

6.046 Spring 2024

LECTURE 1 2/6/24 11AM

Interval Scheduling

- single resource
- requests $R = \{r_1, \dots, r_n\}$ as intervals
- compatibility: $r_i \text{ comp. } r_j \Leftrightarrow r_i \cap r_j = \emptyset$
- goal: find a largest set $S \subseteq R$ of compatible requests

↳ greedy algorithm: earliest finish time

• use exchange/hybrid argument to compute S_{OPT} with
 S_{GREEDY} and do swapping to realize they have
same cardinalities

• sorting takes $O(n \log n)$ → linear scan → $\boxed{O(n \log n) \text{ total}}$

Weighted Interval Scheduling

→ intervals now have weights

Dynamic Programming

$$W(R) = \max \{ W(R \setminus \{r_n\}), w_n + W(R \setminus INC(r_n)) \}$$

Sorted by
end times ↑

weight
 $R \setminus \{r_n\}$ are both prefixes of $\{r_i\}$
 $R \setminus INC(r_n)$

incompatibility
set ↓

$R \setminus INC(r_n) = \{r_1, r_2, \dots, r_i\}$ for some $i < n$

finding i takes $O(\log n)$ via binary search

→ total time $O(n) \times O(\log n) = O(n \log n)$.

Divide & Conquer.

Median Finding, Integer Multiplication

→ split into a problems of size $\frac{n}{b}$

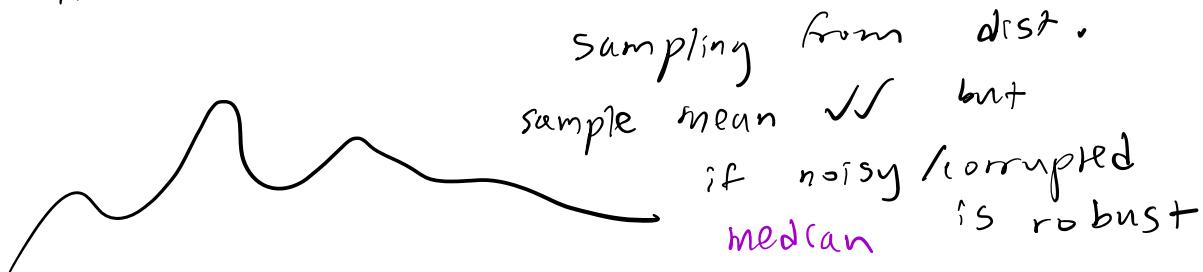
$$\Rightarrow T(n) = a \cdot T\left(\frac{n}{b}\right) + \{ \text{time to combine, etc.} \}$$

Median Findinginput: set S of n distinct #s

$$\text{rank}(x) := \# \underset{\in S}{\text{elements}} \text{ in } S \leq x$$

$$\begin{aligned} \text{upper median} &= \text{rank } \lceil \frac{n+1}{2} \rceil \text{ element } \\ \text{lower median} &= \text{rank } \lfloor \frac{n+1}{2} \rfloor \text{ element } \end{aligned} \quad \begin{cases} \text{same if } n \text{ odd} \\ \text{if } n \text{ even} \end{cases}$$

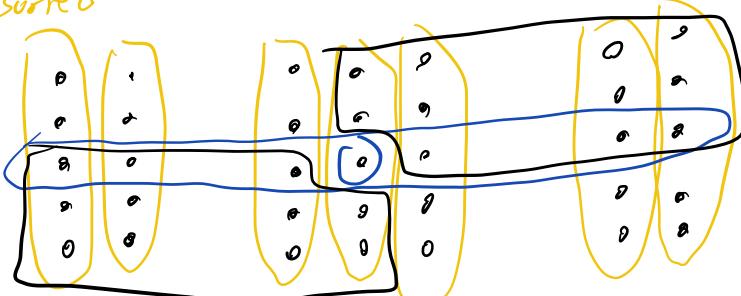
motivation (Robust statistics)

methods sort then find $\rightarrow O(n \log n)$ check all elements $\rightarrow O(n)$ per element $\Rightarrow O(n^2)$.modified Brute Force - compute $L(x) = \{k \mid k < x, k \in S\}$
 $G(x) = \{k \mid k > x, k \in S\}$ for $x \in S$. Then search in bigger bucket for median→ but worst-case $= O(n^2)$ tooDef Say $x \in S$ is c -balanced if $\max \{|L(x)|, |G(x)|\} \leq c \cdot n$ Goal Find some c -balanced elements in time $O(n)$

Subroutine for finding a c-balanced element

- (1) divide S into n/s buckets of size 5.
- (2) compute median of each bucket $\rightarrow O(n)$ time
- (3) compute "median-of-medians" recursively -

assume each
bucket sorted



$$E[G(x)] \sim \frac{3n}{10}$$

$$E[L(x)] \sim \frac{3n}{10}$$

$$\Rightarrow T(n) = T\left(\frac{n}{2}\right) + T\left(\frac{3n}{10}\right) + O(n)$$
$$\Rightarrow T(n) = n$$

Karatsuba - $T(n) = 3T\left(\frac{n}{2}\right) + O(n) \Rightarrow T(n) = n^{\log_2 3}$

Hurley, Van der Hoeven $\leadsto T(n) = O(n \log n)$

LECTURE 3 2/13/24 (Virtual), did not attend)

Monte Carlo

- always efficient
- $P(\text{right}) \rightarrow \text{high}$

Las Vegas

- efficient in expectation
- always right

average-case
(input)

expected-case
(random choice)

Matrix Products

$C = A \times B \rightarrow O(n^3)$ multiplications for $n \times n$.

Strassen's

$[2 \times 2] \times [2 \times 2]$ in 7 multiplications
 $\Rightarrow O(n^{\log_2 7}) \rightarrow O(n^{2.81})$

Matrix Product Verification

Given A, B, C , check $AB \stackrel{?}{=} C$.

↳ can we do better than matrix-multiply??

Friarvald's Algorithm wlog $A, B, C \in \mathbb{F}_2^{n \times n}$

- choose random binary vector $r \in \mathbb{F}_2^n$

If $\begin{cases} ABr = Cr & \text{YES} \\ \text{else} & \text{NO} \end{cases} \Rightarrow O(n^2)$

easy to see if $AB = C \Rightarrow \text{YES}$

$AB \neq C \Rightarrow \text{NO}$ with high probability?

Claim If $C \neq AB$, $\Pr[ABr \neq Cr] \geq \frac{1}{2}$.

Proof If for some r we have

$$(C - AB)r = 0, \text{ we have } (C - AB)(r + \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}),$$

$$(C - AB)\left(r + \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}\right), \dots, (C - AB)\left(r + \begin{bmatrix} 1 \\ \vdots \\ 0 \end{bmatrix}\right) \uparrow$$

one of
these $\neq 0$.

let wlog $(C - AB)\left(r + \begin{bmatrix} 1 \\ \vdots \\ 0 \end{bmatrix}\right) \neq 0$.

Then we have a $|r|$ for each r s.t.

$(C - AB)r = 0, r' = r + \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$ has $(C - AB)r' \neq 0$.

so $\Pr[ABr \neq Cr] \geq \frac{1}{2}$. \square

recall rank-finding from Lecture 2.

worst-case to find rank is $O(n^2)$

↳ how can we "cleverly" pick pivots to

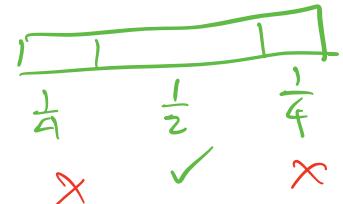
ensure rank of pivot $\geq \frac{1}{4}n, \leq \frac{3}{4}n$??

Quickselect



random pivot is enough, $\frac{1}{2}$ chance!

"Las Vegas Algorithm"



Paranoid Quicksort

keep repicking pivot until rank $\geq \frac{1}{4}n, \leq \frac{3}{4}n$ before divide + conquer
 \Rightarrow expected $\leq 2cn = O(n)$
 $T(n) = O(n) + T(\frac{3}{4}n)$
 $\Rightarrow O(n).$

Quicksort

- ① Pick a pivot
- ② sort elements before and after pivot
- ③ recurse on two halves L, R .

basic - pick 1st or last
randomized
median-find

} $O(n^2)$ worst-case
 $O(n \log n)$ avg case
} $O(n \log n)$ all cases
→ impractical
against mergesort
though

Paranoid Quicksort using Paranoid Quicksort

- $O(n)$ paranoid quickselect

expected

↓
expected

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

$$+ O(n) = T\left(\frac{n}{4}\right) + T\left(\frac{3n}{4}\right) + O(n)$$

$$\Rightarrow O(n \log n) \text{ expected}$$

Union Bound

$$P(A_1 \cup A_2 \cup \dots \cup A_n) \leq \sum P(A_i)$$

Markov's Inequality

$$P[X \geq c \mathbb{E}[X]] \leq \frac{1}{c}$$

Las Vegas
w/ exp time T

run for
 cT

Monte Carlo
w/ success
chance $\geq 1 - \frac{1}{c}$

why does randomized quicksort give sorted w.h.p. in $O(n \log n)$?

OR, with recursion depth L :

let $\gamma_{i,k}$ = indicator of if container with a_i at k^{th}
recursion chooses bad pivot.

$$\mathbb{E} \left[\sum_k \gamma_{i,k} \right] \leq \sum_i \underbrace{0.5 L}_{\substack{\uparrow \\ 0.5 \text{ per level}}} \leq 0.5 L$$

$$\Pr \left[\sum_k \gamma_{i,k} \geq (1+0.2) \cdot 0.5L \right] \leq e^{-0.04 \cdot 0.5L / 3} = e^{-L/150}$$

↓

let $L = 300 \log n$ \longrightarrow

$$\Pr \left[(\text{any } \sum_k \gamma_{i,k} \text{ over all } i) \geq 0.6L \right] \leq \frac{1}{n^2} \cdot n = \frac{1}{n}$$

↳ at least $0.4L$ good pivots

$$\Rightarrow n \left(\frac{3}{4} \right)^{0.4L} = n \cdot \left(\frac{3}{4} \right)^{120 \log n} < n^{-33} \leq 1,$$

so $\Pr[\text{termination}] \geq 1 - \frac{1}{n} \rightsquigarrow 1$ for
 n large \square . ↑
So choice of L is good!

Setting Data structure, member functions called repeatedly.

Can we more accurately capture cost of using the data structure?

Case Study Union-Find / Disjoint Sets

3 Operations

Make-set (x) \rightarrow add $\{x\}$ to collection

Find-set (x) \rightarrow return $\text{REP}[S(x)]$

\uparrow
representative of set of x .

Union (x, y) \rightarrow replace sets $S(x), S(y)$ with
single $S(x) \cup S(y)$ with arbitrary
single representative

Application

connected components
in graphs

Implementation

- doubly-linked lists



make-set $\Rightarrow O(1)$

find-set $\Rightarrow O(n)$, follow pointers
to head

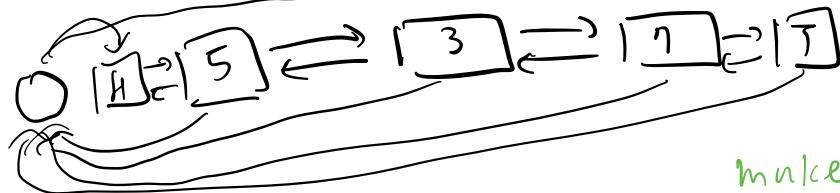


union \Rightarrow find tail of x
head of y , and
link $\Rightarrow O(n)$

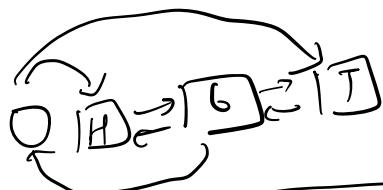
Amortized Cost

An operation has amortized cost T if any sequence of k operations has cost $\leq kT$

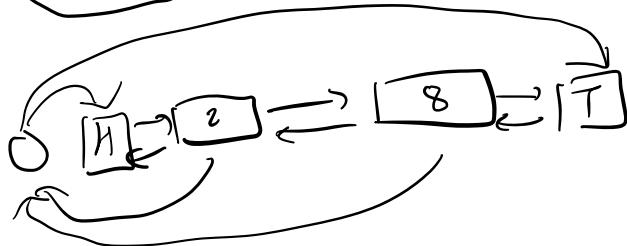
Implementation 2 - doubly-linked lists w/ head/tails



make-set $\Rightarrow \Theta(1)$



find-set $\Rightarrow \Theta(1)$ by finding head
union $\Rightarrow \Theta(n)$ if need to
update all pointers of $S[y]$



no Θ ↓

$\Theta(\min\{IS(x), IS(y)\})$

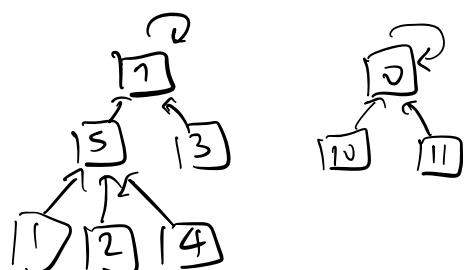
Claim Amortized cost of union is $O(\log n)$

Proof WTS $\{\text{total # of pointer updates}\} \leq O(\log n)$

↓
each element x 's pointer updated $O(\log n)$ times.
 \Rightarrow but pointer only updated if set size more
than doubles $|S(y)| \rightarrow |S(x) \cup S(y)| \geq 2|S(y)|$

□

Implementation 3 - trees



make-set $\Rightarrow \Theta(1)$

worst case

Find-set $\Rightarrow \Theta(\text{height}(S(x))) \rightarrow \Theta(n)$

union \Rightarrow merge two trees by updating
root of shorter one to point to
root of larger

$\Theta(\text{height}(S(x)) + \text{height}(S(y))) \rightarrow \Theta(n)$

worst case

→ union-by-rank.

ensure height of each tree = $O(\log n)$

combine shorter tree into larger tree

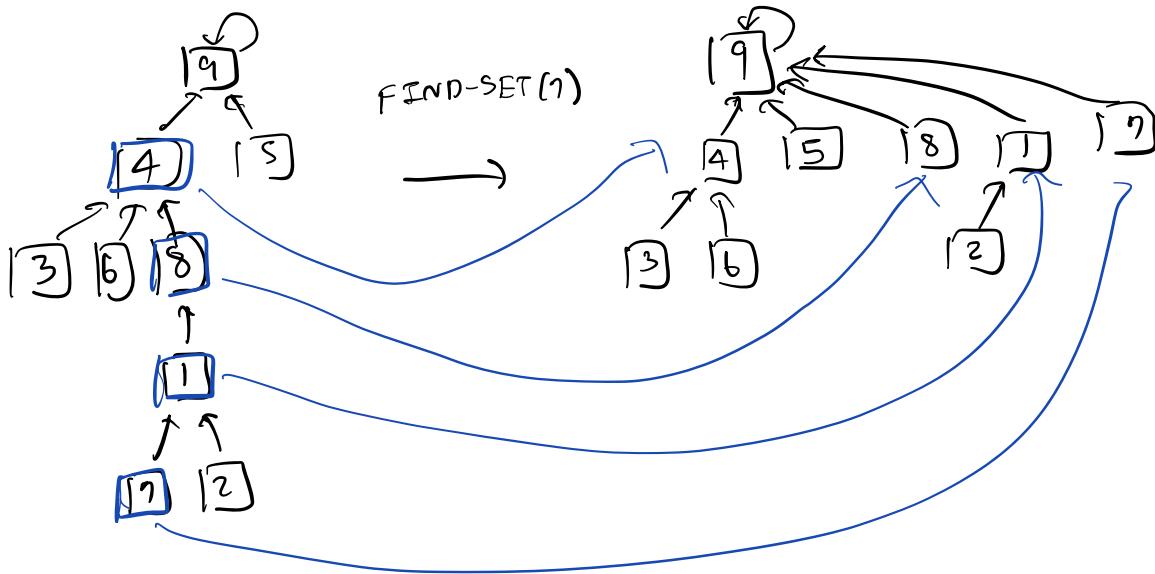
find-set, union $\Rightarrow \Theta(\log n)$

amortized

proof is by induction \Rightarrow show first that all tree's have height $O(\log n)$ after n operations.

then show each element's pointer is updated $\text{depth}(x) = O(\log n)$ times \square .

