

Smooth Sensitivity and Sampling

Sofya Raskhodnikova

Penn State University

Joint work with **Kobbi Nissim** (Ben Gurion University)

and **Adam Smith** (Penn State University)

Our main contributions

- *Starting point:* Global sensitivity framework [DMNS06]
(Cynthia's talk)
- Two new frameworks for private data analysis
- Greatly expand the types of information that can be released

Road map

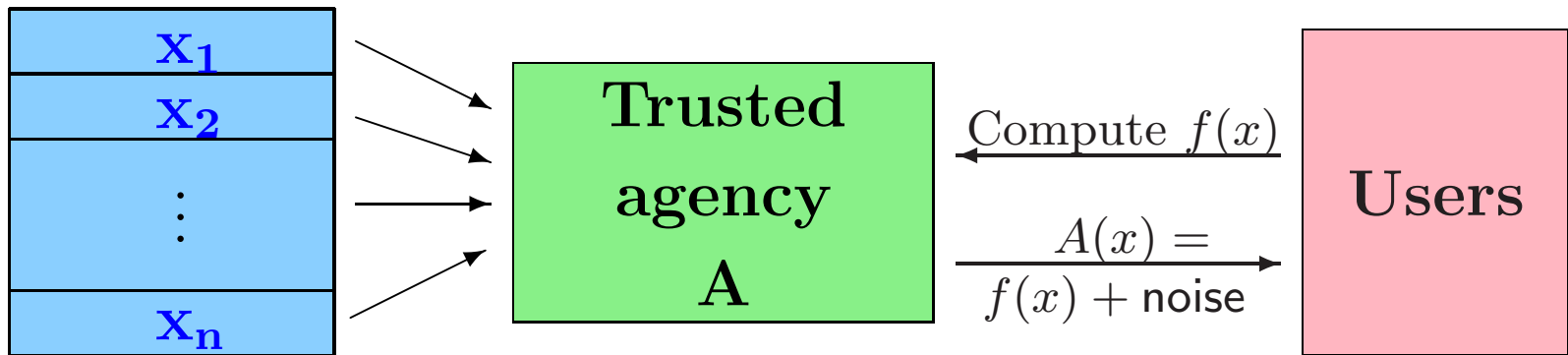
I. Introduction

- Review of global sensitivity framework [DMNS06]
- Motivation

II. Smooth sensitivity framework

III. Sample-and-aggregate framework

Model



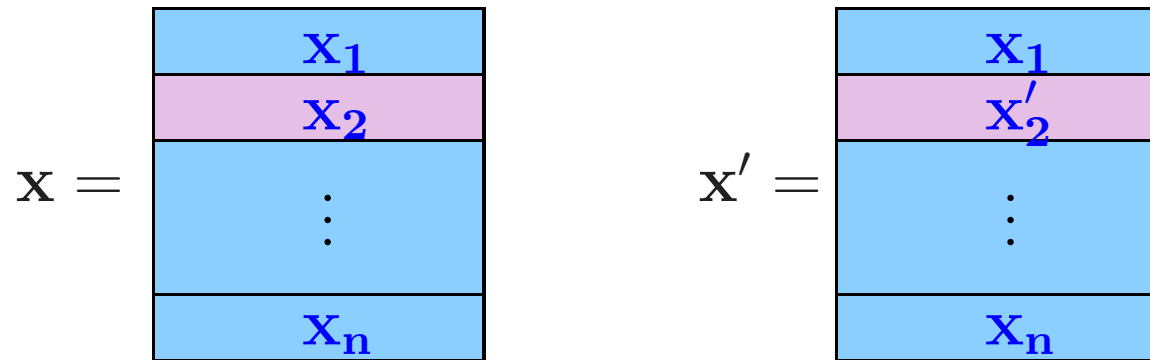
Each row is arbitrarily complex data supplied by 1 person.

For which functions f can we have:

- **utility**: little noise
- **privacy**: indistinguishability definition of [DMNS06]

Privacy as indistinguishability [DMNS06]

Two databases are *neighbors* if they differ in one row.



