

# A Theory of Pricing Private Data

Dan Suciu – U. of Washington

Joint work with: Chao Li, Daniel  
Yang Li, Gerome Miklau

# Motivation

- Private data has value
  - A unique user: \$4 at FB, \$24 at Google [JPMorgan]
- Today's common practice:
  - Companies profit from private data without compensating users
- New trend: allow users to profit financially
  - Industry: personal data locker  
<https://www.personal.com/> , <http://lockerproject.org/>
  - Academia: mechanisms for selling private data [Ghosh11,Gkatzelis12,Aperjis11,Roth12,Riederer12]

# Overview

**This talk:** framework for pricing queries on private data

- **Data owners:** sell their private data
- **Buyer:** buys a query (many buyers, many queries!)
- **Trusted market maker:** facilitates transactions

**What I will address:**

- Consistent prices for arbitrary queries
- Fair compensation of data owners for privacy loss

**What I will not address:**

- Designing truthful, efficient mechanisms
- Prices/payments: at the discretion of market maker

# Challenges

**Perturbation:** is a cost savings mechanism for buyer  
**Price:** computed for each (query, perturbation) pair.

Two extremes:

- **No perturbation**
  - Query returns **raw** data
  - Data owner compensated the full price of data; e.g. **\$10**
  - Buyer pays a high price
- **High perturbation**
  - Query is  **$\epsilon$ -Differentially Private**, for small  $\epsilon$
  - Data owner compensated a tiny price, e.g. **\$0.001**
  - Buyer pays modest price


# Related Work


- Query-based data pricing, Koutris, Upadhyaya, Balazinska, Howe, Suciu, 2012
- Pricing Aggregate Queries in a Data Marketplace, Li and Miklau, 2012
- Selling privacy at auction, Ghosh, A., Roth, A. 2011
- Pricing Private Data, Gkatzelis, Aperjis, Huberman, 2012
- A Market for Unbiased Private Data, Aperjis, Huberman 2011
- Buying Private Data at Auction (...), Roth 2012
- For sale : Your Data By : You, Riederer, Erramilli, Chaintreau, Krishnamurthy, Rodriguez, 2012

# Outline

- Problem Statement
- The Buyer's price:  $\pi$
- Balanced Pricing Framework
- Conclusions

# Main Concepts

- **Database**  $\mathbf{x} = (x_1, \dots, x_n)$ 
  - $x_i = \text{value}$ , owned by some *owner*
- **Buyer's request:**  $\mathbf{Q} = (\mathbf{q}, v)$ 
  - $\mathbf{q} = (q_1, \dots, q_n) = \text{query}$ ;  $\mathbf{q}(\mathbf{x}) = \sum_i q_i x_i$
  - $v = \text{variance}$
- **Randomized answer:**  $\mathcal{K}(\mathbf{x})$ 
  - $E[\mathcal{K}(\mathbf{x})] = \mathbf{q}(\mathbf{x})$ ,  $\text{Var}[\mathcal{K}(\mathbf{x})] \leq v$

Buyer pays  $\pi(\mathbf{Q})$
- **Privacy loss:**
  - $\epsilon_i(\mathcal{K})$  [Ghosh'11]
  - $W(\epsilon_i)$  = its value to the owner

Owner receives  $\mu_i(\mathbf{Q})$

# Example (1/3)

**Data:** 1000 data owners rate two candidates A, B between 0..5:

- Owner 1:  $x_1, x_2$
- Owner 2:  $x_3, x_4$
- ...
- Owner 1000:  $x_{1999}, x_{2000}$

**Price:** \$10 for each raw item  $x_i$

- **Buyer:**
  - Compute rating for candidate A:  $x_1 + x_3 + \dots + x_{1999}$
  - $\mathbf{q} = (1, 0, 1, 0, \dots)$ ,  $v=0$  (raw data)
- **$\mu$ -Payments:** \$10/item
- **Buyer's Price  $\pi$ :** \$10,000

