# Testing Distributions:
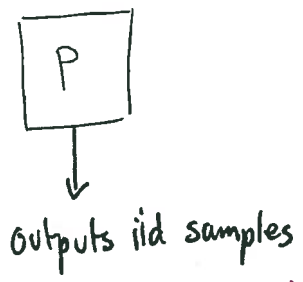
## Uniformity

Turning to a new model : <span style="float:right">prob dists</span>

Probability distributions — get samples of distribution



Domain $D$, $|D| = n$ ← Known

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

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 $\geq 3/4$

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

  $\underbrace{\qquad}$

  which measure of distance?

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

  ↑ today's focus

  good direction for projects!

# Distances

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

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

$\|p-q\|_1 = 2 \cdot \text{TVD dist}(p,q)$

$= \max_{S \subseteq D} \left\{ \sum_{i \in S} |p_i - q_i| \right\}$

"Total variation distance"

**Fact:** $\|p-q\|_2 \le \|p-q\|_1 \le n^{1/2} \|p-q\|_2$

examples:

① $p = (1, 0, 0, \ldots 0)$

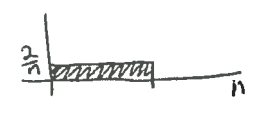$q = (\frac{1}{n}, \frac{1}{n}, \ldots \frac{1}{n})$

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

$l_2$-distance:
$\|p-q\|_2^2 = (1-\frac{1}{n})^2 + (n-1)(\frac{1}{n})^2$
$\approx 1$

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

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

as far as possible ↓

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

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

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

↑ pretty close

$\Rightarrow$ so $l_2$-distance can be weird
(but above fact limits how weird)

## "Plug-in Estimate"

**Algorithm:**

- Take $m$ samples from $p$
- $\forall x$, estimate $p(x)$ by $\hat{p}(x) = \dfrac{\#\text{ times } x \text{ appears in sample}}{m}$
- if $\sum_X |\hat{p}(x) - \frac{1}{n}| > \varepsilon$  reject

  else  accept

**Analysis:** (better analyses exist, e.g. next page)

pick $m$ st. whp $\forall x$  $|p(x) - \hat{p}(x)| < \varepsilon/n$  (assume this holds in following)

$$\Rightarrow \|\hat{p} - p\|_1 < \varepsilon$$

**Correct behavior** {

so if $p = U$, test will accept

to show if $\|p - U\|_1 > 2\varepsilon$, likely to reject, will show contra positive:

If test accepts,   *good approx*   *test accepts*

by $\triangle \neq$ : if $\|p - \hat{p}\|_1 < \varepsilon$ + $\|\hat{p} - U\|_1 < \varepsilon$

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

but how big should $m$ be?

$\Omega(n)$? coupon collector?   ...   :)

Above algorithm gives good approx in $O(n/\varepsilon^2)$ samples:

Thm if $m = \Theta(\frac{n}{\varepsilon^2})$, $\Pr\left[\|\hat{p}-p\|_1 \leq \varepsilon\right] \geq 3/4$

Better analysis : (not done in lecture)

← def of $\hat{p}$ + lin of expectation

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

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

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

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

Jensen's ≠

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

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

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

$$\leq \frac{1}{\sqrt{m}} \cdot \sqrt{n} \leftarrow \text{since} \atop \max_{p \in \text{prob dist} \atop \text{over domain} \atop \text{of size } n} \sum \sqrt{p(x)} \text{ is } \sqrt{n}$$

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

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

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

Let's consider an "easier" problem — $L_2$-distance

## $L_2$ - Distance (squared):

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

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

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

Collision probability of $p$:

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

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 $p$    ① how many samples?
2. let $\hat{c} \leftarrow$ estimate of $\|p\|_2^2$ from sample    ② how?
3. if $\hat{c} < \tfrac{1}{n} + \delta$ pass    ③ what should $\delta$ be?
   else   fail

**Thm** if $s = \theta(\sqrt{n}/\varepsilon^4)$, $\Pr\left[|\hat{c} - \|p\|_2^2| > \tfrac{\varepsilon^2}{2}\right] \leq \tfrac{1}{4}$

$\rightarrow$ Algorithm is property tester for uniformity under $L_2$-dist.