Privacy definition

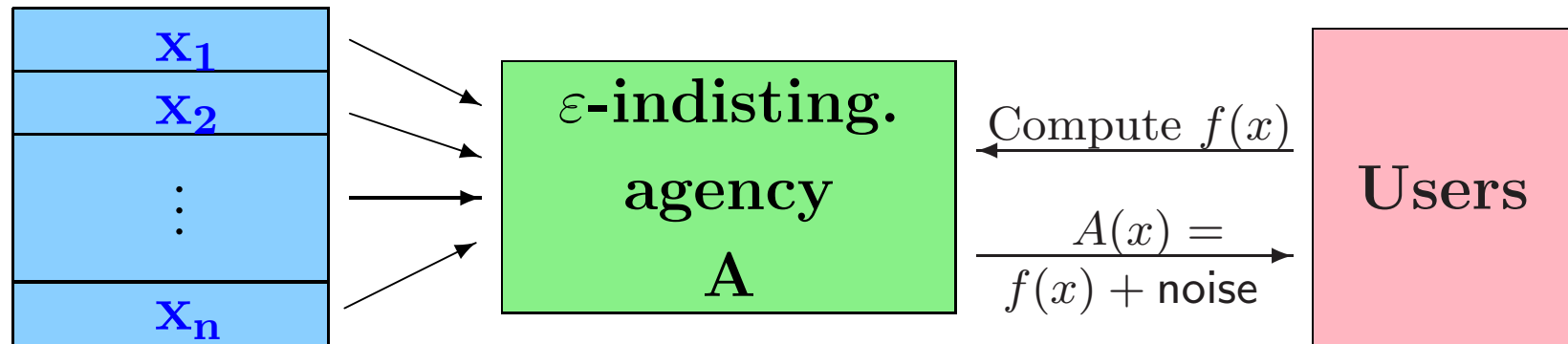
Algorithm A is ϵ -indistinguishable if

- for all neighbor databases x, x'
- for all sets of answers S

$$\Pr[A(x) \in S] \leq (1 + \epsilon) \cdot \Pr[A(x') \in S]$$

Privacy definition: composition

If A is ϵ -indistinguishable on each query,
it is ϵq -indistinguishable on q queries.



Global sensitivity framework [DMNS06]

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

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

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

Theorem

If $A(x) = f(x) + \text{Lap}\left(\frac{\text{GS}_f}{\varepsilon}\right)$ then A is ε -indistinguishable.

Instance-Based Noise

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

Instance-Based Noise

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

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

$$x = \underbrace{0 \dots 0}_{\frac{n-1}{2}} \underbrace{0 \ 1 \dots 1}_{\frac{n-1}{2}}$$

$$\text{median}(x) = 0$$

$$x' = \underbrace{0 \dots 0}_{\frac{n-1}{2}} \underbrace{1 \ 1 \dots 1}_{\frac{n-1}{2}}$$

$$\text{median}(x') = 1$$

$$\text{GS}_{\text{median}} = 1$$

- Noise magnitude: $\frac{1}{\epsilon}$.

Instance-Based Noise

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

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

$$x = \underbrace{0 \dots 0}_{\frac{n-1}{2}} 0 \underbrace{1 \dots 1}_{\frac{n-1}{2}}$$

$$\text{median}(x) = 0$$

$$x' = \underbrace{0 \dots 0}_{\frac{n-1}{2}} 1 \underbrace{1 \dots 1}_{\frac{n-1}{2}}$$

$$\text{median}(x') = 1$$

$$\text{GS}_{\text{median}} = 1$$

- Noise magnitude: $\frac{1}{\epsilon}$.

Our goal: noise tuned to database x

Road map

I. Introduction

- Review of global sensitivity framework [DMNS06]
- Motivation

II. Smooth sensitivity framework

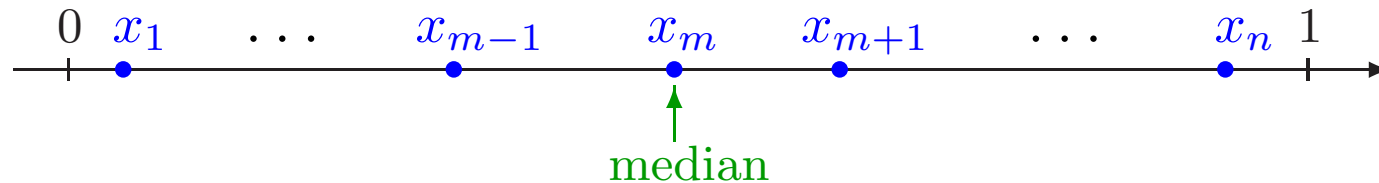
III. Sample-and-aggregate framework

Local sensitivity

$$\text{Local sensitivity } \text{LS}_f(x) = \max_{x': \text{neighbor of } x} \|f(x) - f(x')\|$$

