

2/8/23

HS1090

classes: M, Tu, W, Th (cont'd)

Eval

$$\text{Quiz 1,2} - 15\% \times 2 = 30\%$$

$$- 40\%$$

Endsem =

Spoken assignment

- 10% (task, record video)

[+5 mins.]

Group task

→ presented to class
after quiz 2

- 20% (in English, group of 3-4)

- some aspect of Germany,
culture, lit., philosophy, history

- banned topics: Hitler, Cars,
beer, football

Kaiser - czar - caesar - king

Wechtenstein

2 syllables

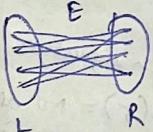
* Oder-Neisse line

2/8/23

CSG170

Perfect matching in bipartite graphs

$G(L, R, E)$



A matching $M \subseteq E$ s.t.

$\forall v \in L \cup R, \exists$ at most one $e \in M$ that is incident on v

A perfect matching is a bijection

$$(|L| = |R|)$$

Q: Given $G(L, R, E)$,

check if G has a perfect matching.

there exist polynomial-time deterministic algos;
but we're interested in looking at a rand. one

Q': Given $G(L, R, E)$, compute a p.m. if one exists.

We can consider some edge & check if it is present in any p.m. by deleting it &

reusing on the smaller graph \Rightarrow property called self-reducibility

We form a soln to this by solving a diff prob.

Polynomial Identity Testing (PIT)

1/p: Program C that computes an n -variate degree d poly.

$$p(x_1, \dots, x_n)$$

$$C(a_1, \dots, a_n) = p(a_1, \dots, a_n)$$

Q: check if $p \equiv 0$.

i.e., $\forall x_1, \dots, x_n, p(x_1, \dots, x_n) = 0$

$$\text{No. of coefficients} = \binom{n+d-1}{n} \approx n^d$$

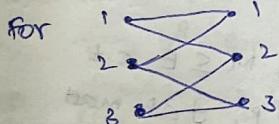
How are these probs. related?

For $G(L, R, E)$, consider the bipartite adj. mat.

$$A \in \{0, 1\}^{n \times n}, n = |L| = |R|$$

$$A(i, j) = \begin{cases} 1, & \text{if } (i, j) \in E \\ 0, & \text{o.w.} \end{cases} \quad [i \in L, j \in R]$$

$$p(i, j) = \underbrace{x_{ij}}_{\text{variable}} A(i, j)$$



$$P = \begin{bmatrix} x_{11} & x_{12} & 0 \\ x_{21} & 0 & x_{23} \\ 0 & x_{32} & x_{33} \end{bmatrix}$$

$$\det(P) = \sum_{\sigma \in S_n} (-1)^{\text{sgn}(\sigma)} \prod_{i=1}^n p(i, \sigma(i))$$

signature
 $\text{sgn}(\sigma) = \text{no. of inversions in } \sigma$

set of bijections from $\{1, \dots, n\}$ to $\{1, \dots, n\}$

$|E|$ -variate degree n poly w/ at most $n!$ terms

Each surviving term corresponds to a valid perfect matching

Obs: G has a P.M. iff $\det(P) \neq 0$

\Rightarrow solving PIT solves P.M. (w/ a n. simple algo)

Solving PPT

DeMillo-Lipton-Schwartz-Zippel lemma:

Let p be a non-zero n -variate deg.d poly. over \mathbb{R} .

Let $S \subseteq \mathbb{R}$. Then

$$\Pr_{\substack{a_1, a_2, \dots, a_n \in S \\ \text{chosen indep.} \\ \& \text{uniformly} \\ \text{at random}}} [p(a_1, a_2, \dots, a_n) = 0] \leq \frac{d}{|S|} \quad (\text{indep. of } n)$$

3/8/23

$a_1, a_2, \dots, a_n \in S$

chosen indep.
& uniformly
at random

$$\text{Here, sample space } \Omega = \underbrace{S \times S \times \dots \times S}_{n \text{ times}} = S^n$$

The event of interest is the set of zeros,

$$\text{i.e., } \emptyset \subseteq S^n = \{(a_1, a_2, \dots, a_n) \in S^n \mid p(a_1, \dots, a_n) = 0\}$$

$$\Pr(\emptyset) = ?$$

The proba. space is $(\Omega, \mathcal{F}, \Pr)$

$$\mathcal{F} = \mathcal{P}(\Omega) = 2^\Omega \quad (\text{set of all events})$$

$\Pr: \Omega \rightarrow [0, 1]$ (proba. of choosing each tuple)

$$\Pr(a_1, \dots, a_n) = \frac{1}{|S|^n} \quad (\text{indep. \& u.a.r.})$$

It satisfies:

$$\textcircled{1} \quad \forall E \in \mathcal{F} \quad 0 \leq \Pr(E) \leq 1$$

$$\textcircled{2} \quad \Pr(\Omega) = 1$$

\textcircled{3} If E_1, E_2, \dots, E_n are pairwise disjoint,

$$\Pr(\cup E_i) = \sum_i \Pr(E_i)$$

Some facts:

* Let E_1, E_2, \dots, E_n be any n events

$$\Pr(\cup E_i) \leq \sum_{i=1}^n \Pr(E_i) \quad [\text{union bound}]$$

* Law of total proba.:

* conditional proba.:

$$\Pr(E|F) = \frac{\Pr(E \cap F)}{\Pr(F)}$$

E_1, \dots, E_n For disjoint events E_1, E_2, \dots, E_n
partition of Ω s.t. $\cup E_i = \Omega$,

* E & F are indep. if $\Pr(E|F) = \Pr(E)$

$$\Rightarrow \Pr(E \cap F) = \Pr(E) \Pr(F)$$

$$\Pr(F) = \sum_{i=1}^n \Pr(F \cap E_i)$$

Proof of DLSZ-lemma

[By induction on n]

Base case: $n = 1$

$$\Pr_{a \in S} [p(a) = 0] \leq \frac{d}{|S|}$$

$$Z = \{a \in S \mid p(a) = 0\}$$

$$|Z| \leq d$$

$$\Rightarrow \Pr(Z) \leq \frac{d}{|S|}$$

Alg0: (A)

sample $a_1, \dots, a_n \in_{\text{i.i.d.}} S$
u.a.r.

If $p(a_1, \dots, a_n) = 0$,
say $P = 0$

else, say $P \neq 0$

Note:

① If $P = 0$, then $\Pr(A \text{ says } P = 0) = 1$

② If $P \neq 0$, then

$$\underbrace{\Pr(A \text{ says } P = 0)}_{\text{error proba.}} \leq \frac{d}{|S|}$$

\Rightarrow one-sided error alg0.

We can boost the success proba. by repeating the sampling.
i.e., A'

repeat K times

- sample $a_1, \dots, a_n \in S$

- If $p(a_1, \dots, a_n) \neq 0$, say $P \neq 0$

Say $P = 0$

If $P = 0$, $\Pr(A' \text{ says } P = 0) = 1$

If $P \neq 0$, $\Pr(A' \text{ says } P = 0)$

= $\Pr(p = 0 \text{ in every iteration})$

$$= \Pr(Z_1 \cap Z_2 \cap \dots \cap Z_K) = \Pr(Z_1) \cdot \Pr(Z_2) \cdot \dots \cdot \Pr(Z_K)$$

$$\leq \left(\frac{d}{|S|}\right)^K$$

Induction step:

$$p(x_1, \dots, x_n) = \sum_{i=0}^d \alpha_i p_i(x_1, \dots, x_n)$$

find largest i s.t. $p_i \neq 0$

$Z_i = \text{set of tuples on which } p_i \text{ evals. to 0}$

$$\deg. p_i \leq d - i$$

$$\Rightarrow \Pr(Z_i) \leq \frac{d-i}{|S|} \quad (\text{induction hypothesis})$$

$$\Pr(Z) = \Pr(Z \cap Z_i) + \Pr(Z \cap \bar{Z}_i) \leq 1$$

$$= \Pr(Z_i) \cdot \Pr(Z \mid Z_i) + \Pr(\bar{Z}_i) \cdot \Pr(Z \mid \bar{Z}_i)$$

we have a lower bound, but we want an upper bound,
 \Rightarrow we can just choose 1

$$\leq \frac{d-i}{|S|} \cdot 1 + 1 \cdot \frac{i}{|S|}$$

$$\leq \frac{d}{|S|}$$

p_i 's are evaluated, now we have a poly. on x_i 's w/ deg. i

18/23

Verifying matrix multiplication (Friordan's algorithm)

Given $n \times n$ matrices A, B, C ,

check if $A \cdot B = C$

- can you do this asymptotically faster than matrix multiplication?

$O(n^3)$ - naive

$O(n^{2.37})$ - complex

conjectured to be $O(n^{2+\epsilon})$ for v. small ϵ

Friordan's algo. \rightarrow randomized $O(n^2)$

Algo

- choose $\bar{r} \in \{0,1\}^n$

- check if $\underbrace{AB\bar{r}}_{A(B\bar{r})} = \underbrace{C\bar{r}}_{O(n^2)}$

$$\Omega = \{0,1\}^n$$

$$\Pr(\bar{a}) = \frac{1}{2^n} \quad \forall \bar{a} \in \Omega$$

Theorem

If $A \cdot B \neq C$,

$$\Pr_{\bar{r} \in \{0,1\}^n} [AB\bar{r} = C\bar{r}] \leq \frac{1}{2}$$

non-zero matrix $D = AB - C$

$$\Pr_{\bar{r} \in \{0,1\}^n} [D\bar{r} = 0]$$

Consider that the first row of D is non-zero (wlog) & $D_{1,1} \neq 0$

$$\leq \frac{1}{2}$$

by DLSZ-lemma,

where variables

are $r_i, i=1, \dots, n$

& identically

zero implies $r \in \{0,1\}^n$

$$AB = C_{1,:}$$

$$\Pr_{r \in \{0,1\}^n} \left(\sum_{i=1}^n D_{1,i} r_i = 0 \right)$$

$$(D_{1,1}, D_{1,2}, \dots, D_{1,n}) \begin{pmatrix} r_1 \\ r_2 \\ \vdots \\ r_n \end{pmatrix}$$

$$= \Pr_{r \in \{0,1\}^n} \left(r_1 = -\frac{\sum_{i=2}^n D_{1,i} r_i}{D_{1,1}} \right)$$

or

$r_1, r_2, \dots, r_n \in \{0,1\}$ \rightarrow we choose r_2, \dots, r_n

[Principle of deferred decisions]

& finally choose r_1

If we want to repeat it until error proba.

is ϵ ,

we want $(\frac{1}{2})^n = \epsilon$

$$n = -\log_2 \epsilon = \log_2 \frac{1}{\epsilon}$$

$$= \sum_{\substack{b_2, \dots, b_n \\ \in \{0,1\}}} \Pr_{r_1 = -\frac{\sum_{i=2}^n D_{1,i} b_i}{D_{1,1}}} \left[r_1 = -\frac{\sum_{i=2}^n D_{1,i} b_i}{D_{1,1}} \right]_{r_2 = b_2, \dots, r_n = b_n}$$

$$\cdot \Pr_{r_2 = b_2, \dots, r_n = b_n} [r_2 = b_2, \dots, r_n = b_n]$$

$$\leq \frac{1}{2}$$

Min-cut in graphs

$G(V, E)$

A cut set is a set of edges whose removal disconnects G
Min-cut is a cut set of smallest size/cardinality

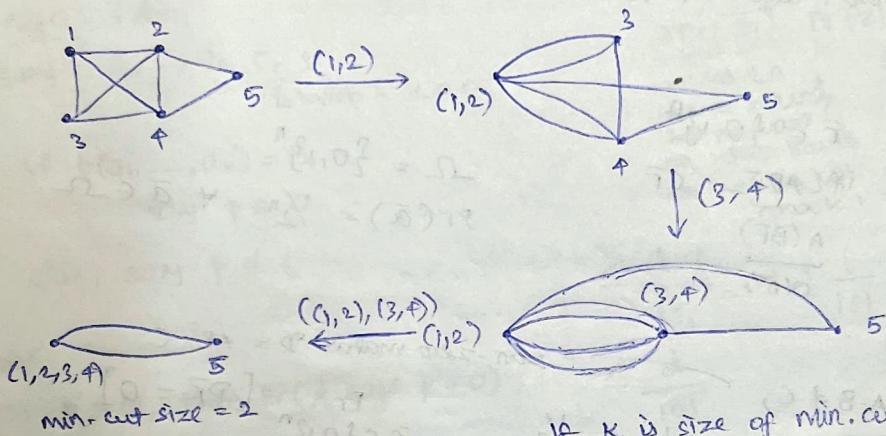
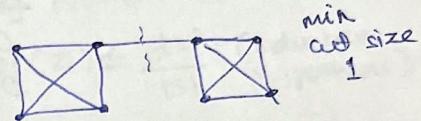
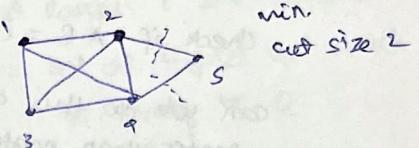
- A rand-algo. (Karger '93)

repeat } - choose $e \in E$ u.a.r.
until 2 vertices remain } - contract e

$$G' \leftarrow G \setminus \{e\}$$

& keep parallel edges

- remaining edges belong to
min cut



If K is size of min.cut,
 $|E| \geq \frac{nK}{2}$ (degree of each vertex $\geq K$)

Obs: If G has min-cut size K ,
then G has at least $\frac{Kn}{2}$ edges

[If min-cut size $\leq K$, then

$\forall u \in V, \deg(u) \geq K$

$$\text{no. of edges} = \frac{1}{2} \sum_{u \in V} \deg(u) \geq \frac{nk}{2}$$

fix a min-cut C of size K

$$\Pr(\text{no edge from } C \text{ is sampled}) \geq 1 - \frac{2}{n}$$

7/18/23

E_i : Event that an edge from C was not sampled
in the i^{th} step

$$\Pr(E_i) \geq 1 - \frac{2}{n}$$

Assume that E_1 occurred

$$G' = G \setminus \{e\}$$

obs: Every cut in G' is

a cut in G

\Rightarrow The min-cut size in
 G' is also K

$F_i = \bigcap_{j=1}^i E_j$ = the event that no edge from C was sampled in the first i iterations

We are interested in $\Pr(F_{n-2})$

$$\Pr(E_2 | F_1) \geq 1 - \frac{2}{n-1} \quad [\text{no of edges in } G \setminus \{e\} \geq \frac{(n-1)k}{2}]$$

$$\Pr(E_i | F_{i-1}) \geq 1 - \frac{2}{n-i+1}$$

$$\Pr(F_{n-2}) = \Pr(F_{n-3}) \cdot \Pr(E_{n-2} | F_{n-3})$$

$$= \Pr(E_{n-2} | F_{n-3}) \cdot \Pr(E_{n-3} | F_{n-4}) \cdot \dots \cdot \Pr(E_1)$$

$$\geq \prod_{i=1}^{n-2} \frac{n-i-1}{n-i+1} = \frac{1 \times 2}{n(n-1)}$$

$$\geq \frac{2}{n(n-1)} \rightarrow \text{v. low success prob.}$$

for some fixed \min cut $\rightarrow O(n^2)$ min cuts possible, more not possible

$$\Pr(C \text{ was not outputted}) \leq 1 - \frac{2}{n(n-1)}$$

Repeat K times

$$\Pr(C \text{ was not output in any iteration}) \leq \left[1 - \frac{2}{n(n-1)}\right]^K$$

$$1-x \leq e^{-x}$$

Max-cut

$$G(V, E)$$

partition $V = V_1 \cup V_2$ s.t.

$|E(V_1, V_2)|$ is maximized
no. of edges across the partition/cut

We analyse a $\frac{1}{2}$ -approx

- $O(2^n)$ brute-force algo.

- NP-hard

- α -approximation

- Algo. that outputs a cut of size $\geq \alpha(\max)$, $\alpha < 1$

- Best known poly. time: $\alpha = 0.8$ for

- $\alpha \geq 0.94$ not possible unless $P = NP$

Algo. A

$\forall v \in V$,
put v in V_i w.p. $1/2$ ind. of other vertices

$A(x, r) \rightarrow$ rand. var.

$A: \Omega \rightarrow \mathbb{R}$
 \downarrow
 n -bit str.
 denoting
 a partition
 cut size

Theorem: $\mathbb{E}(x) \geq \frac{m}{2}$
 \uparrow
 total no.
 of edges

size of the
 rand. cut

Modified algo.

- sample logn bits u.a.r. ~~rand~~
- construct n pairwise indep rand-bits defining a cut
- output the size of the cut

No. of rand. choices = $2^{\text{logn}} = n$

\Rightarrow Brute-forcing over all logn bit strings obtains a poly-time deterministic algo.

same results as before
 $E(X) = m/2$

over the logn bits

10/8/23

Alternate procedure for derandomization

Method of conditional expectation

conditional expectation: Given r.v.s X, Y ,

$$E(X|Y=y) = \sum_{x_i} x_i \cdot \Pr(X=x_i | Y=y)$$

fact: 1) $E(X) = \sum_y \Pr(Y=y) E(X|Y=y)$

2) $E\left(\sum_{i=1}^n x_i | Y=y\right) = \sum_{i=1}^n E(x_i | Y=y)$

$$E(X) = \Pr(V_1 \in V_1) \cdot E(X|V_1 \in V_1) + \Pr(V_1 \in V_2) \cdot E(X|V_1 \in V_2)$$

$$= \frac{1}{2} \cdot E(X|V_1 \in V_1) + \frac{1}{2} \cdot E(X|V_1 \in V_2) = \frac{m}{2} \quad (\text{from analysis before})$$

Put $v_i \in V_1$ arbitrarily

$$E(X|V_1 \in V_1) = \Pr(V_2 \in V_1) \cdot E(X|V_1 \in V_1 \& V_2 \in V_1) \quad \begin{matrix} \nearrow 1/2 \\ \searrow m/2 \end{matrix}$$

$+ \Pr(V_2 \in V_2) \cdot E(X|V_1 \in V_1 \& V_2 \in V_2) \quad \begin{matrix} \nearrow 1/2 \\ \searrow m/2 \end{matrix}$

these must be $\geq m/2$

keep choosing the larger

$$E(X|V_1, V_2, \dots, V_{i-1}) = \Pr(V_i \in V_1) \cdot E(X|V_1, V_2, \dots, V_{i-1}, V_i \in V_1) + \Pr(V_i \in V_2) \cdot E(X|V_1, V_2, \dots, V_{i-1}, V_i \in V_2)$$

$$\hookrightarrow \geq \frac{m}{2}$$

(from prev. steps)

sequential greedy algo:

- Fix some order of vertices v_1, \dots, v_n

- For $i=1, \dots, n$:

- Add v_i to V_i if no. of edges from v_i to v_1, \dots, v_{i-1} is larger

\hookrightarrow Else, add v_i to V_2 .

\hookrightarrow No. of cut edges b/w v_1, \dots, v_{i-1} & v_i
 + No. of cut edges incident on v_i
 $+ \frac{1}{2} \times \text{No. of other edges}$ (at most one vertex incident on v_1, \dots, v_n)

11.10.23 Quicksort

- Randomly choose a pivot
- Partition array based on pivot $\rightarrow O(n)$
- Recursively sort the two parts

Running time:

worst-case $O(n^2)$

Expected $\Theta(n \log n)$

running time \propto No. of comparisons

$$\bar{a} = a_1, a_2, \dots, a_n \xrightarrow{\text{sorted}} b_1, b_2, \dots, b_n$$

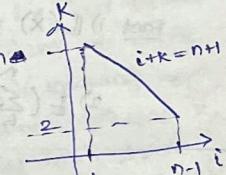
$$X_{ij} = \begin{cases} 1, & \text{if } b_i \text{ & } b_j \text{ were compared} \\ & \text{at some pt. in Quicksort} \\ 0, & \text{o.w.} \end{cases} \rightarrow \text{this occurs when}$$

the first pivot
in $b_i, b_{i+1}, \dots, b_{j-1}, b_j$
appears somewhere in b_i or b_j

$$X = \sum_{i < j} X_{ij}$$

$$E(X_{ij}) = \Pr(X_{ij} = 1) = \frac{2}{j-i+1}$$

$$\begin{aligned} E(X) &= \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{2}{j-i+1} \\ &= \sum_{i=1}^{n-1} \sum_{k=2}^{n-i+1} \frac{2}{k} = \sum_{k=2}^n \sum_{i=1}^{n-k+1} \frac{2}{k} = \sum_{k=2}^n (n-k+1) \frac{2}{k} \\ &= 2(n+1) \sum_{k=2}^n \frac{1}{k} - 2(n-1) \\ &\quad \xrightarrow{\text{log } n + \Theta(1)} \star \text{Harmonic sum} \end{aligned}$$



$$= 2n \log n + \Theta(n)$$

\Rightarrow Expected $\Theta(n \log n)$

we will show later that

$$\Pr(X > cn \log n) \leq \frac{1}{n}$$

A similar analysis holds if we randomly permute the array and then choose the first/last/some fixed index element as the pivot repeatedly

A few standard r.v.s

① Bernoulli r.v. \equiv Indicator r.v.

$$X = \begin{cases} 1, & \text{w.p. } p \\ 0, & \text{w.p. } 1-p \end{cases} \quad E(X) = \Pr(X=1) = p$$

② Binomial r.v.

Let X_1, X_2, \dots, X_n be Bernoulli trials, i.e., rand. expts. w/ success proba. p

A binomial r.v. counts the no. of successes in n Bernoulli trials, i.e., $X = \sum_i X_i$

Range of $X = \{0, 1, \dots, n\}$

$$\Pr(X=i) = \binom{n}{i} p^i (1-p)^{n-i} \rightarrow \text{prob. mass fn.}$$

$$E(X) = np$$

③ Geometric r.v.

