

# CSC2412: Adaptive Data Analysis via Differential Privacy

---

*Sasho Nikolov*

# The adaptive data analysis problem

---

## Estimating population counts

- Unknown distribution  $D$  on  $\mathcal{X}$

↪ models the population

↪ universe of possible data points

- Predicates  $q_1, \dots, q_k : \mathcal{X} \rightarrow \{0, 1\}$

E.g.  $q_1 = ?$  smoker  
 $q_2 = ?$  smoker and male  
 $q_3 = ?$  smoker and PhD  
...

Want to estimate, for all  $i = 1 \dots k$ :

$$q_i(D) = \mathbb{E}_{x \sim D}[q_i(x)].$$

↪ fraction of the population satisfying  $q_i$

## The classical solution

Draw a sample  $X = \{x_1, \dots, x_n\}$  iid from  $D$ .

Hope that  $\forall i: q_i(X) \approx q_i(D)$

$$q_i(X) = \frac{1}{n} \sum_{j=1}^n q_i(x_j)$$

↳ independent, in  $\{q_i\}$

$$\mathbb{E}[q_i(X)] = q_i(D)$$

Hoeffding:  $\forall i: \mathbb{P}(|q_i(X) - q_i(D)| > \alpha) = \mathbb{P}(|q_i(X) - \mathbb{E}q_i(X)| > \alpha)$

$$\leq 2 e^{-2n\alpha^2}$$

$$\mathbb{P}(\exists i: |q_i(X) - q_i(D)| > \alpha) \leq 2k \cdot e^{-2n\alpha^2} \leq \beta \quad \text{if}$$

$$n \geq \frac{\ln(2k/\beta)}{2\alpha^2}$$

## Adaptive queries?

What if  $q_i$  depends on  ~~$q_1, \dots, q_{i-1}$~~ ? the estimates for  $q_1(D), \dots, q_{i-1}(D)$

E.g.  $q_i$  is chosen based on  $q_1(X), \dots, q_{i-1}(X)$

E.g.  $q_1 = ?$  smokers and male }  $\rightarrow$  if even split  
 $q_2 = ?$  smokers and female } ask  $q_3 = ?$  smokers and  $\geq 35$  yrs

---

Suppose we ask  $q_1(X), q_2(X), \dots, q_k(X)$  for  $k \gg n$ ,  $q_i$  random predicates  
and we "invert" to learn  $X$  else stop

$q_{k+1}(x) = \begin{cases} 1 & x \in X \\ 0 & \text{o/w} \end{cases} \Rightarrow q_{k+1}(X) = 1$ . But if  $D$  is uniform on  $\mathcal{X}$  then  $q_{k+1}(D) \approx 0$

## A simple solution

Break  $X = \{x_1, \dots, x_n\}$  into  $X^1 = \{x_1, \dots, x_{n/k}\}$

Answer  $q_1(0)$  by  $q_1(x^1)$

$q_2$  by  $q_2(x^2)$

$\vdots$   
 $q_k$  by  $q_k(x^k)$

$$X^2 = \{x_{n/k+1}, \dots, x_{2n/k}\}$$

$$\vdots$$
$$X^k = \{x_{\frac{(k-1)n}{k}+1}, \dots, x_n\}$$

to get error  $d$  w/ prob  $1-\beta$   
I need

$$\frac{n}{k} \geq \frac{\ln(2k/\beta)}{2d^2} \Leftrightarrow$$

$$n \geq \frac{k \ln(2k/\beta)}{2d^2}$$

Can we do better?

# Transfer theorem

$M$  answers  $q_1$  w/  $M(X)_1$   
 $q_2$  determined from  $M(X)_1 \rightarrow M$  answers w/  $M(X)_2$   
by analyst  $\uparrow$  ...