1. Raw data is too expensive!



# Example (2/3)

**Data:** 1000 data owners rate two candidates A, B between 0..5:

- Owner 1:  $x_1, x_2$
- Owner 2:  $x_3, x_4$
- ...
- Owner 1000:  $x_{1999}, x_{2000}$

**Price:** \$10 for each raw item  $x_i$

- **Buyer:**
  - Can tolerate error  $\pm 300$
  - $\mathbf{q} = (1, 0, 1, 0, \dots)$ ,  $v=0$   $v = 2500^*$  ( $v = \sigma^2 = \text{variance}$ )
- **$\mu$ -Payments:** ~~\$10/item~~ **\$0.001/item** (query is 0.1-DP<sup>\*\*</sup>)
- **Buyer's Price  $\pi$ :** ~~\$10,000~~ **\$1**

\*Probability(error  $< 6\sigma$ )  $> 1/6^2 = 97\%$

\*\*  $\epsilon = \text{Sensitivity}(\mathbf{q})/\sigma = 5/\sigma = 0.1$

2. Perturbed data  
is cheaper.

# Example (3/3)

**Data:** 1000 data owners rate two candidates A, B between 0..5:

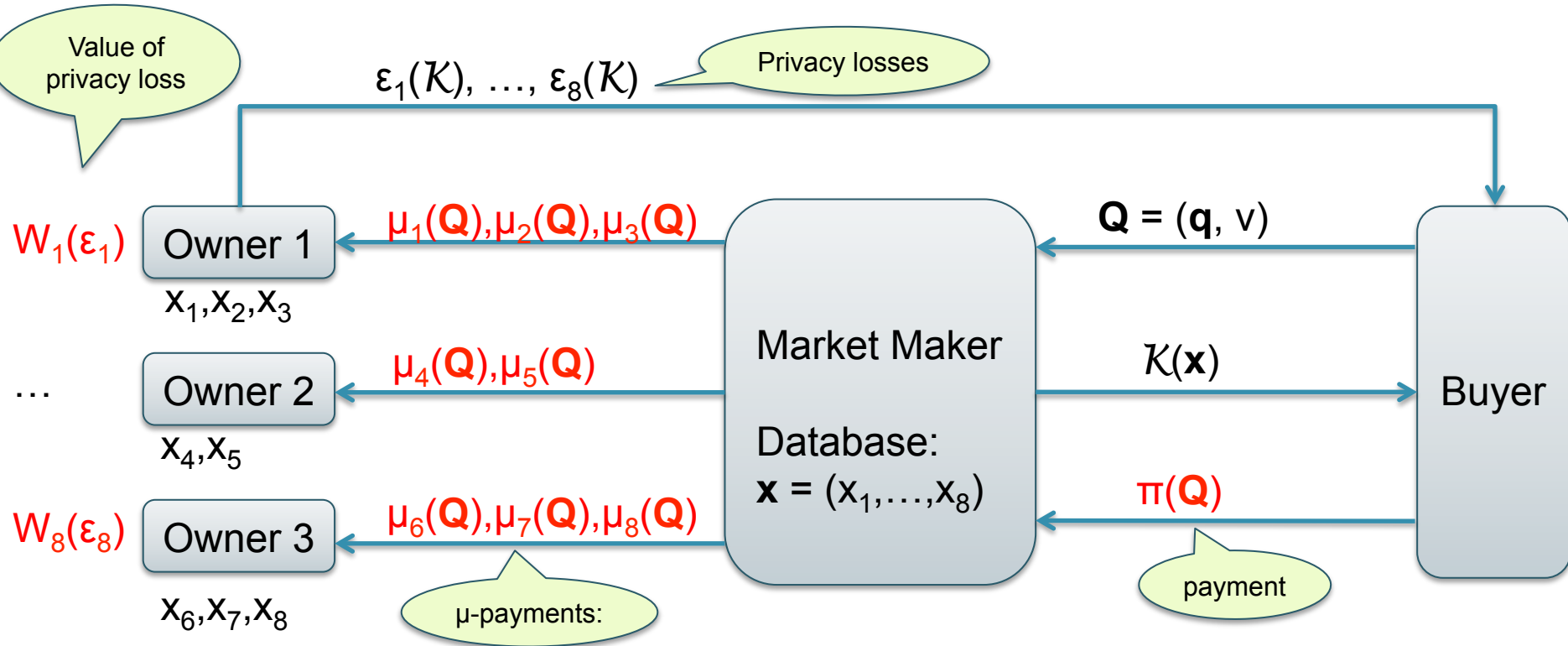
- Owner 1:  $x_1, x_2$
- Owner 2:  $x_3, x_4$
- ...
- Owner 1000:  $x_{1999}, x_{2000}$

