# Pinning Down "Privacy" in Statistical Databases

#### Adam Smith

Computer Science & Engineering Department Penn State

1



Large collections of personal information

- census data
- medical/public health data
- social networks
- recommendation systems
- trace data: search records, etc
- intrusion-detection systems



- larger data sets
- more types of data



- Published "statistics" may be tables, graphs, microdata, decision trees, neural networks, confidence intervals...
- Data may be numbers, categories, tax forms, web searches...
- May be interactive

Drivo ov in	St	U.S. Census Bureau				
FILVACY III		American FactFinder	Main	Search	Feedback	FAQs
		Select Geography				
Individuals	Sei	You are here: Main > Data Sets > Data Sets with Deta Census 2000 Summary File 1 (S	ailed Table: F 1) 100-P	<mark>s ▶ Geogra</mark> j ercent Data,	phy ► Tables ► I Detailed Tables	Results
mannadalo		Choose a selection method				
$x_1$		list name search address search	map	geo wit	hin geo	
$x_2$		Show all geography types   🕧 Explain Census	Geography	۷		
X :	$\rightarrow$	Enter a street address, city and state, or	a street a	address an	d ZIP code. C	lick 'Go'
r		Street Address Quick tips				
$\omega n$		36 leyden st	-			
X		Medford Massachusetts		02155	Go	
	Select one or more geographic areas and click 'Add'					
		Census Tract: Census Tract 3397 Block Group: Block Group 2				
		L., Place: Medford city				

- Published "statistics" may be tables, graphs, microdata, decision trees, neural networks, confidence intervals...
- Data may be numbers, categories, tax forms, web searches...
- May be interactive



- What information can be released?
- Two conflicting goals
  - Utility: Users can extract "global" properties
  - Privacy ("confidentiality"): Individual information stays hidden
- How can these be formalized?



- Variations on model studied in
  - Statistics ("statistical disclosure control")
  - Data mining ("privacy-preserving data mining" \*)
- No coherent theory
- Recently: crypto & theoretical CS
  - Focused on rigorous approach to privacy



- "Privacy" is harder to reason about than "utility"
   Utility is what we're used to
- Existing definitions problematic
  - Many are not specified precisely
  - ➢ Fail in the presence of external information

### **External Information**



## **External Information**



• Users have external information sources

Can't assume we know the sources

Anonymization schemes regularly broken

## **External Information**



- Users have external information sources
  - Can't assume we know the sources
- Anonymization schemes regularly broken
- Example: two hospitals independently release statistics about overlapping populations
  - Combining information "breaks" several current techniques [Ganta, S.]

#### Goal #I: Rigor

> Raise the bar for how we think about privacy

• Especially external information

#### Make clear and refutable statements/conjectures

#### Goal #1: Rigor

Raise the bar for how we think about privacy

• Especially external information

Make clear and refutable statements/conjectures

#### • Goal #2: Interesting science

- (New) Computational phenomenon
- Unify different approaches
- Algorithmic, statistical, cryptographic challenges

"Differential" privacy

Handles arbitrary external information

> What can we compute privately?

Example technique: Output perturbation

Calibrating noise to "sensitivity"

Sample-aggregate methodology

#### "Differential" privacy

Handles arbitrary external information

> What can we compute privately?

Example technique: Output perturbation

Calibrating noise to "sensitivity"

Sample-aggregate methodology

• Intuition:

Changes to my data not noticeable by users

> Output is "independent" of my data



• Data set  $\mathbf{x} = (x_1, ..., x_n) \in D^n$ 

Domain D can be numbers, categories, tax forms

Think of x as **fixed** (not random)

• A = **randomized** procedure run by the agency

> A(x) is a random variable distributed over possible outputs Randomness might come from adding noise, resampling, etc.



x' is a neighbor of x if they differ in one data point



#### x' is a neighbor of x if they differ in one data point

Neighboring databases induce **close** distributions on outputs



x' is a neighbor of x if they differ in one data point

**Definition**: A is  $\epsilon$ -differentially private if,

Neighboring databases induce **close** distributions on outputs

for all neighbors x, x',

for all subsets S of outputs

 $\Pr(\mathsf{A}(\mathsf{x}) \in \mathsf{S}) \le e^{\epsilon} \cdot \Pr(\mathsf{A}(\mathsf{x}') \in \mathsf{S})$ 

- E cannot be too small (think  $\frac{1}{10}$ , not  $\frac{1}{2^{50}}$ )
- Distance measure on distributions matters
- This is a condition on the **algorithm** (process) A
  - Saying "this output is safe" doesn't take into account how it was computed
  - > Common problem in the literature...

