

# Machine Learning and Differential Privacy

Ellen Vitercik

December 5, 2018

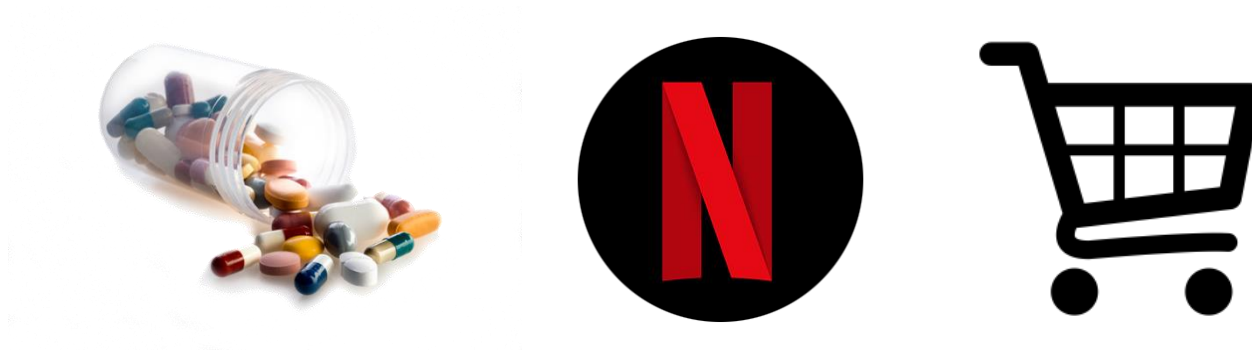
# Today we'll talk about...

1. The importance of privacy in machine learning
2. One way of defining privacy (*differential privacy*)
3. Tools for designing privacy-preserving algorithms

# Learning and privacy

To do machine learning, we need data

What if the data contains sensitive information?





Is it enough to trust the person running the learning algorithm?

**No:** Perhaps algorithm's output reveals sensitive information

# Example: search query completions

why are



why are **manhole covers** round

why are **you** interested in this position

why are **firetrucks** red

why are **cells** so small

why are **cats** afraid of cucumbers

why are **flamingos** pink

why are **we** here

# Example: search query completions

What if we use your friends' search logs to suggest completions?

Might be good for accuracy, but...

why are _	 
why are <b>my feet so itchy?</b>	

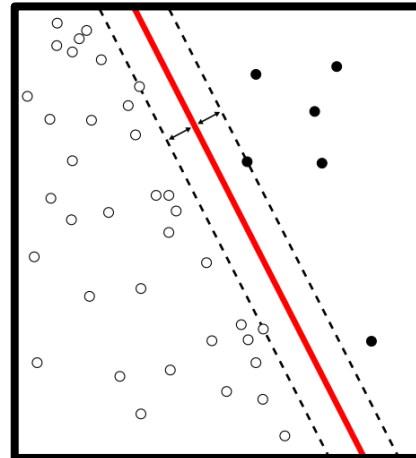


# Privacy leaks can be subtle!

Hospital wants to be able to predict who has condition  $X$

Collect data from residents, use perceptron algorithm

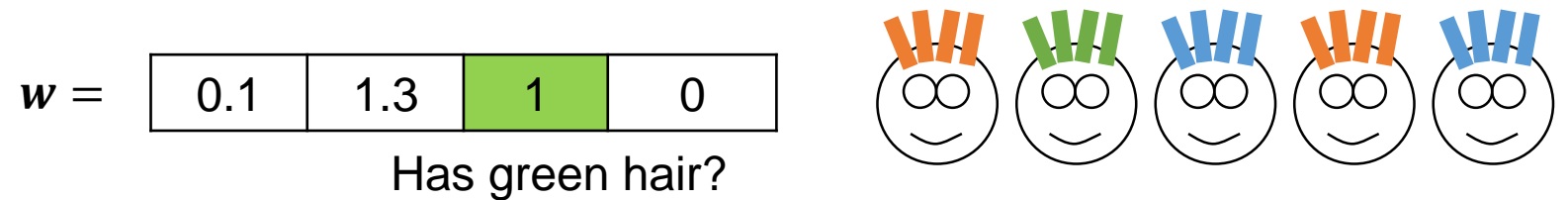
$$w = \begin{array}{|c|c|c|c|} \hline 0.1 & 1.3 & 1 & 0 \\ \hline \end{array}$$



