# *The Privacy of Secured Computations*

Adam Smith

Penn State

Crypto & Big Data
Workshop
December 15, 2015

"Relax – it can only
see metadata."

PennState
College of Engineering

Cartoon: NOISE TO SIGNAL
RobCottingham.com

# *Big Data*

Every <length of time>
your <household object>
generates <metric scale modifier>bytes of data
about <span style="color:red">you</span>

- Everyone handles sensitive data
- Everyone delegates sensitive computations

Crypto &
Big data

# *Secured computations*



- Modern crypto offers powerful tools
  - Zero-knowledge to program obfuscation
- Broadly: specify outputs to reveal
  - … and outputs to keep secret
  - Reveal only what is necessary
- Bright lines
  - E.g., psychiatrist and patient
- Which computations should we secure?
  - Consider average salary in department before and after professor X resigns
  - Today: settings where we must release some data at the expense of others



Crypto  Big data

# *Which computations should we secure?*

- This is a social decision
  - ➤ True, but…

- Technical community can offer tools to reason about security of secured computations

- This talk: privacy in statistical databases

- Where else can technical insights be valuable?

# *Privacy in Statistical Databases*

Individuals          "Curator"          Users



Large collections of personal information
- census data
- national security data
- medical/public health data
- social networks
- recommendation systems
- trace data: search records, etc

# *Privacy in Statistical Databases*

- **Two conflicting goals**
  - Utility: Users can extract "aggregate" statistics
  - "Privacy": Individual information stays hidden

- **How can we define these precisely?**
  - Variations on model studied in
    - Statistics ("statistical disclosure control")
    - Data mining / database ("privacy-preserving data mining" *)
  - Recently: Rigorous foundations & analysis

# *Privacy in Statistical Databases*

- Why is this challenging?
  - ➤ A partial taxonomy of attacks

- Differential privacy
  - ➤ "Aggregate" as insensitive to individual changes

- Connections to other areas

# *External Information*

Individuals          Server/agency          Users



- Users have external information sources
  - Can't assume we know the sources

Anonymous data (often) isn't.

# *A partial taxonomy of attacks*

- ## Reidentification attacks
  - ➢ Based on external sources or other releases

- ## Reconstruction attacks
  - ➢ "Too many, too accurate" statistics allow data reconstruction

- ## Membership tests
  - ➢ Determine if specific person in data set (when you already know much about them)

- ## Correlation attacks
  - ➢ Learn about me by learning about population

# *Reidentification attack example*

[Narayanan, Shmatikov 2008]



**Anonymized** NetFlix data

Public, incomplete **IMDB** data

**Identified** NetFlix Data

On average, four movies uniquely identify user

# *Other reidentification attacks*

- … based on external sources, e.g.
  - ➢ Social networks
  - ➢ Computer networks
  - ➢ Microtargeted advertising
  - ➢ Recommendation Systems
  - ➢ Genetic data [Yaniv's talk]

- … based on composition attacks
  - ➢ Combining independent anonymized releases

[Citations omitted]

# *Is the problem granularity?*

- Examples so far: releasing individual information
  - ➢ What if we release only "aggregate" information?

- Defining "aggregate" is delicate
  - ➢ E.g. support vector machine output reveals individual data points

- Statistics may together encode data
  - ➢ Reconstruction attacks:
    Too many, "too accurate" stats
    $\Rightarrow$ reconstruct the data
  - ➢ Robust even to fairly significant noise

# *Reconstruction Attack Example* [Dinur Nissim '03]

- Data set: $d$ "public" attributes, 1 "sensitive"

people {  } attributes

release → reconstruction → $y' \approx y$

- Suppose release reveals correlations between attributes

  ➤ Assume one can learn $\langle a_i, y \rangle + error$

  ➤ If $error = o(\sqrt{n})$ and $a_i$ uniformly random and $d > 4n$, then one reconstruct $n - o(n)$ entries of y

