## Lecture 13

Testing distributions:
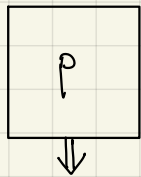
the case of uniformity   (cont)

# A new model:

## Probability distributions: get samples

Discrete Domain $D$ s.t. $|D| = n$ ← Known $n$

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

$\boxed{p}$

↓

this is all we see → { iid samples

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?
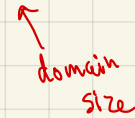
k-histogram?

Methods ?

learn distribution

$\chi^2$- test

plug-in estimate

Max likelihood estimate

Goal : Sample complexity sublinear in n

domain size

# Testing Uniformity

uniform dist on domain $D$

goal: if $p \equiv U_D$ then output PASS — with prob $\geq 3/4$

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

which measure of distance?

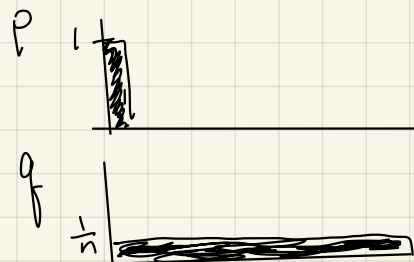$l_1, l_2,$ KL-divergence, Earthmover, Jensen-Shannon ....

today's focus

# Distances

$\ell_1$ - distance :
$$\|p-q\|_1 = \sum_{i \in D} |p_i - q_i|$$

$\ell_2$ - distance :
$$\|p-q\|_2 = \sqrt{\sum_{i \in D} (p_i - q_i)^2}$$

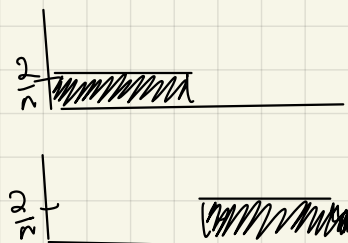$$\|p-q\|_2 \leq \|p-q\|_1 \leq \sqrt{n} \cdot \|p-q\|_2$$

## examples :

① $\quad p = (1, 0, 0, 0, \ldots, 0)$

$\quad q = \left(\frac{1}{n}, \frac{1}{n}, \frac{1}{n}, \ldots, \frac{1}{n}\right)$



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

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

② $\quad p = \left(\frac{2}{n}, \frac{2}{n}, \frac{2}{n}, \ldots \frac{2}{n}, 0, 0, \ldots 0\right)$

$\quad q = \left(0, 0, \ldots 0, \frac{2}{n}, \frac{2}{n}, \ldots \frac{2}{n}\right)$



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

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

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

tiny
even though
$\ell_1$ is big

# Via "Plug-in" Estimate:

- take $m$ samples from $p$

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

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

  else accept

# How many samples?

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

## Let's consider $L_2$-distance (squared):

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

uniform on $1 \ldots n$

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

for $p = U$:
$$\|p\|_2^2 = \tfrac{1}{n}$$

for $p \neq U$:
$$\|p\|_2^2 > \tfrac{1}{n}$$

$$= \underbrace{\sum p_i^2}_{\text{collision prob of}} - \tfrac{1}{n}$$

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

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

we know this since we know $n$

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

## Algorithm to estimate:

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

  else fail

① how big is $s$?

② how to estimate?

③ what should $\delta$ be

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

ie. what should $\delta$ be?

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

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 ✱ holds with $\Delta = \dfrac{\varepsilon^2}{2}$ ?

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

so if use $\delta = \dfrac{\varepsilon^2}{2}$ test should PASS

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

but $\|p\|_2^2 = \|p - U_{[n]}\|_2^2 + \dfrac{1}{n} > \varepsilon^2 + \dfrac{1}{n}$

& ✱ $\implies \hat{c} > \left(\varepsilon^2 + \dfrac{1}{n}\right) - \dfrac{\varepsilon^2}{2} = \dfrac{\varepsilon^2}{2} + \dfrac{1}{n}$

so if we use $\delta = \dfrac{\varepsilon^2}{2}$ test should FAIL

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

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

**Better idea:** "recycle" use <u>all</u> pairs in sample

gives $\Theta(K^2)$ samples of collision prob from $k$ samples of $p$

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

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

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

  $\left. \phantom{\begin{cases} 1 \\ 0 \end{cases}} \right\}$ $b_{ij}$'s are not independent

  $\Rightarrow$ can't use Chernoff

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

Analysis:

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

$$Pr\left[\left|\hat{c} - \|p\|_2^2\right| > \rho\right] \le \frac{Var[\hat{c}]}{\rho^2} \qquad \text{Chebyshev's} \ne$$

recall $Var[x] = E[(x-E(x))^2]$

$$Var[\hat{c}] = \frac{1}{\binom{s}{2}^2} Var\left[\sum_{i<j} b_{ij}\right] \qquad \text{by fact: } Var[aX] = a^2 Var[x]$$

need to bound

difficulty: $b_{ij}$'s not independent

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

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

**Lemma** $\quad Var\left[\sum\limits_{i<j} \sigma_{ij}\right] \leq \binom{s}{2}\|p\|_2^2 + 4\left(\binom{s}{2}\|p\|_2^2\right)^{3/2}$

**Proof**

$\underline{def}\quad \bar{\sigma}_{ij} = \sigma_{ij} - E[\sigma_{ij}]$

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

$\leftarrow$ trick: rewrite variance as $E[\sum \bar{\sigma}_{ij}]$

why? $\quad Var[\sum \bar{\sigma}_{ij}] = E\left[\left(\sum \bar{\sigma}_{ij} - \underbrace{E[\sum \bar{\sigma}_{ij}]}_{=0}\right)^2\right]$

$= E\left[\left(\sum \bar{\sigma}_{ij}\right)^2\right]$

$= E\left[\left(\sum \sigma_{ij} - E\sigma_{ij}\right)^2\right]$

$= Var\left(\sum \sigma_{ij}\right)$

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

**Proof**

$\underline{\text{def}}\quad \bar{b}_{ij} = b_{ij} - E[b_{ij}]$

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

$\Leftarrow$ trick: rewrite variance as $\cancel{E\left[\sum \bar{b}_{ij}^2\right]}^{=0}$

why? $\quad \mathrm{Var}\left[\sum \bar{b}_{ij}\right] = E\left[\left(\sum \bar{b}_{ij} - E\left[\sum \bar{b}_{ij}\right]\right)^2\right]$

$\qquad\qquad\qquad = E\left[\left(\sum b_{ij} - E[b_{ij}]\right)^2\right]$

$\qquad\qquad\qquad = \mathrm{Var}\left[\sum b_{ij}\right]$

So can equivalently bound
$\mathrm{Var}\left[\sum \bar{b}_{ij}\right]$

**Facts:**

- $E\left[\bar{b}_{ij}\,\bar{b}_{k\ell}\right] \leq E\left[b_{ij}\,b_{k\ell}\right]$

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

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

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

(Verify @ home)

## 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

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

$$= E\left[\sum_{\substack{i<j \\ \text{①}}} \bar{\delta}_{ij}^2 + \sum_{\substack{i<j \\ k<l \\ i,j,k,l \\ \text{are all distinct} \\ \text{②}}} \bar{\delta}_{ij}\bar{\delta}_{kl} + \sum_{\substack{i<j \\ i<l \\ i,j,l \\ \text{distinct} \\ \text{③}}} \bar{\delta}_{ij}\bar{\delta}_{il} + \sum_{\substack{i<j \\ k<j \\ i,j,k \\ \text{distinct} \\ \text{④}}} \bar{\delta}_{ij}\bar{\delta}_{kj}\right.$$

$$\left. + \sum_{\substack{i<j<l \\ \text{⑤}}} \bar{\delta}_{ij}\bar{\delta}_{jk} \right]$$

Lets bound each term:

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

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

$\underbrace{\phantom{\Pr[\delta_{ij}=1]}}_{\text{prob of collision}}$

(right margin)

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

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

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

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 $E\left[\displaystyle\sum_{\substack{i<j\\k<l\\ \text{all distinct}}} \bar{\sigma}_{ij}\cdot\bar{\sigma}_{kl}\right] \le \displaystyle\sum_{\substack{i<j\\k<l\\ \text{all distinct}}} E\left[\bar{\sigma}_{ij}\cdot\bar{\sigma}_{kl}\right] = \displaystyle\sum_{\substack{i<j\\k<l\\ \text{all distinct}}} E\left[\bar{\sigma}_{ij}\right]\cdot E\left[\bar{\sigma}_{kl}\right]$

$\overset{\text{independence}}{\checkmark}$

$\underset{0}{\underbrace{\qquad}}$

$= 0$

$\boxed{\text{and}}$ ③ $\boxed{④}$ $\boxed{⑤}$

③ $E\left[\displaystyle\sum_{\substack{i<j\\i<l\\ i,j,l\\ \text{distinct}}} \bar{\sigma}_{ij}\bar{\sigma}_{il}\right] \le E\left[\displaystyle\sum \sigma_{ij}\sigma_{il}\right] = \displaystyle\sum E\left[\sigma_{ij}\sigma_{il}\right]$

