Techniques for Private Data Analysis

Sofya Raskhodnikova

Penn State University

Based on joint work with Shiva Kasiviswanathan, Homin Lee, Kobbi Nissim and Adam Smith

Private data analysis



Collections of personal and sensitive data

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

WHAT INFORMATION CAN BE RELEASED?

- Two conflicting goals
 - utility: users can extract "global" statistics
 - privacy: individual information stays hidden

Related work

Other fields: huge amount of work

- in statistics (statistical disclosure limitation)
- in data mining (privacy-preserving data mining)
- largely: no precise privacy definition (only security against specific attacks)

In cryptography (private data analysis)

- [Dinur Nissim 03, Dwork Nissim 04, Chawla Dwork McSherry Smith Wee 05, Blum Dwork McSherry Nissim 05, Chawla Dwork McSherry Talwar 05, Dwork McSherry Nissim Smith 06, ...]
- rigorous privacy guarantees

Differential privacy [DMNS06]

Intuition: Users learn the same thing about me whether or not I participate in the census

Two databases are *neighbors* if they differ in one row (arbitrarily complex information supplied by one person).



Privacy definition

Algorithm A is ε -differentially private if

- for all neighbor databases x, x'
- for all sets of answers S
 - $\Pr[A(x) \in S] \le (1 + \varepsilon) \cdot \Pr[A(x') \in S]$

Properties of differential privacy



- ε is non-negligible (at least $\frac{1}{n}$).
- Composition: If A_1 and A_2 are ε -differentially private then (A_1, A_2) is 2ε -differentially private
- robust in the presence of arbitrary side information

What can we compute privately?

Research so far:

- Definitions [DiNi,DwNi,EGS,DMNS,DwNa,DKMMN,GKS]
- Function approximation



- Protocols [DiNi,DwNi,BDMN,DMNS,NRS,BCDKMT]
- Impossibility results [DiNi,DMNS,DwNa,DwMT,DwY]
- Distributed protocols [DKMMN,BNiO]
- Mechanism design [McSherry Talwar 07]
- Learning [Blum Dwork McSherry Nissim 05, KLNRS08]
- Releasing classes of functions [Blum Ligett Roth 08]
- Synthetic data [Machanavajjhala Kifer Abowd Gehrke Vilhuber 08]

I. Function approximation

- Global sensitivity framework [DMNS06]
- Smooth sensitivity framework [NRS07]
- Sample-and-aggregate [NRS07]
- II. Learning
 - Exponential mechanism [MT07,KLNRS08]

Function Approximation



For which functions f can we have:

- privacy: differential privacy [DMNS06]
- utility: output A(x) is close to f(x)

Global sensitivity framework [DMNS06]

Intuition: f can be released accurately when it is insensitive to individual entries x_1, \ldots, x_n .

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

Example: $\mathsf{GS}_{\text{average}} = \frac{1}{n} \text{ if } x \in [0, 1]^n.$

	eorem
$If A(x) = f(x) + Lap\left(\frac{GS_f}{\varepsilon}\right)$	then A is ε -diff. private.

Global sensitivity framework [DMNS06]

Intuition: f can be released accurately when it is insensitive to individual entries x_1, \ldots, x_n . Global sensitivity $\mathsf{GS}_f = \max_{\text{neighbors } x, x'} \|f(x) - f(x')\|_1$. Example: $\mathsf{GS}_{\text{average}} = \frac{1}{n}$ if $x \in [0, 1]^n$. Noise $= \mathsf{Lap}(\frac{1}{\varepsilon n})$. Compare to: Estimating frequencies (e.g., proportion of people with blue eyes) from n samples: sampling error $\frac{1}{\sqrt{n}}$.

Theorem

If $A(x) = f(x) + Lap\left(\frac{GS_f}{\varepsilon}\right)$	then A is ε -diff. private.
---	---

Global sensitivity framework [DMNS06]

Intuition: f can be released accurately when it is insensitive to individual entries x_1, \ldots, x_n . Global sensitivity $\mathsf{GS}_f = \max_{\text{neighbors } x, x'} \|f(x) - f(x')\|_1$. Example: $\mathsf{GS}_{\text{average}} = \frac{1}{n}$ if $x \in [0, 1]^n$. Noise $= \mathsf{Lap}(\frac{1}{\varepsilon n})$. Compare to: Estimating frequencies (e.g., proportion of people with blue eyes) from n samples: sampling error $\frac{1}{\sqrt{n}}$.

