Total: $O(n^2)$

↪ not that great ::

$O(n \log n)$ Solution

Sort requests according to Finish time $\rightarrow O(n \log n)$



$M[i] \rightarrow \text{max weight for requests } \{1, 2, 3, \dots, i\}$

$$M[i] = \max \{ w_i + M[p(i)], M[i-1] \}$$

↪ i belongs in
optimal
configuration

↪ does not

subproblems: n
time/problem: $O(1)$
total (for this step): $O(n)$

Overall Total:
 $O(n) + O(n \log n)$
 $= O(n \log n)$

Divide and Conquer

Matrix Multiplication

Square matrices A and B , compute $C = A \cdot B$

$$C_{ij} = \sum_{k=1}^n A_{ik} B_{kj}$$

$\overbrace{\hspace{10em}}$ $O(n)$

Total runtime of naive $O(n^3)$ \nwarrow not good!

Strassen's Algorithm

$$A, B, C \rightarrow 2^k \times 2^k$$

$$2^{k-1} \times 2^{k-1}$$

Divide matrices: $A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$ $\quad \quad \quad C = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}$ $B = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}$

$$C_{11} = A_{11}B_{11} + A_{12}B_{12}$$

$$C_{12} = A_{11}B_{12} + A_{22}B_{22}$$

$$C_{21} = A_{21}B_{11} + A_{22}B_{21}$$

$$C_{22} = A_{21}B_{12} + A_{22}B_{22}$$

} math

$\Theta(n^2)$ since n^2 terms in each of these matrices

$$T(n) = 8T(n/2) + \Theta(n^2)$$

$$T(n) = \Theta(n^{\log_2 8}) = \Theta(n^3)$$

wtf? still the same?

But wait ...

$$\left. \begin{array}{l} M_1 = (A_{11} + A_{22})(B_{11} + B_{22}) \\ M_2 = (A_{21} + A_{22})B_{11} \\ M_3 = A_{11}(B_{12} - B_{22}) \\ M_4 = A_{22}(B_{21} - B_{11}) \\ M_5 = (A_{11} + A_{12})B_{22} \\ M_6 = (A_{21} - A_{11})(B_{11} + B_{12}) \\ M_7 = (A_{12} - A_{22})(B_{21} + B_{22}) \end{array} \right\}$$

magic
math

$$\left. \begin{array}{l} C_{11} = M_1 + M_4 - M_5 + M_7 \\ C_{12} = M_3 + M_5 \\ C_{21} = M_2 + M_4 \\ C_{22} = M_1 - M_2 + M_3 + M_6 \end{array} \right\}$$

more additions

$$\rightarrow T(n) = 7T(n/2) + \Theta(18n^2)$$

$$T(n) = \Theta(n^{\log_2 7}) = \Theta(n^{2.8079})$$

now we only have 7 multiplications!

YAY slightly better

Master's Theorem

$$T(n) = aT\left(\frac{n}{b}\right) + \Theta(f(n))$$

Case I: leaves dominate

$$f(n) = \Theta(n^{\log_b a})$$

$$T(n) = \Theta(n^{\log_b a})$$

To prove, look at recursion tree

Case II: About equal

$$f(n) = \Theta(n^{\log_b a} \cdot \log^k n)$$

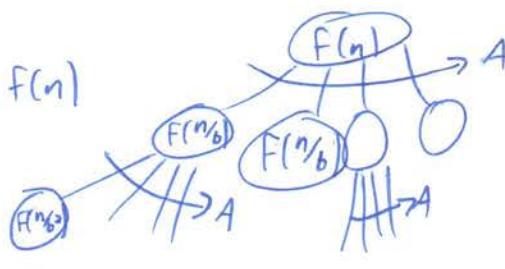
$$T(n) = \Theta(f(n) \cdot n^{\log_b a} \log^{k+1} n)$$

Case III: $f(n)$ dominates

$$f(n) = \Theta(n^{\log_b a} + \epsilon)$$

$$T(n) = \Theta(f(n))$$

$$T(n) = aT\left(\frac{n}{b}\right) + f(n)$$



case II
f(n)

$$aF(n/b)$$

$$a^2 f(n/b^2)$$

$a^h f(n/b^h)$
case I

- Fast Fourier Transform FFT

- polynomial operations vs representations
- divide and conquer algorithms
- collapsing samples / roots of unity
- FFT, IFFT and polynomial multiplication

Polynomial

$$\begin{aligned} A(x) &= a_0 + a_1 x + a_2 x^2 + \dots + a_{n-1} x^{n-1} \\ &= \sum_{k=0}^{n-1} a_k x^k \\ &= \langle a_0, a_1, \dots, a_{n-1} \rangle \quad \text{vector notation just coefficients} \end{aligned}$$

Operations on polynomials

① evaluation: $A(x)$ and $x_0 \rightarrow A(x_0)$

Horner's Rule:

$$A(x) = a_0 + x(a_1 + x(a_2 + \dots + x(a_{n-1} \dots))) \Rightarrow O(n) \text{ time}$$

② Addition: $A(x)$ and $B(x) \rightarrow C(x) = A(x) + B(x) \quad \forall x$

$$c_k = a_k + b_k$$

③ Multiplication: $A(x), B(x) \rightarrow C(x) = A(x)B(x)$

$$(a_0 + a_1 x + a_2 x^2 + \dots) \cdot (b_0 + b_1 x + b_2 x^2 + \dots) = a_0 b_0 + (a_1 b_0 + a_0 b_1) x + \dots$$

$$c_k = \sum_{j=0}^k a_j b_{k-j} \quad \Rightarrow O(n^2) \quad ;$$

TODAY: $O(n \log n)$

Convolution of vectors A and reverse(B)
(inner products of all possible shifts)

Representations

- Ⓐ coeff, vector: $a_0 + a_1 x + a_2 x^2 + \dots$
- Ⓑ roots: r_0, r_1, \dots, r_{n-1} $P = c(x-r_0)(x-r_1) \dots (x-r_{n-1})$
- Ⓒ Samples (x_k, y_k) for $k=0, 1, \dots, n-1$
 $A(x_k) = y_k \quad x_k$'s distinct

Algorithms	① coeffs	② roots	③ samples
① Evaluation	$O(n)$	$O(n)$	$O(n^2)$
② Addition	$O(n)$	∞	$O(n)$
③ Multiplication	$O(n^2)$	$O(n)$	$O(n)$

$O(n \lg n)$

(convert coeffs \longleftrightarrow samples
in $n \lg n$ using FFT)

Matrix View

Vandermonde Matrix

$$\begin{pmatrix} 1 & x_0 & x_0^2 & \dots & x_0^{n-1} \\ 1 & x_1 & x_1^2 & \dots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \dots & x_2^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n-1} & x_{n-1}^2 & \dots & x_{n-1}^{n-1} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_{n-1} \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ y_2 \\ \vdots \\ y_{n-1} \end{pmatrix}$$

$$V_{jk} = x_j^k$$

(coeffs \rightarrow samples = $V \cdot A$ $O(n^2)$)

samples \rightarrow coeffs = $V^{-1}A$

- Gaussian elimination $O(n^3)$

- $V^{-1}A \rightarrow O(n^2)$

Divide and Conquer algorithm

goal: $A(x)$ For $x \in X$

① divide into even and odd coefficients

$$A_{\text{even}}(x) = \sum_{k=0}^{\frac{n}{2}} a_{2k} x^k = \langle a_0, a_2, a_4, \dots \rangle$$

$$A_{\text{odd}}(x) = \sum_{k=0}^{\frac{n}{2}-1} a_{2k+1} x^k = \langle a_1, a_3, a_5, \dots \rangle$$

$$a_{2k} x^k$$

② Conquer: recursively compute $A_{\text{even}}(y)$ and $A_{\text{odd}}(y)$ For $y \in X^2 = \{x^2 \mid x \in X\}$

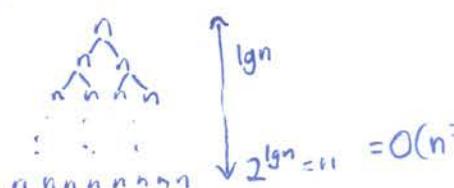
③ Combine

$$A(x) = A_{\text{even}}(x^2) + x A_{\text{odd}}(x^2) \quad \leftarrow \text{Just algebra}$$

For $x \in X$

$$T(n, |X|) = 2 \cdot T(\frac{n}{2}, |X|) + O(n + |X|)$$

$|X|$ never gets smaller



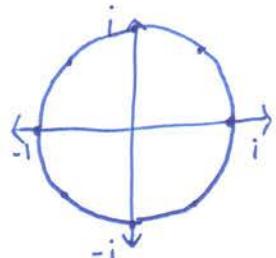
?? ↓

collapsing set X if

$$|X^2| = |X|/2 \text{ and } X^2 \text{ is collapsing}$$

$$\text{or } |X|=1$$

n^{th} roots of unity



$$(\cos \theta, \sin \theta)$$

$$= \cos \theta + i \sin \theta$$

$$= e^{i\theta}$$

$$\text{For } \theta = 0, \frac{T}{n}, \frac{2T}{n}, \dots, \frac{n-1}{n} T$$

$$T = 2\pi$$

$$(e^{i\theta})^2 = e^{i(2\theta)} = e^{i(2\theta \bmod T)}$$

$$e^{i k T/n} = e^{i(2k) T/n}$$

$$\boxed{e^{iT} = 1}$$

$n = 2^k \Rightarrow n^{\text{th}}$ roots of unity collapsing

$$X_k = e^{ik T/n}$$

Fast Fourier Transform (FFT)

= Divide and conquer algorithm for

discrete
DFT

Discrete Fourier Transform (DFT)

$$= V \cdot A \text{ for } X_k = e^{ik T/n}$$

$$V_{jk} = X_j^k = e^{ijk(T/n)}$$

Fast polynomial multiplication

$$A^* = \text{FFT}(A) \quad A \text{ coef} \rightarrow \text{sample}$$

$$B^* = \text{FFT}(B)$$

$$(C_k^*) = A_k^* B_k^* \quad \forall k$$

$$C = \text{IFFT}(C^*)$$

Every # has exactly 2 $\sqrt{2}$

$$|X| = 1 : X = \{1\}$$

$$|X| = 2 : X = \{-1, 1\}$$

$$|X| = 4 : X = \{i, -i, -1, 1\}$$

$$|X| = 8 : X = \left\{ \pm \frac{\sqrt{2}}{2}(1+i), \pm \frac{\sqrt{2}}{2}(1-i), i, -i \right\}$$

Claim $V^{-1} = \bar{V}/n$

$$P = V \cdot \bar{V} = n I = \begin{pmatrix} n & & \\ & \ddots & \\ & & 0 \end{pmatrix}$$

$$P_{jk} = (\text{row } j \text{ of } V) \cdot (\text{col } k \text{ of } \bar{V})$$

$$= \sum_{m=0}^{n-1} e^{i\pi jm/n} - e^{-i\pi mk/n}$$

$$= \sum_{m=0}^{n-1} e^{i\pi m(j-k)}$$

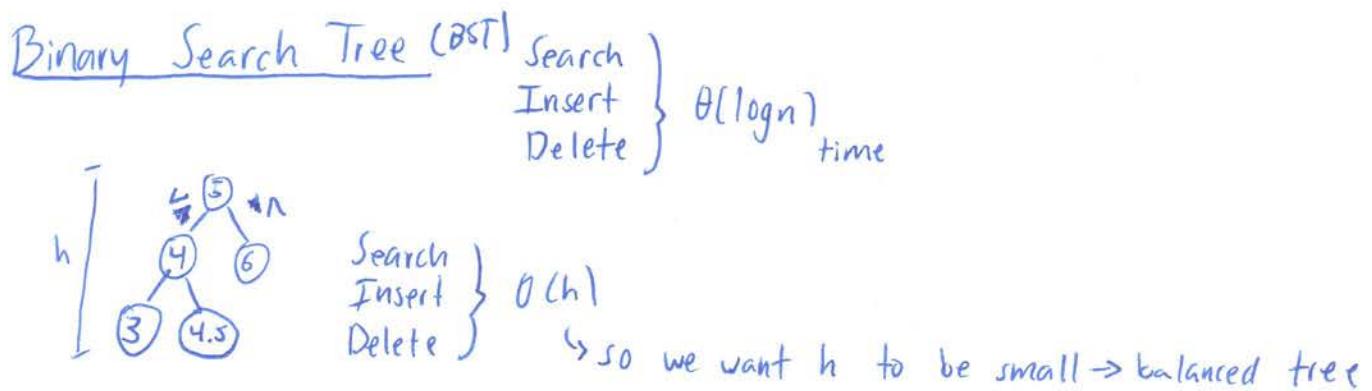
$$\text{If } j \neq k: i = \sum_{m=0}^{n-1} \left(e^{i\pi \frac{j-k}{n}} \right)^m$$
$$= \frac{\left(e^{i\pi \frac{j-k}{n}} \right)^{n+1} - 1}{e^{i\pi \frac{j-k}{n}} - 1}$$

geometric series

$$\sum_{k=0}^{n-1} z^k = \frac{z^n - 1}{z - 1}$$

Applications:

Binary Search Trees
2-3 Trees
B-trees

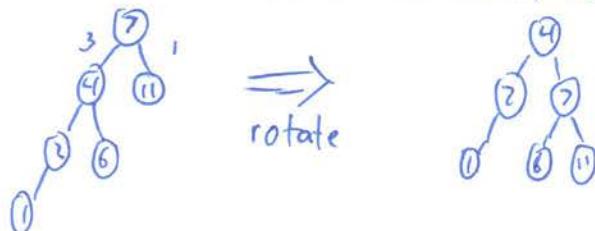


Balanced BST

Invariant: For every node, difference b/w depth of left and right subtree must be at most 1. $\Rightarrow h = O(\log n)$

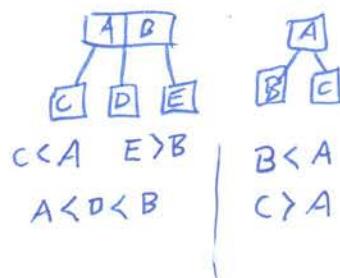
Rotations

Need to make sure invariant holds after insertion

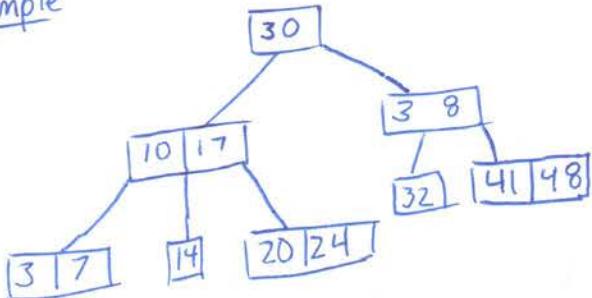


2-3 Trees

- Each node with non-leaf children must have two children or three children
- All leaves are at the same level
- All data is sorted order
- Every leaf node must contain one or two values



Example



elements : n
height = $O(\lg n)$

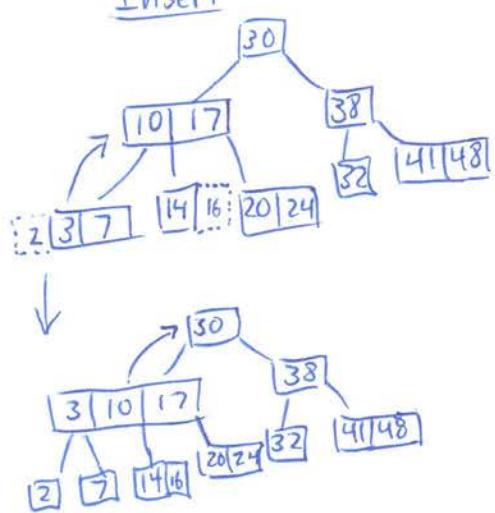
↑
more dense than BST

Operations

Search

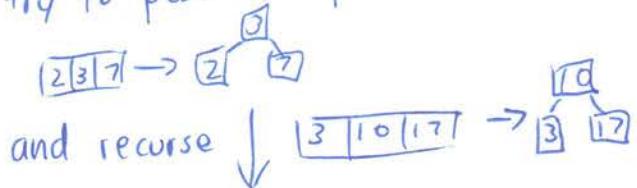
looking for k
like BST with an extra check if 3 node root.
Runtime : $O(\lg n)$

Insert



Insert 16 : Search For 16, insert at spot

Insert 2 : Inserting would create a node with three values → Breaks invariant
Solution: Take median of 3 elements and try to push to parent.

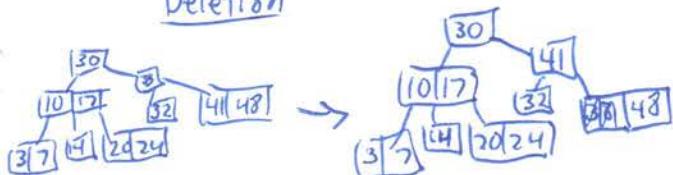


Done :

IF root has three elements
make a new root

(will push an element at most the height of the tree)
Runtime : $O(\lg n)$

Deletion



Delete 38

Successor of any element in an internal node
will be a leaf node.

· Swap with successor
· delete if possible
if not => special case → see book
Runtime $O(\lg n)$

B-trees

Parameterized by constant B

- Each internal node must have $\geq B$ children and $\leq 2B$ children.
- All leaves must be at same depth
- Data is sorted

If a node has i children, it has $i-1$ values.

Insert
Search
Delete $\{ O(\log n) \}$

6.046 Notes

Representing Polynomials

$$A(x) = \sum_{j=0}^{n-1} a_j x^j$$

Coefficient Representation

- Vector of coefficients $\langle a_0, a_1, \dots, a_{n-1} \rangle$

- Evaluating: $\Theta(n)$ using Horner's rule : $A(x_0) = a_0 + x_0(a_1 + x_0(a_2 + \dots + x_0(a_{n-1})))$

- Adding: $\Theta(n)$ = c where $c_j = a_j + b_j$: Add all coeffs w/ same degree

- Multiplying: $\Theta(n^2)$

convolution $c_j = \sum_{k=0}^j a_k b_{j-k}$ $c = a \otimes b$

Point-value representation

- Set of point value pairs $\{(x_0, y_0), (x_1, y_1), \dots, (x_{n-1}, y_{n-1})\}$

where $y_k = A(x_k)$

- Evaluating: $\Theta(n^2)$ using Horner's rule. $\Theta(n \log n)$ if points chosen cleverly.

Interpolation - inverse of evaluating

For any set of n points there is a unique polynomial of degree-bound n

- Adding: $\Theta(n)$

- Multiplying: $\Theta(n)$

★ Choose the complex roots of unity

can take Discrete Fourier Transform of coeff vector \rightarrow point-value

↓
in next section 30.2

inverse DFT in $\Theta(n \log n)$

So: originally

coeffs $\xrightarrow[\Theta(n^2)]{\text{multiply}}$ solution $\Theta(n^2)$

now:

coeffs $\xrightarrow[\Theta(n \log n)]{\text{evaluation}}$ point-value $\xrightarrow[\Theta(n)]{\text{multiply}}$ $\xrightarrow[\Theta(n \log n)]{\text{interpolation}}$ solution $\Theta(n)$

van Emde Boas

- Series of improved DSs
- Insert, Successor
- Delete
- Space

Goal: maintain n elements among $\{0, 1, \dots, u-1\}$

subject to: Insert, Delete, Successor

Successor
in $O(\lg \lg u)$ time if $u = n^{O(1)}$ or $n^{\lg O(1)} \leq u$
then $\lg \lg u = O(\lg \lg n)$

- application: network routers

Intuition: (w)How to get $O(\lg \lg u)$

- binary search on levels of tree

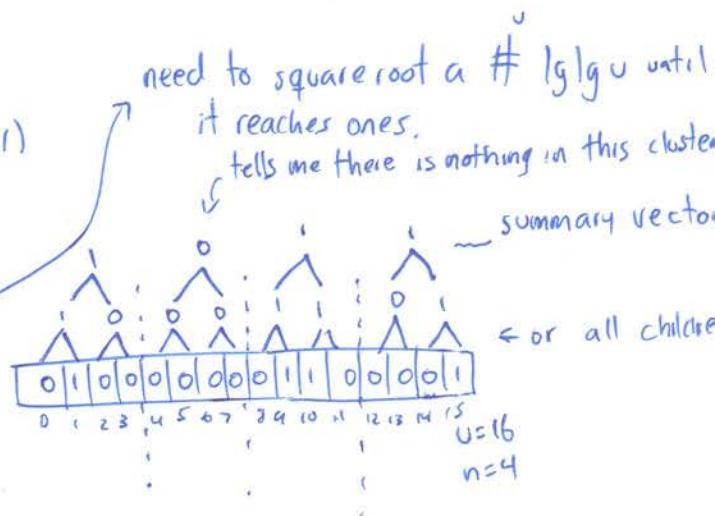
$$\begin{aligned} T(k) &= T(k/2) + O(1) \\ &= O(\lg k) \end{aligned}$$

$$T(\lg u) = T\left(\frac{\lg u}{2}\right) + O(1)$$

$$\begin{aligned} T'(u) &= T(\sqrt{u}) + O(1) \\ &= O(\lg \lg u) \end{aligned}$$

① bit vector = array of size u
 $0 = \text{absent}$ $1 = \text{present}$

Insert/Delete: $O(1)$
successor: $O(u)$



② Split universe into \sqrt{u} clusters of size \sqrt{u}

Insert: $O(1)$

Successor(x): - look in x 's cluster

- look for next one bit in summary vector and find 1st one in that cluster

$O(\sqrt{u})$

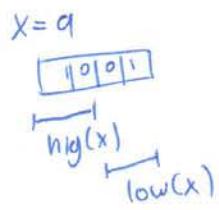
Using recursion we can turn this to $T(\sqrt{u})$!

if $x = i \cdot 5^j + j$
 $0 \leq j \leq 5^j$

$$\text{high}(x) = \lfloor x / 5^j \rfloor$$

$$\text{low}(x) = x \bmod 5^j$$

$$\text{index}(i, j) = i \cdot 5^j + j$$



$$a = 2 \cdot 4 + 1$$

- ③ recurse: $V = \text{size} - u \cdot \text{VEB}$
- $V_{\text{opt}} \text{cluster}[i] = \text{size} - 5^j \cdot \text{VEB}$
- $V_{\text{summary}} = \text{size} \cdot \sqrt{u} \cdot \text{VEB}$

$\text{Insert}(V, x)$:

- Insert into cluster: $\text{Insert}(V_{\text{cluster}}[\text{high}(x)], \text{low}(x))$
- Update summary structure: $\text{Insert}(V_{\text{summary}}, \text{high}(x))$

$\text{Successor}(V, x)$:

$i = \text{high}(x)$
 $j = \text{successor}(V_{\text{cluster}}[i], \text{low}(x))$ \leftarrow try to find in x 's cluster
 if $j = \infty$: // wrong cluster
 $i = \text{successor}(V_{\text{summary}}, i)$ // look for next 1 in summary
 $j = \text{successor}(V_{\text{cluster}}[i], -\infty)$
 return $\text{index}(i, j)$

Runtimes:

These aren't that great because we do more than 1 recursive call
 For each call.

Ex. Insert $T(u) = 2T(5^j) + O(1)$

$T(\lg u) = 2T\left(\frac{\lg u}{2}\right) + O(1) = O(\lg u)$ \leftarrow could be better
 \leftarrow would be fixed if there was only 1 recursive call

Successor:

$$\tilde{O}(\lg u)$$

④ Store min and max

Augment data structure with min and max

Insert(V, x):

```
if  $x < V.\min$ :  
     $V.\min = x$  // Every structure knows minimum  
if  $x > V.\max$ :  
     $V.\max = x$   
:
```

Successor(V, x)

$i = \text{high}(x)$

:

:

:

```
     $j = V.\text{cluster}[i].\min$  // Now could only take 2 recursive calls
```

But we can do better...

Successor(V, x):

$i = \text{high}(x)$

if $\text{low}(x) < V.\text{cluster}[i].\max$:

$j = \text{successor}(V.\text{cluster}[i], \text{low}(x))$

else:

$i = \text{successor}(V.\text{summary}, \text{high}(x))$

$j = V.\text{cluster}[i].\min$

return index(i, j)

// Runtime $O(\lg \lg u)$

Insert not done yet:

First time we insert into cluster still takes two recursive calls...

⑤ Don't store min recursively

```
Insert(V, x):  
    if x < V.min: return V.min  
    if V.min = None: V.min = V.max = x return  
    if x < V.min: swap(x ↔ V.min)  
    if x > V.max: V.max = x  
    if V.cluster[high(x)].min = None:  
        Insert(V.summary, high(x)) // if this gets called  
    Insert(V.cluster[high(x)], low(x)) // this call takes O(1) because of
```

Runtime: $O(\lg \lg u)$

Now for delete:

```
Delete(V, x):  
    i = V.summary.min  
    if i = None: V.min = V.max = None return  
    x = V.min = index(i, V.cluster[i].min)  
    Delete(V.cluster[high(x)], low(x))  
    if V.cluster[high(x)].min = None:  
        Delete(V.summary, high(x))  
    if x = V.max: if V.summary.max = None: V.max = V.min  
    else: i = V.summary.max
```

⑥ Space is $O(u)$

- only store nonempty clusters
- $\Rightarrow V.\text{cluster} = \text{hashtable}$ // instead of array
- $O(n \lg \lg u)$ space

Union-Find and Non-rigorous amortized analysis

Union-Find data structure

Disjoint-set data structure

Collection of pairwise independent sets $\{S_1, S_2, S_3, \dots, S_n\}$
 $\{1, 2\}$
 $\{3\}$
 \uparrow
 \uparrow
 $\{S_1, S_2, S_3, \dots, S_n\}$
each set S_i should have a single representative element $\text{repr}[S_i]$ - MAKE-SET(x) $\rightarrow \{x\}$, $\text{repr}[\{x\}] = x$ - FIND-SET(x) \rightarrow First find set S_x which contains x , and then returns $\text{repr}[S_x]$ - UNION(x, y) \rightarrow Replace S_x and S_y with $S_x \cup S_y$ if $S_x \neq S_y$

Motivation: Keep track of connected components in dynamically changing graph

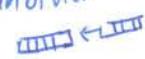
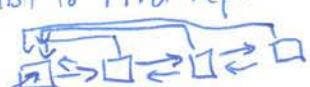
 $G = V, E$ Add node v to G

(v has associated edges)

Linked list solution

Represent each set as a doubly linked list


 $\{3\}$
 $\boxed{3}$

 $\{4, 5, 6\}$
 $\boxed{4} \rightarrow \boxed{5} \rightarrow \boxed{6}$
 $\text{rep}(S_i) = \text{head of list}$
 representing S_i
MAKE-SET(x) \rightarrow Initialize x as a lone node $O(1)$ FIND-SET(x) \rightarrow Hashtable that maps x to node containing x $O(n)$
follow parent pointer to head and return repr . $O(n) \in \text{BAD}$ UNION(x, y) \rightarrow link head of one list and tail of another $\text{repr}(S_x \cup S_y) = \text{head of 1st list}$. $O(n)$ to find head and tail. \curvearrowleft BADProblem: Not efficient to walk through list to find repr Solution: Augment all x with $\text{repr}[S_x]$ 

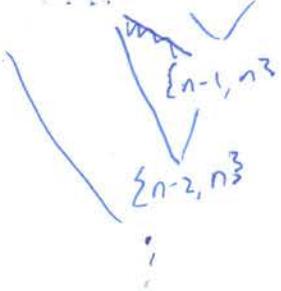
and tail

Find-set: $O(1)$ UNION(x, y): Have to update all pointers in y $O(n)$ \leftarrow still bad

Poorly performing case:

Initially n sets with one element n

$\{1\}$ $\{2\}$ $\{3\}$ $\{4\}$... $\{n-1\}$ $\{n\}$



$$\text{Total cost: } 1+2+3+\dots+n = \frac{n(n+1)}{2} = O(n^2)$$

Solution: Append smaller to the end of larger

↳ then don't have to update too many pointers

$$\text{Total runtime: } 1+1+1+\dots = O(n)$$

Amortized Analysis

Recall table doubling from hashing.

