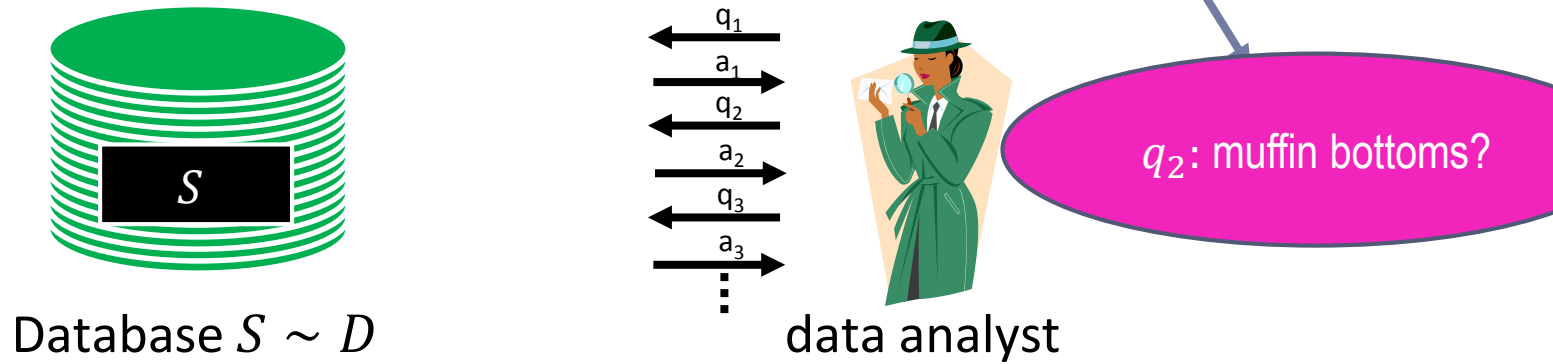


# Universally Adaptive Data Analysis

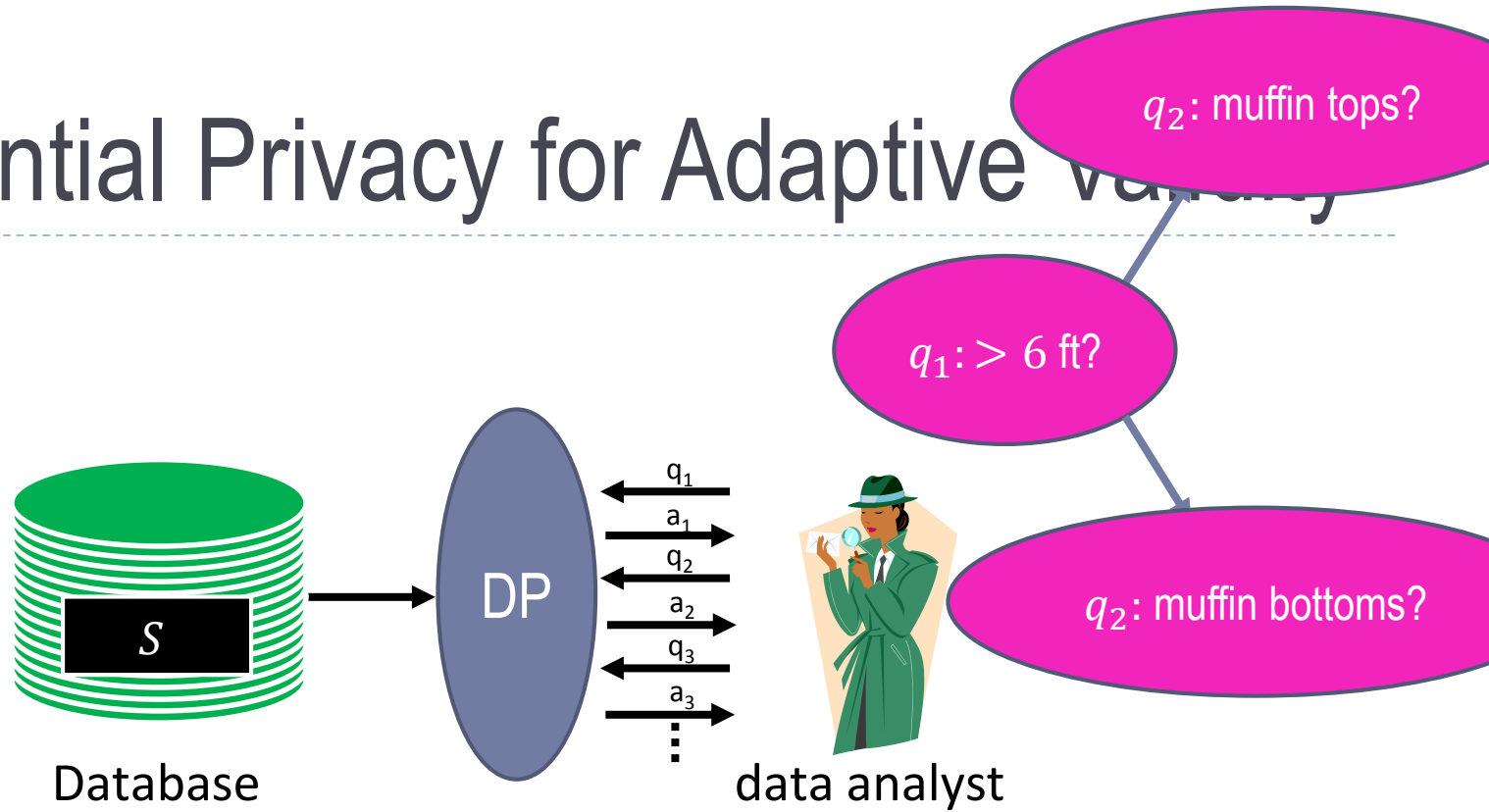
Cynthia Dwork, Microsoft Research

# Adaptive Data Analysis



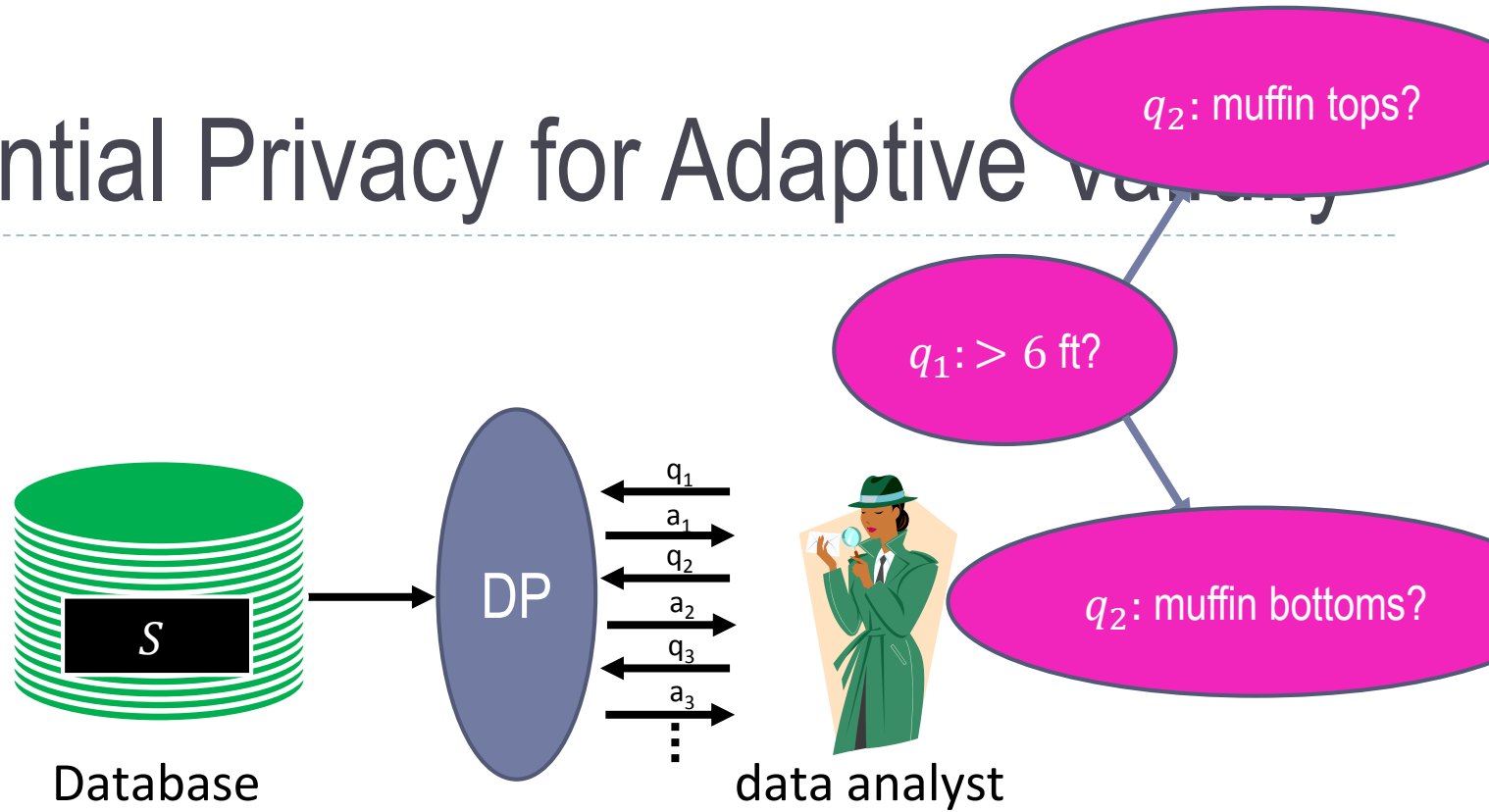
- ▶  $q_i$  depends on  $a_1, a_2, \dots, a_{i-1}$
- ▶ Worry: analyst finds a query for which the dataset is not representative of population; reports surprising discovery