# Privacy leaks can be subtle!

Hospital wants to be able to predict who has condition  $X$

Collect data from residents, use perceptron algorithm



Only one person in town has green hair.

We now know the green-haired person has condition  $X$ !

**How can we be confident that this won't happen?**

# Today we'll talk about...

1. The importance of privacy in machine learning
2. One way of defining privacy (*differential privacy*)
3. Tools for designing privacy-preserving algorithms



# What is privacy?



# What isn't privacy?

Privacy isn't **restricting questions to large populations.**

- “What is the average salary of CMU faculty?”
- “What is the average salary of CMU faculty not named Nina Balcan?”



# What isn't privacy?

Privacy isn't **restricting to “ordinary” facts**

Statistics on Alice's bread buying habits:

*For 20 years she regularly buys bread, then stops.*

Type 2 diabetes?



# What isn't privacy?

Privacy isn't “**Anonymization**”

Case study: Publicly available “anonymized” hospitalization data

Latanya Sweeney re-identified patients by name

---

**The New York Times**

---

Bits

Business, Innovation, Technology, Society

PRIVACY

## With a Few Bits of Data, Researchers Identify ‘Anonymous’ People

By Natasha Singer January 29, 2015 2:01 pm

# What is privacy?

## Attempt 1:

Analysis of dataset D is private if:

Analyst knows no more about Alice after analysis than before.

## Problematic example:

Analysis of dataset D  $\Rightarrow$  West Virginians have high obesity rates

Alice, whose information **isn't** in dataset D, lives in WV

Insurance agency knows Alice lives in WV  $\Rightarrow$  they raise her rates!

Was Alice's privacy violated?

Yes, under this definition...



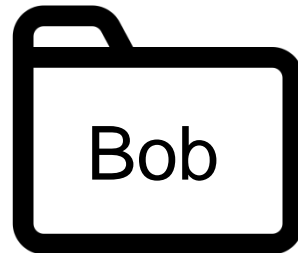
# What is privacy?

## Attempt 2:

Analysis of dataset D is private if:

analyst knows **almost** no more about Alice after analysis  
than he **would have**

had he conducted the **same analysis** on  
an **identical** dataset w/ **Alice's data removed**



# Differential privacy



“Calibrating Noise to Sensitivity in Private Data Analysis.” Dwork, McSherry, Nissim, and Smith. *TCC*. 2006.

“The Algorithmic Foundations of Differential Privacy”. Dwork and Roth. *Foundations and Trends in Theoretical Computer Science*, NOW Publishers. 2014.