Sometimes inserts can be very expensive, but usually really fast

With our improvement to the augmented linked list, n calls to make set followed by n calls to union costs an amortized time of $O(n \log n)$

↓
Union costs $O(\log n)$ on average

(consider one node X . Updating $\text{repr}[S_x]$ costs 1.)

when $\text{UNION}(x, y)$ is called

Either $\text{len}(S_x) > \text{len}(S_y) \rightarrow X's \text{ pointer not updated}$

or $\text{len}(S_y) < \text{len}(S_x) \rightarrow X's \text{ pointer updated}$

$\text{len}(S_x)$ at least doubles post union operation.

$\Rightarrow \text{repr}[S_x]$ is updated at most $\log n$ times.

Amortization

- aggregate method
- accounting method
- charging method
- potential method
- table doubling
- binary counter
- 2-3 trees

Recall: Table doubling [6.066]

- n items in table of size $m \Rightarrow$ expected cost of $O(1 + \frac{m}{n})$
- if n grows to m then double $m \quad m' = 2m \quad O(m)$ work

Total cost: For n insertions $= \Theta(2^0 + 2^1 + 2^2 + \dots + 2^{\lg n}) = \Theta(n)$

\Rightarrow "amortized cost per operation" = $O(1)$

Aggregate Method

total cost of k operations $\frac{1}{k}$ = amortized cost per operation

Amortized bounds

amortized
- assign cost for each operation such that "preserve the sum"

$$\sum_{\text{op}} \text{amortized cost} \geq \sum_{\text{op}} \text{actual cost}$$

Ex: 2-3 trees achieve $O(c)$ worst-case per create-empty $O(i \lg n^*)$ amortized per insert.

\emptyset amortized per delete since can't delete something that was not inserted

$$\text{total cost} = O(c + i \lg n^* + d \lg n^*) = O(c + i \lg n^* + d(O))$$

$\hookrightarrow \leq i$

i^* : max size over all time

Accounting Method

- allow an operation to store credit in bank account (≥ 0)
- allow an operation to pay for time using stored credit. balance ≥ 0 unused time

2-3 Trees

Claim: $O(\lg n)$ per insert NO STAR

- ∅ per delete, amortized so each coin has a different value
- put 1 coin worth $\Theta(\lg n)$ per insert
- delete consumes 1 coin amortized cost = actual cost + deposits - withdrawals
- invariant: 1 coin of $\lg i$ for $i = 1 \dots n$

Table

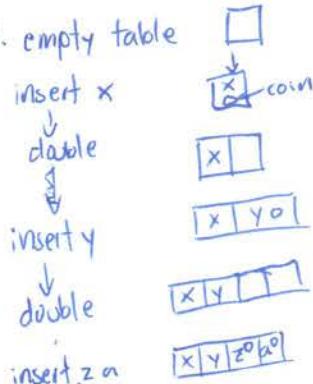
Table doubling (insertions only)

- when insert, add on it worth $c = O(1)$

- when double, last $n/2$ items have coins

$$\Rightarrow \text{amortized cost} = \Theta(n) - c \frac{n}{2} = O(n) \text{ if set } c \text{ large.}$$

Ex: empty table



Charging method:

- allow operations to charge costs retroactively to past (not feature)
- amortized cost = actual cost - total charge to past + total charge in Future

Table doubling (inserts only)

- charge doubling to inserts since last doubling ($n/2$) $\Rightarrow O(1)$ each

- only charge insert once

Table doubling^{and halving} (inserts and deletes): $m = O(n)$

100% Full \rightarrow double
25% Full \rightarrow halve

\rightarrow new array always 50% Full

charge halving to $\sum \frac{m}{4}$ deletions since last resize
 doubling to $\sum \frac{m}{2}$ insertions "

Potential Method

- define potential function Φ mapping data-structure configuration
 → non negative integer

- amortized cost = actual cost + $\Delta\bar{\Phi}$
 $\hookrightarrow \bar{\Phi}(\text{after op}) - \bar{\Phi}(\text{before op})$

$$\sum \text{amortized cost} = \sum \text{actual costs} + \bar{\Phi}(\text{end}) - \bar{\Phi}(\text{beginning})$$

⇒ pay $\bar{\Phi}(\text{beginning})$ at beginning

- should be $O(1)$ or \emptyset

"Potential high costs in Future"

Binary Counter 0011010111 $\xrightarrow{\text{increment}}$ 0011011000

- increment costs $O(\overbrace{\# \text{ trailing } 1's}^+ + 1)$

$$\bar{\Phi} = \# 1 \text{ bits}$$

$\Delta\bar{\Phi}$ should not change that much after each operations

- increment destroys +1s, creates one 1

$$\text{Amortized cost} = 1+1-1+1 = 2$$

2-3 Trees (insertions only)

- how many splits per insert?

$\leq \lg n$ worst case

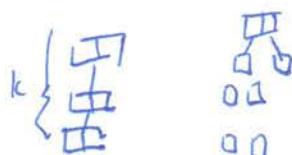
$O(1)$ amortized

$\bar{\Phi} = \# \text{ nodes with 3 children}$

$$\# \text{ splits} = k$$

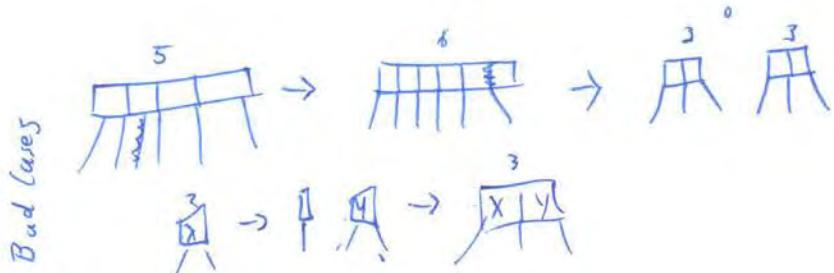
$$\Delta\bar{\Phi} = -k+1$$

$$\text{amortized cost} = k-k+1 = 1$$



2-3 trees

Inserts and Deletes in $(2, 5)$ -trees
 $(2 \leq \# \text{children} \leq 5)$



$$\Phi = \# \text{nodes w/ 2 children} + \# \text{nodes with 5 children}$$

insert w/ k splits $\Delta \Phi = -k + 1$

delete w/ k merges $\Delta \Phi = -k + 1$

$\Rightarrow O(1)$ splits/merges per insert/delete

Randomized Algorithms

- Why randomized
- Matrix multiply verification
- Quick Sort

- Algorithm that generates a random number $r \in \{1, \dots, R\}$ and makes decisions based on r 's value

On the same input on diff executions \rightarrow run for diff # steps
 - could produce diff outputs

Probably Correct

Monte Carlo

Probably Fast

Las Vegas

Probably correct +
probably fast

Atlantic City

Matrix Product

$$C = A \times B$$

Simple Algorithms: $O(n^3)$ multiplications

Strassen - multiply two 2×2 matrices using 7 mults (instead of 8): $O(n^{\log_2 7}) = O(n^{2.81})$

Coppersmith-Winograd: $O(n^{2.376})$

- Get a $O(n^2)$ algorithm

: If $A \times B = C$, then $\text{prob}[\text{output} = \text{YES}] = 1$

: If $A \times B \neq C$, then $\text{prob}[\text{output} = \text{YES}] \leq 1/2$ \leftarrow False positive

entries $\in \{0,1\} \bmod 2$ \hookrightarrow run it k times $\frac{1}{2^k} \Rightarrow$ still $O(n)$

Freval's Algorithm

choose random binary vector $r[1 \dots n]$

show that $\Pr[r_i = 1] = 1/2$ independently for $i=1, \dots, n$

IF $A(B_r) = C_r$ output YES

else output NO

IF $AB = C$, $A(B_r) = (AB)r = C_r$ // No false negatives

Analyzing correctness: $AB \neq C$

Claim: If $AB \neq C \Rightarrow \text{Prob}[ABr \neq Cr] \geq \frac{1}{2}$

Let $D = AB - C$. Our hypothesis is that $D \neq 0$.

We need to show that there are many r such that $Dr \neq 0$

Specifically, $\text{Prob}[Dr \neq 0] \geq \frac{1}{2}$ for randomly chosen r

$Dr = 0$ case // "bad r " \Rightarrow False positive

$D = AB - C \neq 0 \Rightarrow \exists i, j \text{ st } d_{ij} \neq 0$

$$\begin{array}{c} \text{d}_{ij} \neq 0 \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & \text{■} & & \\ \hline & & & \\ \hline \end{array} \end{array} \quad \begin{array}{c} v \\ \vdots \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{array} \quad d_{ij} = \begin{array}{|c|c|} \hline & \\ \hline & \text{■} \\ \hline & \\ \hline \end{array} \quad \begin{array}{c} Du \neq 0 \\ (Du)_{ji} = d_{ij} \neq 0 \\ \neq 0 \end{array}$$

Take any r that can be chosen by our alg s.t. $Dr = 0$

$$r' = r + v \quad \begin{matrix} \text{vector addition} \\ \text{r to r'} is 1 to 1 \end{matrix} \quad Dr' = D(r + v) = 0 + Du \neq 0$$

$$\text{if } r' = r + v$$

$$= r'' + v$$

$$\text{then } r = r''$$

$Dr \neq 0$

Discover an r' s.t. $Dr' \neq 0$ $r \leftrightarrow r'$ 1 to 1 mapping

Number of r' for which $Dr' \neq 0 \geq \frac{\text{number of } r}{\text{for which } Dr = 0}$

$$\Rightarrow \Pr[\text{if } r \neq 0] \geq \frac{1}{2}$$

randomly chosen

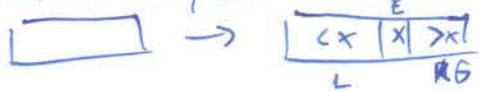
Quick Sort

Divide and Conquer in place no auxiliary space
 $O(n)$
 all work in divide step

n-elmt array A

Divide: 1. Pick a pivot elmt x in A.

Partition array into sub arrays



→ CLRS p 171

Conquer: Recursively sort subarrays L & R

Combine:

Basic Quick Sort

pivot $x = A[1]$ or $A[n]$

- partition given X $O(n)$ time

Worst case analysis

sorted or reverse sorted

one side L or R has $n-1$ elements, and other 0

$$T(n) = T(0) + T(n-1) + \Theta(n)$$

$$\begin{aligned} T(n) &= T(n-1) + \Theta(n) \\ &= \Theta(n^2) \end{aligned}$$

divide step

Pivot selection: intelligently \rightarrow worst case $O(n \lg n)$

guarantee balanced L and R using median selection that runs in $\Theta(n)$ time

$$T(n) = 2T(n/2) + \Theta(n) + \Theta(n)$$

thanks to recursive median selection
 median based pivoting.

Randomized Quicksort

expected time is $O(n \log n)$ for all input array A

CLRS p 181-184

↳ afraid of unbalanced partition

Paranoid Quicksort

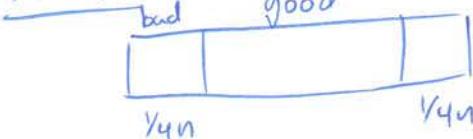
Repeat

choose pivot to be random element of A

Perform partition

Until resulting partition is s.t. $|L| \leq \frac{3}{4}|A|$ and $|R| \leq \frac{3}{4}|A|$

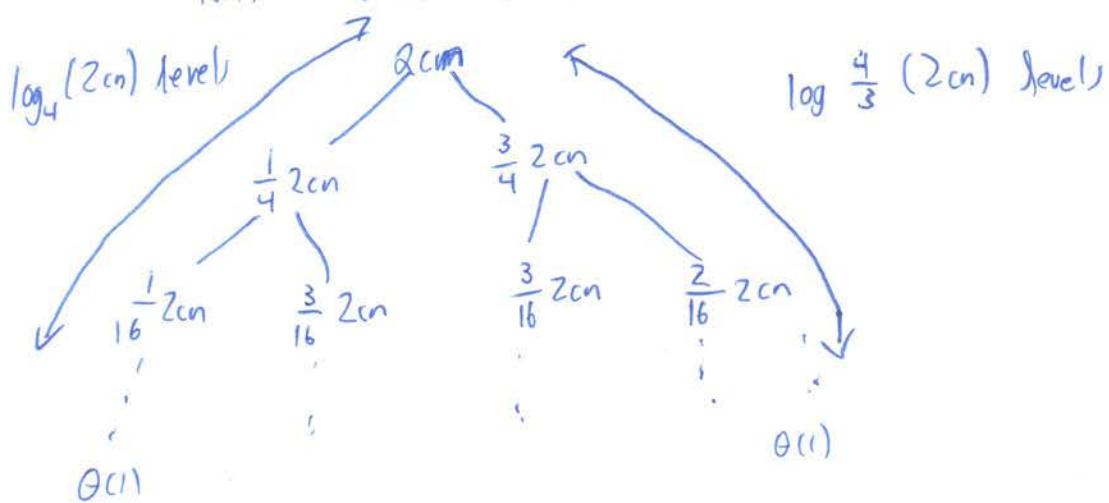
Analysis



A call is good with prob $\geq 1/2$

$$T(n) = T(\frac{n}{4}) + T(\frac{3n}{4}) + E(\# \text{iterations}) \cdot C \cdot n$$

$$T(n) = T(\frac{n}{4}) + T(\frac{3n}{4}) + 2 \cdot cn$$



2cn work at each level
max $\log_{4/3}(2cn)$ levels $\Rightarrow O(n \log n)$

6.046 Recitation

Expected running time of: Randomized Select
Randomized quicksort

Randomized select

Given an array and index. Find the k^{th} largest (or smallest) element in the array.

RANDOMIZED-SELECT (A, p, r, i):

if $p=r$: return $A[p]$

~~q → pivot~~ → Partition(A, p, r)

$k \rightarrow q-p+1$

// Different ways to select pivot

if $i \leq k$: RANDOMIZED-SELECT (A, p, q, i)

else: RANDOMIZED-SELECT ($A, q+1, r, i-k$)

$T(n) \rightarrow$ Runtime of select on input size n

$$\text{Ex}[T(n)] \leq \underbrace{\Theta(n)}_{\substack{\text{from} \\ \text{PARTITION}}} + \text{Ex}[\max(T(1), T(n-1))]$$

RANDOMIZED-PARTITION (A, p, r):

: → RANDOM(p, r)

Exchange $A[p]$ and $A[i]$

return PARTITION (A, p, r)

$$\begin{aligned} T(n) &\leq cn \sum_{i=0}^{n-1} P_i[E_i] (\max(T_i, T(n-1))) \\ &\leq cn + \sum_{i=0}^{n-1} \frac{1}{n} \cdot \max(T_i, T(n-1)) \\ &\leq cn + \frac{2}{n} \sum_{i=0}^{n/2-1} \max(T(i), T(n-1)) \end{aligned}$$

By substitution

$$\text{Guess } T(n) = O(n)$$

By induction:

$$\text{For all } k < n \quad T(n) = O(n) \Rightarrow T(n) \leq cn \quad \text{constant}$$

$$T(n) \leq n + \frac{2}{n} \sum_{i=0}^{n/2-1} \max(T(i), T(n-i))$$

$$\leq n + \frac{2}{n} \sum_{i=0}^{n/2-1} \max(c_i, c_{(n-i)})$$

$$\leq n + \frac{2}{n} \sum_{i=n/2}^{n-1} c_i$$

$$T(n) \leq n + \frac{2}{n} \cdot c \sum_{i=n/2}^{n-1} i \leq n + c\left(\frac{3n}{4}\right) = n\left(1 + \frac{3c}{4}\right)$$

$$T(n) \leq n\left(1 + \frac{3c}{4}\right) = 4n \quad \text{if } c=4 \\ = cn$$

~~can have runtime worse than $O(n)$~~ if we pick multiple bad pivots.
 $T(n) = O(n) + T(n-1) = O(n^2)$

Randomized Quick Sort

RANDOMIZED-QUICK-SORT (A, p, r)

if $p < r$:

$q \leftarrow \text{PARTITION}(A, p, r)$

RANDOMIZED-QS($A, p, q-1$)

RANDOMIZED-QS($A, q+1, r$)

// Recuse on both sides

$$T(n) \leq n + E[X[T(i) + T(n-i)]]$$

$$\leq n + \sum_{i=0}^{n-1} \underbrace{P_r(E_i)}_{1/n} (T(i) + T(n-i))$$

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

Substitution - (like in randomized select)

Guess $T(n) = \Theta(n \log n)$

For all $k < n$ $T(k) = O(k \log k)$

$$T(n) \leq n + \frac{2}{n} \sum_{i=0}^{n-1} i \log i \rightarrow \left[\sum_{i=0}^{n-1} T(i) + T(n-i-1) \right]$$

$$\leq n + \frac{2}{n} \sum_{i=0}^{n-1} i \log i$$

Use:

$$\sum_{i=0}^m i \log i \leq \frac{1}{2} n^2 \log n - \frac{1}{8} n^2$$

$$T(n) \leq n + \frac{2}{n} \sum_{i=0}^{n-1} i \log i \leq n + \frac{2}{n} \left(\frac{1}{2} n^2 \log n - \frac{1}{8} n^2 \right)$$

$$\leq cn \log n + n \left(1 - \frac{c}{4} \right)$$

$$\leq cn \log n \quad \square$$

Expectation vs Amortized

Averaging over random decisions

Averaging over operations

Skip Lists

- Data structure (randomized)
- With high probability (whp) notion
- whp analysis of SEARCH

Skip List

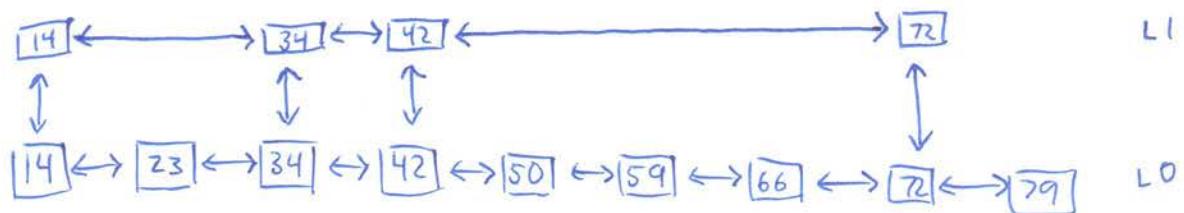
- easy to implement
- Maintain a dynamic set of n elements in $O(\log n)$ time per operation in expectation and whp

One ^{sorted} linked list



Search: $O(n)$

Two sorted linked lists



↳ NYC Subway Express line.

Search(x): (2-linked list)

- Walked right in top linked list (L_1) until going right would go too far
- walk down to bottom list (L_0) until element is found (or not)

Analysis: Search Cost $\approx |L_1| + \frac{|L_0|}{|L_1|}$ ^{only traverse a portion of L_0}

minimize this : when terms are equal

$$|L_1|^2 = |L_0| = n$$

$$|L_1| = \sqrt{n}$$

$$\text{Search cost} \approx O(\sqrt{n})$$

2 sorted lists $\rightarrow 2\sqrt{n}$ dispersed optimally

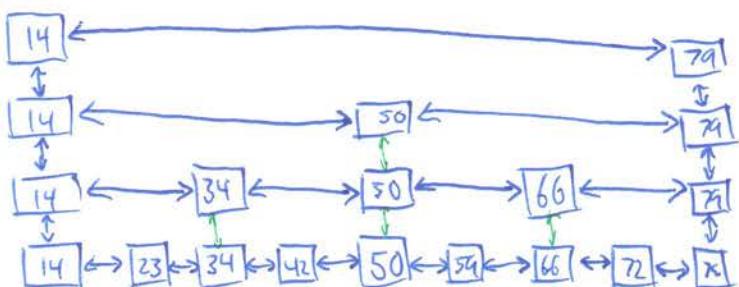
3 sorted lists $\rightarrow 3\sqrt[3]{n}$

k sorted lists $\rightarrow k\sqrt[k]{n}$

$\log n$ sorted lists $\rightarrow \lg n \sqrt{\lg n} = 2 \log n$

* This only works with a static set!

Invariant: If the station i is in L_x . It must also be in $L_{x-1} \dots L_0$



Insert (x)

To insert an element x into skip list:

- Search(x) to see where x fits into bottom list

- Always insert into bottom list (Promote)

Insert into some of the lists above (which ones?)

Delete(x): Search(x) and delete at all levels

Warmup Lemma: # levels in n -element skip list is $O(\lg n)$ whp

Proof: Failure probability (not $\leq c \lg n$ levels)

$$= \Pr(> c \lg n \text{ levels})$$

$$= \Pr(\text{some element } x \text{ gets promoted } > c \lg n \text{ times})$$

$$\leq n \cdot \Pr(\text{element } x \text{ gets promoted } > c \lg n \text{ times})$$

Flip Fair coin
If heads:
promote x to
next level up
and repeat
else: stop

related to $\sum_{k=0}^{\infty} \frac{1}{n^k}$

$$= n \left(\frac{1}{2}\right)^{c \lg n} = \frac{n}{n^c} = \frac{1}{n^{c-1}} = \frac{1}{n^{\alpha}} \quad \alpha = c-1$$

Search: Thm: Any search in an n -element skip list costs $O(\lg n)$ why

Analyze Search backwards

↳ search starts [ends] of node in bottom list

At each node visited

if node was not promoted higher (tails etc)

then we go to came from left

if node was promoted higher

then we go to came from up

- stop [start] when we reach top level (or ∞)

Proof

Backwards search makes "up" moves and "left" moves
each with $Pr = \frac{1}{2}$

moves going up \leq # levels

$\leq c \lg n$ whp (By Warmup Lemma)

Total # moves = # moves till you get $c \lg n$ up moves \uparrow

Claim: Total # moves = # coinflips until $c \lg n$ heads = $O(\lg n)$ whp

□

Chernoff: Let Y be a random variable representing the total # of tails in a series of m independent coin flips where each flip has pr p of heads. Then for all $r > 0$ we have

$$\Pr[Y \geq E[Y] + r] \leq e^{-\frac{2r^2}{m}}$$

Lemma : For any c , there is a constant d
such that whp the #heads in Flipping
 $d \lg n$ fair coins is at least $c \lg n$

$$[d = 8c] \quad E[Y] = m/2$$

$$[r = d \lg n - c \lg n]$$

Randomization: Universal and Perfect hashing

Dictionary Problem: We want an ADT that supports the following operations.

INSERT(item)

DELETE(item)

SEARCH(key)

• items have distinct keys

Hashing

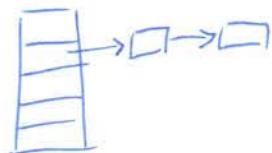
Do these operations in $O(1)$ time and $O(n)$ space

$V = \#$ all possible items

$n = \#$ of current keys in table

$m = \#$ slots in table

Collision Resolution: Hashing with chaining



$\Theta(1 + \alpha)$ time per operation
 $\alpha \leftarrow$ load factor n/m

Simple Uniform Hashing

$$\Pr_{k_1 \neq k_2} [h(k_1) = h(k_2)] = 1/m$$

Like basic quicksort

assuming input keys are random
 \rightarrow solution \rightarrow universal hashing

Universal Hashing

choose random hash function h from H

Universal Hash Family

now we assume h is random
 and keys can be whatever like Randomized quicksort

$$\Pr_{h \in H} [h(k) = h(k') \leq \frac{1}{m} \text{ for all } k \neq k']$$

Theorem:

For n arbitrary distinct keys and for random $h \in H^{\text{universal}}$

$$E[\#\text{keys colliding in a slot}] \leq 1 + \alpha^{n/m}$$

Proof Let $I_{i,j} \stackrel{\text{Indicator Random variable}}{=} \begin{cases} 1 & \text{if } h(k_i) = h(k_j) \\ 0 & \text{otherwise} \end{cases}$

$$\begin{aligned} E[H \text{ keys colliding with } k_i] &= E\left[\sum_{j=1}^n I_{i,j}\right] \\ &= \sum_{j=1}^n E[I_{i,j}] \quad \text{By linearity of expectation} \\ &= \underbrace{\sum_{j \neq i} E[I_{i,j}]}_{\Pr(I_{i,j}=1) = \Pr(h(k_i)=h(k_j)) : V_m} + E[I_{i,i}] \\ &= V_m + 1 = 1 + \alpha \quad \square \end{aligned}$$

Hash Families

Dot product hash family

- assume m prime, $U = m^r$ for int r

view keys in base m : $k = \langle k_0, k_1, \dots, k_{r-1} \rangle$

for key $a = \langle a_0, a_1, \dots, a_{r-1} \rangle$
define $ha(k) = a \cdot k \bmod m$

$$= \sum_{i=0}^{r-1} a_i \cdot k_i \bmod m$$

$$\mathcal{H} = \{ha \mid a \in \{0, 1, \dots, v-1\}^r\}$$

- storing ha in \mathcal{H} requires storing 1 key \Rightarrow computing $ha(k)$ takes $O(1)$

Theorem: dot-product hash family is universal

Proof. Choose any two keys $k \neq k'$, differ in some digit $k_d \neq k'_d$
let not $d = \{0, 1, \dots, r-1\}$

$$\begin{aligned} \Pr_a \{ha(k) = ha(k')\} &= \Pr_a \left\{ \sum_{i=0}^{r-1} a_i k_i = \sum_{i=0}^{r-1} a_i k'_i \pmod{m} \right\} \\ &= \Pr_a \left\{ \sum_{i \neq d} a_i k_i + a_d k_d = \sum_{i \neq d} a_i \cdot k'_i + a_d k'_d \pmod{m} \right\} = \Pr_a \left\{ \sum_{i \neq d} a_i (k_i - k'_i) + a_d (k_d - k'_d) = 0 \pmod{m} \right\} \\ &= \Pr_a \left\{ a_d = -(k_d - k'_d) \mid \sum_{i \neq d} a_i (k_i - k'_i) \pmod{m} \right\} \\ &= \sum_{a_d \in \mathbb{Z}_v} \left[\Pr_{a_d} \{a_d = f(k, k', a_{not d})\} \right] = \sum_x \Pr_{a_d} \{a_d = x\} \Pr_{a_d} \{a_d = f(k, k', x)\} \\ &= E_{a_{not d}} [V_m] = V_m \quad \square \end{aligned}$$

Another universal Hash Family CLRS

- prime $p \geq v$

$$h_{ab}(k) = [(a \cdot k + b) \bmod p] \bmod m$$

$$\mathcal{H} = \{h_{ab} \mid a, b \in \{0, 1, \dots, v-1\}\}$$

Static Dictionary Problem

given n keys to store in table, support $\text{Search}(k)$

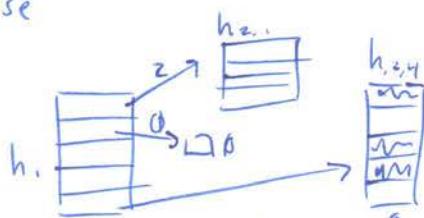
Perfect Hashing - no solutions

- polynomial build time whp

- $O(1)$ time search \rightarrow worst case

- $O(n)$ space \rightarrow worst case

Idea: 2-level hashing



$$\textcircled{1} \quad \text{Pick } h_1 : \{0, 1, \dots, v-1\} \rightarrow \{0, 1, \dots, m-1\}$$