\rightarrow path compression



Claim w/o Union-by-rank, w/ path compression
amortized cost is $O(\log n)$

Pf Let c_i = cost of i^{th} operation.

WTS $\forall k, \sum_{i=1}^k c_i = O(k \log n) \rightarrow$ but c_i can vary.

\Rightarrow Potential Functions

Suppose we define \hat{c}_i s.t. $\sum_{i=1}^k c_i \leq \sum_{i=1}^k \hat{c}_i$

$$\hat{c}_i = c_i + \Delta \Phi_i$$

$$= \Phi_i - \Phi_{i-1}$$

both ≥ 0

let

$$\Phi(D) = \sum_{x \in D} \log \text{size}(x)$$

use $D_0 = \text{initial state of all}$

$$\text{singletons} \Rightarrow \Phi(D_0) = 0,$$

of nodes
in subtree
rooted at x .

$$\text{make-set: } \hat{c}_i = c_i + \Delta \Xi_i \rightarrow \hat{c}_i = O(1) + 0 = O(1).$$

\uparrow
no change

find-set: let $x, u_1, u_2, \dots, u_h = \text{REP}(S(x))$ be the node-to-root path we traverse. $c_i = h$

We have $u_j.\text{SIZE} \rightarrow u_j.\text{SIZE} - u_{j-1}.\text{SIZE}$,
upon path compression

$$\Delta \Xi_i = \sum_{j=1}^{h-1} \log\left(\frac{u_j.\text{SIZE} - u_{j-1}.\text{SIZE}}{u_j.\text{SIZE}}\right)$$

$$\hat{c}_i = h + \sum_{j=1}^{h-1} \log\left(\frac{u_j.\text{SIZE} - u_{j-1}.\text{SIZE}}{u_j.\text{SIZE}}\right)$$

$$= 1 + \sum_{j=1}^{h-1} 1 + \log\left(\frac{u_j.\text{SIZE} - u_{j-1}.\text{SIZE}}{u_j.\text{SIZE}}\right)$$

note:

$$1 + \log\left(\frac{u_j.\text{SIZE} - u_{j-1}.\text{SIZE}}{u_j.\text{SIZE}}\right) \geq 0$$

$$\Leftrightarrow \log\left(1 - \frac{u_{j-1}.\text{SIZE}}{u_j.\text{SIZE}}\right) \geq -1$$

$$\Leftrightarrow \frac{u_{j-1}.\text{SIZE}}{u_j.\text{SIZE}} \leq \frac{1}{2}$$

$$\Leftrightarrow u_j.\text{SIZE} \geq 2u_{j-1}.\text{SIZE}$$

$$= 1 + \log n$$

$$= O(\log n)$$

because u_j can only double at most $\log n$ times
to stay under n
(since u_j is monotonically increasing)

union :

$$\begin{aligned} \hat{\ell}_i &= 1 + \log(x.\text{size} + y.\text{size}) - \log(y.\text{size}) \\ &= 1 + \log\left(1 + \frac{x.\text{size}}{y.\text{size}}\right) \\ &= O(\log n) \quad \square. \end{aligned}$$

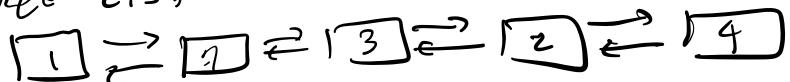
link
 $\xrightarrow{\text{link}}$
 $x \text{ root to } y \text{ root}$

together union-by-rank + path compression $\Rightarrow O(\alpha(n))$.

LECTURE 5 2/22/24 11AM

self-organizing lists

linked-list



stores n elements

w/ distinct keys

one operation - $\text{ACCESS}(x)$ = start from head of L ,
↓ find element with key x .

time / cost = $\text{rank}_L(x)$

suppose we have a sequence of N accesses

→ want to minimize total cost

$\min = N$ $\max = nN$

but, we can also reorder (via transpositions) L after each access → each transposition costs 1

$C_A(L_i)$ = total cost of $\text{access}(x_i)$

$= \text{access cost} + \text{reordering cost}$

$= \text{rank}_L(x_i) + \# \text{ transpositions after } x_i$

offline algorithm - knows all access calls to the array beforehand

$$\text{ex)} L = [1 \ 2 \ 3 \ \dots \ n]$$

$$S = n, n, \dots, n \quad |S| = n.$$

OPT: access n , move n to front, keep accessing n
 $\underbrace{n}_{n-1} \quad \underbrace{\quad \quad \quad}_{|S|-1}$

$$\text{total cost } 2n-1 + |S|-1 = (2n-2) + n.$$

worst case example

$$L = [1 \ 2 \ 3 \ \dots \ n]$$

$$S = n, n-1, n-2, \dots, 1$$

divide L into

$$\begin{array}{c} \boxed{1 \quad 1} \\ \frac{n}{3} \quad \frac{n}{3} \quad \frac{n}{3} \end{array}$$

each of the first $\frac{n}{3}$ calls must either:

- be moved to first $\frac{n}{3}$ spots in list $\rightarrow \text{cost} \geq \frac{n}{3}$
 OR

- be accessed somewhere outside of first $\frac{n}{3}$ spots in list $\rightarrow \text{cost} \geq \frac{n}{3}$

$$\text{so cost is } \geq \frac{n}{3} \cdot \frac{n}{3} = \Omega(n^2)$$

$\uparrow \quad \uparrow$
elements cost

worst case

for online

adversary asks for Last element every time $\Rightarrow O(|S|n)$
 to get access calls one at a time

Move to Front (MTF) Heuristic

after accessing x_i , move x_i to front of L

$$\Rightarrow \text{cost} = \underset{\text{access}}{\underset{\uparrow}{\text{rank}_L(x_i)}} + \underset{\text{move to front}}{\underset{\uparrow}{\text{rank}_L(x_i)}} - 1$$

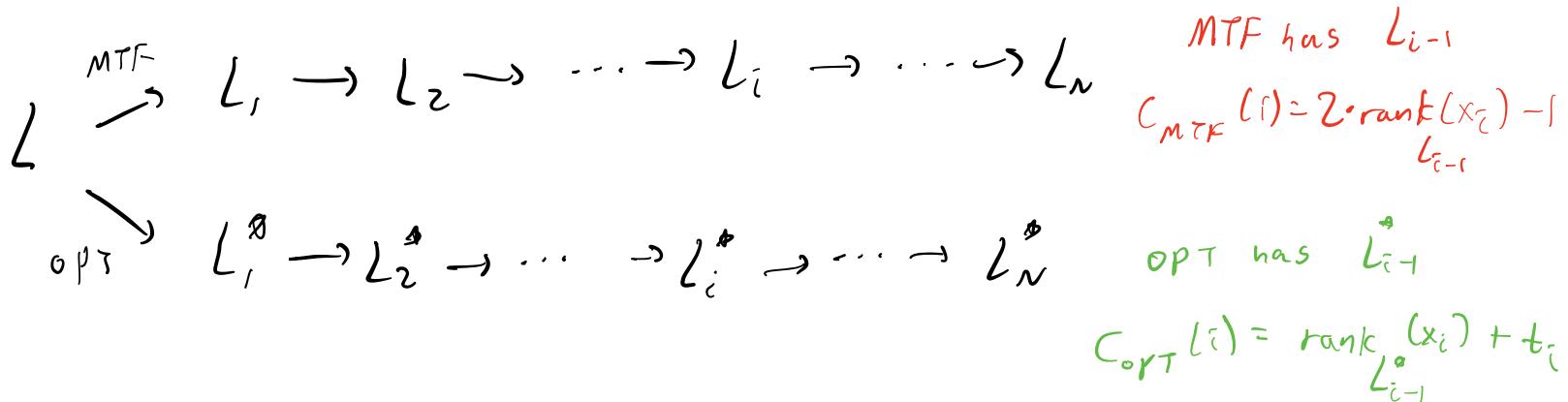
Defn An online algorithm A is α -competitive if there is some cost c s.t. \forall access sequences S ,

$$\text{total cost of } A \text{ on } S \leq \alpha(\text{total cost of OPT}) + c$$

↑
optimal offline

Thm MTF is 4-competitive.

Proof Let $S = x_1, x_2, \dots, x_n$ access sequence
 L_i = list that MTF has after accessing x_i
 L_i^* = first that OPT has after accessing x_i at i^{th} access x_i



we want a potential function to compare L_i^* , L_i .

Def $\text{distance}(L_i, L_i^*) = \min(\# \text{of transpositions to transform } L_i \text{ to } L_i^*)$

proof relies on

fact that $\rightarrow = \# \text{ of inversions}$

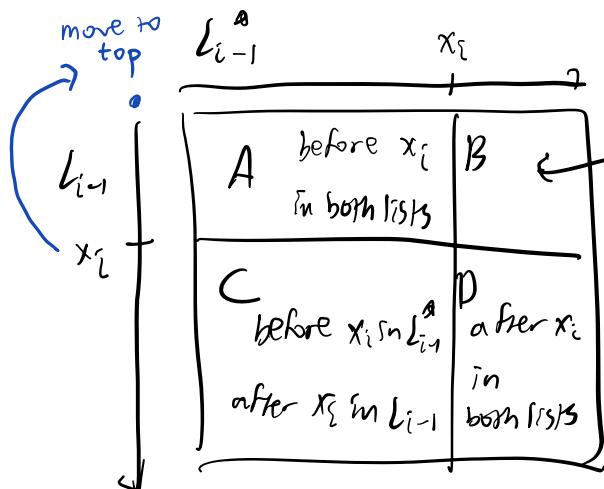
$\# \text{ inversions} \leq \# \text{ transpositions}$ ↗
 $(\# \text{ of pairs in } L_i \text{ that have swapped order in } L_i^*)$

$\# \text{ transpositions} \geq \# \text{ inversions}$
 ↑
 can fix ≥ 1 inversion each time
 every inversion must be fixed

Potential Function $\varphi(i) = 2 \cdot \# \text{ inversions between } L_i \text{ & } L_i^*$

$$\hat{C}_{MTF}(i) = C_{MTF}(i) + \varphi(i) - \varphi(i-1)$$

recall $C_{MTF}(i) = 2 \cdot \text{rank}_{L_{i-1}}(x_i) - 1$

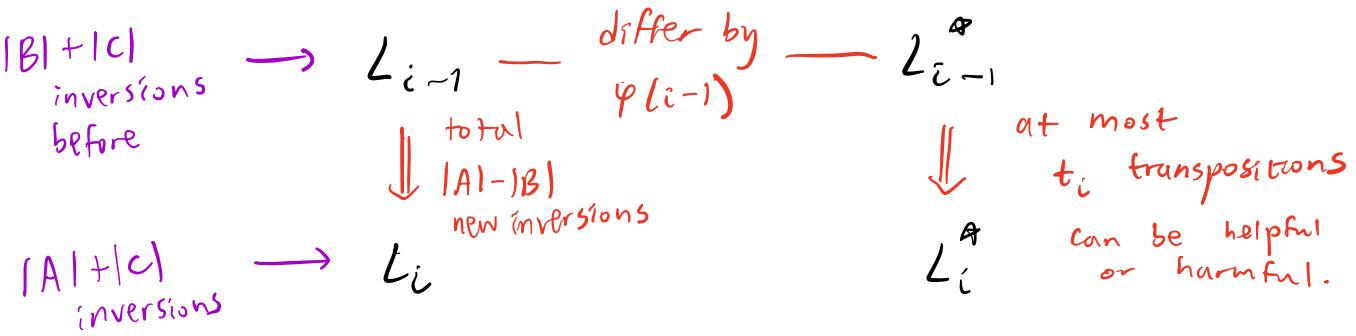


$$L_{i-1} = \boxed{\begin{array}{c|c} A \cup B & C \cup D \end{array}}$$

$$\text{rank} \rightarrow |A| + |B| + 1$$

$$L_i^* = \boxed{\begin{array}{c|c} A \cup C & B \cup D \end{array}}$$

$$\text{rank} \rightarrow |A| + |C| + 1$$



$$\Rightarrow \varphi(i) - \varphi(i-1) \leq 2(|A| - |B| + t_i)$$

$$\begin{aligned}
 \Rightarrow \hat{c}_{MTF}(i) &= c_{MTF}(i) + \varphi(i) - \varphi(i-1) \\
 &\leq 2|A| + 2|B| + 1 + 2|A| - 2|B| + 2t_i \\
 &= 4|A| + 2t_i + 1 \leq 4\text{rank}_{L_{i-1}^*}(x_i) + 4t_i \leq 4c_i^*
 \end{aligned}$$


note $\text{rank}_{L_{i-1}^*}(x_i) = |A| + |C| + 1 \Rightarrow |A| \leq \text{rank}_{L_{i-1}^*}(x_i) - 1$

so $\hat{c}_{MTF} = c_{MTF} + \varphi(|S|) - \varphi(0)$

$c_{MTF} \leq \hat{c}_{MTF} \leq 4c_{OPT}$

□

LECTURE 6 2/27/24

Hashing

Dictionary has $\text{INSERT}(x)$, $\text{DELETE}(x)$, $\text{SEARCH}(x)$

$x = (k, v)$ key-value pair

direct access



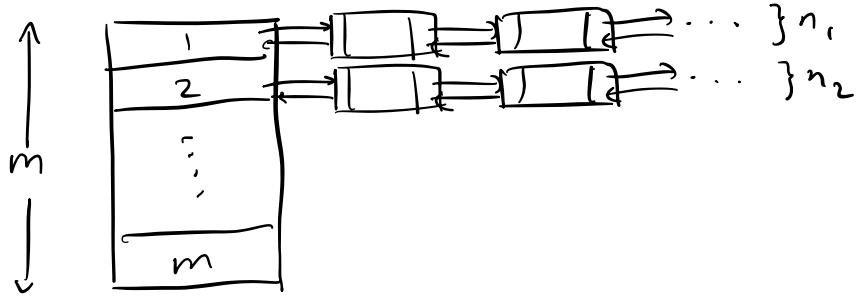
n keys
 m slots

$|U|$ possible keys

SOLUTION 0 - Doubly Linked List $O(n)$ Insert/Delete, $O(n)$ Search

Direct Addressing - $O(1)$ all operations but $|U|$ space

Chaining



INSERT, DELETE $O(1)$

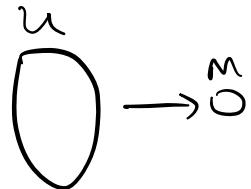
SEARCH $O(n_i)$

best case $n_i = \frac{1}{m} = \alpha$ load factor

worst case $n_i = \Theta(n)$

Hash Functions

$$h: \mathcal{U} \rightarrow M = \{0, 1, \dots, m-1\} \quad m \ll |\mathcal{U}|$$



collisions - k_1, k_2 s.t. $h(k_1) = h(k_2)$

random oracle h

- when given a new key k , choose a random location, save it

↳ can give us expected load $\frac{n}{m} = \alpha$

few collisions

for random oracle h , for any $k_1 \neq k_2$, $P[h(k_1) = h(k_2)] = \frac{1}{m}$

Family of Hash Functions

$$\mathcal{H} = \{h: \mathcal{U} \rightarrow M\}$$

UNIVERSAL HASHING for any $k_1 \neq k_2$, $P[h(k_1) = h(k_2)] \leq \frac{1}{m}$.

how to get such a family \mathcal{H} ? base m representation

$$\mathcal{H} = \{h_a: a \in \mathcal{U}\} \quad a = (a^{(e-1)}, \dots, a^{(1)}, a^{(0)})$$

$$\text{for } k = (k^{(e-1)}, \dots, k^{(0)}), \quad h_a(k) = \sum_{i=0}^{e-1} a^{(i)} k^{(i)} \pmod{m}$$

$$P[h(k_1) = h(k_2)] = P[a^T(k_1 - k_2) = 0] = P[a = 0] = \frac{1}{m} \pmod{m}$$

Open Addressing

→ probe sequence $h(k, p) = p^{\text{th}}$ preferred location of k

→ uniform (permutation) hashing: probe sequences of different keys
are independent uniform permutations
↳ all operations are expected $O(\frac{1}{1-\alpha})$

linear probing: $h_{\text{LIN}}(k, p) = h(k) + p \bmod m$

quadratic probing: $h_{\text{QUAD}}(k, p) = h(k) + c_1 p + c_2 p^2 \bmod m$ for c_1, c_2

double hashing: $h_{\text{DHASH}}(k, p) = h_1(k) + p \cdot h_2(k) \bmod m$
both random hash functions

resizing to lower $\alpha \rightarrow$ low amortized cost !!