- Too many, "too accurate" stats $\Rightarrow$ reconstruct data

  ➤ Cannot release everything everyone would want to know

# *Reconstruction attacks as linear encoding* [DMT'07,...]

- Data set: d "public" attributes per person, 1 "sensitive"



n people $\Big\{$ $\begin{array}{|c|c|} \hline a_i & y \\ \hline \end{array}$

d+1 attributes

release → reconstruction → $y'$ ≈ $y$

- Idea: view statistics as noisy linear encoding $My + e$

$$\begin{array}{|c|} \hline a_i \times a_j \\ \hline M \\ \hline \end{array} \cdot y + e \rightarrow y'$$

- Reconstruction depends on geometry of matrix M

  ➢ Mathematics related to "compressed sensing"

# *Membership Test Attacks*

- **[Homer et al. (2008)]**
  Exact high-dimensional summaries
  allow an attacker
  with knowledge of population
  to test membership in a data set



SNP associations

Genome-wide significance threshold

Chromosomal location

- Membership is sensitive

  ➤ Not specific to genetic data (no-fly list, census data…)

  ➤ Learn much more if statistics are provided by subpopulation

- Recently:

  ➤ Strengthened membership tests
    [Dwork, S., Steinke, Ullman, Vadhan '15]

  ➤ Tests based on learned face recognition parameters
    [Frederiksson et al '15]

# *Membership tests from marginals*

- $X$: set of $n$ binary vectors from distrib $P$ over $\{0,1\}^d$

- $q(X) = \bar{X} \in [0,1]^d$: proportion of 1 for each attribute

- $z \in \{0,1\}^d$: Alice's data

- Eve wants to know if Alice is in X.
  Eve knows
    - $q(X) = \bar{X}$
    - $z$: either in $X$ or from $P$
    - $Y$: $n$ fresh samples from $P$

- [Sankararam et al, '09]
  Eve reliably guesses if $z \in X$
  when $d > cn$

$$X =$$

| 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 |
| 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |

$$\bar{X} =$$

| ½ | ¾ | ½ | ½ | ¾ | ½ | ¼ | ¼ | ½ |
|---|---|---|---|---|---|---|---|---|

$$z =$$

| 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|

# *Strengthened membership tests* [DSSUV'15]

- $X$: set of $n$ binary vectors from distrib $P$ over $\{0,1\}^d$

- $q(X) = \bar{X} \pm \boldsymbol{\alpha}$: approximate proportions

- $z \in \{0,1\}^d$: Alice's data

- Eve wants to know if Alice is in X. Eve knows

  ➢ $q(X) = \bar{X} \pm \boldsymbol{\alpha}$

  ➢ $z$: either in $X$ or from $P$

  ➢ $Y$: $\boldsymbol{m}$ fresh samples from $P$

- [DSSUV'15] Eve reliably guesses if $z \in X$ when $d > c'\left(n + \boldsymbol{\alpha^2 n^2} + \dfrac{\boldsymbol{n^2}}{\boldsymbol{m}}\right)$

$X =$

| 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 |
| 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |

$q(X) \approx$

| ½ | ¾ | ½ | ½ | ¾ | ½ | ¼ | ¼ | ½ |
|---|---|---|---|---|---|---|---|---|

$z =$

| 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|

# *Robustness to perturbation*

- $n = 100$

- $m = 200$

- $d = 5,000$

- Two tests
  - ➢ LR [Sankararam et al'09]
  - ➢ IP [DSSUV'15]



- Two publication mechanisms
  - ➢ Rounded to nearest multiple of 0.1 (red / green)
  - ➢ Exact statistics (yellow / blue)

Conclusion: IP test is robust.
Calibrating LR test seems difficult

# *"Correlation" attacks*

- Suppose you know that I smoke and…
  - ➢ Public health study tells you that I am at risk for cancer
  - ➢ You decide not to hire me

- Learn about me by learning about underlying population
  - ➢ It does not matter which data were used in study
  - ➢ Any representative data for population will do

- Widely studied
  - ➢ De Finetti [Kifer '09]
  - ➢ Model inversion [Frederickson et al '15] *
  - ➢ Many others

- Correlation attacks fundamentally different from others
  - ➢ Do not rely on (or imply) individual data
  - ➢ Provably impossible to prevent **

* Model inversion used two few different ways in [Frederickson et al.]     ** Details later.

