# Lecture 13
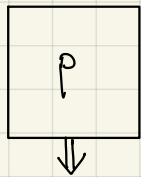
Testing distributions:

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

$\vdots$

What do we need to Know?    is it    uniform?

high entropy?

large support?    (many distinct elts with $\geq 0$ probability)

monotone increasing, K-modal?
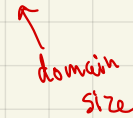
K-histogram?

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 $\quad p \equiv U_D \quad$ then   output   PASS $\quad\leftarrow$ with prob $\geq 3/4$

if $\quad$ dist $(p, U_D) > \varepsilon \quad$ then   output   FAIL

$\underbrace{\qquad\qquad\qquad}$
which measure
of distance?

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

$\uparrow\nearrow$
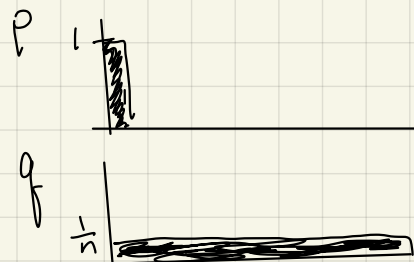today's focus

# Distances

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

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

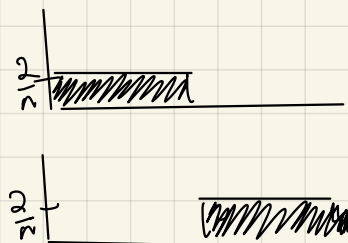$$\|p - q\|_2 \leq \|p - q\|_1 \leq \sqrt{n} \cdot \|p - q\|_2$$

## 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
$l_1$ 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 $\Theta(\frac{n}{\varepsilon^2})$ samples

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

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

*uniform on 1..n*

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

for $p = U$:

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

for $p \neq U$:

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

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

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 $\hat{c} < \frac{1}{n} + \delta$   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 $\maltese$ : $\left| \hat{c} - \|p\|_2^2 \right| < \Delta$

$\underbrace{}$
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 $\maltese$ holds with $\Delta = \frac{\varepsilon^2}{2}$?

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

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

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

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

$\delta \maltese \Rightarrow \hat{c} > \left( \varepsilon^2 + \frac{1}{n} \right) - \frac{\varepsilon^2}{2} = \frac{\varepsilon^2}{n} + \frac{1}{n}$

so if we use $\delta = \frac{\varepsilon^2}{2}$
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 $1 \le i < j \le s$
$$6_{ij} \leftarrow \begin{cases} 1 & \text{if } X_i = X_j \\ 0 & \text{o.w.} \end{cases}$$
- Output $\hat{c} \leftarrow \dfrac{\sum_{i<j} 6_{ij}}{\binom{s}{2}}$

$6_{ij}$'s are not independent $\Rightarrow$ can't use Chernoff

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

$$\Pr\left[\left|\hat{c} - \|p\|_2^2\right| > \rho\right] \le \frac{Var[\hat{c}]}{\rho^2}$$

Chebyshev's ≠

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

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

by fact: $Var[aX] = a^2 Var[x]$

need to bound
difficulty: $b_{ij}$'s not independent

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

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

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

**Proof**

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

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

Facts:

- $E[\bar{\sigma}_{ij} \bar{\sigma}_{k\ell}] \leq E[\sigma_{ij} \sigma_{k\ell}]$

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

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

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

(Verify @ home)

$\Longleftarrow$ trick: rewrite variance as $\cancel{E[\sum \bar{\sigma}_{ij}^2]}^{=0}$

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

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

$\qquad = \text{Var}[\sum \sigma_{ij}]$

So can equivalently bound

$\qquad \text{Var}[\sum \bar{\sigma}_{ij}]$

**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**

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

$$= E\left[\sum_{\substack{i<j}}\bar{\delta}_{ij}^2 + \sum_{\substack{i<j \\ k<\ell \\ i,j,k,\ell \text{ distinct}}}\bar{\delta}_{ij}\bar{\delta}_{k\ell} + \sum_{\substack{i<j \\ i<\ell \\ i,j,\ell \text{ distinct}}}\bar{\delta}_{ij}\bar{\delta}_{i\ell} + \sum_{\substack{i<j \\ k<j \\ i,k,j \text{ distinct}}}\bar{\delta}_{ij}\bar{\delta}_{kj}\right.$$

(1)  (2)  (3)  (4)

$$\left. + \sum_{i<j<\ell}\bar{\delta}_{ij}\bar{\delta}_{j\ell}\right]$$

(5)

Lets bound each term:

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

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

---

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

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

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

**Facts:**

• $E[\bar{\delta}_{ij}\bar{\delta}_{k\ell}] \leq E[\delta_{ij}\delta_{k\ell}]$

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

• $s^2 \leq 3\binom{s}{2}$

• $\binom{s}{3} \leq s^3/6$

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

<span style="color:red">$\underbrace{\quad\quad\quad}$ this is where the trick helps — gets rid of lots of terms</span>

③ (+ similarly ④ + ⑤)