LECTURE 7 2/29/24 11AM

n keys
 m slots
 $|U|$ possible keys

| | Space | Insert | Delete | Search | Build |
|-------------------|----------|------------------------------|------------------------------|------------------------------|-------------|
| SOLUTION ZERO | $O(n)$ | $O(1)$ | $O(1)$ | $O(n)$ | $O(n)$ |
| DIRECT ADDRESSING | $O(U)$ | $O(1)$ | $O(1)$ | $O(1)$ | $O(n)$ |
| CHAINING | $O(m+n)$ | $O(1)$ | $O(1)$ | $\exp O(\alpha)$ | $O(n)$ |
| OPEN ADDRESSING | $O(m)$ | $\exp O(\frac{1}{1-\alpha})$ | $\exp O(\frac{1}{1-\alpha})$ | $\exp O(\frac{1}{1-\alpha})$ | $\exp O(n)$ |
| PERFECT HASHING | $O(n)$ | - | - | $O(1)$ | $\exp O(n)$ |
| CUCKOO HASHING | $O(n)$ | $\exp O(1)$ | $O(1)$ | $O(1)$ | $\exp O(n)$ |

$m=n$ $m > 2n$

STATIC dictionary

- no INSERTS/DELETES

- perfect hashing scheme

↳ want **no** collisions on given set $K \subseteq U$, $|K|=n$.

No Collisions

Take universal hash family $\mathcal{H} = \{h: \mathcal{U} \rightarrow M\}$

find: $P[\text{no collisions on } K \subseteq \mathcal{U}]$

given two keys, $P[\text{collision}] \leq \frac{1}{m}$

over all pairs of keys, $P[\text{collision}] \leq \frac{\binom{m}{2}}{m} \leq \frac{n^2}{2m}$

make $m = n^2 \rightarrow P[\text{no collisions}] \geq \frac{1}{2}$

expected \leqq random calls to \mathcal{H} for no collisions.

FKS Perfect Hashing

① Pick h_1 from universal hash family $\mathcal{H}_1 = \{h_1: \mathcal{U} \rightarrow M\}$

• hash all keys in K with h_1 , (1.5) if $\sum_{i=0}^{m-1} n_i^2 \geq 4n$, resample
(expected 2 samples)

• let $n_i = \# \text{ keys with } h_1(k) = i$ $M_i = \{0, 1, 2, \dots, n_i^2 - 1\}$

② For each i , pick $h_{2,i} \in \mathcal{H}_{2,i} = \{h_{2,i}: \mathcal{U} \rightarrow M_i\}$

• resample until no collisions

↳ Space used = $O(m + \sum_{i=0}^{m-1} n_i^2)$ set $m = n$. \rightarrow expected $O(m+n) = O(n)$

let N_i = random variable for n_i

$$\mathbb{E}\left[\sum_{i=0}^{m-1} N_i^2\right] = \sum_{(k_1, k_2) \in K^2} \mathbb{E}[X_{h_1(k_1) = h_1(k_2)}]$$

indicator for
if $h_1(k_1) = h_1(k_2)$

$$= \sum_{k \in K} 1 + \sum_{(k_1, k_2) \in K^2} \mathbb{P}[h_1(k_1) = h_1(k_2)]$$

$$\leq n + \frac{1}{m} \cdot 2 \binom{n}{2} \leq n + \frac{2n^2}{2m} \leq n + \frac{2n^2}{2n} \leq 2n.$$

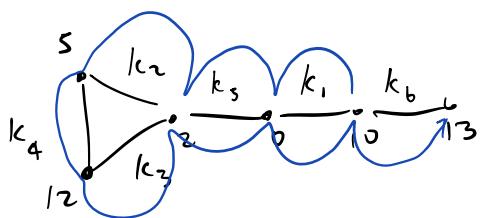
Markov $P\left[\sum_{i=0}^{m-1} N_i^2 > 4n\right] < \frac{1}{2}$, so $P\left[\sum_{i=0}^{m-1} N_i^2 \leq 4n\right] > \frac{1}{2}$

so add step (1.5) above \rightarrow then worst-case build is space is search is $O(1)$. exp $O(n)$ $O(n)$

Cuckoo Hashing

- two hash functions h_1, h_2
 - Malloc list of length $> 2|K|$
 - every key k to $h_1(k)$ or $h_2(k)$
 - option 0
 - option 1
 - add keys 1 at a time and if there's a collision, kick the old one out. the old one cf in option i goes to option 1-i and repeat

we can represent this as a graph, where keys k are edges between $h_1(k)$ and $h_2(k)$.



time for INSERT
= length of path

if we start with
 $k_1 \rightarrow 10$ then adding k_6
 $k_5 \rightarrow 0$ with $h_1(k_6) = 10, h_2(k_6) = 13$
 $k_3 \rightarrow 12$ makes a path on
 $k_4 \rightarrow 5$
 $k_2 \rightarrow 2$ our cuckoo graph

If we bound the number of cycles and lengths of paths between two keys $h_1(k_1)$, $h_2(k_1)$, $h_1(k_2)$, $h_2(k_2)$ we get

$O\left(\frac{1}{m}\right)$ expected probability of a collision
 \Rightarrow expected cost of insert = $O(1)$

→ more hash functions to allow for more slots gives better load factor restrictions than $\frac{1}{2}$.

frees

spanning tree = tree of connected $G = (V, E)$

with vertex set V and edge set $E' \subseteq E$

minimum spanning tree — not necessarily unique

Mega - greedy - MST

$$A = \emptyset \text{ (empty)}$$

for $|V|-1$ iterations:

Find a "safe" e to add to A

Avesel maintains inurlant

$A \leftarrow A \cup \{e\}$

return A

runtime; (greedy, irreversible)

$(n-1) \cdot (\text{cost of finding safe edge})$

e is safe if

e is safe if

$A \cup \{e\}$ is a subset
of edges of some MST.

minimum-weight edge e^* of G is always safe).

A cut of a graph G is a partition $S \cup V \setminus S$.

An edge (u, v) crosses the cut if $u \in S, v \in V \setminus S$ or vice versa.

A cut respects A if no edge of A crosses the cut
(i.e. all on one side)

An edge crossing a cut is a light edge if it has minimum weight among all edges crossing the cut.

Thm Let A be any set of edges of $G = (V, E)$ in some MST. Let $(S, V \setminus S)$ be any cut respecting A . Let e be a light edge for the cut. Then e is safe.

PF Consider an MST T containing A . Let $e^* = (a, b)$ be our light edge.

Let $P_{(a,b)}$ be the path in T from a to b .

must have some edge e crossing the cut.

Take $T' = T \setminus \{e\} \cup \{e'\}$ which is still a spanning tree. □

Kruskal's greedy step

Find minimum weight edge not forming a cycle.

Proof The edge is a light edge between a tree and $G \setminus \{tree\}$.

→ Union - find can give

$$O(m \log m + n \cdot T_{\text{make-set}} + n \cdot T_{\text{union}} + m \cdot T_{\text{find-set}})$$

↑
Sort E
in non-decreasing
order

↑
each vertex
starts as
a "tree"

↑
combine
two
trees

↑
to connect
trees

$$O(\text{sort}(m) + m \alpha(n))$$

↓
 $O(m \log n)$
otherwise

↓
 $O(m \alpha(n))$ if E sorted
OR E can be radix sorted

Prim's Greedy step

Find minimum weight edge incident to growing tree.

Proof The edge is a light edge between growing tree and rest of graph

→ Dijkstra's can give $O(m + n \log n)$ via Fibonacci Heaps.

when to use each??

Kruskal's if E's are sorted / easy to sort

or if m is small enough $\rightarrow m \leq O(n \log n)$

Max-Flow

Input: A network = directed graph $G = (V, E)$

source $s \in V$ sink $t \in V$

edge capacities $c: E \rightarrow \mathbb{R}^{>0}$

let $c(u, v) = 0$ if $(u, v) \notin E$

(CLRS) Flow → inconvenient for flow cycles
we care about net flow



let "gross flow" $g: E \rightarrow \mathbb{R}^{>0}$
 $g(e)$ = flow along $e \in E$

feasibility: $0 \leq g(u, v) \leq c(u, v) \quad \forall (u, v) \in E$

flow conservation: If $v \neq s, t$

$$\sum_{(x, v) \in E} g(x, v) = \sum_{(v, y) \in E} g(v, y)$$

flow in = flow out

(6.1220) Flow - don't assume s, t are source/sink
can be < 0

let "net flow" $f: V \times V \rightarrow \mathbb{R}$ \leftarrow $(u, v), (v, u) \notin E$

feasibility: $f(u, v) \leq c(u, v) \rightarrow \Rightarrow f(u, v) = 0$.
if $(u, v) \notin E$, $f(u, v) \leq 0$

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

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

$\hookrightarrow f(u, u) = 0 \quad \hookrightarrow f(u, v) = \text{net flow from } u \text{ to } v$

\hookrightarrow for every $v \neq s, t \quad \sum_{u: f(u, v) > 0} f(u, v) = \sum_{x: f(v, x) > 0} f(v, x)$

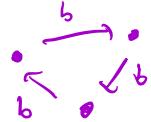
since $\sum_u f(u, v) = 0$ positive flow in = positive flow out

Let value of flow f $|f| = \sum_{v \in V} f(s, v) = \sum_{v \in V} f(v, t)$

Max Flow Problem

Given $G = (V, E, s, t, c)$ we want to find a flow of max value,

Climm Any flow can be decomposed into a collection of $s-t$ flow paths and flow cycles



Formally, let $\text{supp}_f(G) = \text{subgraph of } G \text{ of edges with } f(u, v) > 0$

Proof Induct on #edges of $\text{supp}_f(G)$.

If no edges, then $|f|=0$ and we are trivially done.

Now suppose we have $L > 0$ edges, and that all f' with $|f'| \geq 0$ whose $\text{supp}(G)$ has $< L$ edges can be decomposed.

IF we find a cycle C in $\text{supp}_f(G)$

Let $f_{\min}(C) = \min_{e \in C} f(e)$.

$$\text{new } f'(u, v) = \begin{cases} f(u, v) - f_{\min}(C) & \text{if } (u, v) \in C \\ -f'(v, u) = -(f(v, u) - f_{\min}(C)) & \text{if } (v, u) \in C \\ f(u, v) & \text{otherwise} \end{cases}$$

then $f' \rightarrow$ skew-symmetric by definition

\rightarrow flow conservation still works

$\rightarrow f'(u, v) \leq f(u, v) \leq c(u, v) \wedge (u, v) \text{ w/ } (v, u) \notin C$

$\rightarrow f'(u, v) = f_{\min}(C) - f(v, u) \leq 0 \leq c(u, v) \text{ for } (v, u) \in C$

Consider $\sum_v f'(s, v) = |F'|$.

IF C doesn't touch s , $|f'|$ unaffected

IF C goes into and out of s , $|f'| = |f|$ by skew-symmetry

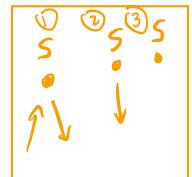
→ we've removed an edge from $\text{supp}_f(G)$ ($\underset{e \in E}{\text{argmin}} f(e)$) !

Finding the cycle

since all vertices except s, t have positive flow in = positive flow out,

- if $\exists (u, v) \in \text{supp}_f(G)$ with $v \neq t$, $\exists (v, u') \in \text{supp}_f(G)$
- if $\exists (u, v) \in \text{supp}_f(G)$ with $v \neq s$, $\exists (v', u) \in \text{supp}_f(G)$

3 cases for $f_1 \geq 0$



DFS from s , we will either
cycle or find $t \rightarrow$

edge out of t

↓ symmetric

① if there's an edge
into s , backward DFS

③ to get cycle or return to t

guarantees

if no edges into or out of s
(and also no edges into/out of t)

then we have
cycle too

then DFS from any other
 $v \in \text{supp}_f(G)$ to get a cycle

② otherwise \exists a path P out of s to t . (and no cycles including s or t)

$$\text{let } f_{\min}(P) = \min_{e \in P} f(e)$$

$$\text{new } f'(u, v) = \begin{cases} f(u, v) - f_{\min}(P) & \text{if } (u, v) \in P \\ -f(v, u) = -(f(v, u) - f_{\min}(P)) & \text{if } (v, u) \in P \\ f(u, v) & \text{otherwise} \end{cases}$$

→ still a flow as before!

no edges into $s \Rightarrow f(s, v) \geq 0 \quad \forall v \in V$

note that $f'(s, v) = f(s, v) \geq 0 \quad \forall v$ except first edge in $P: (s, v')$

then $f'(s, v') = f(s, v') - f_{\min}(P) \geq 0$ so

$$|f'| = \sum_v f(s, v) \geq 0 \quad \text{as desired.}$$

⇒ In both a cycle and path scenario, we've removed edges to induct. □

Now let f^* be a max flow in G (not necessarily unique)

Let $F^* = |f^*|$, How to determine if $F^* > 0$?

Let $G^* =$ subgraph of G with edge of capacities > 0

$$F^* > 0$$

only paths, not
cycles, contribute
to flow
 \iff
can send
positive flow
along the path

G^* has an
 $s \rightarrow t$ path

But how do we show that $F^* = 0$, or $\nexists s \rightarrow t$ paths in G^* ?

\rightarrow A cut?

let $S = \{v \in V \mid \exists s-v \text{ path in } G^*\}$

$t \notin S \iff \text{no } s-t \text{ path} \iff F^* = 0$

Def An $s-t$ cut is a cut $(S, V \setminus S)$ s.t $s \in S$
 $t \in V \setminus S$

let $c(S) = c(S, V \setminus S) = \sum_{u \in S} \sum_{v \in V \setminus S} c(u, v)$

be the capacity of the cut. \uparrow
only from $s \rightarrow v \setminus S$
not other direction

$c(S) = 0 \iff$ cut separates s from t .

$\Rightarrow F^* = 0 \iff \nexists s-t \text{ path in } G^* \iff \exists s-t \text{ cut with } c(S) = 0$

Minimun s-t cut problem

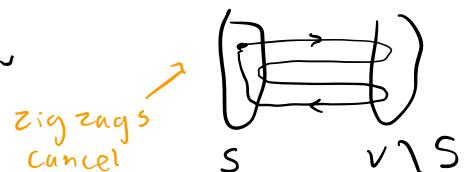
Given $(G = (V, E), s, t \in V, c)$ find an $s-t$ cut of minimum capacity

Trivially $f(S) \leq c(S)$ for any $s-t$ cut S by feasibility.

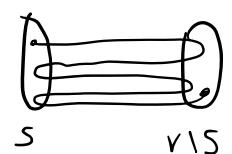
Claim For two $s-t$ cuts S, S' we have $f(S) = f(S')$.

Proof Separate f into flow cycles + $s-t$ flow paths.

→ cycles contribute 0 to flow for both $f(S)$ and $f(S')$



→ paths cancel out except for final crossing which contributes exactly the flow of the $s-t$ path to both $f(S)$ and $f(S')$



Note that $\{s\}, V \setminus \{s\}$ is an $s-t$ cut with $\text{Flow} = |f|$ so A $s-t$ cuts, $f(S) = |f|$.

⇒ Max flow \leq min-cut

$$F^* = |F^*| = f^*(S^*) \leq c(S^*)$$

Given a flow f , how to increase its value? (or conclude $F^* = |f|$?)

Idea Just try to find an $s-t$ path to push more flow through t .