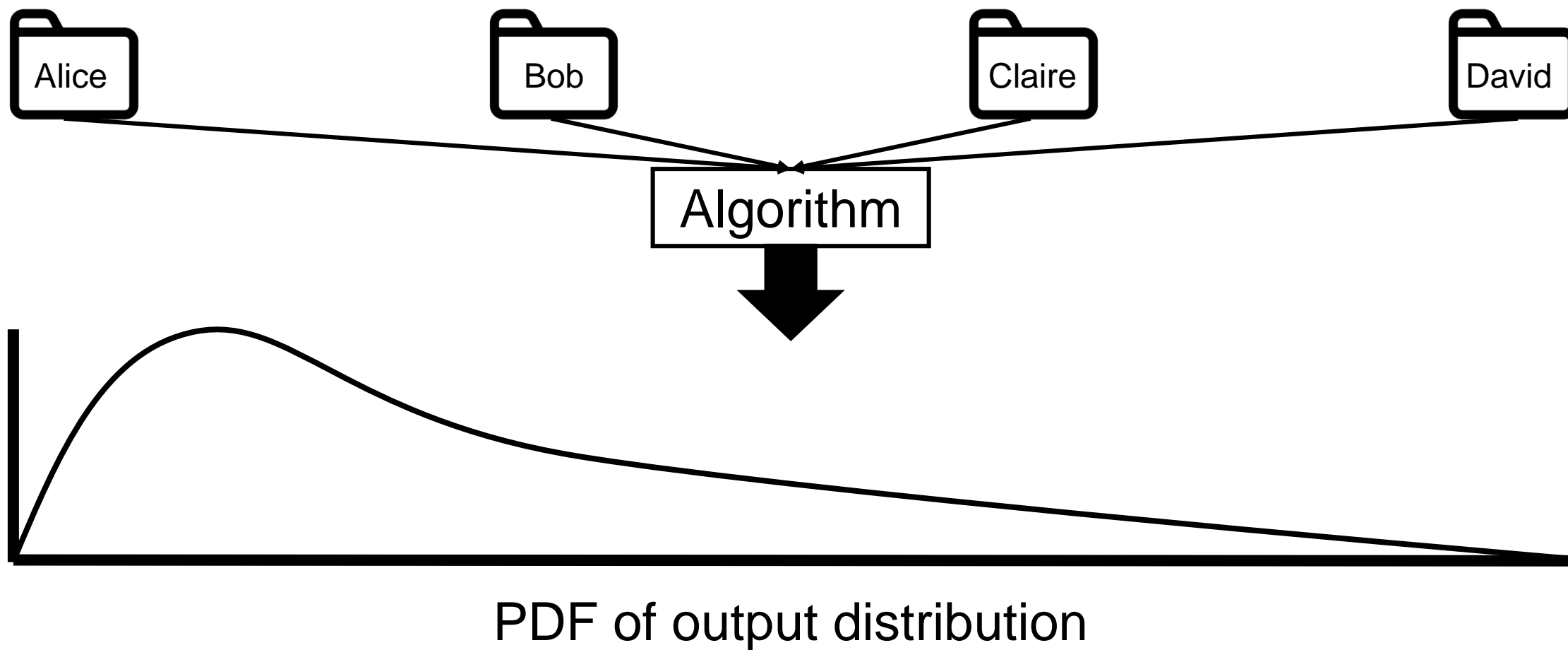


Differential privacy

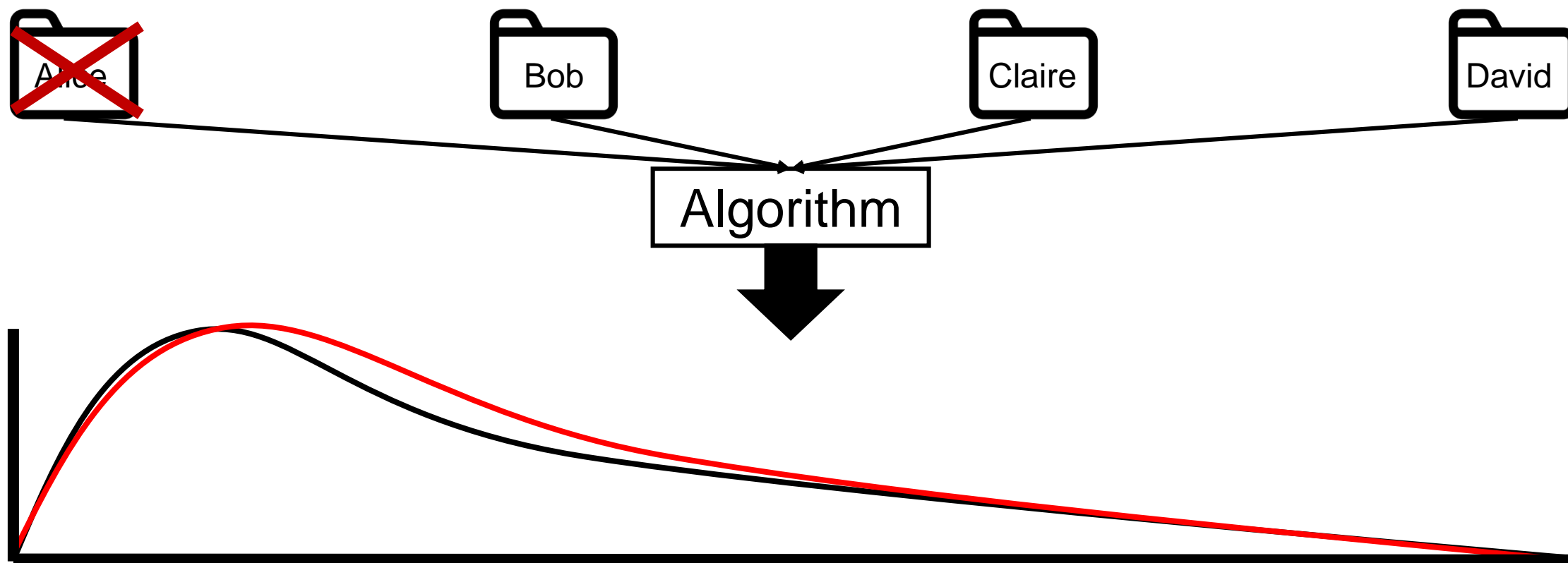




# Differential privacy



# Differential privacy



Can't tell if used Alice's data, let alone what her data was!

# Differential privacy

**Def:** Two datasets  $S, S'$  are **neighboring** if differ on  $\leq 1$  entry  
*1 entry  $\equiv$  1 person*

$S$	$S'$
$x_1$	$x_1$
$\vdots$	$\vdots$
$x_i$	$x'_i$
$\vdots$	$\vdots$
$x_n$	$x_n$

# Differential privacy

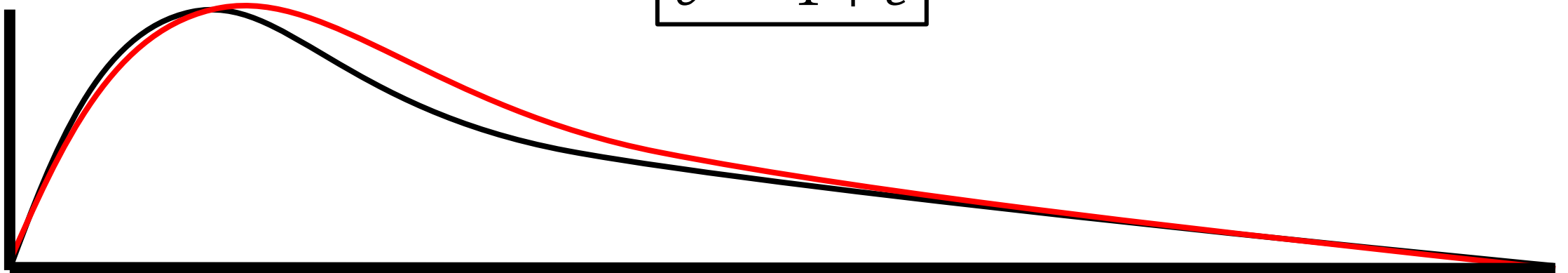
Algorithm  $\mathcal{A}$  is  **$\epsilon$ -differentially private** if:

For all pairs of neighboring sets  $S, S'$  and all sets  $\mathcal{O}$  of outputs,

$$\mathbb{P}[\mathcal{A}(S) \in \mathcal{O}] \leq e^\epsilon \mathbb{P}[\mathcal{A}(S') \in \mathcal{O}]$$

$e^\epsilon \approx 1 + \epsilon$





# DP protects against additional harm

$\mathcal{A} :=$  DP algorithm

$f: \text{Range}(\mathcal{A}) \rightarrow W$  maps  $\mathcal{A}$ 's output to a future world state  $w \in W$

