

Lecture 12

Testing distributions:

the case of uniformity

A new model:

Probability distributions: get samples



this is
all we
see



iid samples

Discrete Domain D st. $|D|=n$

$$P_i = \Pr[p \text{ outputs } i] \leftarrow \text{unknown}$$

Examples:

- lottery data
- Shopping choices
- experimental outcomes
- ⋮

What do we need to know?

is it

uniform?

e.g. lottery

high entropy?

large support?

(many distinct elts
with > 0 probability)

monotone increasing, K -modal?

K -histogram?

Methods ?

learn distribution

χ^2 -test

plug-in estimate

Maxlikelihood estimate

Goal : sample complexity sublinear in n

Testing Uniformity

uniform dist on D

goal: if $p \equiv U_D$ then output PASS

if $\text{dist}(p, U_D) > \epsilon$ then output FAIL

with prob $\geq 3/4$

which measure of distance?

l_1, l_2 , KL-divergence, Earthmover, Jensen-Shannon ...

today's focus

Distances

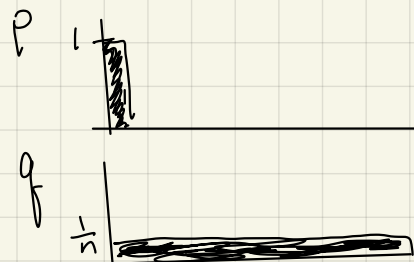
l_1 -distance: $\|p-q\|_1 = \sum_{i \in D} |p_i - q_i|$

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

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

examples:

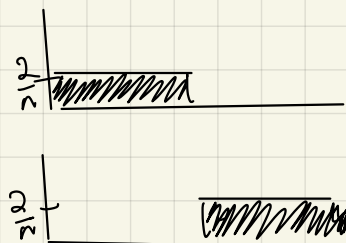
① $p = (1, 0, 0, 0, \dots, 0)$
 $q = (\frac{1}{n}, \frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n})$



$$\|p-q\|_1 = \frac{n-1}{n} + (n-1) \cdot \frac{1}{n} \approx 2$$

$$\|p-q\|_2 = \left(\frac{n-1}{n}\right)^2 + (n-1) \left(\frac{1}{n}\right)^2 \approx 1$$

② $p = (\frac{2}{n}, \frac{2}{n}, \frac{2}{n}, \dots, \frac{2}{n}, 0, 0, \dots, 0)$
 $q = (0, 0, \dots, 0, \frac{2}{n}, \frac{2}{n}, \dots, \frac{2}{n})$



$$\|p-q\|_1 = n \cdot \frac{2}{n} = 2$$

$$\|p-q\|_2 = n \cdot \left(\frac{2}{n}\right)^2 = \frac{4}{n} \text{ so } \|p-q\|_2 = \frac{2}{\sqrt{n}}$$

small
 even
 though
 p & q
 very
 different

Via "Plug-in" Estimate:

- take m samples from p

- estimate $p(x) \forall x$ via $\hat{p}(x) = \frac{\# \text{ times } x \text{ occurs in sample}}{m}$

- if $\sum_x |\hat{p}(x) - \frac{1}{n}| > \varepsilon$ reject

else accept

Naive Analysis: (better analyses exist)

pick m s.t. $\forall x \quad |\hat{p}(x) - p(x)| < \frac{\varepsilon}{n} \Rightarrow \|\hat{p} - p\|_1 < \varepsilon$

so if $p=U$,
likely to pass

by $\Delta \neq$, if $\|p - \hat{p}\|_1 < \varepsilon$ + $\|\hat{p} - U\|_1 < \varepsilon$ then $\|p - U\|_1 < 2\varepsilon$

how big should m be?

maybe need to see each x enough to get $\frac{\varepsilon}{n}$ -additive error on $\hat{p}(x) \Rightarrow \Omega(m)$? $\Omega(n/\varepsilon)$? $\Omega(n/\varepsilon^2)$? $\Omega(\frac{n}{\varepsilon^3})$?