Residual Network given $G = (V, E, s, t, c)$ and flow f

$$G_f = (V, E_f, s, t, c_f)$$

$$\hookrightarrow c_f(u, v) = c(u, v) - f(u, v) \geq 0$$

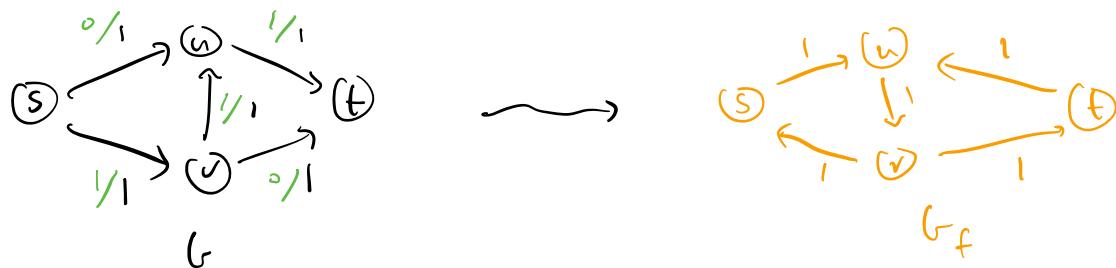
$$E_f = \{(u, v) \mid c_f(u, v) > 0\}.$$

If $(u, v), (v, w) \notin E$, neither has any capacity or flow

\Rightarrow neither is in E_f

$\Rightarrow E_f$ may consist only of edges in G and reverses of edges in G .

ex)



Max Flow Min Cut Theorem

TFAE.

① \exists $s-t$ cut $(S, V \setminus S)$ s.t. $c(S) = f(S) = |f|$

② f is a max flow

③ There is no $s-t$ path in G_f

Pf (① \Rightarrow ②) Let S^* be s.t. $(S^*, V \setminus S^*)$ is a min $s-t$ cut.
let f^* be a max flow

$$\begin{array}{c} |f| \leq |f^*| \leq c(S^*) \leq c(S) = f(S) = |f| \\ \max_{\text{flow}} \quad \text{weak duality} \quad \min_{\text{cut}} \quad \text{all } s-t \\ \end{array} \Rightarrow |f| = |f^*|. \quad \hookrightarrow \text{max flow} = \min \text{ cut}$$

(3) \Rightarrow (1) Suppose no s-t path in G_f .

$\Rightarrow \exists$ a cut $(S, V \setminus S)$ s.t. $c_f(S) = 0$

$$\Rightarrow 0 = c_f(S) = \sum_{u \in S} \sum_{v \in V \setminus S} c_f(u, v) = \sum_{u \in S} \sum_{v \in V \setminus S} c(u, v) - f(u, v) \Rightarrow \sum_{u \in S} c(u, v) = \sum_{u \in S} f(u, v)$$

(NOT 3) \Rightarrow NOT 1) Suppose P is an s-t path in G_f .
contrapositive

Let $c_f(P) = \min_{e \in P} c_f(e) > 0$ by definition of G_f .

We can push flow $c_f(P)$ along P as we've done before!

$$\hookrightarrow f'(x, y) = \begin{cases} f(x, y) + c_f(P) & \text{if } (x, y) \in P \\ -f(y, x) - c_f(P) & \text{if } (y, x) \in P \\ f(x, y) & \text{otherwise} \end{cases}$$

- skew-symmetric
- feasible

(check for all 3 cases
 $(x, y) \in P, (y, x) \in P, \text{ otherwise}$)

now $|f'| = \sum_x f'(s, x) = \sum_{x \neq u} f(s, x) + f'(s, u) = |f| + c_f(P) > |f|$

$(s, u) = \text{first edge of } P$

\hookrightarrow so (1) is false! \square

Thus, we can improve the flow of f iff there is an s-t path in G_f .

Ford-Fulkerson Algorithm

Start with $f = \text{all } 0 \text{ flow}$

$\rightarrow O(|E| \cdot \text{iterations})$ overall!

While \exists s-t path in G_f :

$P \leftarrow$ s-t path in G_f

augment f via $c_f(P)$ along P

$O(|E|)$ via
BFS / DFS

return f] f is a max flow!

Runtime: If all capacities $\in \{0, 1, \dots, C\}$ integers, then

each $c_f(e) \geq 1$ so at most $O(m F^*) \leq O(m n C)$
at most $n-1$ edges from s pseudo-polynomial

Flow Integrality Theorem

If all capacities $\in \mathbb{Z}$, then there is an integral max flow.

\rightarrow rational capacities still finite runtime \rightarrow real capacities may take forever

Two Variants:

1) Max Bottleneck Path Algorithm (MBP)

- find P in G_f from s to t that maximizes $c_f(P)$

↳ doable in $O(m \log n)$

- then augment along P .

$$\text{Runtime } O(m \log |f|) = O(m^2 \log n \log n)$$

2) Edmonds-Karp Algorithm

Find P with minimum # of edges. Augment along P .

Runtime $O(m^2n)$ → each edge part of at most $\frac{n}{2}$ augmenting paths → $O(mn)$
 → BFS to find shortest path in $O(m)$ augmentations

Strongly Poly Runtime → no dependence on edge weights
 (fast, even for real weights)
 ↳ real inputs

Edmonds Karp
 $O(m^2n)$

Best
 $O(mn)$

weakly Poly Runtime → polynomial dependence on $\log C$
 MBP
 $O(m^2 \log n \log n)$ Best ($\text{fast for } C \sim \text{poly}(n)$)
 $O(m^{1+\epsilon} \log C)$ ↳ integer inputs

pseudo Poly Runtime → polynomial dependence on C
 Ford-Fulkerson
 $O(m|f|) \leq O(mnC)$ (only fast for C small)

Linear Programming

$$\begin{array}{ll}
 \max / \min & c_1 x_1 + \cdots + c_n x_n \xleftarrow{\text{(linear)}} \text{objective function} \\
 \text{subject to} & A_{11} x_1 + A_{12} x_2 + \cdots + A_{1n} x_n \left\{ \leq, \geq, = \right\} b_1 \\
 & \vdots \\
 & A_{m1} x_1 + A_{m2} x_2 + \cdots + A_{mn} x_n \left\{ \leq, \geq, = \right\} b_m \\
 \text{mtm} \\
 \text{constraints} & \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{linear} \\
 & \quad \text{constraints} \\
 & x_1, x_2, \dots, x_n \geq 0 \xleftarrow{\text{non-negativity}} \\
 & \quad \underbrace{}_{n \text{ decision variables}} \text{constraints}
 \end{array}$$

Standard Form

can make $\min c_1x_1 + \dots + c_nx_n \rightarrow \max (-c_1)x_1 + \dots + (-c_n)x_n$

can make $A_{i,1}x_1 + \dots + A_{i,n}x_n \geq b_i \rightarrow (-A_{i,1})x_1 + \dots + (-A_{i,n})x_n \leq -b_i$

$$\text{can make } A_{i,1}x_1 + \dots + A_{i,n}x_n = b_i \rightarrow \begin{cases} A_{i,1}x_1 + \dots + A_{i,n}x_n \leq b_i \\ (-A_{i,1})x_1 + \dots + (-A_{i,n})x_n \leq (-b_i) \end{cases}$$

$$\begin{array}{l} \text{Can represent LP as} \\ \rightarrow \max C^T x \\ \text{subject to } Ax \leq b \\ \quad \quad \quad x \geq 0 \end{array}$$

$A = m \times n$
 $b = m \times 1$
 $c = n \times 1$
 $x = n \times 1$

feasibility - feasible region = $\{x \in \mathbb{R}^n \mid Ax \leq b \wedge x \geq 0\}$

infeasible if region = \emptyset , feasible otherwise

bounedness - feasible LP is either bouned or unbouned

feasible region is either bounded or un bounded

bounded feasible region \rightarrow bounded optimum \rightarrow unique optimal value

degeneracy - bounded LP can have multiple optimal solutions

Geometry of an LP

$n=1$: feasible region is $[c_1, \infty]$ or $[-\infty, c_1]$

① feasible region is single continuous interval

lower bound if $c_1 < 0$

0 if $c_1 = 0$

② optimal objective function is at corner point

upper bound if $c_1 > 0$

$n=2$:

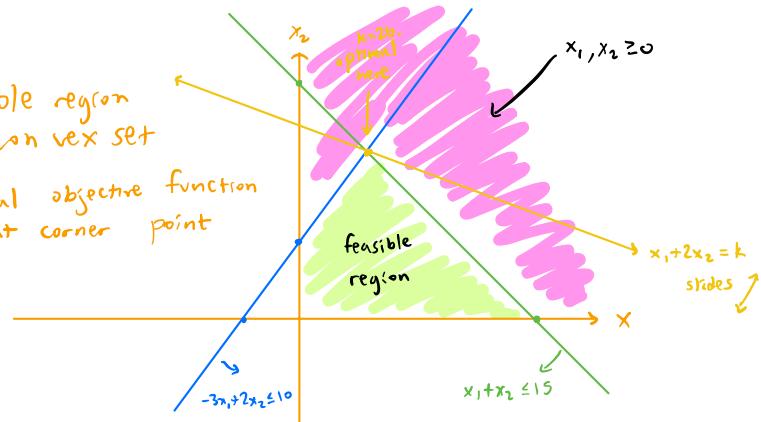
ex) $P_5: \max x_1 + 2x_2$

$$\begin{aligned} -3x_1 + 2x_2 &\leq 10 \\ x_1 + x_2 &\leq 15 \\ x_1, x_2 &\geq 0 \end{aligned}$$

① feasible region is convex set

② optimal objective function is at corner point

$n=3$: planes bound a space.



Solving Linear Programs → double in poly time.

Duality

Primal $P_5: \max x_1 + 2x_2$

$$\begin{aligned} -3x_1 + 2x_2 &\leq 10 \\ x_1 + x_2 &\leq 15 \\ x_1, x_2 &\geq 0 \end{aligned}$$

again. Try

$$x_1 + 2x_2 \leq 2(x_1 + x_2) \leq 30.$$

$$x_1 + 2x_2 \leq \frac{1}{6}x_1 + 2x_2$$

$$= \frac{1}{6}(-3x_1 + 2x_2) + \frac{5}{3}(x_1 + x_2)$$

$$\leq \frac{10}{6} + \frac{25}{3} = \frac{80}{3}$$

$$x_1 + 2x_2 \leq \frac{1}{5}(-3x_1 + 2x_2) + \frac{8}{5}(x_1 + x_2)$$

$$\leq 2 + 24 = 26 \text{ as desired.}$$

LECTURE 12 3/21/24

(Did not attend)

Primal

$$\max c^T x$$

$$\text{subject to } Ax \leq b$$

$$x \geq 0$$

Dual

$$\min b^T y$$

$$\text{subject to } A^T y \geq c$$

$$y \geq 0$$

Note dual of the dual = primal.

ex) $\max x_1 + 2x_2$

$$-3x_1 + 2x_2 \leq 10$$

$$x_1 + x_2 \leq 15$$

$$x_1, x_2 \geq 0$$

$$y_1(-3x_1 + 2x_2) \leq 10y_1$$

$$y_2(x_1 + x_2) \leq 15y_2$$

$$-3y_1 + y_2 \geq 1$$

$$2y_1 + y_2 \geq 2$$

$$\min 10y_1 + 15y_2$$

beyond the standard form...

| Primal | Dual |
|---|---|
| constraint ($\leq, \geq, =$) | decision variable ($\geq, \leq, \text{unrestricted}$) |
| decision variable ($\geq, \leq, \text{unrestricted}$) | constraint ($\geq, \leq, =$) |
| coefficients of objective | constraint bounds |
| constraint bounds | coefficients of objective |

Weak Duality

$$\begin{array}{ll} \max c^T x & \min b^T y \\ \text{subject to } Ax \leq b & \text{subject to } A^T y \geq c \\ x \geq 0 & y \geq 0 \end{array}$$

Let \vec{x}, \vec{y} be primal and dual feasible solutions.

Then

$$c^T \vec{x} \leq b^T \vec{y}$$

PF

$$c^T \vec{x} = \vec{x}^T c \leq \vec{x}^T A^T y \leq b^T y.$$

↳ primal unbounded \rightarrow dual infeasible

↳ dual unbounded \rightarrow primal infeasible

↳ primal bounded \rightarrow dual bounded OR infeasible

↳ dual bounded \rightarrow primal bounded OR infeasible

both primal AND dual
may be infeasible

Strong Duality

primal bounded \leftrightarrow dual bounded.

intuition - ball with coords (y_1, y_2, \dots) dropped into feasible region. moves to optimal solution where normal forces of walls (primal constraints) cancel gravity \rightarrow points normal to objective function

optimal dual solution \rightarrow primal (bordering) constraints maxed
 \rightarrow optimal primal solution

Note - not all walls will be used.

Suppose only some walls i will be used on wall

\hookrightarrow in dual solution this means $\underbrace{A^{(i)}^T}_{\text{column vector}} y = c_i$

\hookrightarrow all other constraints are slack $\Leftrightarrow x_i = 0$

"complementary slackness"

no normal force for non-touching walls

$$\text{so at optimality, } \underbrace{x_i}_{\text{o for no wall o for yes wall}} (\underbrace{A^{(i)}^T y - c_i}_{\text{for yes wall}}) = 0 \quad \forall i \in \{1, 2, \dots, n\}$$

P , NP , NP -Completeness

| poly-time | exp-time, faster unknown |
|---|--|
| unweighted shortest path in $O(E)$ | unweighted longest path in $O(1.66^{ V })$ undirected $O(2^{ V } \cdot E)$ directed |
| maximum bipartite matching in $O(E ^{1+o(1)})$ | maximum tripartite matching in $O(3.52^{ V })$ |
| linear programming in $O(n^{2.372})$ $n \times n$ matrix A | integer programming (integer vector) $n \times n$ integer matrix A → exp-time even if $A, b, c, x \in \{0,1\}$ |

| | | optimization | search | decision |
|------------------------------|---|--------------------------------|-----------------------------------|--|
| s-t shortest path | I | G, s, t | G, s, t, k | G, s, t, k |
| | O | min-weight s-t path | Simple path of weight $\leq k$ | YES if \exists path of weight $\leq k$ NO otherwise |
| s-t longest path | I | G, s, t | G, s, t, k | G, s, t, k |
| | O | max-weight s-t path | Simple path of weight $\geq k$ | YES if \exists path of weight $\geq k$ NO otherwise |
| max bipartite matching | I | G | G, k | G, k |
| | O | maximum pairing of vertices | pairing of $\geq 2k$ vertices | YES if \exists pairing of $\geq 2k$ vertices NO otherwise |

solving optimization → search → decision

↳ no efficient algorithm for decision
 \Downarrow

no efficient algorithm for optimization, search

A decision problem π is solvable in poly time if \exists algorithm A

and constant c s.t. \forall inputs x of size n,

$\hookrightarrow A(x)$ runs in $O(n^c)$

$\hookrightarrow \pi(x) = \text{YES}$ iff $A(x)$ returns 1.

$P = \boxed{\text{All decision problems solvable in poly time.}}$ (deterministic)

$NP = \boxed{\text{decision problems that are easy to verify in poly time.}}$

Formally, $\pi \in NP$ if \exists verification algorithm V_π and constants c, c' s.t.

$\hookrightarrow V_\pi$ takes two inputs x, y where $|y| \leq n^{c'}$

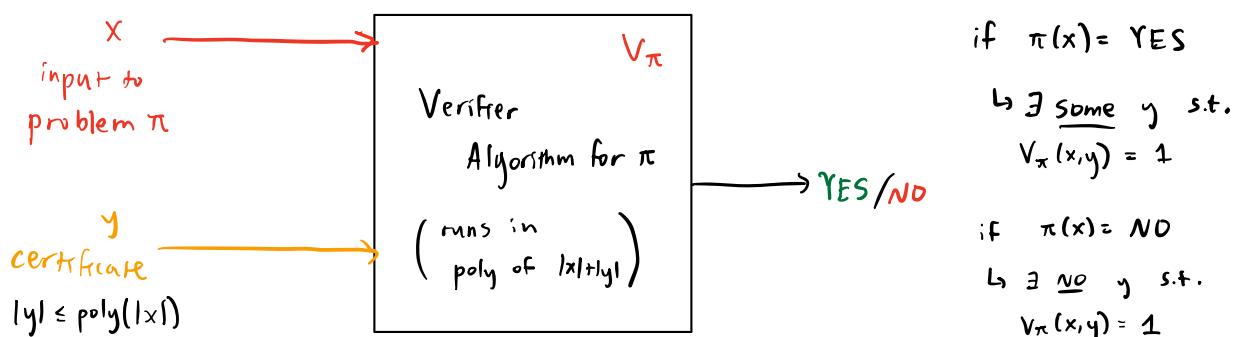
$\hookrightarrow V_\pi(x, y)$ takes $O((|x| + |y|)^c)$ time

\hookrightarrow for every input x of size n,

$\pi(x) = \text{YES}$ iff $\exists y$ of length $\leq n^{c'}$

s.t. $V_\pi(x, y) = 1$.

y is a certificate proving that x is a YES instance of π



y is usually just the decision problem itself

remember, only decision problems can be in P, NP, etc.

no-long path - given $G = (V, E)$, $s, t \in V, k$, is there no path from $s-t$ of length $\geq k$? $\in \text{co-NP}$ unknown if $\subseteq NP$

Thm $P \subseteq NP$.