Theorem

If $A(x) = f(x) + Lap\left(\frac{GS_f}{\varepsilon}\right)$ then A is ε -diff. private.

Functions with low global sensitivity

- Means, variances for data in a bounded interval
- histograms, contingency tables
- singular value decomposition

Big picture for global sensitivity framework:

- add enough noise to cover the worst case for f
- noise distribution depends only on f, not database x

Problem: for some functions that's too much noise

Smooth sensitivity framework [Nissim Smith Raskhodnikova 07]: noise tuned to database x

Local sensitivity $\mathsf{LS}_f(x) = \max_{\substack{x': \text{ neighbor of } x}} \|f(x) - f(x')\|$ *Reminder:* $\mathsf{GS}_f = \max_x \mathsf{LS}_f(x)$ *Example: median* for $0 \le x_1 \le \cdots \le x_n \le 1$, odd n



 $\mathsf{LS}_{\mathrm{median}}(x) = \max(x_m - x_{m-1}, x_{m+1} - x_m)$

Goal: Release f(x) with less noise when $\mathsf{LS}_f(x)$ is lower.

Local sensitivity $\mathsf{LS}_f(x) = \max_{\substack{x': \text{ neighbor of } x}} \|f(x) - f(x')\|$ *Reminder:* $\mathsf{GS}_f = \max_x \mathsf{LS}_f(x)$ *Example: median* for $0 \le x_1 \le \cdots \le x_n \le 1$, odd n



Goal: Release f(x) with less noise when $\mathsf{LS}_f(x)$ is lower.

Local sensitivity $\mathsf{LS}_f(x) = \max_{\substack{x': \text{ neighbor of } x}} \|f(x) - f(x')\|$ *Reminder:* $\mathsf{GS}_f = \max_x \mathsf{LS}_f(x)$ *Example: median* for $0 \le x_1 \le \cdots \le x_n \le 1$, odd n



Goal: Release f(x) with less noise when $\mathsf{LS}_f(x)$ is lower.

- Noise magnitude proportional to $\mathsf{LS}_f(x)$ instead of GS_f ?
- *No!* Noise magnitude reveals information.
- *Lesson:* Noise magnitude must be an insensitive function.

Smooth bounds on local sensitivity

Design sensitivity function S(x)

- S(x) is an ε -smooth upper bound on $\mathsf{LS}_f(x)$ if:
 - for all x: $S(x) \ge \mathsf{LS}_f(x)$
 - for all neighbors $x, x' : \quad S(x) \le e^{\varepsilon} S(x')$



 $\frac{\text{Theorem}}{\text{If } A(x) = f(x) + \text{noise}\left(\frac{S(x)}{\varepsilon}\right) \text{ then } A \text{ is } \varepsilon' \text{-differentially private.}$

Example: GS_f is always a smooth bound on $\mathsf{LS}_f(x)$

Smooth bounds on local sensitivity

Design sensitivity function S(x)

- S(x) is an ε -smooth upper bound on $\mathsf{LS}_f(x)$ if:
 - for all x: $S(x) \ge \mathsf{LS}_f(x)$
 - for all neighbors $x, x' : \quad S(x) \le e^{\varepsilon} S(x')$



 $\frac{\text{Theorem}}{\text{If } A(x) = f(x) + \text{noise}\left(\frac{S(x)}{\varepsilon}\right) \text{ then } A \text{ is } \varepsilon' \text{-differentially private.}$

Example: GS_f is always a smooth bound on $\mathsf{LS}_f(x)$

Smooth Sensitivity

Smooth sensitivity $S_f^*(x) = \max_y \left(\mathsf{LS}_f(y) e^{-\varepsilon \cdot \mathsf{dist}(x,y)} \right)$



Intuition: little noise when **far** from sensitive instances



Smooth Sensitivity

Smooth sensitivity $S_f^*(x) = \max_y \left(\mathsf{LS}_f(y) e^{-\varepsilon \cdot \mathsf{dist}(x,y)} \right)$



Intuition: little noise when **far** from sensitive instances



Computing smooth sensitivity

Example functions with computable smooth sensitivity

- Median & minimum of numbers in a bounded interval
- *MST cost* when weights are bounded
- Number of triangles in a graph

