

Lecture 7: k-wise independence

①

Top 4 probability facts:

{ • "Union Bound": $\Pr[A \text{ or } B] \leq \Pr[A] + \Pr[B]$

{ • "Linearity of Expectation": average of $X+Y = \text{avg. of } X + \text{avg. of } Y$
 $(E[X+Y]) (= E[X] + E[Y])$

*[no matter
if A, B
X, Y independent,
dependent, whatever]*

• Markov's Inequality: If avg. of some positive #'s is 100, at most $\frac{1}{3}$ of them are ≥ 300 .

If $X \geq 0$ and $E[X] = \mu$, $\Pr[X \geq c\mu] \leq \frac{1}{c}$

• $1+x \approx e^x$ if x tiny [not about probability :)] $(c \gg 1)$

[[Why do we care about probability so much in an alg. class?]]

Life/alg. philosophy: if you don't know what choice to make, make a random one. At least this way it's hard for a devious adversary/fair input to force you into bad outcomes. Like when playing R.P.S.: if you play randomly, you know that no element of psychology can make your winning chances $< 50\%$.

[[Let's talk about one of the most famous probability chestnuts of all time:] "Birthday Paradox"]

[[actually play it!]]

Under SUHA

... $m < v$ if you have a good shape
so if you really want no harsh collisions, you're
[Great]

therefore $v > m^2$ is necessary & sufficient

to ensure no collisions "with high probability",
or $m = 1.18 \sqrt{365} \approx 33.5$

$$\bullet \quad m = 1.18 \sqrt{365} \approx 33.5$$

$$m = C \cdot \frac{v}{w} \approx \frac{v}{w}$$

$$m = C \cdot \frac{v}{w} \Leftrightarrow \frac{v}{w} = \frac{m}{C} \Leftrightarrow v \cdot C = m \cdot w$$

$$v \cdot C \ll w \Leftrightarrow v \cdot \frac{m}{w} \ll w \Leftrightarrow v \ll w$$

$$\text{Therefore: } m \ll v \Leftrightarrow \Pr[a \parallel \text{disint}] \approx 1 - \frac{m}{w}$$

$$\frac{m}{w} \approx e^{-\frac{m}{w}}$$

$$\frac{m}{w} \approx \frac{m}{\sqrt{w}} \rightarrow ((1-w) + (1-w) + \dots + (1-w))^{\frac{1}{2}} = e^{-\frac{m}{w}}$$

dangerous to be more precise up to

$$\Pr[\text{Max load } T] \approx e^{-\frac{m}{w}} \cdot e^{-\frac{m}{w}} \cdots e^{-\frac{m}{w}} = e^{-\frac{m}{w} \cdot m}$$

(Topo) if $\Pr[A \cdot B \cdot C \cdot D] = \Pr[A] \cdot \Pr[B] \cdot \Pr[C] \cdot \Pr[D]$

key: Using that all ball locations are independent $(\Pr[A \cdot B \cdot C \cdot D])$

exp [and b1 || does not collide & b2 || does not collide \Rightarrow $\Pr[\text{load } T] \leq \Pr[\text{load } T]$]

$$(1 - \frac{1}{w}) \cdot (1 - \frac{1}{w}) \cdots (1 - \frac{1}{w}) \cdot (1 - \frac{1}{w}) \cdots (1 - \frac{1}{w})$$

which $\Rightarrow \text{Max load } T$.

↑ shared previously

②

Balls & bins vec: m balls (people), v bins (people), $w = n$ cases

Ideas:

Only

Hope

h is "random enough":

Find

a small

set of

hash functions (rows)

all "easy to compute"

choose h randomly from the

only

Can't literally do it. Storing h takes $O(\lg n)$ bits
maybe 2¹⁰⁰⁰ Δ

SUHA = Chooses h at random from all possible h_0, \dots, h_{n-1}

h_{n-1}	h_{n-2}	\dots	h_1	h_0
$1 - v$	$1 - v$	\dots	$1 - v$	$v - 1$
0	1	\dots	0	0
$1 - v$	$1 - v$	\dots	0	0
1	\dots	\dots	0	0
0	0	\dots	0	0
\vdots	\vdots	\vdots	\vdots	\vdots
0	1	\dots	0	0
$1 - v$	$1 - v$	\dots	$1 - v$	$v - 1$

all possible hash fun.