Pf Take any $\pi \in P$.

Let A_π be the poly-time algorithm for π .

Then let V_π just run A_π and ignore y .

Runs in $\text{poly}(|x|) \leq \text{poly}(|x| + |y|)$.

□

Thm $NP \subseteq EXPTIME \leftarrow$ decision problems solvable in $2^{\text{poly}(n)}$

Pf Sketch Just run verifier on all strings y of length n^c
 $\rightarrow \exp(n)$

NP-Complete "hardest problems in NP"

① Must be in NP

② Must be NP-hard

↳ "poly-time algorithm for π can be used to solve any NP problem"

Let Q and π be decision problems.

A many-one polynomial time reduction from Q to π

is an algorithm R that takes input = instance x of Q ,

- runs in $\text{poly}(|x|)$ time

- returns instance y of π

s.t. $Q(x) = \text{YES} \rightarrow \pi(y) = \text{YES}$

$Q(x) = \text{NO} \rightarrow \pi(y) = \text{NO}$

→ y will have length $\text{poly}(|x|)$ since

R runs in $\text{poly}(|x|)$

$Q \leq_p \pi$

"many-one" because many instances x can go to one instance y

Claim If $Q \leq_p \pi$ and $\pi \in P$, $Q \in P$.

Pf Let A be algorithm for π , let R be the many-one poly time reduction of Q to π . Given instance x of Q ,

take $A(R(x))$ for π , which has same answer as x for Q .

Runtime: A, R both run in poly time of input $\rightarrow Q \in P$.

Also known

$P \subsetneq EXPTIME \subsetneq R$

↑
decidable
problems

LECTURE 14 4/14/24 (Did not attend)

Claim $A \leq_p B, B \leq_p C \rightarrow A \leq_p C$

PF

R_1 reduces A_{n_1} to B in $O(n_1^c)$

R_2 reduces B_{n_2} to C in $O(n_2^{c'})$

$A \xrightarrow[\text{YES}]{\quad} B \xrightarrow[\text{YES}]{\quad} C \xrightarrow[\text{YES}]{\quad}$

$\rightarrow R_2 \circ R_1$ reduces $A \rightsquigarrow C$
in $O((n_1^c)^{c'}) = O(n_1^{cc'})$
 $= \text{poly}(n_1)$ \square

So if \mathcal{H} is NP-complete

$\forall Q$ with $\mathcal{H} \leq_p Q$ and $Q \in \text{NP} \rightarrow Q$ is NP-complete.

$\forall Q \in \text{NP} \rightarrow Q \leq_p \mathcal{H}$

$A \leq_p B$ and $B \leq_p A \Leftrightarrow A, B$ equivalent.

\hookrightarrow All NP-complete problems are equivalent

Def CNF Formula on x_1, \dots, x_n (Boolean literals) is of the form
(conjunctive normal form)

$$\bigwedge_{i=1}^m \left(\bigvee_{j=1}^p l_{ij} \right)$$

every $l_{ij} = \neg x_k$ or x_k for $k \in [n]$.

Def k -SAT : $p=k$ for every clause in CNF formula
1SAT, 2SAT $\in \text{P}$ $k \geq 3$ gives NP-complete

Def A vertex cover of $G=(V,E)$ is a subset $S \subseteq V$ s.t.
 $\forall (u,v) \in E$, either u or v or both in S . $\text{size} = |S|$

VERTEX-COVER

Input : $G=(V,E)$, integer k

Output : does there exist a vertex cover of G of size $\leq k$?

VERTEX-COVER \in NP

\hookrightarrow V_{π} checks $y = S \subseteq V$ s.t. $|S| \leq k$
and every edge $\in E$ has at least one
endpoint in S \hookrightarrow poly-time

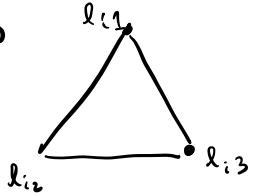
Thm $3\text{-SAT} \leq_p \text{VERTEX-COVER}$.

First construct problem from instance of 3-SAT.

Given $\Phi = (l_1 \vee l_{12} \vee l_{13}) \wedge (l_2 \vee l_{22} \vee l_{23}) \wedge \dots \wedge (l_m \vee l_{m2} \vee l_{m3})$

\uparrow
7 literals
 x_k or $\neg x_k$

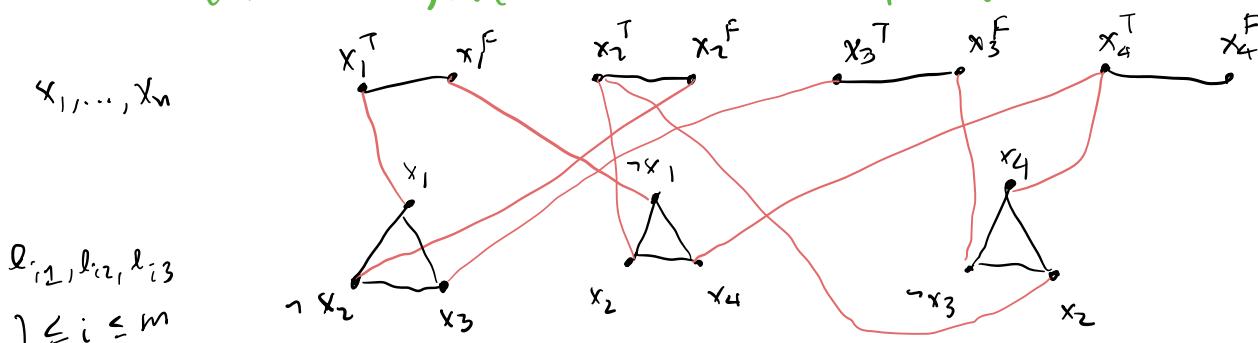
① for every $x_k \rightarrow x_k^T \xrightarrow{\quad} x_k^F$ ③ for each l_{ij} in clause i

② for every clause \rightarrow 

if $l_{ij} = x_p$ else $l_{ij} = \neg x_p$

$\boxed{k = n + 2m}$ 

ex) $(x_1 \vee \neg x_2 \vee x_3) \wedge (\neg x_1 \vee x_2 \vee \neg x_4) \wedge (x_4 \vee \neg x_3 \vee x_2)$



YES instance of 3SAT \rightarrow YES instance of VERTEX-COVER

- (a) \rightarrow put $x_k^T \in S$ for $x_k = 1$ else $x_k^F \in S$
- (b) \rightarrow pick one satisfied literal (from 2 or 3 if needed) per clause,
and put other two literals $\in S$.

obviously ①, ② type edges are covered by S , $|S| = n + 2m$.

③: if $x_k^T \xrightarrow{\quad} x_k$ if neither is covered, $x_k = 1$ by (b)
but then $x_k^T \in S$ by (a) \rightarrow contradiction

$x_k^F \xrightarrow{\quad} \neg x_k$ if neither is covered, $x_k = 0$ by (b)
but then $x_k^F \in S$ by (a) \rightarrow contradiction \square

YES instance of VERTEX-COVER \rightarrow YES instance of SSAT

read off $x_k = 1$ if x_k^T , else $x_k = 0$.

note we need $\exists n \in S$ in the ① type edges

note we need (disjoint) $\geq 2m \in S$ in the (2) type edges

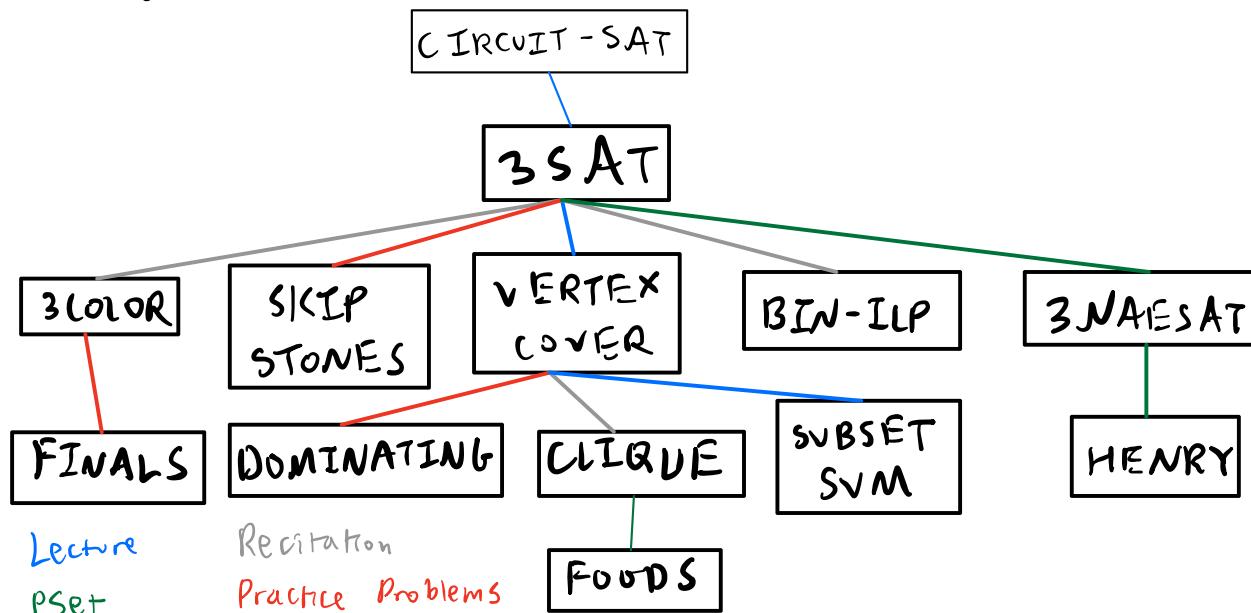
so $|S| = n+2m$ exactly.

↳ exactly one of x_k^T or x_k^F is covered.

for each clause

- each clause
- if none weren't satisfied, then entire 

If all three literals were, then, $\neg c = n+2m \rightarrow$ EXACTLY 2 in S . \square
 $\text{triangle } \in S \rightarrow$ contradiction since $|c| = n+2m$



LECTURE 15 4/9/24 (Did not attend)

Circuits: AND OR NOT



X3 can evaluate assignment in linear time via topological sort then Dynamic Programming.

CIRCUIT-SAT : given circuit C on n inputs x_1, \dots, x_n ,

YES iff \exists assignment of x_i 's to output 1 on C.

CIRCUIT-SAT ∈ NP.

Cook - Levin Theorem CIRCUIT-SAT is NP-complete.

WTS \exists many-one poly-time reduction from every NP problem to CIRCUIT-SAT.

↳ Use existence of verifier algorithm.

↳ make following assumptions on machine model

(1) program is stored in memory as sequence of instructions
 ↳ operations to perform ↓ addresses of operands → addresses to store results

(2) program counter maintains where we are in program
 ↳ branching for loops, etc.

(3) memory holds entire state of program

↳ program counter, program, storage, bookkeeping

(4) runtime $O(p(n))$ means we use $O(p(n) \log p(n))$ space
 ↑
 pointer size
 in RAM

(5) Main Assumption Let L_i be memory at step i of algorithm.
 From each $L_i \rightarrow L_{i+1}$, \exists Boolean circuit taking L_i to L_{i+1} ,
 constructable in $O(p(n))$ time

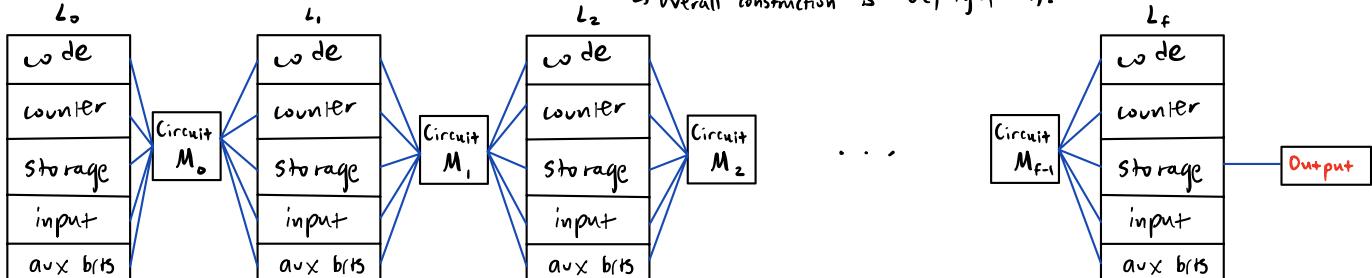
then with these assumptions, we get

Thm Suppose Q is a decision problem that can be solved by a program P in $p(n)$ time on inputs of size n . Then for fixed n \exists Boolean circuit C_n constructable in $O(\text{poly}(p(n)))$ time s.t. if inputs

$$x_1, \dots, x_n, \quad Q(x_1, \dots, x_n) = C_n(x_1, \dots, x_n).$$

PF Idea $f = p(n)$ steps, each M_i constructable in $O(\text{poly}(p(n)))$ time

↳ Overall construction is $O(\text{poly}(p(n)))$.



Proof of Cook-Levin

Let A be any problem NP w/ verifier V_A running in poly time.

Then \exists circuit C_N & input of size N computable in $\text{poly}(N)$ time

if

$V_A(x, y)$ can be represented as $C_N(x_1, \dots, x_n, y_1, \dots, y_{\text{open}})$

where $N = n + p(n)$, just hardcode x_i 's into C_N

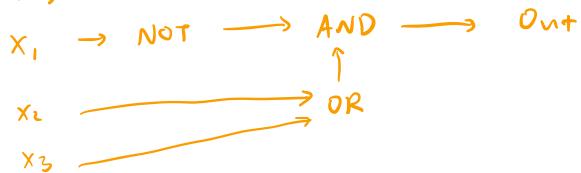
YES-instance of $A \Leftrightarrow \exists y$ for $V_A \Leftrightarrow$ satisfying instance of C_N \square

Now we need $\text{CIRCUIT-SAT} \leq_p \text{3SAT}$.

pf Encode each gate using a variable

(reduction) gate i : $g_i \Leftrightarrow \neg x$ $g_i \Leftrightarrow x \wedge y$ $g_i \Leftrightarrow x \vee y$

ex)



$$(g_1 \Leftrightarrow \neg x_1) \wedge (g_2 \Leftrightarrow x_2 \vee x_3)$$

$$(g_3 \Leftrightarrow g_1 \wedge g_2) \wedge g_3$$

Final out put
should be true.

$$g_i \Leftrightarrow \neg x$$

$$g_i \Leftrightarrow x \wedge y$$

$$g_i \Leftrightarrow x \vee y$$



$$(g_i \vee x) \wedge (\neg g_i \vee \neg x)$$



$$(\neg g_i \vee a \vee b) \wedge (g_i \vee \neg a) \wedge (g_i \vee \neg b)$$

$$(g_i \vee \neg a \vee \neg b) \wedge (\neg g_i \vee a) \wedge (\neg g_i \vee b)$$

can be used to create
exact 3SAT instance!

\square

P - class of minimization/maximization problems

$P \in P$ - one instance

x = solution to some P , $Q_p(x)$ outputs quality of x

$\text{OPT}(P)$ = quality of optimal solution to P

For any algorithm A , $A(P)$ = quality of A 's solution to P

Def P - class of minimization problems

$\alpha \geq 1$ - approximation factor

We say A is an α -approximation for P if and only if $\forall P \in P$

$$[\text{OPT}(P) \leq] A(P) \leq \alpha \cdot \text{OPT}(P)$$

by def

$$\max: \alpha \geq 1 \quad [\text{OPT}(P) \geq] A(P) \geq \frac{1}{\alpha} \cdot \text{OPT}(P)$$

SMTI

Maximum Matching - max nonoverlapping pairs of vertices in G

↙ ↓

max flow for bipartite solvable in poly time in general

connected by edges

greedy algorithm - take edges and add to matching if possible.