# Differential Privacy for Adaptive Querying



- ▶  $q_i$  depends on  $a_1, a_2, \dots, a_{i-1}$
- ▶ **Differential privacy neutralizes risks incurred by adaptivity**
  - ▶ Definition of privacy tailored to statistical analysis of large data sets

# Differential Privacy for Adaptive Querying



- ▶  $q_i$  depends on  $a_1, a_2, \dots, a_{i-1}$
- ▶ **Differential privacy neutralizes risks incurred by adaptivity**
  - ▶  $\exists$  LARGE literature on DP algorithms for data analysis

# Some Intuition

---

- ▶ Fix a query, eg, “What fraction of population is over 6 feet tall?”
- ▶ Almost all large datasets will give an approximately correct reply
  - ▶ Most datasets are representative with respect to this query
- ▶ If, in the process of adaptive exploration, the analyst finds a query for which the dataset is not representative, then she must have “learned something significant” about the dataset.
  - ▶ Preserving the “privacy” of the data may prevent over-fitting.

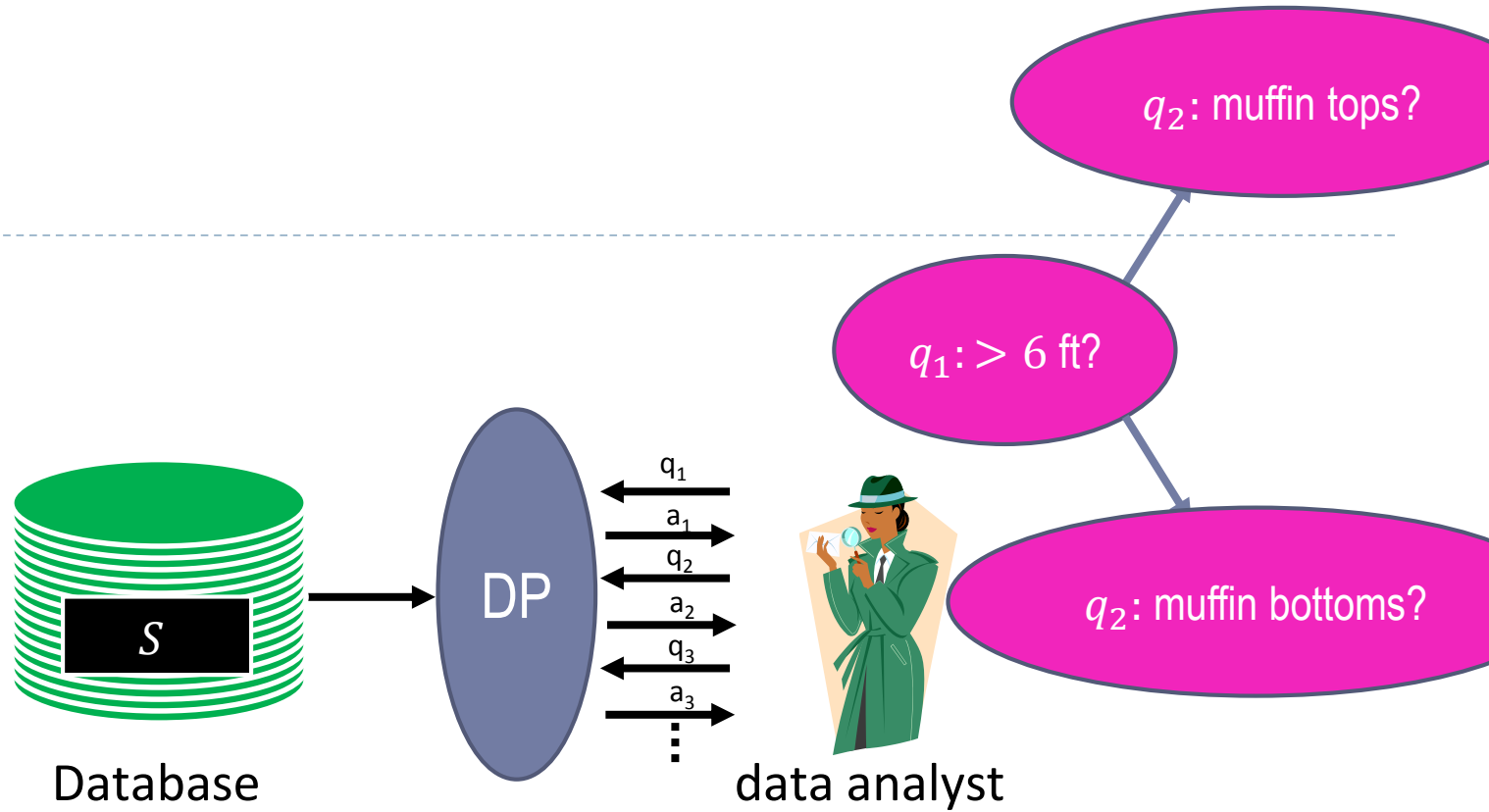


# Intuition After Nati's Talk

---

- ▶ Differential Privacy: The outcome of any analysis is essentially equally likely, independent of whether any individual joins, or refrains from joining, the dataset.
  - ▶ This is a stability requirement.
