

Privacy-Protecting Statistics Computation: Theory and Practice

Rebecca Wright
Stevens Institute of Technology

27 March, 2003

Erosion of Privacy

ì You have zero privacy. Get over it.î

- Scott McNealy, 1999

ì Changes in technology are making privacy harder.

ñ reduced cost for data storage

ñ increased ability to process large amounts of data

ì Especially critical now (given increased need for security-related surveillance and data mining)

Overview

- i Announcements
- i Introduction
- i Privacy-preserving statistics computation
- i Selective private function evaluation

Announcements

- i DIMACS working group on secure efficient extraction of data from multiple datasets. Initial workshop to be scheduled for **Fall 2003**.
- i DIMACS crypto and security tutorials to kick off Special Focus on Communication Security and Information Privacy: **August 4-7, 2003**.
- i NJITES Cybersecurity Symposium, Stevens Institute of Technology, **April 28, 2003**.

What is Privacy?

- i May mean different things to different people
 - ñ seclusion: the desire to be left alone
 - ñ property: the desire to be paid for one's data
 - ñ autonomy: the ability to act freely
- i Generally: the ability to control the dissemination and use of one's personal information.

Different Types of Data

i Transaction data

- ñ created by interaction between stakeholder and enterprise

- ñ current privacy-oriented solutions useful

i Authored data

- ñ created by stakeholder

- ñ digital rights management (DRM) useful

i Sensor data

- ñ stakeholders not clear at time of creation

- ñ growing rapidly

Sensor Data Examples

- i surveillance cameras (especially with face recognition software)
- i desktop monitoring software (e.g. for intrusion or misbehavior detection)
- i GPS transmitters, RFID tags
- i wireless sensors (e.g. for location-based PDA services)

Sensor Data

- i Can be difficult to identify stakeholders and even data collectors
- i Cross boundary between 'real world' and cyberspace
- i Boundary between transaction data and sensor data can be blurry (e.g. Web browsing data)
- i Presents a real and growing privacy threat

Product Design as Policy Decision

- i product decisions by large companies or public organizations become de facto policy decisions
- i often such decisions are made without conscious thought to privacy impacts, and without public discussion
- i this is particularly true in the United States, where there is not much relevant legislation

Example: Metro Cards

Washington, DC

- no record kept of per card transactions
- damaged card can be replaced if printed value still visible

New York City

- transactions recorded by card ID
- damaged card can be replaced if card ID still readable
- have helped find suspects, corroborate alibis

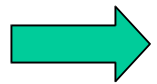
Transactions without Disclosure

→ Don't disclose information in first place!

- i Anonymous digital cash [Chaum et al]
- i Limited-use credit cards [Sha01, RW01]
- i Anonymous web browsing [Crowds, Anonymizer]
- i Secure multiparty computation and other cryptographic protocols
 - ñ perceived (often correctly) as too cumbersome or inefficient to use
 - ñ but, same advances in computing change this

Privacy-Preserving Data Mining

Allow multiple data holders to collaborate to compute important (e.g., security-related) information while protecting the privacy of other information.



Particularly relevant now, with increasing focus on security even at the expense of some privacy.

Advantages of privacy protection

- i protection of personal information
- i protection of proprietary or sensitive information
- i fosters collaboration between different data owners (since they may be more willing to collaborate if they need not reveal their information)

Privacy Tradeoffs?

- i Privacy vs. security: maybe, but doesn't mean giving up one gets the other (who is this person? is this a dangerous person?)
- i Privacy vs. usability: reasonable defaults, easy and extensive customizations, visualization tools

Tradeoffs are to cost or power, rather than inherent conflict with privacy.

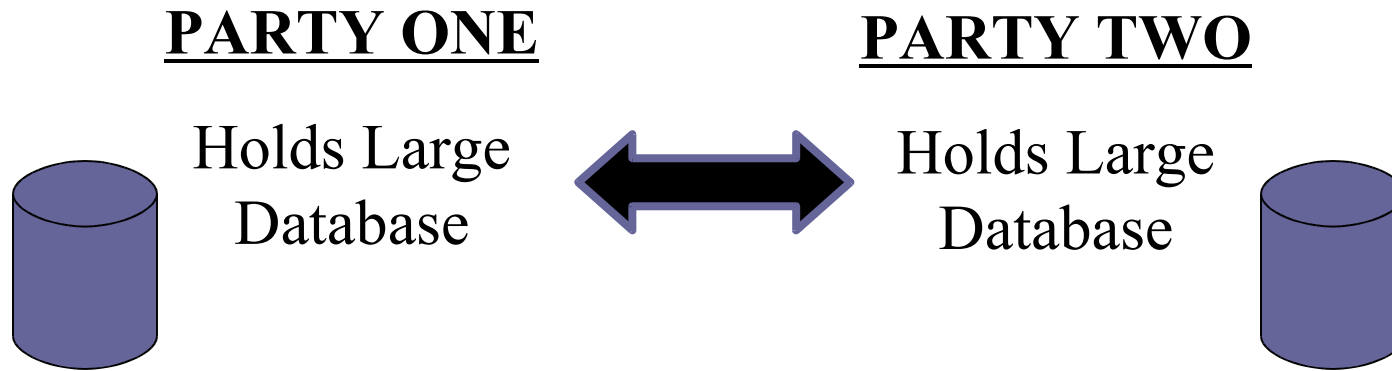
Privacy/Security Tradeoff?