# *A partial taxonomy of attacks*

- Reidentification attacks
  - ➢ Based on external sources or other releases


**Identified** NetFlix Data

- Reconstruction attacks
  - ➢ "Too many, too accurate" statistics allow data reconstruction



- Membership tests
  - ➢ Determine if specific person in data set (when you already know much about them)



- Correlation attacks
  - ➢ Learn about me by learning about population

# *Privacy in Statistical Databases*

- Why is this challenging?
  - A partial taxonomy of attacks

- Differential privacy

- Connections

- "Aggregate" ≈ stability to small changes in input

- Handles arbitrary external information

- Rich algorithmic and statistical theory

# *Differential Privacy*

- Intuition:
  - ➢ Changes to my data not noticeable by users

  - ➢ Output is "independent" of my data

# *Differential Privacy* *[Dwork, McSherry, Nissim, S. 2006]*



- Data set $x = (x_1, ..., x_n) \in D^n$
  - Domain D can be numbers, categories, tax forms
  - Think of x as **fixed** (not random)
- A = **randomized** procedure
  - A(x) is a random variable
  - Randomness might come from adding noise, resampling, etc.

# *Differential Privacy* *[Dwork, McSherry, Nissim, S. 2006]*



- A thought experiment
  - ➢ Change one person's data (or remove them)
  - ➢ Will the distribution on outputs change much?

# *Differential Privacy* *[Dwork, McSherry, Nissim, S. 2006]*



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

Neighboring databases
induce **close** distributions
on outputs

**Definition**:  A is ε-differentially private if,
for all neighbors x, x',
for all subsets S of outputs

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

# *Differential Privacy* *[Dwork, McSherry, Nissim, S. 2006]*



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

**Definition**:  A is (ε,δ)-differentially private if
for all neighbors x, x',
for all subsets S of outputs

$$\Pr(A(x) \in S) \le e^\epsilon \cdot \Pr(A(x') \in S) + \delta$$

Neighboring databases
induce **close** distributions
on outputs

# *Differential Privacy* *[Dwork, McSherry, Nissim, S. 2006]*

- This is a condition on the **algorithm** A
  - ➤ Saying a particular output is private makes no sense
- Choice of distance measure matters
- What is $\varepsilon$?
  - ➤ Measure of information leakage
  - ➤ Not too small (think $\frac{1}{10}$, not $\frac{1}{2^{50}}$ )

Neighboring databases induce **close** distributions on outputs

**Definition**: A is $\varepsilon$-differentially private if, for all neighbors x, x', for all subsets S of outputs

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

# *Example: Noise Addition*



- Say we want to release a summary $f(x) \in \mathbb{R}^p$

  ➢ e.g., proportion of diabetics: $x \in \{0,1\}$ and $f(x) = \frac{1}{n} \sum_i x_i$

- Simple approach: add noise to $f(x)$

  ➢ How much noise is needed?

- Intuition: $f(x)$ can be released accurately when $f$ is insensitive to individual entries $x_1, \dots, x_n$

# *Example: Noise Addition*



function f

$A(x) = f(x) + noise$

local random coins

- Global Sensitivity: $\mathsf{GS}_f = \max\limits_{\text{neighbors } x, x'} \|f(x) - f(x')\|_1$

  ➢ Example: $\mathsf{GS}_{\text{proportion}} = \frac{1}{n}$

# *Example: Noise Addition*



function f

$$A(x) = f(x) + noise$$

local random coins

- Global Sensitivity: $\mathsf{GS}_f = \max_{\text{neighbors } x,x'} \|f(x) - f(x')\|_1$

  ➤ Example: $\mathsf{GS}_{\text{proportion}} = \frac{1}{n}$

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

  ➤ Laplace distribution $\mathsf{Lap}(\lambda)$ has density
$$h(y) \propto e^{-|y|/\lambda}$$
  ➤ Changing one point translates curve

$h(y + \mathsf{GS}_f)$ $h(y)$

# *Example: Noise Addition*



function f

$$A(x) = f(x) + noise$$

local random coins

- Example: proportion of diabetics
  - $\mathrm{GS_{proportion}} = \frac{1}{n}$
  - Release $A(x) = \text{proportion} \pm \frac{1}{\epsilon n}$
- Is this **a lot**?
  - If x is a random sample from a large underlying population, then **sampling noise** $\approx \frac{1}{\sqrt{n}}$
  - A(x) "as good as" real proportion

**proportion**   $A(X)$

# *Useful Properties*

- **Composition:**
  If A₁ and A₂ are **ε**-differentially private,
  then joint output (A₁,A₂) is **2ε**-differentially private.