- * Counts the no. of Bernoulli trials before the first success
- * Range of the r.v. = $\mathbb{N} = \{1, 2, \dots\}$
- * $\Pr(X=i) = (1-p)^{i-1} p$
- * Memorylessness

$$\Pr(X=n+k \mid \underbrace{X > k}_{\substack{\text{no success} \\ \text{in first} \\ k \text{ trials}}}) = \Pr(X=n)$$

$$\begin{aligned} * E(X) &= \Pr(Y=1) \cdot \underbrace{E(X|Y=1)}_{1} + \Pr(Y=0) \cdot \underbrace{E(X|Y=0)}_{x=1+z, z \text{ being an identical geometric r.v.}} \\ &= p + (1-p)(E(1+z)) \\ &= 1 + (1-p)E(z) = 1 + (1-p)E(X) \\ \Rightarrow pE(X) &= 1 \Rightarrow E(X) = \frac{1}{p} \end{aligned}$$

Coupon collector problem

n coupons

A box of chocolate has one of the coupons u.a.r.

Q. How many boxes must be collected to obtain one copy of every coupon?

X_i = No. of boxes that are bought after $i-1$ coupons are obtained before obtaining the i^{th} coupon

X_i is a geometric r.v. with success prob. $\frac{n-i+1}{n}$

No. of boxes, $X = \sum_{i=1}^n X_i$

$$\Pr(X) = \sum_{i=1}^n \Pr(X_i) = \sum_{i=1}^n \frac{n}{n-i+1} = n \sum_{j=1}^n \frac{1}{j}$$

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

Markov's inequality

If X is a non-neg. r.v. & $a > 0$, then

$$\Pr(X \geq a) \leq \frac{E(X)}{a}$$

$$\text{i.e., } \Pr(X \geq kE(X)) \leq \frac{1}{k}$$

Pf:

consider the indicator r.v. $I = \begin{cases} 1, & X \geq a \\ 0, & \text{o.w.} \end{cases}$

$$\Pr(X \geq a) = \Pr(I=1) = E(I)$$

Note that $I \leq \frac{X}{a}$

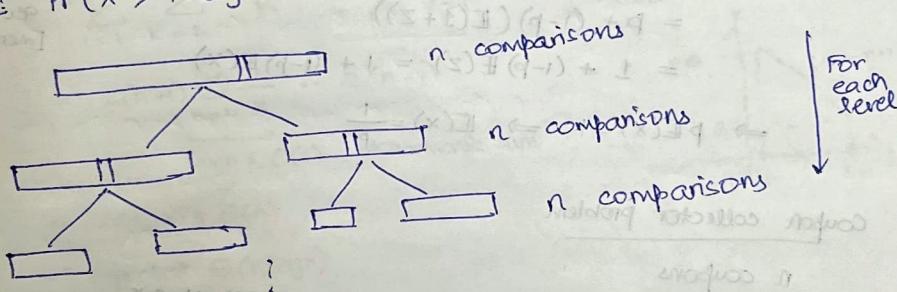
$$\Rightarrow E(I) \leq E\left(\frac{X}{a}\right) = \frac{E(X)}{a}$$

$$\therefore \Pr(X \geq a) \leq \frac{E(X)}{a}$$

A bound for quicksort

$X = \text{no. of comparisons}$, $E(X) = 2n \log n + O(n)$

Thm: $\Pr(X > kn \log n) \leq \frac{1}{n}$, k is a constant indep. of n



Each level performs n comparisons

\Rightarrow Bounding no. of levels bounds no. of comparisons

Fix an element a of the array

Fix an element a of the array containing a in the i th level of recursion

$x_i = \text{size of the array containing } a \text{ in the } i\text{th level of recursion}$

$$x_0 = n$$

$$E(x_i) = \sum_{t>0} \Pr(x_{i-1}=t) \cdot E(x_i | x_{i-1}=t)$$

We want this to be $\propto x_{i-1}$ to bound levels as $O(\log n)$

Choosing an element in 25th-75th percentile bounds the size of the array with a to be at most $\frac{3}{4}x_{i-1}$

$$\begin{aligned} \mathbb{E}(x_i | x_{i-1} = t) &\leq \frac{1}{2} \cdot \frac{3}{4} x_i + \frac{1}{2} \cdot x_i = \frac{7}{8} t \\ \Rightarrow \mathbb{E}(x_i) &\leq \sum_{t>0} \Pr(x_{i-1} = t) \cdot \frac{7}{8} t = \frac{7}{8} \sum_{t>0} \Pr(x_{i-1} = t) \cdot t \end{aligned}$$

$$\begin{aligned} \mathbb{E}(x_i) &\leq \frac{7}{8} \mathbb{E}(x_{i-1}) \\ \Rightarrow \mathbb{E}(x_i) &\leq \left(\frac{7}{8}\right)^n \end{aligned}$$

choose $i = 3\log_{8/7}^n$

$$\begin{aligned} \Rightarrow \mathbb{E}(x_i) &\leq n \left(\frac{7}{8}\right)^{3\log_{8/7}^n} = n \left(\frac{1}{n^3}\right) \\ &\leq \frac{1}{n^2} \end{aligned}$$

we have only considered the bound for a single element a .

By Markov's inequality, for $a=1$,

$$\Pr(x_i > 1) \leq \frac{1}{n^2} \text{ for } i = 3\log_{8/7}^n$$

Let E_a be the event that a exists in a subarray of size > 1

after $i = 3\log_{8/7}^n$ recursive steps.

$$\begin{aligned} \Pr(E_a) &\leq \frac{1}{n^2} \\ \Pr(\exists a \text{ s.t. } E_a) &\leq \sum \Pr(E_a) = n \left(\frac{1}{n^2}\right) \\ &\leq \frac{1}{n} \end{aligned}$$

Variance of a r.v.

$$\begin{aligned} \text{var}(x) &= \mathbb{E}((x - \mathbb{E}(x))^2), \quad \mathbb{E}(x) = \mu \\ &= \mathbb{E}(x^2 - 2\mu x + \mu^2) \\ &= \mathbb{E}(x^2) - \mu^2 \end{aligned}$$

Properties

- ① $\text{var}(kx) = k^2 \text{var}(x)$
 - ② $\text{var}(x+y) = \text{var}(x) + \text{var}(y) + 2\text{cov}(x, y)$,
 $\text{cov}(x, y) = \mathbb{E}((x - \mathbb{E}(x))(y - \mathbb{E}(y))) = \mathbb{E}(xy) - \mathbb{E}(x)\mathbb{E}(y)$
- If x & y are independent,
- $$\text{var}(x+y) = \text{var}(x) + \text{var}(y)$$

* $X \sim \text{Bin}(n, p)$

$$x = x_1 + x_2 + \dots + x_n, \quad x_i \sim \text{Bernoulli}(p)$$

$$\text{var}(x_i) = \mathbb{E}(x_i^2) - \mathbb{E}(x_i)^2 = p - p^2 = p(1-p)$$

$$\text{var}(x) = \sum_{i=1}^n \text{var}(x_i) = np(1-p)$$

* $X \sim \text{Geom}(p)$

$$\begin{aligned} \mathbb{E}(X^2) &= p \cdot 1 + (1-p) \cdot \mathbb{E}((X+1)^2) = p + (1-p) [\mathbb{E}(X^2) + 2\mathbb{E}(X) + 1] \\ &= 1 + \frac{p(1-p)}{p} + (1-p)\mathbb{E}(X^2) \end{aligned}$$

$$p\mathbb{E}(X^2) = \frac{2-p}{p} \Rightarrow \mathbb{E}(X^2) = \frac{2-p}{p^2}$$

$$\text{var}(X) = \mathbb{E}(X^2) - \mathbb{E}(X)^2 = \frac{2-p}{p^2} - \frac{1}{p^2} = \frac{1-p}{p^2}$$

chebyshov's inequality

Let X be any r.v. & $a > 0$

$$\Pr(|X-\mu| \geq a) \leq \frac{\text{var}(X)}{a^2}$$

$$\text{i.e., } \Pr(\mu-a \leq X \leq \mu+a) \geq 1 - \frac{\text{var}(X)}{a^2}$$

Pf: Consider $Y = (X-\mu)^2$

Y is a non-neg r.v.

$$\Rightarrow \text{By Markov's inequality, } \Pr(Y \geq a^2) \leq \frac{\mathbb{E}(Y)}{a^2}$$

$$\Pr(|X-\mu| \leq a) = \Pr((X-\mu)^2 \leq a^2) = \frac{\text{var}(X)}{a^2}$$

Randomness-efficient probability amplification

Consider Frievald's algo.

$$\frac{AB}{n \times n} \stackrel{?}{=} C$$

* We use n rand. bits

* $AB \neq C \Rightarrow \Pr(\text{Error}) \leq \frac{1}{2}$

* $AB = C \Rightarrow \Pr(\text{Error}) = 0$

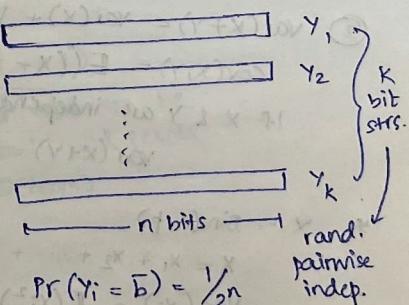
* Repeating the event k times, $\Pr(\text{Error}) \leq \frac{1}{2^k}$
w/ no. of rand. bits used = nk

Alternative:

Repeating k times,

$$\Pr(\text{Error}) \leq \frac{1}{k}, \quad k \leq n$$

No. of rand. bits = $2n$



$$\Pr(Y_i = b) = \frac{1}{2^n}$$

* The construction of κ pairwise indep. bit strings of len. n can be done w/ $\log k$ bits

- For each bit pos., generate $\log k$ purely rand. bits
- Using these bits, gen. κ pairwise indep. rand. bits

$$\Rightarrow \log k \times n$$

~~log k~~

* Here, we can actually do it w/ ~~log k~~ bits

Let y_1, y_2, \dots, y_k be pairwise indep. n -bit strings

$$x_i = \begin{cases} 1, & \text{if } ABY_i \neq CY_i \\ 0, & \text{o.w.} \end{cases}$$

Given $AB \neq C$, $E(x_i) \geq \frac{1}{2}$

$$\text{Consider } X = \sum_{i=1}^k x_i$$

$$E(X) \geq \frac{k}{2}$$

$$\text{Error proba.} = \Pr(X=0)$$

$$\Pr(X=0) \leq \Pr(|X - E(X)| \geq \frac{k}{2}) \leq \frac{\text{var}(X)}{(\frac{k}{2})^2}$$

$$\Pr(X=0) \leq \frac{4}{K^2} \text{ var}(X)$$

$$\text{var}(X) = \sum_{i=1}^k \text{var}(x_i) + 2 \sum_{i < j} \text{cov}(x_i, x_j) \quad (y_i, y_j \text{ are pairwise indep.})$$

$$= \sum_{i=1}^k \text{var}(x_i) \quad x_i = \begin{cases} 1, & \text{w.p. } p \\ 0, & \text{w.p. } 1-p \end{cases}$$

$$\Rightarrow \text{var}(x_i) = p(1-p) \leq \frac{1}{4}$$

$$\Rightarrow \Pr(X=0) \leq \frac{4}{K^2} \cdot (K \times \frac{1}{4}) = \frac{1}{K} \quad \Rightarrow \text{error bound proven}$$

* Consider $\mathbb{Z}_p = \{0, 1, \dots, p-1\}$, $\underbrace{+}_{\text{modulo } p}, \cdot$ ~~1/a finite field~~

If p is prime, the following properties hold:

① associativity

② commutativity

③ additive identity

④ multiplicative identity

⑤ Unique additive inverse ($\forall z \in \mathbb{Z}_p$)

⑥ Unique multiplicative inverse ($\forall z \in \mathbb{Z}_p \setminus \{0\}$) $\rightarrow p$'s primality is reqd. here

A set w/ these operations satisfying the above properties is a field.

\mathbb{Z}_p is a finite field for prime p .

There also exist finite fields \mathbb{Z}_p for non-prime p .

17/8/23

Theorem: Let p be a prime no.

choose $a, b \in_r \mathbb{Z}_p = \{0, 1, \dots, p-1\}$

$$x_i = \underbrace{a \cdot i + b}_{\text{operations}} , i \leq p$$

modulo p → forms a finite field

① x_i is uniformly distributed over \mathbb{Z}_p

② $\{x_i\}$ are pairwise ind.

① i.e., $+c \in \mathbb{Z}_p$

$$\Pr_{a, b \in_r \mathbb{Z}_p} [a \cdot i + b = c] = \frac{1}{p}$$

fix $a \in \mathbb{Z}_p$

$$\Rightarrow b = \underbrace{c - ai}_{\text{additive inverse}}$$

⇒ p of p^2 combos yield c .

② i.e., $+c_1, c_2 \in \mathbb{Z}_p$,

$$\Pr_{a, b \in_r \mathbb{Z}_p} [a \cdot i + b = c_1 \wedge a \cdot j + b = c_2] = \Pr_{a, b \in_r \mathbb{Z}_p} [a \cdot i + b = c_1] \cdot \Pr_{a, b \in_r \mathbb{Z}_p} [a \cdot j + b = c_2]$$

$$= \underbrace{\frac{1}{p}}_{a, b \in_r \mathbb{Z}_p} \cdot \underbrace{\frac{1}{p}}_{a, b \in_r \mathbb{Z}_p}$$

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

yields unique a, b

$$a = \frac{c_1 - c_2}{i - j} \quad \begin{cases} \text{multiplicative} \\ \text{inverse} \end{cases}$$

$$b = \frac{c_{ij} - c_2 i}{i - j}$$

$$\Rightarrow \frac{1}{p^2}$$

Theorem: for every prime p & integer k , ∃ a unique finite field of cardinality p^k , up to isomorphism

[Galois field, GF(p^k)]

construction of GF(\uparrow) [smallest non-prime finite field] ($p=k=2$)

consider polynomials of deg. ≤ 2 over GF(2),
i.e., deg. ≤ 2 & coeffs. from GF(2), i.e., $\{0, 1\}$

We have 8 polynomials:

$$0, 1, x, x+1, x^2, x^2+1, x^2+x, x^2+x+1$$

Irreducible polynomials of deg. 2 over GF(2)

↳ cannot be resolved to a product of lesser polynomials

~~$x^2+x = x(x+1)$~~ not irreducible

x^2+1 has a root of 1, not irreducible

in fact, $(x+1)^2 = x^2+1+(x+x) = x^2+1$

$x^2 + x + 1$ is irreducible

$$GF(4) = \{0, 1, x, x+1\} \quad +, \cdot \text{ performed modulo } \begin{matrix} 1+x+x^2 \\ \text{the irreducible poly.} \\ (\text{of deg. } k) \end{matrix}$$

$x - x \equiv (1+x) \bmod (1+x+x^2)$

$x \cdot (1+x) \equiv 1 \bmod (1+x+x^2)$

$(1+x) \cdot (1+x) \equiv x \bmod (1+x+x^2)$

fact: $\forall l \geq 0$, $x^{2 \cdot 3^l} + x^{3^l} + 1$ is irreducible over $GF(2)$

Chernoff-Hoeffding bounds

$$x_1, x_2, \dots, x_n, \quad x_i = \begin{cases} 1, & \text{w.p. } p_i \\ 0, & \text{w.p. } 1-p_i \end{cases}$$

$$x = \sum_{i=1}^n x_i, \quad \boxed{x_i \text{ are indep.}}, \quad \mu = E[x]$$

$$\textcircled{1} \quad \Pr(|x - \mu| \geq t) \leq e^{-2t^2/n}, \quad t > 0$$

$$\textcircled{2} \quad \Pr(x \geq (1+\delta)\mu) \leq e^{-\delta^2\mu/3} \quad \left. \begin{array}{l} \\ \end{array} \right\} 0 \leq \delta \leq 1$$

$$\textcircled{3} \quad \Pr(x \leq (1-\delta)\mu) \leq e^{-\delta^2\mu/2}$$

Hoeffding's extension:

$$\Pr\left[\left|\frac{1}{n} \sum_{i=1}^n x_i - \mu\right| \geq \epsilon\right] \leq 2e^{-2n\epsilon^2/(b-a)^2}$$

18/8/23
consider a 2-sided error algo.,

$$\text{if } \Pr(\text{error}) \leq \frac{1}{2} - \epsilon$$

To boost the success proba., we

- repeat the algo. k times

- take the majority as the ans.

$$x_1, x_2, \dots, x_k, \quad x_i = \begin{cases} 1, & \text{if } i^{\text{th}} \text{ iteration answers correctly} \\ 0, & \text{o.w.} \end{cases}$$

$$E(x_i) \geq \frac{1}{2} + \epsilon$$

$$x = \sum_{i=1}^k x_i = \text{No. of successful trials}$$

$$E(x) \geq \frac{k}{2} + k\epsilon$$

Proba. of error = $\Pr(x \leq \frac{k}{2})$

$$\Pr(\boxed{E(x) - x} \geq k\epsilon) \leq e^{-2k^2\epsilon^2/k} \quad [\text{Chernoff bound (1)}]$$

When $\epsilon = \frac{1}{n}\epsilon$,
 $k = \frac{1}{\epsilon^2} \log n$ is $\text{poly}(n)$, choosing $k = \frac{1}{\epsilon^2} \log n \leq \frac{1}{\epsilon^2}$

$$\text{Even } k = \frac{n}{\epsilon^2} \text{ is poly}(n),$$

w/ error $\leq e^{-n}$!

Moment generating fn.

$$M_x(t) = E(e^{tx}) = E\left(\sum_{i \geq 0} \frac{t^i x^i}{i!}\right)$$

$$= \sum_{i \geq 0} \frac{t^i}{i!} E(x^i)$$

i-th moment of X

Prf of Chernoff bounds:

$$\Pr(X \geq m) = \Pr\left(\frac{e^{tx}}{e^{tm}} \geq e^{tm}\right), t > 0$$

non-neg. r.v.

$$\leq \frac{E(e^{tx})}{e^{tm}} \quad (\text{Markov's inequality})$$

$$E(e^{tx}) = E\left(e^{t \sum_{i=1}^n x_i}\right) = E\left(\prod_{i=1}^n e^{tx_i}\right)$$

$$= \prod_{i=1}^n E(e^{tx_i}) \quad [\text{indep. r.v.s}]$$

$$E(e^{tx_i}) = p \cdot e^{tp} + (1-p) \cdot 1 = pe^t + q = p(e^{t-1}) + 1$$

$$E(e^{tx}) = (pe^t + q)^n$$

$$\Pr(X \geq m) \leq \frac{(pe^t + q)^n}{e^{tm}}$$

$$\text{Let } m = (p+r)n, \mu = np$$

$$\Pr(X \geq (p+r)n) \leq \frac{(pe^t + q)^n}{e^{t(p+r)n}}$$

calculus ↓

$$\leq \left(\left(\frac{p}{p+r}\right)^{p+r} \left(\frac{q}{q-r}\right)^{q-r}\right)^n$$

strongest version
of the bound
but difficult
to apply

Weakening & simplifying w/ approximations:

$$E(e^{tx_i}) = p(e^{t-1}) + 1 \leq e^{p(e^{t-1})}$$

$$E(e^{tx}) = \prod_{i=1}^n E(e^{tx_i}) \leq e^{n p(e^{t-1})} = e^{\mu(e^{t-1})}$$

$$\Rightarrow \Pr(X \geq (1+\delta)\mu) \leq \frac{e^{\mu(e^{t-1})}}{e^{t\mu(1+\delta)}} \rightarrow \begin{array}{l} \text{Minimized for} \\ t = \ln(1+\delta) \end{array}$$

calculus/
analysis ↓

$$\leq \left(\frac{e^\delta}{(1+\delta)^{1+\delta}}\right)^\mu$$

$$\leq e^{-\delta^2 \mu / 2}, \delta < 1$$

* K servers, n jobs

Allocate jobs to servers in a decentralized way

simple way: Allocate each job to a server uniformly at random.

$$\text{Expected load} = \frac{n}{k} \longrightarrow \text{"balls-in-bins"}$$

21/8/23

Balls & bins

m balls, n bins

Throw the balls into the bins u.a.r.

Coupon collector prob. - what is the val. of m s.t. every bin contains at least one ball?

Hashing - How large should n be compared to m s.t. no bin has a lot of balls?

Birthday problem

What should m be so that the proba. that \exists a bin w/ more than 1 ball is at least $\frac{1}{2}$?

B_i - the event that the i^{th} ball lands in a bin of its own
We want $\Pr\left(\bigcap_{i=1}^m B_i\right)$.

$$\Pr\left(\bigcap_{i=1}^m B_i\right) \rightarrow \text{we want } \Pr(B_i | B_{i-1}, \dots, B_1) = \frac{n-i+1}{n}$$

$$= \prod_{i=1}^m \Pr(B_i | B_{i-1}, \dots, B_1) = \prod_{i=1}^m \frac{n-i+1}{n} = \prod_{i=1}^m \left(1 - \frac{i-1}{n}\right)$$

$$\leq \prod_{i=1}^m e^{-(i-1)/n} = e^{-\sum_{i=1}^m \frac{i-1}{n}} = e^{-m(m-1)/2n} < \frac{1}{2}$$

$$\Rightarrow m = \Theta(\sqrt{n})$$

for hashing, this tells that a hash table of size ~~$\Theta(n^2)$~~

~~$\Theta(n^2)$~~ $\Theta(n^2)$ is likely to find collisions
($\text{prob} > \frac{1}{2}$)

E_i = event that the i^{th} ball lands in ~~one of~~ a separate bin
given that the first $i-1$ balls landed in separate bins

$$\Pr(E_i) = \frac{i-1}{n}$$

$$\Pr\left(\bigcup_{i=1}^m E_i\right) \leq \sum_{i=1}^m \frac{i-1}{n} = \frac{m(m-1)}{2n}$$

If $m < \sqrt{n}$, $\Pr(\text{a ball that doesn't land in a separate bin}) < \frac{1}{2}$

If n balls are thrown into n bins, what is the max. load? logn
loglogn
w.p.
 $1 - \frac{1}{n}$

Thm: If n balls are thrown into n bins w.o.r., then w.p. $\geq 1 - \frac{1}{n}$, the max. load is

at most $\frac{\log n}{\log \log n}$

If: Fix a bin i .

$$\Pr(\text{bin } i \text{ has } k \text{ balls}) \leq \binom{n}{k} \left(\frac{1}{n}\right)^k$$

$$\leq \left(\frac{ne}{k} \cdot \frac{1}{n}\right)^k = \left(\frac{e}{k}\right)^k$$

$$\boxed{\binom{n}{k} \leq \left(\frac{ne}{k}\right)^k}$$

$$\Pr(\text{bin } i \text{ s.t. bin } i \text{ has } k \text{ balls}) \leq n \left(\frac{e}{k}\right)^k \quad (\text{union bound})$$

$$\leq \frac{1}{n} \quad (\text{reqd.})$$

what val. of k do we choose?

