

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

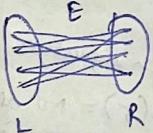
* Oder-Neisse line

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 &

recursing on the smaller graph \Rightarrow property called

self-reducibility

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

Polynomial Identity Testing (PIT)

Input: 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., $\forall x_1, \dots, x_n, p(x_1, \dots, x_n) = 0$

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

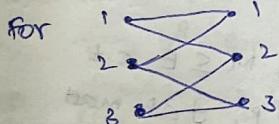
How are these probs. related?

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 PIT

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

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

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

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

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

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

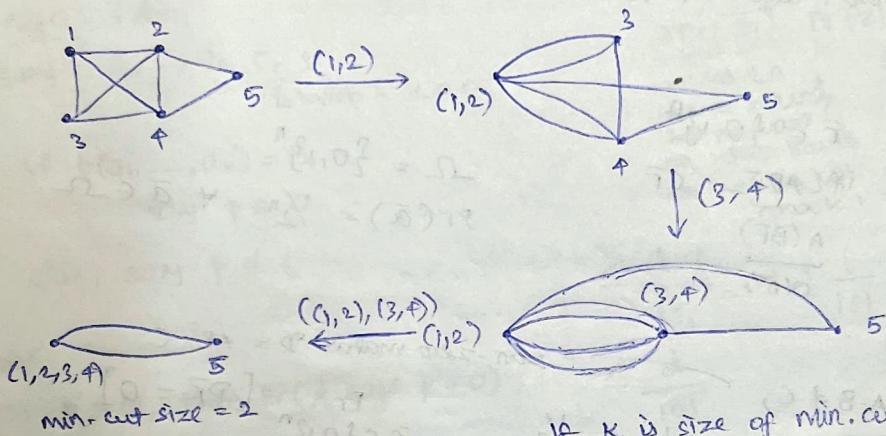
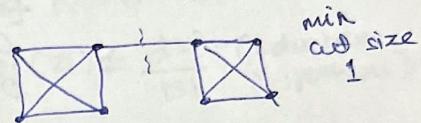
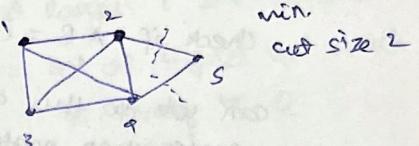
- A rand-algo. (Karger '93)

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: $\mathbb{E}(x) \geq \frac{m}{2}$

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

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

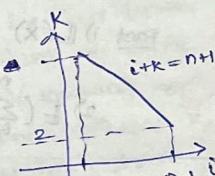
$$X_{ij} = \begin{cases} 1, & \text{if } b_i \text{ & } b_j \text{ were compared} \\ & \text{at some pt. in Quicksort} \\ 0, & \text{o.w.} \end{cases} \rightarrow \text{this occurs when}$$

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

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

$$\begin{aligned} * E(X) &= \Pr(Y=1) \cdot \underbrace{E(X|Y=1)}_{1} + \Pr(Y=0) \cdot \underbrace{E(X|Y=0)}_{x=1+z, z \text{ being an}} \\ &= p + (1-p)(E(1+z)) \quad y: \text{first trial} \\ &= 1 + (1-p)E(z) = 1 + (1-p)E(X) \\ \Rightarrow pE(X) &= 1 \Rightarrow E(X) = \frac{1}{p} \end{aligned}$$

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

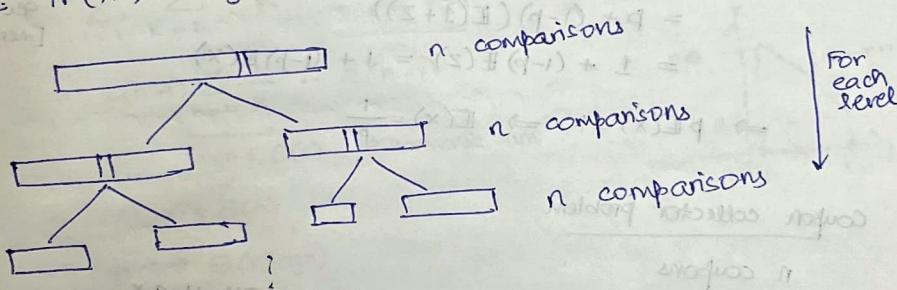
$$\Rightarrow E(I) \leq E\left(\frac{X}{a}\right) = \frac{E(X)}{a}$$