Suppose I have a utility function  $u: W \rightarrow \mathbb{R}$

E.g.,  $u(w)$  = “how happy am I if the world is  $w$ ”



DP guarantees that

$$\mathbb{E}_{w \sim f(\mathcal{A}(s))}[u(w)] \approx e^{\pm \epsilon} \cdot \mathbb{E}_{w \sim f(\mathcal{A}(s'))}[u(w)]$$

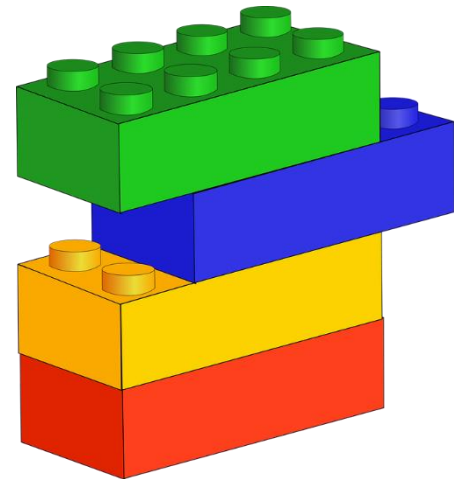
# Today we'll talk about...

1. The importance of privacy in machine learning
2. One way of defining privacy (*differential privacy*)
3. Tools for designing privacy-preserving algorithms
  - a) **Laplace mechanism**
  - b) Exponential mechanism
  - c) Composing private algorithms
  - d) Examples of differentially-private ML tools

# Laplace mechanism

Very useful building block for designing private algorithms.

“Calibrating Noise to Sensitivity in Private Data Analysis.” Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. TCC, 2006.



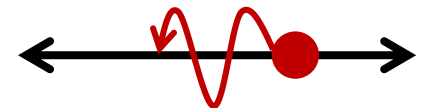
# Laplace mechanism

**Goal:** Evaluate  $f: D \rightarrow \mathbb{R}$  mapping datasets to  $\mathbb{R}$ ; preserve  $\epsilon$ -DP

*Ex.,  $f(S) := \text{mean weight of people in } S$*

**Idea:** Compute  $f(S)$  and add noise to hide any individual's info

How little can we get away with?





# Laplace mechanism

**Goal:** Evaluate  $f: D \rightarrow \mathbb{R}$  mapping datasets to  $\mathbb{R}$ ; preserve  $\epsilon$ -DP

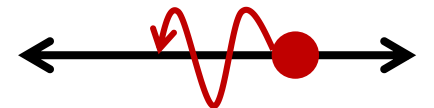
*Ex.,  $f(S) := \text{mean weight of people in } S$*

**Idea:** Compute  $f(S)$  and add noise to hide any individual's info

**Def: Sensitivity** of  $f$  is  $\Delta_f = \max_{S, S' \text{ neighboring}} |f(S) - f(S')|$

**Laplace Mechanism** outputs  $Z_S \sim \text{Lap}\left(f(S), \frac{\Delta_f}{\epsilon}\right)$

$$\text{PDF } p_{Z_S}(z) = \frac{\Delta_f}{2\epsilon} \exp\left(-\frac{\epsilon}{\Delta_f} |z - f(S)|\right)$$



# Laplace mechanism: Privacy guarantees

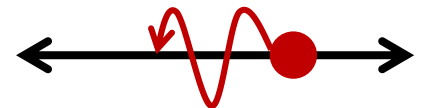
**Def: Sensitivity** of  $f$  is  $\Delta_f = \max_{S, S' \text{ neighboring}} |f(S) - f(S')|$

**Laplace Mechanism** outputs  $Z_S \sim \text{Lap}\left(f(S), \frac{\Delta_f}{\epsilon}\right)$

$$\text{PDF } p_{Z_S}(z) = \frac{\Delta_f}{2\epsilon} \exp\left(-\frac{\epsilon}{\Delta_f} |z - f(S)|\right)$$