"Universe" size: think of as huge, like $U = 2^n$, $n = 1000$, e.g.

→ just putting an upper bound on the size of the universe

numbers → covered by a # (ascii whatever)

hash function: $\{shingles\} \rightarrow \{0, 1, \dots, n-1\}$

Left's challenge is little → [e.g.: no collisions if all 2^m]

what are collisions?

How to choose hash functions that probably achieve properties SUHA

(Now like to "stop cheating" → we'll now try to ask...)

P.I Beyond SUHA

One notion of "random enough":

(4)

\mathcal{H} is "universal" [bad term, but anyway] if...

- If you first fix any two ~~distinct~~ elements $x, y \in U$,
 $x \neq y$
- Then you choose $h \in \mathcal{H}$ at random ..
- Then $\Pr[h(x)=h(y)] \leq \frac{1}{n}$.

Rem: SUHA $\equiv \mathcal{H}$ is all possible ^{"collision"} is universal: $\Pr_h[h(x)=h(y)] = \frac{1}{n}$.

e.g.: $U=3, n=2$.

	0	1	2	objs to hash./balls
h_0	0	0	0	0,1 one "bins"
h_1	0	1	1	weird hash function, hashes everything to 0!
h_2	1	0	1	
h_3	1	1	0	

$\mathcal{H} = \{h_0, h_1, h_2, h_3\}$ is universal. [And only uses 4 out of 8 poss. hash fns.]

↳ for any 2 distinct columns,

pick a random row,

chance of some "bin" is $= \frac{1}{2} \leq \frac{1}{2}$ ✓.

universal

Fact (we'll prove later): For any $U \leq 2^n$, there is a simple family \mathcal{H} of size ~~2^{2n}~~ 2^{2n} , hash functions "named" $h_{a,b}$ for all possible n -bit strings a, b ; $h_{a,b}(x) = \text{simple function of } a, b, x$.

[To implement your hash table, initially pick a, b : just 2^n (a few thousand?) random bits. Then hash with $h_{a,b}$.]

[Let's assume for now...]

Q: Is universality "random enough"?

A: For some properties achievable with SUHA, yes.

Let's go back to birthday paradox / balls & bins. -- ⑤

Say m balls, n bins, $n \gg m^2$. SUHA \Rightarrow max load 1 with high prob.

Proof breaks down without SUHA: $(1 - \frac{1}{n})(1 - \frac{2}{n}) \dots (1 - \frac{m-1}{n})$

assumed independence of

$h(x_1), h(x_2), \dots, h(x_m)$

don't have
any more

Different proof [more or less, up to constant factor] works!

Say you are hashing x_1, \dots, x_m using h drawn from universal ft.
 $h(x_1), \dots, h(x_m)$ are #'s ("bins") in $0 \dots n-1$ randomly

Define a boolean integer (random variable) $C_{1,2} = \begin{cases} 1 & \text{if } h(x_1) = h(x_2) \\ 0 & \text{if } h(x_1) \neq h(x_2) \end{cases}$ "1 & 2 collide"

$$\begin{aligned} E[C_{1,2}] &= \Pr[h(x_1) = h(x_2)] \cdot 1 + \Pr[h(x_1) \neq h(x_2)] \cdot 0 \\ &= \Pr[h(x_1) = h(x_2)] \quad [\text{expectation of "indicator" of event } A = \Pr[A]] \\ &\leq \frac{1}{n} \quad \text{by "universality".} \end{aligned}$$