  - ▶ Gave rise to the folklore that differential privacy yields generalizability.
  - ▶ But we will be able to say something stronger.





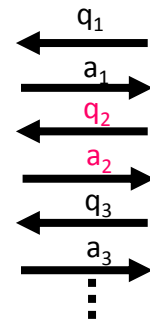
- ▶  $q_i$  depends on  $a_1, a_2, \dots, a_{i-1}$
- ▶ **Differential privacy neutralizes risks incurred by adaptivity**
  - ▶ E.g., for statistical queries: whp  $|E_S[A(S)] - E_P[A(S)]| < \tau$
  - ▶ High probability is important for handling many queries

# Formalization

- ▶ Data sets  $S \in X^n; S \sim D$
- ▶ Queries  $q: X^n \rightarrow Y$
- ▶ Algorithms that choose queries and output results
  - ▶  $A_1 = q_1$  (trivial choice), outputs  $(q_1, q_1(S))$
  - ▶  $A_i: X^n \times Y_1 \times \dots \times Y_{i-1} \rightarrow Y_i$  where
    - ▶  $q_i = C_i(y_1, \dots, y_{i-1})$
    - ▶  $A_i(S, y_1, \dots, y_{i-1}) = (q_i, q_i(S)) = (q_i, a_i)$
- ▶  $H \stackrel{\text{def}}{=} \{(S, q) \mid q(S) \text{ not representative wrt } D\}$ 
  - ▶  $\forall (y_1, \dots, y_{i-1}) \Pr_S [(S, q_i) \in H] \leq \beta_i$

Choose new query based on history of observations

Output chosen query and its response on  $S$



- ▶ We want:  $\Pr[(S, C_i(S)) \in H]$  to be similar
  - ▶  $q_i(S)$  should generalize even when  $q_i$  chosen as a function of  $S$

