

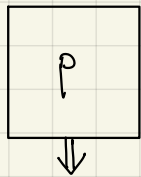
Lecture 12

Testing distributions:

the case of uniformity

A new model:

Probability distributions: get samples



this is
all
we
see

} iid samples

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

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

Examples:

- lottery data
- Shopping choices
- experimental outcomes
- ⋮

What do we need to know?

is it

uniform?

high entropy?

large support?

(many distinct elts
with > 0 probability)

monotone increasing, k -modal?

k -histogram?

Methods ?

learn distribution

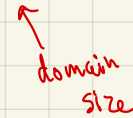
χ^2 -test

plug-in estimate

Maxlikelihood estimate

Goal : sample complexity sublinear in n

domain
size



Testing Uniformity

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

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

which measure
of distance?

$l_1, l_2, \text{KL-divergence, Earthmover, Jensen-Shannon} \dots$

today's focus

with prob
 $\approx 3/4$

Distances

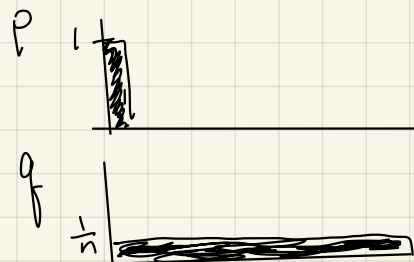
l_1 -distance: $\|p - q\|_1 = \sum_{i \in D} |p_i - q_i|$

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:

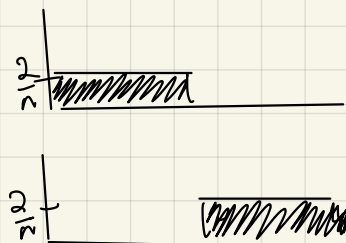
① $p = (1, 0, 0, 0, \dots, 0)$
 $q = (\frac{1}{n}, \frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n})$



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

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

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



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

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

$$\|p - q\|_2 = \frac{2}{\sqrt{n}}$$

tiny even though l_1 is big

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

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

so if $p = u$
likely to pass

if $\|p - u\|_1 > 2\varepsilon$
likely to fail

how big should m be?

do you need to see each x at least once? $\log n$ times?
 $\Omega(n)$? $\Omega(n/\varepsilon)$? $\Omega(n/\varepsilon^2)$? $\Omega(n \log 1/\varepsilon^2)$?

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

Pf of claim

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

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

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

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

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



so pick $m = \Theta(\frac{n}{\epsilon^2})$ gives $E[\|\hat{p} - p\|_1] \leq \frac{\epsilon}{c}$
 & by Markov's \neq with prob $1 - \frac{1}{c}$

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

Jensen's \neq

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

$$\text{Var}(\hat{p}(x)) = \frac{1}{m^2} \cdot \text{Var}[\# \text{ times } x \text{ occurs in sample}]$$

$$= \frac{1}{m^2} \text{Var}\left[\sum_{i=1}^m \mathbb{1}_{i\text{th sample is } x}\right]$$

$$= \frac{1}{m^2} \sum \text{Var}[\mathbb{1}_{i\text{th sample is } x}]$$

$$= \frac{1}{m^2} \cdot m \cdot \underbrace{p(x)(1-p(x))}_{\leq 1} \leq \frac{p(x)}{m}$$

indep for $i \neq j$

$$p(x)(1-p(x))$$

fact $\text{Var}(aX) = a^2 \cdot \text{Var}(X)$

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

$$\|p - u_{[n]}\|_2^2 = \sum_{i \in [n]} (p_i - \frac{1}{n})^2 = \sum (p_i^2 - \frac{2p_i}{n} + \frac{1}{n^2})$$

uniform on $1..n$

$$= \sum p_i^2 - \frac{2}{n} \sum p_i + \sum_{i=1}^n \frac{1}{n^2}$$

$$= \sum p_i^2 - \frac{1}{n}$$

collision prob of uniform distribution = $\|u_{[n]}\|_2^2$

we know this since we know n

collision prob of $p: \|p\|_2^2 = \Pr_{s,t \in P} [s=t] = \sum p_i^2$

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

for $p = u$:

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

for $p \neq u$:

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

Algorithm to estimate:

- take s samples of p
- let $\hat{c} \leftarrow$ estimate of $\|p\|_2^2$ from sample
- if $\hat{c} < \frac{1}{n} + \delta$ pass
else fail

- ① how big is s ?
- ② how to estimate?
- ③ what should δ be

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 $\hat{C} \leq \|U_{[n]}\|_2^2 + \frac{\epsilon^2}{2} \leq \frac{1}{n} + \frac{\epsilon^2}{2}$

so if we use $\delta = \frac{\epsilon^2}{2}$
 test should PASS

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

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

$\star \Rightarrow \hat{C} > \left(\epsilon^2 + \frac{1}{n}\right) - \frac{\epsilon^2}{2} = \frac{\epsilon^2}{2} + \frac{1}{n}$

so if we use $\delta = \frac{\epsilon^2}{2}$
 test should FAIL

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

Naive idea:

- repeat several times;
- take two samples & set $X_i \leftarrow \begin{cases} 1 & \text{if two samples equal} \\ 0 & \text{o.w.} \end{cases}$
- increment i
- output average of X_i 's

$O(k)$ samples of collisions from k samples of p

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

Naive idea: • repeat several times;

$\Theta(k)$ samples of collisions from k samples of p

• 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

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 $1 \leq i < j \leq s$

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

b_{ij} 's are not independent
 \Rightarrow can't use Chernoff

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

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

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

Chebyshev's \neq

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

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

by fact from before

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

so + lemma $\Rightarrow \text{Var}[\hat{c}]$ is $O\left(\frac{\|p\|_2^2}{s^2} + \frac{\|p\|_2^3}{s}\right)$

Proof of lemma in next lecture