We can substitute $\frac{\log n}{\log \log n}$ to verify it works.

x_1, x_2, \dots, x_n - r.v.s corresponding to no. of balls in individual bins dependent, their sum is fixed \rightarrow something indep. would be nicer to work with

$$x_i \sim \text{Bin}(m, 1/n)$$

$$\begin{aligned} \Pr(x_i = k) &= \binom{m}{k} \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{m-k} \\ &= \frac{m!}{k!(m-k)!} \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{m-k} \\ &= \frac{(m-k+1) \cdots (m)}{k! n^k} \left(1 - \frac{1}{n}\right)^{m-k} \\ &\approx \left(\frac{m}{n}\right)^k \left(\frac{1}{k!}\right) e^{-m/n} \quad (\text{assuming } m, n \gg k) \end{aligned}$$

Poisson dist.

$$x \sim \text{Poi}(\lambda)$$

$$\Pr(x = k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

$$E(x) = \lambda$$

$$\lambda = m/n$$

Thm: If $X \sim \text{Bin}(n, p)$ & $\lim_{n \rightarrow \infty} np = \lambda$ (indep. of n),

$$\text{for every fixed } k, \lim_{n \rightarrow \infty} \Pr(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

23/8/22
 $X_1 \sim \text{Poi}(\lambda_1)$, $X_2 \sim \text{Poi}(\lambda_2)$ [indep.]

then $X_1 + X_2 \sim \text{Poi}(\lambda_1 + \lambda_2)$

$$\text{Pr}(X_1 + X_2 = j) = \sum_{k=0}^j \text{Pr}(X_1 = k) \cdot \text{Pr}(X_2 = j - k)$$

$$= \sum_{k=0}^j \frac{e^{-\lambda_1} \cdot \lambda_1^k}{k!} \cdot \frac{e^{-\lambda_2} \cdot \lambda_2^{j-k}}{(j-k)!}$$

$$= \frac{e^{-(\lambda_1 + \lambda_2)}}{j!} \sum_{k=0}^j \underbrace{\lambda_1^k \lambda_2^{j-k}}_{(\lambda_1 + \lambda_2)^j} \underbrace{\frac{j!}{k!(j-k)!}}$$

$$= \frac{e^{-(\lambda_1 + \lambda_2)} (\lambda_1 + \lambda_2)^j}{j!}$$

Chernoff bounds for Poisson r.v.s

$X \sim \text{Poi}(\lambda)$

$$E(e^{tX}) = \sum_{k \geq 0} \text{Pr}(X=k) \cdot e^{tk} = \sum_{k \geq 0} \frac{e^{-\lambda} \cdot \lambda^k}{k!} \cdot e^{tk}$$

$$= e^{-\lambda} \sum_{k \geq 0} \frac{e^{tk} \cdot (e^\lambda)^k}{k!} \cdot e^{-\lambda t} = e^{\lambda(t-1)}$$

$$= e^{\lambda(t-1)} \sum_{k \geq 0} \frac{e^{-\lambda t} \cdot (e^\lambda)^k}{k!}$$

$$= e^{\lambda(t-1)}$$

Let $x \geq \lambda$.

$$\Pr(X \geq x) = \Pr(e^{tX} \geq e^{tx})$$

$$\leq \mathbb{E} \frac{e^{tx}}{e^{tx}}$$

$$\leq \frac{e^{-\lambda} (e^\lambda)^x}{x^n}$$

Find t that minimizes this,
 $t = \log x/\lambda$

Let $x \leq \lambda$

$$\Pr(X \leq x) = \Pr(e^{tX} \geq e^{tx}), \quad t < 0$$

$$\leq \mathbb{E} \frac{e^{-\lambda} (e^\lambda)^x}{x^n}$$

Thm: Let $\{x_i^m\}_{1 \leq i \leq n}$ be the balls-in-bins dbn. w/ m balls & n bins

Let $\{y_i^m\}_{1 \leq i \leq n}$ be indep. Pois. rvs $\sim \text{Poi}(m/n)$

The dbn. $\{y_i^m\}_{1 \leq i \leq n}$ conditioned on $\sum_{i=1}^n y_i^m = k$

is identical to $\{x_i^k\}_{1 \leq i \leq n}$

Pf: $\Pr\left(\bigcap_{i=1}^n (x_i = k_i)\right)$ [multinomial dbn.]
 $\quad \quad \quad \left(\sum_{i=1}^n k_i = k\right)$

$$= \binom{k}{k_1} \cdot \binom{k-k_1}{k_2} \cdots \binom{k - k_1 - k_2 - \cdots - k_{n-1}}{k_n} \cdot \left(\frac{1}{n}\right)^{k_1} \cdot \left(\frac{1}{n}\right)^{k_2} \cdots \left(\frac{1}{n}\right)^{k_n}$$

$$= \frac{k!}{k_1! k_2! \cdots k_n!} \left(\frac{1}{n}\right)^k$$

$$\Pr\left(\bigcap_{i=1}^n (y_i^m = k_i) \mid \sum_{i=1}^n y_i^m = k\right) = \frac{\Pr\left[\bigcap_{i=1}^n (y_i^m = k_i)\right] \cdot \Pr\left[\bigcap_{i=1}^n (y_i^m = k_i) \mid \sum_{i=1}^n y_i^m = k\right]}{\Pr\left(\sum_{i=1}^n y_i^m = k\right)}$$

$\xrightarrow{\text{indep.} \Rightarrow \text{prod.}}$

Thm: $\{x_i^m\}_{1 \leq i \leq n}$ balls-m-bins dbn.
 w/ m balls & n bins

$\{y_i^m\}_{1 \leq i \leq n}$, indep. Pois. (m/n)

E Let f be a non-negative fn.

$$E(f(x_1^m, x_2^m, \dots, x_n^m))$$

$$\leq e^{m\lambda} E(f(y_1^m, y_2^m, \dots, y_n^m))$$

How to apply this? no conditioning!

w/ what proba.
 suppose we want to answer,

no. of bins with no balls is $\geq k$

$$\text{Define } f(x_1, \dots, x_n) = \begin{cases} 1, & \text{if } |\{j \mid x_j = 0\}| \geq k \\ 0, & \text{o.w.} \end{cases}$$

$$E(f(x_1^m, \dots, x_n^m)) = P(\text{no. of bins w/ no balls is } \geq k)$$

2018/23

Pf: $E(f(y_1^m, y_2^m, \dots, y_n^m)) = \sum_{k \geq 0} E(f(y_1^m, \dots, y_n^m) \mid \sum_{i=1}^n y_i^m = k) \cdot \Pr\left(\sum_{i=1}^n y_i^m = k\right)$

$$\geq E(f(y_1^m, \dots, y_n^m) \mid \sum_{i=1}^n y_i^m = m) \cdot \Pr\left(\sum_{i=1}^n y_i^m = m\right)$$

$$= E(f(y_1^m, \dots, y_n^m)) \cdot \frac{e^{-m} m^m}{m!}$$

$$E(f(Y_1^m, \dots, Y_n^m)) \geq E(f(X_1^m, \dots, X_n^m)) \cdot \frac{e^{-\frac{m}{2}} \cdot \frac{m^m}{m!}}{\sqrt{2\pi m}} \cdot \left(\frac{e^m}{m}\right)^m$$

$$E(f(X_1^m, \dots, X_n^m)) \leq e^{\sqrt{m}} E(f(Y_1^m, \dots, Y_n^m))$$

[$e \leq \sqrt{2\pi}$]

Stirling's approx.
 $m! \leq \sqrt{2\pi m} \left(\frac{m}{e}\right)^m$

Thm: If n balls are thrown into n bins w.o.r.y.,
then \exists a bin containing $\approx \frac{\log n}{\log \log n}$ balls

$$n \cdot p \geq 1 - \frac{1}{n}$$

Pf:

$$f(x_1, \dots, x_n) = \begin{cases} 1, & \text{if } x_i \leq k \quad \forall i \in \{1, \dots, n\} \\ 0, & \text{o.w.} \end{cases} \quad \Rightarrow \max\{x_1, \dots, x_n\} \leq k$$

$$E(f(X_1^m, \dots, X_n^m)) = \Pr(\text{Every bin has } \leq k \text{ balls})$$

$$\Pr\left(\bigcap_{i=1}^n (Y_i^m \leq k)\right) = (\Pr(Y_i^m \leq k))^n = (1 - \Pr(Y_i^m \geq k))^n$$

$$\Pr(Y_i^m \geq k) \geq \Pr(Y_i^m = k) = \frac{e^{-1} \cdot \frac{1^k}{k!}}{k!} = \frac{1}{ek!} \quad \hookrightarrow \text{Poi}(1)$$

$$\Pr\left(\bigcap_{i=1}^n (Y_i^m \leq k)\right) = \left(1 - \frac{1}{ek!}\right)^n \leq e^{-\frac{n}{ek!}}$$

$$\Pr(\text{Every bin has } \leq k \text{ balls}) \leq e^{\sqrt{n}} \cdot e^{-n/ek!} \leq \frac{1}{n} \quad (\text{reqd.})$$

Verify that $k = \frac{\log n}{\log \log n}$ works

Static dictionary

Universe U , set $S \subseteq U$

S is static, i.e., known at the start & fixed

$$m = |S| \leq |U| = M \quad \hookrightarrow O(1s)$$

We wish to compactly represent S so that

membership queries ($x \in S$) can be

answered efficiently

$\hookrightarrow O(1)$ [expected or worst-case]

$$\log\left(\frac{e^y - 1}{e^y}\right) - y$$

$$+\frac{y}{e^y - 1} \log(e^y - 1) - \frac{ye^y}{e^y - 1} \quad \log(e^y - 1)(e^y - 1) = ye^y \quad e^y = 1/2 \quad y = -\log 2$$

Bloom filter

Bit-array $A[1, 2, \dots, n]$ (i.e., size is n)

choose k funs. h_1, \dots, h_k u.a.r. [assumes perfectly rand. hash funs.]

$\forall x \in S$, set $A[h_i(x)] = 1 \quad \forall i \in \{1, \dots, k\}$

i.e., $h_j(x) = i$

w.p. n^{-1}

Membership query: Given x , check $A[h_i(x)] \quad \forall i \in \{1, \dots, k\}$

resolution

If all 1's, say $x \in S$

$x \in S \Rightarrow$ correct result guaranteed

$x \notin S \Rightarrow$ false positives may occur

Given $x \notin S$,

$$\Pr[\text{we answer } x \in S] \leq \delta$$

$$\Pr[A[j] = 0] = \underbrace{\left(1 - \frac{1}{n}\right)^{mk}}_{\text{For given } x \notin S} \stackrel{\text{Overall } h_i \text{'s}}{=} p \quad \& \quad x \in \underbrace{m}_{\text{bins}}$$

$\approx e^{-mk/n}$

$$\Pr(Y_i \sim \text{Poi}(mk/n)) = 0$$

$A[j] \& A[j']$
are dependent
(balls-in-bins
w/ km balls
& n bins)

What should be n, k so that proba. for false positives is $\leq \delta$?

Expected no. of empty posns. in the Bloom filter = np

Given $x \notin S$, assuming \uparrow proba. of bin being empty,

$$\Pr(\text{Answer Yes}) = (1-p)^k \quad \xrightarrow{\text{assumes independence}}$$

$$Y_i \sim \text{Poi}(km/n), \quad \Pr(Y_i = 0) = e^{-km/n}$$

$$Z_i = \begin{cases} 1, & \text{if } Y_i = 0 \\ 0, & \text{o.w.} \end{cases} \quad Z_i \sim \text{Ber}(e^{-km/n})$$

$$Z = \sum_{i=1}^k Z_i = \text{No. of } \cancel{\text{posns.}} \text{ in } A \text{ that are } 0$$

$$\Pr[|Z - E(Z)| \geq \epsilon n] \leq 2e^{-ne^2e^{-km/n}/3} \quad [\delta = \epsilon n/\mu]$$

With $f(x)$ as $|x - E(z)| \geq \epsilon n \Rightarrow 1, 0 \text{ o.w.}$

$$\Pr[|X - E(z)| \geq \epsilon n] \leq 2e^{-ne^{km/n}} e^{-ne^{km/n}/3}$$

$$\Pr(\text{false positives}) \approx \left(1 - e^{-km/n}\right)^k \xrightarrow{\min. \text{ for } k=(\ln 2)n} \frac{k \log(1 - e^{-km/n})}{k} + K \left(e^{-km/n}\right) \cdot \frac{n}{1 - e^{-km/n}}$$

$$K = \log(1/\delta) \Rightarrow n = m \log(1/\delta)/k$$

Universal hash families (Carter-Wegman)

κ -universal hash family:

Family of fns $H = \{h: U \rightarrow [n]\}$ s.t.

$\forall x_1, x_2, \dots, x_k$ (all different),

$$\Pr_{h \in H} [h(x_1) = h(x_2) = \dots = h(x_k)] \leq \frac{1}{n^{k-1}}$$

$h \in H$

If $H = \text{set of all fns. from } U \text{ to } [n]$, $(|H| = n^{|U|})$

$$\Pr_{h \in H} [h(x_1) = \dots = h(x_k)] = \frac{1}{n^{k-1}} \quad \forall k \quad (\text{ideally, we want } |H| = \text{poly}(|U|))$$

κ -wise indep. hash family: [strongly κ -universal]

$H = \{h: U \rightarrow [n]\}$ s.t.

$\forall x_1, x_2, \dots, x_k$ (distinct), $\forall y_1, y_2, \dots, y_k$,

$$y_i \in \{1, \dots, n\}$$

$$\Pr_{h \in H} [\bigwedge_{i=1}^k (h(x_i) = y_i)] = \frac{1}{n^k} \quad \text{strongly } \kappa\text{-universal}$$

\Downarrow

$\forall x_1, x_2, \dots, x_k$ (distinct),

$$\Pr_{h \in H} [h(x_1) = h(x_2) = \dots = h(x_k)] = \frac{1}{n^{k-1}} \quad \text{---} \quad \kappa\text{-universal}$$

strongly κ -universal = uniformity, i.e.,

$$\textcircled{1} \quad \Pr_{h \in H} [h(x) = y] = \frac{1}{n} \quad \&$$

$$\textcircled{2} \quad \Pr_{h \in H} \left[\bigwedge_{i=1}^k h(x_i) = y_i \right] = \prod_{i=1}^k \Pr_{h \in H} [h(x_i) = y_i]$$

* m balls into n bins using a 2-universal hash family

$H = \{h: [m] \rightarrow [n]\}$

i.e., choose $h \in H$ & assign balls to bins using h

for any 2 balls $i \neq j$, let

$$x_{ij} = \begin{cases} 1, & \text{if } i \neq j \text{ land in the same bin;} \\ 0, & \text{o.w.} \end{cases}$$

$$\Rightarrow \Pr_{h \in H} [x_{ij} = 1] = \Pr_{h \in H} [h(i) = h(j)] \leq \frac{1}{n}$$

Total no. of collisions, $X = \sum_{i < j} X_{ij}$

$$E(X) \leq \binom{m}{2} \cdot \frac{1}{n} \leq \frac{m^2}{2n}$$

Let y be the max. load

$X \geq \binom{y}{2} = \text{no. of collisions in the bin w/ max. load}$

$$\Pr(X \geq \frac{m^2}{n}) \leq \frac{1}{2} \quad (\text{Markov's ineq.})$$

$$\Pr(\binom{y}{2} \geq \frac{m^2}{n}) \leq \frac{1}{2} \quad (\text{since } \binom{y}{2} \leq X)$$

$$\Pr(Y \geq 1 + \sqrt{\frac{2}{n}}) \leq \frac{1}{2} \quad \left[\binom{y}{2} \geq \frac{(Y-1)^2}{2} \right]$$

If $n = \Theta(m^2)$, then every bin has $\Theta(1)$ balls w.p. $\geq \frac{1}{2}$

for a static dict. on $S \subseteq U$, $|S| = m$

we can obtain a good hash fn. of size $\Theta(m^2)$

- choose a hash fn from a 2-universal hash family
- if no. of collisions is too high, resample

[Expected number of times to sample is 2 from our bound]

- * Dynamic dictionary over U
 - additions (using H , a 2-universal hash family)
 - deletions
 - query

[we want $n = O(m)$ ideally
 \Rightarrow FKS (Fredman-Komlós-Szemeredi)]

Set of all additions, $S \subseteq U$

$x \in U, h \in H$

No. of collisions of $|x|$ elements in S using hash fn. h ,

$$\text{coll}(x, S, h) = |\{y \in S \mid h(x) = h(y)\}|$$

$$E[\text{coll}(x, S, h)] = \frac{1}{|H|} \sum_{h \in H} \text{coll}(x, S, h)$$

$$= \frac{1}{|H|} \sum_{\text{yes}} \sum_{h \in H} \underbrace{\mathbb{1}_{[h(x)=h(y)]}}_{\text{indicator fn.}}$$

$$= \sum_{\text{yes}} \frac{1}{|H|} \sum_{h \in H} \underbrace{\Pr[h(x)=h(y)]}_{\Pr[h(x)=h(y)]} \quad [n = |S|]$$

$$\leq \sum_{\text{yes}} \frac{1}{n} = \frac{|S|}{n} \quad [n \rightarrow \text{size of hash table}]$$

[if $n = \Theta(|S|)$, then $O(1)$ running time for a query, but hash table could be sparse at times]

Explicit 2-universal family

$\mathcal{U}, |\mathcal{U}| = m \rightarrow n$ sized table

Prime $p > m \Rightarrow p = \Theta(m)$ [if prime $p \leq m$ & $2m$]

$$\mathcal{H} = \{h_{a,b} \mid 1 \leq a \leq p-1, 0 \leq b \leq p-1\}$$

$$h_{a,b}(x) = ((ax+b) \bmod p) \bmod n$$

$$|\mathcal{H}| = \Theta(p^2)$$

$$\Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq \frac{1}{n}$$

$$ax + b \equiv z_1 \pmod{p}$$

$$ay + b \equiv z_2 \pmod{p}$$

$$z_1 \equiv z_2 \pmod{n}$$

Then: \mathcal{H} is 2-universal

Pf: To show: $\Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq \frac{1}{n}$

$$ax + b \equiv z_1 \pmod{p}$$

$$ay + b \equiv z_2 \pmod{p}$$

$$z_1 \equiv z_2 \pmod{n}$$

For every $z_1 \neq z_2$, \exists unique soln. for a, b

$$|\{z_2 \mid z_2 \neq z_1, z_1 \equiv z_2 \pmod{n}\}| \leq \frac{p}{n} - 1$$

$$\Rightarrow \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq \frac{p(\frac{p}{n} - 1)}{p(p-1)} = \frac{1}{n} \cdot \frac{(p-n)}{(p-1)}$$

$$\leq \frac{1}{n}$$

Pairwise-indep. hash family

$$|\mathcal{U}| = 2^m \rightarrow 2^n$$

Consider the field $\text{GF}(2^m)$

$$\mathcal{H} = \{h_{a,b} \mid a, b \in \text{GF}(2^m)\}$$

$$h_{a,b}(x) = (ax + b) \text{ in } \text{GF}(2^m)$$

The "pairwise" comes from here
 $\text{GF}(2^m) \Rightarrow m$ -wise indep.

m -bit string \Rightarrow truncate to n bits

Exercise
Verify this is pairwise indep.

Perfect hash family

[Static dict. again]

A family $H = \{h: U \rightarrow [n]\}$ is perfect for sets of size $\leq n$
 if $\forall S \subseteq U, |S| \leq n$,
 $\exists h \in H \text{ s.t. } \forall x \neq y, h(x) \neq h(y)$

[Any n -sized set
in U can be
perfectly hashed
by a fn. in H]

\Rightarrow static dictionary has an $O(n)$ data structure
with $O(1)$ query time

\hookrightarrow Is $h[n(x)] = x$?

Each cell of the table = $O(\log |U|)$

We also want to store each h in a constant
no. of cells & to index H , we need $\log |H|$ bits

$$\Rightarrow \log |H| = O(1) \cdot \log |U| \Rightarrow |H| = |U|^{O(1)}$$

* If $H = \{h: U \rightarrow [n]\}$ is a perfect hash family for sets
of size n , then

$$|H| = 2^{\Omega(n)}$$

$$|U| = n^2 \rightarrow [n]$$

cannot be perfectly hash

FKS hashing (Fredman - Komlos - Szemerédi, '84)

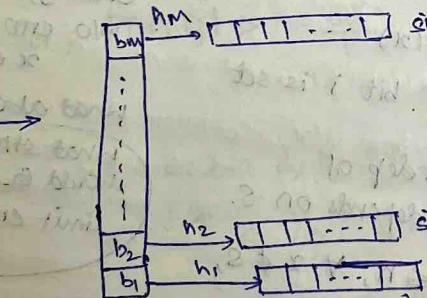
[Cell probe model]

$$U \rightarrow [n], \text{ sets } S, |S| = m$$

choose $h \in H$ (2-universal)

s.t. no. of collisions $\leq m$

Handle collisions using another level of hashing
(2-universal again)



$$\text{Size} = m + \sum_{i=1}^m b_i^2$$

$$= O(m)$$

(Further collisions are
chained, since only
constant collisions now)

$$\sum_{i=1}^m \binom{b_i}{2} \leq m^2$$

$$(m \text{ (no. of collisions)})$$

Bit probe model

Size, query are now in terms of bits rather than cells

FKS: $m \log |U|$ $\log |U|$

What should be n if we

limit ourself to 1 bit query?