**Privacy:** The Laplace mechanism preserves  $\epsilon$ -DP.

*We'll see why on the board.*



# Laplace mechanism: Utility guarantees

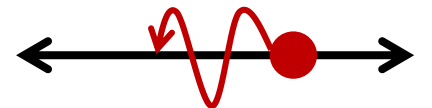
**Def: Sensitivity** of  $f$  is  $\Delta_f = \max_{S, S' \text{ neighboring}} |f(S) - f(S')|$

**Laplace Mechanism** outputs  $Z_S \sim \text{Lap}\left(f(S), \frac{\Delta_f}{\epsilon}\right)$

$$\text{PDF } p_{Z_S}(z) = \frac{\Delta_f}{2\epsilon} \exp\left(-\frac{\epsilon}{\Delta_f} |z - f(S)|\right)$$

**Utility:** With probability at least  $1 - \delta$ ,  $|Z_S - f(S)| \leq \frac{\Delta_f}{\epsilon} \log \frac{1}{\delta}$ .

*Proof idea:* analyze Laplace distribution's CDF.



# Laplace mechanism: Computing means

Given set  $S = \{x_1, \dots, x_n\} \subset [0,1]$ , privately compute  $f(S) = \frac{1}{n} \sum x_i$

**Question:** What is  $\Delta_f = \max_{S, S' \text{ neighboring}} |f(S) - f(S')|$ ?

**Answer:**  $\Delta_f = \frac{1}{n}$

# Laplace mechanism: Computing means

Given set  $S = \{x_1, \dots, x_n\} \subset [0,1]$ , privately compute  $f(S) = \frac{1}{n} \sum x_i$

**Recall:** Laplace mech. outputs  $Z_S \sim \text{Lap}\left(f(S), \frac{1}{n\epsilon}\right)$

# Laplace mechanism: Computing means

Given set  $S = \{x_1, \dots, x_n\} \subset [0,1]$ , privately compute  $f(S) = \frac{1}{n} \sum x_i$

**Utility:** With probability at least  $1 - \delta$ ,  $|Z_S - f(S)| \leq \frac{1}{n\epsilon} \log \frac{1}{\delta}$ .

If  $S \sim P^n$  and goal is to estimate  $\mathbb{E}_{x \sim P}[x]$  using  $f(S)$ , w.p.  $1 - \delta$ ,

$$|\mathbb{E}_{x \sim P}[x] - f(S)| \leq \sqrt{\frac{1}{2n} \ln \frac{1}{\delta}}.$$

**Error due to privacy negligible compared to sampling error!**

# Laplace mechanism: Multi-dim functions

What if function  $f$  maps to  $\mathbb{R}^d$ ? I.e.,  $f: D \rightarrow \mathbb{R}^d$

Example:  $f(S) = \langle \text{mean weight in } S, \text{mean height in } S \rangle$

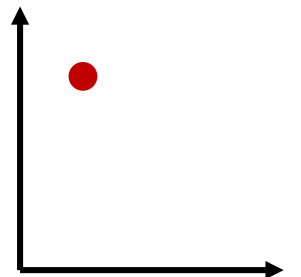
**Def:** The sensitivity of  $f$  is  $\Delta_f = \max_{S, S' \text{ neighboring}} \|f(S) - f(S')\|_1$ .

**Def:** The Laplace Mechanism outputs  $f(S) + \mathbf{Z}$

$\mathbf{Z} \in \mathbb{R}^d$  has components drawn from  $\text{Lap}\left(0, \frac{\Delta_f}{\epsilon}\right)$  distribution

**Privacy:** The Laplace mechanism preserves  $\epsilon$ -DP

**Utility:** With probability at least  $1 - \delta$ ,  $\|\mathbf{Z}\|_\infty \leq \frac{\Delta_f}{\epsilon} \log \frac{d}{\delta}$



# Today we'll talk about...

1. The importance of privacy in machine learning
2. One way of defining privacy (*differential privacy*)
3. Tools for designing privacy-preserving algorithms
  - a) Laplace mechanism
  - b) Exponential mechanism**
  - c) Examples of differentially-private ML tools



# Exponential mechanism

**Goal:** Choose the “best” item from a finite set  $Y$  of items

*E.g., voting in a local election*



Frank McSherry and Kunal Talwar. Mechanism design via differential privacy.  
In Foundations of Computer Science. 2007.