$$\text{Reminder: } \text{GS}_f = \max_x \text{LS}_f(x)$$

Example: median for $0 \leq x_1 \leq \dots \leq x_n \leq 1$, odd n



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

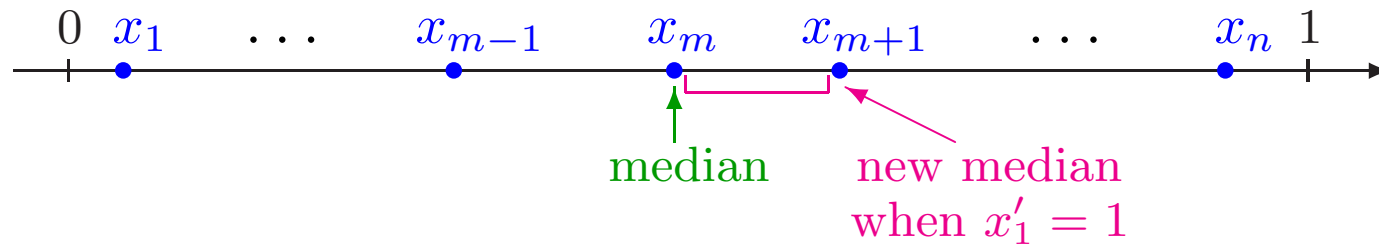
Goal: Release $f(x)$ with less noise when $\text{LS}_f(x)$ is lower.

Local sensitivity

$$\text{Local sensitivity } \text{LS}_f(x) = \max_{x': \text{neighbor of } x} \|f(x) - f(x')\|$$

$$\text{Reminder: } \text{GS}_f = \max_x \text{LS}_f(x)$$

Example: median for $0 \leq x_1 \leq \dots \leq x_n \leq 1$, odd n



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

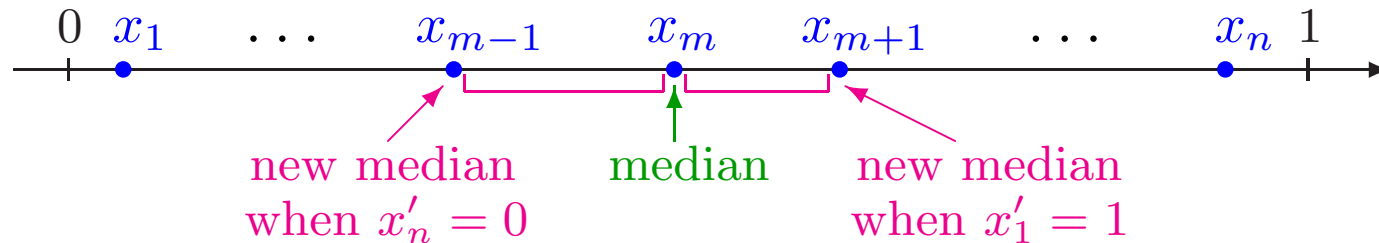
Goal: Release $f(x)$ with less noise when $\text{LS}_f(x)$ is lower.

Local sensitivity

$$\text{Local sensitivity } \text{LS}_f(x) = \max_{x': \text{neighbor of } x} \|f(x) - f(x')\|$$

$$\text{Reminder: } \text{GS}_f = \max_x \text{LS}_f(x)$$

Example: median for $0 \leq x_1 \leq \dots \leq x_n \leq 1$, odd n



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

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

Instance-based noise: first attempt

Noise magnitude proportional to $LS_f(x)$ instead of GS_f ?

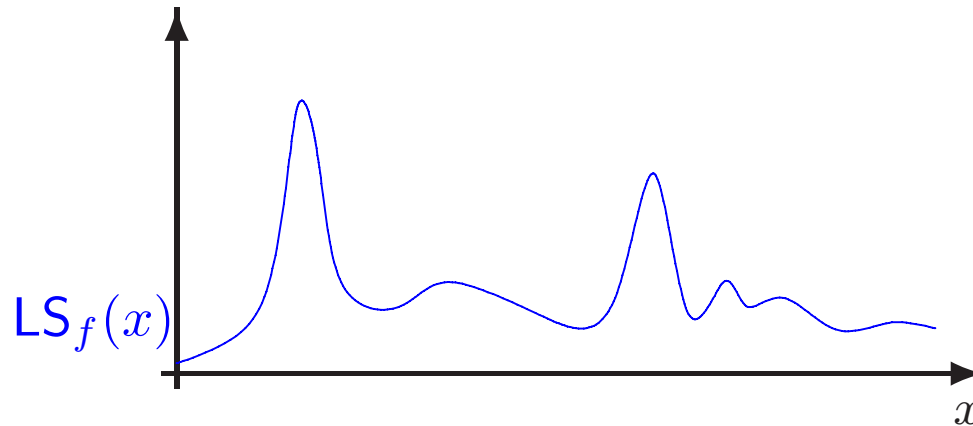
No! Noise magnitude reveals information.

