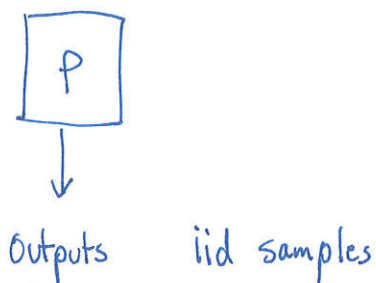


Model



domain D , $|D| = n$ ← known
 $D = \{1 \dots n\} = [n]$ for today
 $p_i = \Pr[P \text{ outputs } i]$
 unknown

Testing Uniformity

- if $P \equiv U_{[n]}$ tester outputs PASS
 - if $\text{dist}(P, U_{[n]}) > \epsilon$ then tester outputs FAIL
- with prob $\geq 3/4$
- distance measures:

$$l_1: \|p - q\|_1 = \sum_{i \in D} |p_i - q_i|$$

$$l_2: \|p - q\|_2 = \sqrt{\sum_{i \in D} (p_i - q_i)^2}$$

$$\|p - q\|_2 \leq \|p - q\|_1 \leq \sqrt{n} \|p - q\|_2$$

Last time: "plug in" estimate seems to need $\Omega(n)$ samples

L₂ - Distance (squared):

$$\begin{aligned} \|p - u\|_2^2 &= \sum_{i \in [n]} (p_i - \frac{1}{n})^2 \\ &= \sum p_i^2 - \frac{2}{n} \sum p_i + \sum \underbrace{\left(\frac{1}{n}\right)^2}_{=\frac{1}{n}} \\ &= \sum p_i^2 - \frac{1}{n} \end{aligned}$$

Collision probability of p :

$$\|p\|_2^2 \equiv \Pr_{s, t \in p} [s = t] = \sum p_i^2$$

for $p = u$, $\|p\|_2^2 = \frac{1}{n}$

for $p \neq u$, $\|p\|_2^2 > \frac{1}{n}$

$$= \|p\|_2^2 - \underbrace{\|u\|_2^2}_{\frac{1}{n}}$$

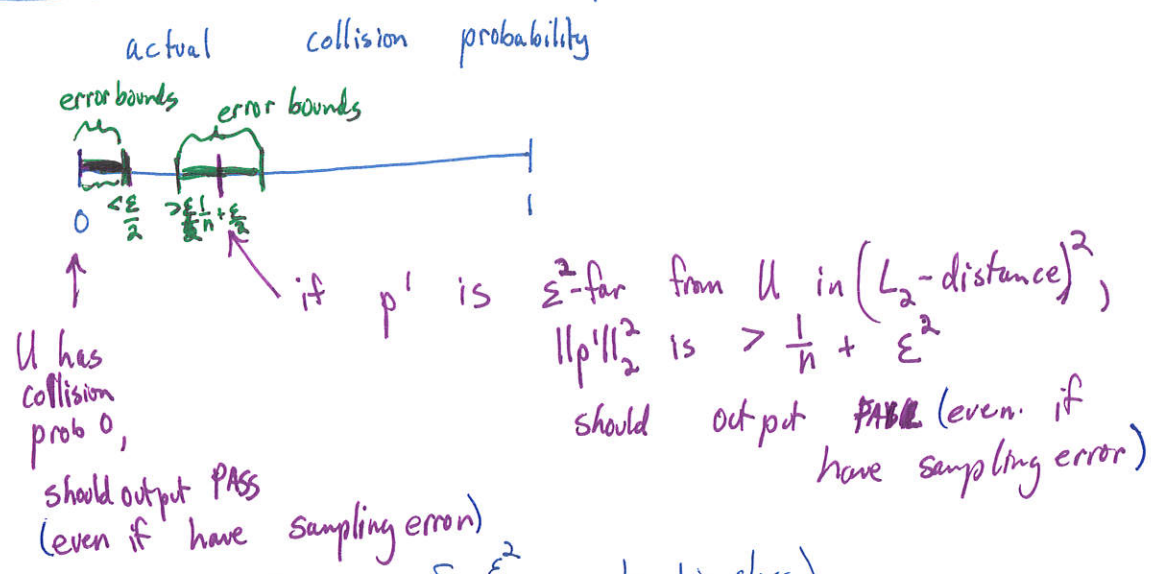
we can estimate this

we know this since we know n

Algorithm

1. take s samples from p ① how many samples?
2. let $\hat{c} \leftarrow$ estimate of $\|p\|_2^2$ from sample ② how?
3. if $\hat{c} < \frac{1}{n} + \delta$ pass ③ what should δ be?
 else fail

How well do we need to estimate $\|p\|_2^2$?