$q_i(S)$  fails to generalize



# Differential Privacy [D.,McSherry,Nissim,Smith '06]

---

$M$  gives  $\epsilon$ -differential privacy if for all pairs of adjacent data sets  $S, S'$ , and all events  $T$

$$\Pr[M(S) \in T] \leq e^\epsilon \Pr[M(S') \in T]$$

Randomness introduced by  $M$



# Differential Privacy [D.,McSherry,Nissim,Smith '06]

---

$M$  gives  $\epsilon$ -differential privacy if for all pairs of adjacent data sets  $S, S'$ , and all events  $T$

$$\Pr[M(S) \in T] \leq e^\epsilon \Pr[M(S') \in T]$$

For random variables  $\mathbf{X}, \mathbf{Y}$  over  $X$ , the max-divergence of  $\mathbf{X}$  from  $\mathbf{Y}$  is given by

$$D_\infty(\mathbf{X}||\mathbf{Y}) = \log \max_{x \in X} \frac{\Pr[\mathbf{X} = x]}{\Pr[\mathbf{Y} = x]}$$

Then  $\epsilon$ -DP equivalent to  $D_\infty(M(S)||M(S')) \leq \epsilon$ .

---



# Differential Privacy [D.,McSherry,Nissim,Smith '06]

---

$M$  gives  $\epsilon$ -differential privacy if for all pairs of adjacent data sets  $S, S'$ , and all events  $T$

$$\Pr[M(S) \in T] \leq e^\epsilon \Pr[M(S') \in T]$$

For random variables  $\mathbf{X}, \mathbf{Y}$  over  $X$ , the max-divergence of  $\mathbf{X}$  from  $\mathbf{Y}$  is given by

$$D_\infty(\mathbf{X}||\mathbf{Y}) = \log \max_{x \in X} \frac{\Pr[\mathbf{X} = x]}{\Pr[\mathbf{Y} = x]}$$

Then  $\epsilon$ -DP equivalent to  $D_\infty(M(S)||M(S')) \leq \epsilon$ .

Closed under post-processing:  $D_\infty(A(M(S))||A(M(S'))) \leq \epsilon$ .

---



# Differential Privacy [D.,McSherry,Nissim,Smith '06]

---

$M$  gives  $\epsilon$ -differential privacy if for all pairs of adjacent data sets  $S, S'$ , and all events  $T$

$$\Pr[M(S) \in T] \leq e^\epsilon \Pr[M(S') \in T]$$

For random variables  $\mathbf{X}, \mathbf{Y}$  over  $X$ , the max-divergence of  $\mathbf{X}$  from  $\mathbf{Y}$  is given by

$$D_\infty(\mathbf{X}||\mathbf{Y}) = \log \max_{x \in X} \frac{\Pr[\mathbf{X} = x]}{\Pr[\mathbf{Y} = x]}$$

Then  $\epsilon$ -DP equivalent to  $D_\infty(M(S)||M(S')) \leq \epsilon$ .

Group Privacy:  $\forall S, S'' D_\infty(M(S)||M(S')) \leq \Delta(S, S'')\epsilon$ .

---



# Properties

---

- ▶ **Closed under post-processing**
  - ▶ Max-divergence remains bounded
- ▶ **Automatically yields group privacy**
  - ▶  $k\epsilon$  for groups of size  $k$
- ▶ **Understand behavior under adaptive composition**
  - ▶ Can bound cumulative privacy loss over multiple analyses
    - ▶ “The epsilons add up”
- ▶ **Programmable**
  - ▶ Complicated private analyses from simple private building blocks

