# Lower Bounds on distributions

last time: sketch of lower bound for uniformity testing

Homework: One way of making it formal (not optimal in all parameters)

Today: Another methodology of showing lower bounds

def Uniformity tester
given samples from $p$ on $[n]$, $\varepsilon$
- if $p = U_n$ output PASS with prob $\geq 3/4$ ⟍ any constant $> 1/2$ to do better than random guess
- if $\|p - U_n\|_1 > \varepsilon$ output FAIL with prob $\gtrsim 3/4$

Thm Uniformity tester needs $\Omega\left(\sqrt{n}/\varepsilon^2\right)$ samples

Proof soon, 1st some observations + basics:

Observation: randomness doesn't help testing algorithms
Pf: h.w.

Information Theory Basics:

$p(y|x) \equiv P(Y=y|X=x)$

Entropy $\qquad H(x) = -\sum_{x \in \text{domain}} p(x) \log p(x)$

Conditional Entropy $\quad H(Y|X) = E_x\left[ \sum_{\substack{y \text{ st.} \\ p(y) \neq 0}} p(y|x) \log \frac{1}{p(y|x)} \right]$

$$= \sum_x p(x) \sum_{\substack{y \text{ st.} \\ p(y) \neq 0}} p(y|x) \log \frac{1}{p(y|x)}$$

Note:
$H(Y|X) = 0$ iff $Y$ determined by $X$
$H(Y|X) = H(Y)$ iff $Y$ independent of $X$

Basic facts:

- $H(x) \geq 0$

- $H(Y|X) \leq H(Y)$

- Chain rule: $H(X,Y) = H(X) + H(Y|X)$
  <span style="color:purple">joint entropy of pair (x,y)</span>

**Mutual Information:**

$$I(X,Y) = H(X) + H(Y) - H(X,Y)$$
$$= H(X) - H(X|Y)$$
$$= H(Y) - H(Y|X)$$

<span style="color:purple">measure of how independent $X,Y$ are or how much $X$ allows you to predict $Y$</span>

Chain rule: $I(X;(Y,Z)) = I(X;Z) + I(X;Y|Z)$

**Main Idea:**

define random var $X$ as fair coin flip

$X$ decides whether pick $K$ samples from ⎡ uniform on $[n]$

<span style="color:purple">⇑ all $K$ from same distribution</span>

uniform on $S$ st. $|S| = \frac{n}{2}$
+ $S$ chosen randomly

<span style="color:purple">we get samples, not $X$. can we figure out what $X$ is from samples?</span>

Will show, if $\mathbf{K}$ small, $I(X, \text{samples}) = o(1)$

<span style="color:green">any improvement over random gives ↓ of c...</span>

<span style="color:green">So what?</span>

**Lemma** if $f$ any fctn (algorithm) s.t. $\Pr_{X, \text{samples}}\left[f(\text{samples}) = X\right] \geq 51\%$

then $I(X;A) \geq 2 \cdot 10^{-4}$

<span style="color:green">So if $I(X, \text{samples}) = o(1) \implies$ no algorithm can solve the testing problem with high enough prob...</span>

$a_i \leftarrow \#$ times elt $i$ appears in sample

$$I(x, \text{samples}) = I(x, \{a_x\}_{i=1}^n)$$

Lets assume $a_i$'s independent  (they are not if $K$ is fixed, but if $K$ chosen as Poisson dist with mean $K_0$, they are independent)

$$I(x, \{a_i\}_{i=1}^n) \leq \sum_{i=1}^n I(x, a_i) \quad \text{by chain rule}$$

↑ drawn identically $\forall i$

$$\stackrel{\leq}{=} n \cdot I(x, a_1) = O\left(\frac{K^2 \varepsilon^4}{n}\right)$$

Lemma  $I(x, a_1) = O\left(\frac{K^2 \varepsilon^4}{n^2}\right)$

Proof: calculations

if $K = \dot{O}\left(\frac{\sqrt{n}}{\varepsilon^a}\right)$

this is $\dot{O}(1)$

## Poissonization

An important way to get rid of dependencies.

why:
if take fixed $K$ # of samples

$Pr[$ see elt $i$ $]$ not independent of $Pr[$ see elt $j]$.

why? if you see elt $i$, you know 1 sample
is not $j$, so less likely you
will see elt $j$ in all $K$
samples (you now only have $K-1$
samples left to "play with").

Poissonization trick:

pick $K$ distributed as Poisson with parameters

def. Poisson dist with parameter $\lambda$ $(\Psi(\lambda))$:

$K$ occurs with prob $\dfrac{\lambda^k e^{-\lambda}}{k!}$ ← Note:
$0! = 1$
$\Psi(0) = 0$

Observe

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

$$E[X] = \lambda \qquad \text{for } X \leftarrow \Psi(\lambda)$$

$$Var[X] = \lambda$$

Poisson Sampling: pick $K \sim \Psi(\lambda)$
take $K$ samples of distribution

Important property of Poisson Sampling:

- \# of occurences of elt $i$ is <u>independent</u> of
  " " " " " $j$    (for $i \neq j$)

- \# of occurences of elt $i \sim \Psi(K \cdot p_i)$
  $$E[\text{"} \quad \text{"} \quad \text{"} \quad \text{"} \quad \text{"} \quad \text{"} \quad ] = K \cdot p_i$$
  $$Var[\text{"} \quad \text{"} \quad \text{"} \quad \text{"} \quad \text{"} \quad \text{"} \quad ] = K \cdot p_i$$

Why does this give us a lower bound?

Suppose you want to show $\geq S_0$ samples are required for a testing problem.