## Theorem

Suppose  $M$  takes a dataset  $X$  and answers  $k$  adaptive queries  $q_1, \dots, q_k$ . If

- $\forall X \in \mathcal{X}^n, \mathbb{P}(\exists i : |q_i(X) - M(X)_i| > \alpha) < \alpha\beta, \rightarrow M$  accurate on the dataset
- $M$  is  $(\alpha, \alpha\beta)$ -DP,

then for a constant  $C$

$$\mathbb{P}_{M, X \sim D^n}(\exists i : |M(X)_i - q_i(D)| > C\alpha) < C\beta.$$

$$X \sim D^n \Leftrightarrow X = \{x_1, \dots, x_n\} \quad x_i \sim D \text{ independently}$$

## Improving on the simple solution

Simple solution; error  $d$  with  $\approx \frac{k \log(k/\beta)}{d^2}$

Can get error  $\alpha$  with  $\approx \frac{\sqrt{k \log k}}{\alpha^2}$  samples.

Gaussian noise + advanced composition

answer  $q_i$  w/  $q_i(x) + z_i$       $z_i \sim \mathcal{N}\left(0, \frac{1}{n^2 \cdot \rho}\right) \approx \frac{k \log(1/\delta)}{n^2 d^2}$

and we get  $(\epsilon, \delta)$ -DP

for any  $\delta$  and  $\epsilon \approx \sqrt{k \rho \ln(1/\delta)}$

Transfer thm: we need  $(d, d\beta)$ -DP

std dev per  $q_i$  is  $\approx \frac{\sqrt{k \ln(1/\delta)}}{n d} = \frac{\sqrt{k \ln(1/(d\beta))}}{n d} \ll d$  if  $n \gg \frac{\sqrt{k \ln(1/(d\beta))}}{d^2}$



## Key Lemma

$$X \in \mathcal{X}^n$$

$q: \mathcal{X} \rightarrow \{0,1\}$   
 $D$  is a distr. on  $\mathcal{X}$

### Lemma

Suppose  $\mathcal{W}$  is  $(\epsilon, \delta)$ -DP, and on input  $X$  outputs a counting query  $q$ . Let  $X \sim D^n$ .

Then

$$|\mathbb{E}[q(D) \mid q = \mathcal{W}(X)] - \mathbb{E}[q(X) \mid q = \mathcal{W}(X)]| \leq \underbrace{e^\epsilon - 1 + \delta}_{\approx \epsilon} \approx \epsilon + \delta$$

↓  
over random choice of  $X \sim D^n$   
and randomness of  $\mathcal{W}$

A DP algorithm cannot find a query that distinguishes  $X$  from  $D$ .

## Proof of Key Lemma

$$q(X) = \frac{1}{n} \sum_{i=1}^n q(x_i) \quad q: \mathcal{X} \rightarrow \{0, 1\}$$

$$\mathbb{E}[q(X) \mid q = \mathcal{W}(X)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[q(x_i) \mid q = \mathcal{W}(X)] = \frac{1}{n} \sum_{i=1}^n \mathbb{P}(q(x_i) = 1 \mid q = \mathcal{W}(X))$$

Take  $x'_i \sim D$  independently from everything else.

$X' = \{x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n\}$   $X, X'$  neighbouring

$$\mathbb{P}(q(x_i) = 1 \mid q = \mathcal{W}(X)) \leq e^\epsilon \mathbb{P}(q(x_i) = 1 \mid q = \mathcal{W}(X')) + \delta$$

$(\epsilon, \delta)$ -DP of  $w$

## Proof part 2

$$X = \{x_1, \dots, x_n\}$$

$$X' = \{x_1, \dots, x'_i, x_{i+1}, \dots, x_n\}$$

Observation:  $(x_i, X')$  has the same distribution as  $(x'_i, X)$

$$\mathbb{P}(q(x_i) = 1 \mid q = w(X)) \leq e^\varepsilon \mathbb{P}(q(x_i) = 1 \mid q = w(X')) + \delta$$

$$q(D) = \mathbb{E}_{x \sim D} q(x) = \mathbb{P}_{x \sim D}(q(x) = 1)$$

$$= e^\varepsilon \mathbb{P}(q(x'_i) = 1 \mid q = w(X)) + \delta$$

$$= e^\varepsilon \mathbb{E}[q(D) \mid q = w(X)] + \delta$$

$$\mathbb{E}[q(X) \mid q = w(X)] \leq e^\varepsilon \underbrace{\mathbb{E}[q(D) \mid q = w(X)]}_{\leq 1} + \delta$$

$$\mathbb{E}[q(X) \mid q = w(X)] - \mathbb{E}[q(D) \mid q = w(X)] \leq e^\varepsilon - 1 + \delta \geq -(e^\varepsilon - 1 + \delta) \text{ analogous } 10$$



