# Lecture 12:

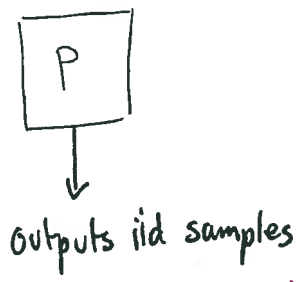## Testing Distributions

- Uniformity

Turning to a new model:

Probability distributions – get samples of distribution

Domain $D$, $|D| = n$ ← known

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

[P]

↓

outputs iid samples

↖ this is all we can learn from

Examples:

Lottery data

Shopping choices

experimental outcomes

⋮

What do we want to know?

is it uniform? eg. lottery

is it high entropy?

large support? (many distinct elements have $>0$ probability

is p monotone increasing, k-modal, monotone hazard rate...?

how can we do it?

$\chi^2$ test

plug in estimate

learn distribution, Maximum likelihood estimates

Goal: sample complexity  SUBLINEAR  in  $n$

# Testing Uniformity

The goal:

↙ Uniform dist on D

- if $P \equiv U_D$ then tester outputs PASS ← with prob ≥ 3/4

- if $\mathrm{dist}(P, U_D) > \varepsilon$ then tester outputs FAIL

which measure of distance?

$\ell_1, \ell_2,$ KL-divergence, Earthmover, Jensen-Shannon

↑ today's focus

good direction for projects!

# Distances

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

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

$$\|p-q\|_2 \leq \|p-q\|_1 \leq n^{1/2} \|p-q\|_2$$

## examples:

① $\quad p = (1,0,0,\dots 0)$



$\quad q = (\tfrac{1}{n}, \tfrac{1}{n}, \dots \tfrac{1}{n})$

$l_1$ distance:
$$\|p-q\|_1 = \left(\tfrac{n-1}{n}\right) + (n-1)\cdot\tfrac{1}{n}$$
$$\approx 2$$
$l_2$-distance:
$$\|p-q\|_2^2 = (1-\tfrac{1}{n})^2 + (n-1)(\tfrac{1}{n})^2$$
$$\approx 1$$

② 

$\quad p = \left(\tfrac{2}{n}, \tfrac{2}{n}, \dots \tfrac{2}{n}, 0,0,\dots 0\right)$

$\quad q = \left(0,0,\dots 0, \tfrac{2}{n}, \tfrac{2}{n}, \dots \tfrac{2}{n}\right)$



as far as
possible ↓

$l_1$ distance:
$$\|p-q\|_1 = n\cdot\left(\tfrac{2}{n}\right) = 2$$

$l_2$-distance: $\|p-q\|_2^2 = n\cdot\left(\tfrac{2}{n}\right)^2 = \tfrac{4}{n}$

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

↑
pretty
close

$\Rightarrow$ so $l_2$-distance can be weird

## "Plug-in" Estimate:

Algorithm:

- 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.

Analysis: (better analyses exist — see next page)

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

by $\triangle \neq$, if $\|p - \hat{p}\|_1 < \varepsilon$ + $\|\hat{p} - U\|_1 < \varepsilon$ ← this happens exactly when test passes

then $\|p - U\|_1 < 2\varepsilon$.

*So, if $p = U_n$ then $p$ passes* (circled note)

*So, if $\|p - U_n\| > 2\varepsilon$ this test is likely to fail* (circled note)

how many samples? $\Omega\left(\frac{n}{\varepsilon}\right)$ maybe even worse ...

*for each $x$, need to see it at least once in order to give non zero estimate.* (circled note)

$\Theta(n)$? Can we do better?

**Uh oh, do we need "coupon collector bound" $\Omega(n \log n)$?**

**Better analysis:**

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

Pf

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

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

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

← Jensen's ≠

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

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

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

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

← since $\max_{p \in \text{prob dist over domain of size } n} \sum \sqrt{p(x)}$ is $\sqrt{n}$

So picking $m = \Omega\left(\frac{n}{\varepsilon^2}\right)$ gives

$$E[\|\hat{p} - p\|_1] \leq \frac{\varepsilon}{2}$$

by Markov's ≠ : with prob $1 - \frac{1}{2}$, $\|\hat{p} - p\|_1 \leq \varepsilon$

Note, this says can "learn" (approximate) any dist w.r.t. $L_1$ distance in $\theta(n/\varepsilon^2)$ samples