Buhman, Miller, Radhakrishnan, Venkatesh, '00
size = $O\left(\frac{m}{\epsilon} \log |U|\right)$ → comparable to FKS result

After seeing
of nice bits
of results
1-bit query w/ 2-sided error of $\leq \epsilon$

If we restrict 1-sided error of $\leq \epsilon$,
size = $O\left(\frac{m^2}{\epsilon^2} \log |U|\right)$

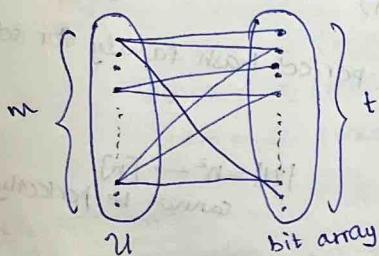
↳ can show existence but can't obtain the data structure
if we want to efficiently find the structure,

$$\text{size} = O\left(\frac{m^2}{\epsilon^2} \log^2 |U|\right)$$

3/18/23

change of notation: $|U| = m$, $|S| = n$
for the proof

1-bit query algo. → bipartite graph



? Randomly
 $x \in S$: choose $i \in [t]$
using the graph
say,
[edge wts. like
prob. of choosing]

Given $G(U, [t], E)$,

$$S \subseteq U, |S| \leq n,$$

$$\text{define } \forall x, N(x) = \{i \mid (x, i) \in E\}$$

set all the bits in $\bigvee_{x \in S} N(x)$

query: $x \in S$

- choose $i \in N(x)$

- Answer "YES" if bit i is set

Note that the graph is indep. of
S. Only the data structure depends on S.

We want: $\forall S \subseteq U, |S| \leq n, \forall x \notin S,$

$$\text{for } N(S) = \bigcup_{y \in S} N(y),$$

$$|N(x) \cap N(S)| \leq \epsilon |N(x)|$$

clearly, this works
w/o error when
 $x \in S$

what about $x \notin S$?
what structure
should G have to
limit error $\leq \epsilon$?

we show a stronger result:

$$\begin{cases} \forall x, |N(x)| = d & [\text{d-regular over } U] \end{cases}$$

$$\begin{cases} \forall x \neq y, |N(x) \cap N(y)| \leq \frac{\epsilon d}{n} & [\text{this implies the above}] \end{cases}$$

such
graphs
exist

An (m, t, d, r) -combinatorial design is a family of sets

T_1, T_2, \dots, T_m s.t. $\forall i \neq j, i \neq j$

① $\forall i, T_i \subseteq [t]$

In our case,

② $|T_i| = d$

$$T_i = N(x_i), x_1, x_2, \dots, x_m \in \mathbb{R}^n$$

③ $|T_i \cap T_j| \leq r$

$$r = \frac{ed}{n}$$

We consider a random collection of sets & show that w.p. > 0 , the properties are satisfied, implying existence.

Pf: $\forall j \in [m], \forall i \in [t]$,
put i in T_j w.p. $\frac{2d}{t}$ independently

$$\Rightarrow E(|T_j|) = 2d$$

$$\Pr(|T_j| < d) \leq e^{-2d\left(\frac{d}{t}\right)^2/2} = e^{-d/4} \quad [\text{Chernoff bound}]$$

$$\Rightarrow \Pr(\exists j |T_j| < d) \leq m e^{-d/4} \quad [\text{Union bound}]$$

< 1 [we want this, otherwise
 \geq no set can have size d]

$$E(|T_i \cap T_j|) = t \left(\frac{2d}{t}\right)^2$$

$$\left[\sum_{k \in [t]} E(|T_i \cap T_j| \mid k \in T_i) \cdot \Pr(k \in T_i) \right] \quad [\text{check}]$$

$$= \frac{4d^2}{t}$$

$$\Pr(|T_i \cap T_j| > \frac{8d^2}{t}) \leq e^{-\frac{4d^2}{st}} \quad [\text{Chernoff bound}].$$

$$\Pr(\exists i, j |T_i \cap T_j| > r = \frac{8d^2}{t}) \leq \binom{m}{2} e^{-\frac{4d^2}{st}}$$

$$\leq m^2 e^{-\frac{4d^2}{st}}$$

We find suitable d, t s.t.

$$m e^{-d/4} \leq \frac{1}{10} \quad \& m^2 e^{-\frac{4d^2}{st}} \leq \frac{1}{10}$$

Then, both probas $\leq \frac{1}{10} \Rightarrow \Pr(\text{Not satisfied}) \leq \frac{1}{5}$

$$\Rightarrow \Pr(\text{satisfied}) \geq \frac{4}{5}$$

$$r = \frac{ed}{n} = \frac{8d^2}{t} \Rightarrow t = \frac{8dn}{e}$$

This is satisfied when

$$d = \Theta\left(\frac{n}{e} \log m\right) \& t = \underbrace{\Theta\left(\frac{n^2}{e^2} \log m\right)}_{\text{Hence proved}}$$

$$r = \Theta(\log m)$$

19/12

PS1 discussion

① Generate n distinct nos. & sort them

Assign each i the $\underbrace{k_1, k_2, \dots, k_n}$
sorted posn. of k_i

choose using $\log n$ bits each

$$\Rightarrow \Pr(k_i = k_j) = \frac{1}{n^2}$$

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

$$\Rightarrow \Pr(j \neq i \mid k_i = k_j) \leq \frac{1}{n}$$

② Intuitively, X increments every 2^x times, so in a sense, it appears to be counting the no. of bits in the representation of n .

Exercise: Find Variance, $\text{Var}(Y)$ & show

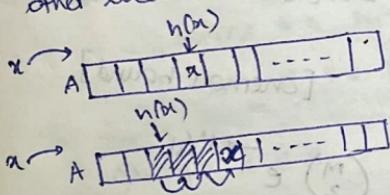
(Morris counter)
L way to count n
times using $\log n$ bits

n can be rep. as $(1 \pm \epsilon)n$
using $\frac{1}{\epsilon^2} \log n$ bits

19/12

Open addressing (Linear probing)

our previous hashing schemes (such as FKS) handled collisions using an auxiliary data structure. This scheme maps collisions into some other location of the primary data structure.



If $A[h(x)]$ is empty

If $A[h(x)]$ contains some other elements
probe linearly for an empty position

Addition: Insertion: Starting from $h(x)$, insert x at the first empty slot

Membership: Starting from $h(x)$, continue until you find x or an empty slot

[doesn't work when there are deletions, in which case we set flags/tombstones at posns. of deleted elements]

works very well

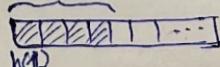
in practice.

Knuth '63: If h is completely random, then expected time is $O(1)$.

2010-11: 5-wise independence is necessary & sufficient for expected $O(1)$ time.

"Run": A run is a maximal contiguous seq. of cells that are filled

Given $i \in S$, what is the length of the run containing $h(i)$?
say, 1



Consider a k -wise indep. hash family H , h.e. H
 for the given \mathcal{I} , consider the interval of pos'g in the hash table
 $I_j = [h(\mathcal{I}) - (2^j - 1), h(\mathcal{I}) + (2^j - 1)]$ (dyadic intervals)

$|\mathcal{I}| = n$, hash table has size t

Expected no. of elements in $I_j = \frac{|I_j| \cdot n}{t}$

$$X_i = \begin{cases} 1, & \text{if } i \in I_j \\ 0, & \text{o.w.} \end{cases}$$

$$X = \sum_{i \in \mathcal{I} \setminus \{\mathcal{I}\}} X_i \Rightarrow E(X) = \frac{n|I_j|}{t}$$

$$\text{If } t = 8n, |I_j| = 2^{j+1} - 1 \Rightarrow E(X) \leq 2^{j-2}$$

Consider R , a run containing \mathcal{I}

$$E(|R|) = \sum_{l=0}^n l \cdot \Pr(|R|=l)$$

$$\leq \sum_{j=0}^{\log t} 2^j \cdot \Pr\left(\underbrace{2^{j-1} \leq |R| \leq 2^j}_{\text{For these, } l \leq 2^j}\right)$$

If $|R| \geq 2^{j-1}$, no. of elements in I_j is $\geq 2^{j-1}$

① Assume h is random $\Rightarrow X$ is a sum of Bernoulli r.v.s

$$\Pr(|R| \geq 2^{j-1}) \leq \Pr(X \geq \underbrace{2^{j-1}}_{2E(X)})$$

$$\leq e^{-2^{j-2}/2}$$

$$\Rightarrow E(|R|) \leq \sum_{j=0}^{\log t} \frac{2^j}{e^{2^j}} = O(1)$$

② Suppose h is picked from a 3-wise indep. hash family

Given $h(\mathcal{I})$, X_i 's are pairwise indep.

$$\Rightarrow \text{var}(X) = \sum_{i \in \mathcal{I} \setminus \{\mathcal{I}\}} \text{var}(X_i) = \sum_{i \in \mathcal{I} \setminus \{\mathcal{I}\}} \underbrace{E(X_i^2)}_{= E(X_i)^2} - [E(X_i)]^2 \quad (\text{Bernoulli})$$

$$\leq \sum_{i \in \mathcal{I} \setminus \{\mathcal{I}\}} E(X_i)$$

$$\Pr(|R| \geq 2^{j-1}) \leq \Pr(X \geq \underbrace{2^{j-1}}_{2\text{var}(X)})$$

$$\leq \frac{1}{\text{var}(X)}$$

$$\leq \frac{1}{2^{j-2}}$$

$$\Rightarrow E(|R|) \leq \sum_{j=0}^{\log t} \frac{2^j}{2^{j-2}}$$

$$= O(\log n)$$

Fourth-moment bound

Let X_1, X_2, \dots, X_n be t -wise indep. st. ~~Exp~~
 $\text{E}(X_i) = p = \mu$ $X_i \sim \text{Ber}(p)$

$$X = \sum_{i=1}^n X_i \Rightarrow \text{E}(X) = np = \mu$$

For $\mu \geq 1, \beta > \mu$

$$\Pr(X > \mu + \beta) \leq \frac{\beta^2}{\beta^t}$$

(3) n is picked from a 5-wise indep. family

\Rightarrow Given $n(1), X_i$'s are t -wise indep

$$\Pr(|X| \leq 2^{j-1}) \leq \Pr(X \geq 2^{j-1})$$

$$= \Pr(X \geq 2^{j-2} + 2^{j-2})$$

$$\leq \frac{4}{\mu^2} = \frac{32}{2^{2j}}$$

$$\mathbb{E}(|X|) \leq \sum_{j=0}^{\log t} \frac{2^j \cdot 32}{2^{2j}} = 32 \sum_{j=0}^{\log t} \frac{1}{2^j} \leftarrow \Theta(1)$$

6/9/23

pf. of fourth-moment bound:

$$\Pr(X \geq \mu + \beta) = \Pr(X - \mu \geq \beta) = \Pr((X - \mu)^+ \geq \beta^+)$$

$$\leq \frac{\mathbb{E}[(X - \mu)^+]}{\beta^+}$$

Define $Y_i = X_i - p$ Since X_i 's are t -wise indep., so are Y_i 's

$$\mathbb{E}(Y_i) = 0$$

$$\mathbb{E}[(X - \mu)^+] = \mathbb{E}\left[\left(\sum_{i=1}^n (X_i - p)\right)^+\right] = \mathbb{E}\left[\left(\sum_{i=1}^n Y_i\right)^+\right]$$

$$= \mathbb{E}\left[\sum_{i=1}^n Y_i^+ + \binom{t}{2} \leq Y_i^2 Y_j^2 + \binom{t}{3} \leq Y_i Y_j^2 \quad i \neq j\right]$$

$$+ \sum_{i+j+k+l} Y_i Y_j Y_k Y_l + \sum_{i+j+k+l} Y_i Y_j Y_k^2 \quad \text{other split}$$

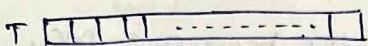
$$= \sum_{i=1}^n \mathbb{E}(Y_i^+) + \sum_{i \neq j} \mathbb{E}(Y_i^2) \geq \mathbb{E}(Y_i^2) \quad \begin{cases} \text{due to } t\text{-wise} \\ \text{indep.} \end{cases}$$

$$\mathbb{E}(Y_i^+) = p(1-p)^+ + (1-p)(0-p)^+$$

$$\mathbb{E}(Y_i^2) = p(1-p)^2 + (1-p)p^2$$

$$\mathbb{E}[(X - \mu)^4] \leq 4np^2 = 4\mu^2 \quad \text{(verify)} \quad \text{Hence, proved.}$$

Cuckoo hashing (Pagh, Rödler - '01)



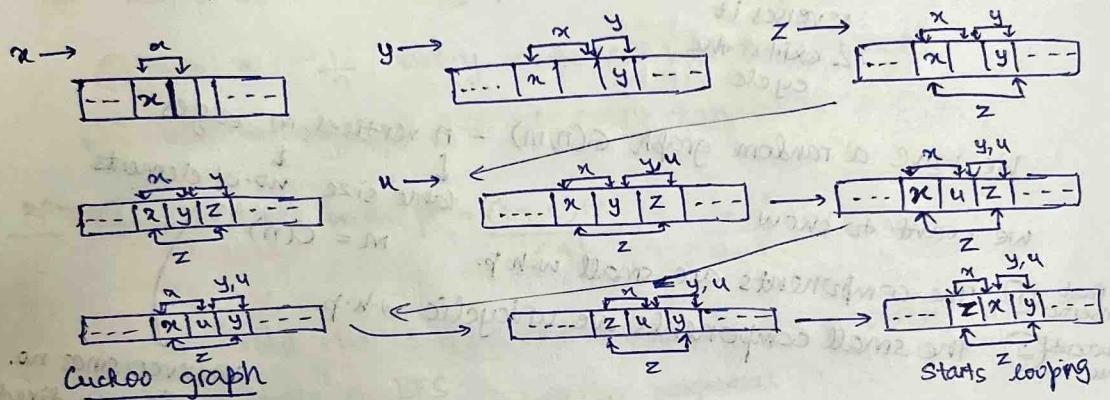
2 hash fn's. h_1, h_2 (completely random)

- * Every x is hashed to either $h_1(x)$ or $h_2(x)$
- ⇒ search for x in $O(1)$ -worst case

Insertions are slightly harder:

Insert (x):

- If one of $h_1(x)$, $h_2(x)$ is empty, then insert in T at the empty loc!
- If both $h_1(x)$ & $h_2(x)$ are occupied,
 - then $h_i(y) = h_1$ for some $y \neq x$, $i \in \{1, 2\}$,
 - place x in $h_1(x)$ & recursively try to insert y
- If this continues indefinitely, choose new h_1, h_2' & rehash the table

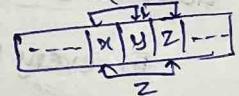


$G(V, E)$

$V \rightarrow$ set of pos's. in the hash table

+ x in the table, add edge $(h_1(x), h_2(x))$

- self loops & 1st edges are possible



=

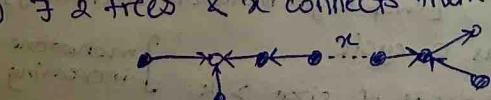
Obs: Let C be a connected component in the cuckoo graph, and x is being inserted. If the new component has > 1 cycle, then insertion is impossible. [without rehashing]

Lemma: If after insertion, the component is unicyclic, then insertion is possible in time $O(|C|)$

↳ new size of the component

Possible cases:

- ① $\exists 2$ trees & x connects them



Results in a DAG, which must have a sink \Rightarrow terminates

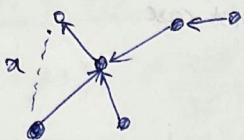
Directions show potential next pos's. for hashed elements

Each vertex has out-degree ≤ 1

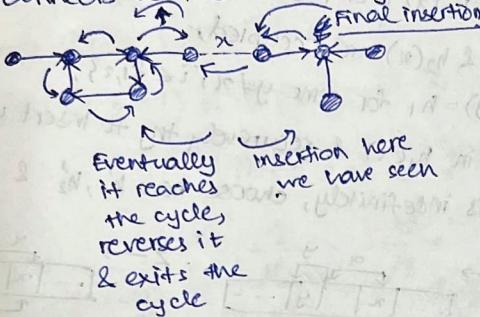
② x lies within a component C

$\Rightarrow C$ must be a tree

Once again, once again, the tree becomes a DAG,
& the sink node gets filled



③ x connects a tree & a ~~cycle~~



We have a random graph $G(n, m)$ - n vertices, m edges

we want to show:

table size \downarrow
no. of elements

$$m = O(n)$$

fast insertion ① The components are small w.h.p.

successful insertion ② The small components are unicyclic w.h.p.

Erdős-Renyi random graph: $G(n, p)$

For every pair of vertices (v_1, v_2) ,

add edge (v_1, v_2) w.p. p indep.

$$\Rightarrow E(\text{edges}) = \binom{n}{2}p = m, \text{ when } p = \frac{m}{\binom{n}{2}}$$

Good h.p. bound too,
since it is binomial r.v.

Instead,
we analyse

However, since no.
of edges is fixed,
existence of edges
is dependent,
making analysis
difficult

Monotone graph property

A property P is monotone increasing if

$G \in P \wedge G \subseteq G' \Rightarrow G' \in P$

G has the
property

of G'

[swap to get monotone decreasing]

(E.g.) P_1 - graph has a cycle

? monotone
increasing

* - If a component of size $> k$

- \mathcal{P} - max. indep-set is of size $\geq k$
- bipartiteness
 - planarity
- $\left. \begin{array}{l} \text{monotone} \\ \text{decreasing} \end{array} \right\}$

Lemma: Let \mathcal{P} be any monotone increasing property.

$$\text{Denote by } P(n, m) = \Pr_{G \sim G(n, m)}(G \in \mathcal{P})$$

$$\& P(n, p) = \Pr_{G \sim G(n, p)}(G \in \mathcal{P})$$

$$\text{If } p^+ = \frac{(1+\epsilon)m}{\binom{n}{2}}, \quad p^- = \frac{(1-\epsilon)m}{\binom{n}{2}}, \text{ then}$$

$$P(n, p^-) - e^{-O(m)} \leq P(n, m) \leq P(n, p^+) + e^{-O(m)}$$

Thm: Let G be a cuckoo graph w/ $m = (1-\epsilon)\frac{n}{2}$, then

- ① w.p. $\geq 1 - \frac{1}{n}$, every component has size $\Theta(\log n)$
 ② w.h.p. the expected size of every component is $O(1)$

Pf: ① $G(n, p^+)$ with $p^+ = \frac{(1+\epsilon)m}{\binom{n}{2}} = \frac{1-\epsilon^2}{n-1}$

To bound the component sizes, we fix a vertex $v \in V$

& start a BFS $\xrightarrow{\text{num. neighbours}}$
 for any vertex $v_i, N_i = |N(v_i)| \sim \text{Bin}(n-i, p^+)$

$$v_1, v_2, \dots, v_{i-1}, v_i \xrightarrow{N_i} \text{Bin}(n-i, p^+)$$

$$\Pr\left(\sum_{i=1}^k N_i > k\right) \xrightarrow{\text{once again dep.}}$$

$$B_i \sim \text{Bin}(n-i, p^+) \Rightarrow \Pr\left(\sum_{i=1}^k B_i > k\right) \xrightarrow{\text{Branching process}}$$

$$B = \sum_{i=1}^k B_i \sim \text{Bin}(k(n-i), p^+)$$

more flexible, so can
bound this to bound
our target

$$E(B) = k(n-i)p^+ = k(1-\epsilon^2)$$

$$\Pr(B > k) = \Pr\left(B \geq \frac{E(B)}{1-\epsilon^2}\right) \leq \Pr(B \leq (1+\epsilon^2)E(B))$$

- ① If $m = (1-\epsilon)\frac{n}{2}$, then the expected size of any connected component in the cuckoo graph is $O(1)$.
- ② w.h.p. all components of size $\Theta(\log n)$ are unicyclic

Fix $v \in V$. Let S be the size of the component containing v in $G(n, m)$.

$$E(S) = \sum_{k>0} k \cdot \Pr(S=k)$$

$$= \Pr(S=1) +$$

$$\Pr(S=2) + \Pr(S=3) +$$

$$\Pr(S=3) + \Pr(S=4) + \Pr(S=5) +$$

\vdots

$$\underbrace{\Pr(S \geq 1)}_{\Pr(S \geq 1)}$$

$$\underbrace{\Pr(S \geq 2)}_{\Pr(S \geq 2)}$$

$$\cdots$$

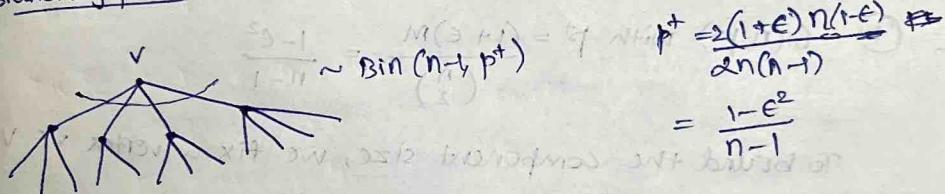
$$\Pr(S \geq n)$$

$$= \sum_{k>0} \Pr(S \geq k) \rightarrow \text{monotone increasing}$$

$$\leq \sum_{k>0} \Pr(S' \geq k) + \underbrace{n\epsilon}_{\text{outside summation, } k \leq n} \rightarrow O(1)$$

$$\leq \sum_{k>0} \Pr(S' \geq k) + \underbrace{n\epsilon}_{\substack{\text{component size} \\ \text{in } G(n, p^+)}} \rightarrow O(1)$$

Branching Process



$y_i = \text{No. of nodes at level } i$

$$Y = \sum_{i \geq 0} Y_i$$

$$E(Y_i) = \sum_{k \geq 0} E(Y_i | Y_{i-1} = k) \cdot \Pr(Y_{i-1} = k)$$

$$= \sum_{k \geq 0} k \cdot (n-1)p^+ \cdot \Pr(Y_{i-1} = k) = (n-1)p^+ E(Y_{i-1})$$

$$= ((n-1)p^+)^i = (1-\epsilon^2)^i$$

$$E(Y) = \sum_{i \geq 0} (1-\epsilon^2)^i = \frac{1}{\epsilon^2} = O(1)$$

$$\Rightarrow E(S) = O(1)$$

② In $G(n, m)$, how can a component not be unicyclic?