# Exponential mechanism

Given utility function  $u(S, y) = \text{“utility of } y \text{ for dataset } S\text{”}$

**Goal:** Find  $y \in Y$  maximizing  $u(S, y)$

**Question:** Why can't we use the Laplace Mechanism?



**Answer:** E.g.,  $Y = \{\text{town welder, town farmer, ...}\}$

How do we add noise to “town mechanic”?

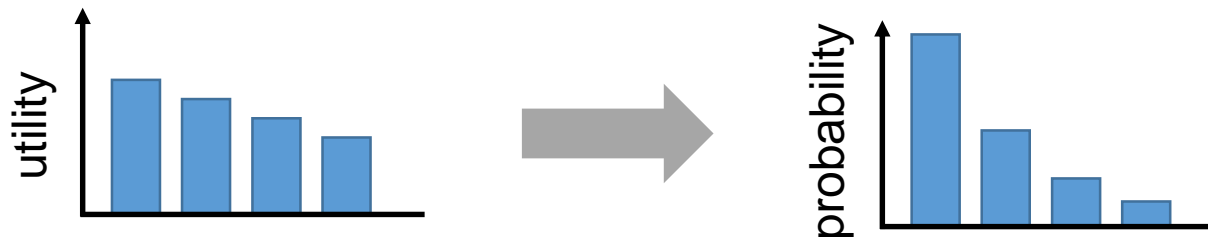
# Exponential mechanism

Given utility function  $u(S, y)$  = “utility of  $y$  for dataset  $S$ ”

**Goal:** Find  $y \in Y$  maximizing  $u(S, y)$

**Def:** The sensitivity of  $u$  is  $\Delta_u = \max_{S, S', y} |u(S, y) - u(S', y)|$

Exponential Mechanism outputs  $y$  with w.p.  $\propto \exp\left(\frac{\epsilon}{2\Delta_u} u(S, y)\right)$



# Exponential mechanism

Given utility function  $u(S, y)$  = “utility of  $y$  for dataset  $S$ ”

**Goal:** Find  $y \in Y$  maximizing  $u(S, y)$

**Def:** The sensitivity of  $u$  is  $\Delta_u = \max_{S, S', y} |u(S, y) - u(S', y)|$

Exponential Mechanism outputs  $y$  with w.p.  $\propto \exp\left(\frac{\epsilon}{2\Delta_u} u(S, y)\right)$

**Privacy:** The exponential mechanism preserves  $\epsilon$ -DP.

*Proof follows from algebraic manipulations of density function.*

# Exponential mechanism

Given utility function  $u(S, y)$  = “utility of  $y$  for dataset  $S$ ”

**Goal:** Find  $y \in Y$  maximizing  $u(S, y)$

**Def:** The sensitivity of  $u$  is  $\Delta_u = \max_{S, S', y} |u(S, y) - u(S', y)|$

Exponential Mechanism outputs  $y$  with w.p.  $\propto \exp\left(\frac{\epsilon}{2\Delta_u} u(S, y)\right)$

**Privacy:** The exponential mechanism preserves  $\epsilon$ -DP.

*We'll see why on the board.*

# Database sanitization



Given dataset  $S$ , produce synthetic dataset  $\hat{S}$ , preserve DP

**Ideally:**  $\hat{S}$  behaves basically the same as  $S$  (for our purposes)

Based on "A learning theory approach to noninteractive database privacy." Avrim Blum, Katrina Ligett, Aaron Roth. *Journal of the ACM (JACM)* 60.2 (2013): 12.