so if $\|p - U\|_2 > 2\varepsilon$
likely to fail

Can do better - don't need good approx for all x

Claim $E[\|\hat{p} - p\|_1] \leq \sqrt{\frac{n}{m}}$

Pf of claim

$$E[\|\hat{p} - p\|_1] = \sum_x E[|\hat{p}(x) - p(x)|]$$

$$\leq \sum_x \sqrt{E[(\hat{p}(x) - p(x))^2]}$$

$$= \sum_x \sqrt{\text{Var}(\hat{p}(x))}$$

$$\leq \sum_x \sqrt{\frac{p(x)}{m}}$$

$$\leq \frac{1}{\sqrt{m}} \cdot \sqrt{n}$$



so pick $m = \Theta(\frac{n}{\epsilon^2})$ gives $E[\|\hat{p} - p\|_1] \leq \frac{\epsilon}{c}$

by Markov's \neq , with prob $1 - \frac{1}{c}$

$$\|\hat{p} - p\|_1 \leq \epsilon$$

$$E[\hat{p}(x)] = \frac{1}{m} E\left[\sum_{i=1}^m \mathbb{1}_{i^{\text{th}} \text{ sample is } x}\right]$$

$$= \frac{1}{m} \sum_{i=1}^m E[\mathbb{1}_{i^{\text{th}} \text{ sample is } x}]$$

$$= \frac{m}{m} \cdot p(x) = p(x)$$

Jensen's \neq

$$\text{Var}(\hat{p}(x)) = \frac{1}{m^2} \cdot m \cdot p(x) \cdot (1 - p(x)) \leq \frac{p(x)}{m}$$

since $\max_p \sum \sqrt{p(x)}$ is \sqrt{n}

p
prob
distributions
over domain
of size n

So can "learn" (approximately) any distribution w.r.t. L_1 distance in $\Theta(\frac{n}{\epsilon^2})$ samples

Let's consider L_2 -distance (squared):

$$\|p - U_{[n]}\|_2^2 = \sum_{i \in [n]} \left(p_i - \frac{1}{n}\right)^2$$

$$= \sum p_i^2 - \frac{2}{n} \sum p_i + \sum \left(\frac{1}{n}\right)^2 = \sum p_i^2 - \frac{1}{n}$$

$\underbrace{\sum p_i}_{=1}$ $\underbrace{\sum \left(\frac{1}{n}\right)^2}_{=\frac{1}{n}}$

$$= \|p\|_2^2 - \|U_{[n]}\|_2^2$$

$\underbrace{\hspace{2em}}$
we can estimate this

$\underbrace{\hspace{2em}}$
we know this since we know n (it is $\frac{1}{n}$)

"Collision" probability of p :

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

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

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

Algorithm to estimate:

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

How well do we need to estimate $\|p\|_2^2$?
 i.e. what should δ be?

Assumption \star : $|\hat{c} - \|p\|_2^2| < \Delta$

will take enough
 samples s.t.
 this holds with
 prob $\geq 3/4$

this is our parameter
 that determines whether
 our approximation is good.

recall:

$$\|p - u_{[n]}\|_2^2 = \|p\|_2^2 - \|u_{[n]}\|_2^2$$

What 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 $\|p\|_2^2 = \|p - u_{[n]}\|_2^2 + \frac{1}{n} > \epsilon^2 + \frac{1}{n}$

$\downarrow \star \Rightarrow \hat{c} > \|p\|_2^2 - \Delta \geq \epsilon^2 + \frac{1}{n} - \frac{\epsilon^2}{2} = \frac{\epsilon^2}{2} + \frac{1}{n}$ so test will FAIL

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

- Naive idea:
- repeat several times;
 - take two samples & set $X_i \leftarrow \begin{cases} 1 & \text{if two samples equal} \\ 0 & \text{o.w.} \end{cases}$
 - output average of X_i 's
- $\Theta(k)$ samples of collisions from k samples of p

Better idea: "recycle" use all pairs in sample
gives $\Theta(k^2)$ samples of collision prob from k samples of p

- Take s samples from p : x_1, \dots, x_s

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

$$b_{ij} \leftarrow \begin{cases} 1 & \text{if } x_i = x_j \\ 0 & \text{o.w.} \end{cases}$$