Approximating smooth sensitivity

- \bullet only smooth upper bounds on LS are meaningful
- simple generic methods for smooth approximations – work for *median* and *1-median* in L_1^d

I. Function approximation

- Global sensitivity framework [DMNS06]
- Smooth sensitivity framework [NRS07]
- Sample-and-aggregate [NRS07]
- II. Learning
 - Exponential mechanism [MT07,KLNRS08]

- Smooth sensitivity framework requires understanding combinatorial structure of f
 – hard in general
- Goal: an automatable transformation from an arbitrary f into an ε-diff. private A
 - A(x) ≈ f(x) for "good" instances x

Example: cluster centers



- Comparing sets of centers: Earthmover-like metric
- Global sensitivity of cluster centers is roughly the diameter of the space. But intuitively, if clustering is "good", cluster centers should be insensitive.
- No efficient approximation for smooth sensitivity

Example: cluster centers



- Comparing sets of centers: Earthmover-like metric
- Global sensitivity of cluster centers is roughly the diameter of the space. But intuitively, if clustering is "good", cluster centers should be insensitive.
- No efficient approximation for smooth sensitivity

Example: cluster centers



- Comparing sets of centers: Earthmover-like metric
- Global sensitivity of cluster centers is roughly the diameter of the space. But intuitively, if clustering is "good", cluster centers should be insensitive.
- No efficient approximation for smooth sensitivity

Sample-and-aggregate framework

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

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



Sample-and-aggregate framework

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

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



Good aggregation functions

• average

- works for L_1 and L_2
- center of attention
 - the center of a smallest ball containing a strict majority of input points
 - works for arbitrary metrics
 - (in particular, for Earthmover)
 - gives lower noise for L_1 and L_2

Sample-and-aggregate method

Theorem

If f can be approximated on xfrom small samples then f can be released with little noise

Sample-and-aggregate method

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 $\approx \frac{r}{\epsilon} + negl(n)$

Sample-and-aggregate method

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 $\approx \frac{r}{\varepsilon} + negl(n)$

- Works in all "interesting" metric spaces
- Example applications
 - k-means cluster centers (if data is separated a.k.a.[Ostrovsky Rabani Schulman Swamy 06])
 - fitting mixtures of Gaussians (if data is i.i.d., using [Achlioptas McSherry 05])
 - PAC concepts (if uniquely learnable,
 - i.e., if learning algorithm always outputs the same hypothesis or something close to it)

I. Function approximation

- Global sensitivity framework [DMNS06]
- Smooth sensitivity framework [NRS07]
- Sample-and-aggregate [NRS07]

II. Learning

• Exponential mechanism [McSherry Talwar 07, Kasiviswanathan Lee Nissim Raskhodnikova Smith 08]

Bank needs to decide which applicants are bad credit risks^{*} Goal: given sample of labeled data (past customers), produce good prediction rule (hypothesis) for future loan applicants

^{*}Example taken from Blum, FOCS03 tutorial

Bank needs to decide which applicants are bad credit risks^{*} Goal: given sample of labeled data (past customers), produce good prediction rule (hypothesis) for future loan applicants

/ go ris	$\underset{\mathrm{inc}}{\mathrm{mmp}}$	other accts	high debt	% down
Y	0.32	Yes	No	10
Y	0.25	No	No	10
N	0.30	No	Yes	5
Y	0.31	Yes	No	20
Y	0.25	No	No	10

^{*}Example taken from Blum, FOCS03 tutorial

Bank needs to decide which applicants are bad credit risks^{*} Goal: given sample of labeled data (past customers), produce good prediction rule (hypothesis) for future loan applicants

	$\frac{\%}{\mathrm{down}}$	$\mathop{ m high}\limits_{ m debt}$	other accts	$\begin{array}{c} \mathrm{mmp}/\\ \mathrm{inc} \end{array}$	good risk?
	10	No	Yes	0.32	Yes
example y_i	10	No	No	0.25	Yes
	5	Yes	No	0.30	No
	20	No	Yes	0.31	Yes
	10	No	No	0.25	Yes

^{*}Example taken from Blum, FOCS03 tutorial

Bank needs to decide which applicants are bad credit risks^{*} Goal: given sample of labeled data (past customers), produce good prediction rule (hypothesis) for future loan applicants

