Today: Continue uniformity testing

Review

Model:



unknown distribution D on $[n]$ → $X_1, X_2, \ldots$ → Algorithm → Solution

independent samples

Task: Uniformity testing

① If $D = U_{[n]}$, output YES w.p. $99/100$

↑ uniform distribution on $[n]$

② If $d_{TV}(D, U_{[n]}) \geq \varepsilon$, output NO w.p. $99/100$

Want to use as few samples as possible

# Algorithm proposed last time

- collect $s = C \cdot \sqrt{n}/\varepsilon^4$ independent samples $X_1, X_2, \ldots, X_s$ from $D$

   *sufficiently large constant* (pointing to $C$)

- count collisions:

$$Y = \sum_{i < j} Y_{ij} = \begin{cases} 1 & \text{if } X_i = X_j \\ 0 & \text{otherwise} \end{cases}$$

- if $Y/\binom{s}{2} \geq \frac{1}{n} + \frac{2\varepsilon^2}{n}$
  output NO
  else output YES

---

Analysis last time:

$$\mathbb{E}\left[ Y/\binom{s}{2} \right] = \|D\|_2^2$$

*distribution treated as $n$-dimensional vector of probabilities* (pointing to $\|D\|_2^2$)

① : $D = U_{[n]}$
$\Rightarrow \|D\|_2^2 = \frac{1}{n}$

② : $d_{TV}(D, U_{[n]}) \geq \varepsilon$
$\Rightarrow \|D\|_2^2 \geq \frac{1 + 2\varepsilon^2}{n}$

Our approach: show that $Y / \binom{s}{2}$ is a good estimator for $\|D\|_2^2$

Lemma: $\mathrm{Var}[Y] \leq 7 \left( \binom{s}{2} \|D\|_2^2 \right)^{3/2}$

Definition: $\bar{Y}_{ij} = Y_{ij} - \mathbb{E}[Y_{ij}]$

$\left( \text{of course, } \mathbb{E}[\bar{Y}_{ij}] = 0 \right)$

Useful facts:

① $\mathbb{E}[\bar{Y}_{ij} \bar{Y}_{kl}] \leq \mathbb{E}[Y_{ij} Y_{kl}]$

Why? $\mathbb{E}[\bar{Y}_{ij} \bar{Y}_{kl}] = \mathbb{E}\left[ (Y_{ij} - \underbrace{\mathbb{E}[Y_{ij}]}_{\|D\|_2^2})(Y_{kl} - \underbrace{\mathbb{E}[Y_{kl}]}_{\|D\|_2^2}) \right]$

$= \mathbb{E}\left[ Y_{ij} Y_{kl} - Y_{ij} \|D\|_2^2 - Y_{kl} \|D\|_2^2 + \left( \|D\|_2^2 \right)^2 \right]$

$= \mathbb{E}[Y_{ij} Y_{kl}] - 2 \left( \|D\|_2^2 \right)^2 + \left( \|D\|_2^2 \right)^2 \leq \mathbb{E}[Y_{ij} Y_{kl}]$

② $\|D\|_3 \leq \|D\|_2$

This in an inequality on norms that follows from Hölder's inequality

20-3

③ $\dfrac{s^2 \leq 3\binom{s}{2}}{}$

$s^2 \leq \dfrac{3}{2} s(s-1)$

⇕

$s \leq \dfrac{3}{2} s - \dfrac{3}{2}$

⇕

$3 \leq s \longleftarrow$ 

$\boxed{\text{true if the constant } C \text{ defining } s \text{ is large enough}}$

④ $\dfrac{\binom{s}{3} \leq \dfrac{s^3}{6}}{}$

$\binom{s}{3} = \dfrac{s(s-1)(s-2)}{6} \leq \dfrac{s^3}{6}$

---

Proof of the variance lemma:

$$Var[Y] = Var\left[\sum_{i<j} Y_{ij}\right] = Var\left[\sum_{i<j} \bar{Y}_{ij}\right]$$

$$= \mathbb{E}\left[\left(\sum_{i<j} \bar{Y}_{ij}\right)^2\right] - \underbrace{\left(\mathbb{E}\left[\sum_{i<j} \bar{Y}_{ij}\right]\right)^2}_{=0}$$

