# Lecture 3

- concentration of measures. (cont.)

## Distribution testing
- uniformity testing

# Useful tools for concentration  (recap)

1) Markov     for   non-negative  r.v. $X$

$$\Pr[\; X > a\;] \leq \frac{E[X]}{a}$$

2)  chebyshev

$$\Pr[\; |X - E[X]| > a\;] \leq \frac{Var[X]}{a^2}$$

3)  chernoff

Sum of $n$   i.i.d   Bernoulli random variables

$$S = \sum_{i=1}^{n} X_i \qquad X_i \sim Ber(p) \quad , \quad \varepsilon \in [0,1]$$

$$\Pr\left[\; \frac{S}{n} > p(1+\varepsilon)\;\right] \leq e^{-np\varepsilon^2/3}$$

$$\Pr\left[\; \frac{S}{n} < p(1-\varepsilon)\;\right] \leq e^{-np\varepsilon^2/2}$$

4) Hoeffding    (same condition as (3))

$$\Pr\left[ \frac{S}{n} > p + \varepsilon \right] \le e^{-2m\varepsilon^2}$$

$$\Pr\left[ \frac{S}{n} < p - \varepsilon \right] \le e^{-2m\varepsilon^2}$$

# distribution testing

An $(\varepsilon, \delta)$ - tester for property $P$

we have an unknown distribution $d$

We aim to design an algorithm $A$
that distinguishes the following w.p. $\geq 1-\delta$:

- if $d \in P$, $A$ outputs accept

- if $d$ is $\varepsilon$-far from $P$, $A$ outputs reject

what is a property?

$P$ = a set of distributions

$P = \{U_n\}$ → a uniform dist. on $[n]$

$P = \{$ a set of unimodal distributions $\}$

$d$ is $\varepsilon$-far iff $\text{dist}(d, P) > \varepsilon$

$$\text{dist}(d, P) = \min_{d' \in P} \text{dist}(d, d')$$

Example distances:

$\ell_1$ - distance: $\|d - d'\|_1 = \sum_{x \in \Omega} |d(x) - d'(x)|$

$\ell_2$ - distance : $\|d - d'\|_2 = \sqrt{\sum_{x \in \Omega} (d(x) - d(x'))^2}$

Total variation distance : $\|d - d'\|_{TV} = \max_{E \subseteq \Omega} |d(E) - d(E')|$

(statistical distance)

$\hookrightarrow$ every event

Turns out $\qquad \|d - d'\|_{TV} = \dfrac{1}{2} \|d - d'\|_1$

---

Today's question : uniformity testing

Design algorithm $A$ that receives $n, \varepsilon, \delta,$ and samples from $d$ and outputs

- accept w.p. $\geq 1 - \delta$ if $d = U_n$

- reject w.p. $\geq 1 - \delta$ if $\|d - U_n\|_1 > \varepsilon$

Q: Which one look like a real dice?

2     3     1     4     6     1

4     6     4     3     4     5

$Q_2$ what did give it away?

$A_2$ repetitions! ⤳ samples from a uniform distribution
looks "less" repeated.

Let's formalize this intuition...

collisions : two samples that are equal to
each other

# collisions in the sample set, tells
us if a distribution is uniform or not.

Algorithm:

Draw $m$ samples from $d$: $X_1, \ldots, X_m$

$\forall \; i < j \in [m]: \quad \sigma_{ij} = \begin{cases} 1 & \text{if } X_i = X_j \\ 0 & \text{o.w.} \end{cases}$

$$Y \leftarrow \sum_{i=1}^{m} \sum_{j > i}^{m} \sigma_{ij} \Big/ \binom{m}{2}$$

if $Y < t$

         output    accept
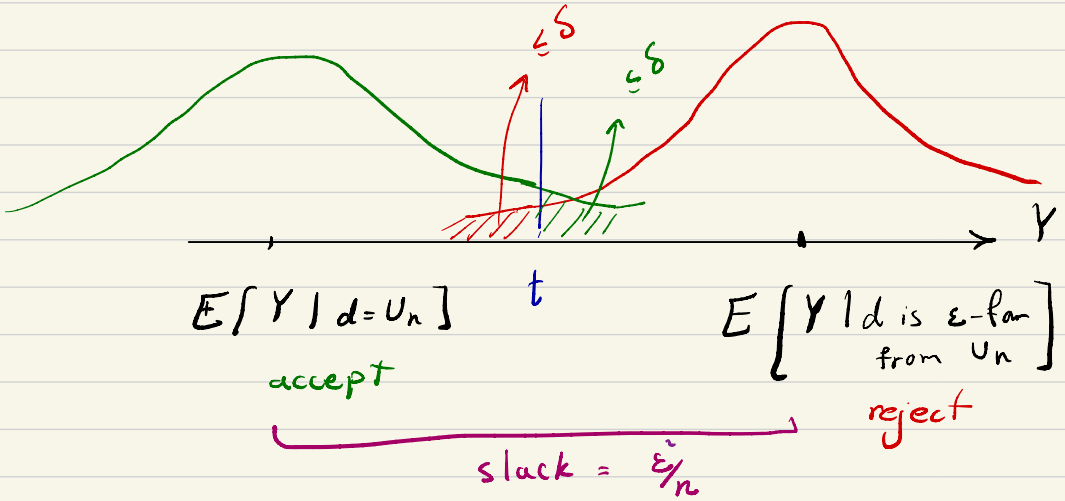
else

         output    reject

Our goal here: what should $m$ & $t$ be?

