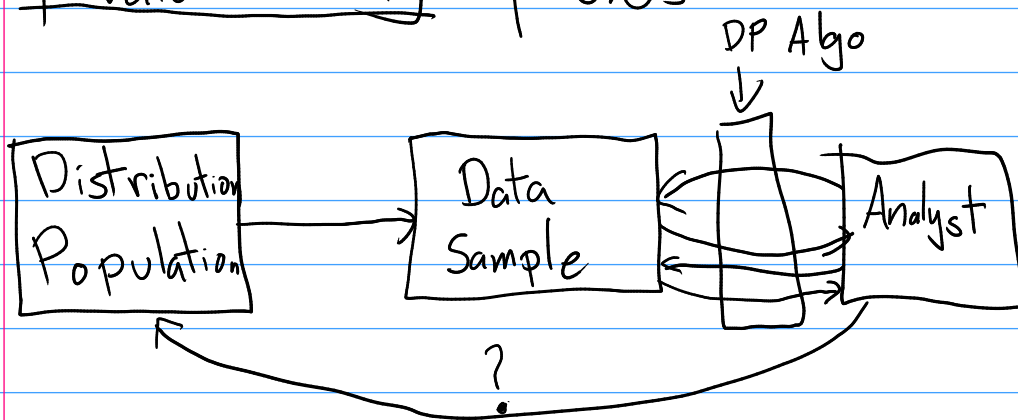# Adaptive Data Analysis
## Setup and Motivation

p-Value hacking   $p < 0.05$



Dist $D$ on $\mathcal{X}$,

$$q_1, \ldots, q_k : \mathcal{X} \to \{0, 1\}$$

$$q_i(D) := \underset{x \sim D}{E}[q_i(x)], \quad \forall i \in [k]$$

$$X = (X_1, \ldots, X_n) \in \mathcal{X}^n, \quad X_j \sim D$$

$$q_i(X) = \frac{1}{n} \sum_{j=1}^{n} q_i(X_j) \qquad q_i(D) \approx q_i(X).$$

$$\forall i \quad \Pr[|q_i(X) - q_i(D)| \geq \alpha] \leq 2\exp(-2\alpha^2 n)$$

$$\Pr[\exists i : |q_i(X) - q_i(D)| \geq \alpha] \leq 2k\exp(-2\alpha^2 n) \leq \beta$$

$$2k/\beta \leq \exp(2\alpha^2 n) \iff n \geq \frac{\log(2k/\beta)}{2\alpha^2}$$

$$n \geq \log k \iff k \leq e^n$$

## Independence

Bad new w/ adaptivity

# What goes wrong adaptively?

$X = \{1, \ldots, k\}$

$q_i : X \to \{0, 1\}$

$q_i(x) = \begin{cases} 1 & \text{if } x = i \\ 0 & \text{else} \end{cases}$

Ask queries $q_1(X), \ldots, q_k(X)$

Learn dataset $X$

$q_{k+1}(x) = \begin{cases} 1 & \text{if } x \in X \\ 0 & \text{else} \end{cases}$

$q_{k+1}(X) = 1$ but $q_{k+1}(D) \leq \eta_k$

Unless $n = \Omega(k)$, no "generalization"

# A naive Solution

Sample splitting

Split $X$ into $k$ parts $X^{(1)}, \ldots, X^{(k)}$

Run $i$th query on $i$th part

$q_i(X^{(i)}) \simeq q_i(D)$ if $n^{(i)} \geq \frac{\log(k/\beta)}{\alpha^2}$

Need $n \geq \frac{k \log(k/\beta)}{\alpha^2}$ (naive)

advanced composition

Today $n \geq \frac{\sqrt{k} \log(2k/\beta) \log(1/\alpha\beta)}{\alpha^2}$ ($\approx$ optimal)

$\log(k)$ (nonadaptive)

($\approx$ optimal)

$n \geq \frac{\sqrt{\log|X| \log(k/\beta) \log(1/\alpha\beta)}}{\alpha^3}$ ($\approx$ optimal)

$\uparrow$ PMW

(McSherry) (Dwork) Feldman Hardt Pitassi Reingold Roth '15

Bassily (Nissim) (Smith) Steinke Stemmer Ullman '16?