	% down	high debt	other accts	$\mathop{\mathrm{mmp}}_{\mathrm{inc}}$	good risk?	
	10	No	Yes	0.32	Yes	
example y_i	10	No	No	0.25	Yes	label
	5	Yes	No	0.30	No	ſ
	20	No	Yes	0.31	Yes	
	10	No	No	0.25	Yes	

 z_i

^{*}Example taken from Blum, FOCS03 tutorial

Bank needs to decide which applicants are bad credit risks^{*} Goal: given sample of labeled data (past customers), produce good prediction rule (hypothesis) for future loan applicants

	% down	high debt	other accts	$\mathop{\mathrm{mmp}}_{\mathrm{inc}}$	good risk?	
	10	No	Yes	0.32	Yes	
example y_i	10	No	No	0.25	Yes	label z_i
	5	Yes	No	0.30	No	
	20	No	Yes	0.31	Yes	
	10	No	No	0.25	Yes	

Reasonable rules given this data:

- Predict YES iff $100 \times \frac{\text{mmp}}{\text{inc}} (\% \text{ down}) < 25$
- Predict YES iff (!high debt) AND (% down > 5)

*Example taken from Blum, FOCS03 tutorial





 \bullet Examples drawn according to distribution ${\cal D}$



 \bullet Examples drawn according to distribution ${\cal D}$



- \bullet Examples drawn according to distribution ${\cal D}$
- A point drawn according to \mathcal{D} has to be classified correctly w.h.p. (over learner randomness and \mathcal{D})

Given distribution \mathcal{D} over examples, labeled by function c, hypothesis h is *good* if it mostly agrees with c:

$$\Pr_{y \sim \mathcal{D}}[h(y) = c(y)] \text{ is close to } 1.$$

Given distribution \mathcal{D} over examples, labeled by function c, hypothesis h is *good* if it mostly agrees with c:

$$\Pr_{y \sim \mathcal{D}}[h(y) = c(y)] \text{ is close to } 1.$$

Definition of PAC learning

Algorithm A PAC learns a concept class C if

- given polynomially many examples, drawn from \mathcal{D} , labeled by some $c \in C$
- A outputs a good hypothesis with high probability in polynomial time

Given distribution \mathcal{D} over examples, labeled by function c, hypothesis h is *good* if it mostly agrees with c:

$$\Pr_{y \sim \mathcal{D}}[h(y) = c(y)] \text{ is close to } 1.$$

Definition of PAC^{*} learning

Algorithm A PAC^{*} learns a concept class C if

- given polynomially many examples, drawn from \mathcal{D} , labeled by some $c \in C$
- A outputs a good hypothesis of polynomial length with high probability in polynomial time

Input: database $x = (x_1, ..., x_n)$ $x_i = (y_i, z_i)$, where $y_i \sim \mathcal{D}$ and $z_i = c(y_i)$ is the label of example y_i

% down	high debt	other accts	$\begin{array}{c} \mathrm{mmp}/\\ \mathrm{inc} \end{array}$	good risk?
10	No	Yes	0.32	Yes
10	No	No	0.25	Yes
5	Yes	No	0.30	No
20	No	Yes	0.31	Yes
10	No	No	0.25	Yes

Input: database $x = (x_1, ..., x_n)$ $x_i = (y_i, z_i)$, where $y_i \sim \mathcal{D}$ and $z_i = c(y_i)$ is the label of example y_i

% down	high debt	other accts	$\begin{array}{c} \mathrm{mmp}/\\ \mathrm{inc} \end{array}$	good risk?
10	No	Yes	0.32	Yes
10	No	No	0.25	Yes
5	Yes	No	0.30	No
20	No	Yes	0.31	Yes
10	No	No	0.25	Yes



Input: database $x = (x_1, ..., x_n)$ $x_i = (y_i, z_i)$, where $y_i \sim \mathcal{D}$ and $z_i = c(y_i)$ is the label of example y_i

% down	high debt	other accts	$\begin{array}{c} \mathrm{mmp}/\\ \mathrm{inc} \end{array}$	good risk?
10	No	Yes	0.32	Yes
10	No	No	0.25	Yes
5	Yes	No	0.30	No
20	No	Yes	0.31	Yes
10	No	No	0.25	Yes



Definition

Algorithm A privately learns concept class C if

- **Utility:** Algorithm A PAC learns class C
- **Privacy:** Algorithm A is differentially private

Input: database $x = (x_1, ..., x_n)$ $x_i = (y_i, z_i)$, where $y_i \sim \mathcal{D}$ and $z_i = c(y_i)$ is the label of example y_i