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 $\text{LS}_f(x)$ if:
 - for all x : $S(x) \geq \text{LS}_f(x)$
 - for all neighbors x, x' : $S(x) \leq e^\varepsilon S(x')$



Theorem

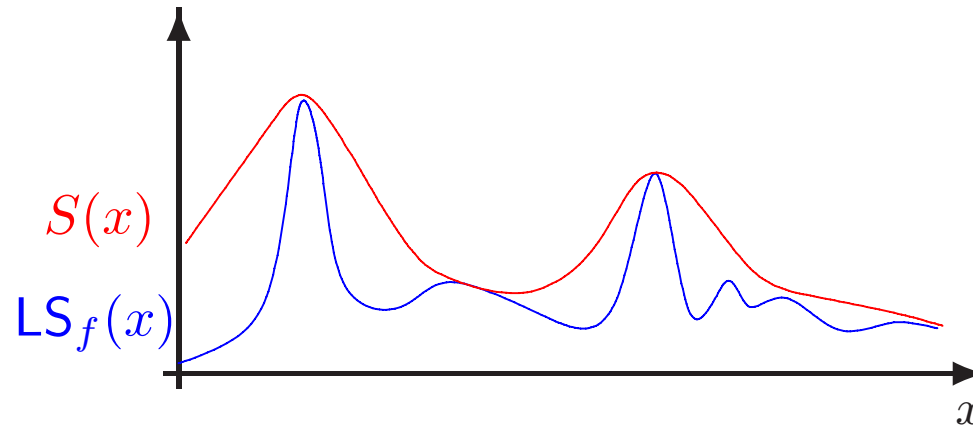
If $A(x) = f(x) + \text{noise} \left(\frac{S(x)}{\varepsilon} \right)$ then A is ε' -indistinguishable.

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

Smooth bounds on local sensitivity

Design sensitivity function $S(x)$

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



Theorem

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

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

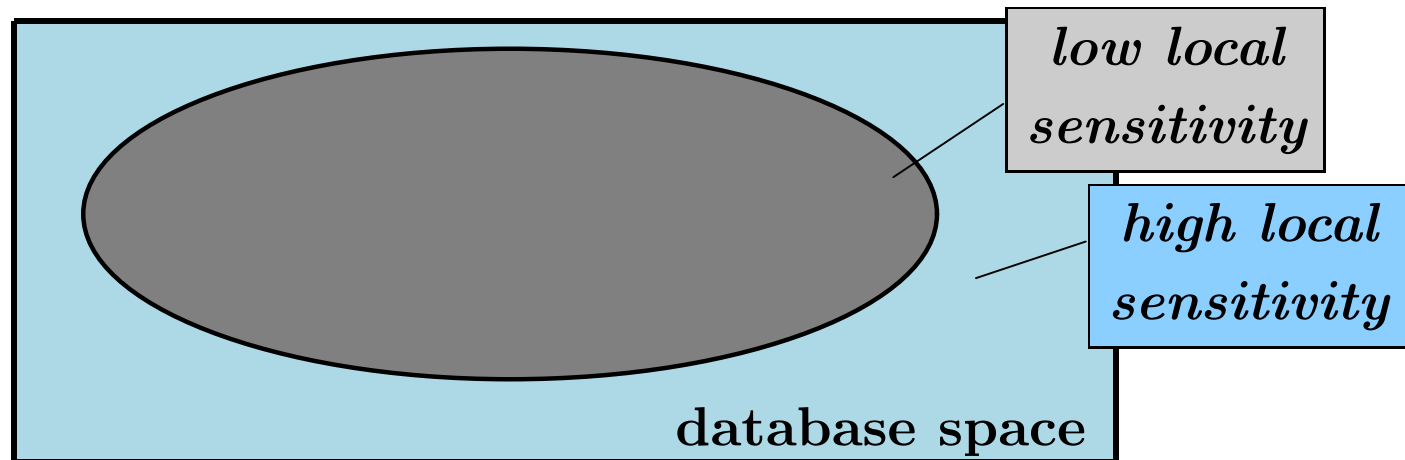
Smooth Sensitivity

Smooth sensitivity $S_f^*(x) = \max_y (LS_f(y) e^{-\varepsilon \cdot \text{dist}(x,y)})$

Lemma

For every ε -smooth bound S : $S_f^*(x) \leq S(x)$ for all x .