20-4

$$= \mathbb{E}\left[ \underbrace{\sum_{i<j} \bar{Y}_{ij}}_{\textcircled{1}} + \underbrace{\sum_{\substack{i<j \\ k<l \\ i,j,k,l \text{ distinct}}} \bar{Y}_{ij} \bar{Y}_{kl}}_{\textcircled{2}} \right.$$

$$+ \underbrace{\sum_{\substack{i<j \\ i<l \\ i,j,l \text{ distinct}}} \bar{Y}_{ij} \bar{Y}_{il}}_{\textcircled{3}} + \underbrace{\sum_{\substack{i<j \\ k<j \\ i,j,k \text{ distinct}}} \bar{Y}_{ij} \bar{Y}_{kj}}_{\textcircled{4}}$$

$$\left. + \underbrace{\sum_{\substack{i<j \\ j<l}} \bar{Y}_{ij} \bar{Y}_{jl}}_{\textcircled{5}} + \underbrace{\sum_{\substack{i<j \\ k<i}} \bar{Y}_{ij} \bar{Y}_{ki}}_{\textcircled{6}} \right]$$

20 - 5

We analyze the expectation term
by term:

△1 : $\mathbb{E}\left[\sum_{i<j} \bar{Y}_{ij}^2\right] \overset{①}{\leq} \mathbb{E}\left[\sum_{i<j} Y_{ij}^2\right]$

$= \mathbb{E}\left[\sum_{i<j} Y_{ij}\right] = \binom{S}{2}\|D\|_2^2$

②2 : $\mathbb{E}\left[\sum_{\substack{i<j \\ k<l}} \bar{Y}_{ij}\bar{Y}_{kl}\right]$

$i,j,k,l$ independent

$\overset{\text{independence}}{=} \sum \underbrace{\mathbb{E}[\bar{Y}_{ij}]}_{=0}\underbrace{\mathbb{E}[\bar{Y}_{kl}]}_{=0} = 0$

③3 : $\mathbb{E}\left[\sum_{\substack{i<j \\ i<l}} \bar{Y}_{ij}\bar{Y}_{il}\right]$

$i,j,l$ distinct

$\overset{①}{\leq} \mathbb{E}\left[\sum Y_{ij}Y_{il}\right]$

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

$p_i = $ probability of drawing $i$ from $D$

$$\leq \underbrace{2\binom{S}{3}}_{\substack{\text{because for} \\ \text{any selection} \\ \text{of three distinct} \\ \text{elements we} \\ \text{have either} \\ i<j<l \\ \text{or } i<l<j}} \underbrace{\sum_{i=1}^{n} p_i^3}_{= \|D\|_3^3}$$

$$\leq \underbrace{2 \cdot \frac{S^3}{6}}_{\text{due to } \boxed{4}} \cdot \underbrace{\|D\|_2^3}_{\text{due to } \boxed{2}}$$

20-7

$$\geqslant \sqrt{3}\left(\binom{s}{2}\|D\|_2^2\right)^{3/2}$$

due to ③, $s^3 = \left(s^2\right)^{3/2} \leqslant \left(3\binom{s}{2}\right)^{3/2}$

④ : same bound as ③

⑤ & ⑥ : $\leqslant \frac{\sqrt{3}}{2}\left(\binom{s}{2}\|D\|_2^2\right)^{3/2}$ each

Almost the same as ③ & ④,
except for each triple selected from
s options, there is a unique assignment
to indices in the summation. So no
factor of 2 needed.

Overall:
$$\mathrm{Var}[Y] \leqslant \binom{s}{2}\|D\|_2^2 + 3\sqrt{3}\left(\binom{s}{2}\|D\|_2^2\right)^{3/2}$$
If C in the definition of s is large
enough, $\binom{s}{2} \geqslant n$. Then $\binom{s}{2}\|D\|_2^2 \geqslant n \cdot \frac{1}{n} \geqslant 1$,
and therefore, $\binom{s}{2}\|D\|_2^2 \leqslant \left(\binom{s}{2}\|D\|_2^2\right)^{3/2}$.

Hence
$$\text{Var}[Y] \le \underbrace{\left(1 + 3\sqrt{3}\right)}_{\le 7} \left(\binom{s}{2} \|D\|_2^2\right)^{3/2}$$

(to be continued in the next lecture)