**Definition**: A is *e*-differentially private if,

for all neighbors x, x',

for all subsets S of outputs

Neighboring databases induce **close** distributions on outputs

```
\Pr(\mathsf{A}(\mathsf{x}) \in \mathsf{S}) \le e^{\epsilon} \cdot \Pr(\mathsf{A}(\mathsf{x}') \in \mathsf{S})
```



$$x_i \in \{0, 1\}$$
$$\bar{x} = \frac{1}{n} \sum_i x_i$$



- Data points are binary responses  $\, x_i \in \{0,1\} \,$
- Server wants to release average  $\bar{x} = \frac{1}{n} \sum_{i} x_{i}$



- Data points are binary responses  $\, x_i \in \{0,1\} \,$
- Server wants to release average  $\bar{x} = \frac{1}{n} \sum_{i} x_{i}$

• (**Claim**: If noise ~ Lap $(\frac{1}{\epsilon n})$  then A is  $\epsilon$ -differentially private









**Definition**: A is  $\epsilon$ -differentially private if, for all neighbors x, x', for all subsets S of transcripts  $\Pr(A(x) \in S) \leq e^{\epsilon} \cdot \Pr(A(x') \in S)$ 

- "Composition": If algorithms A<sub>1</sub> and A<sub>2</sub> are E-differentially private then the outputting results of both algorithms A<sub>1</sub>(x),A<sub>2</sub>(x) is 2E-differentially private
- "Group privacy":  $k \in$ -differential privacy for groups of size  $\leq k$
- Meaningful in the presence of arbitrary external information

**Definition**: A is  $\epsilon$ -differentially private if, for all neighbors x, x', for all subsets S of transcripts  $Pr(A(x) \in S) \leq e^{\epsilon} \cdot Pr(A(x') \in S)$ 

16

• A naïve hope:

Your beliefs about me are the same after you see the output as they were before

- Suppose you know I am the height of average Canadian
  - You could learn my height from database! But it didn't matter whether or not my data was part of it.
  - > Has my privacy been compromised? No!
  - Theorem (Dwork-Naor): Learning things about individuals is unavoidable in the presence of external information

• A naïve hope:

Your beliefs about me are the same after you see the output as they were before

- Suppose you know I am the height of average Canadian
  - You could learn my height from database! But it didn't matter whether or not my data was part of it.
  - > Has my privacy been compromised? No!
  - Theorem (Dwork-Naor): Learning things about individuals is unavoidable in the presence of external information

• A naïve hope:

Your beliefs about me are the same after you see the output as they were before

- Suppose you know I am the height of average Canadian
  - You could learn my height from database! But it didn't matter whether or not my data was part of it.
  - > Has my privacy been compromised? No!
  - Theorem (Dwork-Naor): Learning things about individuals is unavoidable in the presence of external information
- [DM] Differential privacy implies: No matter what you know ahead of time,

You learn the same things about me whether or not I am in the database

Consider an intruder trying to infer personal information
"Background knowledge" = prior distribution on data x
"Conclusions you draw" = posterior p(·|output)
Experiment 0: Run A(x)
Experiment i: Run A(x<sub>-i</sub>) where x<sub>-i</sub> = (x<sub>1</sub>,...,x<sub>i-1</sub>,0,x<sub>i+1</sub>,...,x<sub>n</sub>)
## Why is this a good definition?