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

b_{ij} 's not independent
so can't use Chernoff to argue that \hat{c} is close to $E[b_{ij}]$

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

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

Chebyshev's \neq

recall $\text{Var}[x] = E[(x - E[x])^2]$

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

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

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

} need to bound this but b_{ij} 's not independent

Lemma $\text{Var}\left[\sum_{i < j} b_{ij}\right] \leq \binom{s}{2} \|p\|_2^2 + 4 \cdot \left[\binom{s}{2} \|p\|_2^2\right]^{3/2}$

$$\Rightarrow \text{Var}(\hat{c}) = O\left(\frac{\|p\|_2^2}{s^2} + \frac{\|p\|_2^3}{s}\right)$$

Lemma $\text{Var} \left[\sum_{i < j} \delta_{ij} \right] \leq \binom{s}{2} \|p\|_2^2 + 4 \cdot \left[\binom{s}{2} \|p\|_2^2 \right]^{3/2}$

Proof

def $\bar{\delta}_{ij} = \delta_{ij} - E[\delta_{ij}]$

← trick:
why?

rewrite variance as $E[\sum \bar{\delta}_{ij}^2]$ ↙ = 0

$$\begin{aligned} \text{Var}[\sum \bar{\delta}_{ij}] &= E[(\sum \bar{\delta}_{ij} - E[\sum \bar{\delta}_{ij}])^2] \\ &= E[(\sum \delta_{ij} - E[\delta_{ij}])^2] \\ &= \text{Var}[\sum \delta_{ij}] \end{aligned}$$

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

Facts:

- $E[\bar{\delta}_{ij} \bar{\delta}_{kl}] \leq E[\delta_{ij} \delta_{kl}]$

- $\left(\sum_x p(x)^3 \right)^{1/3} \leq \left(\sum_x p(x)^2 \right)^{1/2}$

- $s^2 \leq 3 \binom{s}{2}$

- $\binom{s}{3} \leq s^3/6$

(Verify @ home)

So can equivalently bound $\text{Var}[\sum \bar{\delta}_{ij}]$

Lemma $\text{Var} \left[\sum_{i < j} \delta_{ij} \right] \leq \binom{s}{2} \|p\|_2^2 + 4 \cdot \left[\binom{s}{2} \|p\|_2^2 \right]^{3/2}$

Proof

$$\begin{aligned} \text{Var} \left[\sum_{i < j} \delta_{ij} \right] &= \text{Var} \left[\sum_{i < j} \bar{\delta}_{ij} \right] = E \left[\left(\sum_{i < j} \bar{\delta}_{ij} \right)^2 \right] \\ &= E \left[\underbrace{\sum_{i < j} \bar{\delta}_{ij}^2}_{(1)} + \underbrace{\sum_{\substack{i < j \\ k < l \\ i, j, k, l \text{ distinct}}}}_{(2)} \bar{\delta}_{ij} \bar{\delta}_{kl} + \underbrace{\sum_{\substack{i < j \\ i < l \\ i, j, l \text{ distinct}}}}_{(3)} \bar{\delta}_{ij} \bar{\delta}_{il} + \underbrace{\sum_{\substack{i < j \\ k < j \\ i, k, j \text{ distinct}}}}_{(4)} \bar{\delta}_{ij} \bar{\delta}_{kj} \right. \\ &\quad \left. + \underbrace{\sum_{i < j < l} \bar{\delta}_{ij} \bar{\delta}_{jl}}_{(5)} \right] \end{aligned}$$

Let's bound each term:

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

$\delta_{ij}^2 = \delta_{ij}$ since indicator var

$$\delta_{ij} \leftarrow \begin{cases} 1 & \text{if } X_i = X_j \\ 0 & \text{o.w.} \end{cases}$$

def $\bar{\delta}_{ij} = \delta_{ij} - E[\delta_{ij}]$

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

Facts:

- $E[\bar{\delta}_{ij} \bar{\delta}_{kl}] \leq E[\delta_{ij} \delta_{kl}]$
- $\left(\sum_x p(x)^3 \right)^{1/3} \leq \left(\sum_x p(x)^2 \right)^{1/2}$
- $s^2 \leq 3 \binom{s}{2}$
- $\binom{s}{3} \leq s^3/6$

$$(2) \quad E \left[\sum_{\substack{i < j \\ k < l \\ \text{all distinct}}} \bar{\delta}_{ij} \bar{\delta}_{kl} \right] \leq \sum E[\bar{\delta}_{ij}] \cdot E[\bar{\delta}_{kl}] = 0$$

this is where the trick helps - gets rid of lots of terms

(3) (+ similarly (4) + (5))

$$E \left[\sum_{\substack{i < j \\ i, j, l \text{ distinct}}} \bar{\delta}_{ij} \bar{\delta}_{il} \right] \leq E \left[\sum_{i, j, l \text{ distinct}} \delta_{ij} \delta_{il} \right] = \sum \Pr[X_i = X_j = X_l]$$

$$\leq \binom{5}{3} \sum_x p(x)^3$$

expected #
3-way collisions

$$\leq \frac{5^3}{6} \left(\sum_x p(x)^2 \right)^{3/2}$$

↳ by facts.

$$\leq \frac{\sqrt{3}}{2} \binom{5}{2}^{3/2} \left(\|p\|_2^2 \right)^{3/2}$$

$$\delta_{ij} \leftarrow \begin{cases} 1 & \text{if } x_i = x_j \\ 0 & \text{o.w.} \end{cases}$$

def $\bar{\delta}_{ij} = \delta_{ij} - E[\delta_{ij}]$

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

Facts:

- $E[\bar{\delta}_{ij} \bar{\delta}_{kl}] \leq E[\delta_{ij} \delta_{kl}]$

- $\left(\sum_x p(x)^3 \right)^{1/3} \leq \left(\sum_x p(x)^2 \right)^{1/2}$

- $5^2 \leq 3 \binom{5}{2}$

- $\binom{5}{3} \leq 5^3/6$

$$\begin{aligned} \text{So, } \text{Var} \left[\sum_{i < j} b_{ij} \right] &= \text{Var} \left[\sum_{i < j} \tilde{b}_{ij} \right] \\ &\leq \binom{s}{2} \|p\|_2^2 + 0 + 3 \cdot \frac{\sqrt{3}}{2} \left(\binom{s}{2} \|p\|_2^2 \right)^{3/2} \\ &\leq \binom{s}{2} \|p\|_2^2 + 4 \cdot \left[\binom{s}{2} \|p\|_2^2 \right]^{3/2} \end{aligned}$$



We have:

$$\text{Var}(\hat{c}) = O\left(\frac{\|p\|_2^2}{s^2} + \frac{\|p\|_2^3}{s}\right)$$

$$b_{ij} \leftarrow \begin{cases} 1 & \text{if } x_i = x_j \\ 0 & \text{o.w.} \end{cases}$$

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

where $s = \# \text{ samples}$

Put into Chebyshev with $p = \frac{\epsilon^2}{2}$:

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

$$\leq \frac{\text{const} \cdot \|p\|_2^2}{\epsilon^4 \cdot s^2}$$

want this ≤ 1
need $s = \Omega(1/\epsilon^2)$

$$+ \text{const} \cdot \frac{1}{\epsilon^4} \cdot \frac{1}{s} \cdot \|p\|_2^3$$

also want this $\ll 1$

$$\text{so pick } s = \Omega\left(\frac{1}{\epsilon^4}\right)$$

\uparrow
BIGGER CONSTRAINT

Note can get better bounds

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

recall:
 $\|p-u_{[n]}\|_2^2 = \|p\|_2^2 - \|u_{[n]}\|_2^2$

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

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

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

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

So either additive estimate of $\|p\|_2^2$ to within $\frac{\varepsilon^2}{2n}$
or mult estimate of $\|p\|_2^2$ to within $(1 \pm \frac{\varepsilon^2}{3})$
suffices

turns out that picking # samples $S \gg \frac{\sqrt{n}}{\varepsilon^4}$ suffices