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



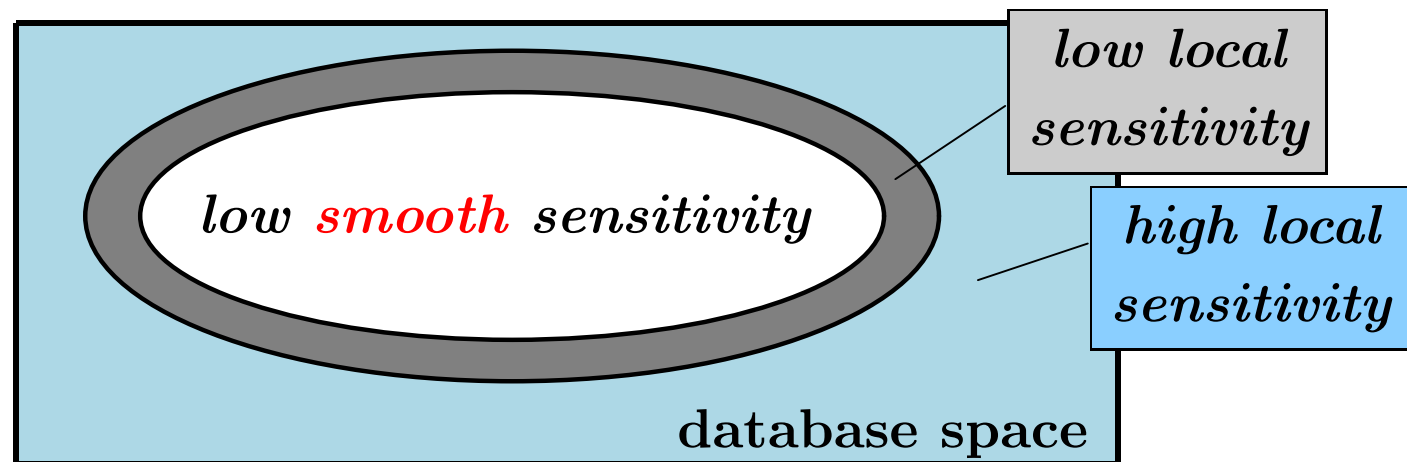
Smooth Sensitivity

$$\text{Smooth sensitivity } S_f^*(x) = \max_y (LS_f(y) e^{-\varepsilon \cdot \text{dist}(x,y)})$$

Lemma

For every ε -smooth bound S : $S_f^*(x) \leq S(x)$ for all x .

