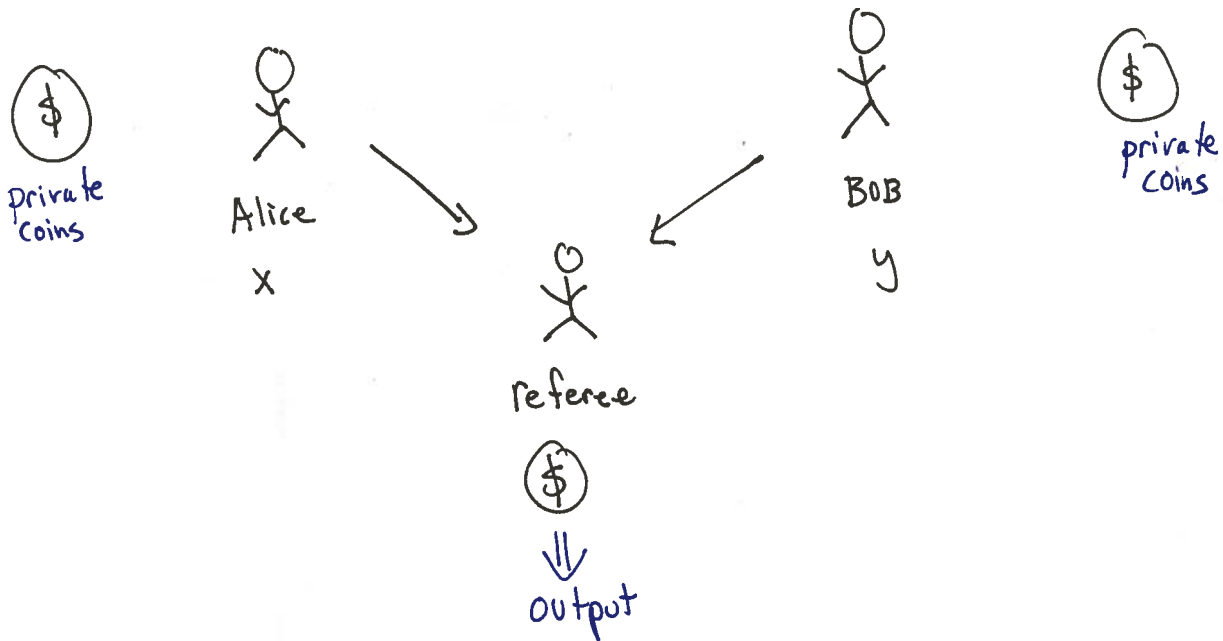


Lecture 21:

More lower bounds for
distribution testing problems

- communication complexity vs. sample complexity
- overview of closeness testing l.b.

Communication Complexity: Simultaneous Message Passing Model (SMP)



note: 1 round

$$\text{let } f(x,y) = \begin{cases} 1 & \text{if } x=y \\ 0 & \text{o.w} \end{cases}$$

Known Thm:

If $|x| = |y| = k$

must communicate $\Omega(\sqrt{k})$ bits to compute f

Yet another $\tilde{\Omega}(\sqrt{n})$ lower bound for testing uniformity
of a distribution :

Useful fact:

\exists error correcting code $C : \{0,1\}^k \rightarrow \{0,1\}^n$ s.t.

- (1) constant rate : $\frac{k}{n} \approx \text{const}$
- (2) relative distance $\delta = \Omega(1)$
- (3) each code word (each y in range of C)
is "balanced" - i.e., $\#1\text{'s} = \#0\text{'s}$

The reduction:

Alice: given x

- computes $C(x)$
- $A \leftarrow \{i \mid C(x)_i = 1\}$
- send s uniformly dist
samples from A to ref

Bob: given y

- computes $C(y)$
- $B \leftarrow \{j \mid C(y)_j = 0\}$
- send s unif dist
samples from B
to ref

Referee:

- receives samples from $A+B$
- constructs S samples of dist $q \equiv \frac{U_A + U_B}{2}$:

sample i :

toss coin
 heads: output next sample from A
 tails: " " " B

- feeds to uniformity tester + outputs pass if uniformity test passes

Why does it work?

if $X=Y$:

$C(X) = C(Y)$
 $A+B$ are partition of $[n]$ s.t. $|A|=|B|$

$$U_{[n]} = \frac{U_A + U_B}{2} \equiv q$$

so uniformity tester accepts q

✓ correct answer for this case

if $X \neq Y$:

$C(X) + C(Y)$ disagree on $\delta = \Omega(\epsilon)$ fraction of domain
 so $|A \cap B| = \frac{\delta}{2}$ + $|[n] \setminus (A \cup B)| = \frac{\delta}{2}$

$\Rightarrow q$ not supported on $\frac{\delta}{2}$ fraction of domain
 + "double weight" on " " " "

$$\text{so } \|q - U_{[n]}\|_1 > \delta$$

if $X \neq y$ (cont.)

so unif test rejects q



correct answer
for this
case

sample complexity:

each sample encoded via $O(\log n)$ bits

\Rightarrow comm. complexity is $O(s \log n)$

but we know that comm complexity for $X \stackrel{?}{=} y$

is $\Omega(\sqrt{k})$

$$\Rightarrow s = \Omega\left(\frac{\sqrt{k}}{\log n}\right) \stackrel{\text{const rate}}{\downarrow} = \Omega\left(\frac{\sqrt{n}}{\log n}\right) = \tilde{\Omega}(\sqrt{n})$$

So what?

why yet another (weaker) proof of uniformity testing?