# $L_2$ - Distance (squared):

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

$$= \sum p_i^2 - \tfrac{2}{n} \underbrace{\sum p_i}_{=1} + \underbrace{\sum (\tfrac{1}{n})^2}_{= \tfrac{1}{n}}$$

$$= \underbrace{\sum p_i^2}_{} - \tfrac{1}{n}$$

Collision probability of $p$:

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

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

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

$$= \underbrace{\|p\|_2^2}_{\substack{\text{we can} \\ \text{estimate} \\ \text{this}}} - \underbrace{\|u_{(n)}\|_2^2}_{\substack{\text{we know this} \\ \text{since we know } n}}$$

# Algorithm

1. take $s$ samples from $p$

2. let $\hat{c} \leftarrow$ estimate of $\|p\|_2^2$ from sample

3. if $\hat{c} < \tfrac{1}{n} + \delta$ pass

   else fail

① how many samples?

② how?

③ what should $\delta$ be?

First:

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

Naive idea:

take two new samples:

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

<span style="color:red">} $\delta_i$'s are independent</span>

" 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 $p$ "

<span style="color:red">} $\delta_{ij}$'s are not independent</span>

$$\delta_{ij} \leftarrow \begin{cases} 1 & \text{if sample i, j are equal} \\ 0 & \text{o.w.} \end{cases}$$

Estimate by recycling:

· Take $s$ samples from $p$: $X_1 \cdots X_s$

· for each $1 \leq i < j \leq s$

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

$\delta_{ij}$'s not independent so can't use Chernoff

· Output $\hat{c} \leftarrow \dfrac{\sum\limits_{i<j} \delta_{ij}}{\binom{s}{2}}$

Analysis: $E[\hat{c}] = \dfrac{1}{\binom{s}{2}} \cdot \binom{s}{2} \cdot E[\delta_{ij}]$

$= \|p\|_2^2$

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

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

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 $\Delta = \frac{\varepsilon^2}{2}$

What happens if ✸ holds with $\Delta = \frac{\varepsilon^2}{2}$ ?

Correct behavior!

- if $p = U_{[n]}$ then $\hat{c} \leq \|U_{[n]}\|_2^2 + \Delta = \frac{1}{n} + \frac{\varepsilon^2}{2}$

  so test will PASS

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

  but $\|p\|_2^2 = \|p - U_{[n]}\|_2^2 + \frac{1}{n}$ ← see p.6

  $\qquad\qquad > \varepsilon^2 + \frac{1}{n}$

  ↓ $\hat{c} > \|p\|_2^2 - \Delta$ ← ✸

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

  so test will FAIL

Remaining Question:

How many samples do we need to estimate $\hat{c}$ to within $\Delta$ ?

## Analysis

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

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

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

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

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

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

**Lemma** $\text{Var}\left[\sum \sigma_{ij}\right] \leq 4\left(\binom{s}{2} \|p\|_2^2\right)^{3/2}$

Fact $\Rightarrow$
$$\text{Var}[\hat{c}] \leq 4 \cdot \frac{\left(\binom{s}{2}\|p\|_2^2\right)^{3/2}}{\binom{s}{2}^2} \leq \theta\left(\|p\|_2^3/s\right)$$

**Why?** (proof...)

$\xleftarrow{}$ trick – will rewrite variance as $E[\sigma_{ij}]$.

$\underset{=}{\text{def}}. \quad \overline{\sigma_{ij}} = \sigma_{ij} - E[\sigma_{ij}]$

So $E[\overline{\sigma_{ij}}] = 0 \quad + \quad \overline{\sigma_{ij}} < \sigma_{ij} \quad (\text{since } E[\sigma_{ij}] > 0)$

Also $\therefore \circ E[\overline{\sigma_{ij}} \, \overline{\sigma_{kl}}] \leq E[\sigma_{ij} \sigma_{kl}]$

verify at home? (or trust...)
- $(\sum p(x)^3)^{1/3} \leq (\sum p(x)^2)^{1/2}$
- $s^2 \leq 3\binom{s}{2}$
- $\binom{s}{3} \leq s^3/6$