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

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

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

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

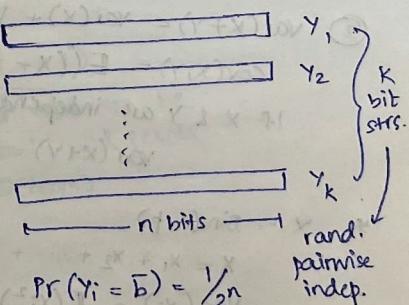
* Repeating the event k times, $\Pr(\text{Error}) \leq \frac{1}{2^k}$
w/ no. of rand. bits used = nk

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

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

$$b = \frac{c_{1j} - c_{2i}}{i - j}$$

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

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

$$GF(4) = \{0, 1, x, x+1\} \quad +, \cdot \text{ performed modulo } 1+x+x^2$$

the irreducible poly.
(of deg. k)
↳ this always
exists, apparently.

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

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

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

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

Chernoff-Hoeffding bounds

$$X_1, X_2, \dots, X_n, \quad X_i = \begin{cases} 1, & \text{w.p. } p_i \\ 0, & \text{w.p. } 1-p_i \end{cases}$$

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

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

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

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

Hoeffding's extension:

18/8/23 X_i 's are indep. & take values in $[a, b]$, $E[X_i] = \mu$

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

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.

$X_1, X_2, \dots, X_k, \quad X_i = \begin{cases} 1, & \text{if } i^{\text{th}} \text{ iteration answers correctly} \\ 0, & \text{o.w.} \end{cases}$

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

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

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

Prob. of error = $\Pr(X \leq \frac{k}{2})$

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

$$\text{When } \epsilon = \frac{1}{n}\epsilon,$$

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

$$\text{choosing } k = \frac{1}{\epsilon^2} \log n, \leq \frac{1}{n^2}$$

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

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

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.

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

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

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

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

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

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

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

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

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

strongest version
of the bound
but difficult
to apply

Weakening & simplifying w/ approximations:

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

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

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

$$\underset{\text{calculus/analysis}}{\leq} \left(\frac{e^\delta}{(1+\delta)^{1+\delta}} \right)^\mu$$

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

* 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 $x \geq \lambda$.

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

$$\leq \mathbb{E} \frac{e^{x(t-1)}}{e^{tx}}$$

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

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

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

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

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

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

$\xrightarrow{\text{indep.} \Rightarrow \text{prod.}}$ $\sum_{i=1}^n y_i^m = k$
 $\xrightarrow{\text{sum of Poisson}} \text{Poi}(m)$

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)

\blacksquare Let f be a non-negative fn.

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

$$\leq e^{m\bar{x}} \underbrace{\mathbb{E}(f(y_1^m, y_2^m, \dots, y_n^m))}$$

How to apply this? no conditioning!

w/ what proba.
suppose we want to answer,

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

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

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

2018/23

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

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

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

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

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

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

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

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 $x \notin S$,

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

$$\Pr[A[j] = 0] = \underbrace{\left(1 - \frac{1}{n}\right)^{mk}}_{\text{For given } x \notin S} \xrightarrow{\text{Overall } h_i \text{'s}} \Pr[x \in m] \approx e^{-mk/n} \Pr(Y_i \sim \text{Poi}(mk/n) = 0)$$

$A[j] \wedge A[j']$
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,

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

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

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

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

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

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

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

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

Just to say they are close
i.e., a good bound

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 $p \leq m$ & $2m$]