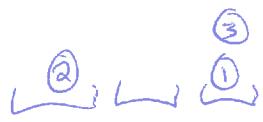
Similarly define $C_{i,j}$ for $1 \leq i < j \leq m$, $\begin{cases} 1 & \text{if "ball } i, j \text{ collide"} \\ 0 & \text{if not.} \end{cases}$

$$\text{Universality } \Rightarrow E[C_{i,j}] \leq \frac{1}{n}.$$

$$\begin{aligned} \text{Define } C &= C_{1,2} + C_{1,3} + \dots + C_{m-1,m} = \sum_{i < j} C_{i,j}. \\ &= \underbrace{\# \text{ of pairs of balls that "collide"}}_{\text{③}} \end{aligned}$$

Note: $C = 0 \Leftrightarrow$ no collisions / max load 1.

$C \geq 1 \Leftrightarrow$ at least one bin has load > 1



$$E[C] = E[C_{1,2}] + E[C_{1,3}] + \dots + E[C_{m-1,m}]$$

$$\begin{aligned} \text{"linearity of expectation!"} &\Rightarrow \leq \frac{1}{n} + \frac{1}{n} + \dots + \frac{1}{n} \quad \underbrace{\binom{m}{2} \text{ times}}_{\text{remember that quantity?}} \\ &= \frac{1}{n} \binom{m}{2} = \frac{1}{n} \frac{m(m-1)}{2} \leq \frac{m^2}{2n} \end{aligned}$$

$$\begin{aligned} C_{1,2} &= 0 \\ C_{1,3} &= 1 \\ C_{2,3} &= 0 \\ \hline C &= 1. \end{aligned}$$

⑥

Say n chosen $\geq 10m^2$.

Then $E[C] \leq \frac{m^2}{2 \cdot n} \leq \frac{m^2}{2 \cdot 10m^2} = \frac{1}{20} = .05$.

[Now think for a sec. C is either $0, 1, 2, 3, \dots$

And supposedly "on avg." it's $\leq .05$. Then it's probably usually $0!$]

"Markov's inequality": $\Pr[C \geq 1] \leq .05$ (otherwise, certainly $E[C] > .05$)
 $\Pr[\text{collisions}]$

\therefore universality is enough to conclude: ~~> 95%~~ chance of max load 1, provided $n \geq 10m^2$.

[What about our "usual" setting, $m=n$?] $\Pr[C \geq 10m] \leq \frac{1}{20} = .05$. ④

If $n=m$, $E[C] \leq \frac{m^2}{2m} \leq \frac{m}{2}$. [not amazing].

Q: If i th bin has load L_i , what will C be (at least)?

$$\begin{array}{c} 0 \\ 0 \\ 0 \\ 0 \\ \hline i \end{array} \xrightarrow{L_i=4} \text{contributes } \binom{4}{2} \text{ to } C. \quad C = \sum_{i=0}^{n-1} \binom{L_i}{2}.$$

$$\binom{L_i}{2} \approx \frac{L_i^2}{2}; \quad \Rightarrow \frac{L_i^2}{4} \quad \text{for } L_i \geq 2. \quad \therefore C \geq \sum_{i=0}^{n-1} \frac{L_i^2}{4}.$$

④ says > 95% of the time, ~~$10m \geq C$~~ $\Rightarrow \sum_i \frac{L_i^2}{4} \leq 10m \Leftrightarrow \sum_i L_i^2 \leq 40m$.

[What does this imply about max load?] \Rightarrow every L_i has $L_i^2 \leq 40m$
 $\Rightarrow L_i \leq \sqrt{40m}$.

> 95% chance of "max load $\leq O(\sqrt{m})$ " is not so great:

SUHA would have $O(\log m)$ ^{log form}.

[Sad thing is, this could happen.] [Good news \rightarrow we'll see a cool trick next time, "double hashing", to work around it.]