- Use $\delta \leq \frac{\epsilon^2}{2}$ (we used $\delta = \frac{\epsilon^2}{4}$ on board in class)
- assume (with prob $\geq 3/4$) estimate of collision prob $\|p\|_2^2$ is within $\pm \delta$ for $\delta \leq \frac{\epsilon^2}{2}$

Proper Behavior {

- Then if $p=U$, $\hat{c} < \frac{1}{n} + \frac{\epsilon^2}{2}$ \Rightarrow PASS
- if $\|p-U\|_2^2 > \epsilon^2$ then $\hat{c} > \frac{1}{n} + \epsilon^2 - \frac{\epsilon^2}{2} = \frac{1}{n} + \frac{\epsilon^2}{2}$ \Rightarrow FAIL

• But how many samples to get prob $\geq 3/4$ of good estimate?

[What if want better estimate of $\|p\|_2^2$, to within $\frac{\epsilon^2}{3n}$]

[see later for estimate of L_1]

How well do we need to estimate $\|p\|_2^2$?

Assumption \star : $|\hat{C} - \|p\|_2^2| < \Delta$
 will take enough samples so that this holds with prob $\geq 3/4$
 this is our parameter that determines whether our approximation is good. Spoiler: will set $\Delta = \frac{\epsilon^2}{2}$

What happens if \star holds with $\Delta = \frac{\epsilon^2}{2}$?

• if $p = U_{[n]}$ then $\hat{C} \leq \|U_{[n]}\|_2^2 + \Delta = \frac{1}{n} + \frac{\epsilon^2}{2}$

so test will PASS

• if $\|p - U_{[n]}\|_2 > \epsilon$ then $\|p - U_{[n]}\|_2^2 > \epsilon^2$

but then $\|p\|_2^2 = \|p - U_{[n]}\|_2^2 + \frac{1}{n}$

$> \epsilon^2 + \frac{1}{n}$

+ $\hat{C} > \|p\|_2^2 - \Delta$

$\geq \epsilon^2 + \frac{1}{n} - \Delta = \epsilon^2 + \frac{1}{n} - \frac{\epsilon^2}{2} = \frac{\epsilon^2}{2} + \frac{1}{n}$

so test will FAIL

How many samples do we need to estimate \hat{C} to within Δ ?

First:

How to estimate $\|p\|_2^2$?

Naive idea:

take two new samples:

$$X_i \leftarrow \begin{cases} 1 & \text{if samples are equal} \\ 0 & \text{o.w.} \end{cases}$$

" gives $\Theta(k)$ samples of collision probability
from k samples of p "

Better idea: recycle - use all pairs in sample

" gives $\Theta(k^2)$ samples of collision probability
from k samples of p "

Estimate by recycling:

- Take s samples from p : X_1, \dots, X_s

- for each $1 \leq i < j \leq s$

$$b_{ij} \leftarrow \begin{cases} 1 & \text{if } X_i = X_j \\ 0 & \text{if } X_i \neq X_j \end{cases}$$

- Output $\hat{c} \leftarrow \frac{\sum_{i < j} b_{ij}}{\binom{s}{2}}$

b_{ij} 's not
independent
so can't use
Chernoff

Analysis: $E[\hat{c}] = \frac{1}{\binom{s}{2}} \cdot \binom{s}{2} \cdot E[b_{ij}]$
 $= \|p\|_2^2$

