



Privacy Issues in Data Publishing

- ▶ Governments and organizations publish anonymous personal information for research, analytics and services

	ZIP Code	Age	Salary	Disease
1	476**	2*	3K	gastric ulcer
2	476**	2*	4K	gastritis
3	476**	2*	5K	stomach cancer
4	4790*	≥ 40	6K	gastritis
5	4790*	≥ 40	11K	flu
6	4790*	≥ 40	8K	bronchitis
7	476**	3*	7K	bronchitis
8	476**	3*	9K	pneumonia
9	476**	3*	10K	stomach cancer

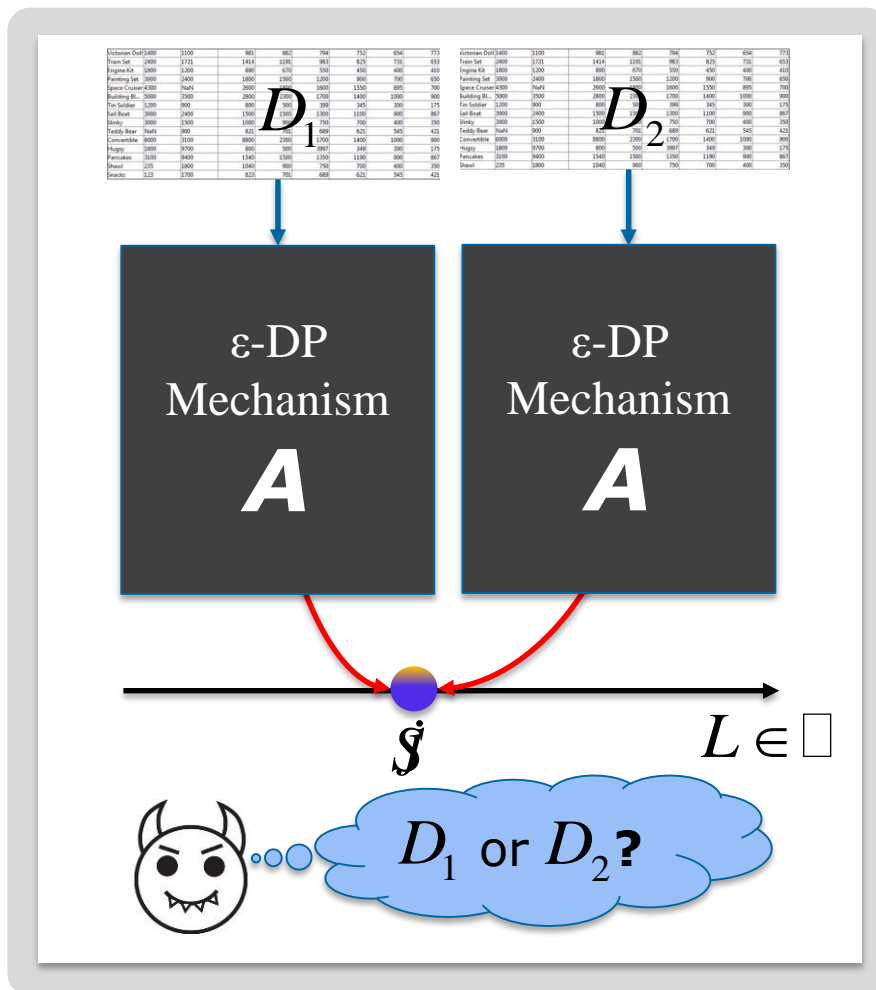
The screenshot shows the 'Movies You've Rated' section on Netflix. It features a list of movies with columns for 'TITLE', 'MPAA', 'GENRE', and 'STAR RATING'. Each movie entry includes an 'Add' button, the movie title and year, the MPAA rating, the genre, a star rating (5 stars), and a 'Clear Rating' button. The movies listed are: 12 Angry Men (1957), The 39 Steps (1935), An American in Paris (1951), The Andromeda Strain (1971), Apollo 13 (1995), The Battle of Algiers (1965), Being There (1979), Big Deal on Madonna Street (1958), The Birds (1963), and Blade Runner (1982).

www.netflix.com

▶ Privacy leak

- identify a person from internet databases
 - de-anonymize Netflix Price dataset [A. Narayanan '08]
- discover an individual's record by comparing databases
 - your record was not in the database last month, but now it is...