% down	high debt	other accts	$\begin{array}{c} \mathrm{mmp}/\\ \mathrm{inc} \end{array}$	good risk?
10	No	Yes	0.32	Yes
10	No	No	0.25	Yes
5	Yes	No	0.30	No
20	No	Yes	0.31	Yes
10	No	No	0.25	Yes



Definition

Algorithm A privately learns concept class C if

- Utility: Algorithm A PAC learns class C (average-case)
- **Privacy:** Algorithm A is differentially private (worst-case)

Designing private learners: baby steps

View non-private learner as function to be approximated

- First attempt: add noise
 - Problem: Close hypothesis
 may mislabel many points



Designing private learners: baby steps

View non-private learner as function to be approximated

- First attempt: add noise
 - Problem: Close hypothesis may mislabel many points
- Second attempt:
 - apply sample-and-aggregate to non-private learning algorithm
 - Works when good hypothesis is essentially unique
 - Problem: there may be many good hypotheses different samples may produce different-looking hypotheses





Each PAC^{*} learnable concept class can be learned privately, using polynomially many samples.







- Output h from C with probability $\sim e^{\varepsilon \cdot \text{score}(x,h)}$
 - may take exponential time



Proof: Adapt the exponential mechanism of [MT07]: score(x, h) = # of examples in x correctly classified by h

• Output h from C with probability $\sim e^{\varepsilon \cdot \text{score}(x,h)}$

– may take exponential time





Each PAC^{*} learnable concept class can be learned privately, using polynomially many samples.

Proof: score(x, h) = # of examples in x correctly classified by h

• Output h from C with probability $\sim e^{\varepsilon \cdot \text{score}(x,h)}$

Utility (learning):

Each PAC^{*} learnable concept class can be learned privately, using polynomially many samples.

Proof: score(x, h) = # of examples in x correctly classified by h

• Output h from C with probability $\sim e^{\varepsilon \cdot \text{score}(x,h)}$

Utility (learning):

Good *h* correctly label all examples: $\Pr[h] \sim e^{\varepsilon \cdot n}$ **Bad** *h* mislabel $\geq 10\%$ of examples: $\Pr[h] \sim e^{\varepsilon \cdot 0.9n}$

Each PAC^{*} learnable concept class can be learned privately, using polynomially many samples.

Proof: score(x, h) = # of examples in x correctly classified by h

• Output h from C with probability $\sim e^{\varepsilon \cdot \text{score}(x,h)}$

Utility (learning):

Good *h* correctly label all examples: $\Pr[h] \sim e^{\varepsilon \cdot n}$ **Bad** *h* mislabel $\geq 10\%$ of examples: $\Pr[h] \sim e^{\varepsilon \cdot 0.9n}$

Sufficient to ensure $n \gg \log(\# \text{ bad hypotheses}) = \text{polynomial}$

 \leq output length of \leq non-private learner

Then w.h.p. output h labels 90% of examples correctly.

Each PAC^{*} learnable concept class can be learned privately, using polynomially many samples.

Proof: score(x, h) = # of examples in x correctly classified by h

• Output h from C with probability $\sim e^{\varepsilon \cdot \text{score}(x,h)}$

Utility (learning):

Good *h* correctly label all examples: $\Pr[h] \sim e^{\varepsilon \cdot n}$ **Bad** *h* mislabel $\geq 10\%$ of examples: $\Pr[h] \sim e^{\varepsilon \cdot 0.9n}$

Sufficient to ensure $n \gg \log(\# \text{ bad hypotheses}) = \text{polynomial}$

 $\stackrel{\rm output \ length \ of}{= non-private \ learner}$

Then w.h.p. output h labels 90% of examples correctly.

By "Occam's razor", if $n \gg \log(\# \text{ hypotheses})$, then h does well on examples $\implies h$ does well on distrib. \mathcal{D}

I. Function approximation

- Global sensitivity framework [DMNS06]
- Smooth sensitivity framework [NRS07]
- Sample-and-aggregate [NRS07]
- II. Learning
 - Exponential mechanism [MT07,KLNRS08]

This talk: partial picture of techniques

- current techniques best for
 - function approximation
 - learning
- New ideas needed for
 - combinatorial search problems
 - text processing
 - graph data (definitions?)
 - high-dimensional outputs