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

$\leq \binom{5}{3} \sum\limits_x p(x)^3$  <span style="color:green">expected # 3-way collisions</span>

$\leq \dfrac{5^3}{6} \left(\sum\limits_x p(x)^2\right)^{3/2}$  <span style="color:green">$\Big\}$ by facts.</span>

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

---

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

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

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

Facts:

- $E[\bar{\delta}_{ij} \, \bar{\delta}_{kl}] \leq E[\delta_{ij} \, \delta_{kl}]$
- $\left(\sum\limits_x p(x)^3\right)^{1/3} \leq \left(\sum\limits_x p(x)^2\right)^{1/2}$
- $5^2 \leq 3\binom{5}{2}$
- $\binom{5}{3} \leq 5^3/6$

So,

$$\text{Var}\left[\sum_{i<j} \delta_{ij}\right] = \text{Var}\left[\sum_{i<j} \delta_{ij}\right]$$

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

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

We have:

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

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

$$\hat{C} \leftarrow \frac{\sum_{i \neq j} b_{ij}}{\binom{s}{2}}$$

where $s = \#$ samples

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

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

$$\leq \frac{\text{const} \cdot \|p\|_2^2}{\varepsilon^4 \cdot s^2} \quad + \quad \text{Const} \cdot \frac{1}{\varepsilon^4} \cdot \frac{1}{s} \cdot \|p\|_2^3$$

want this $\leq 1$   $\leq 1$

also want this $\ll 1$   $\leq 1$

need $s = \Omega(1/\varepsilon^2)$

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

BIGGER CONSTRAINT

Note    can get    better   bounds

$s$ is independent of $n$!

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

recall:
$$\|p-u_{(n)}\|_2^2 = \|p\|_2^2 - \|u_{(n)}\|_2^2$$

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

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

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

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

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

or mult estimate of $\|p\|_2^2$ to within $\left(1 \pm \frac{\varepsilon^2}{3}\right)$

suffices

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

$$s = O(\sqrt{n})$$

# Generalizations:   Given another distribution $q$,

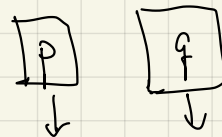is $p = q$ or is $p$ "far" from $q$?

↖ focus on $L_1$ distance

1. "Identity Testing"

$q$ is known to algorithm, no samples of $q$ needed

} focus on sample complexity but runtime can be made similar

2. "Closeness Testing"

$q$ is given via samples



Will see more on these soon
(e.g. Pset, lecture ...)

What is complexity in terms of $n$ ??

# A difficulty in analyzing distribution testers:

typical algorithm:

take $m$ samples $\{s_1 \dots s_m\} = S$

let $X_i = \#$ times $i$ occured in sample

$\vdots$

**problem:**

$X_i$'s are $\underline{not}$ independent

e.g. if $X_1 = \frac{m}{2} + 1$

then $X_2 < \frac{m}{2}$

Can we make the $X_i$'s independent? **Poissonization**

e.g. $S = \{2, 5, 3, 2, 3\}$

$X_2 = X_3 = 2$

$X_5 = 1$

all other $X_i = 0$

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

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

## new algorithm:

① $\hat{m} \leftarrow Poi(m)$

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

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

$\vdots$

equivalent $\iff$

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

$\vdots$

why equivalent?

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

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

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

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

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

use $\lambda = m$

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

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

$$= \Pr[X_i = c \text{ according to (2)}]$$

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

e.g. uniformity test statistic $\qquad Var\left[\hat{C}\right] = O\left(\frac{\|p\|_2^2}{S^2} + \frac{\|p\|_2^3}{3}\right)$

Goal: transform distributions $p, q$ into $p', q'$ ←on new domains st. $\|p'\|_2$ & $\|q''\|_2$ small

give reduction to small $L_2$-norm
$$\begin{cases} & \& \quad p=q \implies p'=q' \\ & \quad \|p-q\|_1 > \varepsilon \implies \|p'-q'\| > \varepsilon \end{cases}$$

Comment: $q$ may be "known" or given via samples

## Transformation of $p$:

$S \leftarrow$ Draw $Poi(k)$ samples from $p$ over domain $[n]$

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

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

$$(i, j) \quad \text{where} \quad j \in [b_i + 1]$$

New distribution $p'$: 

pick $i \in_R p$

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

output $(i, j)$

#samples
$\downarrow$
size $m + n$

Example:    $n = 5$

domain of $p$ is $[5]$

e.g. $S = \{2, 5, 3, 3, 3\}$

$X_2 = X_3 = 2$

$X_5 = 1$

all other $X_i = 0$

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

$p'$ is "mixture" of uniform + observed : $p' \approx \alpha \cdot U + (1 - \alpha) \cdot \hat{p}(i)$
distribution

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

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

$\left. \phantom{\begin{matrix}a\\a\\a\end{matrix}} \right\}$ $\rho'(i,j) = \dfrac{p(i)}{b_i + 1}$

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

Why?

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

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

$$\underset{\star}{\leq} \sum_i \frac{p(i)^2}{k \cdot p(i)} = \frac{1}{k}$$