Consider an intruder trying to infer personal information
"Background knowledge" = prior distribution on data x
"Conclusions you draw" = posterior p(·|output)
Experiment 0: Run A(x)
Experiment i: Run A(x<sub>-i</sub>) where x<sub>-i</sub> = (x<sub>1</sub>,...,x<sub>i-1</sub>,0,x<sub>i+1</sub>,...,x<sub>n</sub>)

(**Lemma**:  $\forall$  prior,  $\forall$  output,  $p_0(\cdot|output) \approx p_i(\cdot|output)$ 

$$\begin{array}{c} \text{prior(x)} \\ \text{output } y \end{array} \left\{ \begin{array}{c} \text{Bayes' rule with} \\ Pr(y \mid x) = Pr(A(x) = y) \\ \hline \text{Bayes' rule with} \\ Pr(y \mid x) = Pr(A(x_{-i}) = y) \end{array} \right\} \xrightarrow{p_0(x \mid y)} \\ \begin{array}{c} \text{close} \\ \text{distributions} \\ \text{distributions} \\ \text{output } y \end{array} \right\}$$

## Why is this a good definition?

- Consider an intruder trying to infer personal information
  "Background knowledge" = prior distribution on data x
  "Conclusions you draw" = posterior p(·|output)
  Experiment 0: Run A(x)
  Experiment i: Run A(x<sub>-i</sub>) where x<sub>-i</sub> = (x<sub>1</sub>,...,x<sub>i-1</sub>,0,x<sub>i+1</sub>,...,x<sub>n</sub>)
- (Lemma:  $\forall$  prior,  $\forall$  output,  $p_0(\cdot|output) \approx p_i(\cdot|output)$

• Proof:  

$$p_{0}(x) = \frac{\Pr(A(x) = \text{output}) * \text{prior}(x)}{\int_{t} \Pr(A(t) = \text{output}) * \text{prior}(t)}$$

$$\approx \frac{\Pr(A(x_{-i}) = \text{output}) * \text{prior}(x)}{\int_{t} \Pr(A(t_{-i}) = \text{output}) * \text{prior}(t)} = p_{1}(x)$$

• Similar lemmas hold for relaxations of definition

### What can we compute privately?

 "Privacy" = change in one input leads to small change in output distribution

What computational tasks can we achieve privately?

- Research so far
  - Function approximation [DN, DN, BDMN, DMNS, NRS, BCDKMT, BLR]
  - Mechanism Design [MT]
  - Learning [BDMN,KLNRS]
  - Statistical estimation [S]
  - Synthetic Data [MKAGV]
  - Distributed protocols [DKMMN,BNO]
  - Impossibility results / lower bounds [DiNi,DMNS,DMT]





"Differential" privacy

Handles arbitrary external information

> What can we compute privately?

Example technique: Output perturbation

Calibrating noise to "sensitivity"

Sample-aggregate methodology

"Differential" privacy

Handles arbitrary external information

> What can we compute privately?

Example technique: Output perturbation

Calibrating noise to "sensitivity"

Sample-aggregate methodology

#### Output Perturbation, more generally



• May be interactive

> Non-interactive: release pre-defined summary stats + noise

Interactive: respond to user requests

• May be repeated many times

 $\succ$  Composition: q releases are jointly qe-differentially private

• How much noise is enough? (How much is too much?)



- Intuition:  $f(\mathbf{x})$  can be released accurately when f is insensitive to individual entries  $x_1, x_2, \ldots, x_n$
- Global Sensitivity:

$$\mathsf{GS}_{f} = \max_{\text{neighbors } x, x'} \|f(x) - f(x')\|_{1}$$

• Example:  $GS_{average} = \frac{1}{n}$ 



- Intuition:  $f(\mathbf{x})$  can be released accurately when f is insensitive to individual entries  $x_1, x_2, \ldots, x_n$
- Global Sensitivity:  $GS_f$

$$= \max_{\text{neighbors } x, x'} \|f(x) - f(x')\|_1$$

• Example:  $GS_{average} = \frac{1}{n}$ 

**Theorem:** If  $A(x) = f(x) + Lap\left(\frac{GS_f}{\epsilon}\right)$ , then A is  $\epsilon$ -differentially private.

**Theorem:** If  $A(x) = f(x) + Lap\left(\frac{GS_f}{\epsilon}\right)$ , then A is  $\epsilon$ -differentially private.

Laplace distribution  $Lap(\lambda)$  has density  $h(y) \propto e^{-\frac{\|y\|_1}{\lambda}}$ 