Claim 2-approximation. \rightarrow tightness (upper bound is optimal) @

PF If M is output of our algorithm, M^* is optimal, then $|M^*| \leq 2|M|$.

For $e \in M, e^* \in M^*$, add label " e^* " to e if e "blocks" e^* i.e. they share a vertex

Then each edge $e \in M$ can have at most two labels. \Rightarrow

But also every edge $e^* \in M^*$ must be blocked by

some $e \in M$ (or just add e^* to M).



$$2|M| \geq \sum_{e \in M} \# \text{labels} \geq |M^*|$$

□

Vertex Cover

S
↓

greedy algorithm - find maximal matching, take all endpoints.

Claim 2-approximation \rightsquigarrow tightness @ 

S is a vertex cover. If \exists some edge e not covered by S, M should have found it.

$|M| \leq |S^*|$ because just consider edges in M $\rightarrow S^*$ must contain one vertex for each $e \in M$.

$$\hookrightarrow |S| = 2|M| \leq 2|S^*|$$

Note: we can't get much better than 2, currently $2 - O\left(\frac{1}{\sqrt{\log|V|}}\right)$
 $\sqrt{2-\epsilon}$ is best if P \neq NP, but no known algorithm.

Weighted Vertex Cover

$G = (V, E, w) \rightarrow$ we want minimum weighted vertex cover $\sum_{v \in S} w(v)$
 $w: V \xrightarrow{T} \mathbb{R}_{\geq 0} \rightarrow$ LP rounding...

LP Rounding

integer program \rightarrow relax to LP \rightarrow solve LP \rightarrow round solution.

Let x_v for $v \in V$ be indicator variables for keeping/excluding v in S

$$\text{want } \min \sum_{v \in V} x_v \quad \text{where} \quad \begin{aligned} x_a + x_b &\geq 1 \quad \forall \{a, b\} \in E \\ x_v &\in \{0, 1\} \quad \forall v \in V \\ \hookrightarrow 0 \leq x_v &\leq 1 \quad \forall v \in V \end{aligned}$$

solve to get $x_v^* \quad \forall v \in V$.

$$x_v := \begin{cases} 0 & x_v^* < \frac{1}{2} \\ 1 & x_v^* \geq \frac{1}{2} \end{cases}$$

Claim S is a cover.

Pf Each $\{u, v\} \in E$ has $x_u^* + x_v^* \geq 1 \rightarrow$ at least one is $\geq \frac{1}{2}$

$$\hookrightarrow \max(x_u, x_v) = 1 \quad \square$$

Now we have

$$\sum_{v \in S} w(v) = \sum_{v \in S} w(v)x_v \leq 2 \sum_{v \in V} w(v)x_v^* \leq 2 \sum_{v \in S^*} w(v) \quad \text{2-approximation?}$$

Max k-SAT Given k -CNF formula ϕ , find an assignment satisfying

↑ the max # of clauses. (Assume literals in each clause distinct)

decision is NP-complete

optimization is harder

Algorithm Independently set each x_i T or F

↳ each clause has $1 - 2^{-k}$ chance of being satisfied

↳ approximation ratio $\alpha = \frac{1}{1 - 2^{-k}}$.

Now derandomize.

$$\mathbb{E}_{x_1, \dots, x_n} [\#SAT(\phi)] = \frac{1}{2} \mathbb{E}_{x_1, \dots, x_n} [\#SAT(\phi) | x_1=T] \quad \leftarrow \text{assign } x_1 \text{ the one with better conditional expectation}$$
$$+ \frac{1}{2} \mathbb{E}_{x_1, \dots, x_n} [\#SAT(\phi) | x_1=F]$$

↳ deterministic $\alpha = \frac{1}{1 - 2^{-k}}$ -approximation.

KNAPSACK n items, (v_i, w_i) need to max $\sum v_i$ with $\sum w_i \leq W$.

Greedy Algo 1

Best $\frac{v_i}{w_i}$ first.

But if we have $K=2C$, items $(2c, 2c), (2, 1)$ our approximation ratio is $\frac{2c}{1} \rightarrow 2C$.

Greedy Algo 2

Take one item with best v_i .

But if we have $W = \frac{C(C+1)}{2}$ and items $(1, 1), (2, 1), \dots, (C, 1)$ then $\alpha = \frac{\frac{C(C+1)}{2}}{C} = \frac{C+1}{2}$.

LP Relaxation

→ Fractional Algorithm

$$\max \sum_{i=1}^n v_i x_i$$

$$\text{s.t. } \sum_{i=1}^n w_i x_i \leq W$$

$$x_i \in \{0, 1\} \quad \forall i \in [n]$$

$$\Downarrow \\ 0 \leq x_i \leq 1 \quad \text{relaxation}$$

| | | | | | |
|-------|-------|-------|-------|-------|-------|
| w_1 | w_2 | w_3 | w_4 | w_5 | w_6 |
|-------|-------|-------|-------|-------|-------|

decreasing order of $\frac{v_i}{w_i}$

↑ take

largest fraction of first non-fitting item and add that to x_i

Final 2-approx Algo

take $\max(\text{Greedy Algo 1}, \text{Greedy Algo 2})$

$$\underline{\text{PF}} \quad \begin{matrix} \text{relaxed} \\ \text{optimal} \end{matrix} \quad \text{LP} \quad \gtrless \quad \begin{matrix} \text{integer} \\ \text{optimal} \end{matrix} \quad \text{LP}$$

but Greedy Algo 1 + Greedy Algo 2 \geq relaxed LP \geq integer LP = OPT
 \Rightarrow contains all non-fractional items compared to relaxed LP \Rightarrow value > fractional item in relaxed LP \therefore our algorithm gets value $= \frac{OPT}{2}$

LECTURE 17 4/16/24 11AM

Online Learning (Part 2)

- LEARNER, set of actions $a \in A$ ← Known to LEARNER & ADVERSARY

→ at each time step t

 - ① pick some probability distribution p_t on A
 - ② pick action according to p_t
 - ③ feedback $c_t(a)$, c_t is cost / loss function

↑

known to adversary

known to learner after a picked

At time step t

- take weighted majority of n experts' predictions

for each expert i . $w_i = 1$ $\forall i$.

$$w_i^{t+1} = (1-\varepsilon) w_i^t$$

$$W^t := \sum_{i=1}^n w_i^t$$

$$W^{t+1} = \sum_{i=1}^n w_i^{t+1} = W^t - \varepsilon \sum_{\substack{i: \text{expert} \\ \text{incorrect}}} w_i^t \leq (1 - \frac{\varepsilon}{2}) W^t$$

↑
 ONLY
 IF LEARNER wrong

Let i^* = best expert, makes m mistakes over T steps

$$W_{\varepsilon}^T = (-\varepsilon)^{m^*}$$

now if our method makes m mistakes over T steps

$$W^T \leq \left(1 - \frac{\epsilon}{2}\right)^m n$$

1-18 have

$$W^T \leq W^T \Rightarrow (1-\varepsilon)^m \leq (1-\frac{\varepsilon}{2})^m$$

$$m^* \ln(1-\varepsilon) \leq m \ln\left(1 - \frac{\varepsilon}{2}\right) + \ln n$$

VI 1A for $0 \leq \varepsilon \leq \frac{1}{2}$

$$m^* (-\varepsilon - \varepsilon^2) \leq -m \frac{\varepsilon}{2} + \ln n$$

$$\text{external regret} \rightarrow m - m^* \leq (1+2\varepsilon)m^* + \frac{2\ln n}{\varepsilon}$$

but $\lim_{T \rightarrow \infty} \frac{1}{T} \left((1+2\epsilon)^m + \frac{2\ln n}{\epsilon} \right) \geq 1$ is possible,

so no vanishing regret

Randomized Multiplicative Weights

At time step t

- pick expert i based on distribution $(w_1^t, w_2^t, \dots, w_n^t)$

for each expert i
correct

$$w_i^{t+1} = w_i^t$$

incorrect

$$w_i^{t+1} = (1-\varepsilon)w_i^t$$

$$W^{t+1} = \sum_{i=1}^n w_i^{t+1} = \sum_{i=1}^n (1-\varepsilon M_i^t) w_i^t = W^t - \varepsilon \mathbb{E}[M^t] W^t \leq e^{-\varepsilon \mathbb{E}[M^t]} W^t$$

$$\text{then } W^T \leq e^{-\varepsilon \mathbb{E}[M^T]} W^0 \Rightarrow (1-\varepsilon)^{m^*} \leq e^{-\varepsilon \mathbb{E}[M^T]}$$

$$w_i^T = (1-\varepsilon)^{m^*}$$

$$\mathbb{E}[M] \leq m^*(1+\varepsilon) + \frac{\ln n}{\varepsilon}$$

$$\hookrightarrow \mathbb{E}[M] - m^* = m^*\varepsilon + \frac{\ln n}{\varepsilon}$$

$$\lim_{T \rightarrow \infty} \frac{1}{T} \left(m^*\varepsilon + \frac{\ln n}{\varepsilon} \right) \leq \varepsilon$$

can control
to get
vanishing regret!!

Define M_i^t = indicator for expert i
making mistake at time t

M^t = indicator for LEARNER
making mistake at time t

$$M = \sum_{t=1}^T M^t = \# \text{ of mistakes of LEARNER by time } T$$

Multiplicative Weights Update

Same as RMW but now $M_i^t \in [-1, 1]$ continuous cost

M = total cost of LEARNER

At time step t

- pick expert i based on distribution

$$(w_1^t, w_2^t, \dots, w_n^t)$$

for each expert i

$$w_i^{t+1} = (1-\varepsilon M_i^t) w_i^t$$

$$\mathbb{E}[M] \leq m_i^* (1+\varepsilon) + \frac{\ln n}{\varepsilon}$$

$$\forall i \in \{1, 2, \dots, n\}$$

$$\text{where } m_i^* = \sum_{t=1}^T M_i^t$$

= total cost of i^{th} expert

LECTURE 18 4/23/24 (Did not attend)

Casino : start w/ \$20, $\pm \$5$ each time w/ $\frac{1}{2}$ chance

Def A Stochastic process is a sequence of random variables $\{X_t : t \in \mathbb{N}\}$
indexed by time

Process is a Markov Chain / has Markov Property if

X_{t+1} only depends on $X_t \quad \forall t \in \mathbb{N}$

* Let D be our state space

ex) $D = \{\$0, \$5, \dots, \$40\}$

$X_t \in D$ = amount of \$ you have after t rounds of play

Graphical Representation

$\forall t \in N$, $G_t = (D, E_t, w_t)$ with $E_t = \{(u, v) \in D^2 \mid \Pr[X_{t+1} = v \mid X_t = u] > 0\}$
 $w_t(u, v) = \Pr[X_{t+1} = v \mid X_t = u] \quad \forall u, v \in D$

time-homogeneous if G_t independent of t .

Walk Matrix Representation

$$W \in \mathbb{R}_{\geq 0}^{D \times D} \quad w_{uv} = \Pr[X_{t+1} = v \mid X_t = u]$$

need W + initial distribution

$$\vec{x}^{(0)} \in [0, 1]^{101} \text{ (row vector)} \\ \text{s.t. } \sum_{i=1}^{101} x_i^{(0)} = 1$$

$$\vec{x}^{(1)} = \vec{x}^{(0)} W, \quad \vec{x}^{(t)} = \vec{x}^{(0)} W^t$$

stationary distribution - $\vec{x} W = \vec{x}$ (usually can just solve!)

DEFINITIONS

strongly connected = irreducible
SCC = communicating class



$$\text{period}(v) = \text{gcd}(\text{len}(\text{all directed cycles in its SCC}))$$

v is periodic if $\text{period}(v) > 1$, else aperiodic

an SCC is periodic if its vertices are periodic, else aperiodic

a graph is periodic if some SCC is periodic, else aperiodic

in undirected graphs, note that v is periodic iff its SCC is bipartite

an SCC is periodic if it is bipartite, else aperiodic

a graph is periodic if some SCC is bipartite, else aperiodic

Given a weighted directed graph we can also build a Markov chain via random walk

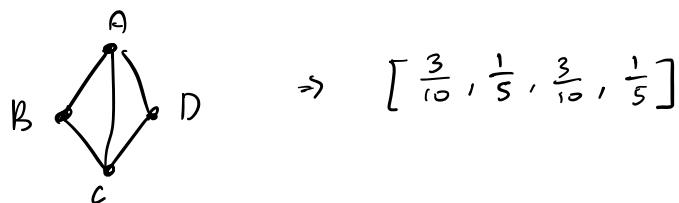
$$w_{uv} = \frac{c(u, v)}{\sum_{x \in D} c(u, x)}$$

now sums to 1

Fundamental Theorem of Markov Chains

Let G_x be the Markov chain for a time-homogeneous Markov process.
 There is always a stationary distribution for G_x (for finite D).
 → if G_x is strongly connected, it is unique
 → if G_x is aperiodic, it converges to some stationary distribution
 intuition: only consider recurrent (no outgoing edges) SCCs since transient (yes outgoing edges) SCCs can lose mass.

Remark If G_x is undirected/unweighted before normalization, stationary distribution is proportional to degree of each vertex



$$\left[\frac{3}{10}, \frac{1}{5}, \frac{3}{10}, \frac{1}{5} \right]$$

α -lazy random walk → do nothing with probability α , normal with probability $1-\alpha$.
 ↳ creates aperiodicity.

LECTURE 19 4/25/24 (Did not attend)

Sampling Problem How do we draw samples from a probability distribution π ?
 e.g. if we want to estimate expectation

example card shuffling - riffle takes ~ 7 shuffles for uniform.

Idea Design walk matrix W s.t. π is its stationary distribution.

Alg: Sample/pick X_0 arbitrary. then
 Then for $t=1, \dots, T$ $X_T \sim \vec{x}^{(0)} W^T \approx \pi$
 sample $X_t \sim W(X_{t-1}, \cdot)$. for big T .

π is usually given to us as a bunch of weights $w(x)$ for $x \in D$ s.t. $\pi(x) = \frac{w(x)}{\sum_{y \in D} w(y)}$, but we can't add up since $|D|$ is huge.

Def W is reversible if it satisfies detailed balance equations

$$\exists \pi \text{ s.t. } \pi(u) W_{uv} = \pi(v) W_{vu} \quad \forall u, v \in D$$

i.e. for $\{X_t : t \in \mathbb{N}\}$, assuming $X_0 \sim \pi$, (X_s, \dots, X_T) has same distribution as (X_T, \dots, X_s)

* For a symmetric walk matrix, we can have a uniform stationary distribution.

$$\pi(y) = \pi(y) \sum_{x \in D} W_{yx} = \sum_{x \in D} \pi(y) W_{xy} = \sum_{x \in D} \pi(x) W_{xy} \text{ as desired.}$$

Motivating Problem - Sampling a K -coloring of G efficiently given that $K >$ maximum degree + 1 (extra flexibility)
AND given a valid K -coloring

Idea → pick a random vertex, change its color to one of the K colors
so that new coloring is still valid.

→ repeat T times

↳ Markov process, states are valid colorings.

① Self-loops since we can pick same color

② Strongly-connected.

Pf take two colorings \vec{c}, \vec{c}' .

First vertex: if $c_i = c'_i$, ok.

Else, set v to an unused color (exists since $K = \max + 1$)

then color all neighbors of v in \vec{c}' that had

color c'_i to a different color

Then set v to c'_i . Induct.

↳ we never need to worry about changing earlier

vertices because \vec{c}' is valid, i.e. desired coloring

has no same-color edges. □

③ Uniform distribution is stationary

Pf Note $W_{\vec{c}, \vec{c}} = W_{\vec{c}', \vec{c}'} \quad \forall \vec{c}, \vec{c}' \in D \rightarrow$ symmetric □

⇒ our idea, if run long enough, gives us a uniformly random coloring!