- **Post processing:** A is **ε**-differentially private,
  then so is g(A) for any function g

- Meaningful in the presence of arbitrary external information

Neighboring databases induce **close** distributions on outputs

**Definition**: A is $\varepsilon$-differentially private if,
for all neighbors x, x',
for all subsets S of outputs

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

# *Interpreting Differential Privacy*

- A naïve hope:

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

- Impossible because of correlation attacks

- **Theorem** [DN'06]: Learning things about individuals is unavoidable in the presence of external information

- Differential privacy implies:

  No matter what you know ahead of time,

  > You learn (almost) the same things about me
  > whether or not my data are used

# *Features or bugs?*

- May not protect sensitive global information, e.g.
  - ➤ Clinical data: Smoking and cancer
  - ➤ Financial transactions: firm-level trading strategies
  - ➤ Social data: what if my presence affects everyone else?

- Leakage accumulates with composition
  - ➤ $\varepsilon$ adds up with many releases
    - Inevitable in some form [reconstruction attacks]
  - ➤ How do we set $\varepsilon$?

# *Variations on the approach*

- Predecessors [DDN'03,EGS'03,DN'04,BDMN'05]

- ($\varepsilon,\delta$)- differential privacy
  - ➢ Require $\Pr(A(x) \in S) \leq e^{\epsilon} \cdot \Pr(A(x) \in S) + \delta$
  - ➢ Similar semantics to ($\varepsilon,0$)- diffe.p. when $\delta \ll 1/n$

- Computational variants [MPRV09,MMPRTV'10,GKY'11]

- Distributional variants [RHMS'09,BBGLT'11,BD'12,BGKS'13]
  - ➢ Assume something about adversary's prior distribution
  - ➢ Deterministic releases
  - ➢ Composition becomes delicate

- Generalizations
  - ➢ [BLR'08, GLP'11] simulation-based definitions
  - ➢ [KM'12, BGKS'13] General language for specifying privacy concerns. Downside: tricky to instantiate

# *What can we **compute** privately?*



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

  What computational tasks can we achieve privately?

- Lots of recent work, interesting questions

  ➢ Across different fields: statistics, data mining, machine learning, cryptography, algorithmic game theory, networking, info. theory

# *A Broad, Active Field of Science*

- Basic Tools and Techniques
- Implemented systems
  - ➢ RAPPOR (Google)
  - ➢ PInQ (Microsoft)
  - ➢ Fuzz (U. Penn)
  - ➢ Privacy Tools (Harvard)
- Theoretical Foundations
  - ➢ Feasibility results: Learning, optimization, synthetic data, statistics
  - ➢ Connections to game theory, robustness, false discovery
- Domain-specific algorithms
  - ➢ Networking, clinical data, social networks, …

# *Basic Technique 1:*
## *Noise Addition*

# *Example: Noise Addition* [Dwork, McSherry, Nissim, S. 2006]

function f

$$A(x) = f(x) + \text{noise}$$

$x_1$
$x_2$
$\vdots$
$x_n$

A

local random coins

- Global Sensitivity: $\text{GS}_f = \max_{\text{neighbors } x, x'} \|f(x) - f(x')\|_1$

  ➤ Example: $\text{GS}_{\text{proportion}} = \frac{1}{n}$

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

➤ Laplace distribution $\text{Lap}(\lambda)$ has density

$$h(y) \propto e^{-|y|/\lambda}$$

➤ Changing one point translates curve

$h(y + \text{GS}_f)$ $h(y)$

# *Example: Histograms*

- Say $x_1, x_2, \ldots, x_n$ in domain D
  - Partition D into d disjoint bins
  - $f(x) = (n_1, n_2, \ldots, n_d)$ where $n_j = \#\{i : x_i \text{ in } j\text{-th bin}\}$
  - $GS_f = 1$
  - Sufficient to add noise $\mathrm{Lap}(1/\epsilon)$ to each count
- Examples
  - Histogram on the line
  - Populations of 50 states
  - Marginal tables
    - bins = possible combinations of attributes