Visual description



First step: slack exists

$$\mathbb{E}[\sigma_{ij}] = \sum_{a=1}^{n} Pr[X_i = a] \cdot Pr[X_j = a]$$

$$= \sum_{a=1}^{n} d_a^2 = \|d\|_2^2$$

$$\mathbb{E}[Y] = \frac{1}{\binom{m}{2}} \sum_{i=1}^{m} \sum_{j=i+1}^{m} \sigma_{ij} = \|d\|_2^2$$

**Case 1:** d is uniform

if $d = U_n$: $\|d\|_2^2 = \sum_{a=1}^{n} d_a^2 = n \times \frac{1}{n^2} = \frac{1}{n}$

**Case 2:** d is $\varepsilon$-far from uniform

if $\|d - U_n\|_1 > \varepsilon$:

$$\|d\|_2^2 = \sum_{a=1}^{n} d_a^2 = \sum_{a=1}^{n} \left( \frac{1}{n} + (d_a - \frac{1}{n}) \right)^2$$

$$= \sum_{a=1}^{n} \frac{1}{n}^2 + \frac{2}{n} \left( d_a - \frac{1}{n} \right) + \left( d_a - \frac{1}{n} \right)^2$$

$$= \frac{1}{n} + \frac{2}{n} \underbrace{\left( \sum_{a=1}^{n} d_a - \frac{1}{n} \right)}_{= 0} + \sum_{a=1}^{n} (d_a - \frac{1}{n})^2$$

$$= \frac{1}{n} + \underbrace{\| d - U_n \|_2^2}_{\text{our slack}}$$

- Our conjecture is correct Y "tends" to be larger when $d$ is $\varepsilon$-far from uniform.

How far ?

we know $\| d - U_n \|_1 > \varepsilon$

Cauchy-schwarz: $\left( \sum \alpha_i^2 \right) \cdot \left( \sum y_i^2 \right) \geq \left( \sum x_i y_i \right)^2$ $\Big\} \Longrightarrow$

$$\left( \sum_a \left( d_a - \frac{1}{n} \right)^2 \right) \cdot \left( \sum_{a=1}^{n} 1^2 \right) \geq \left( \sum \left| d_a - \frac{1}{n} \right| \right)^2$$
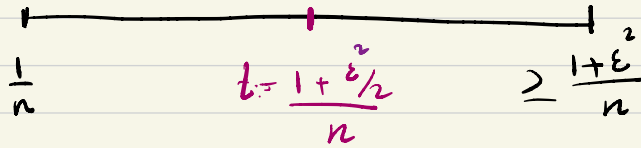
$\Longrightarrow$

$$\| d - U_n \|_2^2 = \sum_{a=1}^{n} \left( d_a - \frac{1}{n} \right)^2 \geq \frac{\left( \sum \left| d_a - \frac{1}{n} \right| \right)^2}{n}$$

$$= \frac{\| d - U_n \|_1^2}{n} > \frac{\varepsilon^2}{n}$$

$E[Y \mid d = U_n]$                   $E[Y \mid d \text{ is } \varepsilon\text{-far}]$

$$\frac{1}{n} \qquad\qquad t := \frac{1 + \varepsilon^2/2}{n} \qquad\qquad \geq \frac{1 + \varepsilon^2}{n}$$

**Next step :** Concentration

Let set $t$ to be in the middle : $t \leftarrow \frac{1 + \varepsilon^2/2}{n}$

If we show the following, we get an

$(\varepsilon, \delta)$ — tester

① $\Pr\left[ Y \geq \frac{1 + \varepsilon^2/2}{n} \;\middle|\; d = U_n \right] \leq \delta$     $\delta = 0.1$

② $\Pr\left[ Y \leq \frac{1 + \varepsilon^2/2}{n} \;\middle|\; d \text{ is } \varepsilon\text{-far from } U_n \right] \leq \delta$    $\delta = 0.1$

$$Y = \frac{1}{\binom{m}{2}} \sum_{i<j} \sigma_{ij}$$

not a great candidate
for chernoff. bound

(why ?)

Our plan : Using chebyshev's

Lets compute the variance of Y

Lemma 1 Var ( Y ) $= \frac{1}{\binom{m}{2}^2} \cdot \left( \binom{m}{2} \|d\|_2^2 + 6\binom{m}{3} \|d\|_3^3 \right)$

proof is deferred for now .