Metropolis-Hastings

estimate $v: D \rightarrow \mathbb{R}$ given $\bar{v} = k \cdot v$ for large (nonenumerable) K

idea pick a known Markov chain $g: D^2 \rightarrow \mathbb{R}$

sample based on g but accept/reject using info in \bar{v}

$$\text{ex)} \quad v(\vec{c}) = \begin{cases} \frac{1}{N} & \text{if } \vec{c} \text{ valid} \\ 0 & \text{otherwise} \end{cases} \quad \bar{v}(\vec{c}) = \begin{cases} 1 & \text{if } \vec{c} \text{ valid} \\ 0 & \text{otherwise} \end{cases}$$

$k=N = \# \text{valid colorings}$

$g(\vec{c}', \vec{c})$ is nonzero iff \vec{c}' differs from \vec{c} by one vertex or fewer in color.

Formally we design a walk matrix W

$$W_{xy} = g(y|x) \cdot p_{\text{acc}}(y, x) \quad p_{\text{acc}}(y, x) = \min \left\{ 1, \frac{\bar{v}(y)}{\bar{v}(x)} \cdot \frac{g(x|y)}{g(y|x)} \right\}$$

\uparrow \uparrow

$\Pr[X_{t+1}=y | X_t=x]$ $\frac{v(y)}{v(x)}$

$$W_{xx} = 1 - \sum_{y \neq x} W_{xy} = \Pr[\text{no move/reject}]$$

\rightarrow strongly connected and aperiodic if

$$g: D^2 \rightarrow \mathbb{R}^{>0} \text{ and } \bar{v}: D \rightarrow \mathbb{R}^{>0}$$

Claim W has stationary distribution v .

Pf We prove W is reversible: $v(x) W_{xy} = v(y) W_{yx}$.

So we need to show

$$v(x) g(y|x) \min \left\{ 1, \frac{v(y)}{v(x)} \cdot \frac{g(x|y)}{g(y|x)} \right\} = v(y) g(x|y) \min \left\{ 1, \frac{v(x)}{v(y)} \cdot \frac{g(y|x)}{g(x|y)} \right\}$$

casework on \uparrow vs. \uparrow shows this indeed. \square

Rate of Convergence

$\{X_t : t \in \mathbb{N}\}$ $X_0 \sim \tilde{x}^{(0)}$. want to understand how fast
 $X_t \sim \tilde{x}^{(0)} w^{t \rightarrow \pi}$

Now imagine another Markov process $\{Y_t : t \in \mathbb{N}\}$ s.t. $Y_0 \sim \pi$

and $\Pr[Y_{t+1}=v | Y_t=u] = W_{uv} \forall u, v, t.$

Coupling: Correlate X_t and Y_t s.t. $\mathbb{E}[T]$ is minimized.

Here, $T = \min_{\tau} \{t : X_t = Y_t\}$ (random walk collides with stationary distribution)

↳ after $O(\mathbb{E}[T])$ steps the distributions are "close"

LECTURE 20 4/30/24 (Did not attend)

Multiplying polynomials \rightarrow

$$p(x) = a_0 + \dots + a_{n-1}x^{n-1}$$

$$q(x) = b_0 + \dots + b_{n-1}x^{n-1}$$

takes $\Theta(n^2)$ to do naively.

Assume we can multiply real/complex numbers precisely in $O(1)$ time.

↳ More on this later.

coefficient representation

we can also describe polynomials of degree $n-1$ as
 a set of n values at n inputs (then interpolate)

then $p(x_i)q(x_i)$ for $0 \leq i \leq n-1$ is easy to compute.

→ but is this useful? how to convert \Rightarrow coefficient representation?

if we carefully choose x_i , turns out that we can
 do this all in $O(n \log n)$.

Fast Fourier Transform

- can take $\Theta(n)$ to evaluate degree- n polynomial at x_i , but we can choose related x_i to do less work

observe that (for n even)

$$p(x) = a_0 + a_2 x^2 + \dots + a_{n-2} x^{n-2} + a_1 x + a_3 x^3 + \dots + a_{n-1} x^{n-1}$$

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

$$\hookrightarrow p(-x) = \underbrace{p_{\text{even}}(x^2)}_{\text{even}} - x \underbrace{p_{\text{odd}}(x^2)}_{\text{even}}$$

- define size of degree- $(d-1)$ polynomial as $d = \# \text{coefficients}$
- let $S^k = \{x^k \mid x \in S\}$ for sets S .

idea

- compute $\{p_{\text{even}}(x^2) \mid x \in S\} \Leftrightarrow \{p_{\text{even}}(y) \mid y \in S^2\} \quad T(\frac{n}{2})$
- compute $\{p_{\text{odd}}(x^2) \mid x \in S\} \Leftrightarrow \{p_{\text{odd}}(y) \mid y \in S^2\} \quad T(\frac{n}{2})$
- combine $p(x) = p_{\text{even}}(x^2) + x p_{\text{odd}}(x^2)$ for $x \in S \quad O(n)$

$$T(\text{size } n \text{ poly @ } n \text{ pts}) = 2 T(\text{size } \frac{n}{2} \text{ poly @ } \frac{n}{2} \text{ pts}) + O(n) \rightarrow O(n \log n).$$

for this we need $|S^{\frac{n}{2}}| = \frac{|S^n|}{2} \nearrow$ to reduce computations.

\hookrightarrow use roots of unity! (first part n s.t. $n = 2^k$)

then $S = \{w^i \mid w^n = 1\}$ works

so we can now apply...

Discrete Fourier Transform (coefficient \rightarrow point-value)

start with coefficients $(a_0, \dots, a_{n-1}) \rightarrow$ values at $(1, \omega, \omega^2, \dots, \omega^{n-1})$

$$y_k = \sum_{i=0}^{n-1} a_i \omega^{ik}$$

$$\hookrightarrow (y_0, y_1, \dots, y_{n-1})$$

$\text{FFT}((a_0, \dots, a_{n-1}), n) :=$ $n=2^k$ But how do we go
if $n=1 \rightarrow$ return a_0 backward, i.e. from
else:

values $\rightarrow y_{\text{even}} \leftarrow \text{FFT}((a_0, a_2, \dots, a_{n-2}), \frac{n}{2})$
at $1, \omega, \dots, \omega^{n-2} \rightarrow y_{\text{odd}} \leftarrow \text{FFT}((a_1, a_3, \dots, a_{n-1}), \frac{n}{2})$
for $j \in \{0, 1, \dots, \frac{n}{2}-1\}$:
 $y[j] \leftarrow y_{\text{even}}[j] + \omega^j y_{\text{odd}}[j]$
 $y[j + \frac{n}{2}] \leftarrow y_{\text{even}}[j] - \omega^j y_{\text{odd}}[j]$
return y

$(y_0, \dots, y_{n-1}) \rightarrow (a_0, a_1, \dots, a_{n-1})?$

- $O(n^3) \sim O(n^{2.373})$ by solving system of equations on $p(a_i) = y_i$
- $O(n^3) \sim O(n^2)$ using Lagrange Interpolation.

L E C T U R E 21 5/2/24 (Did not attend)

Note that

$$\begin{pmatrix} y_0 \\ y_1 \\ y_2 \\ \vdots \\ y_{n-1} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & \omega & \omega^2 & \cdots & \omega^{n-1} \\ 1 & \omega^2 & \omega^4 & \cdots & \omega^{2(n-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \omega^{n-1} & \omega^{2(n-1)} & \cdots & \omega^{(n-1)^2} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_{n-1} \end{pmatrix}$$

\vec{y} W \vec{a}

And that $\frac{1}{\sqrt{n}} W$ is unitary. So $\vec{a} = \frac{1}{\sqrt{n}} W^T \vec{y} = \bar{W} \frac{1}{\sqrt{n}} \vec{y} = \overline{W \frac{\vec{y}}{\sqrt{n}}}$.
 ↳ apply FFT on $\frac{\vec{y}}{\sqrt{n}}$ then take conjugate.

Convolution: ignore polynomials, just think of sequences (a_0, a_1, \dots, a_m)
 (b_0, b_1, \dots, b_n)

$$(c_0, \dots, c_{m+n}) := c_k = \sum_{i+j=k} a_i b_j$$

↳ solvable in $O((m+n) \log(m+n))$.

Minkowski Sum given $X, Y \subseteq \{0, 1, \dots, m-1\}$, compute $X+Y = \{x+y \mid x \in X, y \in Y\}$.

↳ naively takes $O(|X||Y|)$. → ok if $m \gg n$.

↳ let (x_0, \dots, x_{m-1}) s.t. $x_i = 1$ if $i \in X$, else 0. same for (y_0, \dots, y_{m-1}) .

→ convolve to get (c_0, \dots, c_{2m-2}) , then $i \in X+Y$ iff $c_i > 0$. $\Rightarrow O(m \log m)$

String Matching in $\{0,1\}$ for now. Let $m \in \mathbb{N}$.

search string $S = (s_0, s_1, \dots, s_{n-1})$, target string $P = (p_0, \dots, p_{m-1})$

want to find $M := \{k \mid s_{k+i} = p_i \ \forall i \in [m]\}$

Naively takes $O(mn)$ → quadratic if $m = \frac{n}{2}$

→ construct $P^{\text{rev}} = (p'_{m-1}, p'_{m-2}, \dots, p'_1, p'_0)$ where

convolve $P^{\text{rev}} * S'$ to get $L = (l_0, \dots, l_{m+n-2})$

↳ $l_{i+m-1} = m$ iff $i \in M$

$$\begin{aligned} s'_i &= \begin{cases} 1 & s_i = 1 \\ -1 & s_i = 0 \end{cases} \\ p'_i &= \begin{cases} 1 & p_i = 1 \\ -1 & p_i = 0 \end{cases} \end{aligned}$$

Proof

$$l_{i+m-1} = \sum_{j+k=i+m-1} p'_{m-1-k} s'_j = \sum_{j-i=x} p'_x s'_j = \sum_{j-i=x} \begin{cases} 1 & p_x = s_j \\ -1 & p_x \neq s_j \end{cases} = m - 2(\text{mismatches})$$

let $m-k-1 = x$

can also easily check how many substrings have $\geq m-k$ characters in common: $l_{i+m-1} \geq m-2k$

LECTURE 22 5/7/24 (did not attend)

Sublinear Algorithms - read a fraction of input, hope to get an approximately correct answer

① classical approximation

→ answer within δ of correct one

→ optimization, not useful for decision problems.

② property testing

YES if it is or is close to a YES-instance

NO if it is far from a YES-instance (and sometimes if close)



③ diameter of point set - given $N = n \times n$ distance matrix D

$$\hookrightarrow D_{ij} = D_{ji} \quad \hookrightarrow D_{rj} \leq D_{rk} + D_{kj} \quad \forall i, j, k,$$

estimate $D^* = \max_{i,j} D_{ij}$

Algorithm
→ pick 1 point p , return $\max_i D_{ip} \rightarrow O(\sqrt{n})$

if D_{ip} is best, then

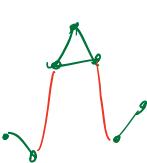
$$2 \max(D_{ip}, D_{jp}) \geq D_{ip} + D_{jp} \geq D_{pj}$$

can't do better,
e.g. $D = \text{all } 1's \text{ vs. } D' = \text{all } 1's \text{ except } D_{ii} = D_{jj} = 2 \text{ for some } i, j$
need $\Omega(n^2)$ to find i, j

(2a) graph connectedness (use adjacency list representation)

$\epsilon > 0$ small number.

Def Graph G with n vertices, max degree d is ϵ -close to connected if adding $< \epsilon \cdot d \cdot n$ edges will connect it

ex) 
 $n = 7$
 $d = 2$
 $\epsilon < 7$ gives ϵ -close to connected.

if $\epsilon \leq \frac{1}{dn}$, $\epsilon dn \leq 1$ so G is already connected! So let $\epsilon > \frac{1}{dn}$.

Algorithm (d, n, G, ϵ given)

Repeat $\frac{c}{\epsilon d}$ times:

pick a random point, DFS from it.

if $< \frac{2}{\epsilon d}$ vertices in connected component, output NO.

if $\geq \frac{2}{\epsilon d}$ vertices, continue.

$$\epsilon > \frac{1}{dn} \Rightarrow \frac{1}{\epsilon d} \ll n$$

output YES.

Runtime is $\frac{c}{\epsilon d} \cdot \frac{2}{\epsilon d} \cdot d = O\left(\frac{1}{\epsilon^2 d}\right) = O\left(\frac{1}{\epsilon}\right) \cdot O\left(\frac{1}{\epsilon d}\right) \Rightarrow$ sublinear.

↑
constant

→ Obviously YES if connected.

→ if NOT ϵ -close to connected, at least $\epsilon \cdot d \cdot n$ edges needed to connect G

↪ $\geq \epsilon \cdot d \cdot n$ CC's.

↪ $\geq \frac{\epsilon \cdot d \cdot n}{2}$ CC's with $< \frac{2}{\epsilon d}$ vertices (else too many vertices)

↪ $\geq \frac{\epsilon \cdot d \cdot n}{2}$ vertices in CC's of $< \frac{2}{\epsilon d}$ vertices ← each CC has ≥ 1 vertex

$$\begin{aligned} \text{so } \Pr[\text{No}] &= 1 - (\text{never find "small" CC}) \geq 1 - \left(1 - \frac{\epsilon d}{2}\right)^{\frac{c}{\epsilon d}} \\ &\geq 1 - \left(\left(1 - \frac{1}{2/\epsilon d}\right)^{\frac{2}{\epsilon d}}\right)^{\frac{c}{2}} \\ &\geq 1 - e^{-\frac{c}{2}} \\ &\geq \frac{3}{4} \end{aligned}$$

for large enough c .

(if ϵ -close to connected but not connected, we don't care)

(2b) List sortedness

Def A list of length n is ϵ -close to sorted if $\leq \epsilon n$ elements can be removed to sort the list

Idea 1

Random x_i, x_{i+1} check if sorted.

↳ NO, ex) 1, 3, 5, 7, 9, ..., 2, 4, 6, ...

only one pair tells us NO.

$\Pr \approx \frac{1}{n}$ of discovery

Idea 2

Random $i < j$, check if $x_i < x_j$.

↳ NO ex) 2 1 4 3 6 5 8 7 ...

again only $\frac{n}{2}$ pairs tell us NO.

$\Pr \approx \frac{1}{n}$ of discovery

Algorithm

Repeat $\frac{c}{\epsilon}$ times.

pick random i

binary search for x_i in L as if it was sorted

if you find inconsistency or don't end up in location i , output NO

output YES.

Runtime is $\frac{c}{\epsilon} \cdot O(\log n) = O\left(\frac{\log n}{\epsilon}\right) \ll n$ for large enough n .

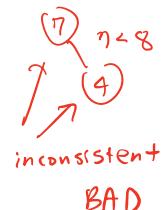
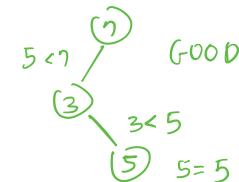
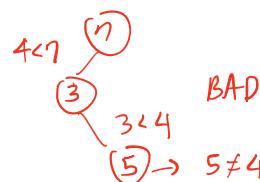
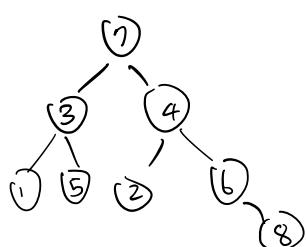
bad index examples

ex) $L = 1 3 5 7 2 4 6 8$

1) $i=5, x_i=4$

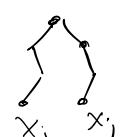
2) $i=2, x_i=5$

3) $i=7, x_i=8$



Claim If $i < j$ and i, j are both GOOD, then $x_i < x_j$.

Pf Consider just the paths from root to x_i and x_j .
these two paths must be consistent so $x_i < x_j$.



So removing all BAD nodes will leave us with a sorted list!

" $\leq \epsilon n$ BAD indices \Rightarrow list is ϵ -close to sorted"



"List isn't ϵ -close to sorted $\Rightarrow \geq \epsilon n$ BAD indices"

Algorithm Analysis:

→ obviously YES if sorted.

→ if not ϵ -close to sorted \Rightarrow > en BAD indices.

$$\Pr[\text{NO}] = 1 - (\text{we never find BAD index}) \geq 1 - (1 - \epsilon)^{\frac{c}{\epsilon}} \geq 1 - e^{-c} \geq \frac{3}{4}$$

for $c \geq 2$.