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

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

$$|H| = \Theta(p^2)$$

$$\Pr_{(a,y) \in H} [h(x) = h(y)] \leq 1/n$$

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

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

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

Then: H is 2-universal

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

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

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

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

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

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

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

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

Pairwise-indep. hash family

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

Consider the field $GF(2^m)$

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

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

The "pairwise" comes from here

any $(k-1) \Rightarrow k$ -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)

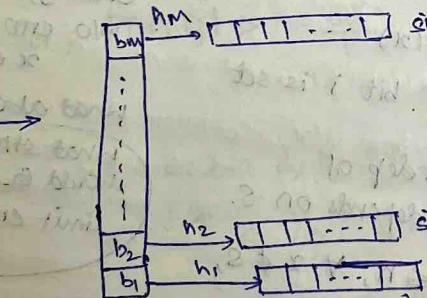
[Cell probe model]

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

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)

$$\sum_{i=1}^m \binom{b_i}{2} \leq m^2 \quad (\text{No. of collisions})$$

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

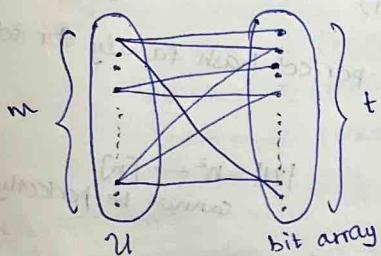
↳ can show existence but can't obtain the data structure
if we want to efficiently find the structure,

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

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

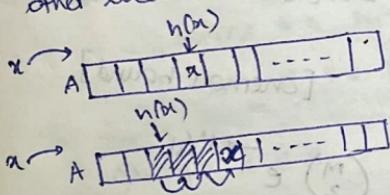
(Morris counter)
L way to count n
times using $\log n$ bits

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

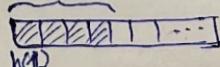
in practice.

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 4-wise indep. st. ~~Exp~~
 $x_i \sim \text{Ber}(p)$

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

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

$$\Pr(X > \mu + \beta) \leq \frac{\beta^2}{\beta^4}$$

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

\Rightarrow Given $n(1), X_i$'s are 4-wise indep

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

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

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

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

6/9/23

Pf. of fourth-moment bound:

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

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

Define $y_i = x_i - p$ Since x_i 's are 4-wise indep., so are y_i 's

$$E(y_i) = 0$$

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

$$= E\left[\sum_{i=1}^n y_i^+ + \binom{4}{2} \leq y_i^2 y_j^2 + \binom{4}{3} \leq y_i y_j^2 \quad i \neq j\right]$$

$$+ \sum_{i+j+k+l} y_i y_j y_k y_l + \sum_{i+j+k+l} y_i y_j y_k^2 \quad \text{other split}$$

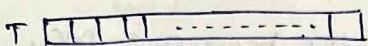
$$= \sum_{i=1}^n E(y_i^+) + \sum_{i \neq j} E(y_i^2) \geq E(y_i^2) \quad \begin{cases} \text{due to 4-wise indep.} \\ \text{2 go to 0} \end{cases}$$

$$E(y_i^+) = p(1-p)^+ + (1-p)(0-p)^+$$

$$E(y_i^2) = p(1-p)^2 + (1-p)p^2$$

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

Cuckoo hashing (Pagh, Rödler - '01)