**Theorem:** If  $A(x) = f(x) + Lap\left(\frac{GS_f}{\epsilon}\right)$ , then A is  $\epsilon$ -differentially private.

Laplace distribution  $Lap(\lambda)$  has density  $h(y) \propto e^{-\frac{\|y\|_1}{\lambda}}$ 



Sliding property of 
$$\operatorname{Lap}\left(\frac{\operatorname{GS}_{f}}{\varepsilon}\right)$$
:  $\frac{h(y)}{h(y+\delta)} \leq e^{\varepsilon \cdot \frac{\|\delta\|}{\operatorname{GS}_{f}}}$  for all  $y, \delta$ 

**Theorem:** If  $A(x) = f(x) + Lap\left(\frac{GS_f}{\epsilon}\right)$ , then A is  $\epsilon$ -differentially private.

Laplace distribution  $Lap(\lambda)$  has density  $h(y) \propto e^{-\frac{\|y\|_1}{\lambda}}$ 

 $h(y+\delta)$ 

Sliding property of  $\operatorname{Lap}\left(\frac{\operatorname{GS}_{f}}{\varepsilon}\right)$ :  $\frac{h(y)}{h(y+\delta)} \leq e^{\varepsilon \cdot \frac{\|\delta\|}{\operatorname{GS}_{f}}}$  for all  $y, \delta$  *Proof idea:* A(x): blue curve A(x'): red curve  $\delta = f(x) - f(x') \leq \operatorname{GS}_{f}$ 

## Examples of low global sensitivity

- Many natural functions have low GS, e.g.:
  - ➤ Sample mean
  - Histograms and contingency tables
  - Covariance matrix
  - > Estimators with uniformly bounded sensitivity curve
  - Distance to a property
  - > Functions that can be approximated from a random sample
- [BDMN] Many data-mining and statistical algorithms access the data via a sequence of low-sensitivity questions
  - > e.g. perceptron, some EM algorithms, "SQ" learning algorithms

• Average: 
$$A(x) = \bar{x} + Lap(\frac{1}{\epsilon n})$$
  
> Suppose X<sub>1</sub>, X<sub>2</sub>, X<sub>3</sub>, ..., X<sub>n</sub> are i.i.d. random variables  
>  $\bar{X}$  is a random variable, and  $\sqrt{n} \cdot (\bar{X} - \mu) \xrightarrow{\mathcal{D}} Normal$   
>  $\underbrace{\frac{A(X) - \bar{X}}{StdDev(\bar{X})}}_{P} 0$  if  $\epsilon \sqrt{n} \to \infty$  with  $n$ 

> No "cost" to privacy:

- A(X) is "as good as"  $\bar{X}$  for statistical inference\*



Theorem: For any exponential family, can release
 "approximately sufficient" statistics
 > Suff. stats T(X) are sums, add noise d/(εn) for dimension d
 > A(X) - T(X) P/(StdDev(T(X)) → 0

• **Theorem:** For any exponential family, can release "approximately sufficient" statistics • Suff. stats T(X) are sums, add noise  $\frac{d}{\epsilon n}$  for dimension d •  $\frac{A(X) - T(X)}{\operatorname{StdDev}(T(X))} \xrightarrow{P} 0$ 

**Theorem:** For any well-behaved parametric family, one can construct a private efficient estimator A, if  $\epsilon \sqrt[4]{n} \to \infty$ > A(X) converges to MLE

Requires additional techniques

- **Theorem:** For any exponential family, can release "approximately sufficient" statistics • Suff. stats T(X) are sums, add noise  $\frac{d}{\epsilon n}$  for dimension d •  $\frac{A(X) - T(X)}{\operatorname{StdDev}(T(X))} \xrightarrow{P} 0$ 
  - **Theorem:** For any well-behaved parametric family, one can construct a private efficient estimator A, if  $\epsilon \sqrt[4]{n} \to \infty$ 
    - > A(X) converges to MLE

Requires additional techniques

Bounds gets worse as dimension increases

> What is the "best" private estimator?