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

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

Road map

I. Introduction

- Review of global sensitivity framework [DMNS06]
- Motivation

II. Smooth sensitivity framework

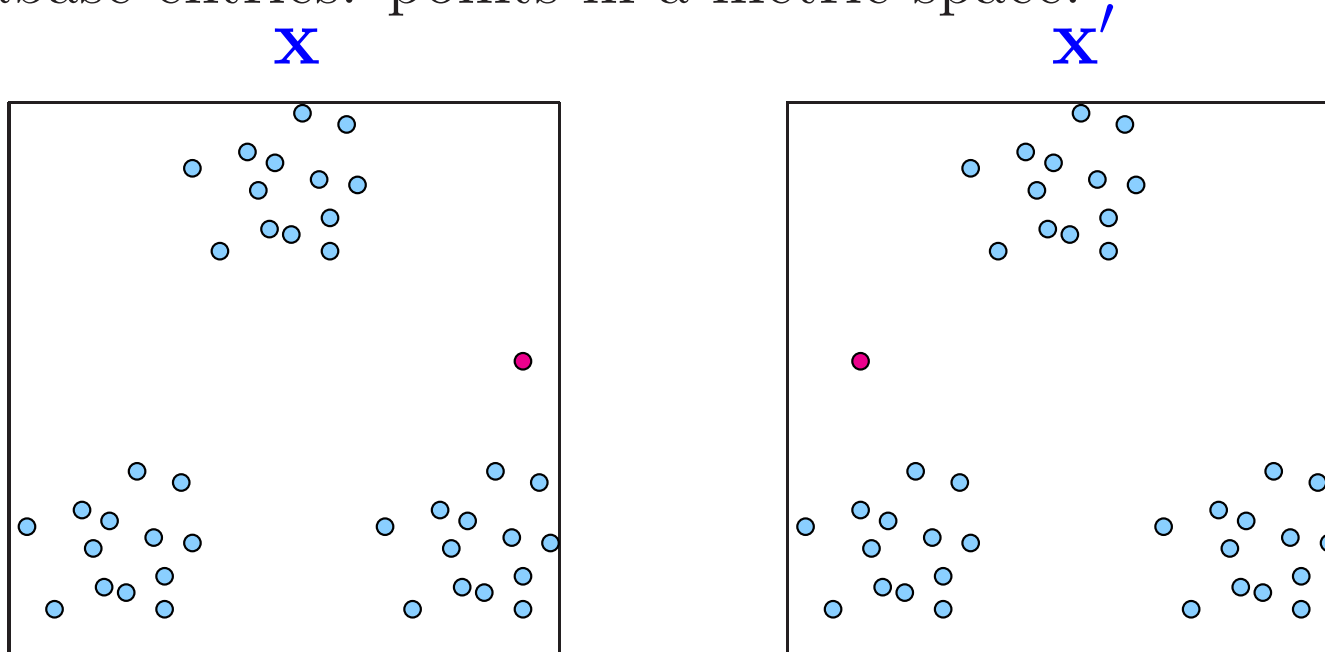
III. Sample-and-aggregate framework

New goal

- Smooth sensitivity framework requires understanding combinatorial structure of f
 - hard in general
- **Goal:** an automatable transformation from an arbitrary f into an ε -indistinguishable A
 - $A(x) \approx f(x)$ for "good" instances x

Example: cluster centers

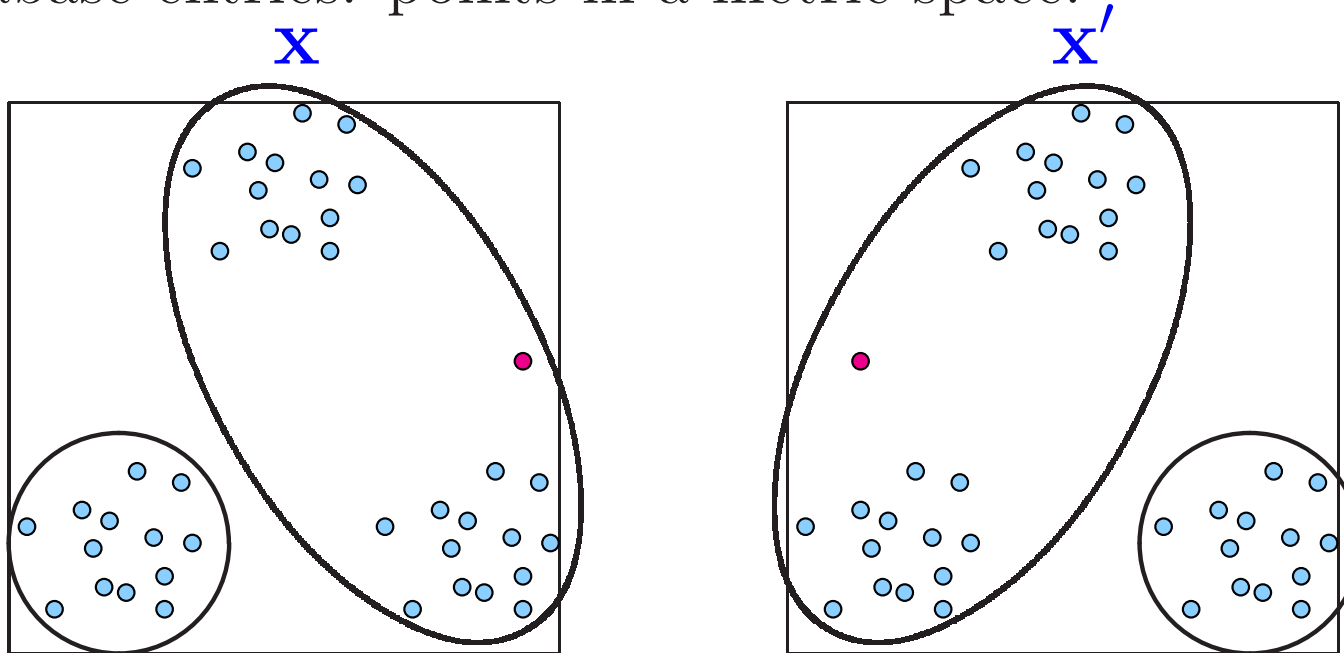
Database entries: points in a metric space.