From a universal hash family

For $m = \Theta(n)$ (nearby prime): hash all items with chaining with h_1

\textcircled{2} For each slot $j \in \{0, 1, \dots, n-1\}$

- let $\ell_j = \# \text{items in slot } j = |\{i \mid h_1(k_i) = j\}|$

- pick $h_{2,j} : \{0, 1, \dots, v-1\} \rightarrow \{0, 1, \dots, m_j\}$
From universal hash family

For $\ell_j^2 \leq m_j \leq O(\ell_j^2)$ (nearby prime)

- replace chain in slot j with hashing-with-chaining using $h_{2,j}$

$$\text{Space} = O(n + \sum_{j=0}^{m-1} \ell_j^2)$$

\textcircled{3} If $\sum_{j=0}^{m-1} \ell_j^2 > cn$ redo step \textcircled{1} to guarantee space = $O(n)$
some constant k

Search time = $O(1)$ For first table h_1 ,
+ $O(\max \text{chain size in second table})$

(2.3) while $\overbrace{h_{2,j}(k_i)}^{\text{ANY collision}} = h_{2,j}(k_{i'})$ for any $i \neq i', j$
repick $h_{2,j}$ and rehash those ℓ_j items
↑ guarantees search $O(1)$

\Rightarrow no collisions in second level

3/6/15

6.046 Recitation

Hashing and Dynamic Programming

Universal HashingHash Family $\mathcal{H} = \{h_1, h_2, \dots\}$ Pick h at random from \mathcal{H}
guarantee that

$$\Pr(h(k) = h(k')) \leq 1/m \quad \forall k \neq k'$$

No assumption about input keys

$$\mathcal{H} = \left\{ h_a : h_a(k) = a \cdot k \bmod m \right\}_{\substack{\text{random } 0 \leq a \leq m-1 \\ a \text{ prime}}}$$

Proof that \mathcal{H} is universal: $\Pr(h_a(k) = h_a(k'))$ where $k \neq k'$

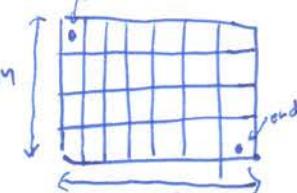
$$k = \{k_0, k_1, \dots, k_{r-1}\} \quad \text{There must exist some index } d \text{ s.t. } k_d \neq k'_d$$

$$k' = \{k'_0, k'_1, \dots, k'_{r-1}\}$$

$$\begin{aligned} \Pr(h_a(k) = h_a(k')) &\equiv \Pr\left(\sum_{i=0}^{r-1} a_i \cdot k_i \equiv \sum_{i=0}^{r-1} a_i \cdot k'_i \pmod{m}\right) \\ \Pr\left(\sum_{i \neq d} a_i \cdot k_i + a_d k_d \equiv \sum_{i \neq d} a_i \cdot k'_i + a_d k'_d \pmod{m}\right) \\ &= \Pr\left(a_d(k_d - k'_d) \equiv \left(\sum_{i \neq d} a_i(k_i - k'_i)\right) \pmod{m}\right) \end{aligned}$$

Dynamic Programming

- ① Define subproblems
- ② Determine relationship b/w subproblems (guessing)
- ③ Determine how to compute final answer given solutions to subproblems

Example

Goal: Get robot to end
can only move
now many ways

 $S[i, j] \rightarrow \# \text{ways of getting from (1,1) to (i,j)}$

$$S[i, j] = S[i-1, j] + S[i, j-1]$$

Base case

$$S[1, 1] = 1$$

$$S[M, N] = \text{solutions}$$

subproblems: $M \times N$
time / subproblem: 1
Total: $O(M \times N)$

Example

n rectangular blocks \rightarrow length l_i , height h_i , breadth b_i .

want to maximize height of a tower subject to constraints
constraint:

IF B_i is stacked on top of B_j , then $l_i \leq l_j$ and $b_i \leq b_j$.

1) Sort blocks in non-increasing order of length and then breadth

$H[i] \rightarrow$ Height of tallest tower with B_i on top

$$H[i] = h_i + \max_{j < i, l_j \geq l_i \text{ and } b_j \geq b_i} H[j]$$

\hookrightarrow guess over all blocks that come before B_i :

Example

Counting Boolean Parenthesizs

True, False, And, OR, XOR

True AND False XOR TRUE

$$\begin{array}{c} \Downarrow \\ (\text{True AND False}) \text{ XOR TRUE} \rightarrow \text{True} \\ \text{True AND (False XOR True)} \rightarrow \text{True} \end{array}$$

Solution: 2

Goal: Find # ways to parenthesize string to evaluate to true.

$T[i, j] \rightarrow$ # ways parenthesizing substring $s[i, j]$ s.t. the expression evaluates to true,

$F[i, j] \quad \text{True} \quad \text{False}$

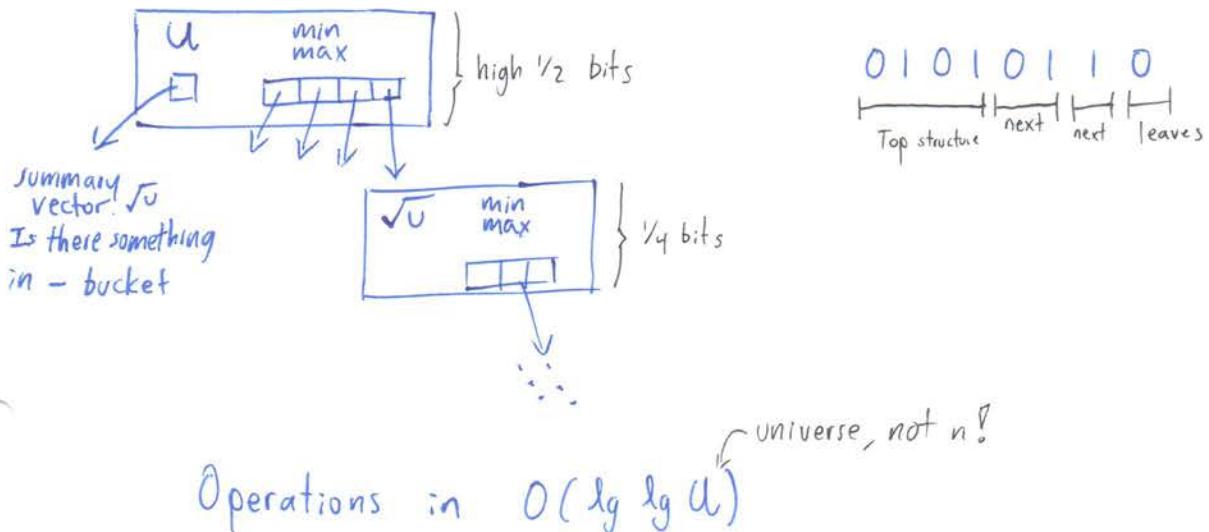
$$T[i, j] = \sum_{k=1}^{j-i+1} \left\{ \begin{array}{ll} T[i, k] \cdot T[k+1, j] & \text{if } s[k] = \text{AND} \\ T[i, k] \cdot F[k+1, j] + F[i, k] \cdot T[k+1, j] & \text{if } s[k] = \text{XOR} \\ T[i, k] \cdot T[k+1, j] + T[i, k] \cdot F[k+1, j] + F[i, k] \cdot T[k+1, j] & \text{if } s[k] = \text{OR} \end{array} \right.$$

Solution $T[1, n]$

6.046 Quiz 1 Review

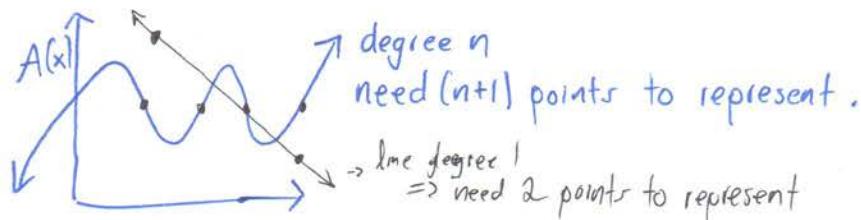
- Topics:
- Divide and Conquer
 - vEB trees
 - FFT
 - B-Trees
 - Amortization and Union Find
 - Augmentation
 - Randomized Algo
 - Hashing
 - Skip lists

vEB Trees



FFT

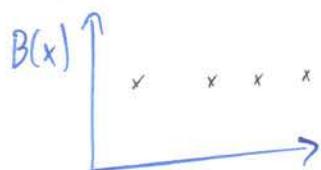
Polynomials: $A(x) = a_0 + a_1x + a_2x^2 + \dots + a_{n-1}x^{n-1}$ (degree $n-1$)



$$B(x) = (\text{degree } n-1)$$

$$A(k) = ??$$

can evaluate in $O(n)$ times
with Horner's Rule.



$$A(x) = \text{degree } n$$

$$B(x) = \text{degree } n$$

$$C(x) = O(2n) \Rightarrow 2n+1 \text{ points}$$

$$A(x) = A_{\text{even}}(x^2) + x A_{\text{odd}}(x^2)$$

where

A_{even} = even terms of $A(x)$

$A_{\text{odd}} = \dots$

Derivation:

$$\sum_{k=0}^{\lceil \frac{n}{2}-1 \rceil} a_{2k} x^k$$

$$A(x)/x \sum_{k=0}^{\lfloor \frac{n}{2}-1 \rfloor} a_{2k+1} x^k$$

FFT picks the $(2n+1)$ th roots of unity as points to evaluate

$$\left\{ e^{j \frac{2\pi i k}{2n+1}} \right\} \text{ for } 0 \leq k \leq 2n+1$$

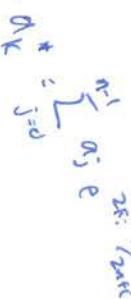
\downarrow squared term $|x|^2$

half the elements

$\frac{2n+1}{2}$ roots of unity

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

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



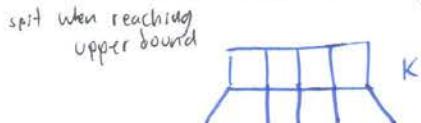
IFFT \rightarrow inverse FFT to go back!

Example Problem: SP 2012 Q1 P4

B-Trees

- Each node can have multiple keys. k
- Each ^{non-leaf} node has $k+1$ children -?
- non root node $B-1 \leq \# \text{keys} \leq 2B-1$
- root node $1 \leq \# \text{keys} \leq 2B-1$
- non leaf node $B \leq \# \text{children} \leq 2B$

- achieves $O(\lg n)$
- Insert, delete, search
- node maintains key in sorted order.



Amortization

Union Find

- Make-Set(x) $O(1)$
- Find-Rep(x) $O(1)$
- Union(x, y) $O(n)$

\downarrow
 $O(\lg n)$ amortized



Randomization

- Quicksort
- Skip Lists
- Hashing
- Matrix Multiplication

Randomized quicksort

- pick a random pivot
- compare to pivot, partition
- Recurse

Skip Lists

- $O(\lg n)$ sorted layer
- every element appears at bottom layer

Insert(x)

- Search for item
- insert at bottom layer
- Flip a coin
 - { Heads \Rightarrow promote to next level and repeat}
 - { Tails \rightarrow return}

Want to ensure whp that there aren't more than $c \log n$ layers

Sample Problems

Worst case running time search: $O(n)$
 \rightarrow every element promoted or none

Hashing

Universal Hashing : choose function $h \in \mathcal{H}$

\mathcal{H} - Universal hash family

$$\Pr_{\substack{h \in \mathcal{H} \\ k \neq k'}} \{ h(k) = h(k') \} \leq \frac{1}{m}$$

set of hash functions

examples

Dot Product Hash Function:

Mult Mod Prime.

- Pick prime p . pick $a, b < p$

$$h_{a,b,p}(k) = [ak + b \bmod p] \bmod m$$

Example Prob: Consider family hash functions $H = \{f, g, h\}$
map $\{0, 1, 2\} \rightarrow \{0, 1\}$

$$f(0) = 1 \quad f(1) = 0 \quad f(2) = 0$$

$$g(0) = 0 \quad g(1) = 1 \quad g(2) = 0$$

$$h(0) = 0 \quad h(1) = 0 \quad h(2) = 1$$

Is this a universal Hash Family?

Yes

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

$$\Pr\{\text{collision}\} = \frac{1}{3}$$

Perfect Hashing

- Know all elements in advance
- $O(1)$ lookups **WORST CASE**
- $O(n)$ space
- $O(n^2 \log n)$ whp create table

Sample Question:

When achieving perfect hashing, there are no collisions at 1st and 2nd levels.

FALSE

Can have collisions at first level but NOT at second.

Monte Carlo Algorithm - Always polynomial - expected correct ^{whp}

Las Vegas Algorithm - Expected polynomial - Always correct ^{whp}

Maximum Subsequence

- 1) Define subproblem
 - 2) Relate subproblem to base case w/ recurrence + base case
 - 3) Recurse and memoize or bottom up DP table
 - 4) Solve original problem
- Analyze runtime
- Input: An strictly increasing subsequence
Output: $A[1..n]$
- Goal: Find longest sequence (not necessarily contiguous)
- Runtime: $O(n^2)$

Building Bridges

Input: $A[1..n]$
Output: Build as many (non-crossing) bridges as possible

Solution = $\max_j \{ L(j) \}$ (like longest increasing subsequence)
 $L(j)$: index of corresponding city on Northern bank.

Runtime: $O(n^2)$

Edit Distance

Input: String $A[1..n]$ & $B[1..m]$
Costs: C_d, C_i, C_m ↪ replacement insertion deletion

Goal: Min cost of transforming $A \rightarrow B$ to transform $A[1..i]$ into $B[1..j]$

Algorithm: $T(i,j) = \min \left\{ \begin{array}{l} T(i-1, j-1) + C_d \\ T(i-1, j) + C_i \\ T(i, j-1) + C_m \end{array} \right\}$

For $j=1$ to c # bottom up approach
compute $(M(j))$

Goal: Pack items in knapsack maximize value.
Input: $A[1..n]$ (repeated items ok)

$M(j) = \max_{i=1..n} \{ M(j-i) + V_i \}$

Input: $A[1..n]$ & $B[1..m]$
Goal: Construct stack of boxes with box i on top.

Solution: Store a back pointer to $M[i-1]$ that was will give items used.

Input: n integers $\{A_1, A_2, \dots, A_n\}$
Goal: Partition integers into S_1 and S_2

Minimize $|S_1| - \sum(S_2)$
 $P(i,j) = \begin{cases} 1 & \text{if some subset of } \{A_1, \dots, A_i\} \\ 0 & \text{otherwise} \end{cases}$

Input: n boxes each with depth d_i
Goal: Construct stack of boxes with box i on top.

$S = (\sum A_i) / 2$
Find min val of S_i :
Junk front $P(n,i) = 1$
 $\Rightarrow \min \left\{ \begin{array}{l} P(1,i) \\ P(2,i) \end{array} \right\} = P(n-1,i) = 1 \right\}$

Input: n boxes with box i on top.
Goal: Max height of stack of boxes with box i on top.

$H(i) = \max \left\{ \begin{array}{l} H(i-1) + h_i \\ \sum_{j=1..i} w_j \end{array} \right\}$
Solution: $\max_j \{ H(j) \}$

Input: n boxes with box i on top.
Goal: Max height of stack of boxes with box i on top.

$H(i) = \max \left\{ \begin{array}{l} H(i-1) + h_i \\ \sum_{j=1..i} w_j \end{array} \right\}$
Solution: $\max_j \{ H(j) \}$

Input: n boxes with box i on top.
Goal: Max height of stack of boxes with box i on top.

$H(i) = \max \left\{ \begin{array}{l} H(i-1) + h_i \\ \sum_{j=1..i} w_j \end{array} \right\}$
Solution: $\max_j \{ H(j) \}$

Input: n boxes with box i on top.
Goal: Max height of stack of boxes with box i on top.

$H(i) = \max \left\{ \begin{array}{l} H(i-1) + h_i \\ \sum_{j=1..i} w_j \end{array} \right\}$
Solution: $\max_j \{ H(j) \}$

Input: n boxes with box i on top.
Goal: Max height of stack of boxes with box i on top.

$H(i) = \max \left\{ \begin{array}{l} H(i-1) + h_i \\ \sum_{j=1..i} w_j \end{array} \right\}$
Solution: $\max_j \{ H(j) \}$

- ea

Augmentation

- easy tree augmentation
- order statistic trees
- level-linked 2-3 trees
- range trees

Idea: Take "off the shelf" data structure and modify to store extra info.

Easy tree augmentation:

- goal: store $f(\text{subtree rooted at } x)$ at each node x in $X.F$
- suppose $X.F$ can be computed in $O(1)$ time from x , children, and $\text{children}.F$
- if modify set S of nodes then costs $O(\# \text{ancestors of } S)$ to update F fields.
 $\Rightarrow O(lgn)$ updates in AVL tree \rightarrow 2-3 trees

* isolated to subtree

Order statistic trees

- ADT / interface:

- $\text{insert}(x)$, $\text{delete}(x)$, $\text{successor}(x)$
- $\text{rank}(x)$: index of x in sorted order
- $\text{select}(i)$: return key of $\text{rank}(i)$

$O(lgn)$ /operation

Augmentation

- use easy tree augmentation with $F(\text{subtree}) = \#\text{nodes in subtree}$

$$x.F_{\text{size}} = 1 + \sum (c.F \text{ for } c \text{ in } x.\text{children})$$

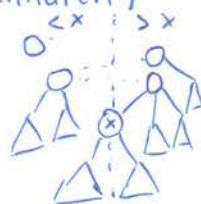
Implementation

$\text{rank}(x)$:

- $\text{rank} = x.\text{left.size}$

- walk from x to root of tree

: If we go left $x \rightarrow x'$: $\text{rank} += x'.\text{left.size} + 1$



constant work/step
 $O(lgn)$ steps

$\text{select}(i)$:

$x = \text{root}$

$\text{rank} = x.\text{left.size} + 1$

→ ; If $i = \text{rank} < \text{return } x$

; If $i < \text{rank} = x = x.\text{left}$

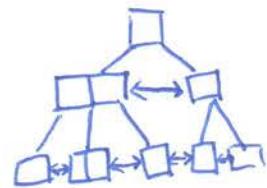
→ ; If $i > \text{rank} = x = x.\text{right}$ $i -= \text{rank}$

repeat

Pro-tip: Only augment things that can be easily maintained.
 NOT
 depth
 rank

Level-linked 2-3 trees

Idea: store level-left and level-right pointers for each node
 - easy to maintain in $O(1)$ overhead

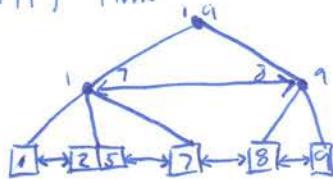


Finger Search Property

search(x from y):

$O(\log |\text{rank}(x) - \text{rank}(y)|)$ time

Store the data in the leaves:



Search - augment to store min and max of subtrees

Insert $\boxed{a|b|c} \Rightarrow \begin{cases} a & b \\ & c \end{cases}$

$\leftrightarrow \begin{cases} | & | \\ | & | \end{cases} \Rightarrow \begin{cases} & | \\ & | \end{cases} \leftrightarrow$

Search(x from y)

$v = \text{leaf. containing } y$

if $v.\text{min} \leq x \leq v.\text{max}$

return regular search in v 's subtree

else if $x < v.\text{min}$: $v = v.\text{level-left}$

else if $x > v.\text{max}$: $v = v.\text{level-right}$

~~*=~~ $v = v.\text{parent}$ // go up

↓??

Orthogonal range search

goal: pre process n points in $d-1$ to support query: given box

Find - # points in box

- k points in box $\rightarrow \Theta(k)$ output

$$O(\lg^d n + |\text{output}|)$$

1D:



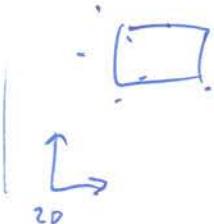
- sorted array

- BST $O(\lg n + k \lg n)$

- level linked $O(\lg n + kc)$



3D
↑



2D
↑

1D range trees

range-query $[(a, b)]$

- search(a)

- search(b)

+ trim common prefix



$\rightarrow O(\lg n)$

nodes and rooted subtrees

2D range search

- 1D range tree on all points by \times

for each node v in x-tree:

store 1D range tree by y on all points in v 's routes sub-tree

$O(\lg^2 n)$ time

Space:

each key lives in $\lg n$ subtrees (all ancestors |

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

B-trees

B-trees

minimum degree t - Each node must have at least $t-1$ keys. \rightarrow all node at least t children
each node can have at most $2t-1$ keys \rightarrow " at most $2t$ children
so if $t=2$ we get a 2-3-4 tree!

$$\text{Height: } h \leq \log_+ \frac{n+1}{2}$$

All leaves are at the same level

3/16/15

6.046 Lecture

Dynamic Programming

- Longest Palindrome sequence
- Optimal Binary Search Trees
- Alternating coin game

DP Notions

- ① Characterize structure of optimal solution
- ② Recursively define the value of an optimal solution based on optimal solutions of subproblems
- ③ Compute value of optimal solution
- ④ Construct optimal solution from computed information

Longest Palindromic Sequence

Given a string : $A[1 \dots n]$ $n \geq 1$

Find longest palindrome that is a subsequence
 Ex: character → character turboventilator → rotator

Strategy $L(i, j)$: length of longest palindromic subsequence $x[i, \dots, j]$ $i \leq j$

```
def L(i, j):
    if i == j: return 1
    if x[i] == x[j]:
        if i+1 == j: return 2
        else: return 2 + L(i+1, j-1)
    else: return max(L(i+1, j), L(i, j-1))
```

→ look at $L[i, j]$ //memoizing
 //don't recurse if $L(i, j)$ already
 //computed

$$\# \text{subproblems} \times \text{time to solve each subproblem given smaller ones are solved}$$

$$\Theta(n^2) \quad \Theta(1) = \Theta(n^2)$$

Optimal BST

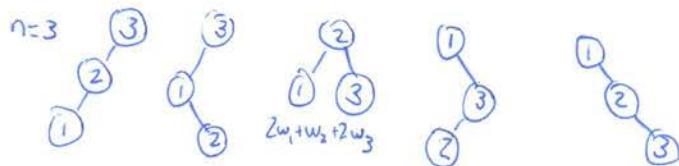
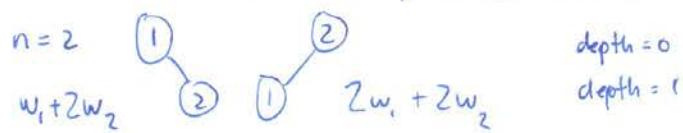
Keys k_1, k_2, \dots, k_n $k_1 < k_2$ wlog $k_i = i$

Weights w_1, w_2, \dots, w_n (search probability)

Find BST T that minimized $\sum_{i=1}^n w_i (\text{depth}_T(k_i) + 1)$

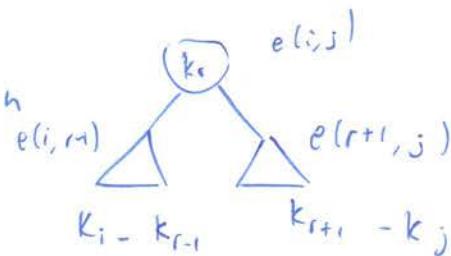
minimize expected search cost

Enumeration exponentially many trees



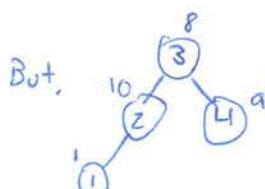
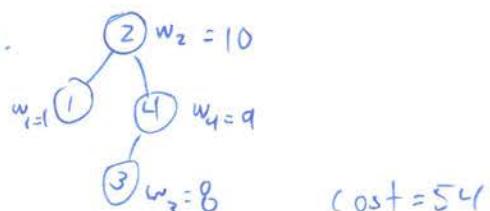
Greedy Solution

Pick k_r in some greedy fashion
↳ highest weight



$e(i, j)$ = cost of optimal BST on k_i, k_{i+1}, \dots, k_j
Counterexample

Greedy algorithm would produce.



is better; cost = 49

So, Fully greedy alg won't work

DP Solution

Have to guess which node will be root

$$e(i, j) = \begin{cases} w_i & \text{if } i=j \\ \min_{1 \leq r \leq j} (e(i, r-1) + e(r+1, j) + w_r + w(i, j)) \end{cases}$$

, sum of all weights from i to j

Alternating coin game

Row of n coins of values V_1, \dots, V_n n even

select either the first or last coin from
remove permanently, receive the value

Basic Strategy

V_1	V_2	V_3	...	V_{n-1}	V_n	n even
-------	-------	-------	-----	-----------	-------	--------

First player: compare $V_1 + V_3 + \dots + V_{n-1}$ to $V_2 + V_4 + \dots + V_n$
 If greater go first
 else go second \Rightarrow never lose

Goal: maximize value AND assuming you move first.