Fix k vertices - $\binom{n}{k}$

Cayley's formula: There are K^{K-2} labelled trees on k vertices

Edges are placed into the tree w.p. $\frac{m}{n^2} \cdot \binom{m}{k-1} \cdot (k-1)! = \left(\frac{1}{n^2}\right)^{k-1}$
 choosing that edge in $G(n, m)$

Add 2 edges $\binom{m-k+1}{2} \cdot 2! \underbrace{\left(\frac{1}{n^2}\right)^k}_{\text{within the tree}}$

We don't want any more edges to the tree (from other vertices)

$\Rightarrow \left(1 - \frac{k(n-k)}{n^2}\right)^{m-k-1} \rightarrow$ can add more cycles
 inside k -vertex-set, which double counts \Rightarrow lower bound

\Rightarrow Prob. that a component of size k isn't unicyclic:

$$\leq \binom{n}{k}^{k-2} \cdot \binom{m}{k-1} \cdot (k-1)! : \binom{m-k+1}{2} \cdot 2! k^k \left(\frac{1}{n^2}\right)^{k+1} \left(\frac{k(n-k)}{n^2}\right)^{m-k}$$

$\leq \frac{1}{n}$ when $k = \Theta(\log n)$, after heavy calculation

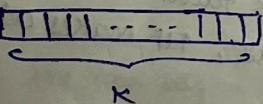
Midsem portions complete

Online algorithms

- Data arrives as a stream & not all at once
- At each step, we make irrevocable decisions
- We want to know how good the decisions (made with incomplete information) are

Paging problem

Cache



seq. of page reqs. r_1, r_2, \dots, r_m

If r_i does not exist in the cache, \Rightarrow cache miss / page fault
 evict some item & bring r_i in

Goal: Minimize page faults

Deterministic online strategies:

① LRU - Least recently used

② FIFO - First-in first-out

③ LFU - Least frequently used

An adversary can design the requests to forcefully trigger faults at every step

\Rightarrow what does it mean to minimize page faults?

How do we know if some strategy is good?

Optimal offline strategy

farthest-in-the-future

\Rightarrow compare against this

Competitive ratio

An algo. A is c -competitive if

(indep. of m)

Edges are placed into the tree w.p. $\frac{1}{n^2}$. choosing edge in GMA

Add 2 edges $\binom{m-k+1}{2} \cdot 2! \cdot \underbrace{\left(\frac{k^2}{n^2}\right)^2}_{\text{within the tree}}$

We don't want any more edges to the tree (from other vertices)

$$\Rightarrow \left(1 - \frac{k(n-k)}{n^2}\right)^{m-k-1} \rightarrow \text{can add more cycles inside } k\text{-vertex-set, which double counts} \Rightarrow \text{lower bound}$$

\Rightarrow Prob. that a component of size k isn't unicyclic:

$$\leq \binom{n}{k}^{k-2} \cdot \binom{m}{k-1} \cdot (k-1)! : \binom{m-k+1}{2} \cdot 2! \cdot k^2 \left(\frac{1}{n^2}\right)^{k+1} \left(\frac{1}{n^2}\right)^{k-1}$$

$\leq \frac{1}{n}$ when $k = \Theta(\log n)$, after heavy calculation

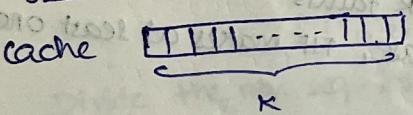
Midsem portions complete

11/9/23

Online algorithms

- Data arrives as a stream & not all at once
- At each step, we make irrevocable decisions
- We want to know how good the decisions (made with incomplete information) are

Paging problem



seq. of page reqs. r_1, r_2, \dots, r_m

If r_i does not exist in the cache, \Rightarrow cache miss / page fault
evict some item & bring r_i in

Goal: Minimize page faults

Deterministic online strategies:

① LRU - Least recently used

② FIFO - First-in first-out

③ LFU - Least frequently used

An adversary can design the requests to forcefully trigger faults at every step

\Rightarrow what does it mean to minimize page faults?

Competitive ratio

An algo. A is c -competitive if \forall seqs. of reqs. r_1, r_2, \dots, r_m , $f_A(r_1, r_2, \dots, r_m) \leq c f_{OPT}(r_1, r_2, \dots, r_m)$

(indep. of m)

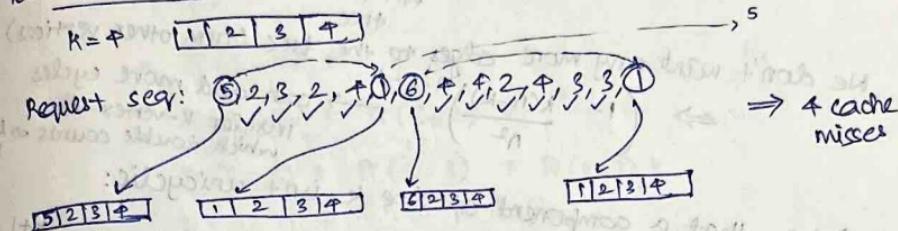
$f \approx \text{loss fn.}$

for $m \rightarrow \infty$

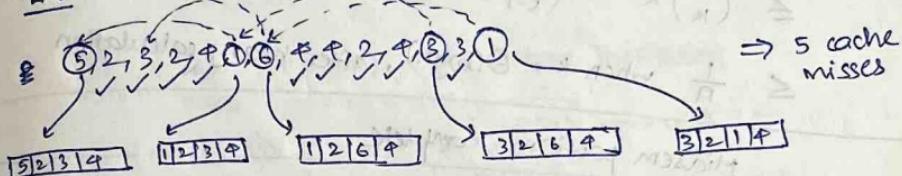
optimal algo., possibly offline

$\frac{f_A(r_1, r_2, \dots, r_m)}{f_{OPT}(r_1, r_2, \dots, r_m)}$ is bounded

Farthest-in-the-Future (FIF)



LRU



Thm: LRU is K -competitive

i.e., For any seq. r_1, r_2, \dots, r_m , if LRU makes n faults,
then FIF makes $\geq \frac{n}{K}$ faults

Idea:

Divide r_1, r_2, \dots, r_m into phases

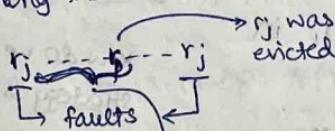
where in each phase, LRU makes K faults

$$r_1 \dots r_i \underbrace{r_{i+1} \dots r_j}_{K \text{ faults}} \dots \underbrace{r_{j+1} \dots r_m}_{K \text{ faults}}$$

We want to show that in each phase, FIF makes at least one fault.

Phase i

① Among the K faults, the same page r_j faulted twice



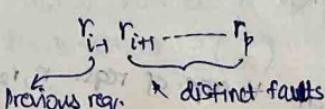
All other $K-1$ elements
in the cache were requested

\Rightarrow No. of distinct requests $\geq K-1$

In the interval $(r_j \dots r_p)$, there were $\geq K+1$ distinct requests

\Rightarrow FIF must also have faulted in that interval

② All the K faults are distinct



At the start of phase i,

r_{i-1} was in the cache

(a) r_{i-1} was one of the faults

\Rightarrow FIFO faults, w/ a similar analysis as in ①

(b) r_{i-1} wasn't one of the faults

From the $K-1$ "free" posns., there were K distinct reqs.

\Rightarrow FIFO faults

\therefore In each of the $\frac{n}{K}$ phases, FIFO faults at least once

\Rightarrow FIFO makes $\geq \frac{n}{K}$ faults

\therefore LRU is K -optimal

Thm: If a deterministic paging algo. is c -competitive,
then $c \geq K$.

Idea: Fix an algo. A.

Construct a seq. r_1, r_2, \dots, r_m s.t.

if $f_A(r_1, r_2, \dots, r_m) = n$,

then $f_{\text{OPT}}(r_1, r_2, \dots, r_m) \leq \frac{n}{K}$

1 2 3 ... K-1 K

1, 2, 3, ..., K-1, K, (k+1, e_1, e_2, ..., e_m), $e_i \in \{1, 2, \dots, K+1\}$

we can construct

by running A

element evicted

at prev. req.

$\Rightarrow f_A(r_1, r_2, \dots, r_m) = m$

Divide the req. seq. into phases of size K

A faults K times.

If OPT encounters a fault, then the $(k+1)^{\text{th}}$ element (not in cache) has been requested, instead of an element (that is not requested in the phase at all) in the cache.

Hence, post this fault, there will be no more faults in this phase

$\Rightarrow f_{\text{OPT}}(r_1, r_2, \dots, r_m) \leq \frac{m}{K} \Rightarrow$ Best possible det. algo. is K -competitive, i.e., for example, LRU

13/9/23

Randomized paging

Analysis will differ based on the power of the adversary.

We will analyse our algorithms assuming the weakest possible adversary:

Oblivious adversary \rightarrow knows only the source code of the algorithm

\rightarrow unaware of the internal coin tosses

\rightarrow non-adaptive; (input seq. decided at first, essentially)

- More powerful adversaries can be adaptive:
- offline adversary \rightarrow strongest form, knows the optimal offline algorithm
 - online adversary & final comparison is essentially against this
[not too clear, read up a bit more]

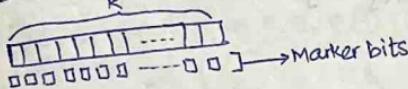
Competitive ratio

An algo. A has cr. C if $\forall r_1, r_2, \dots, r_M$,

$$E(f_A(r_1, r_2, \dots, r_M)) \leq C \cdot f_{\text{OPT}}(r_1, r_2, \dots, r_M)$$

\downarrow
over the internal randomness of the algo. A

Marker algorithm



The algo. proceeds in rounds. At the start of each round, marker bits are

for each req. r_i ,

If r_i is in the cache, set the marker bit to 1

Else, choose a pos. with marker bit 0 u.a.r.,
place r_i in that pos & set the marker bit to 1.

When all marker bits are 1, begin the next round.

(Each round has K distinct reqs.)

\Rightarrow If a cache entry is req'd., it is not evicted in that round

Thm: Marker algo. is dH_K -competitive

($H_K = k^{\text{th}}$ Harmonic sum)

Pf: for the optimal algo:

Round i :

The cache has the following types of items:

- stale items \rightarrow req'd. in round $i-1$ & round i (still in the cache)

- new items \rightarrow req'd. in round i , not in round $i-1$
denote by n_i

Let S_i be the items in the cache for OPT

(at the start)
(of round i)

$$d_i = |S_i \setminus S'_i|$$

may contain n_i 's

OPT will fault on at least $n_i - d_i$ reqs.

$$d_{i+1} = |S_{i+1} \setminus S'_{i+1}|$$

OPT will have faulted on $\geq d_{i+1}$ reqs.

[each item thrown away would've been a fault think about this]

$$\Rightarrow \text{OPT faults on } \geq \max\{n_i - d_i, d_{i+1}\}$$

$$\Rightarrow n_i - d_i + d_{i+1} - n_i \leftarrow$$

S'_{i+1} \curvearrowright queries in round $i \Rightarrow$ evictions come up as missing non in S'_{i+1}

Total no. of faults for OPT over m rounds $\geq \sum_{i=1}^m n_i$

Now, we want to show that expected no. of faults for Marker is $\leq H_K \sum_{i=1}^m n_i$

for Marker:

Round i :

New elements will ~~cause~~ page faults

Stale elements may cause page faults, if a new element evicts it before the stale element is queried

assume

\Rightarrow For the worst case, all the n_i new queries come first.

At least one of the new queries evicts ^{the first} stale element w.p. $\leq \frac{n_i}{K}$

($\frac{1}{K}$, union bound)

For the second stale element, prob. of

there are 2 possibilities

eviction before query $\leq \frac{n_i}{K-1}$

First stale query faulted

First stale query didn't fault

One new element in stale #1's pool,
other $(n_i - 1)$ new queries & stale #1 can pick from $(K-1)$ locations

stale #1 in own pool,
 n_i new queries have searched & picked from $(K-1)$ locations

$\frac{n_i}{K-1}$

$\frac{n_i}{K-1}$

For the whole round, no. of faults $\leq \sum_{i=1}^m n_i + \frac{n_i}{K} + \frac{n_i}{K-1} + \dots + \frac{n_i}{K-(K-1)}$

$\leq n_i H_K$

\Rightarrow Over all rounds, no. of faults $\leq H_K \sum_{i=1}^m n_i$

\Rightarrow c.r. dH_K

Yao's minimax principle \rightarrow std. technique to lower bound rand. algos.

The process is like a game b/w algo. designer & adversary

\Rightarrow 2-player zero-sum game

payoff matrix M

consider rock-paper-scissors

pay the cost
= no. of pg. faults

| | | R | P | S |
|---|--|----|----|----|
| | | 0 | -1 | 1 |
| | | 1 | 0 | -1 |
| R | | 0 | -1 | 1 |
| P | | 1 | 0 | -1 |
| S | | -1 | 1 | 0 |

entries: How much col. player pays row player
(payoff to row player)

Consider a slightly diff. version:

| | R | P | S |
|---|----|----|---|
| R | 0 | 1 | 2 |
| P | -1 | 0 | 5 |
| S | -2 | -1 | 0 |

Row player: $\max_i \min_j M_{ij}$

= Maximise the minimum payoff
(of a row)

Col-player $\equiv \min_j \max_i M_{ij}$

= Minimise the maximum amt. given away

If $\max_i \min_j M_{ij} = \min_j \max_i M_{ij}$, the game has a value

Here, both are 0 for (R, R) \rightarrow pure strategy, action (R) is fixed

For simple RPS, $\max_i \min_j M_{ij} = -1$, $\min_j \max_i M_{ij} = 1$

\rightarrow The game doesn't have a value

$\forall M$, $\max_i \min_j M_{ij} \leq \min_j \max_i M_{ij}$

Mixed strategies

The action isn't fixed, rather they are chosen using a proba. dbn.

$p \sim$ dbn. over rows, or $q \sim$ dbn. over cols.

Here, we consider the expected payoff = $\sum_{i=1}^m \sum_{j=1}^n p_i q_j M_{ij}$

$$= p^T M q,$$

Row player: $\max_p \min_q p^T M q$

~~= best strategy for well~~ P

Col player: $\min_q \max_p p^T M q$

Fixing the rand. bits for a rand. algo. makes it deterministic.

\Rightarrow A rand. algo. \sim dbn. over deterministic algos.

Von-Neumann minimax theorem

For every payoff matrix M (for 2-person zero-sum games),

$$\max_p \min_q p^T M q = \min_q \max_p p^T M q$$

The safest option for the row player is to play R, as it guarantees 0 payoff, while P & S could lead to a payoff of -1.

For the col. player, safest option is again R, as they would have to pay at most 0.

Suppose we fix p .

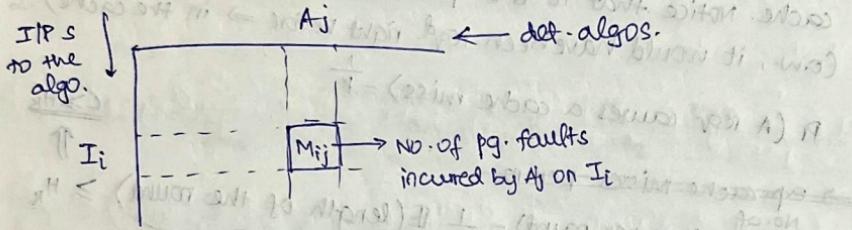
$$p^T M = \left[\sum_{i=1}^m p_i M_{i,1} \quad \sum_{i=1}^m p_i M_{i,2} \quad \dots \quad \sum_{i=1}^m p_i M_{i,n} \right] = [x_1 \ x_2 \ \dots \ x_n]$$

We want to find q that minimizes $p^T M q$, s.t. $\sum_{i=1}^n q_i = 1$

$$p^T M q = \sum_{i=1}^n x_i q_i \geq \min_i x_i \sum_{i=1}^n q_i = \min_i x_i \Rightarrow \text{choose the action that minimises}$$

Corollary:

$$\forall M, \max_{p^T} \min_j p^T M e_j = \min_q \max_i e_i^T M q$$



We want to show,

~~VP (for every rand. algo) s.t. dbn over A~~
~~VP (for every rand. algo) s.t. dbn. over A~~
~~FI (some bad input) s.t.~~

$E(c(A, I_i))$ is large \Rightarrow some row "sum" is large

remember,
pay off is
still not det.,
because of \oplus

It suffices to show that all col. "sums" are large.

The principle:

$$\forall M, p, q, \min_j p^T M e_j \leq \max_i e_i^T M q$$

exp. cost of
a det. algo.

when IP is
sampled
acc. to p

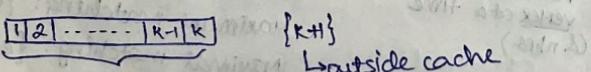
↓
Find a worst-case
IP dbn for all det. algos.

expected cost
of a rand. algo.

↓
we want to show
this is large

Thm: # for any c -competitive online paging algo., $C \geq H_k$

PF:



(P)
 L/P segs.
 $r_i = K+1$
 $\& e_r \{1, 2, \dots, K\}$

$r_i \in \{1, 2, \dots, K+1\} \setminus \{r_{i-1}\}$

Printers, I/O devices, etc.

Divide the reqs. into rounds

- In each round, there are k distinct page reqs.

OPT: At most 1 fault per round (At first faults exist element that won't be reqd this round)

Fix a det algo. A

It can be uniquely characterized at each step by the element outside the cache.

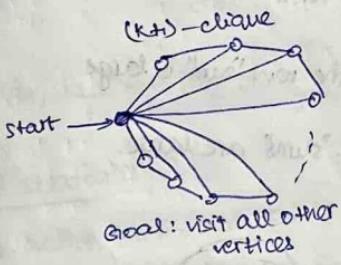
A cache miss occurs if the req. is the element outside the cache. Notice that it is a part of the current valid reqs. (otherwise, it would have been reqd. right before \Rightarrow in the cache)

$\Pr(\text{A req causes a cache miss}) = \frac{1}{k}$

$C \geq H_k$
↑

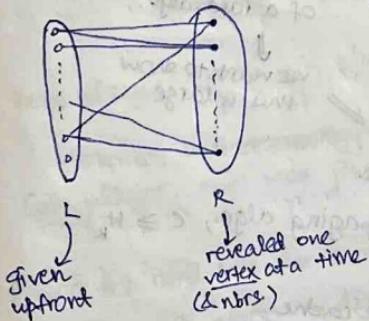
~~No. of cache misses per round~~
 $E(\text{cache misses per round}) = \frac{1}{k} E(\text{length of the round}) \geq H_k$
each round reqs requires k distinct reqs.

$\geq k H_k$, like from coupon collector's



~ coupon collector's prob., but we only want to visit k of $K+1$

Online bipartite matching



Computing maximal matching online is easy, just see if there is an edge to a free vertex in L & add if so

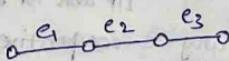
How good are maximal matchings, w.r.t. maximum matchings?