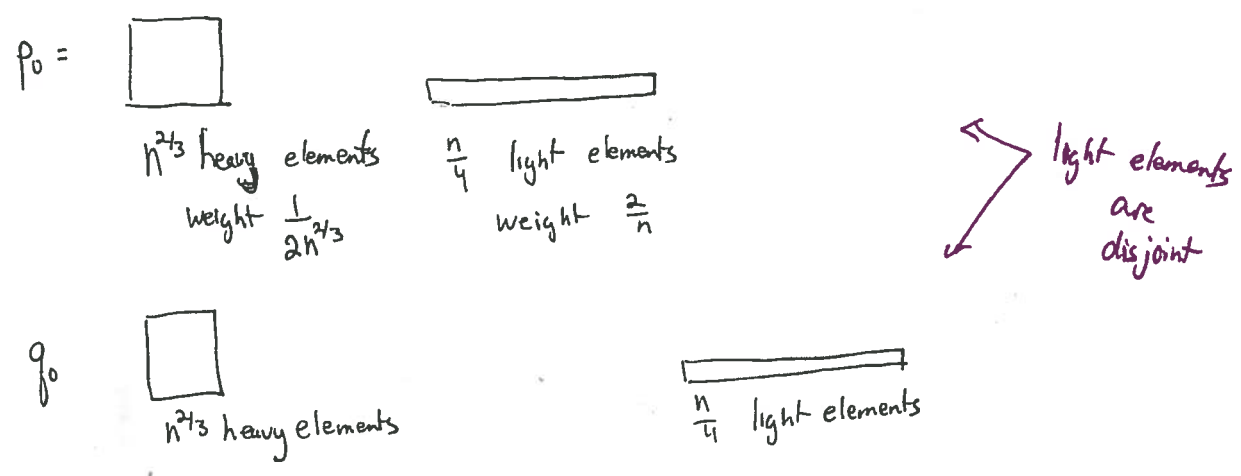
methodology works for many other problems in

dist testing (via fancier codes, reductions)

Sketch of l.b. for p, q given by samples \Leftarrow "closeness testing"

Thm Closeness testing requires $\Omega(n^{2/3})$ samples

Proof idea:



Positive pairs

Negative pairs

$l, \text{dist}=0 \Rightarrow (\pi(p_0), \pi(p_0)) \forall \pi$ $(\pi(p_0), \pi(q_0)) \nexists \pi \Leftarrow l, \text{dist}=1$

where $\pi(p)$ relabels domain elts randomly

$\pi(p_0), \pi(p_0)$ applies same relabeling to both

Main idea: Only Collision Statistics matter!

for positive pairs have collisions in both heavy + light elts

for negative pairs have collisions only in heavy elts

when see a collision, usually can't tell if it was a heavy or light element!

After $o(n^{2/3})$ samples:

probability see any small element twice really small \leftarrow
 probability see any heavy element 3X is small \leftarrow happens, but not too often
 probability see any small elt 3X is tiny \leftarrow
 heavy " 4X is tiny \leftarrow unlikely to happen

So, what collision statistics could we have?

how many elts in domain appear n_p times, n_q times in p, q ?

P	0	0	1	0	2	1	0	3	1	2	4	0	3	1	2
q	0	1	0	2	0	1	3	0	2	1	0	4	1	3	2

#domain elts

will happen less in pos pairs than in neg pairs?

will happen more in pos pairs than in neg pairs?

only heavy elements - same distribution for pos + neg pairs

unlikely - can ignore

When you see collision, you don't know if it came from heavy or light element

$2m = \#$ samples (pretend m heavy, m light) * cheat #1
 $H = \#$ heavy collisions
 $L = \#$ light collisions (1 from each dist)

\leftarrow same distribution for pos + neg pairs

$\leftarrow = 0$ when neg pair

$$E[\# \text{ collisions in pos pair}] = E[H] + E[L] = \frac{m^2}{2n^{2/3}} + \frac{m^2}{n} \sim \frac{m^2}{2n^{2/3}}$$

$$E[\# \text{ collisions in neg pair}] = E[H] = \frac{m^2}{2n^{2/3}}$$

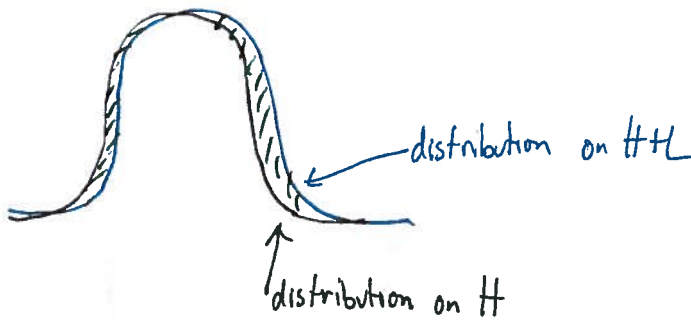
Need to show something a bit stronger - can't distinguish the random variables!

$$E[H] = \frac{m^2}{8n^{2/3}} \approx m^2 \text{ pairs, each collides with prob } \frac{1}{2n^{2/3}}$$

$$\text{Var}[H] \approx \frac{m^2}{n^{2/3}}$$

$$E[L], \text{Var}[L] \approx \frac{m^2}{n} \approx m^2 \text{ pairs, each collides with prob } \frac{1}{n}$$

L_1 distance = small
 \Downarrow
 almost same distribution
 \Downarrow
 hard to distinguish!



how do we show L_1 dist is small?

if they were gaussian,
 could show that $\sqrt{\text{Var}(H)} \leq E[L]$

$$\Leftrightarrow \frac{m}{n^{2/3}} \leq \frac{m^2}{n}$$

$$\Leftrightarrow m \geq n^{2/3}$$

* cheat #2
 \uparrow
 bigger difficulty!
 \Leftarrow they aren't quite, so it's more difficult.