$v(i, j)$ max value we can definitely win if it is our turn and only
coins V_i, \dots, V_j remain

$v(i, i)$ one coin left \Rightarrow pick i

$v(i, i+1)$ pick max of two

$$v(i, j) = \max \left\{ \begin{array}{l} \text{range is } i, j \\ \text{pick } V_i \end{array}, \begin{array}{l} \text{range is } i, j-1 \\ \text{pick } V_j \end{array} \right\}$$

opponent moves
 Need to assume opponent will make
the best move for them

Solution $v(i+1, j)$ subproblem with opponent picking

$$\Rightarrow \text{we are guaranteed the min} \left\{ \begin{array}{l} v(i+1, j-1) \\ \text{opponent picks } V_i \end{array}, \begin{array}{l} v(i+2, j) \\ \text{opponent picks } V_{i+1} \end{array} \right\}$$

$$v(i, j) = \max \left\{ \min \left\{ \begin{array}{l} v(i+1, j-1) \\ v(i+2, j) \end{array} \right\} + V_j, \min \left\{ \begin{array}{l} v(i, j-2) \\ v(i+1, j-1) \end{array} \right\} + V_j \right\}$$

subproblems $\Theta(n^2)$
time / subproblems $\Theta(1)$

- dynamic programming
- matrix multiplication
- Floyd-Warshall algorithm
- Johnson's algorithms
- different constraints

Recall: single-source shortest paths

- given directed graph $G = (V, E)$, vertex $s \in V$, edge weights $w: E \rightarrow \mathbb{R}$
- Find $\delta(s, v) = \text{shortest-path weight } s \rightarrow v \quad \forall v \in V$

<u>situation</u>	<u>algorithm</u>	<u>time</u>
unweighted ($w=1$)	BFS	$O(V+E)$
nonneg edge weights	Dijkstra	$O(V \lg V + E)$
general	Bellman-Ford	$O(VE)$
acyclic (DAGs)	Topological Sort + Bellman-Ford	$O(V+E)$ \hookrightarrow Relax edges in topo sort

All-pairs shortest paths

given $G = (V, E, w)$ Find $\delta(u, v)$ for all $u, v \in V$

<u>situation</u>	- Run SSSP algorithm V times
general	$ V \times \text{Bellman-Ford}$
general	Johnson's

Dynamic Programming

- ① Subproblems
- ② Guessing
- ③ Recurrence
- ④ Acyclic (topo order)
- ⑤ Solve original problem

For m in range $[1, n]$:

For u in V :

For v in V :

For x in V

if $d_{uv} > d_{ux} + d_{xv}$

$w(x, v)$

\nearrow

\nwarrow

\swarrow

\searrow

\nearrow

\swarrow

\nearrow

\swarrow

// Relaxation step

$\Rightarrow O(V^4)$

$$d_{uv}^{(m)} = \text{weight of shortest path } u \rightarrow v \text{ using } \leq m \text{ edges}$$

$$\text{last edge } (x, v)$$

$$d_{uv}^{(m)} = \min_{\substack{\text{for } m=0, \dots, n-1 \\ \text{for } u, v \in V}} \{ d_{uv}^{(m-1)} + w(x, v) \text{ for } x \in V \}$$

$$d_{uv}^{(0)} = \begin{cases} \infty & \text{if } u \neq v \\ 0 & \text{if } u = v \end{cases}$$

$$d_{uv}^{(n-1)} = d_{uv}^{(n)}$$

$$\text{if no negative weight cycles} \rightarrow \text{should not repeat vertices}$$

$$\Leftrightarrow \text{no neg if } d_{vv}^{(n-1)} \text{ is negative}$$

Matrix Multiplication

$$A, B \rightarrow C = AB$$

$O(n^3)$ $O(n^{2.8})$ $O(n^{2.37})$

$$C_{ij} = \sum_{k=1}^n a_{ik} \cdot b_{kj}$$

define $\odot = +$, $\oplus = \min$ (semiring)

$$D^{(m)} = D^{(m-1)} \odot W = W^{(m)}$$

$$W^\oplus = \begin{pmatrix} 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix}$$

$$\hat{D}^{(m)} = \left(d_{ij}^{(m)} \right)_{ij} \quad W = \left(w_{ij} \right)_{ij} \quad V = \{1, 2, 3, \dots, n\}$$

Repeated Squaring: $W^\odot W^\oplus \rightarrow W^\oplus$
 $W^\oplus W^\oplus \rightarrow W^\oplus$ $\Rightarrow O(V^3 \lg V)$ ↳ better than $O(V^4)$
 $W^\oplus W^\oplus \rightarrow W^\oplus$

Transitive Closure: $t_{ij} = \begin{cases} 1 & \text{if } \exists \text{ path } i \rightarrow j \\ 0 & \text{otherwise} \end{cases}$
 $\odot = \text{OR}$ $\oplus = \text{AND}$

ring $O(n^{2.3728})$

Floyd-Warshall algorithm

DP II

① Subproblems $C_{uv}^{(k)} = \text{weight of shortest path } u \rightarrow v \text{ whose intermediate vertices } \in \{1, 2, \dots, k\}$

② guessing: is $k \in \text{path?}$ $E = O(V^2)$

③ $C_{uv}^{(k)} = \min \left\{ C_{uv}^{(k-1)}, C_{uk}^{(k-1)} + C_{kv}^{(k-1)} \right\} \Rightarrow \text{W} \ O(V^3)$ τ_{yay}

$$C_{uv}^{(0)} = w(u, v)$$

```

C = {w(u, v)}
for k = 1, 2, ..., n
  for u in V
    for v in V
      if C_{uv} > C_{uk} + C_{kv}
        C_{uv} = C_{uk} + C_{kv}
  
```

Johnson's algorithm

- ① Find function $h: V \rightarrow \mathbb{R}$ such that $w_h(u, v) = w(u, v) + h(u) - h(v) \geq 0 \quad \forall u, v$
- ② Run dijkstra on $(V, E, w_h) \rightarrow s_h(u, v)$
- ③ claim: $s(u, v) = s_h(u, v) - h(u) + h(v)$

Proof of claim:

consider path p from u to v

$$u = v_0, v_1, v_2, \dots, v_k = v$$

$$w_h(p) = \sum_{i=1}^k w_h(v_{i-1}, v_i)$$

$$= \sum_{i=1}^k [w(v_{i-1}, v_i) + h(v_{i-1}) - h(v_i)] \quad \text{telescope}$$

$$= w(p) + h(v_0) - h(v_k)$$

$$= w(p) + h(u) - h(v)$$

\Rightarrow shortest paths preserved \square

- ① Finding H $w(u, v) + h(u) - h(v) \geq 0$
 $h(v) - h(u) \leq w(u, v) \quad \forall u, v$
- \nearrow
System of difference constraints
 if neg weight cycle

Theorem: IF (V, E, w) has neg weight cycle \Rightarrow no solution to difference constraints

Proof: Consider negative weight cycle $v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_k \rightarrow v_0$

$$\begin{aligned} h(v_1) - h(v_0) &\leq w(v_0, v_1) \\ + \quad h(v_2) - h(v_1) &\leq w(v_1, v_2) \\ &\vdots \\ h(v_k) - h(v_{k-1}) &\leq w(v_{k-1}, v_k) \\ h(v_0) - h(v_k) &\leq w(v_k, v_0) \end{aligned}$$

everything cancels

$$0 \leq w(c) < 0 \quad \square$$

Theorem: if (V, E, w) has no neg-weight cycle then \exists solution

Proof:

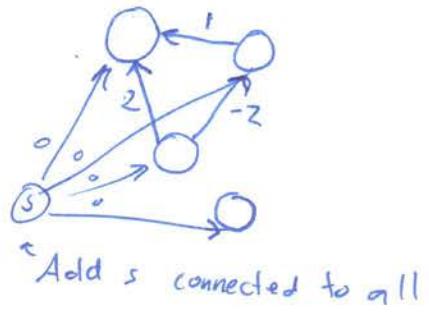
constructive proof

Add s to V

add (s, v) to $E \forall v \in V$

$w(s, v) = 0 \quad \forall v \in V$

$h(v) = f(s, v)$ finite $\forall v \in V$



$$w(u, v) + h(u) - h(v) \geq 0$$

$$w(u, v) + f(s, u) + f(s, v) \geq 0$$

$$f(s, v) \leq f(s, u) + w(u, v) \leftarrow \text{Triangle Inequality}$$

□

$$O(V^2 \lg V + VE)$$

Greedy Algorithms

- MST (recap)
- Process scheduling
- Event overlap
- Fractional knapsack change

Greedy Algorithms

Myopic

Best decision now is the best global decision as well. ← how to prove this?

Minimum Spanning Tree

Graph $G \rightarrow V, E$

Tree that contains all V vertices, and minimizes edge weight sums

Given a vertex subset $S \subset V$

The edge (yu) with min weight s.t. $y \in S$ and $u \in V - S$
must belong to the MST of G .

Proof (by contradiction)

let the MST of graph G by T .

\hookrightarrow Min weight edge between
vertex $y \in S$ and
 $u \in V - S$

Assume, for purposes of contradiction, let T not

contain edge e .

Let's say T contains the edge e_1 instead

see notes

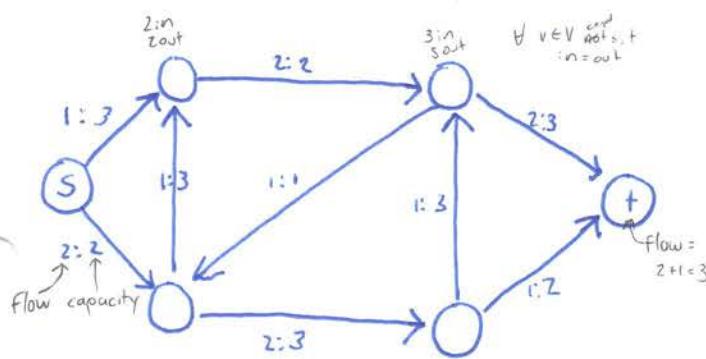
Network Flow

- Flow Networks
- Max Flow problem
- Cuts
- Residual networks
- Augmenting Paths
- Max Flow min cut theorem ↗ Ford Fulkerson algorithm

Flow Network $G(V, E)$ directed graph

- Two distinguished vertices source s , sink t
- each edge $(u, v) \in E$, non-negative capacity $c(u, v)$
- If $(u, v) \notin E$, $c(u, v) = 0$

↪ bandwidth / amount of traffic that can go through edge.



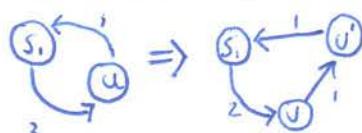
Goal: Increase flow at t by changing flow on other edges. Cannot exceed node capacity

Maximum Flow problem

Given a flow network, & find a flow with maximum value on G

Assumptions

No self-loop edges allowed



Defs

Flow - A flow on G is a function $F : V \times V \rightarrow \mathbb{R}$ satisfying:

capacity constraint: For all $u, v \in V$ $F(u, v) \leq c(u, v)$

Flow conservation: For all $u \in V - \{s, t\}$ $\sum_{v \in V} F(u, v) = 0$

Skew-symmetry: for all $u, v \in V$, $F(u, v) = -F(v, u)$ ↗ implicit summation

The value of Flow F , denoted $|F|$ $|F| = \sum_{v \in V} F(s, v) = F(s, V)$

simple problems

$$f(x, x) = 0$$

$$f(x, y) = -f(y, x) \quad \rightarrow \text{don't double count}$$

$$f(x \cup y, z) = f(x, z) + f(y, z) \quad ; F \times \cap Y = \emptyset$$

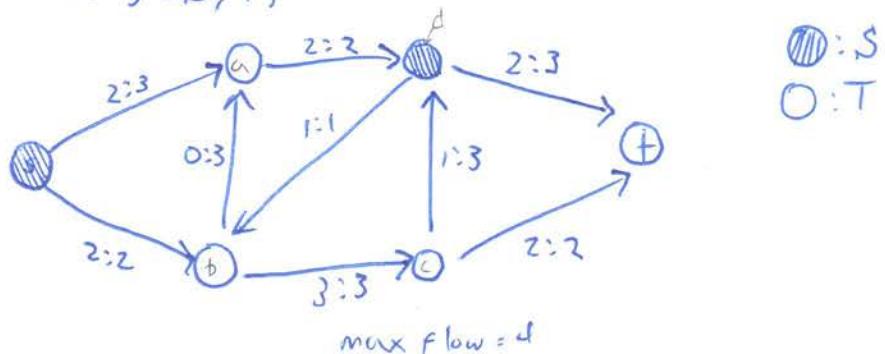
Theorem $|F| = f(V, +)$

$$\begin{aligned}
 |F| &= f(s, V) = f(V \setminus s) - f(V - s, V) \\
 &= f(V, V - s) \quad \text{O} \\
 &= f(V, +) + f(V, V - s - +) \\
 &\quad \quad \quad \text{--- } f(V - s - +, V) \\
 &\quad \quad \quad \text{O by flow conservation} \\
 &= f(V, +) \quad \square
 \end{aligned}$$

Cuts

A cut (S, T) of a flow network $G = (V, E)$ is a partition of V s.t. $s \in S$ and $t \in T$

If f is a flow on G , then the flow across the cut is $f(S, T)$



$$f(S, T) = (2+2) + (-2+1 -1+2)$$

$$\text{(Capacity of cut } c(S, T) = [3+2] + [1+3] = 9)$$

Value of any flow is bounded by the capacity of any cut

Another characterization of flow value.

Lemma: For any Flow F and any cut (S, T) we have
 $|F| = F(S, T)$

Proof

$$\begin{aligned} F(S, T) &= F(S, V) - F(S \setminus S) \\ &= F(S, V) \quad \text{--- } S \text{ does not contain } t \\ &= F(S, V) + F(S \setminus S, V) \\ &= F(S, V) = |F| \quad \text{--- } S \text{ does not contain } t \end{aligned}$$

(Can figure out maximum flow through network by making arbitrary cuts?)

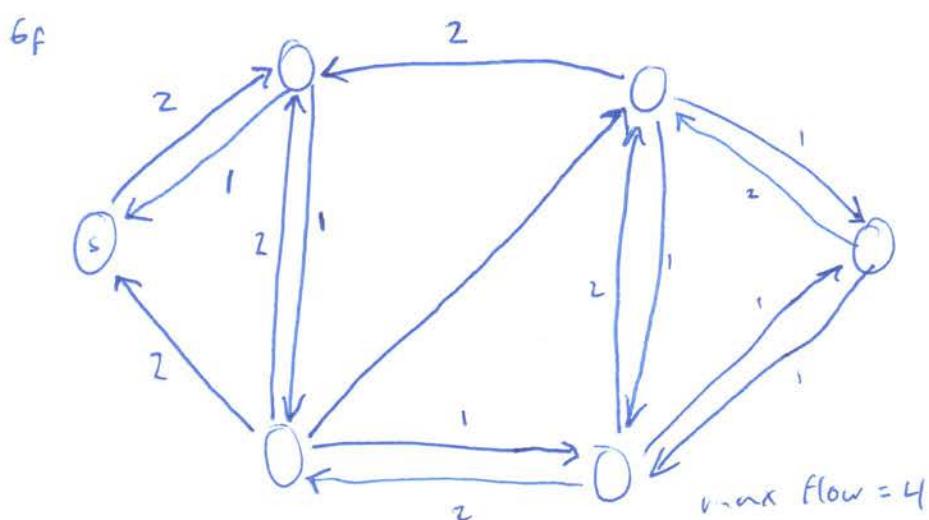
Residual Network $G(V, E)$

$G_F(V, E_F)$: strictly positive residual capacity

$$c_F(u, v) = c(u, v) - f(u, v) > 0$$

edges in E_F admit more flow

If $(v, u) \notin E$, $c(v, u) = 0$ but $f(v, u) = -f(u, v)$



Flow Networks

Capacity constraint: For all $u, v \in V$ $0 \leq f(u, v) \leq c(u, v)$

Flow from one vertex to another, non negative and not exceed capacity

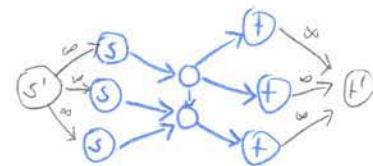
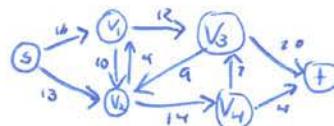
Flow conservation: For all $u \in V - \{s, t\}$

$$\sum_{v \in V} f(v, u) = \sum_{v \in V} f(u, v)$$

↑ source ↑ sink

Flow into vertex must come out of vertex

Flow Problem:

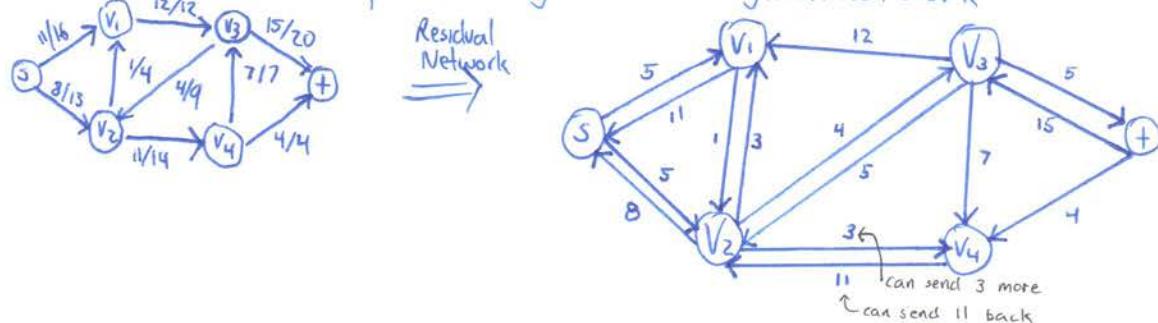


Multiple sources and sinks:

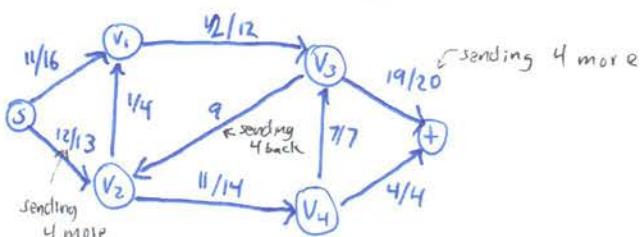
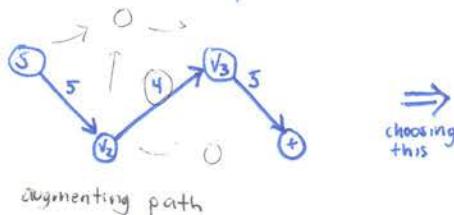
- Create a supersource s' that connects to all sources and $c(s', s) = \infty$
- Create a supersink t' that connects to all sinks and $c(t', t) = \infty$

Ford Fulkerson Method

Residual Network - Roadmap for adding Flow to original network



Augmenting path - path from s to t in the residual network



FORD-FULKERSON(G, s, t)

$$f = 0$$

while there exists an augmenting path P in residual network G_f
augment Flow F along P

return F

Linear Programming

- Examples: politics, flow, shortest paths
- General/Standard Form
- Simplex algorithm - Iterating over slack forms

Politics

How to campaign to win an election?

estimate votes obtained / dollar spent advertising in support of a particular issue

<u>Policy</u>	<u>Demographic</u>			
	<u>Urban</u>	<u>Suburban</u>	<u>Rural</u>	
Build roads	-2	5	3	Want a majority in EACH demographic by spending the minimum amount of money.
Gun control	8	2	-5	
Farm subsidies	0	0	10	
Gasoline tax	10	0	2	
population	\$100,000	200,000	50,000	
majority	50,000	100,000	25,000	

Algebraic Setup

Let x_1, x_2, x_3, x_4 denote money spent / issue

Minimize $x_1 + x_2 + x_3 + x_4$ (n variables)

$$(1) \text{ subject to } -2x_1 + 8x_2 + 0x_3 + 10x_4 \geq 50,000 \quad \leftarrow \text{urban column}$$

$$(2) \quad 5x_1 + 2x_2 + 0x_3 + 0x_4 \geq 100,000$$

$$(3) \quad 3x_1 - 5x_2 + 10x_3 + 2x_4 \geq 25,000$$

$$(\text{m constraints}) \quad x_1, x_2, x_3, x_4 \geq 0$$

optimal solution

$$x_1 = 205,000 / 111$$

$$x_2 = 425,000 / 111$$

$$x_3 = 0$$

$$x_4 = 625,000 / 111$$

x_i : real number

$$x_1 + x_2 + x_3 + x_4 = 310,000 / 111$$

Standard Form for Linear programming

- minimize or maximize linear objective function subject to linear inequalities
 (or equations)

variables: $\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$ objective function: $\vec{c} \cdot \vec{x} = c_1x_1 + c_nx_n$
 inequalities: $A\vec{x} \leq \vec{b}$
 $\max \vec{c} \cdot \vec{x}$, st $\vec{x} \geq 0$

Certificate of Optimality

Is there a short certificate that shows LP soln is optimal?

consider:

$$\frac{25}{222}x(1) + \frac{46}{222}x(2) + \frac{14}{222}(3) \quad \text{equations from page 1}$$

$$\Rightarrow x_1 + x_2 + \frac{140}{222}x_3 + x_4 \geq \frac{3100000}{111}$$

$$x_1 + x_2 + x_3 + x_4 \geq x_1 + x_2 + \frac{140}{222}x_3 + x_4 \geq \frac{3100000}{111}$$

LP Duality

	<u>primal form</u>	<u>dual form</u>
Theorem:	$\max \vec{c} \cdot \vec{x}$ s.t $A\vec{x} \leq \vec{b}$ $\vec{x} \geq 0$	$\min \vec{b} \cdot \vec{y}$ s.t $A^T\vec{y} \geq \vec{c}$ $\vec{y} \geq 0$

Converting to standard form

- 1) Minimize $-2x_1 + 3x_2$ Negate to $2x_1 - 3x_2$ and maximize
- 2) Suppose x_j does NOT have a non-negativity constraint x_{j*}
 x_j replace with $x_j' - x_j''$ $x_j' \geq 0$ $x_j'' \geq 0$
- 3) equality constraint $x_1 + x_2 = 7$ $x_1 + x_2 \overset{e}{=} 7$ $-x_1 - x_2 \overset{e}{=} -7$
- 4) \geq constraint translates to \leq by (-1) multiply

Maximum Flow

$$\max \sum_{v \in V} f(s, v) = |F|$$

skew-symmetry s.t. $f(u, v) = -f(v, u) \quad \forall u, v \in V$

conservation $\sum_{v \in V} f(u, v) = 0 \quad \forall u \in V - \{s, t\}$

capacity $f(u, v) \leq c(u, v) \quad \forall u, v \in V$

Linear constraints

f_1, c_1, f_2, c_2 (two distinct disjoint qts)

f_1, f_2 single capacity

$$f_1(u, v) + f_2(u, v) \leq c(u, v)$$

Shortest path From vertex s

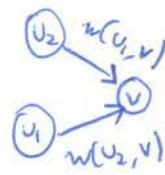
$d[v] = \frac{\text{weight of a path from } s \rightarrow v}{\text{path}}$

$$\min \sum_v d[v]$$

$$\text{s.t. } d[v] - d[u] \leq w(u, v) \quad \forall (u, v) \in E$$

$$d[s] = 0$$

~~strong~~ inequality



$$d[v] - d[u_1] \leq w(u_1, v)$$

$$d[v] = \min ()$$

$$d[v] - d[u_2] \leq w(u_2, v)$$

Simplex Algorithm

Flow: Represent LP in slack form

worst case $\binom{m+n}{n}$ iterations

Convert one slack form into an equivalent slack form whose objective value has not decreased and has likely increased.

Keep going till the opt solution becomes obvious

See notes

Reductions

NP-completeness

Definitions:

P: Set of decision problems D that are polynomial time solvable

NP: Set of decision problems that are verifiable in polynomial time

NP-Hard: " " " " " reducible from any problem in NP

NP-Complete: " " " " " NP and NP-hard

Reductions: Compare the hardness of different problems

Karp Reduction

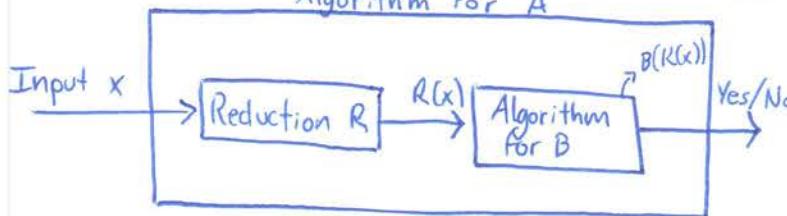
Let A be a problem that takes $x \in X$ as input and spits out $\{0,1\}$ as output

Let B " " " " " "

A is poly-reducible to B if there exists a function

Algorithm for A

R: $X \rightarrow Y$ s.t $A(x) = B(R(x))$



IF B can be solved in poly time \Rightarrow A can be solved in poly time

IF A is "hard" then B is "hard" as well

Reducing Hamiltonian Cycle to Hamiltonian Path

HAM-CYCLE is NP complete. Want to prove HAM-PATH is NP complete too

HAM-CYCLE - Given a directed graph $G = (V, E)$ is there a cycle that visits every vertex exactly once

HAM-PATH - " " " " " is there a path " "

Step 1: Show HAM-PATH is in NP. (can verify in polynomial time)

Input $x \rightarrow$ graph G . Want to prove that \exists $V(x, y)$ that runs in poly time
 $y \rightarrow$ path P that we want to check.

Step 1 \rightarrow Does y touch each vertex exactly once? $O(V^2)$

Step 2 \rightarrow Is every edge in P in G ? $O(VE)$

\Rightarrow HAM PATH IS IN NP

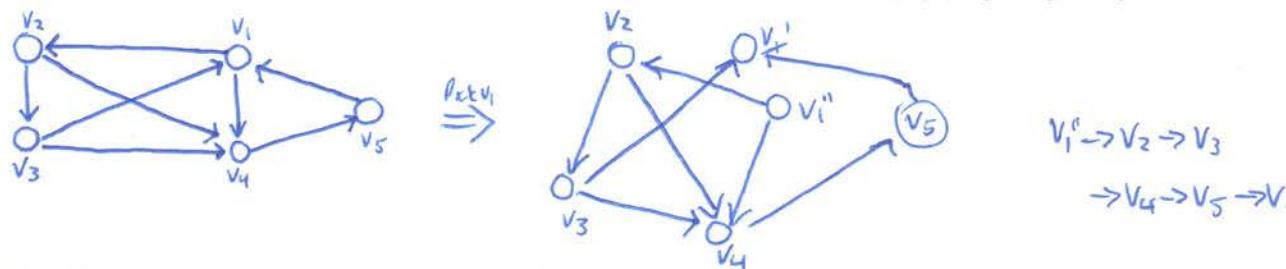
Step 2 - Show HAM-PATH is NP-HARD

Karp Reduction from HAM-CYCLE to HAM-PATH

a) Given an instance of HAM-CYCLE G . Want to create input to HAM-PATH.

Pick some arbitrary vertex $v \in G$ and split it into v^l and v^{ll}

All edges into v are now into v^l and all edges out of v are out of v^{ll}



→ If G has a ham cycle then G' has ham path. Because we can use edges in our original cycle but now our path starts at v_1^{ll} and ends at v_1^l .

← If G' has a ham path then G must have a ham cycle.

The hamiltonian path p in G' must start at v_1^{ll} and end at v_1^l

⇒ HAM PATH is NP-HARD

□

Approximation Algorithms

Definitions

Vertex cover
Set cover
Partition } NP-complete or NP-hard

Approximation Algorithms and Schemes

An algorithm for a problem of size n has an approximation ratio $\rho(n)$ if for any input, algo produces a solution with cost c s.t.

$$\max\left(\frac{c}{c_{\text{opt}}}, \frac{c_{\text{opt}}}{c}\right) \leq \rho(n) - \text{approx}_{\text{algo}}$$

An approximation scheme takes as input $\epsilon > 0$ and for any fixed ϵ , the scheme is a $(1 + \epsilon)$ -approximation algorithm

$O(n^{2/\epsilon})$ running time poly in n , but not necessarily in ϵ

Polynomial time approx
scheme PTAS

Fully PTAS:

poly in n and $\epsilon^{-2} \Rightarrow O(n/\epsilon^2)$

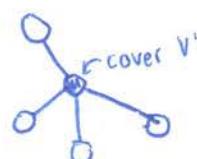
Vertex Cover

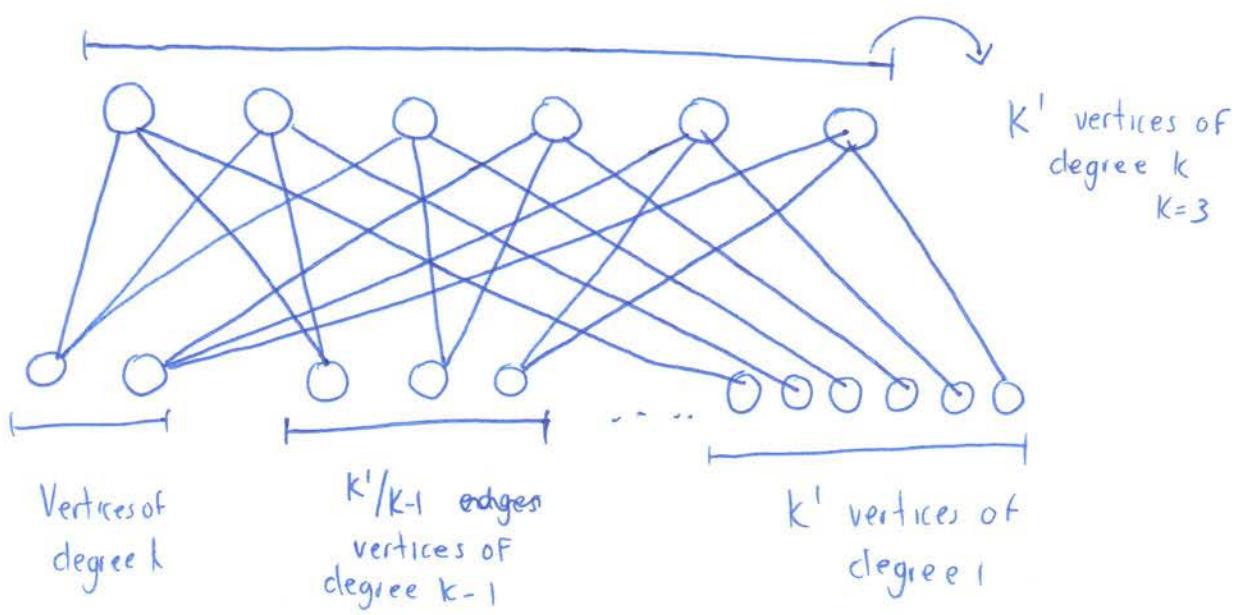
- Set of vertices that cover all edges

Undirected graph $G(V, E)$

Find subset $V' \subseteq V$ s.t. if (u, v) is an edge of G then either $u \in V'$ or $v \in V'$ or both

- Find a V' so $|V'|$ is minimum





Algo could pick all bottom vertices BUT
top: optimal: $= k! = n$

Sol $k! \left(\frac{1}{k} + \frac{1}{k-1} + \dots + 1 \right) \approx k! \log k$
 ↳ still dependent on n !