**Price:** \$10 for each raw item  $x_i$

- **Another buyer:**
  - $\mathbf{q} = (1,0,1,0,\dots)$ , ~~variance = 0~~, ~~variance = 2500~~ variance = 500
- **$\mu$ -Payments:** ~~\$10/item, \$0.001/item~~ \$0.1/item? \$1/item?
- **Buyer's Price  $\pi$ :** ~~\$10000, \$1~~ \$100? \$1000?
- **Buyer will refuse to pay more than \$5!**
  - Instead purchases 5 times variance=2500, for \$5, takes avg.

3. Multiple queries: must be consistent, compensate owners for privacy loss.

# Pricing Framework

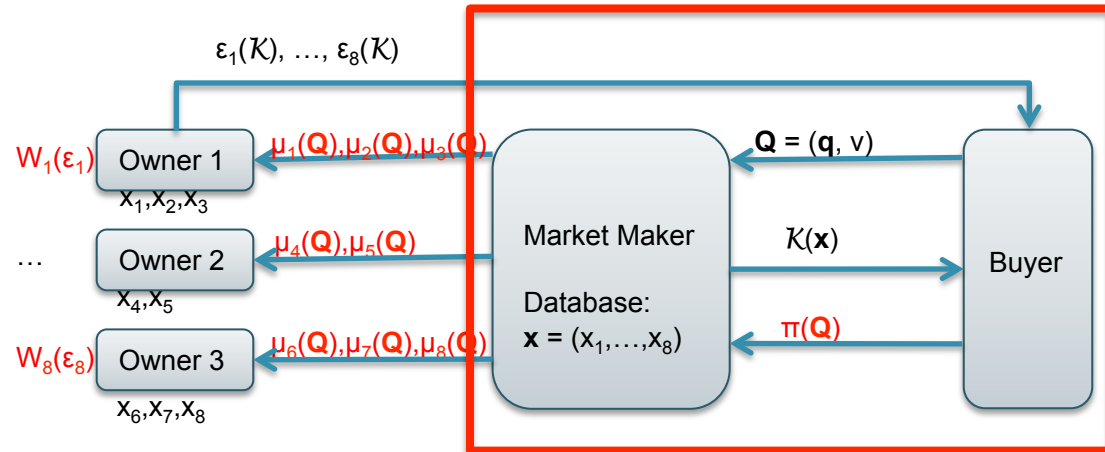


Market maker needs to **balance** the pricing framework

- Satisfy the buyer: use  $\mathcal{K}$  to answer  $\mathbf{Q}$ , charge him  $\pi(\mathbf{Q})$
- Satisfy the owner: pay her  $\mu_i(\mathbf{Q}) \geq W_i(\epsilon_i)$
- Recover cost:  $\mu_1 + \dots + \mu_n \leq \pi$

# Outline

- Problem Statement
- The Buyer's price:  $\pi$
- Balanced Pricing Framework
- Conclusions



# Designing a Pricing Function

For any query/variance request  $Q = (q, v)$

define a price:  $\pi(Q) \in [0, \infty]$

What can go wrong?