Maximum matching:

- Matching having the largest cardinality over all possible matchings

Maximal matching:

- Matching that cannot be extended



Maximum matching: $\{e_1, e_3\}$

Maximal matching: $\{e_1, e_3\}, \{e_2\}$

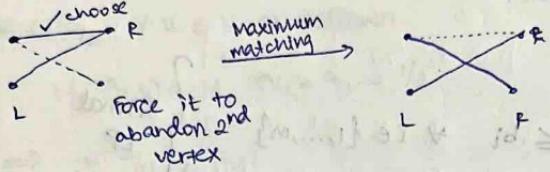
Det-greedy algo

- Add ~~each~~ an edge if possible (maximal matching)

⇒ $\Omega(1)$ size \Rightarrow online α -competitive algo. for maximum matching

2. competitive is the best we can get (deterministically)

M: Greedy maximal matching
OPT: Max. matching
 $\text{TO show: } \text{OPT} \leq 2\text{M}$



There exists a randomized online matching algo w/ c.r. $(1 - \frac{1}{e})$

2019/23 [rand. online $\geq (1 - \frac{1}{e})$ max. matching]

Midsem discussion

2. a) $n_1 = \# \text{ of voters for G}_1$

$$n_2 = n - n_1$$

$$E(X) = \sum_{k=0}^m k \cdot \frac{\binom{n_1}{k} \binom{n_2}{m-k}}{\binom{n_1+n_2}{m}} = \frac{k \cdot m}{n_1+n_2} = p_1 \cdot m$$

Or, show that $\Pr(X_i = 1) = p_1$

$$\Pr(X_{i+1} = 1) = \sum_{k=0}^{i-1} \frac{\binom{n_1}{k} \binom{n_2}{i-k}}{\binom{n_1+n_2}{i}} \cdot \frac{n_1 - k}{n_1 + n_2 - i}$$

A linear program for matching

for every edge e , define $x_e = \begin{cases} 1, & \text{if } e \text{ is in the matching} \\ 0, & \text{o.w.} \end{cases}$

$$\max \sum_{e \in E} x_e$$

$$\text{s.t. } \sum_{u \in N(v)} x_{(u,v)} \leq 1 \quad \forall v \in V \quad \left. \begin{array}{l} \text{defines all possible} \\ \text{matchings} \end{array} \right\} \quad \left. \begin{array}{l} \text{i.e., Every feasible soln.} \\ \text{is a matching} \end{array} \right.$$

In general, integer programs are v. hard to solve

\Rightarrow we usually relax constraints to reals rather than integers

A relaxed linear program: \rightarrow solvable in poly. time

$$\max \sum_{e \in E} x_e$$

$$\text{s.t. } \sum_{u \in N(v)} x_{(u,v)} \leq 1 \quad \forall v \in V$$

$$x_e \geq 0 \quad \forall e \in E$$

for example, $(0.5, 0.5, 0.5)$ is feasible

However, this may beget non-integral solns., which cannot correspond to a true soln!

$$V \ni u \cdot 0 \leq u$$

21/9/23

A general linear program:

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

$$\text{s.t. } \sum_{j=1}^n a_{ij} x_j \leq b_i \quad \forall i \in \{1, \dots, m\}$$

$$x_i \geq 0 \quad \forall i \in \{1, \dots, n\}$$

Primal
LP

Finding upper bounds on the obj-fn. allows us to estimate the soln. & tighter bounds bring us closer to the maximum. We can do this in a principled way using another LP.

$$y_1, y_2, \dots, y_m \geq 0$$

$$y_1 \left(\sum_{j=1}^n a_{1j} x_j \right) + y_2 \left(\sum_{j=1}^n a_{2j} x_j \right) + \dots + y_m \left(\sum_{j=1}^n a_{mj} x_j \right) \leq \sum_{i=1}^m b_i y_i$$

$$\Rightarrow x_1 \underbrace{\left(\sum_{i=1}^m a_{1i} y_i \right)}_{\geq c_1} + x_2 \underbrace{\left(\sum_{i=1}^m a_{2i} y_i \right)}_{\geq c_2} + \dots + x_n \underbrace{\left(\sum_{i=1}^m a_{ni} y_i \right)}_{\geq c_n} \leq \sum_{i=1}^m b_i y_i$$

If we force

$$\sum_{i=1}^m c_i x_i \leq$$

upper bound
on the
obj-fn.)

Dual LP

$$\min \sum_{i=1}^m b_i y_i$$

$$\text{s.t. } \sum_{i=1}^m a_{ii} y_i \geq c_i$$

$$\sum_{i=1}^m a_{ni} y_i \geq c_n$$

$$y_i \geq 0 \quad \forall i \in \{1, \dots, m\}$$

Every feasible soln.
for the dual is
an upper bound
on the optimum
(maximum)
of the primal

Weak-duality

$$\text{OPT(Primal)} \leq \text{OPT(Dual)}$$

Strong-duality

$$\text{OPT(Primal)} = \text{OPT(Dual)}$$

Dual program for "matching"

$$\min \sum_{v \in V} y_v$$

$$\text{s.t. } y_u + y_v \geq 1 \quad \forall (u, v) \in E$$

$$y_u \geq 0 \quad \forall u \in V$$

[Each edge shows up in the
constraints for the vertices
it is incident on]

Given the maximal matching M , we construct a feasible soln. for the dual.

Define $q_u + v \in V$ as follows:

$$\text{If } (u,v) \in M, \quad q_u = q_v = 1/2$$

$$\Rightarrow \sum_{u \in V} q_u = |M|$$

However, this is not a feasible assignment to q_u ($q_u = q_v$)
since edges in the matching do not have their constraints satisfied.

What about $q_u = 2q_v$?

This is feasible. Any vertex not in the matching has a neighbor in M .
 \rightarrow The constraint is satisfied for all edges

$$\Rightarrow \text{OPT} \leq \sum_{u \in V} q_u = 2 \sum_{u \in V} q_u = 2|M|$$

weak-duality

Primal-dual analysis (in general)

- Formulate the primal LP
- Use some algo. to construct a dual feasible soln.
- Apply weak-duality

Online fractional matching

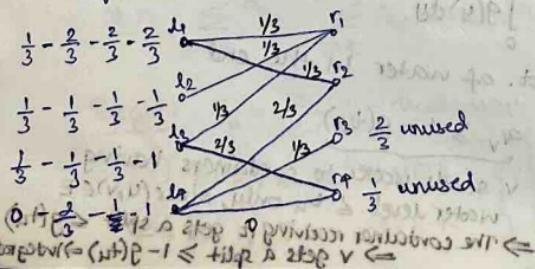
Essentially solving the relaxed LP for matching, in an online fashion.
 $x_e \geq 0$ instead of $x_e \in \{0,1\}$

f - fractional matching

$$\Rightarrow f \geq (1 - \frac{1}{e}) \text{OPT}$$

Waterlevel algo.

The vertices in L are imagined as containers that can hold one unit of water. The vertices in R are like reservoirs containing one unit of water. Edges are like pipes. Each new vertex brings with it some pipes & transfers water to the least full container (equally divided, if not unique) until the containers fill or the reservoir empties (if water levels equalise in between, the supply is further divided).



\Rightarrow Matching size: 3

Maximum matching: 4

$(l_1, r_1), (l_2, r_2), (l_3, r_3), (l_4, r_4)$

$(u, v) \in E$

Fix $l_w + u' \in L \setminus \{u\}$

Run the algo. on $G \setminus \{u\}$. Let u^* be the vertex matched to v .

Fix $l_{u^*} = 1$ if v is unmatched.

claim: If $l_u < l_{u^*}$, then u is matched in G .

Pf: ① u is already matched by the time v arrives

② Else, u gets matched to v

since u hasn't been matched yet, the run of the algo. on G & $G \setminus \{u\}$ are identical before v arrives.

$\forall v' \in N(v) \setminus \{u\}, l_{v'} > l_{u^*} \Rightarrow l_u < l_{u^*} + u' \in N(v) \setminus \{u\}$
 $\Rightarrow u$ gets matched to v

[If v was not matched in $G \setminus \{u\}$, now it gets matched to u]

If u is matched, $q_u = e^{l_u - 1}$

$$E(q_u) \geq \int_0^{l_u} e^{y-1} dy = e^{l_u - 1} - \frac{1}{e}$$

conditioned on $l_u + u' \in L \setminus \{u\}$

$$\text{claim: } q_v \geq \frac{\sum_{u' \in L \setminus \{u\}} q_u \Pr(l_u < l_{u'})}{1 - e^{l_u - 1}} \quad [l_u < l_{u^*} \text{ is a lower bound (smaller event) for } u \text{ being matched in } G]$$

+ values of l_u

Pf: Run the algo. on G & $G \setminus \{u\}$ in parallel

The set of unmatched vertices in G is a superset of those in $G \setminus \{u\}$

$\Rightarrow u^*$ is present in G as well when v arrives.

\Rightarrow The vertex chosen by G , say, w , is s.t.

$$l_w < l_{u^*}$$

$$\Rightarrow q_v \geq 1 - e^{l_{u^*} - 1}$$

$$\Rightarrow E(q_u + q_v) \geq 1 - \frac{1}{e} + (u, v) \in E$$

$$\therefore \text{set } y_u = \frac{e}{e-1} q_u$$

$$E(y_u) + E(y_v) \geq 1 + (u, v) \in E$$

$$\Rightarrow \text{OPT} \leq \frac{e}{e-1} \sum_{u \in L} E(y_u) + \sum_{v \in V} E(y_v) = \frac{e}{e-1} \sum_{u \in L} E(q_u)$$

$$\frac{1}{e-1} \leq (M + NP) \quad \rightarrow \leq (M + NP) = \frac{e}{e-1} E(MI)$$

27/9/23

Multiplicative weights update

Task: Make Yes/No decisions

wrong decisions incur unit cost

Goal: Minimize cost

We also have n experts e_1, e_2, \dots, e_n

No. of timesteps = T , m_i^T = Loss incurred by expert i

$$\Rightarrow \min. (m^T - \min_i m_i^T) \quad] \text{ regret}$$

Mistakes made
by best expert

① Omnipotent expert

- Never makes a mistake

S: set of potential omnipotent experts

$$S \leftarrow \{e_1, e_2, \dots, e_n\}$$

Consider the majority answer of all $e \in S$.

If the answer is correct, then no cost is incurred.

Else, S becomes halved.

$$\Rightarrow m^T - \min_i m_i^T \leq \log_2 |S| \quad [\text{Assume assuming that there is an omnipotent expert}]$$

② Weighted majority algo.

$w_i^{(t)}$: weight assigned to expert i at the t^{th} timestep.

$$w_i^{(0)} = 1 \quad \forall i \in \{1, \dots, n\}$$

If $\sum_{i \in S} w_i^{(t)} \geq \sum_{i \notin S} w_i^{(t)}$, answer Yes
 says yes says no

else All experts i with wrong answer have $w_i^{(t+1)} = (1-\epsilon)w_i^{(t)}$

Thm: Let $m_i^{(T)}$ be the loss incurred by expert i ,

$m^{(T)}$ be the loss incurred by the algo.

$$\Rightarrow m^{(T)} \leq 2(1+\epsilon)m_i^{(T)} + \frac{2\ln n}{\epsilon} \quad \forall i \in \{1, \dots, n\}$$

Thm: Every det. algo. will make at least twice as many mistakes as the best expert.

Pf: Suppose there are 2 experts e_1, e_2

e_1 always answers Yes, e_2 always answers No.

$$\sum_{i=1}^n m_i^{(T)} = q \cdot \sum_{i=1}^n m_i^{(T)} = \text{total mistakes}$$

Adversary can always give the opposite answer as the algo.

$$M^T = T, \text{ but } \min_i M_i^T \leq \frac{1}{2} T \quad \left[M_1^T + M_2^T = T \right] \\ \Rightarrow M_1^T > \frac{T}{2}, M_2^T > \frac{T}{2} \\ \Rightarrow \text{impossible}$$

PF (of prev. thm.):

$$\text{Potential fn, } \phi(t) = \sum_{i=1}^n w_i^{(t)}$$

$$\phi(0) = n$$

$$\phi(t+1) = \sum_{i, e_i} w_i^{(t+1)} + \sum_{i, e_i} w_i^{(t+1)} \\ \text{said Yes} \quad \text{said No}$$

Let correct
answer at
time step t
be No.

$$= (1-\epsilon) \sum_{i, e_i} w_i^{(t)} + \sum_{i, e_i} w_i^{(t)} \\ \text{said Yes} \quad \text{said No}$$

Algo. made error
in the t^{th} step

$$\sum_{i, e_i} w_i^{(t)} \geq \sum_{i, e_i} w_i^{(t)} \\ \text{said Yes} \quad \text{said No} \\ \geq \frac{\phi(t)}{2}$$

$$= \phi(t) - \epsilon \sum_{i, e_i} w_i^{(t)} \\ \text{said Yes}$$

$$\leq \phi(t) \left(1 - \frac{\epsilon}{2}\right) \quad \text{no. of mistakes} \\ = M(t)$$

Also, $w_i^{(t)} = (1-\epsilon)^{M_i^{(t)}}$

$$\Rightarrow \sum_{i=1}^n w_i^{(t)} = \phi(t) \leq \phi(0) \left(1 - \frac{\epsilon}{2}\right) = n \left(1 - \frac{\epsilon}{2}\right)$$

$$\Rightarrow (1-\epsilon)^{M_i^{(t)}} \leq \phi(t) \leq \left(1 - \frac{\epsilon}{2}\right) \cdot n$$

$$\therefore M^{(t)} \leq M_i^{(t)} \frac{\ln(1-\epsilon)}{\ln(1-\frac{\epsilon}{2})} - \frac{\ln n}{\ln(1-\frac{\epsilon}{2})} \quad \rightarrow \text{complete to get the bound}$$

4/10/23

Randomized version

Cost incurred, $m_i^{(t)} \in [-1, 1]$

Dfn. $p^{(t)}$ over the experts [which expert to go by]

\Rightarrow choose e_i acc. $p^{(t)}$ & give the same answer as, e_i

$M_1^{(t)}, M_2^{(t)}, \dots, M_n^{(t)}$ - costs incurred by the experts

\Rightarrow update $p^{(t)} \rightarrow p^{(t+1)}$

Expected cost for the algo. in time $t = \sum_{i=1}^n M_i^{(t)} p_i^{(t)}$

Expected total cost = $\sum_{t=1}^T M^{(t)} \cdot p^{(t)}$

Regret = $\sum_{t=1}^T M^{(t)} \cdot p^{(t)} - \min_i \sum_{t=1}^T M_i^{(t)}$

MWU algo.

$$\eta \leq \frac{1}{2}, w_i^{(1)} = 1 \quad \forall i \in \{1, \dots, n\}$$

$$\text{For each } t, \phi(t) = \sum_{i=1}^n w_i^{(t)}$$

$$p^{(t)} = \left\{ \frac{w_1^{(t)}}{\phi(t)}, \frac{w_2^{(t)}}{\phi(t)}, \dots, \frac{w_n^{(t)}}{\phi(t)} \right\}$$

Observe $m_1^{(t)}, m_2^{(t)}, \dots, m_n^{(t)}$.

$$\forall i \in \{1, \dots, n\}, w_i^{(t+1)} = w_i^{(t)} (1 - \eta m_i^{(t)})$$

$\begin{cases} \text{Large } m_i^{(t)} \rightarrow \text{large loss} \\ \Rightarrow \text{make } w_i^{(t+1)} \text{ a lot smaller} \\ \text{Neg. } m_i^{(t)} \rightarrow \text{inc. } w_i^{(t+1)} \end{cases}$

Thm: $\forall i \in \{1, \dots, n\}$,

$$\underbrace{\sum_{t=1}^T m_i^{(t)} \cdot p^{(t)}}_{\text{Total expected cost}} \leq \sum_{t=1}^T m_i^{(t)} + \eta \sum_{t=1}^T |m^{(t)}| + \frac{\ln n}{\eta}$$

[Similar to prev. result w/o a factor of 2]

Total expected cost.

* If $m_i^{(t)}$ corresponds to the gain instead of the loss/cost,

$$\text{the update rule becomes: } w_i^{(t+1)} = w_i^{(t)} (1 + \eta m_i^{(t)})$$

& the bound on the total expected gain becomes:

$$\underbrace{\sum_{t=1}^T m_i^{(t)} \cdot p^{(t)}}_{\text{Total expected gain}} \geq \sum_{t=1}^T m_i^{(t)} - \eta \sum_{t=1}^T |m^{(t)}| - \frac{\ln n}{\eta}$$

$\begin{cases} \text{Simply make } m_i^{(t)} \rightarrow -m_i^{(t)} \\ \text{or } m_i^{(t)} \rightarrow 0 \end{cases}$

Learning a linear classifier (Winnow algo.)

[Littlestone, 1983]

[Winnowing \rightarrow separating line]

$$\text{datapoints } \{(\bar{a}_i, l_i)\}_{1 \leq i \leq m} \quad \bar{a}_i \in \mathbb{R}^n \quad l_i \in \{\pm 1\}$$

Goal: Find $\bar{w} \in \mathbb{R}^n$ s.t.

$$\forall i, \text{sgn}\left(\sum_{j=1}^n a_{ij} w_j\right) = l_i \quad \rightarrow \quad \sum_{j=1}^n a_{ij} w_j \geq 0, \text{ by absorbing } l_i \text{ into } \bar{a}_i$$

Assumptions: ① $|a_{ij}| \leq 1$

$$\text{Additional constraints: } \begin{cases} \text{② } \sum_{i=1}^n w_i = 1 \\ \text{③ } w_i \geq 0 \quad \forall i \in [n] \end{cases}$$

w.l.o.g., though it may increase soln. space from

$$\text{④ } \exists w^* \text{ s.t. } \forall i \in [m] \quad \sum_{j=1}^n a_{ij} w_j^* \geq \epsilon \quad \mathbb{R}^n \text{ to } \mathbb{R}^m, \text{ by adding } -\bar{a}_i \text{ to get an extended vector } \bar{a}_i^* \text{ in } \bar{a}_i$$

Algo:

$$w_i^{(1)} = \frac{1}{n} \quad \forall i \in [n]$$

for $t = 1$ to T :

check if $\forall j \quad \bar{a}_j \cdot \bar{w}^{(t)} \geq 0$ \rightarrow normalized

If $\exists j$ s.t. $\bar{a}_j \cdot \bar{w}^{(t)} < 0$,

$$\text{then } w_i^{(t+1)} = w_i^{(t)} \left(1 + \frac{\epsilon}{2} \sum_j a_{ij}\right)$$

T is yet to be fixed

$$w^{(t)} \approx p^{(t)}$$

$$a_{ij}^{(t)} \approx m_i^{(t)} \quad \text{for } \bar{a}_j \cdot \bar{w}^{(t)} < 0$$

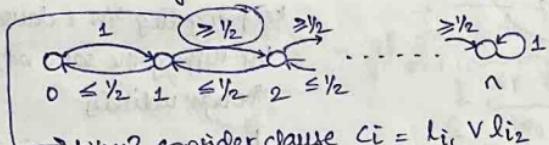
$$\text{Algo.: MWU}$$

Rand. 2-SAT as a Markov chain

No. of states = $n+1$

state i is associated w/ assignments that are at a Hamming dist. of $n-i$ from a sat. assignment

Each bit flip either increases or decreases the Hamming distance by 1. [Assuming there is 1 sat. assignment, or we only analyse for one]

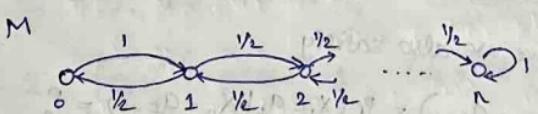


why? consider clause $C_i = l_{i_1} \vee l_{i_2}$

Both l_{i_1} & l_{i_2} cannot match the sat. assignment, otherwise $C_i = 1$.

⇒ w.p. $\frac{1}{2}$ (or 1), the flipping moves closer to the sat. assignment $\neg C_i$ (that are 0)

However, this is not a Markov chain, since the transition probas. are unknown. Instead, let us fix the $\{\geq \frac{1}{2}, \leq \frac{1}{2}\}$ edge wts. as $\frac{1}{2}$. If this reaches state n in t steps in expectation, so will the actual MDP process.



Expected time for M to go from 0 to n?

γ_i = TIME for M to go from i to n

$$E(\gamma_i) = \frac{1}{2}(1 + E(\gamma_{i-1})) \rightarrow \frac{1}{2}(1 + E(\gamma_{i+1}))$$

$$= \frac{1}{2}E(\gamma_{i-1}) + \frac{1}{2}E(\gamma_{i+1}) + 1$$

$$E(\gamma_0) = 1 + E(\gamma_1)$$

$$E(\gamma_n) = 0$$

Verify that $E(\gamma_i) = n^2 - i^2$

⇒ By running for $2n^2$ steps, we find a sat. assignment w.p. $\geq \frac{1}{2}$

Note that this won't work for 3-SAT,

since the trans. probas. forward can only be lower-bounded by $\frac{1}{3}$.

9.10.03

3-SAT

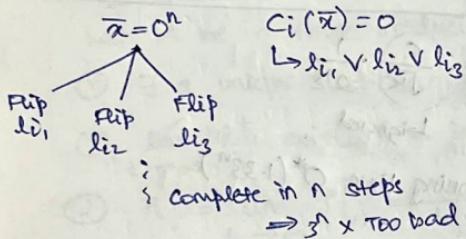
$$\phi = c_1 \wedge c_2 \wedge c_3 \wedge \dots \wedge c_m$$

$$c_i = l_{i1} \vee l_{i2} \vee l_{i3}$$

?

$$\exists \bar{x} \in \{0,1\}^n \text{ s.t. } \phi(\bar{x}) = 1$$

* NP-complete

* Trivial algo: $O^*(2^n)$ [$O(2^n \text{ poly}(n))$]Deterministic $O^*(1.732^n)$ 

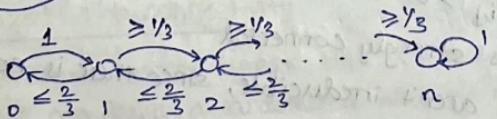
Obs: Every assignment $\bar{x} \in \{0,1\}^n$ is at a Hamming distance of $\leq \frac{n}{2}$ from 0^n or 1^n $\Rightarrow \frac{n}{2}$ steps on 0^n , if not satisfied, start over from 1^n & do $\frac{n}{2}$ steps

Randomized 3-SAT ①

- choose $\bar{x} \in \{0,1\}^n$ → only change (random instead of arbitrary)
- If $\phi(\bar{x}) \neq 1$

$$\text{Let } c_i = l_{i1} \vee l_{i2} \vee l_{i3} \text{ s.t. } c_i(\bar{x}) = 0$$

choose l_{i1}, l_{i2}, l_{i3} w.r.t. & flip



If ϕ is satisfiable,

$$\Pr(\text{success}) \geq \sum_{k=0}^n \underbrace{\binom{n}{k}}_{\text{# chosen } \bar{x} \text{ has Hamming dist. } k \text{ from } \bar{x}} \cdot \frac{1}{2^n} \cdot \left(\frac{1}{3}\right)^k$$

Prob. of moving forward until soln. r. weak!

$$\geq \frac{1}{2^n} \cdot \left(\frac{4}{3}\right)^n = \left(\frac{4}{3}\right)^n$$

\Rightarrow repeating $\left(\frac{4}{3}\right)^n$ gives a soln. w.h.p

$$\therefore O^*(1.5^n)$$

Schöning's algo.