Approx Vertex cover is a 2-approximation algorithm

Let A denote the edges that are picked

$C = 2|A|$ vertices picked.

Show that $C_{\text{opt}} \geq |A|$ Need to cover every edge including all edges in A .

$C \leq 2C_{\text{opt}}$ Need to pick a different vertex from each edge in A

Approx Vertex Cover

$$C \leftarrow \emptyset \quad E' \leftarrow E$$

while $E' = \emptyset$

Pick $(u, v) \in E$ arbitrarily

$$C \leftarrow C \cup \{u\} \cup \{v\}$$

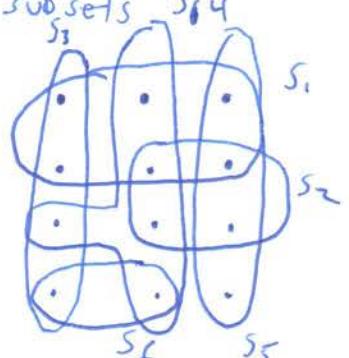
Delete C from all edges incident on u or v

return C

Set cover:

Given a set X and a family of (possibly overlapping) subsets $S_1, S_2, \dots, S_m \subseteq X$ s.t. $\bigcup_{i=1}^m S_i = X$

Find $C \subseteq \{1, 2, \dots, m\}$ while minimizing $|C|$
 s.t. $\bigcup_{i \in C} S_i = X$



opt: S_3, S_4, S_5
 $|S'| = 3$

Approximate-Set-Cover is a $(\ln(n) + 1)$ approx algo

Pick largest S_i , remove all elements in Section X and others S_j

Proof: Assume there is a cover C_{opt} $|C_{\text{opt}}| = t$ let X_k be set of elements in iteration k ($X_0 = X$) $\forall k$, X_k can be covered by t sets

✓ Pigeon Hole principle

\Rightarrow one of them covers at least $\frac{|X_k|}{t}$ elements

\Rightarrow algo picks a set of (current) size $\geq \frac{|X_k|}{t}$

$\Rightarrow \forall k \quad |X_{k+1}| \leq \left(1 - \frac{1}{t}\right) |X_k|$

$|X_k| = 0$, done $\Rightarrow \forall k, |X_k| \leq 1 - \left(\frac{1}{t}\right)^k$, constant rate of shrinking,

cost = $k \quad e^{1/t} n < 1$

(so note)

- Dynamic Programming
- Greedy Algorithms
- Network Flow
- Linear Programming
- Complexity and Approx Algorithms
- All-pair shortest paths

DP Vertex Cover Problem

A vertex cover of an undirected graph is a subset of its vertices s.t. for every edge (u, v) in the graph, either v or u or both are in vertex cover.

- Find minimal vertex cover for a tree

Step 1: $d[v]$ → Minimum vertex cover for a tree rooted at v

Step 2: Guessing: Is vertex v in minimal vertex cover or not?

Step 3: $d[v] = \min \left\{ 1 + \sum_{u \in v.\text{children}} d[u] \right\}$

Step 4: Solution = $d[\text{tree.root}]$

Assume binary tree

subproblems: $O(n)$ Time / Subproblem = $O(1)$ ⇒ Total time: $O(n)$

All-pair shortest paths

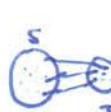
DP $d_{uv}^{(m)}$: Distance to go from $u \rightarrow v$ in at most m edges

$$d_{uv}^{(m)} = \min \left(d_{ux}^{(m-1)} + w(x, v) \right) \Rightarrow O(V^4) \quad \text{or} \quad d_{uv}^{(m)} = \min \left(d_{ux}^{m/2} + d_{xv}^{m/2} \right) \Rightarrow O(V^3 \log V)$$

Ex: Find diameter of weighted graph in $O(V^3)$

Sol: Run APSP, look through all pairs and pick max

Greedy Algorithms

MST: cut property:  → lightest edge that connects S and T will be in an MST

Know Prim's and Kruskal's

Ex: All edges in graph G have unique weights ⇒ MST is unique

Ex: Cycle → edge with the lowest weight is in some MST

→ edge with highest weight is not in an MST

F us

$$\textcircled{1} F(u, v) \leq c(u, v)$$

\textcircled{2} Conservation of flow

\textcircled{3} Flow leaving S = flow entering T

$$\text{max flow} = c(\text{min cut})$$

Ford Fulkerson - slow $O(|E|F)$

Edmon Karp - $O(VE^2)$ ← bound # augmenting paths

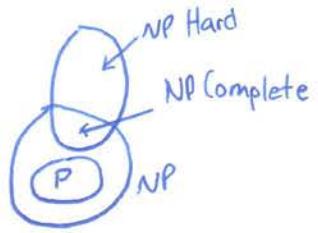
BFS to find augmenting paths to $\min(VE, |F|)$

NP Hard / Complete Problems

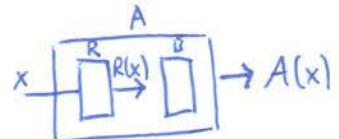
$P = \{ \text{problems solvable in poly-time} \}$

$NP = \{ \text{decision problems with w/ efficient verification algorithm} \} \xrightarrow{\text{poly-time}}$
 $= \{ \text{decision problems solvable in non-deterministic poly-time} \}$

NP-hard = Problem D is NP-Hard if $A \leq_p D \wedge A \in NP$



Reduction: Problem A poly-time reduces to B. $A \leq B \quad A \leq_p B$
 If \exists a poly-time alg R s.t. $A(x) = B(R(x))$



NP-complete: Problem D is NP complete if $D \in NP$, and D is NP Hard

Important Stuff

→ 3 SAT :

- 3DM

- Supermario, 4 port

→ Hamiltonian Cycle / Path

Example

Ham Path = p Ham Cycle

$G = (V, E) \Rightarrow$ Create G'

1) new vertex v to G

2) Add edges to/from all vertices from/to v

Ham Cycle in $G' \Rightarrow$ Ham Path in G

consider $v_i \rightarrow v_{i-1} \rightarrow v \rightarrow v_{i+1} \rightarrow \dots \rightarrow v_n \rightarrow v_i$

rearrange $v \rightarrow v_{i-1} \rightarrow \dots \rightarrow v_n \rightarrow v_i \rightarrow v_{i-1} \rightarrow v$

↑
Hamiltonian Path in G

TODO show other direction

Ham path in G \Rightarrow Ham Cycle in G'

$v_i \rightarrow v_n$

$v \rightarrow v_i \rightarrow v_n \rightarrow v$

Approximation Algorithms

Algorithm for a problem of size n has an approximation ratio of $\varphi(n)$ if it produces a solution s.t. the cost C of the solution is at most $\varphi(n)$ Factor away from optimal.

PTAS - Poly Time Approximation Scheme

Takes in ϵ and a problem instance and produces $(1+\epsilon)$ approx of optimal

$O(n^\epsilon)$

FPTAS - Fully PTAS - Same thing but runtime is poly in $n, \frac{1}{\epsilon}$

$O(\frac{n}{\epsilon})$

Vertex Cover - Repeat:

take any $(u, v) \in E$

add u, v to cover-set

remove all edges from u and v

Distributed Algorithms

Distributed Algorithms - run on multiple machines/processes. Complicated and error prone

- Infinite state interactive state machines

Distributed Networks

$\Gamma(u)$ = set of neighbors of vertex u .

vertex = process

Synchronous Network Model

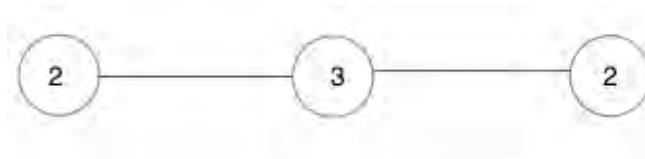
- Each process has out/in ports.
 - processes don't know who its ports channels connect to.
- Processes may be indistinguishable from each other

Synchronous rounds

- in each round
 - each process uses current state to send message to all processes
 - receives message and change state based

6.046 Problem 1-1Collaborators: *Erika Lu, Benjamin Gunby*

- a)** In the following graph, the optimal solution would be to pick both of the restaurants with value of 2 for a total value of 4. This greedy algorithm will however, pick the middle restaurant with value 3, and remove all the neighbors for an incorrect solution with a total.



- b)** We use a dynamic programming approach. Where for each subproblem we consider two scenarios. We either have the restaurant in the final solution or we don't. We'll call these $R(i)$ and $R'(i)$ respectively. The base cases are as follows:

$$\begin{aligned} R(i) &= p(i) \\ R'(i) &= 0 \end{aligned}$$

The recursive solution for each of the equations is:

$$R(i) = p(i) + \sum_{j \in \text{children}(i)} R'(j)$$

$$R'(i) = \sum_{j \in \text{children}(i)} \max(R'(j), R(j))$$

For the first equation, we assume we are including restaurant i so we add its profit plus the sum of the profit we could get if we don't include all of i 's children. For the second equation, we assume we are not including restaurant i so we don't include its profit and add the max profit we could get from each children (we could add each one or not).

We return the max of $R(n)$ and $R'(n)$, where n is the total number of restaurants.

Proof: Since each sub problem (by definition) produces the correct output, our final solution will be correct.

Analysis: Number of subproblems: $O(n)$

Time per subproblem: $O(1)$

Total Runtime: $O(n)$

c) We begin by making an array N of size n where n is the total number of vertices in the graph G . We go through each vertex and count the number of neighbors of that vertex (this is the degree). We place each vertex v in $A[\text{degree}(v)]$. We then check if there are any restaurants in $N[0]$. If there are, we "pick" these restaurants and remove them from the array. Otherwise, we look at all of the vertices in $N[1]$ and pick one at random. We "pick" this restaurant, remove it's neighbor (a) and update all of a's neighbors by shifting them one to the left in the array (effectively updating their neighbor counts). We repeat this process, each time "picking" a new restaurant. The algorithm terminates when there are no elements in $N[0]$ or $N[1]$. At this point, we return all of the picked restaurants.

Proof: Every time we pick a restaurant, we pick the one that will remove the least number of restaurants from consideration. We are essentially always removing leaves from the tree, since we know that the most number of nodes will reside at the deepest level of the tree.

Analysis: The initial construction of N takes $O(n)$ time since we only visit each restaurant once. Since the graph is a tree, we know that there are at most $n - 1$ edges. Since we can only remove each edge once, we can update at most $n - 1$ times. Therefore, the total running time of this algorithm is $O(n)$

d) For each node, we consider two cases and recursively calculate the maximum profit for each. In one case, we don't include this node in our final solution, we have to consider the $n - 1$ remaining nodes. If we do include the node, we have to consider at most $n - 2$ nodes, since we removed the node and its neighbor. The recurrence is as follows $T(n) = T(n - 1) + T(n - 2)$. This is the fibonacci recurrence and it solves to $T(n) = \varphi^n$

6.046 Problem 1-2Collaborators: *Erika Lu, CLRS*

- a)** If a set of requests is on or on the boundary of the same square this means that their Euclidian distance $d \leq 1$. This is because the squares are $\frac{1}{2} \times \frac{1}{2}$ and the maximum distance two points can be apart in this square is equal to 1.
- b)** This algorithm will use a divide and conquer algorithm to correctly compute if any two points are within distance one of each other in $O(n \log n)$ time. In the trivial case where there are only two points, we simply check if their distances are greater than one and return them if they are. Otherwise, we find a vertical line that bisects the set of all points into R_L and R_R such that $R_L = |L|/2 = R_R$ (we can break parity ties by taking the ceil of $|L|/2$ for R_R) where the points in R_R are to the right of the line and the points in R_L are to the left.

For each one of these subsets, we make a recursive call to the algorithm making sure to pass in the correct (updated) L_X and L_Y for each of the subsets. The "combine" part of this algorithm can assume that both the left and right subsets don't have any points distance 1 from each other, if that was the case, the algorithm would have returned before reaching this step. With that assumption, we only need to check the points of both sets that are near our vertical bisecting line, since these are the only points that could possibly be closer than 1 in the entire combined set. In order for two points to be closer than 1 to each other they must be a distance of 2 away from the vertical bisecting line, knowing this, we only have to check points residing in this area. We can create an array consisting of only the points that lie in this area and attempt to find another point that is 1 distance or less apart. If we find a pair, we return it, else, we continue with the merge.

Proof: Since we only compare a finite amount of points on each merge (only the ones that lie one distance away from either side of the dividing line). Since both sides have already been checked, for each point close to the line, we only need to check 7 adjacent squares (using the logic from a). If the algorithm never returns a point each sub problem has been verified to have no close points and the entire set composed of all subsets has also been checked. At this time we can correctly conclude that there are no pairs of points a distance of 1 apart.

Analysis: This algorithm gives the following recurrence: $T(n) = 2T(n/2) + O(n)$ since we split the problem into two problems of half the size and each comparison takes $O(n)$ time. This recurrence, is identical to merge sort, and it solves to $O(n \log n)$

- c)** The algorithm described above can be modified to detect 3 pairs that are a distance of 1 or less apart. We begin by changing our base (trivial) case to when there are only 3 points in the input. In this case, we compare each point against each other to detect if they are

1 or less distance apart. We continue the algorithm as before, but when merging instead of checking 7 adjacent squares we will have to check 15 (2 for each square plus yourself) of them. If we detect three points that are close together, we return them, else we continue with the algorithm. The running time changes only by a constant factor so asymptotically it remains $O(n \log n)$

6.046 Problem 2-1 and 2-2Collaborators: *Erika Lu, Amol Bhave***Problem 2-1)**

a) Description: This naive algorithm checks every possibility of the source string to see if it matches the pattern string. We iterate through the entire source string, starting at index 0, and checking if the substring from 0 to m of the source string matches the pattern string (taking into account that * can be either an 'a' or 'b'). If it matches, we add our current index to a list. We then increase the index i by 1 and check the substring from i to $i+m$. We repeat this process until the algorithm terminates when i is equal to the length of the source string - length of pattern string.

Correctness: Since we are checking every single possibility, the algorithm will always return the correct output.

Analysis: Each check of the algorithm requires $O(m)$ checks, since (in the case of a match), we will have to check every index of the pattern string. We are doing $O(n)$ checks (one for each index of the source string). Therefore, this algorithm will run in $O(nm)$ time.

b) This algorithm takes inspiration from cross-correlation used in signal processing to determine how similar two signals are to each other. We begin by mapping each string to a polynomial. Our chosen mapping, which will be explained later is as follows:

$$A = 1, B = -1, * = 0$$

We chose to represent the coefficient of the term with degree i as the index of the string. For example, the string "aba" would map to $1 - x + x^2$ or $\langle 1, -1, 1 \rangle$ as a vector. To get the cross-correlation, we map each string S and $\text{reverse}(P)$ to their corresponding polynomials using the method described above and get S' and P' respectively. Notice that we reverse the pattern string, before the multiplication. This is because of the inherent nature of convolution where the second polynomial is essentially reversed. By reversing the polynomial before the multiplication, we guarantee that we will get the cross-correlation that we are looking for (after the multiplication). This approach also compares parts of each substring with each other, but we only want the entire pattern string compared to an equally long substring of the source string. Since we want to make sure we compare the entire pattern string, we need to strip away the first and last $m - 1$ terms of our product. Each coefficient (for a term of degree i) in this polynomial represents how similar the two strings are after the pattern has been shifted i times). For example, the polynomial represented by the vector $\langle 2, 0, 2 \rangle$ with a calculated max of 2 suggests that the pattern string matched at index 0 and 2.

It is easy to calculate what the maximum coefficient will be. This happens in the case where we have an exact match between the two compared strings and we get a coefficient equal to $m - \text{num(wildcards)}$. This is because for every character exactly matched we get a $+1$ for the coefficient since $1^2 = 1$ and $(-1)^2 = 1$. So, in the event of an exact match the coefficient will be m . If the pattern string has a wild card (*) we need to account for this by updating our max to m minus the number of wild cards in the pattern string.

Once we have calculated the maximum, we can simply take the "peaks" of the cross-correlation. This is where the strings are most similar, in fact, this is when the strings are an exact match. So we just need to scan through our cross-convolution for coefficient values equal to the max and take note of the term degree they were in. After scanning through, we simply return a list of all of the matches. We chose the following mapping in order to maximize the value of our cross-correlation when two strings are similar and minimize the value when they are not. For example, the cross-correlation of the two strings 'a' and 'b' would be -1 . 'a' and 'a' (or 'b' and 'b') would have a value of 1 because of our chosen mapping. '*' is chosen to map to 0 so that it neither increases or decreases the cross-correlation and essentially remains a 'neutral' state.

Example: $S = "ababbab"$, $P = ab*$. The first step in the algorithm is to convert both of the strings into the polynomial using the mapping described above. S maps to $1 - x + x^2 - x^3 + x^4 + x^5 - x^6$ in vector form, this is $\langle 1, -1, 1, -1, 1, -1 \rangle$. Similarly P maps to $\langle 1, -1, 0 \rangle$. For reasons described above, we need to reverse one of the polynomials, let's say $P' = \langle 0, -1, 1 \rangle$. Multiplying P' and S yields our cross-correlation $C = -x + 2x^2 - 2x^3 + 2x^4 - 2x^6 + 2x^7 - x^8$. In vector form we have $C = \langle 0, -1, 2, -2, 2, -2, 2, -1 \rangle$. The comparison starts comparing the last character of the pattern string with the first character of the source string, leaving the beginning of the pattern string out. Since we are only focused on exact matches with the pattern string, we need to discard the first and last $m - 1$ elements. This leaves us with $C' = \langle 2, -2, 2, -2 \rangle$. The maximum we are looking for is $3 - 1 = 2$ since $|m| = 3$ and there is one wild card. All that is left to do is to scan C' and see when we see a coefficient greater than or equal to this max. This clearly happens at 0 and 2 so the algorithm will return $[0, 2]$.

Runtime: The runtime is equivalent to however long it takes to multiply the polynomials. This is $O(n^2)$ using a simple multiplication approach or $O(n \log n)$ using FFT as described in lecture. See c) for a more descriptive answer.

c) Constructing each of the polynomials takes $O(n)$ and $O(m)$ respectively. Reversing one of the polynomials (m is shorter so we'll reverse that), will take $O(m)$. Since $m \ll n$ this is $O(n)$. Multiplying these polynomials using the FFT algorithm takes $O(n \log n)$ time. We then strip the resulting polynomial in $O(m)$ time and scan through the polynomial in $O(n - 2m)$ time. Asymptotically, this algorithm will run in $O(n \log n)$.

d) We can use the algorithm described in b) with a few modifications in order to use 4

symbols (A,C,G,T). We need to come up with another mapping for these four values that results in a positive number when two matches are multiplied together and a negative (or invalid (to be explained later)) when two mismatches are multiplied. The mapping is as follows: For the source string, we use:

$$A = 1, C = -1, G = i, T = -i$$

and for the pattern string we use

$$A = 1, C = -1, G = -i, T = i, * = 0$$

Notice that the mapping for G and T is interchanged for each of the strings. This is to achieve the desired addition when two matches are multiplied together since $i \cdot (-i) = 1$ and $i \cdot i = -1$. We keep the wild card as 0 for similar reasons as described in b). The algorithm follows as described in b) with the only difference being at the end when we check the final polynomial C' . We know that any coefficient that is not a real number cannot be a match since it was at some point multiplied by a value that was not its match and thus resulted in a non real answer. Since only the mapping of the initial polynomials is different, the algorithm has the same running time and can be proved like b).

Example: $D = ACGACCAT$ and $P = AC * A$. The first step is to match each string to its corresponding polynomial using the specific mappings described above. Doing this we get $D = 1 - x + ix^2 + x^3 + -x^4 + -x^5 + x^6 - ix^7$ and $P = 1 - x - x^3$. In vector form this is $D = < 1, -1, i, 1, i, i, 1, -i >$ and $P = < 1, -1, 0, 1 >$. Multiplying D and the reverse of P yields $C = -ix^{10} + (1+i)x^9 - 2x^8 - ix^7 + 3x^6 - (2-i)x^5 - (2+i)x^4 + 3x^3 - (1-i)x^2 - x + 1$. For this example the max is $4 - 1 = 3$ because $m = 4$ and there is one wildcard. C in vector form is $< 1, -1, (1-i), 3, (-2-i), 3, -i, -2, (1+i), -i >$. Trimming C from the first and last $m - 1$ coefficients gives us $C' = < 3, (-2-i), (-2-i), 3, -i >$ from here, we can scan through C' and see that we have matches at 0 and 3 so the algorithm will return $[0, 3]$.

Problem 2-2)

a) Description: When h_1 and h_2 are the same height, the combine algorithm is as follows: In essence, we join the root nodes of both trees and do an insertion of the key at the new root node. This will take care of any splitting that may be necessary if the new root has too many keys.

Correctness: Both of the initial trees were valid B trees to begin with and since all of the keys in T_1 are strictly smaller than all of the keys in T_2 we are allowed to join the two root nodes and maintain all of the B tree properties with one exception. This occurs when joining both of the root nodes creates a node with more than $2t - 1$ keys. In this case, the regular inserting algorithm will fix the issue, by correctly splitting the node, increasing the height of the entire tree by one. The correctness of this operation follows from the correctness of B tree insertion described in the textbook.

Runtime: This algorithm takes constant time since it is completely independent of the height of the tree. This is because we are essentially "inserting" the new key at the top level and splitting the root node at most once if necessary. The runtime is $O(1)$

b)

Description: This approach is very similar to A except we insert the left b tree and the key at a different level. We start by adding k to the left child (as another key of this node) of the root of T_2 . Since k is less than any key in T_2 we should add this key to the leftmost of the node. We can then combine that node with the root of T_1 and proceed with splitting if necessary.

Correctness: Since both T_1 and T_2 are valid B-tree and we are inserting the shorter tree on the second level of the taller tree, we will end up with all of the leaves at the same level. Additionally, since we are adding a key k to the affected node, we assure that we have at least one more children than we have keys (since T_2 was already a valid B tree)

Runtime: Inserting k and T_1 both take constant time and we are at most going to do a constant amount of splits (1 or two) so the total runtime of this algorithm remains $O(1)$

c)

Insert: Notice that the height of a subtree only changes when a node is moved up a level after a split. This is the basis for our slight modification to the insert algorithm. Anytime we insert a node, it will be inserted at the leaf, so we should set the height of its subtree to 0. Then, if at any point the root goes up a level in the tree (after a split) we create a new node with its height equal to $1 + \text{height of children}$. Since the height of a tree only changes on this occasion, as long we update the augmentation any time this happens, we will have an accurate height for all of the subtrees after any number of inserts.

Delete: Since deletions on B trees only happen at the root (never at the leaves) the height of every sub tree will remain constant regardless of the number of deletions. This is because when the height of the tree will only change at the root of the tree, which does not affect the height of any of the sub trees. It is important to notice that during a deletion we don't swap the entire node, but only the key of that node, so we don't have to worry about updating its height (since the node itself did not move).

d)

Description: This algorithm uses a very similar approach as b), but makes use of the tree augmentation defined in c) to correctly combine T_1 and T_2 of arbitrary height (h_1 and h_2 respectively). The algorithm is as follows: Find h_1 and h_2 by looking at the height augmentation of the respective roots of each tree. We can easily find the larger tree by comparing both heights, we'll refer to this as T_L and the smaller tree as T_S . The next step

is very similar to b), but instead of placing the smaller tree at the second level, we need to insert T_S on the $|h_1 - h_2|$ th level of T_L . We do this by walking down T_L $|h_1 - h_2|$ times. If $T_L = T_2$ we walk to the left every time otherwise, we walk to the right. This is because all of the keys in T_1 are strictly smaller than all of the keys in T_2 . The algorithm follows as described in b).

Correctness: The correctness of this algorithm follows from the correctness of b) with the only difference being that we are inserting the smaller tree at a different level on the larger tree. Since we insert the smaller tree at the $|h_1 - h_2|$ th level of the larger tree, we guarantee that all of the leaves remain at the same level. This is maintained even if the merging requires more than one split, since the height of the entire tree will change by at most one.

Runtime: Looking and comparing the height augmentations takes constant time. Walking down the tree to the appropriate level takes $O(|h_1 - h_2|)$ since we need to take $|h_1 - h_2|$ steps to get to that level. The insertion will take $O(|h_1 - h_2|)$ time since in the worst case we will have to split at every level up to the root node. Overall the runtime is $O(|h_1 - h_2|)$

6.046 Problem 3-1Collaborators: *Erika Lu***a)**

Since we are splitting the data structure into $u^{1/3}$ groups of $u^{2/3}$ numbers each, the original recurrence from CLRS changes to:

$$T(2^m) = T(2^{2m/3}) + O(1)$$

where m is $\lg u$. We will keep the original indexing (changes explained later), so we do a constant amount of work per level. Like the original recurrence, we use changing variables: Let $S(m) = T(2^m)$, now we have:

$$S(m) = S(2m/3) + O(1)$$

By case 2 of the master method, the recurrence solves to $S(m) = O(\lg m)$. Substituting back we see that $T(2^m) = O(\lg \lg u)$. Notice that the asymptotic runtime remains the same as when we divided the structure into $u^{1/2}$ groups of $u^{1/2}$ numbers each.