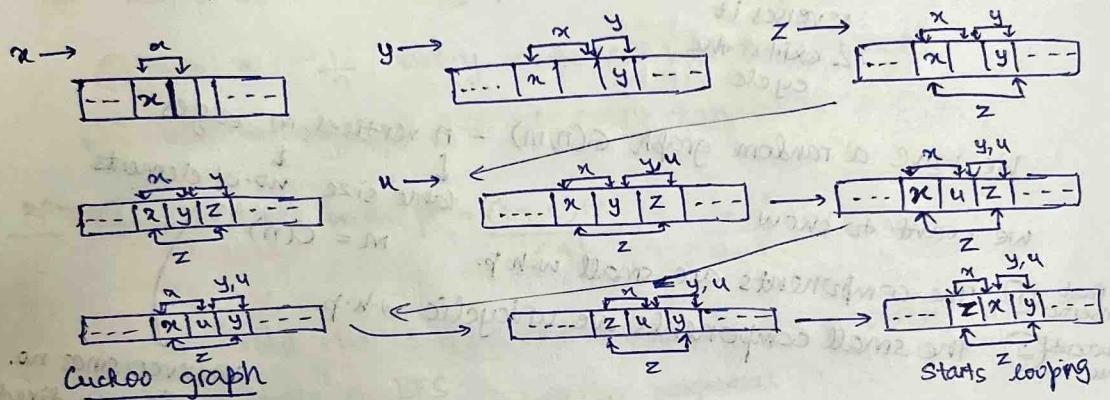
2 hash fn's. h_1, h_2 (completely random)

- * 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

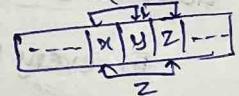


$G(V, E)$

$V \rightarrow$ 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



=

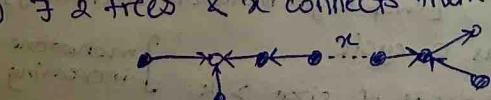
Obs: Let C be a connected component in the cuckoo graph, and x is being inserted. If the new component has > 1 cycle, then insertion is impossible. [without rehashing]

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

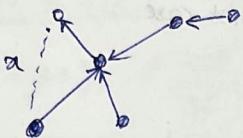
Directions show potential next pos's. for hashed elements

Each vertex has out-degree ≤ 1

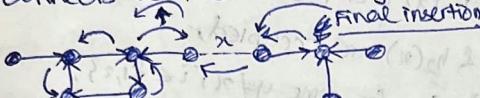
② x lies within a component C

$\Rightarrow C$ 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~~



Eventually insertion here
it reaches we have seen
the cycle, reverses it
& exits the cycle

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

table size \downarrow
no. of elements

we want to show:

$$m = O(n)$$

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.

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

Instead,
we analyse

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

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

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'

G' has the
property

[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

} monotone decreasing

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 $v_i, N_i = |N(v_i)| \sim \text{Bin}(n-i, p^+)$

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

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

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

$$B = \sum_{i=1}^k B_i \sim \text{Bin}(k(n-i), p^+) \quad \begin{matrix} \text{more flexible, so can} \\ \text{bound this to bound} \\ \text{our target} \end{matrix}$$

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

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

- 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

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

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

$$\cdots$$

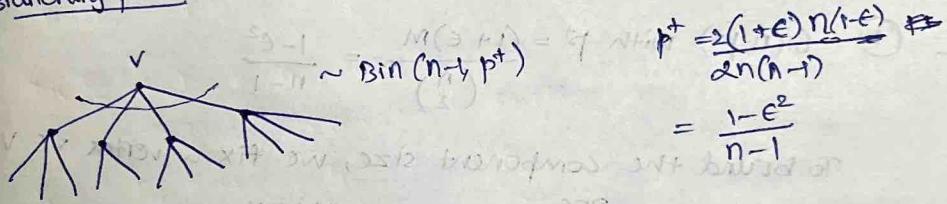
$$\Pr(S \geq n)$$

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

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

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

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:

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

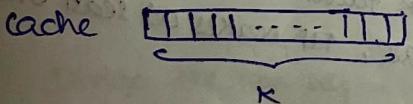
$\leq \frac{1}{n}$ when $k = \Theta(\log n)$, after heavy calculation

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



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 competitive if

How do we know if some strategy is good?

Optimal offline strategy

farthest-in-the-future

\Rightarrow compare against this

(indep. of m)