LECTURE 23 5/9/24 (Did not attend)

Sketching!

Counting - string of n 1's, fed 1 at a time
 - want to find # of 1's (i.e. n)

↳ OPTION 1: just count $\rightarrow O(\log n)$ memory
 optimal space for exact answer.

but can we estimate n in $O(\log n)$ space?? (excluding output?)

Estimation of θ with $\tilde{\theta}$ (ϵ, δ) - estimate \rightarrow close to 0 as possible

\rightarrow want $\tilde{\theta} \in [(\mathbf{1}-\epsilon)\theta, (\mathbf{1}+\epsilon)\theta]$ with $\Pr \geq 1-\delta$.

Morris's Algorithm $\sim O(\log \log n)$ space

"count fewer as you go further along"

$$X = 1^n = n 1's \quad f(X) = n \quad c(X) \leftarrow \text{sketch.}$$

\uparrow
function to
be computed

\rightarrow initialize $c(X) = 0$

\rightarrow on input of 1, $c(X) += 1$ with $\Pr = \frac{1}{2^{c(X)}}$

\rightarrow output $\tilde{f}(X) = 2^{c(X)} - 1$.

toss fair coin
 $c(X)$ times
 and ONLY if
 all heads
only $O(c(X))$ space.

induct on $n=k$.

$$\mathbb{E}[\tilde{f}(X)] = \mathbb{E}[2^{c(X)}] - 1$$

$$= \sum_{j=0}^{k-1} \Pr[c(1^{k-1}) = j] \cdot \mathbb{E}[2^{c(1^k)} | c(1^{k-1}) = j] - 1$$

$$\begin{aligned}
&= \sum_{j=0}^{k-1} \Pr[C(1^{k-1}) = j] \cdot \left(2^{j+1} \cdot \frac{1}{2^j} + 2^j \cdot \left(1 - \frac{1}{2^j}\right) \right) - 1 \\
&= \sum_{j=0}^{k-1} \Pr[C(1^{k-1}) = j] \cdot (2^j + 1) - 1 \\
&= \sum_{j=0}^{k-1} \Pr[C(1^{k-1}) = j] \cdot 2^j + 1 - 1 \\
&= (k-1) + 1 = k
\end{aligned}$$

induction on $n=k$ to get $\mathbb{E}[2^{c(x)}] = \frac{3}{2}k^2 + \frac{3}{2}k + 1$

$$\begin{aligned}
\mathbb{V}[\tilde{f}(x)] &= \mathbb{E}[(2^{c(x)} - 1)^2] - \mathbb{E}[2^{c(x)} - 1]^2 \\
&= \frac{3}{2}k^2 + \frac{3}{2}k + 1 - 2(k+1) + 1 - k^2 \\
&< \frac{1}{2}k^2 \quad \Rightarrow \mathbb{E}[\tilde{f}] = n, \mathbb{V}[\tilde{f}] < \frac{n^2}{2}.
\end{aligned}$$

$$\text{Chebyshov.} \hookrightarrow \mathbb{P}[\tilde{f} - n \geq \epsilon n] \leq \frac{\frac{n^2}{2}}{(\epsilon n)^2} \leq \frac{1}{2\epsilon^2}$$

$$\mathbb{P}[\text{in interval}] \geq 1 - \frac{1}{2\epsilon^2} \quad \delta = \frac{1}{2\epsilon^2}$$

can we do better? (For general algorithm A with $\mathbb{E}=\Theta, \mathbb{V}=\epsilon^2, \text{space}=S$)

→ run several times and average (in parallel)

$$A^m(\epsilon, \delta) : \text{let } R = \lceil \frac{\epsilon^2}{\epsilon^2 \Theta^2 \delta} \rceil$$

run A R times and take average $\tilde{\theta} = \frac{\sum_{i=1}^R \tilde{\theta}_i}{R}$.

$$\text{we get } \mathbb{E}[\tilde{\theta}] = \Theta \quad \mathbb{V}[\tilde{\theta}] = \frac{1}{R} \sum_{i=1}^R \sigma^2 = \frac{\sigma^2}{R} \leq \epsilon^2 \Theta^2 \delta$$

$$\text{so } \Pr[|\tilde{\theta} - \theta| \geq \epsilon \theta] \leq \frac{\epsilon^2 \Theta^2 \delta}{\epsilon^2 \Theta^2} = \delta.$$

but this takes $O(RS) = O\left(\frac{\epsilon^2}{\epsilon^2 \Theta^2 \delta} \cdot S\right)$ space. ↪ $O\left(\frac{1}{\epsilon^2 \delta} \cdot \log \log n\right)$ for Morris.

→ MEDIAN of means of estimates.

$$A^{\text{mom}} : \text{let } R = \lceil 36 \ln(\frac{1}{\delta}) \rceil$$

run $A^m(\varepsilon, \frac{1}{3})$ R times and take their median.

Claim If $\geq \frac{1}{2}$ of A^m 's are in $[(1 \pm \varepsilon)\theta]$ then so is our result.

$$\Pr[\geq \frac{1}{2} \text{ of } A^m \text{ are out of range}] = \Pr[\text{fails} \geq (1 + \frac{1}{2}) \mathbb{E}[\text{fails}]]$$

$$\stackrel{\text{"fails"}}{\leq} e^{-\frac{1}{4} \cdot \frac{1}{3} R \cdot \frac{1}{3}} \leq e^{-\ln(\frac{1}{\delta})} = \delta.$$

$$\#[\text{fails}] = \frac{1}{3}R \Rightarrow \Pr[\text{overall success}] \geq 1 - \delta.$$

overall space is now

$$O(S \cdot \frac{3\sigma^2}{\varepsilon^2 \theta^2} \cdot \ln(\frac{1}{\delta})) = O(\frac{\sigma^2 S}{\varepsilon^2 \theta^2} \ln(\frac{1}{\delta})) \text{ for } A^m.$$

\nearrow $O(\frac{1}{\varepsilon^2} \log(\frac{1}{\delta}) \log \log n)$

\uparrow

we actually get

$\log \log(\frac{n}{\varepsilon \delta})$ by showing no counter exceeds $O(\log \frac{n}{\varepsilon \delta})$ with $\Pr > 1 - \delta$.

N := total # of counters in our algorithm

$$= \lceil \frac{\sigma^2}{\varepsilon^2 \theta^2 \cdot \frac{1}{3}} \rceil \cdot \lceil 36 \ln(\frac{1}{\delta}) \rceil = O\left(\frac{1}{\varepsilon^2} \log\left(\frac{1}{\delta}\right)\right)$$

\hookrightarrow if any counter exceeds $\log(\frac{Nn}{\delta})$

then $\Pr[\text{increase}] \leq \frac{\delta}{Nn} \leftarrow$ over n time steps and N counters

$\Pr[\text{any counter increases beyond } \log(\frac{Nn}{\delta})] \leq \delta \quad \hookrightarrow O(\log \log(\frac{n}{\delta \varepsilon}))$ space

$$\log(\frac{Nn}{\delta}) = O\left(\log\left(\frac{1}{\varepsilon^2} \log\left(\frac{1}{\delta}\right) \cdot n \cdot \frac{1}{\delta}\right)\right) = O\left(\log\left(\frac{n}{\varepsilon \delta}\right)\right)$$

\hookrightarrow overall $O\left(\frac{1}{\varepsilon^2} \log\left(\frac{1}{\delta}\right) \log \log\left(\frac{n}{\delta \varepsilon}\right)\right)$

time complexity.

\rightarrow each $A^m(\varepsilon, \frac{1}{3})$ takes $O(\log(\frac{n}{\varepsilon \delta}) \cdot n)$ \Rightarrow run all R in parallel

$\uparrow \quad \uparrow$

max counter value numbers to read

$\left(\begin{array}{l} \rightarrow \text{max flips of} \\ \text{coin per step} \end{array} \right)$

\rightarrow median of $O(\log(\frac{1}{\delta}))$ items double in $O(\log(\frac{1}{\delta}))$

\rightarrow overall $O(n \log(\frac{n}{\varepsilon \delta}) + \log(\frac{1}{\delta}))$

Distinct Elements

count # of distinct elements in a stream.

FM85 Algorithm - suppose we have a truly random $h: \mathcal{U} \rightarrow [0, 1]$

↳ initialize $c(x) = 1$

↳ after x_i , compute $h(x_i)$, update $c(x) = \min \{ c(x), h(x_i) \}$

↳ output $\tilde{f}(x) = \frac{1}{c(x)} - 1$.

$$\mathbb{E}[c(x)] = \frac{1}{d+1} \quad \text{based on } d \text{ uniform values } \in [0, 1]$$

$$\mathbb{E}[c(x)^2] = \int_0^1 \Pr[c(x)^2 \geq y] dy = \int_0^1 \Pr[c(x)^2 \geq \sqrt{y}] dy$$

$$y = x^2 \quad = \int_0^1 (1 - \sqrt{y})^d dy = \int_0^1 2 \times (1-x)^d dx$$

$$dy = 2x dx \quad = \int_0^1 2(1-x)x^d dy = \frac{2}{d+1} - \frac{2}{d+2} = \frac{2}{(d+1)(d+2)}$$

$$\mathbb{V}[c(x)] = \frac{2}{(d+1)(d+2)} - \frac{1}{(d+1)(d+1)} = \frac{d}{(d+1)(d+2)}$$

but we don't have a truly random h at our disposal!.

↳ t-wise independent hashing: \forall distinct $k_1, \dots, k_t \in \mathcal{U}$
not-necessarily distinct $i_1, \dots, i_t \in \mathcal{M}$

$$\Pr_{h \in \mathcal{H}} \left[\bigwedge_{j=1}^t h(k_j) = i_j \right] = \frac{1}{m^t}$$

can achieve with $\mathcal{H} = \{h_a : a \in \mathcal{M}^t\}$, $h_a(k) = \sum_{\ell=0}^{t-1} a_\ell k^\ell \pmod{m}$

$h_a(k)$: time, space both $O(t \log m)$, $m = |\mathcal{M}|$, to store a .

KMV Algorithm

↳ pick random $h: \mathcal{U} \rightarrow M = \{0, 1, \dots, n^3\}$, $|M| = n$, from 2-wise independent family

↳ $k = \lceil \frac{24}{\epsilon^2} \rceil$, initialize $C(X) = \emptyset$, $f(X) = d$, input stream (x_1, \dots, x_n)

↳ at each step let $C(X) = \min_k \{C(X) \cup \{h(x_i)\}\}$

$$\text{output } \tilde{f}(X) = \begin{cases} |C(X)| & \text{if } |C(X)| < k \\ \frac{k n^3}{\max \{C(X)\}} & \text{otherwise} \end{cases}$$

old $C(X) \approx \frac{\max \{C(X)\}}{k n^3} \leftarrow \begin{matrix} \text{protect from} \\ \text{high variance} \end{matrix}$
 $\uparrow \quad \uparrow$
 $k^{\text{th}} \text{smallest scale to } [0, 1]$

$$\Pr [\tilde{f} \in [(1 \pm \epsilon)d]] \geq \frac{2}{3}$$

↳ any $1 - \delta$ if using mom.

time = n evaluations of h + n updates of $C(X)$

$$= n \log n + n \log k = O(n \log n k) = O\left(n \log \frac{n}{\epsilon^2}\right)$$

space = storing h + storing $C(X)$

$$= \log n + k \log n = O(k \log n) = O\left(\frac{\log n}{\epsilon^2}\right)$$

Similarity Search - n sets $A_1, A_2, \dots, A_n \subseteq \mathcal{U}$, assume $|A_i| \leq d$

Given $A \subseteq \mathcal{U}$, detect if $\exists A_i$ s.t. $J(A, A_i) = \frac{|A \cap A_i|}{|A \cup A_i|} \geq s$.

Linear Scan takes $\Omega(nd)$

Fast Approximate Linear Scan

Modify so that

\rightarrow if $J(A, A_i) \geq s$ for some A_i , YES } both satisfied
 \rightarrow if $J(A, A_i) < s'$ for all A_i , NO } with $> \frac{1}{2}$ chance.

use randomness, t hash functions $h_1, h_2, \dots, h_t \rightarrow [0, 1]$

$$\text{let } \sigma_h(A_i) = (\sigma_{h_1}(A_i), \dots, \sigma_{h_t}(A_i)), \quad \sigma_{h_j}(A_i) = \min_{a \in A_i} h_j(a)$$

note $\Pr[\sigma_h(x) = \sigma_h(y)] = J(x, y)$ for random $h \in [0, 1]$

estimate $\text{Sim}(\sigma_h(A), \sigma_h(A_i)) = \frac{|\{j : \sigma_{h_j}(A) = \sigma_{h_j}(A_i)\}|}{t}$

YES if $\text{Sim}_i \geq s - \epsilon$, NO otherwise.

if $t = \frac{c}{\epsilon^2} \log(\frac{1}{\delta})$ $\Pr_h [\left| \text{Sim}(\sigma_h(x), \sigma_h(y)) - J(x, y) \right| > \epsilon] \leq \delta$

using Chernoff.

$$\delta = \frac{1}{2n} \rightarrow t = O\left(\frac{\log n}{\epsilon^2}\right), \text{ then with probability } \geq 1 - \frac{n}{2n} = \frac{1}{2}$$

all A_i
simultaneously

we get YES if $\exists J(A, A_i) \geq s$ and NO if no $J(A, A_i) \geq s - 2\epsilon$

runtime $O(ndt)$ preprocessing of $\sigma_h(A_i)$'s

$$O((n+o)t) \begin{cases} O(dt) \text{ computing } \sigma_h(A) \\ O(nt) \text{ finding how many } j : \sigma_{h_j}(A) = \sigma_{h_j}(A_i) \text{ for all } 1 \leq j \leq t \\ 1 \leq i \leq n \end{cases}$$

Sublinear Nearest Neighbor Search

→ initialize h_1, h_2, \dots, h_L s.t.

$$\sigma_{h_L}(A) = (\sigma_{h_{L,1}}(A), \dots, \sigma_{h_{L,t}}(A)) \quad \text{where } \sigma, h_{L,j} \rightarrow [0, 1]$$

are from before

→ initialize L hash tables.

→ store set A_i in hash table l with key $\sigma_{h_L}(A)$ for all $1 \leq l \leq L$

for A , calculate $\sigma_{h_L}(A)$ then compute $J(A, A_i)$ $\forall i$ where

$\sigma_{h_L}(A) = \sigma_{h_L}(A_i)$ in some hash table. If $\geq s'$, return YES. Else NO.

If $J(A, A_i) \geq s$ for some $1 \leq i \leq n$, then

$$\Pr[\sigma_{h_L}(A) = \sigma_{h_L}(A_i)] \geq s^t \quad \text{for each } l.$$

→ probability of failure $\leq (1-s^t)^L$. always NO if $\neq A_i$
YES with $\Pr \geq 1 - (s-s^t)^L$ if $\exists A_i$

Expected runtime: if all $J(A, A_i) \leq s'$

preprocessing by computing $\sigma_{h_L}(A_i) \forall i, l \rightarrow O(Lndt)$

$\exp O(Ld(t+ns't))$ $\left\{ \Pr [\sigma_{h_L}(A) = \sigma_{h_L}(A_i)] \leq s'^t, \text{ so } O(Lnd s'^t)$ check time
 computing key $\sigma_{h_L}(A)$ $\forall 1 \leq L \leq L$ takes $O(Ldt)$

ex) $s = \frac{1}{2}, s' = \frac{1}{4} \rightarrow t = \frac{\log n}{2}; L = \sqrt{n}$

$$\Pr[\text{fail}] \leq (1 - \left(\frac{1}{2}\right)^{\frac{\log n}{2}})^{\sqrt{n}} = \left(1 - \frac{1}{\sqrt{n}}\right)^{\sqrt{n}} \leq e^{-1} \approx 0.37,$$

expected time

Sublinear!

$$O(\sqrt{n} d (\log n / 2 + n \cdot \left(\frac{1}{4}\right)^{\frac{\log n}{2}})) = O(\sqrt{n} d \log n + n \sqrt{n} d \cdot \frac{1}{n}) = O(\sqrt{n} d \log n)$$