#### $f(x) = (n_1, n_2, ..., n_d)$ where $n_j = \#\{i : x_i \text{ in } j\text{-th interval}\}$

#### $Lap(1/\epsilon)$



### Example: Histograms

- Say x<sub>1</sub>,x<sub>2</sub>,...,x<sub>n</sub> in [0,1]
  - $\succ$  Partition [0,1] into d intervals of equal size
  - $> f(x) = (n_1, n_2, ..., n_d)$  where  $n_j = #{i : x_i \text{ in } j\text{-th interval}}$ > GS<sub>f</sub> = 2
  - $\blacktriangleright \mbox{ Sufficient to add noise } \mbox{ Lap}(1/\epsilon)$  to each count
    - Independent of the dimension



## Example: Histograms

• Say x<sub>1</sub>,x<sub>2</sub>,...,x<sub>n</sub> in [0,1]

 $\geq$  Partition [0,1] into d intervals of equal size

 $> f(x) = (n_1, n_2, ..., n_d)$  where  $n_j = #{i : x_i \text{ in } j\text{-th interval}}$ > GS<sub>f</sub> = 2

 $\blacktriangleright$  Sufficient to add noise Lap $(1/\epsilon)$  to each count

• Independent of the dimension

For any smooth density h, if X<sub>i</sub> i.i.d. ~ h, noisy histogram converges to h
▷ Expected L<sub>2</sub> error O(1/(3/n)) if ε ≥ 1/(3/n)
▷ Same as non-private estimator

## Example: Histograms

- Say  $x_1, x_2, ..., x_n$  in  $\{0, 1\}$  arbitrary domain D
  - ➢ Partition [0,1] into d intervals of equal size into d disjoint "bins"  $F(x) = (n_1, n_2, ..., n_d)$  where  $n_j = \#\{i : x_i \text{ in } j\text{-th } interval\}$  bin  $GS_f = 2$
  - $\blacktriangleright \mbox{ Sufficient to add noise } \mbox{ Lap}(1/\epsilon)$  to each count
    - Independent of the dimension



# **Contingency** Tables

- Work horse of releases from US statistical agencies
   Frequencies of combinations of set of categorical attributes
- Treat as a "histogram"
  - Eight bins (O+,O-,...,AB+,AB-)
  - Can add constant noise to counts
  - $\succ$  Change to proportions is  $O(\frac{1}{n})$
  - Below sampling noise if n >> #bins
- Problem for practice:

ABO and Rh Blood Type Frequencies in the United States

ABO Type	Rh Type positive	How Many Have It	
0		38%	450/
0	negative	7%	45%
Α	positive	34%	40%
A	negative	6%	
В	positive	9%	11%
В	negative	2%	
AB	positive	3%	- 4%
AB	negative	1%	

(Source: American Association of Blood Banks)

- Some entries may be negative. Multiple tables inconsistent.
- [BCDKMT] Multiple noisy tables can be "rounded" to a consistent set of tables without increasing noise

#### Example: Distance to a Property

- Say P = set of "good" databases
   > e.g. well-clustered databases
- Distance to P = # points in x that must be changed to make <u>x</u> in P
  - $\succ$  Always has GS = 1
- Examples:
  - Distance to good clustering
  - Weight of minimum cut in graph



#### **Global Sensitivity Summary**

- Simple framework for output perturbation with strong privacy guarantees
  - > Noise levels small enough to allow meaningful analysis
- Improved in several respects
  - Worst case definition: even if f is sensitive on only one input, must add lots of noise
    - [NRS] Add less noise on "good" instances
  - One function at a time: To answer q queries, naive analysis suggests making noise increase linearly with q
    - [BLR] Simultaneously answer many "simple" questions

Focus on function approximation: many tasks not so simple

• Auction design [MT], supervised learning [KLNRS]

"Differential" privacy

Handles arbitrary external information

> What can we compute privately?