# Database sanitization



## More formally:

- Let  $S \subseteq \{0,1\}^d$  be a dataset of  $d$ -dimensional binary vectors
- Let  $H$  be a set of functions  $h: \{0,1\}^d \rightarrow \{0,1\}$  with VC-dim  $D$
- Let  $h(S) = \frac{1}{|S|} \sum_{x \in S} h(x)$  be the fraction of  $x \in S$  with  $h(x) = 1$

**If  $|S| \geq \tilde{O}\left(\frac{dD}{\alpha^3 \epsilon}\right)$ , can find  $\hat{S} \subset \{0,1\}^d$  while preserving  $\epsilon$ -DP s.t. w.h.p., for all  $h \in H$ ,  $|h(S) - h(\hat{S})| \leq \alpha$ .**

*Proof uses VC dimension guarantees and probabilistic method*

# Today we'll talk about...

1. The importance of privacy in machine learning
2. One way of defining privacy (*differential privacy*)
3. Tools for designing privacy-preserving algorithms
  - a) Laplace mechanism
  - b) Exponential mechanism
  - c) **Examples of differentially-private ML tools**

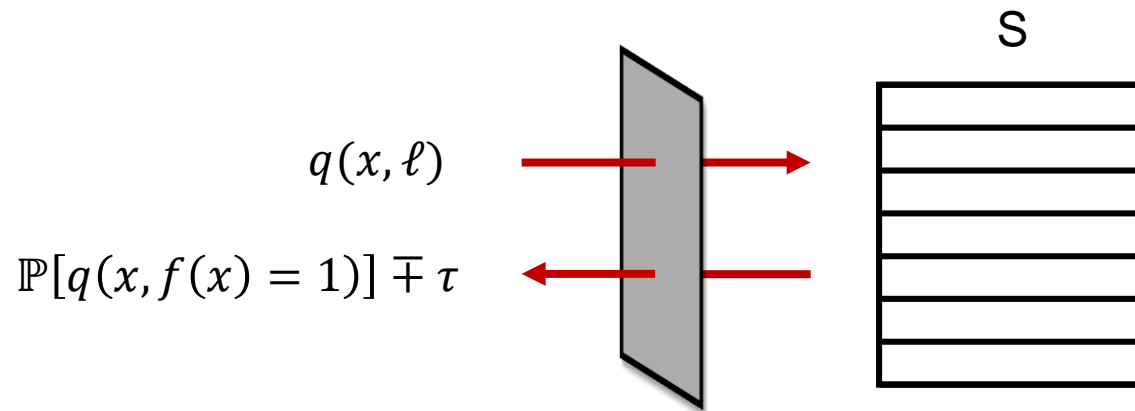


# DP + ML using statistical queries

Anything learnable using *statistical queries* is privately learnable.

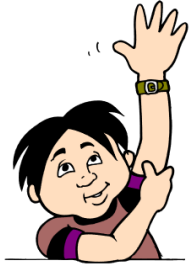
Practical Privacy: The SuLQ Framework. Blum, Dwork, McSherry, Nissim. *PODS*. 2005.

**Statistical query model** [Kearns, '98]:



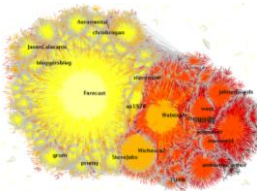
Many algorithms (e.g., ID3, Perceptron, SVM, PCA) can be re-written to interface via statistical queries.

# DP+ML more generally



## Active learning

Balcan and Feldman. "Statistical active learning algorithms." *NeurIPS*. 2013.



## Clustering

Balcan, Dick, Liang, Mou, and Zhang. "Differentially Private Clustering in High-Dimensional Euclidean Spaces." *ICML*. 2017.



## Distributed learning

Blacan, Blum, Fine, and Mansour. "Distributed Learning, Communication Complexity and Privacy." *COLT*. 2012.