

Lecture 14

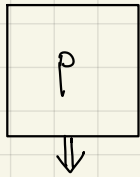
More on testing distributions:

- Poissonization
- Dealing with large l_2 -norm

(+ by the way, ... Testing Closeness)

Recall our setting:

Probability distributions: get samples (only)



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

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

this is
all
we
see

} iid samples

Distances

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

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

$$\|p\|_2 = \sqrt{\sum p_i^2}$$

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

Last time

Testing Uniformity

uniform dist
on domain D

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

with prob
 $\geq 3/4$

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

$l_1 + l_2$ distance measures

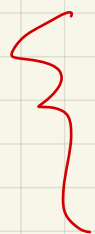
Generalizations:

Given another distribution q ,

is $p=q$ or is p "far" from q ?

today: focus on L_1 -distance

$q = \text{uniform}$
 $O(\sqrt{n})$
for all q
 $O(\sqrt{n})$



1. "Identity Testing"

q is known to algorithm, no samples of q needed

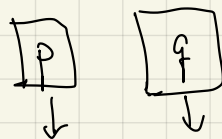


$\Theta(n^{2/3})$



2. "Closeness Testing"

q is given via samples



Tolerant version:

$$\|p - q\|_1 < \epsilon$$

$$\|p - q\|_1 \geq \epsilon' \text{ for } \epsilon' \gg \epsilon \quad \left. \vphantom{\|p - q\|_1 \geq \epsilon'} \right\} \Theta\left(\frac{n}{\log n}\right)$$

What is the sample complexity of these problems in terms of n

Recall: "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

How many samples?

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

A difficulty in analyzing distribution testers:

typical algorithm:

take m samples $\{S_1, \dots, S_m\} = S$
 let $X_i = \#$ times i occurred in sample
 \vdots
 \circ

Problem:
 X_i 's not indep
 e.g. if $X_1 = \frac{m}{2} + 1$
 then $X_2 < \frac{m}{2}$

Can we make the X_i 's independent? Poissonization

Poi(λ): $\Pr[X=k] = \frac{e^{-\lambda} \lambda^k}{k!}$
 $E[X] = \text{Var}[X] = \lambda$

new algorithm:

①

$\hat{m} \leftarrow \text{Poi}(m)$
 Take \hat{m} samples to get \hat{S}
 let $X_i = \#$ times i occurred in \hat{S}
 \vdots
 \circ

equivalent \longleftrightarrow

For each $i \in [n]$
 $X_i \leftarrow \text{Poi}(m \cdot p_i)$ ②
 add X_i copies of i to sample
 Randomly permute the sample
 \vdots
 \circ

Why equivalent?

$$\Pr[X_i = c \text{ according to (1)}] = \sum_{k=c}^{\infty} \Pr[\hat{m} = k] \cdot \binom{k}{c} p_i^c (1-p_i)^{k-c}$$

$$= \sum_{k=c}^{\infty} \frac{e^{-m} m^k}{k!} \cdot \frac{k!}{c!(k-c)!} p_i^c (1-p_i)^{k-c}$$

$$= \frac{e^{-m} m^c p_i^c}{c!} \sum_{k=c}^{\infty} \frac{m^{k-c} (1-p_i)^{k-c}}{(k-c)!} = \sum_{k'=0}^{\infty} \frac{m^{k'} (1-p_i)^{k'}}{(k')!} = e^{m(1-p_i)}$$

use $\lambda = m$

$$= \frac{e^{-m} m^c p_i^c e^{m(1-p_i)}}{c!} = \frac{e^{-mp_i} (mp_i)^c}{c!} = \Pr[X_i = c] = \Pr[X_i = c]$$