Operations on the structure remain relatively unchanged, we the only difference being in the *HIGH*, *LOW*, and *INDEX* functions since these are the only that are dependent on the size of the groups. The functions are as follows:

$$\text{high}(x) = x/u^{2/3}$$

$$\text{low}(x) = x \bmod u^{2/3}$$

$$\text{index}(x, y) = x \cdot u^{2/3} + y$$

b) To implement the desired changes, we need to make sure to modify the pseudocode to exclude the maximum element from lower-level vEB structures. In essence, we can accomplish this by mirroring (in each function) any operations that have to do with the max, the rest remains the same (as explained in CLRS 20.3) with some minor edge cases that will require special attention. The functions are as follows:

INSERT: Insert remains the same but we mirror $x < V.\min$ with a new case for $x > V.\max$. This means that after line 4 (in the CLRS pseudo code), if $x \notin V.\max$ we want to exchange x with $V.\max$ (mirroring what we do if $x \in V.\min$). If $V.\min$ is equal to $V.\max$ (this happens when there is only one node in the tree) we set y equal to the current $V.\min$ and then update $V.\min$ a $V.\max$ with the min and max of x and y respectively. We need to return after this, effectively skipping lines 3-11 on CLRS Pseudo Code

SUCCESSOR: Also remains the same, but we mirror when min is not found in the predecessor. The original implementation of **successor** assumes that the max is included in the lower levels. Since this is no longer the case, we need to do a check to see if the max is the successor before returning NIL in line 13.

DELETE: Again we use the same implementation, mirror cases from the min to the max and account of some edge cases. In line 13, we add the following check: if $x \in V.\max$: exchange x with $V.\max$. Additionally, if x not equal to NIL after line 11. we need to exchange x with $v.\max$. We also need to consider the case when the element we are trying to delete is the current maximum. In this case, we need to replace the max with the predecessor of the current max before deletion.

Since for each function, we “mirror” anything we do to the minimum to the maximum our runtime will at most double. Additionally, one of the additional checks for edge cases will be large enough to affect the overall runtime. Asymptotically, the runtime for all of the functions will remain the same.

6.046 Problem 4-1, 4-2Collaborators: *Erika Lu***Problem 4-1**

a) The proposed data structure consists of a regular FIFO queue call it Q , and a doubly linked list with pointers to both the head and the tail, call this M . Q keeps track of the $ENQUEUE$ and $DEQUEUE$ operations and M is used to keep track of the minimum,

b) The operations are as follows:

$ENQUEUE(x)$ First, push x unto Q as any FIFO queue. Next we insert x onto the end of M , and check the element adjacent to it, call this j . If $x < j$ we pop off j from M by setting $j.next = j.next.next$. We repeat this process, checking the next adjacent element until we reach the end of the list (no adjacent element) or $x > j$. Note that x will never equal j since all elements are distinct.

$DEQUEUE(x)$ First, pop the first element from Q and return it. To update M we need to look at the head of the list and check if it is equal to x . If so, we remove the element from M . If the head of the element is not equal to x , no updating of M is required.

$FIND-MIN(x)$ Return the end of M by following the pointer. No updating of M or Q is required since this operation does not mutate the data structure.

c) $ENQUEUE$ and $DEQUEUE$ both rely on the original FIFO queue implementation of Q with minor changes to correctly modify M . Therefore, this proof will be focused on why $FIND-MIN(x)$ will always return the correct answer.

Invariant 1: When an element x is QUEUED all elements greater than x will never be the minimum in the queue.

Proof: When we add an element x , we add it to the end of M and then "bubble" it up until any elements that were bigger than x are popped from M . It is safe to remove these elements, because they will never be the minimum of the queue. This is because Q is a FIFO queue and there are only two ways for the min to be updated:

Case 1: A smaller element is QUEUED. If an inserted element x is smaller than all elements in Q it will bubble up to the head of M through the $ENQUEUE$ operation. The minimum is updated.

Case 2: The current minimum is DEQUEUED. Since Q is a FIFO queue, the only way for the

m to be dequeued is for all of the other elements inserted before m to be dequeued first. Because of this, we can safely pop any elements bigger than m from M when ENQUEUEing m since once m is inserted the other elements will never be the min of the queue.

Invariant 2: The minimum of the queue is always located at the head of M .

Proof: This follows from the implementation of insert and is maintained through DEQUEUE() operations. Assume DEQUEUE returns x . There are two cases:

Case 1: x is at the head of M . This means that we are dequeuing the current minimum so M will need to be updated. We do this by removing the head (min) of M effectively setting the new head (new minimum) to the next element in the doubly linked list. This element will be the second smallest element in the queue because of INSERT's implementation (it bubbled all the way to the second spot, but was larger than the absolute min). The new minimum of the queue will still be located at the head of M .

Case 2: x is not at the head of M . This means that x was not the minimum of the queue. In this case, DEQUEUE simply removes the element from Q . The only way for x to not be at the head of M is that a smaller element was inserted after x . Because of the way insert is implemented, this operation would have already removed x from M . Therefore, we can conclude that x is not in M and M is correct as is, with the minimum still at the head.

In both cases, the minimum remains at the head of M , maintaining the invariant.

The correctness of all operations follow from the above mentioned invariants.

d) Overview: The worst case for *ENQUEUE* occurs when we add m elements in ascending order and then insert x where x is smaller than all elements. This last operation will take $O(m)$ time since x will have to be "bubbled" up all the way to the head m times. This operation however, must be preceded by the addition of m elements in $O(1)$ time so the total amortized time remains $O(1)$. In addition, the next operation (after the worse case) will take constant time since all of the previous elements were removed. The worst case for *DEQUEUE* is constant time since we only pop an element from Q , compare it to the head of M and remove it if necessary. This takes constant time. Similarly, *FIND-MIN* will take constant time regardless of the input since we simply look at the head of M .

Analysis: We use the charge method to complete the analysis. When we enqueue an element we charge 1 unit to the present (for the insertion) and reserve 1 unit for the future for a total of 2 unit cost of insertion. When we dequeue an element, any popping that is necessary by the insertion can be charged to the past of the element being removed. This will only happen once, since an element can only be removed once. Find-Min takes constant time regardless of the previous inputs. Overall, the running time of a combination of m of these operations will be $O(m)$. With this analysis we conclude that *ENQUEUE*, *DEQUEUE* and

FIND-MIN all take O(1) amortized time.

Problem 4-2

- a) Imagine an array of size m divided into 4 equal segments. Assuming a random pivot there is a probability of $1/2$ that this pivot will lie on either of the middle segments. If we pick a pivot in these segments, the size of a sub-array can be at most $3m/4$ (this would occur when the pivot point is at the edge of either of these segments, creating sub-arrays of size $3m/4$ and $m/4$). The probability that x_i is a pivot is $1/m$ since each element has an equal chance of being chosen as the pivot. The probability of either of these events happening is at least $1/2$ since the first event has a probability of $1/2$
- b) Using the claim in the prompt we know that if we flip a coin $3(\alpha + c)lgn$ times, then with a probability of $1 - \frac{1}{n^\alpha}$ we will get at least $c lgn$ heads. For this problem, we'll say that heads means that we picked a pivot in the middle two segments which creates a sub array of size no greater than $3m/4$. As proven in a, the probability of this is at least $1/2$ so choosing heads is appropriate. Now, each heads essentially takes the subarray containing x_i and divides into two subarrays, one of which will have size of at most $3m/4$. We are interested in the number of times we would need to get heads in order for quick sort to reach its last level of recursion, that is, when the subarray is of size 1 and only contains x_i . This would take $\log_{4/3} n$ "levels" or "heads" to accomplish. The total number that x_i will be compared with the pivot is $d lgn$ for some constant d to be derived next. Again, from the claim we know that we can get at least $c lgn$ heads if we flip a coin $3(\alpha + c)lgn$ times, then with a probability of $1 - \frac{1}{n^\alpha}$. Since we need at least $\log_{4/3} n$ heads we can set $\log_{4/3} n = c lgn$. We now have $c = \log_{4/3} n / lgn$. Using change of base, we find that $\log_{4/3} n = \frac{\log n}{\log(4/3)}$. Substituting this back into c we get:

$$c = \frac{1}{\log(4/3)}$$

Since we are interested in a probability of at least $1 - \frac{1}{n^2}$ we let $\alpha = 2$. Plugging this back into the equation in the claim we get a total number of $3(2 + \frac{1}{\log(4/3)})lgn = d lgn$. We conclude that

$$d = 3(2 + \frac{1}{\log(4/3)})$$

- c) If we take into account all $x_i \in A$, the probability that the number of comparisons will be at most $d'n lgn$ is upperbounded by the sum of the individual probabilities by the Union bound. This is equal to the opposite, 1 minus the probability that at least one event is not at most $d'n \log n$ comparisons. We have:

$$P(\text{not}) = 1 - (1 - 1/n^2) = 1/n^2$$

Adding all of these n individual event gives us (by the Union Bound) an upperbound of $n(1/n^2) = 1/n$. We are interested in the probability that all of the events are at most

$d' n \log n$ comparisons, we can get this using the above calculations and get $P = 1 - (1/n)$. Since we used this same answer in b) to calculate the probability we can conclude that

$$d = d' = 3\left(2 + \frac{1}{\log(4/3)}\right)$$

d) We can generalize our results from b and c to obtain a bound on the number of comparisons made by QUICKSORT that holds with probability $1 - \frac{1}{n^\beta}$. Notice that from b) to c) the lower bound of our probabilities changed from $1 - \frac{1}{n^2}$ to $1 - \frac{1}{n}$. Continuing this pattern if we want our probability to be $1 - \frac{1}{n^\beta}$ we need to begin with the individual probabilities being:

$$1 - \frac{1}{n^{\beta+1}}$$

. Plugging these into our original solution for be, we get our bounds equal to:

$$3\left((\beta + 1) + \frac{1}{\log(4/3)}\right)$$

6.046 Problem 5-1, 5-2Collaborators: *Erika Lu***Problem 5-1**

- a) The following is modeled after a simple Touring state machine. We start by comparing x to k . If $x < k$ we will be trying to move to the right. Otherwise, we will be moving left. Without loss of generality, this description will assume that we are moving to the right. The first state tries to go up if it can, otherwise it goes to the right until going to the right would be too far. At this point we switch over to the second state.

In this state we begin to look for k and follow the standard search algorithm except we start from our current node. We do this by trying to go to the right, if this is to far we go down and repeat. When we are at level 0, we go to the right until we find k . If we see a key greater than k at the bottom level, we can return that k does not exist.

The code is as follows:

```
FINGER-SEARCH(x, k):
    if (x.value == k) return x
    if (x.value < k) direction = 1 //Move right
    else direction = -1 //Move left

    state1 = true; //Start in state1
    state2 = false;
    currentNode = x;

    while (state1)
        if (currentNode.upParent)
            currentNode = currentNode.upParent //Go up if possible
        elif (currentNode.next.value > k)
            state1 = false
            state2 = true
            break
        elif( currentNode.value == k) return currentNode
        elif (currentNode.value == -INFINITY or currentNode.value == INFINITY)
            return None
        else currentNode = currentNode.next(direction)
    while (state2) //Go down
        if (currentNode.value == k) return currentNode
```

```

        elif (currentNode.bottomLevel) currentNode = currentNode.next(direction)
        elif (currentNode.next(direction) < k)
            currentNode = currentNode.next(direction)
        elif (currentNode.next(direction) > k) return None
        elif (currentNode.value == -INFINITY or currentNode.value == INFINITY)
            return None
        else currentNode = currentNode.downChildren
    
```

This algorithm is essentially a search backwards followed by a search forward so it's correctness follows from the correctness of the original search procedure for a skip list. The runtime can be improved however if m is close to k . This is because we only touch nodes towards k and never have to touch any nodes to the left/right of m . This reduces the runtime to $O(lgm)$ where m is the difference in ranks of m and k . This runtime is guaranteed to run in $O(lgm)$ steps with high probability since (as previously mentioned) it is simply a backwards search followed by a forwards search. The original skip list runtimes and probabilities remain the same.

b) To support efficient rank search on skip list we augment each node with its number of left and right children as defined below: Left children: All nodes below the parent node and to the left of the node all the way until the next node to the left of the parent on the same level.

Right children: All nodes below the parent node and to the right of the node all the way until the next node to the right of the parent on the same level.

Slight modifications are needed to maintain the augmented data. The changes are as follows:

SEARCH(): Search remains the same as no data is mutated. The usual high-probability order-of-magnitude time bounds remain the same since no changes are made.

INSERT(x): Insert remains mainly the same with some additional book keeping on each step to maintain the augmented data. When we reach the bottom level we start counting the number of spots traveled, once we find a spot for x , this count will be the number of left children for x . We keep going to the right (and count) until we find y where y is the last node that does not have an "up pointer". This count will be the number of right children for x . We now follow the next points from x in both directions (after each of x 's random promotions) and update its neighbors children by subtracting $x.leftChildren$ and $x.rightChildren$ from each neighbor respectively. The algorithm remains relatively the same, with the extra step of counting the right children. This is however, still $O(lgm)$ whp because it is essentially the same amount of work (asymptotically) as **INSERT(y)** where we have to walk $x.leftChildren_x.rightChildren$ spots to the right before getting to Y . The high probability runtime follows from the original Insert procedure in Skip lists.

DELETE(x): The modification to delete is similar to that of Insert, but when we delete a node, we need to make sure we add all of its children to their respective neighbors. This is done the same way we subtracted from the neighbors when inserting an element. Like Insert(x) the high probability runtime remains unchanged since the amount of extra work hides under a constant factor in the asymptotic notation of the runtime.

- c) Rank search is very similar to range search, but we want to search for an element that is r elements away from x in faster than $O(r)$ in fact, we want this algorithm to run in $O(\ln(r + 1))$ with high probability.

```
RANK-SEARCH(x, r)
    state1 = true //Go to the right and up
    state2 = false //Start going down
    currentNode = x;
    remainingRank = r
    while (state1)
        numRightChildren = currentNode.next(1).children.length
        if (currentNode.upParent)
            currentNode = currentNode.upParent // move up if possible
        elif (currentNode.next(1).children().length > remainingRank )
            //Going to the right would be too far!
            state1 = false;
            state2 = true;
            break;
        elif ( numRightChildren < remainingRank)
            //update remaining rank and move the the right
            remainingRank = remainingRank - numRightChildren
            currentNode = currentNode.next(1)
    while (state2)
        if (remainingRank == 0) return currentNode
        else
            currentNode = currentNode.next(1) //Move right
            remainingRank = remianingRank - 1
```

Problem 5-2

- a) The following greedy algorithm uses a min heap of pointers to the original set of n prizes. The elements in the min heap represent the elements with the highest prize value. We start by adding the first m elements of the set to a min heap, call this M . We then scan through the remaining $n - m$ elements of the list. We compare each of these elements i to the min of M . If $i > \min(M)$ we pop the min from M and insert i into M . We do this because a larger element has been found, and M is supposed to contain the m prizes with largest

value. Otherwise, we just move to the next element and repeat. In the end, we will have a heap of size m (since we started with m elements and every time we add a new one we also remove the current min), these are the prizes with maximum value, and we should chose these prices.

To maintain the original order, we can maintain a heap of pointers to the original set and use those values for the heap operations. At the end, we can follow the points and "mark" each element in the original set. We can then go through the set (in order) and add all elements that are marked to a new "final" set. We then return this final set. These operations take $O(n)$ and won't affect the asymptotic running time analyzed below.

Creating the min heap of size m will take $O(m)$ time. After that, we need to do $n - m$ min operations and $O(n - m)$ inserts for a total of $O((n - m)lgm)$ time. In total we get $O(m + (n - m)lgm) = O(nlgm)$. For space, the original set takes $O(n)$ while the heap takes $O(m)$ for a total of $O(m + n)$

It's worth noting that a simple DP algorithm could also be used, but it's running time would be $O(mn)$.

Another greedy solution involves a variation of the median-finding algorithm and runs in $O(n)$. We begin by modifying the media finding algorithm dicussed in lecture to find the element with $rank = m - 1$ call this x . Instead of $m/2$. We then scan through the original set and add any prize with a higher value than x . After the scan completes we will have selected the m prizes with the highest value from the set. This is guaranteed to maximize the total prize value. This algorithm runs in $O(n)$ as the median finding algorithm runs in $O(n)$ and the scan takes $O(n)$.

b) Since a simply greedy algorithm is not guaranteed to produce the most optimal result, the following algorithm uses Dynamic programming and will run in $O(mn)$ time. The sub problem $S_{j,q,A}$ means the max amount of value you can get if there are a total of j prices and you are allowed to pick at most q . Additionally, the last prize taken was of type A. Similarly $S_{j,q,B}$ means the max amount of value you can get if there are a total of j prices and you are allowed to pick at most q . Additionally, the last prize taken was of type B. The DP recurrence is as follows:

if $p_j.type == A$:

$$S_{j,q,A} = \max(p_j.val + S_{j-1,q-1,A}, S_{j-1,q,A})$$

$$S_{j,q,B} = S_{j-1,q,B}$$

elif $p_j.type == B$:

$$S_{j,q,A} = S_{j-1,q,A}$$

$$S_{j,q,B} = \max(p_j.val + S_{j-1,q-1,A}, p_j.val + S_{j-1,q-1,B}, S_{j-1,q,B})p_j$$

For equation 1, we can either take p_j or not take it. For equation 2 we can't take p_j since we currently picked a prize of type B. Similarly for Equation 3, we cannot take the prize. For the last equation we can either take the B for the first time, take it (not for the first time), or don't take it.

The base cases are: $S_{j,j,(A|B)} = \text{sum}(pj)$ - add all prizes since we can pick all of them.

Our final solution will be $\max(S_{n,m,A}, S_{n,m,B})$. The overall runtime is $O(nm)$ since there are $O(n)$ subproblems and each takes $O(m)$.

c) We use the same $S_{j,q}$ subproblems as in b) (ignoring type) but now we have to consider more sub sequences when making a selection. The recurrence is as follows. The base cases are: $S_{j,j} = \text{sum}(pj)$ - add all prizes since we can pick all of them. The recurrence becomes:

$$S_{j,q} = \max(S - j - 1, q - 1, S - j - 2, q - 2, \dots, S - 1, 1)$$

with a very important exception that we can only consider a case if the previous considered case is less than the current value under consideration. Our final solution will be $S_{n,m}$. The overall runtime of this dynamic programming solution is $O(mn^2)$ since each sub problem takes $O(n)$ to solve and there are $O(nm)$ subproblems.

d) This approach uses DP to find an optimal solution in $O(mn^2)$. We start by the node of the tree since (by the constraints of the problem) we have to take. We then consider the following cases:

Take m elements from the left, 0 elements from the right Take m-1 elements from the left, 1 elements from the right take 0 elements from the left, m elements from the right

We start at the bottom and memorize, working our way up the tree until we have our solution. The base cases are trivial, where the max value you can get at a leaf is equal to the value of the leaf. Define $S(m, i)$ as the max prize amount when we take m elements from the left and m-i elements from the right. The recurrence is as follows:

$$S(m, i) = \max(S(m - i, i), S(m - i - 1, i - 1), \dots, S(1, 1))$$

The final solution will be the the max of all $S(m,i)$ for all i from 0 to m.

The overall runtime of this dynamic programming solution is $O(mn^2)$ since each sub problem takes $O(n)$ to solve and there are $O(nm)$ subproblems.

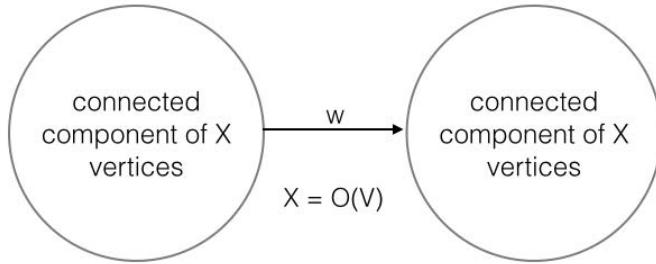
6.046 Problem 6-1, 6-2Collaborators: *Erika Lu, Patrick Liscio***Problem 6-1**

a) Case 1: $r = w_{i,j}$ Replacing an edge with another one of equal weight will not make any noticeable changes to D or Π . Therefore, no action is necessary. In this case the algorithm runs in $O(1)$ time.

Case 2: $r < w_{i,j}$ In this case we have to look into the distance matrix D and compare $D(v \text{ to } i) + r + D(j \text{ to } v)$ to the $D(u \text{ to } v)$. If an update is required (when the distance is less) we update the distance matrix accordingly. Furthermore, we can update the predecessor matrix with the predecessor matrix of j to v . The running time of this algorithm is $O(V^2)$ since we look at all V^2 pairs.

Case 3: $r > w_{i,j}$ In this case, we will have to run Floyd Warshall again to update D and Π . This will take $O(V^3)$ time.

b) In the following example, changing the center edge from weight w to weight r where $r \neq w_{i,j}$ would force the algorithm to look at all pairs of points taking $\Omega(V^2)$ time. In the example below, the triple would be (i, j, r) where i is the vertex with an outgoing edge to the other set of connected components and j is the vertex this edge is adjacent to. r can be any non negative integer as long as $r \neq w_{i,j}$. G and D follow from the picture below:



c) The modification to Floyd-Warshall algorithm is as follows:

Subproblem $c(u, v, k, p) =$ weight of shortest path u to v whose intermediate vertices $\in 1, 2, \dots, k$ in p number of hops.

Recursion:

$$c(u, v, k, 0) = 0$$

$$c(u, v, 0, 1) = w(u, v)$$

$$c(u, v, k, p) = \min(c(u, v, k-1, p), c(u, k, k-1, 0) + c(k, v, k-1, 1), \dots, c(u, k, k-1, p) + c(k, v, k-1, p))$$

The rest of the algorithm follows from the Floyd Warshall algorithm described in lecture.

The final solution will be the min of all $c(u, v, k, h)$. Since the recursion now considers the all of the number of hops, the time to solve each subproblem increases by h . Since the number of sub problems also increases this algorithm will run in $O(V^4h)$

We can reduce the number of subproblems to V^2h by implementing an algorithm very similar to the Dynamic program 1 described in Lecture 11 (APSP). Subproblems: $d(u, v, m) =$ weight of shortest path from u to v using $j = m$ edges Recurrence:

$$d(u, v, m) = \min(d(u, x, m - 1) + w(x, v) \text{ for } x \in V)$$

$$d(u, v, 0) \text{ if } u = v; \infty \text{ else}$$

Topological order: for $m = 0, 1, 2, \dots, h$: for u and v in V : Final Solution: $d(u, v, h)$.

Number of subproblems = V^2h . Time per subproblem $O(V)$. Overall running time: $O(V^3h)$.

d) The following modification to matrix multiplication solves the bounded-hop APSP problem in $O(V^3\log h)$. In regular matrix multiplication we multiply elements and then add them together. We keep most of the process the same but instead of multiplying we add the entries and instead of adding we take the min of the entries. With this modified version of "multiplying" matrices we can make use of successive squaring to solve the bounded-hop ASPS problem. Using successive squaring, we can calculate $\log h$ powers of 2 for the weight matrix and multiply the corresponding matrices to get to h . For example if h is 14 we would multiply $W^2 * W^4 * W^8$. The resulting matrix will be the distance matrix D giving the shortest distances for at most h hop paths because of the way we are "multiplying" matrices.

The time required to "multiply" the matrices remains the same $O(V^3)$ since the changes do not change the asymptotic running time of the original algorithm. Since we use successive squaring of the weight matrix, we will only need to square the matrix at most $\log h$ times. Therefore, this algorithm runs in $O(V^3\log h)$.

e) Case 1: $r = w_{i,j}$ Replacing an edge with another one of equal weight will not make any noticeable changes to D or Π . Therefore, no action is necessary. In this case the algorithm runs in $O(1)$ time.

Case 2: $r < w_{i,j}$ In this case we will have to consider and keep track of the least weight path of distance $< H$ and the least weight path of distance $< H - 1$ all the way to 1 for a total of H matrices. The rest is analogous to case 2 a). This operation takes $O(h)$ per matrix, since

there are $O(h)$ of these matrices and we need to check $O(V^2)$ pairs of points, this algorithm will run in $O(h^2V^2)$.

Case 3: $r > w_{i,j}$ Case three is analogous to part a) where we will have to rerun the algorithm to update D and Π . We can rerun the matrix multiplication algorithm explained in d) to update these in $O(V^3 \log h)$ time.

Problem 6-2

a) Theorem: The undirected graph $G = (V, E)$ with a weight function providing nonnegative real valued weights, where all weights are different, has a unique MST.

Proof: (by contradiction). Consider G' , an MST of G . Now, for the purposes of contradiction, assume that there is another MST for G , call it G'_2 . Since an MST of G with no non-negative edges will have $v - 1$ edges G' and G'_2 must have the same number of edges. Since all the edges are different there is only one set of $v - 1$ edges in an spanning tree that will minimize the weight.

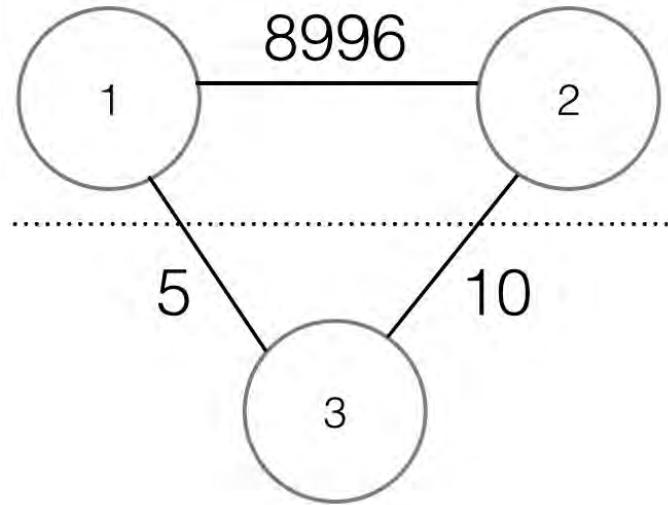
This can also be proved by the correctness of the algorithm described in b). Since b is a completely deterministic algorithm that makes no arbitrary decisions like picking the first element in a tie, it will always pick the unique MST of G.

b) The algorithm is correct. Theorem: For each node n the edge with the least weight coming out of n must be in the MST of G . Proof: (by contradiction). Assume, for the purposes of contradiction, that G' is an MST of G and $n \in V$ does not contain its original adjacent edge with the least weight. Now, we could replace any edge from n with its edge with least weight and still keep G' connected. With this change we have created an tree with a lower edge value sum than the original, contradicting the original claim that G' is an MST of G . We can generalize this theorem to connected components instead of nodes

For each iteration of the algorithm, we add at least a number of edges equal to half of the number of connected components. This happens on the worst case where we only connect each connected component to another. Since the number of connected components halves on each iteration we will need to do this $O(\log V)$ times to get an MST. On each iteration, we need to check $O(E)$ edges. Therefore, the overall running time is $O(E \log(V))$.

c) The algorithm is incorrect. Counterexample: In the graph below, the MST has a total weight of $10 + 5 = 15$. With the algorithm described above, assume that the graph is partitioned arbitrarily into two sections as denoted by the dashed line. The algorithm will find the MST of both sections (for the top section this is simply vertices 1 and 2 directly connected to each other and just vertex 3 for the bottom part) and join the two with the least weight edge between them (in this case the edge between vertex 1 and vertex 3 worth 5). The total edge weight for the "MST" returned by this algorithm is over 9000. By the

initial analysis the tree returned by this algorithm is not the MST of the graph.



d) The algorithm is correct. Begin by performing DFS on the graph. Each time you visit a node n you add it to a list, call this L and a hash set, call this H . If n is in H it means we have seen this node before. We can backtrack using L until we find n again. All of the nodes we pass in this operation are part of a cycle. We take this cycle and remove the edge with the heaviest weight by looking at each edge in the cycle. We repeat this process of finding cycles and removing the most costly edge until no cycles are found. This happens when we transverse all V edges without finding one that we have already visited.

