



# 6.S896 - Algorithmic Statistics

## Lecture 1: Introduction

### Uniformity Testing

#### Recent News Titles

- ① "Biden's Job Approval 42%" - Gallup  
8/25/23
- ② "Obesity spreads through Social Networks"  
- Harvard Medical School  
7/25/07
- ③ "Better Cybersecurity w/a New Quantum  
Random Number Generator"  
- SciTechDaily  
9/4/23
- ④ "EHR-Safe: generating high-fidelity & privacy  
preserving synthetic electronic health records"  
- Nature Digital Medicine  
8/11/23

Q: How much stock to put on such claims?

Challenge: all these claims involve high-dimensional distributions which are, in general, hard to understand

This Class. Develop  $\left\{ \begin{array}{l} \text{models} \\ \text{mathematical tools} \\ \text{algorithmic tools} \\ \text{complexity tools} \end{array} \right\}$  that

are useful to make sense of high-dimensional data

### Application 1: Polling

underlying assumption:  $n$  people are asked whether they approve Biden & their answers,  $X_1, X_2, \dots, X_n$ , are **independently & identically sampled** from  $\text{Bernoulli}(p)$   
 $\nearrow$  estimand

natural estimator:  $\hat{p} = \frac{1}{n} \sum_i X_i$

$$\mathbb{E}[\hat{p}] = p$$

$$\text{Var}[\hat{p}] = \frac{1}{n} p(1-p) \leq \frac{1}{4n}$$

$$\hookrightarrow \text{std}(\hat{p}) \leq \frac{1}{2\sqrt{n}}$$

so I expect  $\hat{p} \approx p \pm \frac{1}{\sqrt{n}}$  w pr  $\geq 95\%$

$$\hookrightarrow 2 \times \text{std}(\hat{p})$$

a.k.a. to approximate unknown  $p$  to within  $\pm \epsilon$   
it suffices to ask  $\approx \frac{1}{\epsilon^2}$  people

Question: is  $\frac{1}{\epsilon^2}$  necessary?

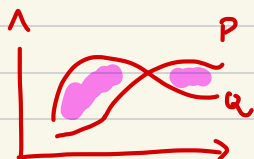
Le Cam's two-point method

setting: two distributions  $P, Q$  over  $\mathcal{X}$

goal: find a test  $T: \mathcal{X} \rightarrow \{P, Q\}$   
such that error probability of  $T$  is small

$$\max \left\{ \Pr_{x \sim P}^{|||} [T(x) = Q], \Pr_{x \sim Q} [T(x) = P] \right\}$$

Lemma: error prob.  $\geq \frac{1}{2} (1 - \text{TV}(P, Q))$



$$\text{TV}(P, Q) = \frac{1}{2} \sum_x |P(x) - Q(x)|$$

$$= \max_{A \subseteq \mathcal{X}} |P(A) - Q(A)|$$

$$\text{TV} = \frac{1}{2} \cdot \text{area of pink shaded region}$$

let's use Le Cam's method to show

Lemma:  $\Omega(\frac{1}{\epsilon^2})$  iid samples from a Bernoulli are necessary to estimate its mean to within  $\pm \epsilon$ .

Stronger Lemma:  $\Omega(\frac{1}{\epsilon^2})$  iid samples are necessary to distinguish Bernoulli( $\frac{1}{2}$ ) vs Bernoulli( $\frac{1}{2} + \epsilon$ )

Proof: Suppose  $P = \text{Ber}(\frac{1}{2})^{\otimes n}$

$$Q = \text{Ber}(\frac{1}{2} + \epsilon)^{\otimes n}$$

$n$  iid samples from  $\text{Ber}(\frac{1}{2}) \leftrightarrow 1$  sample from  $P$

$-||-$   $\text{Ber}(\frac{1}{2} + \epsilon) \leftrightarrow -||-$   $Q$

distinguisher between  $\text{Ber}(\frac{1}{2})$  &  $\text{Ber}(\frac{1}{2} + \epsilon)$   
is a function:

$$T: \{0,1\}^n \rightarrow \{P, Q\}$$

Le Cam: Probability of error  $\geq \frac{1}{2} (1 - \text{TV}(P, Q))$   
↪ ??