# Arbitrage!

## Def.

- $Q=(q, v)$  is *answerable* from  $Q_1, \dots, Q_k$  ( $=(\mathbf{q}_1v_1), \dots, (\mathbf{q}_kv_k)$ ) if there exists a function  $f$  s.t. whenever  $K_1, \dots, K_k$  answer  $Q_1, \dots, Q_k$ ,  $f(K_1, \dots, K_k)$  answers  $Q$
- $Q$  is *linearly answerable* from  $Q_1, \dots, Q_k$  if  $f$  is a linear function;  
notation:  $Q_1, \dots, Q_k \rightarrow Q$

**Examples:**  $(\mathbf{q}_1, v_1), (\mathbf{q}_2, v_2), (\mathbf{q}_3, v_3) \rightarrow (\mathbf{q}_1 + \mathbf{q}_2 + \mathbf{q}_3, v_1 + v_2 + v_3)$

$(\mathbf{q}, v) \rightarrow (c \mathbf{q}, c^2 v)$

$(\mathbf{q}, v), (\mathbf{q}, v), (\mathbf{q}, v), (\mathbf{q}, v), (\mathbf{q}, v) \rightarrow (\mathbf{q}, v/5)$

**Def.** *Arbitrage* happens when  $Q_1, \dots, Q_k \rightarrow Q$  and  $\pi(Q_1) + \dots + \pi(Q_k) < \pi(Q)$

**Example:** If  $5 \times \pi(\mathbf{q}, v) < \pi(\mathbf{q}, v/5)$ , then we have arbitrage

# Arbitrage-Free Pricing

**Def.** The pricing function  $\pi$  is *Arbitrage-Free* if:  
 $Q_1, \dots, Q_k \rightarrow Q$  implies  $\pi(Q_1) + \dots + \pi(Q_k) \geq \pi(Q)$

Do AF-pricing functions exist?

Remark: AF generalizes the following known property of  $\varepsilon$ -DP:

If  $Q_1$  is  $\varepsilon$ -DP, and  $Q = f(Q_1)$ , then  $Q$  is also  $\varepsilon$ -DP

Indeed: if  $\pi(Q_1) \leq \$0.001$  then  $\pi(Q) \leq \$0.001$

# Designing Arbitrage-Free Pricing Functions

$$\pi(\mathbf{q}, v) = (q_1^2 + q_2^2 + \dots + q_n^2) / v \quad \text{is AF}$$

Price of raw data  $\pi(\mathbf{q}, 0) = \infty$

More generally:

$$\pi(\mathbf{q}, v) = \|\mathbf{q}\|^2 / v \quad \text{is AF, where } \|\mathbf{q}\| \text{ is any } \underline{\text{semi-norm}}$$

$$\pi(\mathbf{q}, v) = 20,000 / 3.14 \times \arctan[(q_1^2 + q_2^2 + \dots + q_n^2) / v]$$

Price of raw data  $\pi(\mathbf{q}, 0) = 10,000$

More generally:

If  $f$  is sub-additive, non-decreasing and  $\pi_1, \dots, \pi_k$  are AF then  $\pi = f(\pi_1, \dots, \pi_k)$  is AF