$X_i \sim \text{Poi}(mp_i)$

process II

$$X \sim \text{Poi}(\lambda)$$

$$\Pr[X=k] = \frac{e^{-\lambda} \lambda^k}{k!}$$

$$E[X] = \text{Var}[X] = \lambda$$

Need to check joint distributions samp

Another difficulty: $\|p\|_2$ can be large

eg. uniformity test statistic $\text{Var}[\hat{c}] = O\left(\frac{\|p\|_2^2}{s^2} + \frac{\|p\|_2^3}{s}\right)$

Goal: transform distributions p, q into p', q' st. $\|p'\|_2 + \|q'\|_2$ small

"reduction"
to
small
 L_2 -norm
case

$$\begin{aligned} p=q &\Rightarrow p'=q' \\ \|p-q\|_1 > \varepsilon &\Rightarrow \|p'-q'\|_1 > \varepsilon \end{aligned}$$

will work when q known & when given via samples

Transformation of p : $m = \# \text{ samples "expected" by original alg}$

Diakonikolas Kone
FOCS 2016
See also Sublinear Sp 2019
course notes by
Maryam Alakbarpour

$S \leftarrow \text{Draw } m \stackrel{!}{=} \text{Poi}(m) \text{ samples from } p \text{ over domain } [n]$

$b_i \leftarrow \# \text{ times } i \text{ appears in } S \quad \forall i \in [n]$

$\forall i$ add $b_i + 1$ elements to new domain
 (i, j) where $j \in [b_i + 1]$

\leftarrow equivalently:
 $b_i \sim \text{Poi}(p(i) \cdot m)$

New distribution p' :
pick $i \in_R P$
pick $j \in_R [b_i + 1]$
output (i, j)
 $p'(i, j) = \frac{p(i)}{b_i + 1}$

\swarrow size = $m + n$

Example:
domain of p is $[5]$
e.g. $S = \{2, 5, 3, 2, 3\}$
 $b_2 = b_3 = 2$
 $b_5 = 1$
all other b_i 's = 0

domain of p' :
 $\{ (1, 1)$
 $(2, 1) (2, 2) (2, 3)$
 $(3, 1) (3, 2) (3, 3)$
 $(4, 1)$
 $(5, 1) (5, 2) \}$

Prob

	$p(1)$		
$\frac{p(2)}{3}$	$\frac{p(2)}{3}$	$\frac{p(2)}{3}$	$\frac{p(2)}{3}$
$\frac{p(3)}{3}$	$\frac{p(3)}{3}$	$\frac{p(3)}{3}$	$\frac{p(3)}{3}$
$\frac{p(4)}{2}$	$\frac{p(5)}{2}$	$\frac{p(5)}{2}$	

$b_i \leftarrow \# \text{ times } i \text{ appears in } S \quad \forall i \in [n]$

ϕ' :
 pick $i \in_R P$
 pick $j \in_R [b_i+1]$
 output (i, j)

$$p'(i, j) = \frac{p(i)}{b_i+1}$$

Claim: $E[\|\phi'\|_2^2] \leq \frac{1}{m}$

Problem?
 we don't know if q 's l_2 norm gets small

Why? $E[\|\phi'\|_2^2] = E\left[\sum_{i=1}^n \sum_{j=1}^{b_i+1} p'(i, j)^2\right] = E\left[\sum_{i=1}^n \sum_{j=1}^{b_i+1} \frac{p(i)^2}{(b_i+1)^2}\right]$

$$= E\left[\sum_{i=1}^n \frac{p(i)^2}{b_i+1}\right] = \sum_{i=1}^n E\left[\frac{p(i)^2}{b_i+1}\right] = \sum_{i=1}^n E\left[\frac{1}{b_i+1}\right] \cdot p(i)^2$$

$$\leq \sum_i \frac{p(i)^2}{m \cdot p(i)} \leq \frac{1}{m} \cdot \sum p(i) \leq \frac{1}{m} \quad \square$$

$X \sim \text{Poi}(\lambda)$
 $\Pr[X=k] = \frac{e^{-\lambda} \lambda^k}{k!}$
 $E[X] = \text{Var}[X] = \lambda$