**Naive idea:** (pair off samples)

take two new samples:

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

$\}$ $\sigma_i$'s are independent

" 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 $p$ "

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

$\}$ $\sigma_{ij}$'s are not independent

**Estimate by recycling:**

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

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

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

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

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

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

$= \|p\|_2^2$

Question ③

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

(& how to pick $\delta$ ?)

Pick $\delta = \Delta = \varepsilon^2/2$

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 $\hat{c} \overset{by ✱}{\leq} \|U_{[n]}\|_2^2 + \Delta = \frac{1}{n} + \frac{\varepsilon^2}{2}$

  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

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

  & $\hat{c} > \|p\|_2^2 - \Delta$ ← ✱

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

  so test will FAIL

Remaining Question: (Question 1)

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

# Question ①:

**Analysis:**
$$E[\delta_{ij}] = Pr[\delta_{ij} = 1] = \|p\|_2^2$$

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

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

$$\underline{Fact} \quad Var[aX] = a^2 Var[X]$$

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

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

**Lemma** $\quad Var\left[\sum_{i<j} \delta_{ij}\right] \leq O\left(s^3 \cdot \|p\|_2^3\right)$

**Corr** $\quad Var[\hat{c}] \leq O\left(\|p\|_2^3 / s\right)$

**Proof of lemma** $\quad \underline{def} \quad \overline{\delta_{ij}} = \delta_{ij} - E[\delta_{ij}] \qquad$ (nice trick)

$$\text{note} \quad E[\overline{\delta_{ij}}] = 0 \quad \& \quad \overline{\delta_{ij}} < \delta_{ij} \quad \text{since } E[\delta_{ij}] > 0$$

$$\text{also} \quad E[\overline{\delta_{ij}}\,\overline{\delta_{k\ell}}] \leq E[\delta_{ij}\,\delta_{k\ell}]$$

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

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

$$= E\left[\underbrace{\sum_{\substack{i<j \\ k<l}} \bar{b}_{ij}\,\bar{b}_{kl}}_{\substack{|\{i,j,k,l\}|=2 \\ \textcircled{1}}} + \underbrace{\sum_{\substack{i<j \\ k<l}} \bar{b}_{ij}\,\bar{b}_{kl}}_{\substack{|\{i,j,k,l\}|=3 \\ \textcircled{2}}} + \underbrace{\sum_{\substack{i<j \\ k<l}} \bar{b}_{ij}\,\bar{b}_{kl}}_{\substack{|\{i,j,k,l\}|=4 \\ \textcircled{3}}}\right]$$

## Bounding $\textcircled{1}$:

$$E\left[\sum_{\substack{i<j \\ k<l \\ |\{i,j,k,l\}|=2}} \bar{b}_{ij}\,\bar{b}_{kl}\right] = E\left[\sum_{i<j} \bar{b}_{ij}^{\,2}\right] \leq E\left[\sum_{i<j} b_{ij}^{2}\right] = \binom{s}{2}\|p\|_2^2$$

note $b_{ij}^2 = b_{ij}$ since indicator variable

## Bounding $\textcircled{3}$:

$$E\left[\sum_{\substack{i<j \\ k<l \\ |\{i,j,k,l\}|=4}} \bar{b}_{ij}\,\bar{b}_{kl}\right] \underset{\substack{\text{independence} \\ \text{+linearity of expectations}}}{=} \sum_{\substack{i<j \\ k<l \\ |\{i,j,k,l\}|=4}} E[\bar{b}_{ij}]\,E[\bar{b}_{kl}] = 0$$

note: moving to $\bar{b}_{ij}$ means all of these terms drop out!