## Aside: Generalization from DP

Almost the same proof  
as the lemma (exercise)

### Theorem

For any non-negative loss  $\ell(\theta, (x, y))$ ,  $X = \{(x_1, y_1), \dots, (x_n, y_n)\} \sim D^n$ , and

$$L_X(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\theta, (x_i, y_i)) \quad L_D(\theta) = \mathbb{E}_{(x,y) \sim D}[\ell(\theta, (x, y))],$$

if  $\theta$  is computed by an  $(\epsilon, \delta)$ -DP algorithm, then

$$\mathbb{E}[L_D(\theta)] \leq e^\epsilon \mathbb{E}[L_X(\theta)] + \delta \max_{\theta, x, y} \ell(\theta, (x, y)).$$

Population loss is not much more  
than empirical loss for  
DP algo.

DP  
implies  
generalization

## A simpler transference theorem

### Theorem

If the mechanism  $\mathcal{M}$  satisfies that

1.  $\forall X \in \mathcal{X}^n$ , and all sequence of adaptive queries  $q_1, \dots, q_k$ ,  
 $\mathbb{E}[\max_i |q_i(X) - \mathcal{M}(X)_i|] \leq \alpha$
2.  $\mathcal{M}$  is  $(\epsilon, \delta)$ -DP,

$1-q_1$        $1-q_k$

then

$$\mathbb{E}[\max_i |q_i(D) - \mathcal{M}(X)_i|] \leq \alpha + e^\epsilon - 1 + \delta \approx \alpha + \epsilon + \delta$$

$q_1, \dots, q_k$  are adaptively chosen based on  $\mathcal{M}(X)_1, \dots, \mathcal{M}(X)_k$

$X \sim D^n$

# Proof

$$q_i(x) = 1 \Leftrightarrow 1 - q_i(x) = 0$$

$$q_i(x) = 0 \Leftrightarrow 1 - q_i(x) = 1$$

**Trick:** Suppose that if  $q_i$  is asked, so is  $1 - q_i$ , and is answered by  $1 - \mathcal{M}(X)_i$ .

Then  $\max_{i=1}^k |q_i(D) - \mathcal{M}(X)_i| = \max_{i=1}^k q_i(D) - \mathcal{M}(X)_i$ .

$$|q_i(D) - \mathcal{M}(X)_i| = \max \left\{ q_i(D) - \mathcal{M}(X)_i, \underbrace{\mathcal{M}(X)_i - q_i(D)}_{1 - q_i(D) - (1 - \mathcal{M}(X)_i)} \right\}$$

Define  $\mathcal{W}$  s.t. it \* Simulates  $\mathcal{M}$  on the adaptive

$\mathcal{M}$  is  $(\epsilon, \delta)$ -DP

$\Rightarrow \mathcal{W}$  is  $(\epsilon, \delta)$ -DP

post-processing

\* Outputs  $q_i$  s.t.  $q_i$  has max error

queries  $q_1, \dots, q_k$

$$q_i \text{ s.t. } q_i(D) - \mathcal{M}(X)_i = \max_{j=1}^k q_j(D) - \mathcal{M}(X)_j$$

## Proof pt 2

$$\mathbb{E} \max_{i=1}^k q_i(D) - \mathcal{U}(X)_i = \mathbb{E} [q_i(D) - \mathcal{U}(X)_i \mid q_i = \omega(X)]$$

$$= \mathbb{E} [q_i(D) - q_i(X) \mid q_i = \omega(X)]$$

$$+ \mathbb{E} [q_i(X) - \mathcal{U}(X)_i \mid q_i = \omega(X)]$$

$e^{\varepsilon-1} + \delta$   
by lemma

$$\mathbb{E} \left[ \max_{j=1}^k \underbrace{q_j(X) - \mathcal{U}(X)_j}_{\leq \alpha} \right]$$

$$\leq e^{\varepsilon-1} + \delta + \alpha.$$