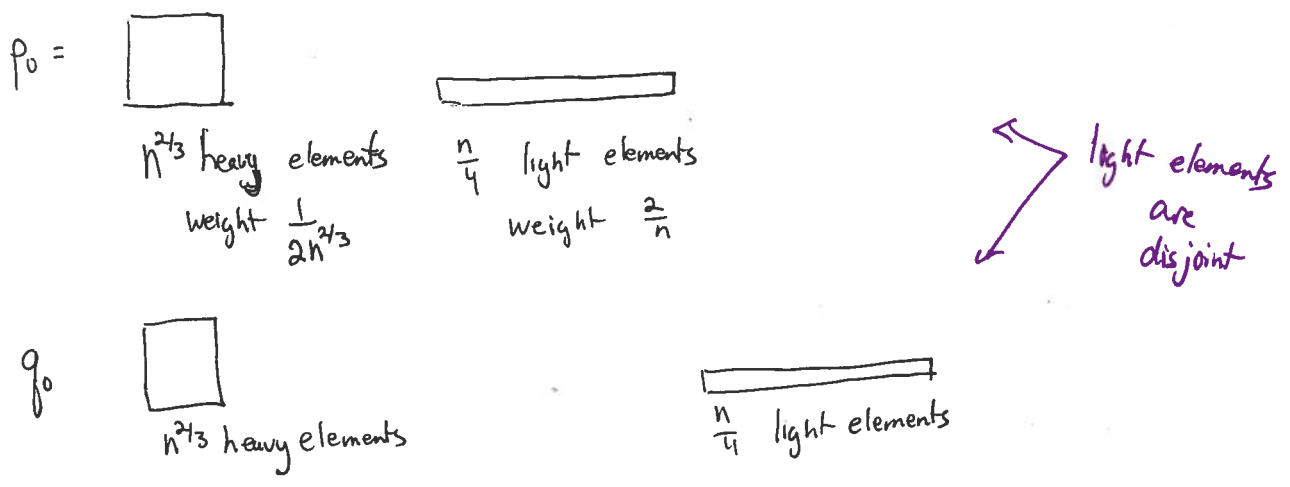i.e. $\forall \mathcal{A}$ taking $S_0$ samples, $\mathcal{A}$ correct with probabilty $\geq 2/3$

$$\Downarrow$$

$\forall \mathcal{A}'$ taking $\Psi(c \cdot S_0)$ samples, $\mathcal{A}'$ correct with prob $\geq 2/3 - \text{"tiny"}$

<span style="color:green">$\underbrace{\quad}_{<>1}$
expectation $c \cdot S_0$
prob \# samples $< S_0$
is "tiny"</span>

<u>Contrapositive</u>: if $\mathcal{A}'$ needs $\geq \Psi(c \cdot S_0)$ samples
then $\mathcal{A}$ needs $\geq S_0$ samples

Sketch of l.b. for $p, q$ given by samples ⟸ "closeness testing"

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

<u>Proof idea</u>:

$p_0 = \square$

$n^{2/3}$ heavy elements
weight $\frac{1}{2n^{2/3}}$

$\frac{n}{4}$ light elements
weight $\frac{2}{n}$

light elements are disjoint

$q_0$ $\square$

$n^{2/3}$ heavy elements

$\frac{n}{4}$ light elements

Positive pairs

$\ell_1 \, dist = 0 \Rightarrow$ $(\pi(p_0), \pi(p_0)) \; \forall \pi$

Negative pairs

$(\pi(p_0), \pi(q_0)) \quad \forall \pi$ ⟸ $\ell_1 \, dist = 1$

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

$\pi(p_0), \pi(p_0)$ applies <u>same</u> relabeling to both

Main idea: ⌐Only Collision Statistics matter!¬ for positive pairs have collisions in both heavy + light elts

for negative pairs have collisions <u>only</u> 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 ⟸

    probability see any heavy element 3X is small ⟸ happens, but not too often

    probability see any small elt 3X is tiny ⟸

    heavy " 4X is tiny ⟸ 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

$m$ = # samples

$H$ = # heavy collisions

$L$ = # light collisions (1 from each dist) ⟵ = 0 when neg pair

⟵ same distribution for pos & neg pairs

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

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

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

$$E[H] = \frac{m^2}{4n^{2/3}}$$

$\binom{m}{2}$ pairs, each collides with prob $\frac{1}{2n^{2/3}}$

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

$$E[L], Var[L] \approx \frac{m^2}{n}$$

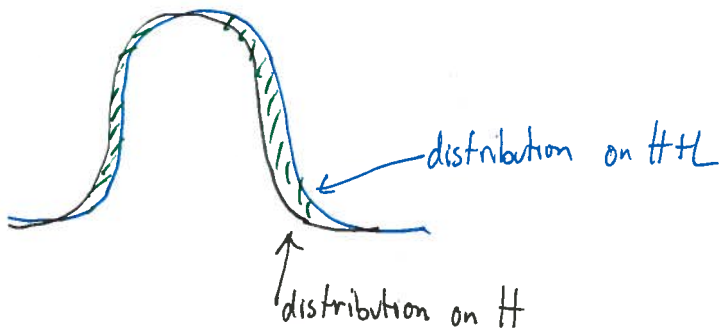$\binom{m}{2}$ pairs, each collides with prob $\frac{1}{n}$

$L_1$ distance small
⇓
almost same distribution
⇓
hard to distinguish!



—distribution on H+L

↑ distribution on H

how do we show $L_1$ dist is small?

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

⇐ they aren't quite, so it's more difficult.

$$\Leftrightarrow \frac{m}{n^{1/3}} \leq \frac{m^2}{n}$$

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