## Bounding $\textcircled{2}$:

$$E\left[\sum_{\substack{i<j \\ k<l \\ |\{i,j,k,l\}|=3}} \bar{b}_{ij}\,\bar{b}_{kl}\right] \leq E\left[\sum_{\substack{i<j \\ k<l \\ |\{i,j,k,l\}|=3}} b_{ij}\,b_{kl}\right] = \sum_{\substack{i<j \\ k<l \\ |\{i,j,k,l\}|=3}} E[b_{ij}\,b_{kl}] = \sum_{\substack{|\{a,b,c\}|=3}} \Pr[X_a = X_b = X_c]$$

e.g. $i=k, j\neq k, i\neq l, j\neq l\ldots$

\# ways to pick $i<j$ $k<l$ s.t. $|\{i,j,k,l\}|=3$: (pick 3 indices, pick one to be repeated twice leaves 2 options)

$$\leq 6 \cdot \binom{s}{3} \cdot \sum_x p(x)^3$$

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

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

$$\leq O\left(s^3 \cdot \left(\|p\|_2^2\right)^{3/2}\right) = O\left(s^3 \|p\|_2^3\right)$$

So: $\quad \mathrm{Var}\left[\sum_{i<j} b_{ij}\right] = O\left( \binom{s}{2} \|p\|_2^2 + 0 + s^3 \cdot \|p\|_2^3 \right)$

$$= O\left( s^3 \|p\|_2^3 \right)$$

So how many samples?

for property tester wrt $L_2$-distance

need to estimate $\|p\|_2^2$ to within (additive) $\Delta = \frac{\varepsilon^2}{2}$

$$\Pr\left[ \left| \hat{c} - \|p\|_2^2 \right| > \frac{\varepsilon^2}{2} \right] \leq \frac{\mathrm{Var}[\hat{c}]}{\varepsilon^4/4} = \frac{1}{\binom{s}{2}^2} \cdot \frac{s^3 \cdot \|p\|_2^3}{\varepsilon^4} \cdot \mathrm{Const}$$

$$\leq O\left( \frac{1}{s} \cdot \frac{1}{\varepsilon^4} \cdot \underbrace{\|p\|_2^3}_{\leq 1} \right)$$

$\underbrace{\phantom{xxxxxx}}$
want this
to be small

Pick $\quad s = \Omega\left(\frac{1}{\varepsilon^4}\right)$ $\qquad\qquad$ (can do better)

what about $L_1$-distance?

**Now:** Distinguish $\|p - U\|_1 \geq \varepsilon$ from $p = U$

via $L_2$ - testing

**Thm** there is distribution testing algorithm

which tests uniformity (in $L_1$) & outputs

correct answer w/ prob $\geq 1 - \delta$

using $O\left(\frac{1}{\varepsilon^4} \sqrt{n} \, \log(1/\delta)\right)$ samples

**why?**

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

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

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

need to get **multiplicative** estimate of $\|p\|_2^2$ to w/in $\left(1 \pm \frac{\varepsilon^2}{4}\right)$ $\quad \overset{\gamma}{\underset{\|}{=}} \overset{do}{\underset{\swarrow \text{this}}{}}$

$\quad\quad$ or **additive** " " " " " $\frac{\varepsilon^2}{2n}$

$\Pr\left[\,|\hat{c} - \|p\|_2^2| > \gamma \cdot \|p\|_2^2\right] \leq \frac{\mathrm{Var}[\hat{c}]}{\gamma^2 \|p\|_2^4} \leq \frac{\text{Const} \cdot \|p\|_2^3 / s}{\gamma^2 \|p\|_2^4} = \frac{\text{Const}}{s \gamma^2 \cdot \|p\|_2}$

$\quad\quad\quad\quad\quad\quad \uparrow \atop \geq \frac{1}{\sqrt{n}}$

$\quad\quad\quad\quad\quad\quad\quad\quad \leq O\left(\frac{\sqrt{n}}{s \gamma^2}\right)$

$\quad\quad\quad$ So pick $s = \Omega\left(\frac{\sqrt{n}}{\varepsilon^4}\right)$