Example technique: Output perturbation

Calibrating noise to "sensitivity"

Sample-aggregate methodology

## High Global Sensitivity: Median

Example 1: median of  $x_1, \ldots, x_n \in [0, 1]$ 



- Noise magnitude:  $\frac{1}{\varepsilon}$ . Too much noise!
- But for most neighbor databases x, x', |median(x) - median(x')| is small.
- Can we add less noise on "good" instances?

# High Global Sensitivity: MST Cost

Example 2: the weight of a minimum spanning tree Database entries: edge weights in the range [0, 1].



 $\mathsf{GS}_{\mathrm{MST-weight}} = 1$ 

# High Global Sensitivity: MST Cost

Example 2: the weight of a minimum spanning tree Database entries: edge weights in the range [0, 1].



 $\mathsf{GS}_{\mathrm{MST-weight}} = 1$ 

## High Global Sensitivity: Cluster centers



Global sensitivity of cluster centers is roughly the diameter of the space.

• But intuitively, if clustering is "good", cluster centers should be insensitive.

## High Global Sensitivity: Cluster centers



Global sensitivity of cluster centers is roughly the diameter of the space.

• But intuitively, if clustering is "good", cluster centers should be insensitive.

## High Global Sensitivity: Cluster centers



Global sensitivity of cluster centers is roughly the diameter of the space.

• But intuitively, if clustering is "good", cluster centers should be insensitive.

## Getting Around Global Sensitivity

- Local sensitivity measures variability in neighborhood of specific data set [Nissim-Raskhodnikova-S, STOC 2007]
  - Connections to robust statistics
    - Bounded influence function implies expected local sensitivity is small
  - Local sensitivity needs to be smoothed
    - Interesting algorithmic/geometric problems
  - $\succ$  Not this talk
- Instead: Generic framework for smoothing functions so they have low sensitivity

## Sample-and-Aggregate Methodology

Intuition: Replace f with a less sensitive function  $\tilde{f}$ .

 $\tilde{f}(x) = g(f(sample_1), f(sample_2), \dots, f(sample_s))$ 



## Example: Efficient Point Estimates

- Given a parametric model  $\{f_{\theta}: \theta \in \Theta\}$
- $\mathsf{MLE} = \operatorname{argmax}_{\theta}(f_{\theta}(x))$
- Converges to Normal
   > Bias(MLE) = O(1/n)
  - ➤ Can be corrected so that
    bias( $\hat{\theta}$ ) = O(n<sup>-3/2</sup>)



**Theorem**: If model is well-behaved, then sampleaggregate using  $\hat{\theta}$  gives efficient estimator if  $\epsilon n^{1/4} \to \infty$ 

• Question: What is the best private estimator?

> Error bounds degrade with dimension...

# Sample-and-Aggregate Methodology

#### Theorem

If f can be approximated on xfrom small samples

then f can be released with little noise

# Sample-and-Aggregate Methodology

#### Theorem

If f can be approximated on x within distance r from small samples of size  $n^{1-\delta}$ then f can be released with little noise  $\approx \frac{r}{\varepsilon} + negl(n)$
## Sample-and-Aggregate Methodology

## Theorem

If f can be approximated on x within distance r from small samples of size  $n^{1-\delta}$ then f can be released with little noise  $\approx \frac{r}{\epsilon} + negl(n)$ 

- Works in several different metric spaces
- Example application: clustering

> I.i.d. random inputs: parametric estimation of mixture models

Arbitrary inputs: approximate optimal k-means clustering if data is "separated" à la [OstrovksyRabaniSchulmanSwamy'06]

## Conclusions

- Define privacy in terms of my effect on output
  - > Meaningful despite arbitrary external information
  - > I should participate if I get benefit
- What can we compute privately?
  - Lots of recent work
  - Existing techniques work best for highly structured computations. What about graph data, text, searches, ...?
- Data privacy is now (even) more challenging than in past
  - Data vastly more varied and valuable
  - External information more available
  - > How should we think about data privacy? (This is one example.)