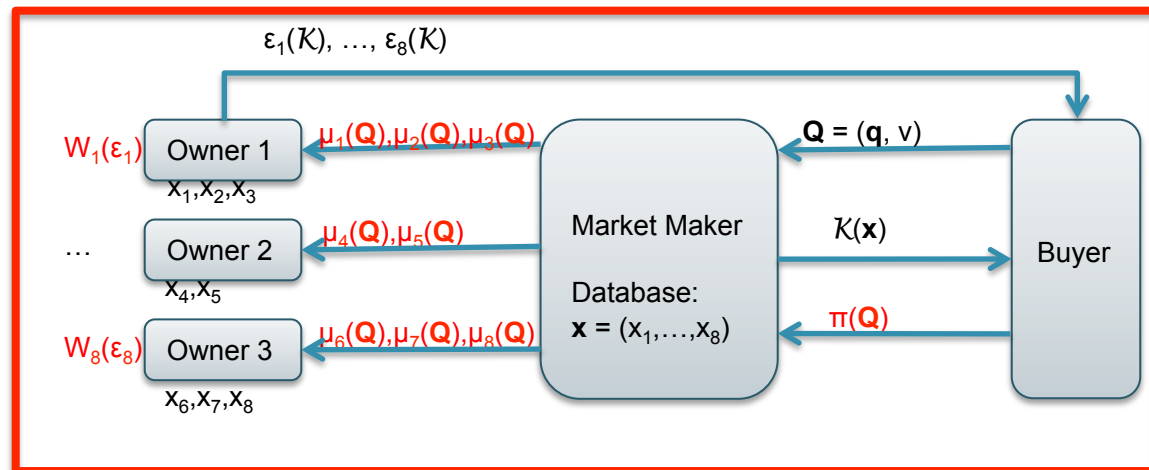


# Discussion

- Query answerability is well studied for relational queries (**no noise!**) [Nash'2010]
  - Checking answerability: NP ... undecidable
- New for linear queries **with noise**:
  - Checking linear answerability is in PTIME
  - Checking general answerability is open

# Outline

- Problem Statement
- The Buyer's price:  $\pi$
- **Balanced Pricing Framework**
- Conclusions



# The Perspective of the Data Owner

- Micropayment to owner  $i$ :  
 $\mu_i(\mathbf{Q})$  = what the market maker pays her

- Must compensate for her privacy loss: [Ghosh'11]

$$\varepsilon_i(\mathcal{K}) = \sup_{S, \mathbf{x}} \left| \log \frac{\Pr(\mathcal{K}(\mathbf{x}) \in S)}{\Pr(\mathcal{K}(\mathbf{x}^{(i)}) \in S)} \right|$$

$W_i(\varepsilon_i)$  = the owner's value for the privacy loss

$W_i(\infty)$  = price for her raw data; e.g. = \$10

# Properties of $\mu_i$

**Assumptions:** the pricing framework is defined by  $\mu_i$ ,  $W_i$ , plus:

- $\mathcal{K}$  = Laplacian answering mechanism:  
$$\mathcal{K}(\mathbf{x}) = \mathbf{q}(\mathbf{x}) + \text{Lap}(\text{sqrt}(v/2))$$
- $\pi = a(\mu_1 + \dots + \mu_n) + b$ , for some  $a \geq 1$ ,  $b \geq 0$

$\varepsilon_i(\mathcal{K})$  derived  
from sensitivity

market maker  
recovers the costs

**Def.** The pricing framework is ***balanced*** if is

- (1)  $\mu_i$  is arbitrage free,
- (2) compensates owner:  $\mu_i(\mathbf{Q}) \geq W_i(\varepsilon_i(\mathcal{K}))$
- (3) is fair:  $q_i = 0$  implies  $\mu_i(\mathbf{q}, v) = 0$

Market maker must design a ***balanced*** pricing framework

# Designing Balanced Pricing Frameworks

The pricing-frameworks below are **balanced** (assume  $x_i \in [0,5]$ )

$$\begin{aligned}\mu_i(\mathbf{q}, v) &= 5c_i |q_i| / \sqrt{v/2} \\ W_i(\varepsilon_i) &= c_i \varepsilon_i\end{aligned}$$

$c_i$  is any constant

Price of raw data:  
 $\mu_i(\mathbf{q}, 0) = W_i(\infty) = \infty$

$$\begin{aligned}\mu_i(\mathbf{q}, v) &= 20 / 3.14 \times \arctan(5c_i |q_i| / \sqrt{v/2}) \\ W_i(\varepsilon_i) &= 20 / 3.14 \times \arctan(c_i \varepsilon_i)\end{aligned}$$

Raw data:  
 $\mu_i(\mathbf{q}, 0) = W_i(\infty) = \$10$

More generally:

If  $\mu_{i1}, \dots, \mu_{ik}$  and  $W_{i1}, \dots, W_{ik}$  are balanced and  $f_i$  is non-decreasing, subadditive then  $\mu_i = f(\mu_{i1}, \dots, \mu_{ik})$ ,  $W_i = f(W_{i1}, \dots, W_{ik})$  are balanced