0    1/d    1

**ABO and Rh Blood Type Frequencies in the United States**

| ABO Type | Rh Type | How Many Have It | |
|---|---|---|---|
| O | positive | 38% | 45% |
| O | negative | 7% | |
| A | positive | 34% | 40% |
| A | negative | 6% | |
| B | positive | 9% | 11% |
| B | negative | 2% | |
| AB | positive | 3% | 4% |
| AB | negative | 1% | |

(Source: American Association of Blood Banks)

# *Using global sensitivity*

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

- Many natural functions have low sensitivity
  - e.g., histogram, mean, covariance matrix, distance to a function, estimators with bounded "sensitivity curve", strongly convex optimization problems

- Laplace mechanism can be a programming interface [BDMN '05]
  - Implemented in several systems [McSherry '09, Roy et al. '10, Haeberlen et al. '11, Moharan et al. '12]

# *Variants in other metrics*

- Consider $f : \mathcal{D}^n \rightarrow \mathbb{R}^d$

- Global Sensitivity: $GS_f = \max_{\text{neighbors } x,x'} \|f(x) - f(x')\|_2$

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

$$N\left(0, \left(\frac{GS_f \cdot 3 \cdot \sqrt{\ln(1/\delta)}}{\epsilon}\right)^2\right)$$

- Example                    cates

  - $f(x)$ = vector of counts.

  - 

  - Add noise $GS_f = \sqrt{d}$            per entry instead of

$$\frac{\sqrt{d \ln(1/\delta)}}{\epsilon} \qquad\qquad \frac{d}{\epsilon}$$

# *Global versus local [NRS07]*



- Global sensitivity is worst case over inputs

- Local sensitivity:

$$\mathsf{LS}_f(x) = \max_{x' \text{ neighbor of } x} \| f(x) - f(x') \|_1$$

- Reminder:

- [NRS'07,DL'09, ...] Techniques with error ≈ local sensitivity

$$\mathsf{GS}_f(x) = \max_x \mathsf{LS}_f(x)$$

# *Basic Technique 2:*
# *Exponential Sampling*

# *Exponential Sampling* *[McSherry, Talwar '07]*

- Sometimes noise addition makes no sense
  - mode of a discrete distribution
  - minimum cut in a graph
  - classification rule

- [MT07] Motivation: auction design

- Subsequently applied very broadly

# *Example: Popular Sites*

- Data: $x_i$ = {websites visited by student i today}
- Range: Y = {website names}
- "Score" of y:     $q(y; x) = | \{i : y \subseteq x_i\} |$
- Goal: output the most frequently visited site

**Mechanism:** Given x,
- Output website $y_0$ with probability $r_x(y) \propto \exp(\epsilon q(y; x))$

- Utility: Popular sites exponentially more likely than rare ones
- Privacy: One person changes websites' scores by ≤1

$q(y; x)$

$r_x(y)$

$r_{x'}(y)$

# *Analysis*

**Mechanism:** Given x,

- Output website $y_0$ with probability $r_{\mathsf{x}}(y) \propto \exp(\epsilon q(y; \mathsf{x}))$

- **Claim:** Mechanism is 2ε-differentially private

- Proof: $\dfrac{r_{\mathsf{x}}(y)}{r_{\mathsf{x}'}(y)} = \dfrac{e^{\epsilon q(y; \mathsf{x})}}{e^{\epsilon q(y; \mathsf{x}')}} \cdot \dfrac{\sum_{z \in Y} e^{\epsilon q(z; \mathsf{x}')}}{\sum_{z \in Y} e^{\epsilon q(z; \mathsf{x})}} \le e^{2\epsilon}$

- **Claim:** If most popular website has score *T*, then

$$\mathbb{E}[q(y_0; x)] \ge T - (\log |Y|)/\epsilon$$

- Proof: Output y is bad if q(y;x) < T - k

  - $\Pr(\text{bad outputs}) \le \dfrac{\Pr(\text{bad outputs})}{\Pr(\text{best output})} \le \dfrac{|Y| e^{\epsilon(T-k)}}{e^{\epsilon T}} \le e^{\log |Y| - \epsilon k}$