$\underset{\substack{1 \text{ iff saw same element in}\\ i\text{th}, j\text{th} + l\text{th sample}\\ \text{"3-way collision"}}}{\underbrace{\qquad\qquad}}$

$= \displaystyle\sum_{\substack{i,j,l\\ \text{distinct}}} \Pr[X_i = X_j = X_l]$

$= \binom{s}{3}\displaystyle\sum_x p(x)^3$

$\le \dfrac{s^3}{6}\cdot\left(\displaystyle\sum_x p(x)^2\right)^{3/2}$

$\le \dfrac{\sqrt{3}}{2}\left(\dfrac{s}{2}\right)^{3/2}\left(\|p\|_2^2\right)^{3/2}$ ⟵ by facts

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

$\underline{\text{def}}$ $\bar{\sigma}_{ij} = \sigma_{ij} - E[\sigma_{ij}]$

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

Facts:

• $E[\bar{\sigma}_{ij}\bar{\sigma}_{kl}] \le E[\sigma_{ij}\sigma_{kl}]$

• $\left(\displaystyle\sum_x p(x)^3\right)^{1/3} \le \left(\displaystyle\sum_x p(x)^2\right)^{1/2}$

• $s^2 \le 3\binom{s}{2}$

• $\binom{s}{3} \le s^3/6$

$$S_0 \quad Var\left( \sum_{i<j} b_{ij} \right) = Var\left[ \sum_{i<j} \overline{b_{ij}} \right]$$

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

$$\leq \binom{S}{2} \|p\|_2^2 + 4 \left( \binom{S}{2} \|p\|_2^2 \right)^{3/2}$$

We have:

$$\mathrm{Var}\left[\hat{C}\right] = O\left(\frac{\|p\|_2^2}{S^2} + \frac{\|p\|_2^3}{S}\right)$$

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

$$\hat{C} \leftarrow \frac{\sum_{i<j} b_{ij}}{\binom{S}{2}}$$

where $S = \#$ Samples

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

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

$$\leq \frac{const}{\varepsilon^4 \cdot S^2} \cdot \underbrace{\|p\|_2^2}_{\leq 1} \qquad + \qquad const \cdot \frac{1}{\varepsilon^4} \cdot \underbrace{\frac{1}{S}}_{\substack{\text{want to} \\ \text{be} \ll 1}} \cdot \underbrace{\|p\|_2^3}_{\leq 1}$$

$$\underbrace{\phantom{\frac{const}{\varepsilon^4 \cdot S^2}}}_{\substack{\text{want this} \\ \text{to be} \leq 1}}$$

need $S = \Omega\left(\frac{1}{\varepsilon^2}\right)$

need $S = \Omega\left(\frac{1}{\varepsilon^4}\right)$

# Samples $S$ to be $O\left(\frac{1}{\varepsilon^4}\right)$

note can get better bounds

$$S = O\left(\frac{1}{\varepsilon^2}\right)$$

$S$ is independent of $n$ !!!!!

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

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

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

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

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

So either additive estimate of $\|p\|_2^2$ to within $\frac{\varepsilon^2}{2n}$

or mult " " " to within

$(1 \pm \frac{\varepsilon^2}{3})$

suffices

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

$S = \sqrt{n}$ suffices

($\dagger$ actually $s = \frac{\sqrt{n}}{\varepsilon^2}$ sufficient)

Generalizations: Given another distribution $q$,

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

1. "Identity Testing"

$q$ is known to algorithm, no samples of $q$ needed $\}$ focus on sample complexity but runtime can be made similar

ik "DNA"

2. "Closeness Testing"

$q$ is given via samples

$\boxed{p}$  $\boxed{q}$

samples

Will see more on these ...

(e.g. pset, lecture ...)

What is complexity in terms of $n, \varepsilon$?

# A difficulty in analyzing distribution testers:

typical algorithm:

take $m$ samples $\{S_1, \ldots S_m\} = S$

let $X_i = $ # times $i$ occurred in sample

$\vdots$

Can we make the $X_i$'s independent?

Poissonization

$$Poi(\lambda): \Pr[x=k] = \frac{e^{-\lambda} \lambda^k}{k!}$$

$$E[x] = Var[x] = \lambda$$

## new algorithm 1

$\hat{m} \leftarrow Poi(m)$

Take $\hat{m}$ samples to get $\hat{S}$

let $X_i = $ # times $i$ occured in $\hat{S}$

$\vdots$

## new algorithm 2

For each $i \in [n]$

$X_i \leftarrow Poi(m \cdot p_i)$

add $X_i$ copies of $i$ to sample

Randomly permute the sample

$\vdots$