- 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

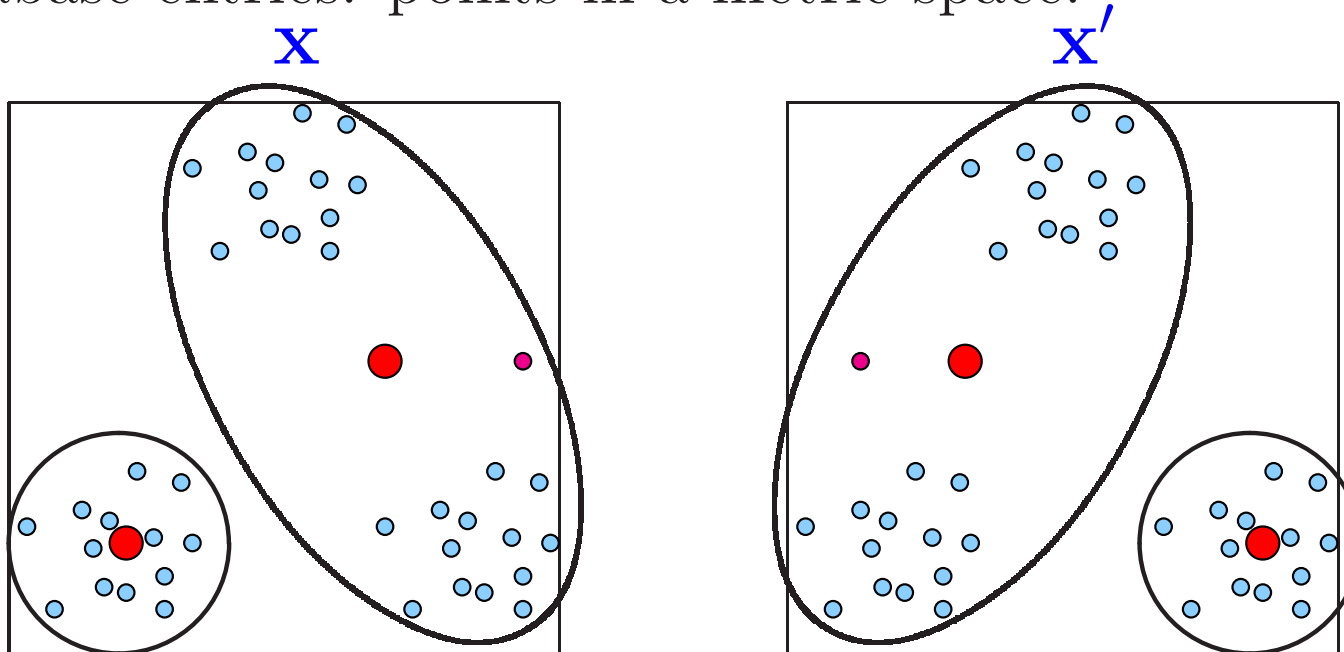
Database entries: points in a metric space.



- 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

Database entries: points in a metric space.

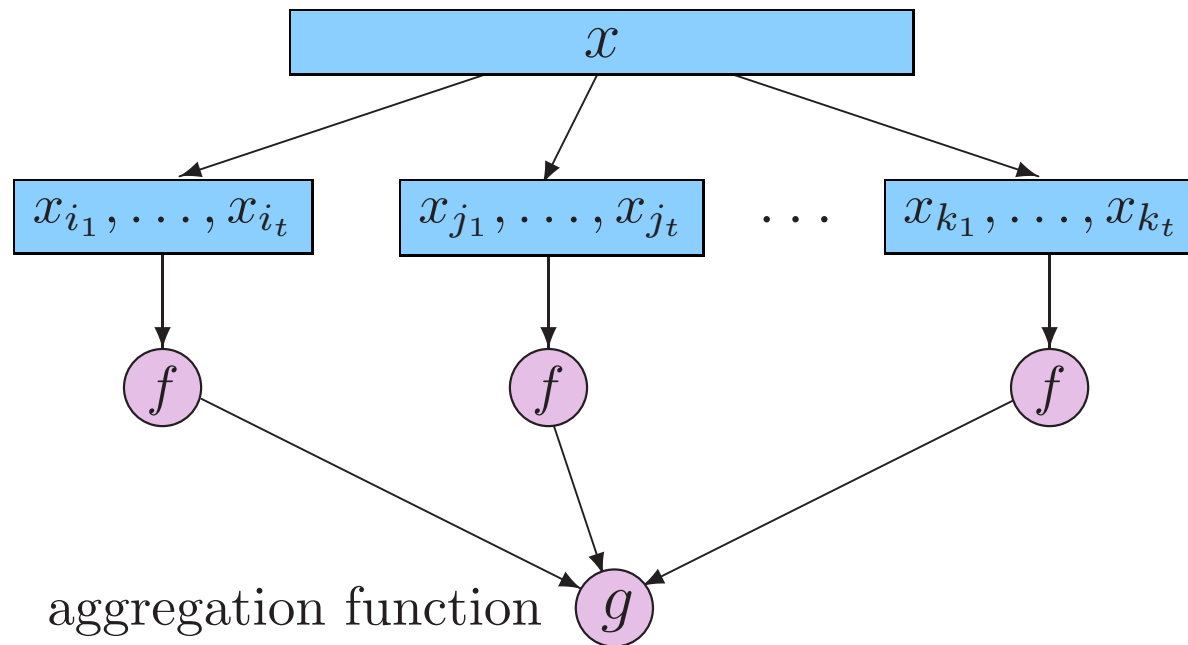


- Comparing sets of centers: **Earthmover-like** metric
- Global sensitivity of cluster centers is roughly the diameter of the space. But intuitively, if clustering is "good", cluster centers should be insensitive.
- No efficient approximation for smooth sensitivity