Claim:  $\text{TV}(P, Q) \leq O(\sqrt{n} \cdot \epsilon)$

$$\geq \frac{1}{2} (1 - O(\sqrt{n} \cdot \epsilon))$$

$$\geq \frac{1}{4} \text{ unless } n = \Omega\left(\frac{1}{\epsilon^2}\right)$$

□

## Proof of Claim: Tensorization of Kullback-Leibler Divergence

$$\left[ \begin{array}{l} \text{Def: If } P, Q \text{ distn's over (finite) } \mathcal{X} \\ \\ KL(P \parallel Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)} \\ \\ \text{Fact: } TV(P, Q) \leq \sqrt{\frac{1}{2} KL(P \parallel Q)} \end{array} \right]$$

Now suppose  $P = \text{Ber}(p)^{\otimes n}$  &  $Q = \text{Ber}(q)^{\otimes n}$

then

$$KL(P \parallel Q) = \sum_{x \in \{0,1\}^n} \prod_i P(x_i) \log \frac{\prod_i P(x_i)}{\prod_i Q(x_i)}$$

$$= \sum_x \prod_i P(x_i) \sum_i \log \frac{P(x_i)}{Q(x_i)}$$

$$= \sum_i \underbrace{\sum_{x_{-i}} \prod_{j \neq i} P(x_j)}_1 \cdot \underbrace{\sum_{x_i} P(x_i) \log \frac{P(x_i)}{Q(x_i)}}_{KL(\text{Ber}(p) \parallel \text{Ber}(q))}$$

$$= n \cdot KL(\text{Ber}(p) \parallel \text{Ber}(q))$$

now if  $p = \frac{1}{2}$   $q = \frac{1}{2} + \epsilon$

$$\begin{aligned}
KL(\text{Ber}(\frac{1}{2}) \parallel \text{Ber}(\frac{1}{2} + \epsilon)) &= \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{2} + \epsilon} + \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{2} - \epsilon} \\
&= \frac{1}{2} \log \frac{1}{1 + 2\epsilon} + \frac{1}{2} \log \frac{1}{1 - 2\epsilon} \\
&= -\frac{1}{2} \log(1 - 4\epsilon^2) \\
&= -\frac{1}{2} (-4\epsilon^2 + O(\epsilon^4)) \\
&= 2\epsilon^2 + O(\epsilon^4)
\end{aligned}$$

$$\begin{aligned}
\Rightarrow KL(P \parallel Q) &= 2n \cdot \epsilon^2 + O(n \cdot \epsilon^4) \\
&= O(n \cdot \epsilon^2)
\end{aligned}$$

$$\Rightarrow TV(P \parallel Q) = O(\sqrt{n} \cdot \epsilon)$$



Proof of Le Cam:

$$\begin{aligned}
&\max \left\{ \Pr_{x \sim P}(T(x) = Q), \Pr_{x \sim Q}(T(x) = P) \right\} \\
&\geq \frac{1}{2} \left( \Pr_{x \sim P}(T(x) = Q) + \Pr_{x \sim Q}(T(x) = P) \right) \\
&= \frac{1}{2} \left( \sum_x P(x) \cdot \mathbb{1}_{\{T(x) = Q\}} + \sum_x Q(x) \cdot \mathbb{1}_{\{T(x) = P\}} \right)
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \left( \sum_x P(x) \cdot (1 - \mathbb{1}_{\{T(x)=P\}}) + \sum_x Q(x) \cdot \mathbb{1}_{\{T(x)=P\}} \right) \\
&= \frac{1}{2} \cdot \left( 1 - \sum_x (P(x) - Q(x)) \cdot \mathbb{1}_{\{T(x)=P\}} \right) \\
&\geq \frac{1}{2} \left( 1 - \sum_{x: P(x) > Q(x)} (P(x) - Q(x)) \right) \\
&= \frac{1}{2} (1 - \text{TV}(P, Q)) \quad \square
\end{aligned}$$

## Application 2: Uniformity Testing

- Big domain  $\mathcal{X}$ ,  $|\mathcal{X}| = N$
- Distribution  $P$  over  $\mathcal{X}$

(e.g. that Quantum random number generator)

Question: is  $P$  uniform?

More precise question:

distinguish using iid samples from  $P$  between

$$P = U(\mathcal{X}) \quad \text{vs} \quad \text{TV}(P, U(\mathcal{X})) \geq \epsilon$$



How many samples do we need?

Thm:  $\Theta(\frac{\sqrt{N}}{\epsilon^2})$  samples are necessary & sufficient  
[Paninski]

trouble: if  $\mathcal{X} = \{0,1\}^N$ , #samples  $\frac{2^{N/2}}{\epsilon^2} !!!$

Today: proof of lower bound of  $\Omega(\frac{\sqrt{N}}{\epsilon^2})$

Intuition: unless  $\gtrsim \sqrt{N}$  samples are drawn  
how can we distinguish between  $P$   
uniform over full or uniform  
over half of the domain?

Proof (of lower bound) use Le Cam!

w.l.o.g.  $\mathcal{X} = [N] = \{1, \dots, N\}$

$P = \mathcal{U}([N])^{\otimes n}$  ( $P$  is dist'n of  $n$  samples from  $\mathcal{U}([N])$ )

Ok and what is a bad dist'n  $Q$ ?

if  $Q$  is uniform over a specific half of the domain (e.g. all even numbers in an array it would be easy to distinguish from  $P$ ...