Suppose \bar{x} has Hamming wt. k

Prob. that the algo. makes l moves
in the wrong direction = $\binom{k+2l}{l} \left(\frac{2}{3}\right)^l \left(\frac{1}{3}\right)^{k+l}$

(for the Markov chain)

$$\Pr(\text{Algo. reaches sat. assignment}) \geq \max_l \binom{k+2l}{l} \left(\frac{2}{3}\right)^l \left(\frac{1}{3}\right)^{k+l}$$

$$\geq \underbrace{\binom{3k}{k} \left(\frac{2}{3}\right)^k \left(\frac{1}{3}\right)^{2k}}_{l=k} \geq \frac{c}{jk} \left(\frac{1}{2}\right)^k$$

$$\boxed{\sqrt{2\pi k} \left(\frac{k}{e}\right)^k \leq k! \leq 2\sqrt{\pi k} \left(\frac{k}{e}\right)^k}$$

$$\Rightarrow \Pr(\text{success}) \geq \sum_{k=0}^n \binom{n}{k} \cdot \frac{1}{2^n} \cdot \frac{c}{jk} \cdot \frac{1}{2^k}$$

Lognated

$$\geq \left(\frac{3}{4}\right)^n$$

\Rightarrow running time $O^*(1 \cdot 3^n)$

More about Markov chains

* vector Π - dbn. over Ω

(P = transition matrix)

ΠP - dbn. over Ω after one timestep

ΠP^n - dbn. over Ω after n timesteps of the random walk

P^n : n -step transition matrix

* If $\exists n \quad P_{ij}^n > 0$, \exists a way to reach j from i (n steps)

Defn: A Markov chain M is irreducible if

$\forall i, j, \exists n \quad P_{ij}^n > 0$

\Leftrightarrow The digraph is strongly connected

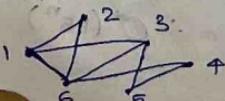
Our satisfiability MCs aren't irreducible, since there is no
way out from state n

* Defn: State j in a MC M is periodic if ~~s.t.~~ \rightarrow probably $\exists \Delta > 1$

$$\Pr(X_{s+\Delta} = j \mid X_s = j) = 0 \text{ unless } \Delta \mid t$$

A MC M is aperiodic if it does not have any periodic states.

* $G(V, E)$ - undirected



Random walk on G :

$$P = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 & 6 \\ 0 & Y_2 & Y_3 & 0 & 0 & 1/3 \\ 2 & 1/2 & 0 & 0 & 0 & 1/2 \\ 3 & & & & & \\ 4 & & & & & \\ 5 & & & & & \\ 6 & & & & & \end{bmatrix}$$

If G is connected, the random walk MC is irreducible

If G is bipartite, the random walk MC is periodic ($\Delta=2$)

Exercise: show that the random walk on an undirected graph is aperiodic iff the graph is not bipartite

[For non-bipartite graphs, we use the odd cycle to show it]

* Stationary dbn.

A dbn. π on the MC is stationary if

$$\pi P = \pi$$

* Fundamental thm. of MCs

Let M be a finite irreducible aperiodic MC.

- ① \exists a unique stationary dbn. π
- ② $\lim_{t \rightarrow \infty} P_{ij}^t$ exists & is indep. of i [prob. of reaching j , indep. of where we start]
- ③ $\pi_j = \lim_{t \rightarrow \infty} P_{ij}^t = \frac{1}{h_{ij}}$ [exp. time of first visit]

h_{ij} = Expected time for the MC to return to j starting from i

* Undirected, non-bipartite, connected G

↓
aperiodic

↓
irreducible

What is $\pi^{st} \cdot \pi P = \pi$?

We know, $P_{u,v} = \frac{1}{\deg(u)}$

Verify: $\pi(u) = \frac{\deg(u)}{2m}$, $m = |E|$

π is uniform when
 G is regular!



Randomized s-t connectivity

- BFS/DFS takes linear time & space

- we can design a randomized algo. via a random walk, using no extra space (besides the curr. index)

Algo:

Perform a random walk for l steps from s & check if t was visited

- $O(\log n)$ space \rightarrow RL space (randomized logspace)

- $O(l)$ time
 $= O(mn)$

\hookrightarrow In 2004, shown to be

solvble in det. $O(\log n)$ space

Cover time: Expected time for a random walk to visit all vertices

claim: $\forall (u, v) \in E, h_{u,v} + h_{v,u} \leq 2|E|$ $h_{u,v}$: Expected time for a random walk starting from u to reach v

Pf: Consider M' over directed edges

by changing $(u, v) \rightarrow \underbrace{u \rightarrow v, v \rightarrow u}_{\text{states of } M'}$

Assume M is bipartite
(can do so by adding self-loops)

No. of states of M' is $2|E|$

M' is aperiodic & irreducible.

$$P_{u \rightarrow v, v \rightarrow u} = \frac{1}{\deg(v)}$$

Exercise: Stationary dist. is uniform $\frac{1}{2|E|}$.

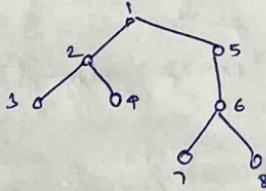
$$\Rightarrow h_{u,u} = 2|E|$$

$$\Rightarrow \cancel{\forall (u, v) \in E, h_{u,v} + h_{v,u} \leq 2|E|}$$

Lemma: Cover time of the random walk $\leq 2m(n-1)$

Expected time to visit all vertices in the component

Pf: Spanning tree of $G \setminus T$



DFS traversal:

$$1, 2, 3, 2, 1, \\ 5, 6, 7, 6, 8, 6, 5, 1$$

Consider a random walk, forcing a visiting in this DFS order, i.e., we are first interested in $1 \rightarrow 2$, then $2 \rightarrow 3$, then $3 \rightarrow 2$, etc.

Expected len. of the random walk $\leq h_{u,v} + h_{v,u}$ $(u, v) \in T$

$$= 2m(n-1)$$

$$< 2n^3 \Rightarrow \text{can random walk for } 8n^3 \text{ steps & get error rate of } \leq \frac{1}{4}$$

Monte Carlo methods - counting & sampling \rightsquigarrow Idea: "throw darts", based on sampling & identifying properties

Defn: Let V be a qty. to be estimated. X is an (ϵ, δ) -approx. of V if $\Pr(|X - V| \geq \epsilon V) \leq \delta$. cr.v.)

Fully Poly-time Randomized Approx Scheme (FPRAS)

Task: Given x ; calculate $V(x)$

An algo. A is an FPRAS if it takes x, ϵ, δ as I/P s.t.

$A(x, \epsilon, \delta)$ is an (ϵ, δ) -approx. of $V(x)$ & runs in time $\text{poly}\left(\frac{1}{\epsilon}, \ln\frac{1}{\delta}, \log n\right)$

DNF counting

$$\phi = T_1 \vee T_2 \vee \dots \vee T_m \longrightarrow \text{DNF}$$

$$T_i = l_{i1} \wedge l_{i2} \wedge \dots \wedge l_{ir}$$

* DNF-SAT $\in P$

* $\# \text{-DNF} = |\{\bar{a} \mid \phi(\bar{a}) = 1\}| \in \#P\text{-complete}$
 (sharp DNF) $\underbrace{|\{\bar{a} \mid \phi(\bar{a}) = 1\}|}_{\text{No. of satisfying assignments}}$

Trivial algo: Sample $\bar{a} \in \{0,1\}^n$ } Regs. $n^{O(1)}$ repeats
 & check if $\phi(\bar{a}) = 1$ if no. of sat. assignments
 is $\text{poly}(n)$

Expressing the problem differently:

Sets $S_1, S_2, \dots, S_m \equiv \text{sat. assignments of } T_1, T_2, \dots, T_m$

Given x , check $x \in S_i$

Computing $| \bigcup_{i=1}^m S_i |$

$$- S_i = \{ \bar{a} \mid T_i(\bar{a}) = 1 \}$$

$$T_i = \overline{x_1} \wedge x_2 \wedge x_3 \wedge \overline{x_4} \text{ (say)}$$

0 1 1 0 \longrightarrow necessary to satisfy

other variables can be set arbitrarily

\Rightarrow computing $|S_i|$ is

- can sample u.a.r. from S_i

easy enough

$$|S_i| = 2^{n-k}, \text{ if } T_i \text{ is satisfiable}$$

& T_i has k literals

$$C_\phi = \{(i, \bar{a}) \mid T_i(\bar{a}) = 1\}$$

$$|C_\phi| = \sum_{i=1}^m |S_i|$$

$$\text{Obs: } |C_\phi| \leq m \left| \bigcup_{i=1}^m S_i \right|$$

$$S_\phi = \{(i, \bar{a}) \mid T_i(\bar{a}) = 1 \wedge \forall j < i, T_j(\bar{a}) = 0\}$$

$\equiv T_i$ is the first term satisfied by \bar{a}

$$\text{Obs: } |S_\phi| = \left| \bigcup_{i=1}^m S_i \right| \Rightarrow S_\phi \text{ is similar in size to } C_\phi, \text{ which avoids the prob we faced before}$$

Algo:

- Sample i w.p. $|S_i| / \sum |S_i|$ } sample (i, \bar{a}) u.a.r. from C_ϕ

- sample \bar{a} u.a.r. from S_i

- check if $(i, \bar{a}) \in S_\phi$

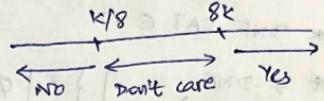
18/10/23 #SAT → generic boolean formula satisfiability

Thm: If $SAT \in P$, then there is an FPRAS for #SAT.

GapsAT

IP: ϕ, k

O/P: Yes, if $\#(\phi) > 8k$
No, if $\#(\phi) < k/8$



Lemma 1: If $SAT \in P$, GapsAT has an efficient randomized algo.

Lemma 2: If \exists an efficient rand.-algo for GapsAT, then \exists algo. A

$$\& i \in \mathbb{N} \text{ s.t. } \Pr\left[\frac{\#(\phi)}{2^i} \leq A(\phi) \leq \#(\phi).2^i\right] \geq 1 - \frac{1}{n}$$

Lemma 3: If \exists algo A as described in lemma 2, then #SAT has an FPRAS.

Pf. of lemma 3:

Given ϕ , define $\phi' = \bigwedge_{j=1}^j \phi(\bar{x}^{(k)})$

$$\#(\phi') = \#(\phi)^j$$

- choose $j = i/e$

- construct ϕ'

- compute $A(\phi')^j$

$$\Pr\left[\frac{\#(\phi')}{2^i} \leq A(\phi') \leq \#(\phi').2^i\right] \geq 1 - \frac{1}{n}$$

$$\Rightarrow \Pr\left[\frac{\#(\phi')^j}{2^{ij}} \leq (A(\phi'))^j \leq \frac{\#(\phi')^j \cdot 2^{ij}}{\#(\phi).2^{ij}}\right] \geq 1 - \frac{1}{n}$$

FPRAS O/P

Pf. of lemma 2:

Find smallest i s.t.

$$\text{GapsAT}(\phi, 8^i) = \text{Yes} \quad \& \quad \text{GapsAT}(\phi, 8^{i+1}) = \text{No}$$

$$\#(\phi) \geq 8^{i-1} = \frac{A(\phi)}{8} \quad \#(\phi) \leq 8^{i+2} = 64A(\phi)$$

choosing $A(\phi) = 8^i$

$$\Rightarrow \frac{\#(\phi)}{8^i} \leq A(\phi) \leq 8\#(\phi) \leq 64\#(\phi), \text{i.e., } i=6$$

Lemma: Let $\mathcal{H}_{n,m}$ be a family of pairwise-indep. hash fns.

$$h: \{0,1\}^n \rightarrow \{0,1\}^m, \text{ let } \epsilon > 0.$$

Let $S \subseteq \{0,1\}^n$ s.t. $|S| \geq 2^m/\epsilon^3$.

$$\Pr_{h \in \mathcal{H}_{n,m}} \left[\frac{(1-\epsilon)|S|}{2^m} \leq |\{x \in S \mid h(x) = 0^m\}| \leq (1+\epsilon) \frac{|S|}{2^m} \right] \geq 1 - \epsilon$$

Pf. of lemma 1:

Given ϕ, k ,

choose $m = \log k$

define $\phi'(x) = \phi(x) \wedge (h(x) = 0^m)$, $h \in \mathcal{H}_{n,m}$

If $\#(\phi) \geq 8k$, then ϕ' is satisfiable.

If $\#(\phi) \leq \frac{k}{8}$, then ϕ' is not satisfiable.

$$S = \{x \mid \phi'(x) = 1\}, |S| \geq \frac{8 \cdot 2^m}{8k}$$

choose $\epsilon = 1/2$,

$$\Pr(|\{x \in S \mid h(x) = 0^m\}| > 4) \geq 1/2$$

$$|S| \leq \frac{2^m}{8}$$

$$\begin{aligned} \Pr(\exists x \in S, h(x) = 0^m) &\geq \Pr(\exists x \in S \mid h(x) = 0^m \geq 1) \\ &\leq \underbrace{\Pr(h(x) = 0^m)}_{\substack{x \in S \\ \leq \frac{2^m}{8}}} \end{aligned}$$

$$\leq \frac{1}{8}$$

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

DNF-sampling

Given $\phi = T_1 \vee T_2 \vee \dots \vee T_m$

Goal: Sample \bar{x} u.a.r. s.t. $\phi(\bar{x}) = 1$

If $m=1$, i.e., $\phi = T_1 = (\text{say}) \bar{x}_1 \wedge x_2 \wedge \bar{x}_3 \wedge x_4$

$$\phi(\bar{x}) = 1 \Rightarrow x_1 = 0, x_2 = 1, x_3 = 0, x_4 = 1$$

\Rightarrow The remaining $n-4$ bits can be set u.a.r.

$$|\{\bar{x} \mid \phi(\bar{x}) = 1\}| = 2^{n-4}$$

For a more general ϕ , we sample a T_i & then sample $x \in \{0,1\}^n$

- sample i w.p. $\frac{|S_i|}{\sum_{i=1}^m |S_i|}$ ($S_i = \{\bar{a} \mid T_i(\bar{a}) = 1\}$)

- sample u.a.r. from S_i

Fix $\bar{a} \in \{0,1\}^n$ s.t. $\phi(\bar{a}) = 1$.

$$\begin{aligned} \Pr(\bar{a} \text{ is sampled}) &= \sum_{\substack{i \\ |T_i(\bar{a})|=1}} \frac{|S_i|}{\sum_{i=1}^m |S_i|} \cdot \frac{1}{|S_i|} \\ &= \frac{m(\bar{a})}{\sum_{i=1}^m |S_i|} \quad (m(\bar{a}) = |\{i \mid T_i(\bar{a}) = 1\}|) \end{aligned}$$

This is not yet uniform. In order to make it uniform, we need another sampling step w.p. $\frac{1}{m(\bar{a})}$, which is

- Accept \bar{a} w.p. $\frac{1}{m(\bar{a})}$. Else continue
(rejection sampling)

Fully Polynomial Almost Uniform Sampler (FPAUS)

I/P : x

Associated w/ x is a set $\Omega(x)$, which we want to sample from. x may be a graph & $\Omega(x)$ may be the set of all indep. sets or the set of all matchings in x .

O/P : $w \in \Omega(x)$ (ϵ -sampler)

A sampling algo. generates an ϵ -uniform sample of Ω

If w satisfies $\forall S \subseteq \Omega, |\Pr(w \in S) - \frac{|S|}{|\Omega|}| \leq \epsilon$

we also want our algo. to run

in $\text{poly}(\log(\frac{1}{\epsilon}), |\Omega|)$

actually uniform

This can also be expressed as follows (not proved in class yet)

Uniform dbn. over Ω , $\Pr(w) = \frac{1}{|\Omega|}, w \in \Omega$

sampling algo. A defines a dbn. over Ω , $A_\Omega(w), w \in \Omega$

s.t. $\frac{1}{2} \sum_{w \in \Omega} \left| \frac{1}{|\Omega|} - A_\Omega(w) \right| \leq \epsilon$

Thm: (Jerrum-Sinclair)

FPAVS \Leftrightarrow FPRAS, for self-reducible problems
- satisfiability,
- matching
- indep. set

Thm: If \exists an FPAVS for
indep. sets, then \exists
an FPRAS for indep. sets

(sampling \Rightarrow counting)

Idea: Given G , $\Omega(G) =$ set of all indep. sets in G

Suppose G has m edges e_1, e_2, \dots, e_m .

$G_i =$ Graph with e_1, e_2, \dots, e_i .

$$|\Omega(G)| = \underbrace{\frac{|\Omega(G_m)|}{|\Omega(G_{m-1})|}}_{\substack{2 \\ \downarrow}} \cdot \underbrace{\frac{|\Omega(G_{m-1})|}{|\Omega(G_{m-2})|}}_{\substack{2 \\ \downarrow}} \cdot \dots \cdot \underbrace{\frac{|\Omega(G_2)|}{|\Omega(G_1)|}}_{\substack{2 \\ \downarrow}} \cdot |\Omega(G_1)|$$

this is where
self-reducibility
plays its role

estimated by sampling from
denominator using FPAVS &
checking membership in numerator

We need to show that these ratios
aren't too small

$$r_i = \frac{|\Omega(G_i)|}{|\Omega(G_{i-1})|}$$

estimate \tilde{r}_i

Monte-Carlo

$$\text{Output } 2^n \prod_{i=1}^m \tilde{r}_i \approx 2^n \prod_{i=1}^m r_i$$

$$R = \frac{m}{n} \prod_{i=1}^m \tilde{r}_i$$

we want

Lemma: If $(\tilde{r}_i)_i$ is an $(\frac{\epsilon}{2m}, \frac{8}{m})$ -approx. of r_i , $\Pr(|\tilde{r}_i - r_i| \leq \frac{\epsilon}{2m} r_i) \geq 1 - \frac{8}{m}$

$$\text{then } \Pr(|R - 1| \leq \epsilon) \geq 1 - 8$$

we need a lower bound for $\frac{|\Omega(G_i)|}{|\Omega(G_{i-1})|}$

$$* \quad \Omega(G_i) \subseteq \Omega(G_{i-1})$$

Let (u, v) be the i^{th} edge e_i .

Consider $\Omega(G_{i-1}) \setminus \Omega(G_i)$.

Every $I \in \Omega(G_{i-1}) \setminus \Omega(G_i)$

contains both u & v .

[If it contains at most one of u, v , its presence in $\Omega(G_{i-1})$ \Rightarrow presence in $\Omega(G_i)$]

Claim: $|\Omega(G_{i-1}) \setminus \Omega(G_i)| \leq |\Omega(G_i)|$

$I \in \Omega(G_{i-1}) \setminus \Omega(G_i)$

I contains $u \& v$

remove v : $I \setminus \{v\} \in \Omega(G_i)$

$$\Rightarrow \frac{|\Omega(G_i)|}{|\Omega(G_{i-1})|} \geq \frac{1}{2}$$

20/10/23

- Sample from $\Omega(G_{i-1})$ using FPAUS \rightsquigarrow How many times?
- Count the fraction of samples in $\Omega(G_i)$

Say we sample M ~~several~~ times using an $\frac{\epsilon}{6m}$ -sampler for $\Omega(G_{i-1})$

Let $X_k = \begin{cases} 1, & \text{the } k^{\text{th}} \text{ sample } \in \Omega(G_i) \\ 0, & \text{o.w.} \end{cases}$

By FPAUS defn., $(\Omega = \Omega(G_{i-1}), S = \Omega(G_i))$

$$\left| \Pr(X_k = 1) - \frac{|\Omega(G_i)|}{|\Omega(G_{i-1})|} \right| \leq \frac{\epsilon}{6m}$$

$$\Rightarrow \left| E(X_k) - \frac{|\Omega(G_i)|}{|\Omega(G_{i-1})|} \right| \leq \frac{\epsilon}{6m}$$

$$\Rightarrow \left| E\left(\frac{1}{M} \sum X_k\right) - \frac{|\Omega(G_i)|}{|\Omega(G_{i-1})|} \right| \leq \frac{\epsilon}{6m}$$

$$\Rightarrow \left| E(\tilde{r}_i) - \frac{r_i}{M} \right| \leq \frac{\epsilon}{6m}$$

for large enough M , $E(\tilde{r}_i) \geq \frac{1}{3}$

what should M be so that

$$|\tilde{r}_i - E(\tilde{r}_i)| \leq \frac{\epsilon}{12m} E(\tilde{r}_i)$$

We bound $\Pr[|\tilde{r}_i - E(\tilde{r}_i)| \leq \frac{\epsilon}{12m} E(\tilde{r}_i)] \geq 1 - \frac{8}{m}$

Using Chernoff bounds, $M = O\left(\frac{m^2}{\epsilon^2} \ln\left(\frac{2m}{8}\right)\right)$

$$\Rightarrow 1 - \frac{\epsilon}{12m} \leq \frac{\tilde{r}_i}{E(\tilde{r}_i)} \leq 1 + \frac{\epsilon}{12m}, \text{ w.p. } \geq 1 - \frac{8}{m}$$

$$\Rightarrow 1 - \frac{\epsilon}{6mr_i} \leq \frac{E(\tilde{r}_i)}{r_i} \leq 1 + \frac{\epsilon}{6mr_i}$$

Using $r_i \geq \frac{1}{2}$, we get

$$1 - \frac{\epsilon}{2m} \leq \frac{\tilde{r}_i}{r_i} \leq 1 + \frac{\epsilon}{2m}$$

This is not yet a complete algo., since we don't have an algo. for FPAVS yet, and we have only constructed a redⁿ. from FPAVS to FPRAS. To construct an FPAVS, we want to design a Markov chain with states corresponding to the set we sample from & a uniform stationary dbn.

Markov chain Monte Carlo

Let Ω be a finite state space.