Finding the cycle takes $O(V)$ times since when visiting each node, it can either be an already visited node, or a new node. Since there are only V nodes, this can only happen $O(V)$ time. Similarly, looping through the edges to find the one with the heaviest weight will take $O(V)$ time since we can have at most V nodes in a cycle to loop though. We will do this no more than $O(E - (v - 1))$ times since we can only remove enough edges to be left with $v-1$ for the solution to be a tree. Therefore, the overall running time of this algorithm is $O(VE - V^2 + V)$

This algorithm works because we know that the largest edge weight in cycle cannot be in the final MST, by as many high weighted edges as we can, we guarantee to end up with an MST since the only edges that we can remove are those that make up a cycle.

6.046 Problem 7-1, 7-2, 7-3Collaborators: *Erika Lu***Problem 7-1**

a) 1) By the max-flow min-cut theorem (from CLRS), we know that flow is maximum if the residual network contains no more augmenting paths. By increasing the capacity of an edge by 1, we are allowing for (potentially) one more augmenting path in the residual graph. We know that there may be at most 1 augmenting path (after the update) because of the following: 1) The original residual network had no augmenting paths (by max-flow min-cut theorem). 2) If an augmenting path is found, it must go through the updated edge (since there were no augmenting paths before). After the path is "used" it will be updated according to the Ford Fulkerson algorithm (effectively "removing" that augmenting path). The correctness of this operation follows from the correctness of Ford Fulkerson. Since the capacity can only increase by one, no more augmenting paths will exist. Since there can only be one augmenting path, the max flow can only increase by one.

2) Increasing the capacity of an edge by k is essentially the same as increasing the capacity of an edge by 1 k times. Therefore, this proof follows from the proof in 1). If we increase an edge by one k times, where each time we can increase the max flow by at most 1, then in total, we can increase the max flow by at most k . In other words, after updating the edge on the residual network, there will be at most k more augmenting paths.

3) In the worst case, decreasing the capacity of an edge by 1 will decrease the max flow by 1. This is because, after the update, there will be at most 1 less augmenting path in the residual network. This is a path that made use of (u, v) before, but using it now would exceed the capacity of the edge. Assuming the edge was already at its full capacity, reducing its capacity by 1 would force one of those augmenting paths to not exist. The max flow cannot decrease by more than 1 since the capacity of the edge only changed by one, and can therefore only affect 1 augmenting path. Similarly, the max flow cannot go up since there is no way to create new augmenting paths by decreasing the capacity of an edge. Therefore, the max flow can decrease by at most 1.

4) Decreasing the capacity of an edge by k is essentially the same as decreasing the capacity of an edge by 1 k times. Therefore, this proof follows from the proof in 3). If we decrease an edge by one k times, where each time we can decrease the max flow by at most 1, then in total, we can decrease the max flow by at most k .

b) In order to update the max flow without rerunning all of Ford Fulkerson, we will need to use the current flow and residual network and correctly update the residual edge affected by

the change. From here, we can simply run another iteration of Ford Fulkerson by checking to see if an augmenting path exists in the residual graph. If a path exists we augment the flow along that path and continue looking for paths.

Correctness: If we increase the capacity of one of the current edges by k , then the max flow can only increase a max of k times. Once we update the graph, this is essentially an "unfinished" Ford Fulkerson algorithm with at most k iterations to go. Therefore, the correctness follows from the correctness of Ford Fulkerson as described in detail in CLRS.

Runtime: For this algorithm, on the worst case, we will increase the flow by k (as proven in the first part of this problem). This means we will have to find at most k augmenting paths on the residual graph. Assuming we run BFS to do this, this will take $O(E' + V)$ time per iteration where E' is the number of edges in the residual graph. Notice that the number of vertices remains the same. Since each edge in the original graph can have a reversed edge $|E'| \leq |E|$. Therefore the total time it takes to find a single augmenting path is $O(2E + V)$. Since we might have to find up to k path, the total running time of this algorithm is $O(k(2E + V))$. Notice that the initial step of updating the residual edge accordingly takes constant time.

c) In order to update the max flow given that a single edge (u, v) has reduced its capacity we will need to do the following: We know that the max flow can only decrease (or stay the same) as proven in part a). First, we want to check if the reduction of the edge will even affect the max flow. If the current flow through that edge is still less or equal than the new (reduced) capacity, then the max flow will remain unchanged. If the capacity of the updated edge is less than the current flow, we would be violating max-capacity property and need to correct is as follows:

First, reduce the flow on (u, v) by k . Since the capacity of that edge also reduced by k and we started with a valid flow network, we will now have a valid flow network. At this point, with the capacity updated, and the flow through that edge updated, we will have to add back as much flow as we can. We can do this by running the same algorithm as described in b).

Correctness: First, if $f(u, v) <= r$ then no additional computation is needed since we still have a valid flow network, which we were told was the max flow of the graph. By our results from a) reducing the capacity of an edge by k can only decrease the total flow by k . With this in mind, we can (for a moment) assume that this will occur and subtract k from the flow to that edge. So, $f(u, v) = f(u, v) - k$. At this point, we can add up to $k - 1$ flow back into an edge (implying that the capacity change only reduced the max flow by 1). This is analogous to the algorithm described in b), where the flow can only go up. We find augmenting paths and update the flow accordingly. The correctness of this procedure follows from the correctness of b).

Runtime Updating the edge capacity and checking if the flow is affected takes constant time. Reducing the edge capacity also takes constant $O(1)$ time. In the worst case, we will remove k flow and add $k - 1$ flow back. Notice that we cannot add more than $k - 1$ flow because we did that initial check at the beginning (and the total flow can only go down (from the original) since we reduced the capacity of an edge). This means that we might have to run $k - 1$ iterations of the algorithm described in b), yielding the same asymptotic running time of $O(k(2E + V))$.

Problem 7-2 We can model this problem as a max flow problem by setting the capacity of all edges to 1. $c(u, v) = 1$ for all u, v . Since an edge represents a road, a capacity of 1 represents that that road can only be taken once. Additionally, we connect all of the companies sources to a super sink S with an edge of capacity of 1, this ensures that we only find (at most) one path for each company.

With the graph built, we can run the Ford Fulkerson method to find the max flow of the graph. If the max flow is less than k $|f| < k$ there is no way to route k trucks so we return impossible. Otherwise we can returned the k paths as explained below: At this point we know that we can send k trucks through the graph given the constraints. Therefore, we can go through each of the source vertices and arbitrarily follow edges until we get to the sink. We return the path taken as the path for that source, remove the edges from the graph and continue to the next source.

On each iteration of Ford Fulkerson, when we find an augmenting path in the residual network, we will have increased the flux by exactly one. This is because all of the edge weights have a capacity of at most 1.

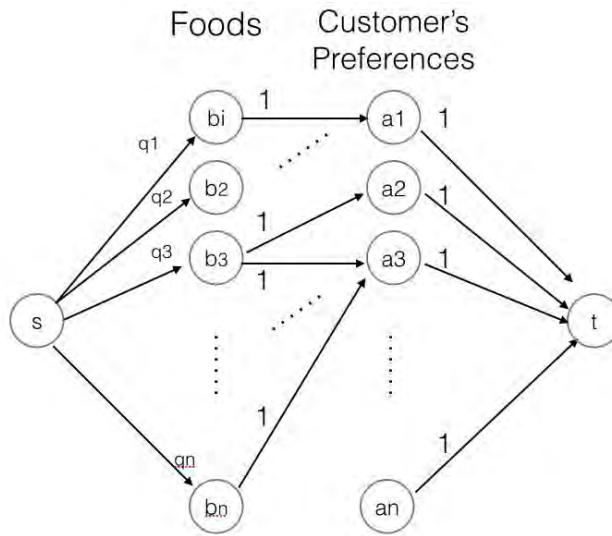
Correctness For this problem, we modified the given graph with capacities and a supersink and then model is as a max flow problem solved by Ford Fulkerson. The "flow" in this problem are the trucks, which take edges from their company source to the common destination/ the sink. An edge in the graph represents a road and by setting the capacity of this edge to 1, we represent a road that can only be taken once. Furthermore, we connect all company sources to a supersource by an edge with a capacity of 1, ensuring that each company can only send at most 1 truck to the sink. The correctness of the max-flow problem follows from the correctness of Ford Fulkerson as explained in CLRS. After running Ford-Fulkerson, we will know whether we can "send" k trucks through the graph or not. If the max flow is less than k , then there is no way to send all k trucks through the graph. Otherwise since it edge can only be used once, we can pick any path for each company, remove the used edges, and still have available edges for the rest of the company.

Runtime Building the graph will take $O(k + E)$ time since we have to set the capacity of all edges to 1 ($O(E)$) and link all of the k companies to the super source ($O(k)$). Once the graph is built we need to run the Ford Fulkerson method. This will take $O(E|f|)$ time. The max possible flow on any graph constructed this way is k because each of the k sources can

only be used once (since they are connected by an edge of capacity 1 to the super source). Therefore, running Ford Fulkerson will take $O(Ek)$. Note that finding all of the paths will take $O(E)$ since we will go through each edge exactly once (because we remove an edge after using it). Overall, the algorithm will run in $O(k + E) + O(Ek) = O(Ek)$

Problem 7-3 This problem can be modeled as a max flow problem if we construct a flow graph G as follows:

We begin with a source s that connects to all of the vertices on the first column. This column represents the different available foods b_1, b_2, \dots, b_n . Also, $q_i c(s, b_i) = q_i$ this means that an edge from s to b_i will have a capacity equal to the quantity of that food. The third column of the graph represents customers and an edge from b_i to a_i implies that customer a_i can eat food b_i . All of the elements on the third column connect to the drain node t . We assign each one of these edges a capacity of 1 to make sure that each customer only gets 1 food item assigned to him/her. Below is a graphical representation of this description:



After building the graph, we can use the Ford Fulkerson method to find the max flow. From here, we can scan through the third column and look at the flow of incoming and outgoing edges to a_i . If customer a_i has an incoming flow of 1, this means that he/she got a food item. To find this item, we have to see where the flow is coming from and follow the edge to appropriate vertex on the second column. If customer a_i has an incoming flow of 0 he/she did not get a food item and should get a voucher. We continue to do this for all customers. At the end we will know what food each customer got, or if they got a voucher. Since this is a max flow problem, we tried to maximize the number of food items distributed to customers which minimized the number of vouchers to be given out.

Runtime: This algorithm is essentially two parts. The first part consists on building the graph and the second consists on running the Ford Fulkerson method on that graph. Since the overall running time will be driven by the latter, I will analyze that first.

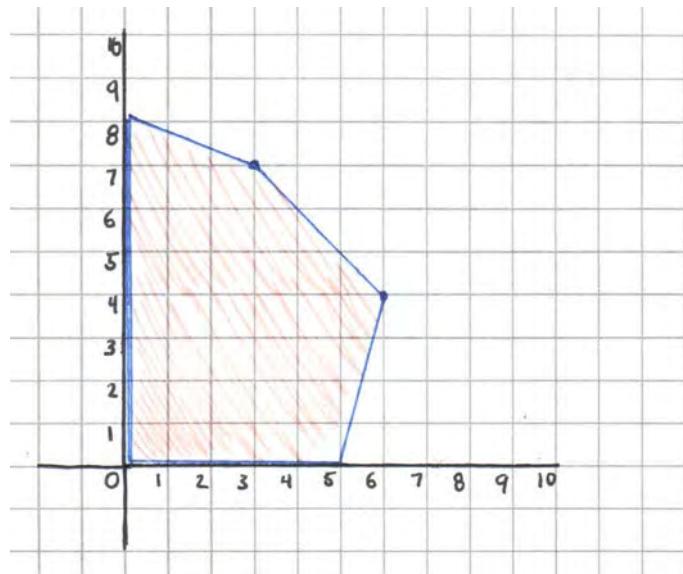
Ford Fulkerson runs in $O(E|f|)$. In the worst case we will have every customer getting something so the max flow $|f| = O(n)$. In the worst case, we will have a very dense graph where every customer sets every food item in their preferences, this forces the number of edges $E = O(nm)$. Combining these two, we get that the Ford Fulkerson method runs in $O(mn^2)$.

Building the graph is upperbounded by mn , which occurs when every customer sets every food item in their preferences. Combining both parts of the algorithm we get a runtime of $O(mn^2) + O(mn) = O(mn^2)$

Correctness: Our model of the problem as a max flow problem is correct because of the way the graph was constructed. In this problem we can think of the flow as the food, which travels through a graph of customers till the destination (sink). We connect all foods to the source, since initially all foods are available with a limited quantity equal to the capacity of the edge that connects to it from s . This ensures that we don't try to give out more food than we have. These foods connect to customers that stated that they can eat that food with an edge of capacity one. This ensures that each customer only receives one of each item. To ensure that each customer gets at most 1 item, we set the edges from the customer vertices to the drain to have a capacity of 1. The correctness of the max flow problem follows from the correctness of Ford Fulkerson. Once we are done with Ford Fulkerson, for each customer vertex, there are two cases: Either the customer received an item which means there is an incoming edge with flow of 1 into that vertex (and subsequently an edge with flow one leaving the vertex by the flow conservation property) or the customer did not receive any food and should get a voucher. Again, since this is a max flow problem, we tried to maximize the number of food that was distributed to customers which minimized the number of vouchers to be given out.

6.046 Problem 8-1, 8-2Collaborators: *Erika Lu***Problem 8-1**

- a) In the following diagram of the feasible region, the x axis corresponds to x_1 while the y axis corresponds to x_2



b)

Standard form:

$$x_1 + x_2 \leq 10$$

$$4x_1 - x_2 \leq 20$$

$$x_1 + 3x_2 \leq 24$$

$$x_1, x_2 \geq 0$$

Slack Form:

$$10 - x_1 - x_2 = x_3$$

$$20 - 4x_1 + x_2 = x_4$$

$$24 - x_1 - 3x_2 = x_5$$

$$Z = 4x_1 + x_2$$

c) First Pivot: $x_1 = 5$ and is then swapped with x_4

$$x_3 = 5 - \frac{5x_2}{4} + \frac{x_4}{4}$$

$$Z = 20 + 2x_2 - x_4$$

$$x_1 = \frac{20 + x_2 - x_4}{4}$$

$$x_5 = 19 - 13x_2/4 + x_4/4$$

2nd Pivot: $x_2 = 4$ and is then swapped with x_3

$$x_2 = 1/5(20 + x_4 - 4x_3)$$

$$Z = 28 - 3x_4/5 - 4x_3/5$$

$$x_1 = 6 - x_4/5 - x_3/5$$

$$x_5 = 6 - 2x_4/5 + 13x_3/5$$

Since the coefficients are all negative at this point in the modified LP, we can conclude that the max solution will be 29. We can then plug this in for the original objective function and solve for x_1 and x_2 given that x_3 and x_4 are zero. Therefore $x_1 = 6$ and $x_2 = 4$.

In our diagram from a), on the first iteration, we got our objective function to intersect at point (5,0). On the second iteration we got it at (6,4). At this point, we found a maximum and we can calculate our solution.

d) The dual LP of the standard form LP from part b) can be proved by obtaining a certificate of optimality. By the LP Duality theorem, we show that this solution is optimal.

$$\vec{b} \cdot \vec{y} = [10y_1 \quad 20y_2 \quad 24y_3]$$

$$A^T \vec{y} \geq \vec{c}$$

$$\begin{bmatrix} 1 & 4 & 1 \\ 1 & -1 & 3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \geq \begin{bmatrix} 4 \\ 1 \end{bmatrix}$$

Problem 8-2

- a) Step 1: Show TRIPLE-SAT is NP

Given a solution to Φ , I can plug in the variables into the original equation and verify that they simplify to True. This can be done in polynomial time, for each of the 3 (different) proposed solutions.

Step 2: Reduce from known NP-Complete and show both ways

If we have a black box that can solve any SAT problem, we can modify our original TRIPLE-SAT problem in poly time, and use SAT to solve the problem. We modify the original problem p as follows. We take our original problem p and OR it with the result of XORing two new variables A and B . The modified problem looks like this:

$$p \parallel (A \oplus B)$$

By constructing our problem this way, we guarantee that three solutions exists, (when A is true and B is false, A and B are false, A is false and B is true). We can then use SAT on this modified problem to generate three solutions to the original problem. Therefore,

To prove from TRIPLE-SAT to SAT, we just need to pick the solution where A and B are both false from TRIPLE-SAT and use that in SAT. We can do this in poly time. Since TRIPLE-SAT will return true if SAT returns true and viceversa, and TRIPLE-SAT is in NP, we can conclude that TRIPLE-SAT is NP Complete.

- b) For this problem, we will reduce DONUT from 3SAT in order to prove that DONUT is NP Hard. To begin the reduction we will need to devise a way to map a 3SAT problem into its corresponding map for a DONUT problem. We do this as follows. For each clause, we create individual nodes and link them to each other. We also link each one of these nodes to a new node (call these variables) corresponding to their complement. Finally, we connect each complement to a new node corresponding to the original value (or the complement of the complement). By connecting node A to its complement, we assure that they both can't be true (since you can't select two neighbors). Saying there is a donut shop at node A equates to saying A is true on the 3SAT problem.

We can build this graph from a 3SAT problem in polynomial time, we can then calculate k: $k = \text{num(clauses)} + \text{num(variables)}$.

Since DONUT can be reduced from 3SAT, and 3SAT is NP complete. A polynomial time

algorithm for DONUT would automatically imply that 3SAT is solvable via a polynomial time algorithm (with a reduction). This would then show that $P = NP$.

c) We can turn this problem into a decision problem by asking whether it is possible to produce a schedule that yields a profit of k . With this decision problem, we could solve the original problem by simply binary searching for the max profit. Note that the min profit is 0 and the max profit is the sum of all p_i : $k = \sum_{i=0}^n p_i$.

Step 1: Show that SCHEDULE is in NP Given a schedule, we could simply loop through the events, check that they are valid (not conflicting with each other) and then add up all the profits. From there we could determine if the the profits are greater than or equal to k .

Step 2: Reduce from known NP-Complete and show both ways

For this problem, we will reduce PARTITION to a specific case of scheduling to show that any problem in NP can be solved if scheduling can be solved. To understand the logic behind the reduction, I will first explain some properties of PARTITION. In PARTITION, we want to take a set of integers as an input, and output two sets with equal sums. We can easily prove that each of these sums must be equal to $S/2$ where S is the sum of all elements in the original set. If we are "picking" different numbers to add to a set, we need to stop adding numbers once our sum reaches $S/2$.

With this in mind, we begin our reduction by setting the profit of a task equal to the corresponding integer value on the input to partition. So an input of (1,2,4) to partition would imply that there is a task with 1 profit, another with 2 and another with 4. Now we want k to be equal to $S/2$, since we are interested in knowing if we can get at least k profit with the scheduled tasks. We also say that deadline = $S/2$ since we want to stop picking tasks once we reach this point (similarly to how PARTITION stops picking integers once the first subset reaches $S/2$). Lastly, we set the duration of each task equal to the profit.

With this, we are able to reduce PARTITION to a specific case of scheduling to show that scheduling is NP Hard.

Since SCHEDULE is in NP and is NP hard, SCHEDULING is NP Complete.

6.046 Problem 9-1, 9-2Collaborators: *Erika Lu, Ben Gunby, Patrick Liscio, 18.434 Notes***Problem 9-1**

- a) To prove that *Alg1* does not guarantee any constant approximation ratio, we first need to show that picking the most dense object may not be the best decision. This could be the case when picking the most dense object, prevents you from picking two less dense objects with a larger total value.

More generally, we want the density of a higher valued object to be less than the density of a lower valued object and have the sizes of both objects big enough so we can only pick one of them. The algorithm will pick the denser lower valued object instead of the optimal less dense one. Consider two objects a_1 and a_2 . To meet the scenario described above, we want the following to hold:

$$\begin{aligned}\frac{v_1}{v_2} &= \frac{1}{k} \\ \frac{1}{k} &> \frac{s_1}{s_2}\end{aligned}$$

To satisfy these, we can set the values to the following

$$\begin{aligned}v_1 &= \frac{B}{k+1} = s_1 \\ v_2 &= \frac{Bk}{k+1} \\ s_2 &= B\end{aligned}$$

- b) The first part of this algorithm is trivial. If we can fit all of the items, we take them all. This is the optimal solution. Otherwise, we sort the items by density and find the smallest index i that makes the total size of the first i items exceed the max capacity B . At this point, we consider taking either the first $i - 1$ items, or item i .

We can establish an upperbound on the optimal solution, call this V^* . By taking the first i items would be taking i of the most dense items, and would also exceed our capacity, this is clearly an upper bound. In math terms:

$$\sum_{j=1}^i v_j \geq V^*$$

Since the algorithm picks the max value of taking either the first $i - 1$ items, or item i , and together these two (added) form an upperbound to the optimal solution, each of these must be at worst half of the optimal value. Therefore, Alg_2 always yields a 2-approximation to the optimal solution.

- c) The following Dynamic Programming solution solves the knapsack problem. Define the subproblem to be $S_{i,v}$ which equals the min subset $a_1 \dots a_i$ whose total value is exactly v . The base cases are:

$$S_{0,v} = \inf$$

$$S_{i,0} = []$$

For the first case, no solution exists, so we return infinity. For the second case, the solution is empty since no objects have 0 or negative value. The recurrence is as follows:

$$S_{i,v} = \min(S_{i-1,v-v_i} \cap [a_i], [S_{i-1,v}])$$

The left side (inside the min), represents choosing item i and the left side represent not choosing item i . We will need to calculate all $S_{i,v}$ and then chose the optimal form there. For the runtime analysis, note that we can obtain an upperbound on the optimal value by looking at V , the largest value of the all items in the set. An upperbound is then nV where n is the size of the set of the items in consideration. With this in mind, the total runtime will be $O(n^2V)$ since we will have to go through at most n different values for each $S_{i,v}$ that we calculate.

- d) Notice that Alg_3 is not fully polynomial because the values of the items may not be polynomially bounded in n . Alg_4 achieves FPTAS by scaling the values of each item to be polynomially bounded by n . It then runs Alg_3 and returns the solution S' .

To show that Alg_4 is FPTAS, we will begin by showing $V(S') \geq (1 - \epsilon)OPT$ for the returned solution S' where OPT is the optimal solution. First, we must note that any object i will have $Kv'(i) \leq v(i)$ because we scaled down the value of each item by $K = \frac{\epsilon v_1}{n}$. This means that the max value for an optimal solution S^* can decrease by at most K for each item (since we scaled down the value of each by K). Since we have n items, the total value can decrease by at most nK . We then run Alg_3 with the scaled values and are guaranteed to get an optimal set for the scaled values and we know that this must be at least as good as picking S^* with the original values. At this point, we know that the value of the returned set S' is at least as good as $KV'(S^*)$. As shown above this is the same as $V(S^*) - nK = OPT - \epsilon V = (1 - \epsilon) * OPT$. The running time of the algorithm follows from the running time of Alg_3 but V is now polynomial in n . More specifically, this algorithm runs in $O(n^2V) = O(n^2\text{floor}(\frac{n}{\epsilon}))$ Since Alg_4 is also within a factor of $(1 - \epsilon)$ from the OPT Alg_4 is FPTAS.

Problem 9-2

a) A minimal cycle cover A of T will include exactly one edge u, v from each cycle on T . This is because a cycle cover includes at least one edge from each cycle and a minimal would include at most one. Knowing that we only have one edge from each cycle in A , reversing all edges in A is equivalent to reversing one edge in every cycle of T , effectively breaking every cycle.

To finish the proof, we need to consider the case where reversing an edge would introduce a new cycle. This case however, cannot happen. The proof is by contradiction. Assume that there is an edge that when reversed, created a new cycle in the tournament graph T . Then, that edge must have been the only edge in A on some directed cycle T , this means there was another cycle not accounted in the cycle cover which is a contradiction.

b) This problem makes use of some important observations about Triangles and cycle covers. First, two triangles cannot share two edges, they can only share one. This can be proved by contradiction: If two distinct triangles share two edges, they share three vertices and are thus not distinct triangles. Next, we note that if any edge is shared among k or more triangles, it must be in the cycle cover. The last observation is that given an edge in a cycle we can construct a triangle containing that edge.

With these observations in mind, we can begin. The first step is to reverse all edges shared among more than k triangles. These are edges that would be in the set (per observation 2). For each edge that we find, we decrement k by one and continue until we can't find another edge or $k = 0$. On the latter case, we just check if the remaining graph is acyclic and then we are done. Otherwise, we will end up with a graph where every edge is in at most k triangles. It is important to note that at this point, we have not removed anything from the graph, we've only reversed some edges.

The next step can be thought of as a reversal of part a). In part a) we showed that if we have a k cover and reverse those k edges we end up with an acyclic graph. Now, we will be considering an acyclic graph and try to find which edges to reverse to get to it. Recall that each edge has is shared with at most k triangles. We now consider how many vertieces there are in these k triangles. For each triangle, we may have 2 vertices + k from the other k triangles. Since we have k of these we end up with $k(k + 2) = k^2 + 2k$.

c) The first step in this algorithm is to perform the reduction/kernel described in part b). This step will take $O(n^3)$ since we will need to count the number of triangles and all other operations are $O(n^3)$ At this point, we can perform any brute force algorithm (check every combination) to solve the problem, and return the appropriate decision. the runtime of this would be $O(k^{f(k)})$ which is still a FTP algorithm since $k^{f(k)} = 2^{f(k)\log(k)}$

Dynamic Programming

Knapsack

$T(n) = 2T(n-1) + \Theta(\log n)$
 1) Define subproblems
 2) Compute subproblem solutions
 3) Reuse and memoize
 4) Solve original problem
 Analyze runtime

Example: Increasing Subsequence
 $L(j) = L(j-1)$ ending at position j
 Goal: Find longest strictly increasing
 (not necessarily contiguous)
 Solution: $\max\{L(j)\}$

Runtime: $O(n^2)$
 Edit Distance
 Input: Strings $A[1..n], B[1..m]$
 Costs: $C[i,j]$, $C[i,j] = \min\{C[i-1,j-1], C[i-1,j], C[i,j-1]\}$

$L(j) = \max\{L(j-1), 1\} + 1$
 Goal: Build array (row-crossing)
 indices are possible
 like longest increasing subsequence
 i.e., index of corresponding city on
 northern bank in O(n log n) + i.e.
 find longest subsequence
 of $k(1), \dots, k(n)$ $\Theta(n^2)$

Integer Knapsack Problem
 Input: n types of items each with
 size s_i , $i=1..k$ and value v_i , knapsack C
 Goal: Pack items in knapsack maximize
 (repeated items ok)
 $M(j) = \max$ val pack in knapsack $\leq j$
 $T(i,j) = \max\{T(i-1,j) + v_i, T(i-1,j-s_i)\}$
 Note: s_i is not explicitly in N_p .
 If complete, $\leq N_p$ will do.