# A "Transfer" Theorem

**Theorem 1** *(Transfer Theorem)* *Suppose that the mechanism $\mathcal{M}$ takes $n$ iid samples $X$ from $\mathcal{D}$, and answers $nk$ (adaptive) queries $q_1, \ldots, q_k$ on $X$ such that,*

*i)* $\forall X \in \mathcal{X}^n, \mathbb{P}_{\mathcal{M}}(\exists i : |q_i(X) - \mathcal{M}(X)_i| > \alpha) \leq \beta$, *where $\mathcal{M}(X)_i$ is the answer given by $\mathcal{M}$ to the $i$-th query;* ← $\mathcal{M}$ is accurate on dataset

*ii)* $\mathcal{M}$ *is* $(\alpha, \alpha\beta)$-*DP.* ← $\mathcal{M}$ is private

*Then,*

$$\mathbb{P}_{\mathcal{D},\mathcal{M}}(\exists i : |\mathcal{M}(X) - q_i(\mathcal{D})| > C\alpha) \leq C\beta, \qquad (1)$$

$\mathcal{M}$ generalizes

*for some constant $C$.*

Jung Ligett Neel Roth Sharifi-Malvajerdi Shenfeld '20

# A Simpler Transfer Theorem

**Theorem 3 (Easier Transfer Theorem)** *Let $X \sim \mathcal{D}^n$ and let $\mathcal{M}$ be $(\epsilon, \delta)$-DP such that for every adaptive $q_1, \ldots, q_k$ and for all $X \in \mathcal{X}^n$,*

↖ Private

$$\mathbb{E}_{\mathcal{M},X}[\max_i |q_i(X) - \mathcal{M}(X)_i|] \leq \alpha.$$

← Accurate

*Then,*

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

# Proof of Easier Transfer Theorem

**Theorem 3 (Easier Transfer Theorem)** *Let $X \sim \mathcal{D}^n$ and let $\mathcal{M}$ be $(\epsilon, \delta)$-DP such that for every adaptive $q_1, \ldots, q_k$ and for all $X \in \mathcal{X}^n$,*

$$\mathbb{E}_{\mathcal{M}, X}[\max_i |q_i(X) - \mathcal{M}(X)_i|] \leq \alpha.$$

*Then,*

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

# Key Lemma

**Lemma 2** *Suppose $\mathcal{W}$ is $(\epsilon, \delta)$-DP and on input $X \in \mathcal{X}^n$, it outputs a counting query $q$. Let $X \sim \mathcal{D}^n$ (independent rows). Then,*

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

$$E[q(X) | q = W(x)] = \frac{1}{n} \sum E[q(X_i) | q = W(x)]$$
$$= \frac{1}{n} \sum Pr[q(X_i) = 1 | q = W(x)]$$

$$X_i' \sim D, \quad X' = (X_1, \ldots, X_i', \ldots, X_n)$$

$$Pr[q(X_i) = 1 | q = W(x)] \leq e^\epsilon Pr[q(X_i) = 1 | q = W(x')] + \delta$$

$$Pr[q(X_i) = 1 | q = W(x')] = Pr[q(X_i') = 1 | q = W(x)]$$
$$= E[q(D) | q = W(x)]$$

$$E[q(X) | q = W(x)] \leq e^\epsilon \cdot E[q(D) | q = W(x)] + \delta$$

$$|E[q(X) | q = W(x)] - E[q(D) | q = W(x)]|$$
$$\leq (e^\epsilon - 1) E[q(D) | q = W(x)] + \delta \leq e^\epsilon - 1 + \delta \quad \square$$

# Proof

**Theorem 3 (Easier Transfer Theorem)** *Let $X \sim \mathcal{D}^n$ and let $\mathcal{M}$ be $(\epsilon, \delta)$-DP such that for every adaptive $q_1, \ldots, q_k$ and for all $X \in \mathcal{X}^n$,*

$$\mathbb{E}_{\mathcal{M},X}[\max_i |q_i(X) - \mathcal{M}(X)_i|] \leq \alpha.$$

*Then,*

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

$W$: Choose $q_i$ maximizing $|q_i(x) - M(x)_i|$

Double # of queries. $2k$ queries

$q_1, \ldots, q_k \qquad (1-q_1)_1, \ldots, (1-q_k)$

$M_1(x) \qquad\qquad 1 - M_1(x)$

$\max_i |q_i(D) - M(x)_i| = \max_i q_i(D) - M(x)_i$

$$\mathbb{E}[\max_i q_i(D) - M(x)_i] = \mathbb{E}[q_i(D) - M(x)_i \mid q_i = W(x)] \qquad {}^{+q_i(x) - q_i(x)}$$

$$= \mathbb{E}[q_i(D) - q_i(x) \mid q_i = W(x)] + \mathbb{E}[q_i(x) - M(x) \mid q_i = W(x)]$$

$$\leq (e^\epsilon - 1 + \delta) + (\alpha) \qquad \boxtimes$$