- i **Claim:** No inherent tradeoff between security and privacy, though the cost of having both may be significant.
- i Experimentally evaluate the practical feasibility of strong (cryptographic) privacy-preserving solutions.

Examples

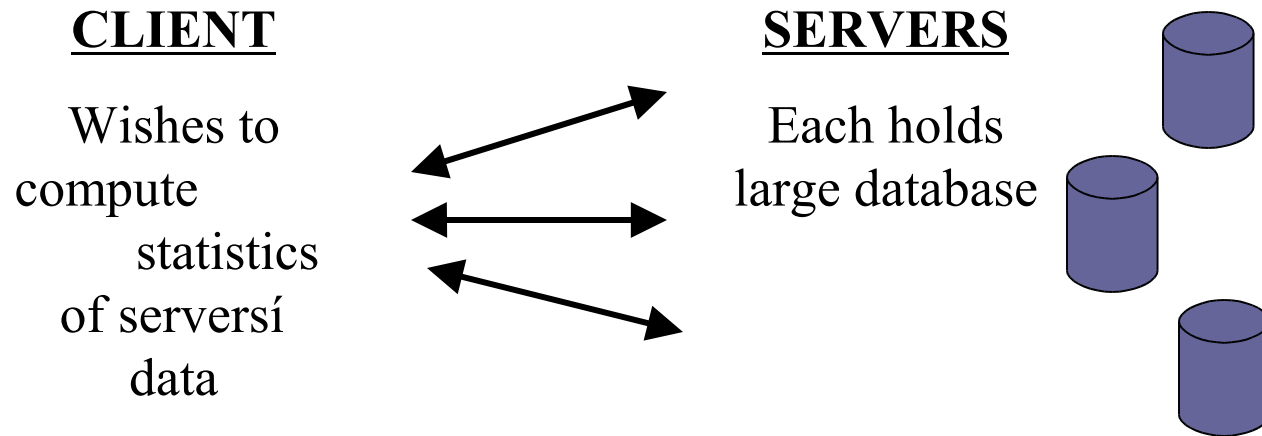
- i Privacy-preserving computation of decision trees [LP00]
- i Secure computation of approximate Hamming distance of two large data sets [FIMNSW01]
- i Privacy-protecting statistical analysis [CIKRRW01]
- i Selective private function evaluation [CIKRRW01]

Similarity of Two Data Sets



- i Parties can efficiently and privately determine whether their data sets are similar
- i Current measure of similarity is approximate Hamming distance [FIMNSW01]
- i Securing other measures is topic for future research

Privacy-Protecting Statistics [CIKRRW01]



i Parties communicate using cryptographic protocols designed so that:

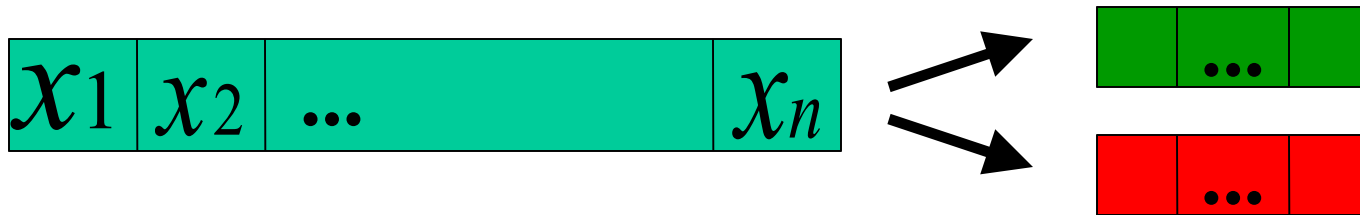
- ñ Client learns desired statistics, but learns nothing else about data (including individual values or partial computations for each database)
- ñ Servers do not learn which fields are queried, or any information about other servers' data
- ñ Computation and communication are very efficient