Union-Find
 MAKE-SET(x): x is one node out
 FIND-SET(x): Return rep of $x out$

Master Theorem: Solve recurrence $T(n) = a(n/b)T(n/b) + f(n)$
 ① calculate $\log_b n$ and compare w/ $f(n)$
 Case 1: $f(n) \geq \Theta(n^{\log_b a}) \Rightarrow T(n) = \Theta(n^{\log_b a})$
 Case 2: $f(n) = \Theta(n^{\log_b a}) \Rightarrow T(n) = \Theta(n \log n)$
 Case 3: $f(n) \leq \Theta(n^{\log_b a - \epsilon}) \Rightarrow T(n) = \Theta(f(n))$

Building Bridges
 Input: $C[1..n], l[1..n]$
 Goal: Build as many (row-crossing)
 bridges as possible

Runtime: $O(n^2)$
 Edit Distance

$T(i,j) = \min$ cost to transform $A \rightarrow B$
 $A[i-1..j] \rightarrow B[i..j]$
 $T(i,j) = \min\{T(j,i-1) + AC[i] \neq B[i]\}, T(i-1,j)$
 $M(j) = \max$ val pack in knapsack $\leq j$

For $j=1$ to n # bottom up approach
 $M(0) = 0$

For $j=1$ to n # bottom up approach
 $M(j) = \max\{M(j-1), M(j-1) + v_j\}$

Solution: $T(n,m)$
 Input: n integers $[a_1, \dots, a_n]$
 goal: partition integers into S_1 and S_2
 Requests compatible if no overlap
 $P(i,j) = \begin{cases} 1 & \text{if same subset of } S_A \text{ and } S_B \\ 0 & \text{otherwise} \end{cases}$

Balanced Partition ("fairly dividing sets")
 Input: n integers $[a_1, \dots, a_n]$
 goal: partition integers into S_1 and S_2
 Requests compatible if no overlap
 $P(i,j) = \begin{cases} 1 & \text{if same subset of } S_A \text{ and } S_B \\ 0 & \text{no sum } j \end{cases}$

$P(i,j) = 1 \text{ if } P(i-1,j) = 1 \text{ or } P(i-1,j-A) = 1$
 $= \max\{P(i-1), P(i-1,j-A)\}$
 $= \min\{S_{i-1}, P(i-1)\}$

Weighted Interval Scheduling
 Resources and requests $1..n$
 Start time s_i End time e_i
 End time e_i

Request compatible if no overlap
 Goal: select compatible subset of n requests
 of max weight. w_i weight of

subproblems:
 $R_i = \{ \text{request } j \in R \mid s_j \geq x_i \}$
 $\text{opt}(R_i) = \max\{ \sum_{j \in R_i} w_j + \text{opt}(R_{i'}), 0 \}$
 Runtime: $O(n^2)$

Horners Rule
 $A(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$
 Evaluation: $\frac{a_0}{x} + \frac{a_1}{x} + \dots + a_n$ Now: coeffs $\frac{a_0}{x^n} + \dots + \frac{a_n}{x}$
 (can be original) if present first

Change of Variables:
 $T(n) = 2T(\frac{n}{2}) + \Theta(\log n)$
 Let $n = 2^m$ and substitute
 $T(2^m) = 2T(\sqrt{2^m}) + \Theta(\log 2^m)$
 $= 2T(\sqrt{2^m}) + \Theta(m)$
 Divide $\sqrt{2^m} = (\sqrt{2})^m = T(n)$
 $\sqrt{2^m} = T(2^m) = 2 + (2^m) + \Theta(m)$
 $= 2 + (m^2) + \Theta(m)$

Solve with Nested Recurrences:
 $T(n) = 2T(\sqrt{n})$
 $4T(\sqrt{n})$
 $4T(\sqrt{n}) = \Theta(n)$

Linked List Solution:
 repr Initial: x as one node out
 MAKE-SET(x): one node out
 FIND-SET(x): return $repr$ of $x out$

$UNION(X, Y)$: Join X and Y

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

Overall can get almost constant amortized time / operation

Random Algorithms

Monte Carlo - Always polynomial expected correct with high probability

Las Vegas - Expected polynomial, always correct

Quicksort:
 $\Theta(n \log n)$ pivot element x Partition

Quicksort with pivots L, R combine $T(L) + T(R)$

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

Amortized Analysis

Aggregate Method: $\frac{\text{total cost}}{k}$ operations per operation

Accounting Method: operations can store credit

can pay for time using stored credit. $\Theta(n \log n)$

Example: $2-3$ trees insert \Rightarrow add 1 in worth $\log n$. $\&$ cost that can always pay since we can't delete something wasn't inserted

Method: operations can charge to the past amortized cost = actual cost - charge to past + charge by future

Example: Table doubling charges $\Theta(n)$ to $\Theta(n^2)$ inserts since last doubling

\hookrightarrow Inserts only charged once!

Potential Method: Define potential function Φ mapping structure configuration \rightarrow non-negative integer.
 Φ - Potential high costs in the Future $\Phi(0) = 0$

Example: Binary counter \Rightarrow $\Phi(1) = 1$ after much work
 $\Phi(0) = 0$ after work $\Rightarrow \Phi(1) - \Phi(0) = 1$

Charging Method:
 amortized cost = actual cost - charge to past + charge by future

Example: B-trees: $\Phi = \# \text{ nodes w/ } B \text{ children}$

Amortized cost $= \Phi(B-1) + 1$ for insert
 $= \Phi(B) - \Phi(B-1) + 1$

Search(x): Walk right on top list L_n until going right would be too far.
 Insert(x): Search(x) and insert into bottom list L_0 tails \rightarrow stop

Failure Prob: $\Pr[\text{No level}]$
 $\Pr[\text{some level got promoted}] \leq \Pr[\text{at least one level got promoted}]$
 $\leq n \cdot \Pr[\text{level got promoted}]$

Search(x): $\Pr[x \in L_n]$

Chernoff Bounds:

$\Pr[Y \geq E[Y] + r] \leq e^{-r^2/2E[Y]}$

Friodell's algorithm $\Pr[A \text{ true}] = \frac{1}{2}$

Choose random binary vector A and B

$\Pr[AB = C] = \Pr[AB \neq C] \geq \frac{1}{2}$

IF $AB = C$ $\Pr[A \text{ true}] = 1$

IF $AB \neq C$ $\Pr[A \text{ true}] = 0$

Return YES/NO

Verify Matrix multiplication

$A = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$

$B = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$

$C = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$

$\Pr[AB = C] = 1$

$\Pr[AB \neq C] = 0$

$\Pr[A \text{ true}] = 1$

Runtimes

Merge Sort (worst): $T(n) = 2T(\frac{n}{2}) + \Theta(n) = \Theta(n \lg n)$
 Binary Search (worst): $T(n) = T(\frac{n}{2}) + \Theta(1) = \Theta(\lg n)$
 Randomized Quicksort (Expected)
 $T(n) = \frac{1}{n} \sum_{k=1}^{n-1} T(k) + \Theta(n) = \Theta(n \lg n)$

Randomized Quicksort (worst): $T(n) = T(n-1) + \Theta(n) = \Theta(n^2)$

Strassen's Algorithm (worst): $T(n) = 7T(\frac{n}{2}) + \Theta(n^2) = \Theta(n^{\lg 7})$

Selection (worst): $T(n) = T(\lceil n/2 \rceil) + T(\lceil \frac{n}{2} \rceil) + \Theta(n) = \Theta(n)$

Randomized Selection (worst): $T(n) = T(n-1) + \Theta(n) = \Theta(n^2)$

Randomized Selection (Expected): $T(n) = \frac{1}{n} \sum_{k=1}^n T(k) + \Theta(n) = \Theta(n)$

Median Finding: $\Theta(n)$
 Weighted Interval Scheduling: $\Theta(n \lg n)$
 Polynomial Evaluation (Horner's): $\Theta(n)$
 Polynomial Multiplication (FFT): $\Theta(n \lg n)$
 BST Search Delete: $\Theta(h)$, $h = \Theta(\lg n)$
 Van Emde Boas balanced: $\Theta(\lg \lg U)$
 Union Find make set: $\Theta(\alpha)$ Amortized
 Skip List: $\Theta(\lg n)$ with high probability

SSSP: DFS $\Theta(V+E)$
 Dijkstra $\Theta(V \lg V + E)$
 Bellman-Ford $\Theta(VE)$
 Topo + D for DAG $\Theta(V+E)$
 APSP: Floyd-Warshall $\Theta(V^3)$
 Johnson's $\Theta(V^2 \lg V + VE)$
 Flow: Ford-Fulkerson: $\Theta(VE)$
 Edmond-Karp: $\Theta(VE^2)$
 Leader Election-Bidirectional: $\Theta(\lg n)$
 Ring of n processes: $\Theta(\lg n)$
 Luby(MIS): $\Theta(V \lg \lg V)$

Hash Family H is universal if for each pair of distinct keys $x, y \in U$, the # of hash functions $h \in H$ for which $h(x) = h(y)$ is exactly $1/m$.

It is possible to compute the smallest \tilde{m} elements of an n -element array in sorted order in $\tilde{O}(n)$ time: SELECT \tilde{m} elements and partition around it, then sort those \tilde{m} elements in $\tilde{O}(n)$ time.

* $\tilde{O}(\lg \lg(n^m)) = \tilde{O}(\lg(n \lg(m))) = \tilde{O}(\lg(n) + \lg \lg(n)) = \tilde{O}(\lg(n))$

* If the deterministic median finding algorithm had two recursive subproblems of sizes $\frac{3}{5}n$ and $\frac{31}{40}n$ (instead of $\frac{1}{2}n$ and $\frac{1}{2}n$) then algo still runs in $\tilde{O}(n)$.

Total subproblem size = $\frac{1}{5} + \frac{31}{40}n = \frac{36}{40}n$ on each level

Sum of geometric series \Rightarrow linear

Perfect hashing can use more than one hash function, one for first level, and one for every bucket of the 2nd

* If Φ is a potential function. For S , 2Φ is also a function for S and the amortized running time of each operation $\geq 2\Phi$, i.e. at most twice Φ .

* skip lists can perform all operations in $\tilde{O}(\lg n)$ as long as prob of promoting p < 1

* Amortized cost = actual cost + $\Phi(\text{after}) - \Phi(\text{before})$

NP Hard / Complete Problems

P-set of problems solvable in poly-time
NP decision problems with poly-time verification / decision problems solvable in non-deterministic polynomial time

NP-Hard: Problem D is NP-hard if $A \leq_p D$ & A is NP. Reduces to D

NP Complete: Problem D is NP-complete if D is NP and D is NP-hard

Reduction: Problem A poly-time reduces to B aka $A \leq_p B$ or $A \leq_p B$ if \exists a poly-time alg R s.t. $A(x) = B(R(x))$

Prove X is NP-Complete:

- \exists NP via nondeterministic algo or certificate + verifier
- Reduce from known NP-complete problem Y to X
 - Poly-time conversion from Y inputs to X inputs.
 - If $Y_{\text{YES}} \rightarrow X_{\text{YES}}$
 - If $X_{\text{YES}} \rightarrow Y_{\text{YES}}$

SAT - Given boolean formula of the form: $(x_1 \vee x_2 \vee x_3) \wedge (\bar{x}_2 \vee x_3 \vee x_4)$

Is there a variable x_{true} clause such that formula is true?

NP-hard: convert algo into circuit

NP-Verifiable in poly-time

Sample Reduction: Hamiltonian Cycle \Rightarrow Traveling Salesman Problem

Step 1: Show Ham-PATH is NP-hard
Reduction. Given an instance of HAM-Cycle b, want to create input to HAM PATH.

Step 2: Show Ham-PATH is in NP

a) Does y touch each vertex once?

b) Is every edge in p in G even that runs in poly-time. Graph \Rightarrow Path

Flow Networks

Flow Network: Directed graph with source and sink + Each edge (u, v) has capacity $c(u, v)$

Flow \rightarrow Rate

Capacity Constant: For all $uv \in E$, $0 \leq f_{uv} \leq c(u, v)$

Flow Conservation: For all $u \in V$, $f_{u,u} = \sum_{v \in V} f_{uv}$

Flow into vertex must come out

$|F| = \sum_{u \in V} f_{(s, u)} = f_{(s, V)}$ $f_{(X \cup Y, Z)} = f_{(X, Z)} + f_{(Y, Z)}$

Properties: $f(X, X) = 0$ $f(X, Y) = -f(Y, X)$

$|F| = F(V, V) = F(V, U) + F(U, V)$

Optimal Substructure: optimal solution to problem incorporate optimal solution(s) to subproblem(s)

Greedy Choice Property: locally optimal choices lead to globally optimal solution.

Minimum Spanning Tree (MST)

Spanning tree $T \subseteq E$ of min weight

Greedy choice property: cut-property

- lightest edge that connects S to T will be in an MST.

Kruskal's Algorithm:

start with random vertex \rightarrow visit all edges out of visited that don't go to already visited vertex

pick smallest one + visit new edge \rightarrow augment path

$= \alpha(V, E) + \beta(E)$ using Fibonacci heap?

Kruskal's Algorithm:

Keep adding lowest-weight edges that Pseudo-code using Union-Set

$T = \emptyset$ // Build MST

For $v \in V$ - Make-Set(v)

Sort E by weight

For $e = (u, v) \in E$ // in order

if $\text{Find-set}(u) \neq \text{Find-set}(v)$

add e to T

Correctness: when adding edge, we make use of greedy choice property

Ford-Fulkerson Algorithm:

$f[u, v] \leftarrow 0$ for all $u, v \in V$

From every source vertex $s \rightarrow$ get $\delta_n(u, v)$ for all $u, v \in V$

③ Claim $\delta_n(u, v) = \delta_n(u, u) - h(u) + h(v)$

④ $\delta_n(u, v) \geq 0$ distance to go from u to v in at most m edges

$\delta_n(u, v) = \min(\delta_n(u, u) + \delta_n(u, v)) \Rightarrow \delta_n(u, v) \geq 0$

$\delta_n(u, v) = \min(\delta_n(u, u) + \delta_n(v, v)) \Rightarrow \delta_n(u, v) \geq 0$

Flow Programming

Politics example: How to campaign to win election. Votes obtained per dollar:

Demographic distribution Rural

x_1	7	5	3
x_2	2	5	5
x_3	0	0	2
x_4	10	0	0
majority	8000	100000	25000

Want to win majority in each demographic by spending minimum amount of money et x_1, x_2, x_3, x_4 denote dollars spent/issue.

Minimize: x_1, x_2, x_3, x_4 subject to:

Subject to:

$$\begin{aligned} x_1 + x_2 + x_3 + x_4 &\geq 8000 \\ 5x_1 + 2x_2 + 10x_3 + 0x_4 &\geq 100000 \\ 3x_1 + 5x_2 + 10x_3 + 2x_4 &\geq 25000 \end{aligned}$$

Certificate of Optimality - shows that LP solution is indeed optimal

Primal: $\max \tilde{c}^T \tilde{x}$ min $b^T y$

$\begin{cases} \tilde{c}_1 \leq \tilde{b}_1 \\ \tilde{c}_2 \leq \tilde{b}_2 \\ \vdots \\ \tilde{c}_m \leq \tilde{b}_m \end{cases} = \begin{cases} \tilde{c}_1^T \tilde{x} \geq \tilde{b}_1 \\ \tilde{c}_2^T \tilde{x} \geq \tilde{b}_2 \\ \vdots \\ \tilde{c}_m^T \tilde{x} \geq \tilde{b}_m \end{cases}$

Inverting to Standard Form:

Want to minimize $-2x_1 + 3x_2 \Rightarrow$

Negate coefficients and minimize.

IFF x_j does not have a non-negativity constraint x_j replaced by $x_j'' - x_j$ $x_j'' \geq 0$

Equality constraint: $x_1 + x_2 = 7$ translates to $x_1'' - x_1 + x_2'' - x_2 = 7$

\Rightarrow constraint translates to \leq by multiplication of -1 .

Worst case running time for simplex is exponential, but it is usually efficient.

IFF constraint graph of linear system is acyclic, solution exists and can be found in $O(|V| \cdot |E|)$ time using DAG shortest paths.

Approximation Algorithms

4 Reducing Clique to Independent Set

Clique: Given graph $G = (V, E)$ and integer k , is there a set of vertices $C \subseteq V$ with $|C| = k$ that form a complete graph?

Independent Set: Given graph $G = (V, E)$ and integer k , is there a set of vertices $I \subseteq V$ with $|I| = k$ such that for any $u, v \in I$, $(u, v) \notin E$? Given that Clique is NP-Complete, we prove that Independent Set is NP-complete.

1. Show Independent Set \in NP

To prove this, we need to prove there exists a verifier $V(x, y)$. Let $x = (G, k)$ be a "yes" input. Let y be I that satisfies the condition.

It takes $O(|I|^2)$ to check that $|I| = k$. It takes $O(|I|^2)$ to check that for every $u, v \in V$, $(u, v) \notin E$. This checks in polynomial time that the certificate y proves that x is a valid input. Therefore, Independent Set is in NP.

2. Show Independent Set \in NP-hard

We prove this by giving a Karp-reduction of Clique to Independent Set.

(a) Given an input $x = (G, k)$ to Clique, create input G' which has the same vertices, but has edge (u, v) if and only if $(u, v) \notin E$. This takes $O(|E|)$ time, so this reduction takes polynomial time.

(b) If I is a set of vertices that form a k -Independent Set for G' , then $C = I$ is a k -Clique for G because for $u, v \in I$, Independent Set says that $(u, v) \notin E'$, but this implies that $(u, v) \in E$ for the Clique problem due to the method of construction. This shows that there are edges between every pair of nodes in C . In addition $|C| = k$, and so C is a k -Clique.

(c) If C is a set of vertices that form a k -Clique in G , then $I = C$ is a k -Independent set for G' . This is because $u, v \in C$ implies that $(u, v) \in E$ for Clique, and this implies that $(u, v) \notin E'$ for Independent Set. Since $|I| = |C| = k$, this shows that there are k elements in the construction G' that are not adjacent to each other.

3. This proves Clique reduces to Independent Set in polynomial time, which means that Independent Set is at least as hard as Clique, so k -Independent Set is NP-hard.

TAS - Polynomial time approx scheme

Takes in ϵ and problem instance α^k and produces $(1+\epsilon)$ approx of optimal solution S and ϵ poly in n and ϵ

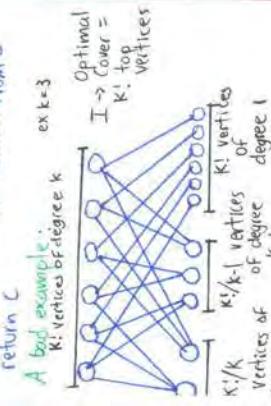
Greedy Center: Repeat - 2-approx Alg to cover take any $(u, v) \in E$. Add (u, v) to center remove all edges from u and v

Approx Vertex Cover (Natural)

```

 $C \leftarrow \emptyset$ 
 $E' \leftarrow E$ 
while  $E' \neq \emptyset$ :
    pick  $v$  with max degree
     $C = C \cup \{v\}$ 
    remove  $v$  and all incident edges from  $E'$ 
return  $C$ 

```



But algorithm may end up picking all the bottom vertices:

$$SOL = K! \left(\frac{1}{k} + \frac{1}{k-1} + \dots + 1 \right) \approx K! \log k$$

$\Rightarrow \log k$ worse than optimal!

to prove fully PTAS show that:

$$\text{OPT} \leq (1+\epsilon) \cdot R \quad \text{or} \\ \text{OPT} (1-\epsilon) \leq R$$

(in)approximation algorithm: cost C

produced by alg is within a Factor of $\rho(n)$ of the cost C^* of optimal solution

$$\max \left(\frac{L}{C^*}, \frac{C^*}{L} \right) \leq \rho(n)$$

Problem 4. Red-Blue Network Flow [30 points] (2 parts)

You are given a flow network $G = (V, E)$ with distinguished source vertex s and distinguished sink vertex t . As usual, the edges in E are directed edges. The edges in E are now, however, of two types: blue and red. The blue edges are of the usual type: each blue edge e has a capacity $c(e)$ that is the *maximum possible* amount of flow that can pass through that edge. The red edges are built out of strange alien technology: each red edge e must have at least a certain amount $d(e)$ of flow passing through that edge. That is, $d(e)$ is the *minimum required* amount of flow that must pass through that edge.

Other than the fact that some edges are now red, everything is as usual (e.g., flow must be conserved at all intermediate vertices, etc.). Each edge is either red or blue. We call such a network problem a "red-blue network flow problem". As usual, the goal is to find a flow of maximum possible value.

- (a) Show that it is possible to compute a maximum red-blue network flow, or determine that no legal flow exists, in polynomial time.
- (b) You now discover how to "turn off" red edges. If a red edge e is on, at least $d(e)$ units of flow must pass over it; if it is off, no flow can pass over it. You can turn some red edges on and others off. Show that it is now NP-complete to determine the maximum flow possible in a red-blue network.

Solution: (a) The problem can be formulated as a linear programming problem. Let $f(u, v)$ be the flow on edge $(u, v) \in E$. B and C be the sets of blue and red edges respectively. The LP is

$$\begin{aligned} \max \sum_e f(s, v) & \text{ s.t.} \\ f(u, v) \leq c(u, v) & \text{ for any } (u, v) \in B, f(u, v) \geq d(u, v) \text{ for any } (u, v) \in C \\ f(u, v) \geq 0 & \text{ for any } u \in V - s, t \end{aligned} \quad (1)$$

The LP can be solved in weakly polynomial time by ellipsoid or interior point algorithm.

(b) Since the NP-complete is for decision problems, the question can be interpreted as determining if there is a legal k flow in the flow network. To prove the problem is NP-hard, we can reduce the SUBSET-SUM problem to the red-blue flow network problem. The SUBSET-SUM is the decision problem that determines if there is a subset in $S = \{x_1, x_2, \dots, x_n\}$ whose sum equal to k . We first create a red-blue flow network with source s , sink t and another n vertices. The edges in the flow network are as follows,

(s, u_i) for $i = 1, 2, \dots, n$ are red edges that connects s to n vertices that represent n variables in the subset, $c(s, u_i) = x_i$.

(u_i, t) for $i = 1, 2, \dots, n$ are blue edges that connects n vertices that connect n to sink t , $d(u_i, t) = x_i$. The way we create the graph guarantees that the flow from s to t through u_i can either be x_i if the red edge (s, u_i) is on or 0 if the red edge is off. The flow in the red-blue network is therefore equal to the sum of the x_i whose corresponding red edge is on. Determine if there exists a flow k is the same as determine if there is a subset sum k . Creating the red-blue network takes polynomial running time. Since SUBSET-SUM is NP-complete, the red-blue flow network problem is NP-hard. Given the assignment of flow on each edge and the on/off assignment of red edges, it takes polynomial time to check if flow k exists, so the problem is in NP. Combining with the NP-hardness proof, the red-blue network-flow problem is NP-complete.

DP Vertex Cover To Blue: A vertex cover of an undirected graph is a subset of E that minimal vertex cover. Find minimal vertex cover for every edge (u, v) in the graph either v or u or both are in vertex cover.

$\sim O(d(C)) \Rightarrow$ size of min vertex cover for tree rooted at v if v is vertex v in min vertex cover or not

3) $d[C] = \min \left\{ 1 + \sum_{u \in C} d[u] \text{ over } \begin{array}{l} \text{such that } \\ \text{either } v \text{ or } u \text{ or both are in } \\ \text{vertex cover} \end{array} \right\}$

4) Solution:

$d[[k \times root]] \# \text{subproblems } O(n) \text{ Time } \propto O(1) \Rightarrow O(n)$

* If A is NP-hard and can be reduced to B, Then B is also NP Hard.

* Approximation factor not necessarily carried over in reductions.

Distributed Algorithms

Algs that run on networked/multiple processors and share memory.

Distributed Networks associate process w/ vertex Edge = two directed communication channel

Synchronous - Rounds

• Communicate using messages

Processes need not be distinguishable

Each Round:

Processes determine messages to send to all

Processes complete new state based on message

Leader Election. Goal: one process elects itself leader

Breaking Symmetry: No algorithm with indistinguishable processes work

This is because all nodes start the same, send the same message, and make same change.

Solved by VID (Unique ID) For each process randomly

large set: $\{1, \dots, r\}$ $r = n^2 / 2e$ n processes

w w w $\Pr[1 - \epsilon]$ all IDs are unique

Algorithm:

Choose random ID

Exchange IDs; if max unique > select, else repeat

Expected Time: $\leq 1/\epsilon$

w w P $\geq 1 - \epsilon$ algo finishes in one round

Maximal Independent Set (MIS)

Every node either both in set in S or has neighbor

can't add any more nodes in S

Each process should output "IN" or "OUT"

Lazy MIS Algorithm

Round 1: choose random val from 1 ... n. send to all

neighbours

If $r >$ all received join MIS. output "IN"

Round 2: If you joined, announce to all neighbors

If you receive announcement, don't join = quit dev

Repeat these rounds until all nodes decided

Correctness: Independence: when a node joins all

Max: only becomes inactive if it joined at join

termination: show that graph gets smaller on each phase.

Breadth First Spanning Tree

Distinguished vertex v0 = root knowledge about G

processes know their VID and neighbor's VID

Processes must produce BFS tree rooted at v0

Vertex of distance from v0 appears at depth

Places shade output parent.

marked! boolean

parent: vid or undefined

send: boolean

Cryptography Hashing

OW - one way

Infeasible, given $y \in \{0,1\}^d$ to find any $x \in \mathbb{R}^d$ s.t. $h(x) = y$

(R - collision resistance) Infeasible to find $x \in \mathbb{R}^d$ s.t. $x \neq x'$ and $h(x) = h(x')$

PRF - Pseudo-randomness

Behavior indistinguishable from Random Oracle

NH - Non-malleability

Infeasible given $h(x)$ to produce $h'(x')$ where x' and x are related

Random Oracle ideal but not achievable in practice.

• tries something

random every time, except as previous answers. (deterministic)

OW $\not\rightarrow$ CR, TCR (neither holds)

Collisions can be found in $O(2^{d/2})$

Inversion can be found in $O(2^d)$

Algorithm: each round:

Send distance (dist) message to all

Receive messages

Relaxation step \rightarrow dist = $\min_j (c_j + w_{j,j})$

If dist decreased set parent = j

Time Complexity: $O(n \cdot |E|)$

Termination w/ converge cast

Asynchronous Distributed Algorithms

"no rounds", "steps" executed one by one in any order FIFO queue for messages.

Synchronous BFS algo won't work

Strategy: Each process keeps track of hop distance, changes parent on new shorter path

and propagates improved distances

State: parent, dist, send(v) \rightarrow FIFO of N

Transitions: receive(m) in v
if $m_{i,j} < \text{dist}_i$, $\text{dist} = m_{i,j}$; parent = v
for every w : add $m_{i,w}$ to send(w)

remove head of send(v)

Message head(s) send(s)

Time: $O(|E| \cdot n \cdot d)$ \leftarrow Ignore local processing complexity. Instead traverse each link

For any node v, and any r <= diam, there is an at most r hop path from v0 to v

by time r node v is dist $\leq r$

Shortest Path Tree: Bellman Ford - correct for small paths, asymptotic

Transitions: receive(m) in v, parent = v
if $m_{v,w} < \text{dist}_w$, $\text{dist} = m_{v,w}$
for every w : add $m_{v,w}$ to send(w)

remove head of send(v)

Simple BFS Algorithm:

Processes mark themselves as they became part of tree. Initially s marked

Processes must produce BFS tree rooted at v0

Vertex of distance from v0 appears at depth

Places shade output parent.

marked! boolean

parent: vid or undefined

send: boolean

Aim for exact solution but isolate exponential term to a parameter

small parameter = fast solution

Kenneth Levenson simplifies self reduction Poly time alg converting input into small equivalent input

Ex: remove vertices you know will be stale and add them at the end

Counting # of nodes:

Set up a spanning tree and sum up the nodes at the end during conversion

BRSTree is complete!

Counting # of nodes:

Set up a spanning tree and sum up the nodes at the end during conversion

BRSTree is complete!