Analysis

$$E [b_{ij}] = \Pr [b_{ij}] = 1 \\ = \|p\|_2^2$$

$$E [\hat{c}] = \frac{1}{\binom{s}{2}} \binom{s}{2} E [b_{ij}] = \|p\|_2^2$$

$$\Pr [|\hat{c} - \|p\|_2^2| > \rho] \leq \frac{\text{Var} [\hat{c}]}{\rho^2}$$

Chebyshev \neq

Fact $\text{Var} [aX] = a^2 \text{Var} [X]$

$$\text{So } \text{Var} [\hat{c}] = \text{Var} \left[\frac{1}{\binom{s}{2}} \sum_{i < j} b_{ij} \right] \\ = \frac{1}{\binom{s}{2}^2} \text{Var} \left[\sum_{i < j} b_{ij} \right]$$

Lemma $\text{Var} \left[\sum b_{ij} \right] \leq 2 \left(\binom{s}{2} \|p\|_2^2 \right)^{3/2}$

Why? (proof...)

def. $\bar{b}_{ij} = b_{ij} - E[b_{ij}]$

so $E[\bar{b}_{ij}] = 0$

Also: $E[\bar{b}_{ij} \bar{b}_{kl}] \leq E[b_{ij} b_{kl}]$

Verify at home? (or trust...)

- $\left(\sum p(x)^3 \right)^{1/3} \leq \left(\sum p(x)^2 \right)^{1/2}$
- $s^2 \leq 3 \binom{s}{2}$
- $\binom{s}{3} \leq s^3/6$

e.g. $(a^3 + b^3)^2 \leq (a^2 + b^2)^3$
 $a^6 + 2a^3b^3 + b^6 \leq a^6 + b^6 + 3a^4b^2 + 3a^2b^4$

So

$$\text{Var} \left[\sum_{i < j} \delta_{ij} \right] = E \left[\left(\sum_{i < j} \delta_{ij} - E \left[\sum_{i < j} \delta_{ij} \right] \right)^2 \right]$$

$$= E \left[\left(\sum_{i < j} \bar{\delta}_{ij} \right)^2 \right]$$

$$= E \left[\sum_{i < j} \bar{\delta}_{ij}^2 + \sum_{\substack{i < j \\ k < l \\ i, j, k, l \text{ distinct}}} \bar{\delta}_{ij} \bar{\delta}_{kl} + \sum_{\substack{i < j \\ k = l \\ i, j, l \text{ distinct}}} \bar{\delta}_{ij} \bar{\delta}_{kl} + \sum_{\substack{i < j \\ k < l \\ i, k, j \text{ distinct}}} \bar{\delta}_{ij} \bar{\delta}_{kl} \right]$$

①
②
③
④

① $E \left[\sum_{i < j} \bar{\delta}_{ij}^2 \right] \leq E \left[\sum \delta_{ij}^2 \right] = \binom{s}{2} \|p\|_2^2$

$E[\delta_{ij}] = E[\delta_{ij}^2]$ since δ_{ij} is indicator var

② \swarrow independent

$$E \left[\sum_{\substack{i < j \\ k < l \\ \text{all 4 distinct}}} \bar{\delta}_{ij} \bar{\delta}_{kl} \right] \leq \sum E[\bar{\delta}_{ij}] E[\bar{\delta}_{kl}] = 0$$

③ $E \left[\sum \bar{\delta}_{ij} \bar{\delta}_{il} \right] \leq E \left[\sum_{\substack{i, j, l \\ \text{distinct}}} \delta_{ij} \delta_{il} \right] = \sum_{\substack{i, j, l \\ \text{distinct}}} \text{pr}[x_i = x_j = x_l]$

$$\leq \binom{s}{3} \sum_x p(x)^3 \quad \text{expected \# 3-way collisions}$$

$$\frac{1}{6} \binom{s}{3}^{3/2} < \frac{\left(3 \binom{s}{2} \right)^{3/2}}{6} = \frac{\sqrt{3}}{2} \binom{s}{2}^{3/2}$$

$$\leq \frac{s^3}{6} \left(\sum_x p(x)^2 \right)^{3/2}$$

$$\leq \frac{\sqrt{3}}{2} \binom{s}{2}^{3/2} \left(\|p\|_2^2 \right)^{3/2} \quad \text{by the facts}$$