# Finding Out the Owner's Valuation $W_i$

Mechanisms proposed [Ghosh'11, Gkatzelis'12, Riederer'12]

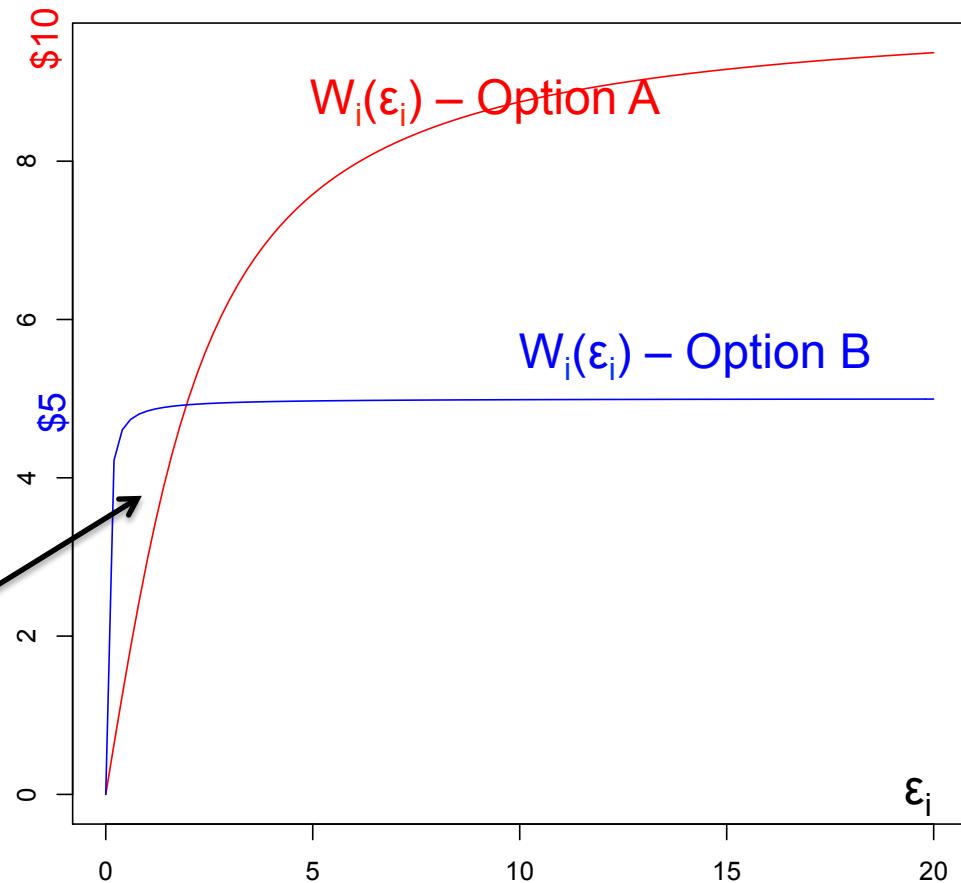
We use an idea from [Aperjis&Huberman'11]:

Market Maker

gives users 3 options

- **Option A**: risk neutral
- **Option B**: risk averse
- **Option C**: opt-out

“Typical” query has  
small privacy loss



# Outline

- Problem Statement
- The Buyer's price:  $\pi$
- Balanced Pricing Framework
- Conclusions

# Conclusions

- **The Contract in differential-privacy:**
  - *Privacy loss  $\epsilon_i$  = bounded by a fixed, small  $\epsilon$*
  - **Privacy budget** (defined by  $\epsilon$ ) = limit on the number of queries
- **The Contract in private data markets:**
  - *Privacy loss  $\epsilon_i$  = arbitrary; compensated by micro-payment  $\mu_i$*
  - **Cash-and-carry** = unlimited queries
- **Special case 1:** Answer contains raw data
- **Special case 2:** Answer is  $\epsilon$ -DP
- **Challenge:** Designing a balanced pricing framework