Idea: take  $Q$  to be uniform on a random half of the domain!

$Q$ : - sample  $Z_1, \dots, Z_{N/2}$  i.i.d. from  $\{\pm 1\}$

importantly  
note that  
any  $q$   
thus defined  
satisfies that  
 $TV(q, U[N]) \geq \epsilon$

- define dist'n  $q$  over  $[N]$  s.t.  $q_{2i-1} = \frac{1+Z_i \cdot \epsilon}{2}$

$$q_{2i} = \frac{1-Z_i \cdot \epsilon}{2}$$

& generate  $n$  samples from  $q$

Le Cam:

$$Pr[\text{error}] \geq \frac{1}{2} (1 - \underbrace{TV(P, Q)})$$

↳ how small is  
this as a function  
of  $N, n, \epsilon$ ?

TL;DR Lemma:  $TV(P, Q) \asymp \sqrt{e^{O(n^2 \epsilon^4 / N)} - 1}$

so, unless  $n \geq \Omega\left(\frac{\sqrt{N}}{\epsilon^2}\right)$ ,  $Pr[\text{error}] \geq \frac{1}{4}$

Proof of TL;DR lemma: will sketch this for  
a slightly different setting  
where rather than  $n$  samples  
we draw  $\tilde{n} \sim \text{Poisson}(n)$  samples

• So now:  $P$ : - sample  $\tilde{n} \sim \text{Poisson}(n)$   
- draw  $\tilde{n}$  samples from  $\mathcal{U}([N])$

$Q$ : - sample  $\tilde{n} \sim \text{Poisson}(n)$

- sample  $Z_1, \dots, Z_{N/2}$  i.i.d. from  $\{\pm 1\}$

- define dist'n  $q$  over  $[N]$  s.t.  $q_{2i-1} = \frac{1+Z_i \cdot \varepsilon}{N}$

$$q_{2i} = \frac{1-Z_i \cdot \varepsilon}{N}$$

- draw  $\tilde{n}$  samples from  $q$

• Suppose  $X_i$  is how many times  $i \in [N]$  appeared in the sample

under  $P$ :  $X_1, \dots, X_N$  independent

&  $X_i \sim \text{Poisson}\left(\frac{n}{N}\right)$

$$P(X_1 = n_1, \dots, X_N = n_N) = e^{-n} \prod_{i=1}^N \frac{\left(\frac{n}{N}\right)^{n_i}}{n_i!}$$

under  $Q$ :  $(x_{2i-1}, x_{2i})$  independent from  $(x_{2j-1}, x_{2j})$   
if  $j \neq i$

&  $\Pr[x_{2i-1}, x_{2i}]$  is a mixture

$$\begin{aligned} Q(x_1=n_1, x_2=n_2) &= \frac{1}{2} e^{-\frac{n(1+\epsilon)}{N}} \frac{\left(\frac{n(1+\epsilon)}{N}\right)^{n_1}}{n_1!} \cdot e^{-\frac{n(1-\epsilon)}{N}} \frac{\left(\frac{n(1-\epsilon)}{N}\right)^{n_2}}{n_2!} \\ &\quad + \frac{1}{2} e^{-\frac{n(1+\epsilon)}{N}} \frac{\left(\frac{n(1+\epsilon)}{N}\right)^{n_1}}{n_2!} \cdot e^{-\frac{n(1-\epsilon)}{N}} \frac{\left(\frac{n(1-\epsilon)}{N}\right)^{n_2}}{n_1!} \\ &= \frac{1}{2} e^{-\frac{2n}{N}} \frac{\left(\frac{n}{N}\right)^{n_1} \left(\frac{n}{N}\right)^{n_2}}{n_1! n_2!} \left( (1+\epsilon)^{n_1} (1-\epsilon)^{n_2} + (1-\epsilon)^{n_1} (1+\epsilon)^{n_2} \right) \end{aligned}$$

$$Q(x_1=n_1, \dots, x_N=n_N) = e^{-n} \prod_{i=1}^N \frac{\left(\frac{n}{N}\right)^{n_i}}{n_i!} \prod_{i=1}^{N/2} \left( (1+\epsilon)^{n_{2i-1}} (1-\epsilon)^{n_{2i}} + (1-\epsilon)^{n_{2i-1}} (1+\epsilon)^{n_{2i}} \right)$$

next we use that  $Tv(P, Q) \leq \sqrt{\frac{1}{2} KL(Q \| P)}$

$$\leq \sqrt{\frac{1}{2} \chi^2(Q \| P)}$$

$$\chi^2(Q \| P) = \mathbb{E}_{x \sim Q} \left[ \frac{Q(x)}{P(x)} \right] - 1$$

$$\hookrightarrow \left( \frac{e^{2n\epsilon^2/N} + e^{-2n\epsilon^2/N}}{2} \right)^{N/2} - 1$$