Consider a directed graph on Ω that is strongly connected.

$$N(x) = \{y \mid \xrightarrow{(x,y)} \in E\}$$

$$M \geq \max_x |N(x)|$$

Define the trans? matrix P as

$$P_{x,y} = \begin{cases} \frac{1}{M}, & x \neq y, y \in N(x) \\ 1 - \frac{|N(x)|}{M}, & x = y \\ 0, & \text{o.w.} \end{cases}$$

P has a uniform stationary dbn.

* We want to construct such an MC for indep. sets, but we cannot do it explicitly by finding all indep. sets, which may be exponential in number & hence, impractical.

Let X_i be an indep. set.

$\forall v \in V$, \rightarrow defines the transition

- If $v \in X_i$, then $X_{i+1} = X_i \setminus \{v\}$

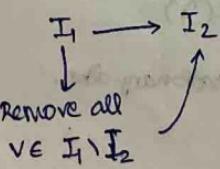
- If $v \notin X_i$ & $\{v\} \cup X_i$ is an indep. set,

then $X_{i+1} = X_i \cup \{v\}$

- Else, $X_{i+1} = X_i$

If we show that this MC is irreducible & aperiodic, its stationary dbn. is uniform & we are done.

Irreducible:



Aperiodic:

If the graph has an edge, the MC has a self-loop, allowing arbitrarily many timesteps at the state w/ a self-loop \Rightarrow aperiodic

We also want to bound the time to converge to the stationary dbn., since our FPAVS should run in poly(n), $\Omega = O(2^n)$

25/10/23
Monomer-Dimer problem (Mark Jerrum & Alistair Sinclair)

Graph $G(V, E)$, $|V| = 2n$

Molecules are placed on the graph s.t.

2 types of molecules
 Monomers occupy one vertex each,
 Dimers occupy adjacent vertices (i.e., an edge)



A config. is a placement of monomers & dimers s.t.
 every vertex has only one molecule.

M_K = set of configs. w/ K dimers ($\Rightarrow 2(n-K)$ monomers)

$$|M_K| = m_K$$

Note that M_n = set of all perfect matchings

for a config. M , $|M| = \text{No. of dimers in } M$
 Our example config. has $|M| = 1$

statistical
 physicists study these & make use of a partition fn.

$$Z(\lambda) = \sum_M w(M)$$

$$w(M) = \lambda^{|M|}$$

$$Z(\lambda) = \sum_{K=0}^{m_K} m_K \cdot \lambda^K \quad \rightarrow \text{A generating fn. for the no. of matchings}$$

The task of computing $Z(\lambda)$ given λ is #P-complete.

Jerrum-Sinclair provided an FPRAS for the problem.

The idea is similar to the previous, writing as a telescopic product.

$$Z(\lambda) = \frac{Z(\lambda_r)}{Z(\lambda_{r-1})} \cdot \frac{Z(\lambda_{r-1})}{Z(\lambda_{r-2})} \cdots \frac{Z(\lambda)}{Z(\lambda_0)} \cdot Z(\lambda_0)$$

$$\text{s.t. } 0 = \lambda_0 < \lambda_1 < \lambda_2 < \cdots < \lambda_r = \lambda$$

Sampling from a Gibbs dist. $\Pi_\lambda(M) = \frac{\lambda^{|M|}}{Z(\lambda)}$
 allows us to deal w/ this

product well \rightsquigarrow we need to construct
 a Markov chain w/ such a stationary dist.

$$\text{Using, } \lambda_i = \left(1 + \frac{1}{n}\right)^{i-1} \lambda_1, \quad \lambda_1 = \frac{1}{|E|}$$

$$\text{Estimating } \frac{z(\gamma_i)}{z(\gamma_{i-1})}$$

After a few times - sample M acc. to the Gibbs dbn. (via the MC)

- Estimate $\hat{x} = \left(\frac{\gamma_{i-1}}{\gamma_i} \right)^{IMI}$

We claim this is a good estimator

$$E(x) = \sum_M \left(\frac{\gamma_{i-1}}{\gamma_i} \right)^{IMI} \cdot \frac{\gamma_i^{IMI}}{z(\gamma_i)} = \underbrace{\frac{1}{Z(\gamma_{i-1})} \sum_M (\gamma_{i-1})^{IMI}}_{z(\gamma_{i-1})}$$

$$= \frac{z(\gamma_{i-1})}{z(\gamma_i)}$$

MC w/ a

Now, we need a method to convert an uniform stationary dbn. to one w/ an arbitrary stationary dbn.

Metropolis algo.

Given finite, irreducible, aperiodic MC w/ uniform stationary dbn. & given a dist. π over the states (M) \rightarrow A new MC w/ the same state space w/ stationary dbn. π

A MC for matchings

M_i = current matching

What are the transitions, M_{i+1} ?

- w.p. $\frac{1}{2}$, $M_{i+1} = M_i$

w.p. $\frac{1}{2}$, $\# \text{edges } (u, v) \in E$,

- e chosen w.p. $\frac{1}{|E|}$
 - $M_{i+1} = M_i \setminus \{e\}$, if $e \in M_i$
 - $M_{i+1} = M_i \cup \{e\}$, if both u & v are unmatched in M_i
 - $M_{i+1} = M_i \setminus \{e'\} \cup \{e\}$, if exactly one of u & v is matched in M_i via $e' \in E$
 - $M_{i+1} = M_i$, o.w.

The MC is irreducible (can remove all matching edges & then add new ones) & aperiodic (self loops exist) \Rightarrow it has a stationary dbn.

There are $|E|$ out edges in the second transition

$\Rightarrow |E|$ more from the first (self-loops)

$\Rightarrow 2|E|$ out edges per state of MC

\Rightarrow The MC is uniform over the set of all matchings

Recall the uniform MC we earlier constructed.
 state space Ω , $N(x) = \text{set of nbrs. of } x, x \in \Omega$

$$M \geq \max_x |N(x)|$$

$$P_{x,y} = \begin{cases} \frac{1}{M}, & \text{if } y \in N(x) \\ 0, & \text{if } y \notin N(x) \\ 1 - \sum_{y \in N(x)} P_{x,y}, & \text{if } y = x \end{cases}$$

↓ To get a dbn. Π

$$P_{x,y} = \begin{cases} \frac{1}{M} \min \left\{ 1, \frac{\pi_y}{\pi_x} \right\}, & \text{if } y \in N(x) \\ 0, & \text{if } y \notin N(x) \\ 1 - \sum_{y \in N(x)} P_{x,y}, & \text{if } y = x \end{cases}$$

This needs the dbn. Π . But, for matchings, notice that the transitions do not change the size of the matching by more than 1.

$$\Rightarrow \frac{\pi_y}{\pi_x} = \frac{\pi_y}{\pi_x} \frac{\pi_x(M')}{\pi_x(M)} = \lambda \in \{-1, 0, 1\}$$

$$\in \left\{ \frac{1}{\lambda}, 1, \lambda \right\}$$

⇒ A simple additional step in the MC is sufficient
 (for matchings)

$$- w.p. \min \left\{ 1, \frac{\pi_x(M_{i+1})}{\pi_x(M_i)} \right\}, \text{ move to } M_{i+1}$$

Lemma: Let M be a finite, irreducible, aperiodic MC w/ transition matrix P . Let Π be some dbn. over the state space Ω .

$$\text{If } \forall x, y \quad \pi_x P_{x,y} = \pi_y P_{y,x}$$

then Π is the stationary dbn.

Pf: $\sum_{x \in \Omega} \pi_x P_{x,y} = \sum_{x \in \Omega} \pi_y P_{y,x} = \pi_y \underbrace{\sum_{x \in \Omega} P_{y,x}}_1 = \pi_y \Rightarrow \Pi \text{ is the stationary dbn.}$

To verify that the above MC has Π as its stationary dbn,

notice that if $\pi_x < \pi_y$,

$$P_{x,y} = \frac{1}{M} \cdot \frac{\pi_x}{\pi_y} \quad \& \quad P_{y,x} = \frac{1}{M} \cdot \frac{\pi_y}{\pi_x}$$

$$\Rightarrow \pi_x P_{x,y} = \pi_y P_{y,x}$$

Now, we are interested in using our estimate of $Z(\bar{\lambda})$ to approximate m_n , the no. of perfect matchings.
 i.e., given algo. for (ϵ, δ) -approx. of $Z(\bar{\lambda}) \approx \hat{Z}(\bar{\lambda})$,
 approximate m_n .

- sample M from the Gibbs dbn.
- count the frac. of matchings that are perfect matchings (say, x)

$$E(x) = \frac{m_n}{Z(\bar{\lambda})} \cdot \frac{x \cdot \hat{Z}(\bar{\lambda})}{x^n} = \frac{x \cdot \hat{Z}(\bar{\lambda})}{x^n} = (\bar{\lambda})^n$$

$$\text{Estimate of } m_n, Y = \frac{x \cdot \hat{Z}(\bar{\lambda})}{x^n}$$

Lemma: For a graph G w/ no. of matchings of size K being m_K ,

$$\forall K, m_{K-1} \cdot m_{K+1} \leq m_K^2 \quad (\text{log-concavity})$$

$$\Rightarrow \frac{m_{K-1}}{m_K} \leq \frac{m_K}{m_{K+1}} \leq \dots \leq \frac{m_{n-1}}{m_n}$$

$$\begin{aligned} \frac{m_K}{m_n} &= \frac{m_K}{m_{K+1}} \cdot \frac{m_{K+1}}{m_{K+2}} \cdot \dots \cdot \frac{m_{n-1}}{m_n} \\ &\leq \left(\frac{m_{K-1}}{m_n} \right)^{n-K} \end{aligned}$$

$$E(X) = \frac{\bar{\lambda}^n}{\sum_k \left(\frac{m_k}{m_n} \right) \bar{\lambda}^k} \geq \frac{\bar{\lambda}^{n-K}}{\left(\frac{m_{K-1}}{m_n} \right)^{n-K}}$$

$$\begin{aligned} \frac{m_{K-1}}{m_n} &= \frac{m_{n-1}}{m_n} \Rightarrow E(X) \geq \frac{1}{n+1} \\ &\leq \left(\frac{m_K}{m_n} \right) \cdot \left(\frac{m_{n-1}}{m_n} \right)^K \leq (n+1) \cdot \bar{\lambda}^K \end{aligned}$$

$$\Rightarrow E(X) \geq \frac{1}{n+1}$$

27/10/23

Mixing times of Markov chains

Total variation distance of dbns. D_1 & D_2 supported on Ω ,

$$\|D_1 - D_2\| = \frac{1}{2} \sum_{x \in \Omega} |D_1(x) - D_2(x)|$$

$$= \max_{A \subseteq \Omega} |D_1(A) - D_2(A)| \quad (\text{can be shown})$$

$$\sum_{x \in A} D_1(x)$$

Let π be the stationary dist. of an MC M .

$P_\alpha^t = \text{Dist. over } \Omega \text{ after } t \text{ steps of } M \text{ starting from } \alpha$

$$\Delta_\alpha(t) = \|P_\alpha^t - \pi\|$$

$$\Delta(t) = \max_{x \in \Omega} \Delta_x(t)$$

Mixing time for a state $x \in \Omega$,

$$T_x(\epsilon) = \min \{t \mid \Delta_x(t) \leq \epsilon\}$$

Mixing time of M , $T(\epsilon) = \max_{x \in \Omega} T_x(\epsilon)$ Bounding this helps bound our algo's running time

M is rapidly mixing if $T(\epsilon)$ is $\text{poly}(\ln(1/\epsilon), n)$

\hookrightarrow IP size depends on context, not necessarily Ω

Coupling of Markov chains

A coupling of an MC M_t on the

state space Ω is a MC Z_t on the state space $\Omega \times \Omega$ s.t.

$$\hookrightarrow (x_t, y_t)$$

$$\Pr(X_{t+1} = x' \mid Z_t = (x, y)) = \Pr(M_{t+1} = x' \mid M_t = x)$$

$$\& \Pr(Y_{t+1} = y' \mid Z_t = (x, y)) = \Pr(M_{t+1} = y' \mid M_t = y)$$

A simple coupling could be just 2 indep. copies of an MC.

* Shuffling of cards

M_t : Pick $i \in \{1, \dots, n\}$, move the card in posⁿ. i to the top

Coupling (X_t, Y_t)

$X_t : M_t$ Let 'C' be the card that was chosen

$Y_t : \text{Move the card } C \text{ to the top}$

$X_{t+1} \sim M_{t+1}$ is clearly satisfied

$Y_{t+1} \sim M_{t+1}$ is because C is chosen u.a.r.

The fact that the choice was done by X_t doesn't affect the probab.

Notice that,

In expectation, after $n \log n + \Theta(n)$ steps, since this is just the coupon collector's problem.

Coupling lemma

Let $Z_t = (X_t, Y_t)$ be a coupling of an MC M_t on Ω .

Suppose that $\exists T$ s.t. $\forall x, y \in \Omega$,

$$\Pr(X_T \neq Y_T | X_0 = x, Y_0 = y) \leq \epsilon$$

Then, $\overline{\tau(\epsilon)} \leq T$.

Mixing time for M_t

Pf: Let π_0 be chosen acc. to the stationary dbn. π of M_t .

Let $A \subseteq \Omega$. Given T, ϵ ,

$$\begin{aligned} \Pr(X_T \in A) &\geq \Pr((X_T = Y_T) \cap (Y_T \in A)) \\ &\geq 1 - \Pr(X_T \neq Y_T) - \Pr(Y_T \notin A) \quad [\text{union bound}] \\ &= \Pr(Y_T \in A) - \underbrace{\Pr(X_T \neq Y_T)}_{\leq \epsilon} \\ &\geq \pi(A) - \epsilon \end{aligned}$$

Replacing A w/ $\Omega \setminus A$,

$$\Pr(X_T \in \Omega \setminus A) \geq \pi(\Omega \setminus A) - \epsilon$$

$$1 - \Pr(X_T \in A) \geq 1 - \pi(A) - \epsilon \Rightarrow \Pr(X_T \in A) \leq \pi(A) + \epsilon$$

$$\Rightarrow \pi(A) - \epsilon \leq \Pr(X_T \in A) \leq \pi(A) + \epsilon \Rightarrow \tau(\epsilon) \leq \epsilon$$

Independent sets of a fixed size K

$$K \leq \frac{n}{3\Delta+3}, \text{ where } n = |V|, \Delta = \text{max. deg.}$$

X_i - indep. set, $|X_i| = K \rightarrow |X_{i+1}| = K$

$v \in X_i, w \in V$

If $w \notin X_i \& X_i \setminus \{v\} \cup \{w\}$, $X_{i+1} = X_i \setminus \{v\} \cup \{w\}$
is an indep. set

Else, $X_{i+1} = X_i$

The MC is aperiodic ($w=v$)

It is irreducible because there are enough vertices outside X (src.) & Y (dest) to make appropriate transitions.

$$|X \cup Y| \leq 2K$$

$$\Rightarrow |X \cup Y \cup N(X \cup Y)| \leq 2K(\Delta+1) \leq 2n/3$$

Hence, it has a stationary dbn.

30/10/22

CS6170

Coupling

A coupling of an MC M_t on the state space Ω is an MC

$$Z_t = (X_t, Y_t) \text{ over } \Omega \times \Omega \text{ s.t.}$$

$$\Pr(X_{t+1} = x' | Z_t = (x, y)) = \Pr(M_{t+1} = x' | M_t = x)$$

$$\& \Pr(Y_{t+1} = y' | Z_t = (x, y)) = \Pr(M_{t+1} = y' | M_t = y)$$

Coupling lemma

Let $Z_t = (X_t, Y_t)$ be a coupling of an MC M_t on Ω .

Suppose $\exists T$ s.t. $\forall x, y \in \Omega$,

$$\Pr(X_T \neq Y_T | X_0 = x, Y_0 = y) \leq \epsilon_T.$$

Then, $T(\epsilon) \leq T_{\max}$

Independent sets of size k

$$\text{Max. deg.} \leq \Delta, k \leq \frac{n}{3(\Delta+1)}$$

MC: X_i Ind. set of size k

-choose $v \in X_i, w \in V - X_i$

-if $X_i \setminus \{v\} \cup \{w\}$ is an ind. set,

$$X_{i+1} = X_i \setminus \{v\} \cup \{w\}$$

$$\text{Else, } X_{i+1} = X_i$$

This is aperiodic & irreducible. We want to show it mixes well.

Coupling (X_t, Y_t)

$$X_t: v \in X_t, w \in V$$

& same as M_t

$$Y_t: \text{if } v \in Y_t \text{ (same } w \text{ as } X_t) \text{ move acc to } M_t$$

~~else, $v \in Y_t \setminus X_t$, move acc to v, w~~

~~else, $v \in Y_t \setminus X_t$, move acc to v, w~~

We want a bound on the no. of steps before $X_t = Y_t$,

$$\text{i.e., } d_t = |X_t - Y_t| = 0$$

Notice that $d_{t+1} \in \{d_t - 1, d_t, d_t + 1\}$

[Think about aperiodic signature again]

$$\Rightarrow \text{To bound: } \Pr(d_{t+1} = d_t^+ | d_t > 0) \leq ?$$

$$\Pr(d_{t+1} = d_t^- | d_t > 0) \geq ?$$

When is $d_{t+1} = d_t + 1$? Something must be removed from $X_t \cap Y_t$
 i.e., $v \in X_t \cap Y_t$

$$\Pr(d_{t+1} = d_t + 1 | d_t > 0)$$

$$\leq \underbrace{\frac{k-d_t}{k}}_{r \in X_t \cap Y_t} \cdot \frac{2d_t(\Delta+1)}{n}$$

$$(\& w \in (X_t \setminus Y_t) \cup (Y_t \setminus X_t))$$

$$\cup N(X_t \setminus Y_t) \cup N(Y_t \setminus X_t)$$

When is $d_{t+1} = d_t - 1$? something must be added to $X_t \cap Y_t$
 i.e., $v \in X_t \setminus Y_t, v' \in Y_t \setminus X_t$

$$\Pr(d_{t+1} = d_t - 1 | d_t > 0) \Rightarrow (\& w \in X_{t+1} \cap Y_{t+1})$$

$$\geq \frac{d_t}{k} \cdot \frac{n - (d_t+k)(\Delta+1)}{n} \Rightarrow \text{it is sufficient that}$$

$$w \notin \underbrace{(X_t \cup Y_t)}_{(d_t+k) \text{ size}} \cup N(X_t \setminus Y_t)$$

$$\mathbb{E}(d_{t+1} | d_t) = d_t \cdot \Pr(d_{t+1} = d_t) + (d_t + 1) \cdot \Pr(d_{t+1} = d_t + 1) + (d_t - 1) \cdot \Pr(d_{t+1} = d_t - 1)$$

$$= d_t + \Pr(d_{t+1} = d_t + 1) - \Pr(d_{t+1} = d_t - 1)$$

$$\leq d_t \left[1 + \frac{k-d_t}{k} \cdot \frac{2(\Delta+1)}{n} + \frac{n - (d_t+k)(\Delta+1)}{nk} \right]$$

$$\leq d_t \left[1 - \frac{n - (3k-1)(\Delta+1)}{nk} \right] \quad (\text{verify})$$

$$\Rightarrow \mathbb{E}(d_{t+1}) \leq \mathbb{E}(d_t) \left[1 - \frac{n - (3k-1)(\Delta+1)}{nk} \right]^{t+1}$$

$$\mathbb{E}(d_{t+1}) \leq \underbrace{\mathbb{E}(d_0)}_{\leq k} \left[1 - \frac{n - (3k-1)(\Delta+1)}{nk} \right]^{t+1}$$

$$\Pr(d_t \geq 1) \leq \mathbb{E}(d_t) \leq \epsilon$$

verify that we get $t = \text{poly}(n, \ln(1/\epsilon))$

Vertex coloring

$c: V \rightarrow \zeta$ Goal: color vertices s.t. adj. vertices

get diff. colors

max. deg. $\leq \Delta$

No. of colors $\geq 4\Delta + 1$

Task is to sample colorings

MC M : Choose, $v \in V$, color v w.r.t.,
try to color v with ℓ
Picking the same color allows aperiodicity.

Ex: Show that it is irreducible.

Each state has $|V| \cdot |\mathcal{C}|$ out edges \Rightarrow uniform stationary dist.

Coupling ($X_t \rightarrow Y_t$)

X_t & Y_t will behave exactly the same, i.e., same v & l

$d_t = \text{No. of vertices w diff. colors}$

Once again, $d_{t+1} \in \{d_t - 1, d_t, d_t + 1\}$

$$\Pr(d_{t+1} = d_t - 1 | d_t > 0) \geq \frac{d_t}{n} \cdot \frac{c - 2\Delta}{c}$$

v should
be colored
differently,
i.e., it is one
of the d_t

New color
must not be
used in its
nbrs $\rightarrow 2\Delta$
may be used

$$\Pr(d_{t+1} = d_t + 1 | d_t > 0) \leq \frac{n - d_t}{n} \cdot \frac{2\Delta}{c}$$

Must be
colored
the same
 \downarrow
color should be
in one of the
two colorings

We can now bound the expectations & consequently, the no. of steps, using a similar analysis as for ind. sets.

1/1/23

canonical paths

MC M over Σ - irreducible, aperiodic

$$\text{s.t. } \forall x, y \quad \pi(x) p_{x,y} = \pi(y) p_{y,x} = Q(x, y)$$

stationary dist.

Graph \tilde{G} : $V = \Omega$, $E = \{(x, y) | Q(x, y) > 0\}$
(undirected)

$\gamma_{x,y}$ = canonical path from x to y in \tilde{G} .

$$\Gamma = \{\gamma_{x,y} | x, y \in \Sigma\}$$

Load for an edge $e \in \tilde{G}$,

$$P(e) = \frac{1}{Q(e)} \sum_{(x,y): e \in \gamma_{x,y}} \pi(x) \pi(y) | \gamma_{x,y} |$$

\hookrightarrow len. of canonical path