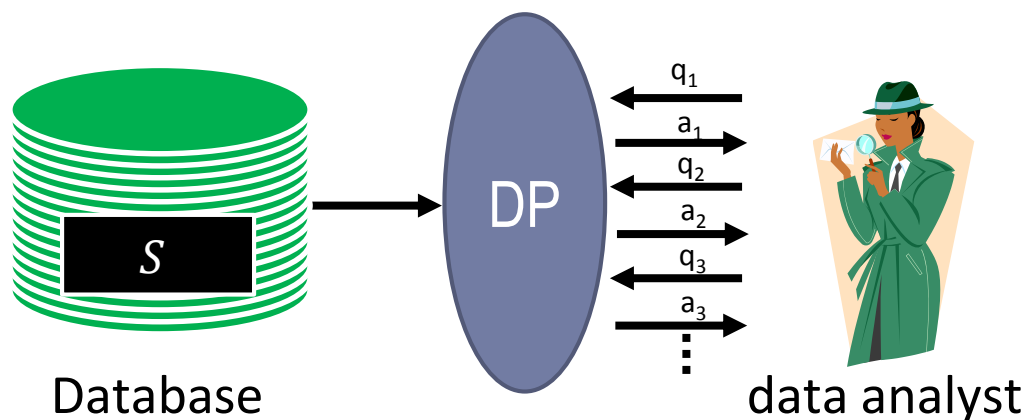


# The Power of Composition

- ▶ Lemma: The choice of  $q_i$  is differentially private.
  - ▶ Closure under post-processing.
- ▶ Inductive step (key): If  $q$  is chosen in a differentially private fashion with respect to  $S$ , then

$\Pr[(S, C(S)) \in H] \text{ is small}$

- ▶ Sufficiency: union bound.



# Description Length

---

▶ Let  $A: X^n \rightarrow Y$ .

▶ Description length of  $A$  is the cardinality of its range

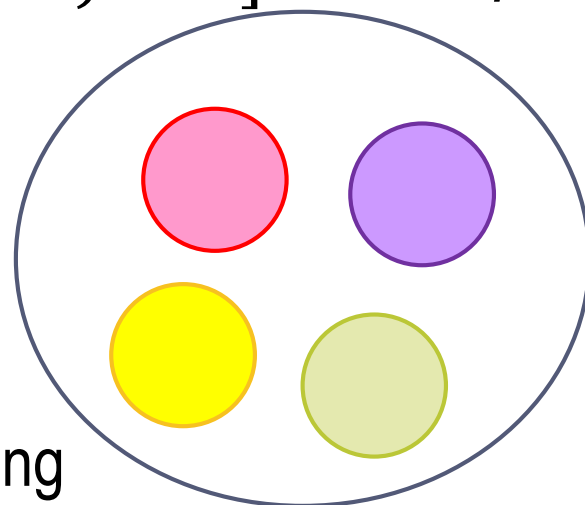
If  $\forall y \Pr_S [(S, y) \in H] \leq \beta$ , then  $\Pr_S [(S, A(S)) \in H] \leq |Y| \cdot \beta$

▶ Description length composes too.

▶ Product:  $\beta \cdot \prod_i |Y_i|$

▶ And, morally, it is closed under post-processing

▶ Once you fix the randomness of the post-processing algorithm



# Approximate max-divergence

---

$\beta$ -approximate max-divergence of  $\mathbf{X}$  from  $\mathbf{Y}$

$$D_{\infty}^{\beta}(\mathbf{X}||\mathbf{Y}) = \log \max_{T \in \mathcal{X}, \Pr[\mathbf{X} \in T] > \beta} \frac{\Pr[\mathbf{X} \in T] - \beta}{\Pr[\mathbf{Y} \in T]}$$





# Approximate max-divergence

---

$\beta$ -approximate max-divergence of  $\mathbf{X}$  from  $\mathbf{Y}$

$$D_{\infty}^{\beta}(\mathbf{X}||\mathbf{Y}) = \log \max_{T \in \mathcal{X}, \Pr[\mathbf{X} \in T] > \beta} \frac{\Pr[\mathbf{X} \in T] - \beta}{\Pr[\mathbf{Y} \in T]}$$

We are interested in (with  $\beta$ , but too messy)

$$D_{\infty}((\mathbf{S}, A(\mathbf{S}))||\mathbf{S} \times A(\mathbf{S})) = \log \max_T \frac{\Pr[(\mathbf{S}, A(\mathbf{S})) \in T]}{\Pr[\mathbf{S} \times A(\mathbf{S}) \in T]}$$



# Approximate max-divergence

---

$\beta$ -approximate max-divergence of  $\mathbf{X}$  from  $\mathbf{Y}$

$$D_{\infty}^{\beta}(\mathbf{X}||\mathbf{Y}) = \log \max_{T \in \mathcal{X}, \Pr[\mathbf{X} \in T] > \beta} \frac{\Pr[\mathbf{X} \in T] - \beta}{\Pr[\mathbf{Y} \in T]}$$

We are interested in (with  $\beta$ , but too messy)

$$D_{\infty}((\mathbf{S}, A(\mathbf{S}))||\mathbf{S} \times A(\mathbf{S})) = \log \max_T \frac{\Pr[(\mathbf{S}, A(\mathbf{S})) \in T]}{\Pr[\mathbf{S} \times A(\mathbf{S}) \in T]}$$

How much more likely is  $A(S)$  to relate to  $S$  than to a fresh  $S'$ ?