Sample-and-aggregate framework

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

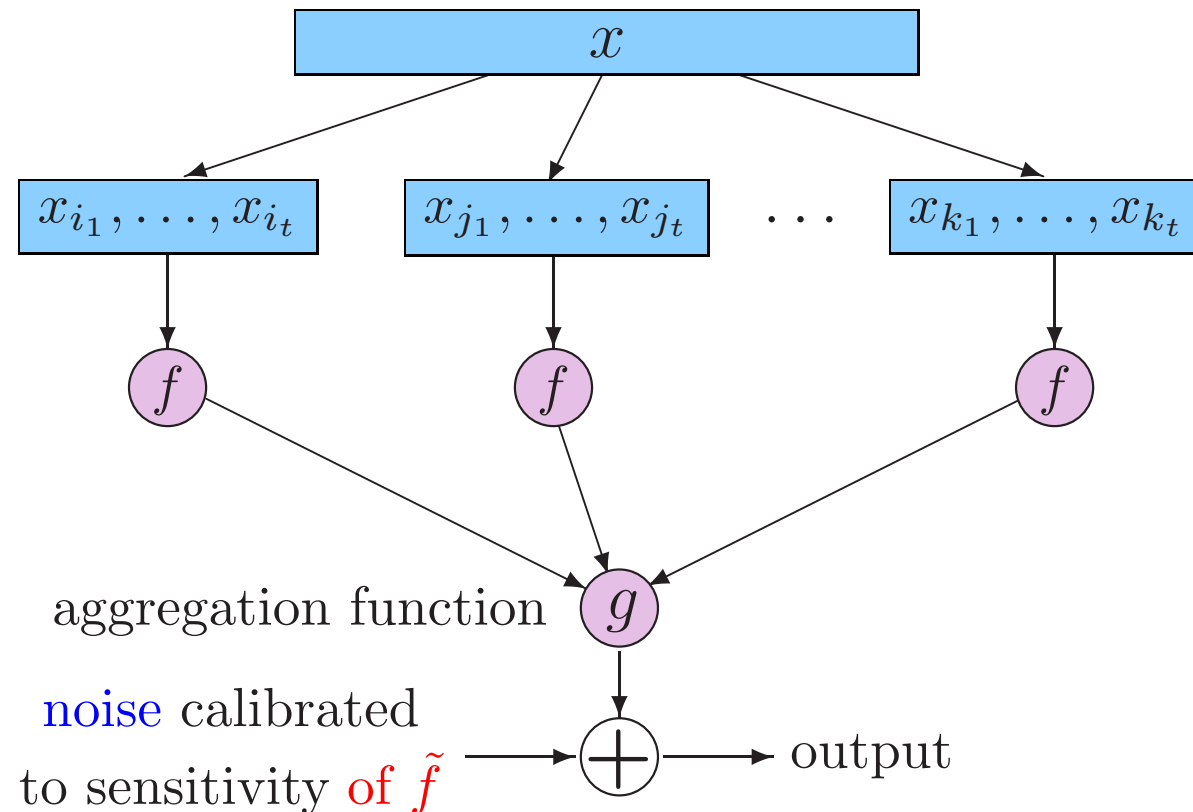
$$\tilde{f}(x) = g(f(\text{sample}_1), f(\text{sample}_2), \dots, f(\text{sample}_s))$$



Sample-and-aggregate framework

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

$$\tilde{f}(x) = g(f(\text{sample}_1), f(\text{sample}_2), \dots, f(\text{sample}_s))$$



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 results

Theorem

*If f can be approximated **on x**
from small samples
then f can be released with little noise*

Sample-and-aggregate results

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

Sample-and-aggregate results

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

- 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 [Vempala Wang 04, Achlioptas McSherry 05])
 - **PAC concepts** (Adam Smith's talk)

Road map

I. Introduction

- Review of global sensitivity framework [DMNS06]
- Motivation

II. Smooth sensitivity framework

III. Sample-and-aggregate framework

Conclusion: fundamental question

Which computations are not too sensitive to individual inputs?

Which functions f admit ε -indistinguishable approximation A ?