  - Get expectation bound via formula $E(Z) = \sum_{k > 0} \Pr(Z \ge k)$

# *Exponential Sampling*

**Ingredients:**

- Set of outputs Y with prior distribution p(y)
- Score function q(y;x) such that
  for all outputs y, neighbors x,x':   $|q(y;x) - q(y;x')| \leq 1$

**Mechanism:** Given x,

- Output $y_0$ from Y with probability

$$r_x(y) \propto p(y) e^{\epsilon q(y;x)}$$

- Basis for first synthetic data results [Blum, Ligett, Roth '08]
  ➤ Preserve *k* linear statistics about data set with domain D

$$\frac{(\log^{1/2} k)(\log^{1/4} |D|)}{n^{1/2}}$$

# *Using Exponential Sampling*

- Mechanism above very general
  - Every differentially private mechanism is an instance!
  - Still a useful design perspective
- Perspective used explicitly for
  - Learning discrete classifiers [KLNRS'08]
  - Synthetic data generation [BLR'08,…,HLM'10]
  - Convex Optimization [CM'08,CMS'10]
  - Frequent Pattern Mining [BLST'10]
  - Genome-wide association studies [FUS'11]
  - High-dimensional sparse regression [KST'12]
  - …

# *Digital Good Auction* *[McSherry, Talwar '07]*

- 1 seller with a digital good

  Cite me maybe

- n potential buyers
  - Each has a secret value $v_i$ in [0,1] for song
  - Setting price p will get revenue rev(p) = p|{i: vi ≥ p}|
  - How can seller set p to get revenue ≈ OPT = max rev(p)?
- Straightforward bidding mechanism
  - Each player reports vi'
  - Lying can drastically change best price
- Instead, sample p* from density r(p) ∝ exp(ε . rev(p))
  - Expected revenue ≥ OPT - O( ln( ε n ) / ε )

# *A Broad, Active Field of Science*

- Basic Tools and Techniques
- Implemented systems
  - ➢ RAPPOR (Google)
  - ➢ PInQ (Microsoft)
  - ➢ Fuzz (U. Penn)
  - ➢ Privacy Tools (Harvard)
- Theoretical Foundations
  - ➢ Feasibility results: Learning, optimization, synthetic data, statistics
  - ➢ Connections to game theory, robustness, false discovery
- Domain-specific algorithms
  - ➢ Networking, clinical data, social networks, …

# *Implications for other areas*

- Game theory & economics
  - Differentially private mechanisms are automatically "approximately truthful"
  - Participating in a DP mechanism doesn't hurt me

- Statistical analysis: Differential privacy is a strong type of stability or robustness
  - Regularization techniques from optimization help design DP algorithms
  - Control **false discovery** in adaptive data analysis

# *Ongoing Work*

- Practical implementations

- Efficient algorithms

- Relaxed definitions
  - ➢ Exploit adversarial uncertainty

- Differently-structured data
  - ➢ E.g., social network data: which data is "mine"?

# *Conclusions*

- Define privacy in terms of my effect on output
  - ➢ Meaningful despite arbitrary external information
  - ➢ I should participate if I get benefit
- Rigorous framework for private data analysis
  - ➢ Rich algorithmic literature (theoretical and applied)
  - ➢ There is no competing theory

- What computations can we secure?
  - ➢ Differential privacy provided a surprising formalization for a previously ad hoc area
  - ➢ What other areas need formalization?
    - How should we think about correlation attacks?

# *Further resources*

- Tutorial from CRYPTO 2012
  - ➢ http://www.cse.psu.edu/~asmith/talks/2012-08-21-crypto-tutorial.pdf
- Courses:
  - ➢ http://www.cis.upenn.edu/~aaroth/courses/privacyF11.html
  - ➢ http://www.cse.psu.edu/~asmith/privacy598
- DIMACS Workshop on Data Privacy (October 2012)
  - ➢ Videos of tutorials
  - ➢ http://dimacs.rutgers.edu/Workshops/DifferentialPrivacy/
- Simons Institute Big Data & DP Workshop (Dec 2013)
  - ➢ Talk videos online