Captures the maximum amount of information that an output of an algorithm might reveal about its input

---



# Unifying Concept: Max-Information

---

- ▶  $I_{\infty}^{\beta}(\mathbf{X}; \mathbf{Y}) = D_{\infty}^{\beta}((\mathbf{X}, \mathbf{Y}) || \mathbf{X} \times \mathbf{Y})$
- ▶ We are interested in  $I_{\infty}^{\beta}(\mathbf{S}; A(\mathbf{S}))$
- ▶ Theorem: If  $I_{\infty}^{\beta}(\mathbf{S}; A(\mathbf{S})) \leq k$  then for any  $T \subseteq X^n \times Y$ 
  - ▶  $\Pr[(\mathbf{S}, A(\mathbf{S})) \in T] \leq 2^k \Pr[\mathbf{S} \times A(\mathbf{S}) \in T] + \beta$
  - ▶ So  $\Pr[(\mathbf{S}, A(\mathbf{S})) \in H] \leq 2^k \max_{y \in Y} \Pr[(\mathbf{S}, y) \in H] + \beta !$

# Unifying Concept: Max-Information

- ▶  $I_{\infty}^{\beta}(\mathbf{X}; \mathbf{Y}) = D_{\infty}^{\beta}((\mathbf{X}, \mathbf{Y}) || \mathbf{X} \times \mathbf{Y})$
- ▶ We are interested in  $I_{\infty}^{\beta}(\mathbf{S}; A(\mathbf{S}))$
- ▶ Theorem: If  $I_{\infty}^{\beta}(\mathbf{S}; A(\mathbf{S})) \leq k$  then for any  $T \subseteq X^n \times Y$ 
  - ▶  $\Pr[(\mathbf{S}, A(\mathbf{S})) \in T] \leq 2^k \Pr[\mathbf{S} \times A(\mathbf{S}) \in T] + \beta$
  - ▶ So  $\Pr[(\mathbf{S}, A(\mathbf{S})) \in H] \leq 2^k \max_{y \in Y} \Pr[(\mathbf{S}, y) \in H] + \beta !$
- ▶ Max-Information composes and is closed under post-processing
- ▶ For  $\epsilon$ -DP  $A$ :  $I_{\infty}(A, n) \leq \epsilon n \log_2 e$ . Better bounds for  $I_{\infty}^{\beta}(A, n)$ .
- ▶  $I_{\infty}^{\beta}(A, n) \leq \log \left( \frac{|Y|}{\beta} \right)$

Bound on worst case approximate max info for any distribution on  $n$ -element databases

# Abstract is Good

---

- ▶ Focusing on *properties* is powerful
  - ▶ Completely universal approach to validity of adaptive analysis
    - ▶ DP, small description length, low max-information
  - ▶ Large numbers of arbitrary adaptively chosen computations
    - ▶ Closure under post-processing and composition



# Long Live the Dataset!

---

- ▶ Leaking information slowly prolongs the lifetime of the system
- ▶ Similar to the situation with privacy for the sake of privacy
  - ▶ To avoid too much cumulative loss, answer with smaller values of  $\epsilon$
  - ▶ Essential: Fundamental Law of Information Leakage
    - ▶ Overly accurate estimates of too many statistics is blatantly non-private.
    - ▶ Dealer's choice
- ▶ **Conjecture: The same is true for adaptivity.**
  - ▶ **Failure to control cumulative max-info leads to failure to generalize**
  - ▶ **Important policy Implications!**
  - ▶ Supporting evidence: Hardt-Ullman queries





# Thank you!



NIPS Workshop on Adaptive Data Analysis, Montreal, 12/11/15