

- Closeness Testing (p, q unknown)
- Learning & Testing monotone distributions

Some other extensions:

What if p, q both unknown? "Closeness testing"

L_2 distance is similar, but what does it say?

$$L_2 \text{ distance: } \|p - q\|_2^2 = \sum_i (p_i - q_i)^2$$

$$= p_i^2 - 2 \sum_i p_i q_i + q_i^2$$

\uparrow \uparrow \uparrow
 self-collision prob of p cross-collision probability of $p+q$ self-collision prob of q

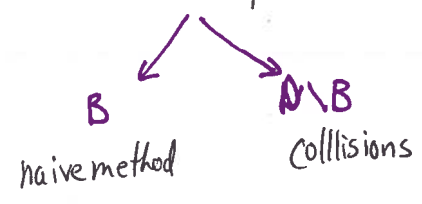
• Can bound variance of $\|p\|_2^2, \sum p_i q_i + \|q\|_2^2$ estimators if max prob element is bounded by b

• what about other case?

Use naive method on elements whose prob $\geq b$
 $\leq \frac{1}{b}$ of these

Oneway: Filtering algorithm:

learn B = domain elements with prob $\geq b \leftarrow O(\frac{1}{b})$ samples
filter rest of samples



Note strange dependence on $n!$

$n^{2/3}$ is fight!! \rightarrow Turns out

$$O\left(\frac{1}{\epsilon^2} \cdot \frac{1}{b}\right) \text{ samples}$$

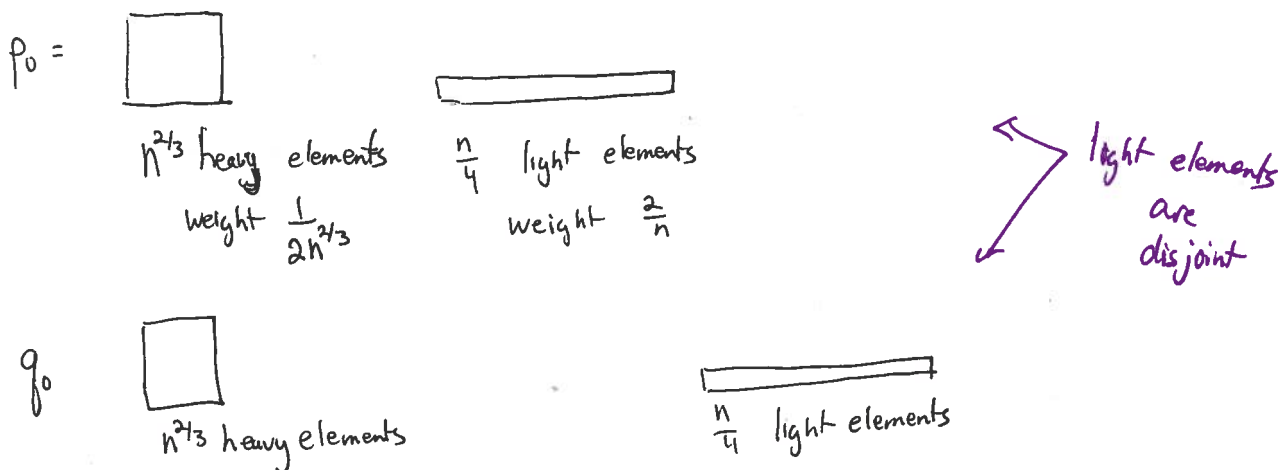
$$O\left(\frac{1}{\epsilon^4} n^{2/3}\right)$$

samples suffice [recent improvements on ϵ known]

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

Thin closeness testing requires $\Omega(n^{2/3})$ samples

Proof idea:



Positive pairs

Negative pairs

$l_1 \text{ dist} = 0 \Rightarrow (\pi(p_0), \pi(p_0)) \forall \pi$

 $(\pi(p_0), \pi(q_0)) \forall \pi \Leftarrow l_1 \text{ dist} = 1$

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

$m = \#$ samples

$H = \#$ heavy collisions

$L = \#$ light collisions (1 from each dist)

\leftarrow same distribution for pos + neg pairs

$\leftarrow = 0$ when neg pair

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

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

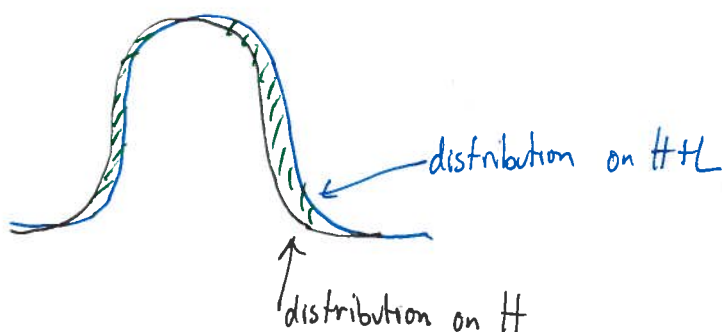
Need to show something a bit stronger - can't distinguish the random variables!