Claim for $Z \sim \text{poi}(\lambda) \quad E\left[\frac{1}{Z+1}\right] \leq \frac{1}{\lambda}$

$$E\left[\frac{1}{Z+1}\right] = \sum_{z=0}^{\infty} \frac{e^{-\lambda} \lambda^z}{(z+1)z!} = \frac{1}{\lambda} \sum_{z=0}^{\infty} \frac{e^{-\lambda} \lambda^{z+1}}{(z+1)!} = \frac{1}{\lambda} \sum_{z=1}^{\infty} \frac{e^{-\lambda} \lambda^z}{(z)!}$$

prob of disjoint events so sum to ≤ 1

After transform $p \rightarrow p'$ using same S :
 $f \rightarrow f'$

$$\begin{aligned} \|p - q\|_1 &= \sum_x |p(x) - q(x)| \\ &= \sum_x \sum_{y=1}^{b_{x+1}} \frac{|p(x) - q(x)|}{b_{x+1}} \\ &= \sum_x \sum_{y=1}^{b_{x+1}} |p'(x) - q'(x)| \\ &= \|p' - q'\|_1 \end{aligned}$$

Transformation of p :

$S \leftarrow$ Draw $\text{poly}(m)$ samples from p over domain $[n]$

$b_i \leftarrow$ # times i appears in $S \quad \forall i \in [n]$

$\forall i$ add $b_i + 1$ elements to new domain

(i, j) where $j \in [b_i + 1]$

New distribution p' : pick $i \in_R P$

$p'(i, j) = \frac{p(i)}{b_i + 1}$
pick $j \in_R [b_i + 1]$

output (i, j)

Claim: $E[\|p'\|_2^2] \leq \frac{1}{m}$

L_2 distance estimation between two distributions p & q :

easier when both $\|p\|_2^2 + \|q\|_2^2$ are small

Thm* given samples of p, q , distributions on $[n]$, s.t. $b \geq \max\{\|p\|_2, \|q\|_2\}$
can distinguish $p \neq q$ from $\|p - q\|_1 > \varepsilon$ in $O(bn/\varepsilon^2)$ samples

Corr if $b = \min\{\|p\|_2, \|q\|_2\}$

can distinguish $p \neq q$ from $\|p - q\|_1 > \varepsilon$ in $O(bn/\varepsilon^2)$ samples

PF idea:

1. estimate $\|p\|_2 + \|q\|_2$ to mult factor of c (can do this in $O(\sqrt{n})$ samples)
 2. if differ by $> c$ mult factor infer $p \neq q$ & reject
 3. else use Thm* with $b' = c \cdot b$
-

Testing Closeness

Corr if $b \geq \min \{ \|p\|_2^2, \|q\|_2^2 \}$

can distinguish $p \neq q$ from $\|p - q\|_1 \geq \epsilon$ in $O(bn/\epsilon^2)$ samples

1. let $k = n^{2/3} \epsilon^{-4/3}$

2. $S \leftarrow$ multiset of $\text{Poi}(k)$ samples from q

3. run tester of Corr on p, q' wrt. S



Samples of $p+q$

Why does it work?

distinguishing $p+q \neq p', q'$ are equivalent

how many samples needed?

• whp $|S|$ is $\Theta(k)$

• $E[\|q'\|_2^2] = O(1/k)$ so whp $\|q'\|_2 = O(1/\sqrt{k})$

$$\begin{aligned} & O\left(k + \frac{1}{\sqrt{k}} \cdot n \cdot \frac{1}{\epsilon^2}\right) = O\left(n^{2/3} \cdot \epsilon^{-4/3} + \frac{1}{n^{1/3}} \cdot \epsilon^{-2/3} \cdot n \cdot \frac{1}{\epsilon^2}\right) \\ & \quad \uparrow \text{picking } S \quad \uparrow \text{run tester on } p', q' \\ & = O\left(n^{2/3} \epsilon^{-4/3}\right) \end{aligned}$$