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

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

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

$$= E\left[\underbrace{\sum_{i<j}\bar{\sigma}_{ij}^2}_{①} + \underbrace{\sum_{\substack{i<j\\k<\ell\\i,j,k,\ell \text{ distinct}}}\bar{\sigma}_{ij}\bar{\sigma}_{k\ell}}_{②} + \underbrace{\sum_{\substack{i<j\\i<\ell\\i,j,\ell\\ \text{distinct}}}\bar{\sigma}_{ij}\bar{\sigma}_{i\ell}}_{③} + \underbrace{\sum_{\substack{i<j\\k<j\\i,k,j\\ \text{distinct}}}\bar{\sigma}_{ij}\bar{\sigma}_{kj}}_{④}\right]$$

$$⑤ \quad + \sum\bar{\sigma}_{ij}\bar{\sigma}_{j\ell}$$
$$⑥ \quad + \sum\bar{\sigma}_{ij}\bar{\sigma}_{ki}$$

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

$\boxed{E[\sigma_{ij}] = E[\sigma_{ij}^2] \text{ since } \sigma_{ij} \text{ is indicator var}}$

② $E\left[\sum_{\substack{i<j\\k<\ell\\ \text{all 4 distinct}}}\bar{\sigma}_{ij}\bar{\sigma}_{k\ell}\right] \overset{\text{independent}}{\leq} \sum E[\bar{\sigma}_{ij}]E[\bar{\sigma}_{k\ell}] = 0$

③ $E\left[\sum\bar{\sigma}_{ij}\bar{\sigma}_{i\ell}\right] \leq E\left[\sum_{\substack{i,j,\ell\\ \text{distinct}}}\sigma_{ij}\cdot\sigma_{i\ell}\right] = \sum_{\substack{i,j,\ell\\ \text{distinct}}}\Pr[X_i=X_j=X_\ell]$

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

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

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

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

④   same   as  3

⑤

⑥

In total:

$$\mathrm{Var}\left[\sum_{i<j} b_{ij}\right] \leq \mathrm{Var}\left[\sum_{i<j} \bar{b}_{ij}\right]$$

$$\leq \binom{5}{2}\|p\|_2^2 + 0 + 4\cdot\frac{\sqrt{3}}{2}\left(\binom{5}{2}\|p\|_2^2\right)^{3/2}$$

$$\leq 4\left[\binom{5}{2}\|p\|_2^2\right]^{3/2}$$

∎

Putting lemma into Chebyshev :

use $\rho = \frac{\varepsilon^2}{2}$

$$\Pr\left[\,|\hat{c} - \|p\|_2^2| \,>\, \frac{\varepsilon^2}{2}\,\right] \;\leq\; \frac{\mathrm{Var}[\hat{c}]}{\varepsilon^4} \cdot 4$$

recall this crude from proof? const in proof?

$$\leq\; \frac{4\left[\binom{5}{2}\|p\|_2^2\right]^{3/2}}{\binom{5}{2}^2 \,\varepsilon^4}\cdot 4 \;\leq\; \frac{32}{\varepsilon^4}\cdot\frac{1}{5}\cdot\|p\|_2^{3}$$

note $\frac{1}{\binom{5}{2}^{1/2}} \leq \frac{1}{\sqrt{\frac{5^2}{2}}} \leq \frac{2}{5}$

also want this to be $\leq 1$   $\leq 1$

So Pick $s \geq \Omega\left(\frac{1}{\varepsilon^4}\right)$

<u>Note</u>: Can get better bnd

1) Testing closeness to any known
   distribution — reduce to uniform case!

2) lower bound

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

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

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

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

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

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

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

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

$s \geq \dfrac{const \cdot \sqrt{n}}{\varepsilon^2}$ samples

since $\Pr\left[|\hat{c} - \|p\|_2^2| \geq \gamma \|p\|_2^2\right] \leq \dfrac{k \cdot \|p\|_2^3}{s \cdot \gamma^2 (\|p\|_2^2)^2} \leq \dfrac{k}{s \cdot \gamma^2 \cdot \|p\|_2}$  ↙ const

$\left[\text{note} \quad \|p\|_2^2 > \frac{1}{n} \text{ so } \|p\|_2 > \frac{1}{\sqrt{n}} \text{ so } \frac{1}{\|p\|_2} < \sqrt{n}\right]$

$\leq \dfrac{k \cdot \sqrt{n}}{s \cdot \gamma^2}$   $\left[\text{note: we need } \gamma \approx \frac{\varepsilon^2}{3}\right]$

so picking $s >> \frac{\sqrt{n}}{\varepsilon^4}$ will give small probability error $\implies$ $\approx \dfrac{k \cdot \sqrt{n}}{s} \cdot \frac{1}{\varepsilon^4}$