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

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

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

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 $\sqrt{\text{Var}(H)} \leq E[L]$

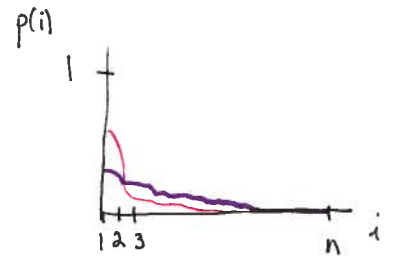
\Leftarrow they aren't quite, so it's more difficult.

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

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

Testing & Learning Monotone Distributions (over totally ordered domain)

Def. p over $[n]$ is "monotone decreasing"
if $\forall i \in [n-1] \quad p(i) \geq p(i+1)$



Monotonicity Tester:

- if p monotone increasing, Pass with prob $\geq 3/4$
- if p ϵ -far in L_1 dist from mon increasing, Fail with prob $\geq 3/4$

Useful tool: "Birge Decomposition"

(note: this is a different decomposition than in homework
in particular, it is oblivious!)

decompose domain $1..n$ into $l = \Theta\left(\frac{\log \epsilon n}{\epsilon}\right) \approx \Theta\left(\frac{\log n}{\epsilon}\right)$ intervals

$$I_1^\epsilon, I_2^\epsilon, \dots, I_l^\epsilon \quad \text{s.t.}$$

$$|I_{k+1}^\epsilon| = \lceil (1 + \epsilon/2) \cdot |I_k^\epsilon| \rceil$$

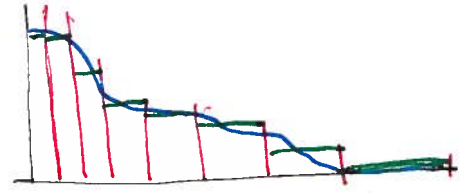
← will drop ϵ
in notation
once it is fixed

so $|I_1^\epsilon| = 1$
 $|I_2^\epsilon| = 2$
 $|I_3^\epsilon| = 3$
but then at some point the sizes grow
exponentially

define "flattened distribution"

$$\forall 1 \leq j \leq l$$

$$\forall i \in I_j \quad \tilde{q}_\epsilon(i) = \frac{q(I_j)}{|I_j|}$$



← assign all elements in same interval the same probability

note: $q(I_j) = \tilde{q}_\epsilon(I_j)$

Thm if q mon decreasing then $\|\tilde{q}_\epsilon - q\|_1 < \epsilon$

Coroll if q ϵ -close to mon decreasing then $\|\tilde{q}_\epsilon - q\|_1 < O(\epsilon)$

Testing Algorithm:

Take samples of q
do uniformity test for each partition (using samples that fell in it)
(if not enough samples then pass) fail if any partition fail

$w_j \leftarrow$ # samples that fell in partition j
use LP to verify w close to monotone

* note this is LP on $O(\log n)$ vars

How many samples?

for each partition with enough weight, say $\frac{\epsilon}{\log n}$, need $\frac{\sqrt{n}}{\epsilon^2}$ samples

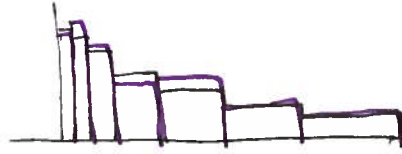
$$\approx O(\sqrt{n} \text{ polylog } n \cdot \text{poly } \frac{1}{\epsilon})$$

← need $\frac{\sqrt{n} \cdot \log n}{\epsilon}$ for each ϵ
need another $\log \log n$ for union bound

(note: this can be improved !!)

Last step:

difficulty



purple is not monotone
but is close

good thing: only $\frac{\log n}{\epsilon}$ variables!

can be solved via brute force
LP (actually quite efficient)
⋮

Slightly changing perspective...

What if we know dist q is monotone, can we learn it?

Yes! use sampling to estimate $\tilde{q}_\epsilon(I_j)$'s

Birge's Thm Can learn monotone distributions to w/in ϵ L_1 error
in $\Theta(\frac{1}{\epsilon^3} \log n)$ samples.