**Case 1 :**    $d = U_n$

$$\Pr\left[\,|Y - E[Y]| \geq \frac{\varepsilon^2}{2n}\,\right] \leq \frac{Var\,(Y)}{(\varepsilon^2/2n)^2}$$

$$\leq \frac{1}{\binom{m}{2}^2} \cdot \left(\binom{m}{2}\|d\|_2^2 + 6\binom{m}{3}\|d\|_3^3\right) \cdot \frac{4n^2}{\varepsilon^2}$$

$$= \Theta\left(\frac{n^2}{m^4\varepsilon^4} \cdot \left(m^2 \cdot \frac{1}{n} + \frac{m^3}{n^2}\right)\right)$$

$$= \Theta\left(\frac{n}{m^2\,\varepsilon^4} + \frac{1}{m\varepsilon^4}\right) \leq 0.1$$

if  $m = c \cdot \left(\frac{1}{\varepsilon^4} + \frac{\sqrt{n}}{\varepsilon^2}\right)$

for sufficiently large $c$

**Case 2:** $\|d - U_n\|_1 > \varepsilon$

The bound on the variance can be large.

$$\binom{m}{2} \|d\|_2^2 + 6\binom{m}{3} \|d\|_3^3$$

Could be problematic if we require $|Y - E[Y]| \leq \frac{b}{n}$

$\hookrightarrow$ adjust the length accordingly

$$\Pr\left[ Y - E[Y] \geq \frac{\varepsilon^2}{2} E[Y] \right] \leq 4 \frac{\text{Var}[Y]}{\varepsilon^4 E[Y]^2}$$

$$\leq \frac{1}{\binom{m}{2}^2} \cdot \frac{\binom{m}{2} \|d\|_2^2 + 6\binom{m}{3} \|d\|_3^3}{\varepsilon^4 \|d\|_2^4} =$$

$$= \theta\left( \frac{1}{m^2 \cdot \varepsilon^4 \|d\|_2^2} + \frac{\|d\|_3^3}{m \, \varepsilon^4 \|d\|_2^4} \right) \leq 0.1$$

$$m = c \cdot \frac{\sqrt{n}}{\varepsilon^4}$$

using $\|d\|_3^3 \leq \|d\|_2^3$

$\Uparrow$

$\ell_p$-norm inequality $\|d\|_3 \leq \|d\|_2$

**Lemma    1**  $\mathrm{Var}(Y) = \frac{1}{\binom{m}{2}^2} \cdot \left( \binom{m}{2} \|d\|_2^2 + 6\binom{m}{3}\|d\|_3^3 \right)$

proof:

$$\mathrm{Var}(Y) = \mathrm{Var}\left( \frac{1}{\binom{m}{2}} \sum_{i<j} \sigma_{ij} \right)$$

$$= \frac{1}{\binom{m}{2}^2} \mathrm{Var}\left( \sum_{i<j} \sigma_{ij} \right)$$

$$= \frac{1}{\binom{m}{2}^2} \left( E\left[ \left( \sum_{i<j} \sigma_{ij} \right)^2 \right] - \underbrace{\left( \sum_{i<j} E[\sigma_{ij}] \right)^2}_{\|d\|_2^2} \right)$$

$$= \frac{1}{\binom{m}{2}^2} E\left[ \sum_{i<j} \sum_{\ell<k} \sigma_{ij}\, \sigma_{\ell k} \right]$$

$$- \|d\|_2^4$$

$$E\left[\ d_{ij}^{2}\ \right]\ =\ \|d\|_{2}^{2}$$

$$E\left[\ d_{ij}\ d_{lk}\ \right]\ =\ \|d\|_{3}^{3}$$

$\hookrightarrow$ Pr [ three samples are equal]

$$E\left[\ d_{ij}\ d_{lk}\ \right]\ =\ E\left[\ d_{ij}\right]\cdot E\left[d_{lk}\right]$$

$$=\ \|d\|_{2}^{4}$$

$\nearrow\ \binom{3}{2}\cdot(\tfrac{3}{2}-1)$

$$\Rightarrow\ \text{Var}\ [Y]\ =\ \frac{1}{\binom{m}{2}^{2}}\left[\ \binom{m}{2}\cdot\|d\|_{2}^{2}\ +\ 6\binom{m}{3}\|d\|_{3}^{3}\right.$$

$$\left.+\ \binom{m}{2}\binom{m-2}{2}\|d\|_{2}^{4}\ -\ \binom{m}{2}^{2}\|d\|_{2}^{4}\right]$$

$$\leq\ \frac{1}{\binom{m}{2}^{2}}\left[\ \binom{m}{2}\|d\|_{2}^{2}\ +\ 6\binom{m}{3}\|d\|_{3}^{3}\right]\ \square$$

$$\binom{m}{2}\ +\ 6\binom{m}{3}\ +\ \binom{m}{2}\binom{m-2}{2}\ =\ \binom{m}{2}^{2}$$