Privacy Concerns

- i Protect clients from revealing type of sample population, type of specific data used
- i Protect database owners from revealing unnecessary information or providing a higher quality of service than paid for
- i Protect individuals from large-scale dispersal of their personal information

Privacy-Protecting Statistics (single DB)

- i Database contains **public** information (e.g. zip code) and **private** information (e.g. income):



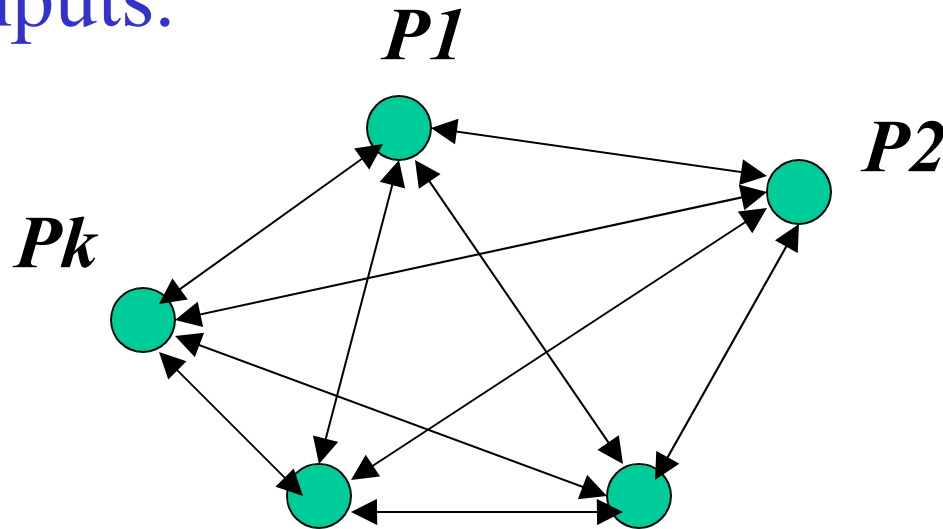
- i Client wants to compute statistics on private data, of subset selected by public data. Doesn't want to reveal selection criteria or private values used.
- i Database wants to reveal only outcome, not personal data.

Non-Private and Inefficient Solutions

- i Database sends client entire database (violates database privacy)
- i For sample size m , use SPIR to learn m values (violates database privacy)
- i Client sends selections to database, database does computation (violates client privacy, doesn't work for multiple databases)
- i general secure multiparty computation (not efficient for large databases)

Secure Multiparty Computation

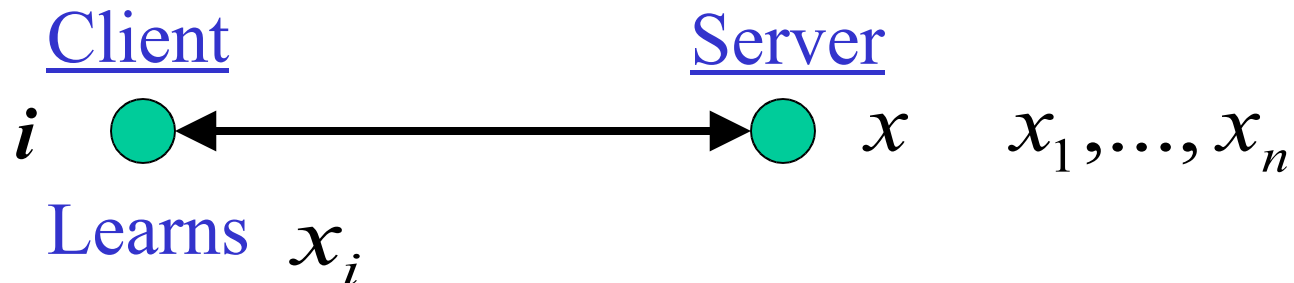
- i Allows k players to privately compute a function f of their inputs.



- i Overhead is polynomial in size of inputs and complexity of f [Yao, GMW, BGW, CCD, ...]

Symmetric Private Information Retrieval

- ii Allows client with input i to interact with database server with input x to learn (only) x_i



- ii Overhead is polylogarithmic in size of database x
[KO,CMS,GIKM]

Homomorphic Encryption

i Certain computations on encrypted messages correspond to other computations on the cleartext messages.

i For additive homomorphic encryption,

$$\tilde{\text{ñ}} \ E(m_1) \text{ i } E(m_2) = E(m_1 + m_2)$$

$$\tilde{\text{ñ}} \ \text{also implies } E(m)^x = E(mx)$$

Privacy-Protecting Statistics Protocol

i To learn mean and variance: enough to learn sum and sum of squares.