(4) same as 3

In total:

$$\begin{aligned} \text{Var} \left[\sum_{i < j} \delta_{ij} \right] &\leq \text{Var} \left[\sum_{i < j} \bar{\delta}_{ij} \right] \\ &\leq \binom{s}{2} \|\phi\|_2^2 + 0 + 2 \cdot \frac{\sqrt{3}}{2} \left(\binom{s}{2} \|\phi\|_2^2 \right)^{3/2} \\ &\leq 2 \left[\binom{s}{2} \|\phi\|_2^2 \right]^{3/2} \end{aligned}$$



Putting lemma into Chebyshev

1) use $p = \frac{\epsilon^2}{2}$

$$\Pr[|c - \|p\|_2^2| > \frac{\epsilon^2}{2}] \leq \frac{\text{Var}[\hat{c}] \cdot 4}{\epsilon^4}$$

$$\leq \frac{2 \left[\binom{s}{2} \|p\|_2^2 \right]^{3/2}}{\binom{s}{2}^2 \epsilon^4} \cdot 4 = \frac{8}{\epsilon^4} \cdot s^3 \cdot \|p\|_2^3$$

$\underbrace{\hspace{10em}}_{\text{also want this to be } \leq 1}$

So pick $s \geq \frac{1}{\epsilon^{4/3}}$

Note: Can get better bound if have bound on $\|p\|_\infty$
 \uparrow
 max prob element

In homework:

- 1) Testing closeness to any known distribution — reduce to uniform case!
- 2) lower bound

How to estimate $\|p-u\|_1$?

$$1) \|p-u\|_1 = 0 \Leftrightarrow \|p-u\|_2^2 = 0 \Leftrightarrow \|p\|_2^2 = \frac{1}{n}$$

$$2) \text{ if } \|p-u\|_1 > \varepsilon \Rightarrow \|p-u\|_2 > \frac{\varepsilon}{\sqrt{n}}$$

$$\Rightarrow \|p-u\|_2^2 > \frac{\varepsilon^2}{n}$$

$$\Rightarrow \|p\|_2^2 > \frac{1}{n} + \frac{\varepsilon^2}{n}$$

either additive estimate with error $\leq \frac{\varepsilon^2}{2n}$

or mult error $\leq (1 \pm \frac{\varepsilon^2}{3n})$

suffices

would have this
 if have
 additive error $\leq \frac{\varepsilon^2}{3n} \cdot \|p\|_2^2$

to get additive error $\leq \frac{\varepsilon^2}{3n} \|p\|_2^2$

suffices to have

$$s \geq \frac{\text{const} \cdot \sqrt{n}}{\varepsilon^2}$$

samples

Some other extensions:

What if p, q both unknown?

L_2 distance is similar, but what does it say?

near tests

$$\begin{aligned}
 L_2 \text{ distance: } \|p - q\|_2^2 &= \sum_i (p_i - q_i)^2 \\
 &= \sum_i p_i^2 - 2 \sum_i p_i q_i + \sum_i q_i^2 \\
 &\quad \uparrow \qquad \qquad \qquad \uparrow \\
 &\quad \|p\|_2^2 \qquad \text{cross-collision probability of } p+q \qquad \|q\|_2^2 \\
 &\quad \text{self-collision prob of } p \qquad \qquad \qquad \text{self-collision prob of } q
 \end{aligned}$$

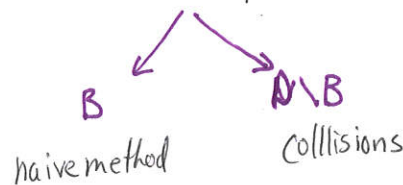
• can bound variance of $\|p\|_2^2, \sum p_i q_i + \|q\|_2^2$ estimators if max prob element is bounded by b

• what about other case?

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

Filtering algorithm:

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



Note strange dependence on $n!$

$n^{2/3}$ is tight!! \rightarrow Turns out $O(\frac{1}{\epsilon^4} n^{2/3} \log n)$ samples suffice [recent improvements on $\epsilon, \log n$ known]