i Server stores:

x_1	x_2	...	x_n
-------	-------	-----	-------

$(z_i \quad x_i^2)$

z_1	z_2	...	z_n
-------	-------	-----	-------

and responds to queries from both

i efficient protocol for sum \rightarrow efficient protocol for mean and variance

Weighted Sum

Client wants to compute selected linear combination of m items: $\sum_{j=1}^m w_j x_{i_j}$

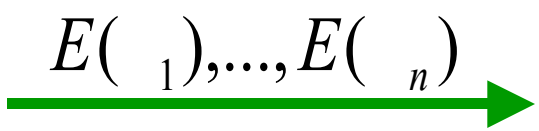
Client

Server

Homomorphic encryption E, D

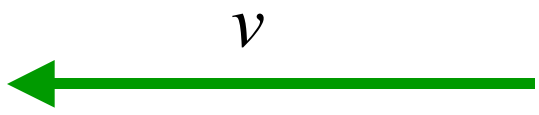
w_j if $i = i_j$
 0 o/w

computes



$$v = \sum_{i=1}^n (E(x_i))^{w_i}$$

decrypts to obtain



$$E\left(\sum_{i=1}^n w_i x_i\right)$$

$$\sum_{i=1}^n x_i \quad \sum_{j=1}^m w_j x_{i_j}$$

Efficiency

- i Linear communication and computation (feasible in many cases)
- i If n is large and m is small, would like to do better

Selective Private Function Evaluation

- i Allows client to privately compute a function f over m inputs x_{i_1}, \dots, x_{i_m}
- i client learns **only** $f(x_{i_1}, \dots, x_{i_m})$
- i server does **not** learn i_1, \dots, i_m

Unlike general secure multiparty computation, we want communication complexity to depend on m , not n . (More accurately, polynomial in m , polylogarithmic in n).

Security Properties

- i **Correctness**: If client and server follow the protocol, client's output is correct.
- i **Client privacy**: malicious server does not learn client's input selection.
- i **Database privacy**:
 - ñ **weak**: malicious client learns no more than output of some m -input function g
 - ñ **strong**: malicious client learns no more than output of specified function f

Solutions based on MPC

i Input selection phase:

ñ server obtains blinded version of each x_{i_j}

i Function evaluation phase

ñ client and server use MPC to compute f on the m blinded items

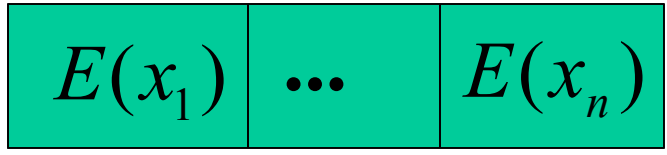
Input selection phase

Client

Server

Homomorphic encryption D, E
Computes encrypted database

Retrieves
 $E(x_{i_1}), \dots, E(x_{i_m})$
using SPIR



Picks random
 c_1, \dots, c_m
computes
 $E(x_{i_j} \parallel c_j)$



Decrypts received values:
 $s_j \parallel x_{i_j} \parallel c_j$

Function Evaluation Phase

i Client has $\mathcal{C} = \mathcal{C}_1, \dots, \mathcal{C}_m$

i Server has $\mathcal{S} = \mathcal{S}_1, \dots, \mathcal{S}_m$ $\mathcal{S}_j = x_{i_j} \cdot \mathcal{C}_j$

Use MPC to compute:

$$g(\mathcal{C}, \mathcal{S}) = f(\mathcal{S}, \mathcal{C}) = f(x_1, \dots, x_m)$$

ii Total communication cost polylogarithmic in n , polynomial in $m, |f|$

Distributed Databases

- i Same approach works to compute function over distributed databases.
 - ñ Input selection phase done in parallel with each database server
 - ñ Function evaluation phase done as single MPC
 - ñ only final outcome is revealed to client.

Performance

	<i>Complexity</i>	<i>Security</i>
1	$m \text{ SPIR}(n,1,k) + O(k f)$	Strong
2	$m \text{ SPIR}(n,1,1) + \text{MPC}(m, f)$	Weak
3	$\text{SPIR}(n,m,\log n) + \text{MPC}(m, f) + km^2$	Weak
4	$\text{SPIR}(n,m,k) + \text{MPC}(m, f)$	Honest client only

→ Current experimentation to understand whether these methods are efficient in real-world settings.

Conclusions

- i Privacy is in danger, but some important progress has been made.

- i Important challenges ahead:
 - ñ Usable privacy solutions
 - ñ Sensor data

 - ñ better use of hybrid approach: decide what can safely be disclosed, use cryptographic protocols to protect critical information, weaker and more efficient solutions for the rest

- i Technology, policy, and education must work together.