Differential Privacy (DP)



▶ ϵ -Differential Privacy [C. Dwork '06]

$$e^{-\epsilon} \leq \frac{P_r[A(D_1) \in S]}{P_r[A(D_2) \in S]} \leq e^{\epsilon}$$

– privacy \rightarrow information loss

▶ ϵ -DP Mechanisms

– DP Noise-adding mechanisms

- *Laplacian, Geometric*

– other DP mechanisms

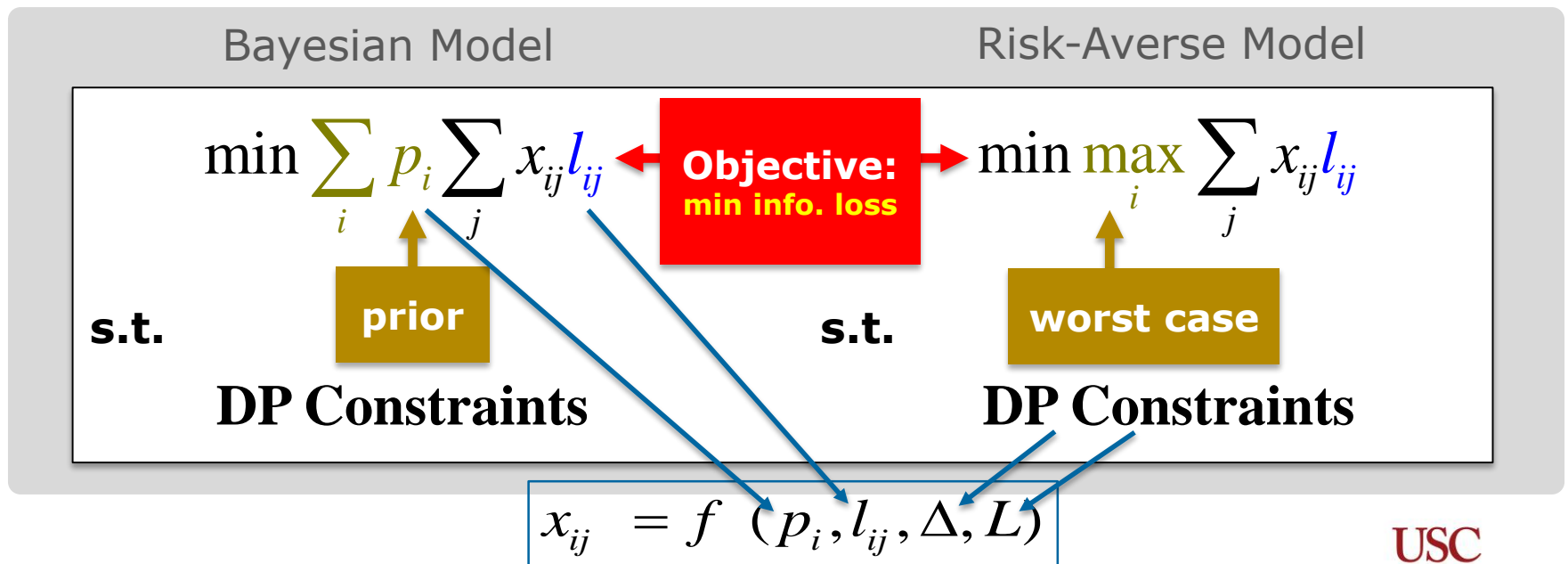
- *Matrix* [C. Li '10], *K-norm* [M. Hardt '09]

– non-numeric DP mechanism

- *Exponential* [F. McSherry '07]

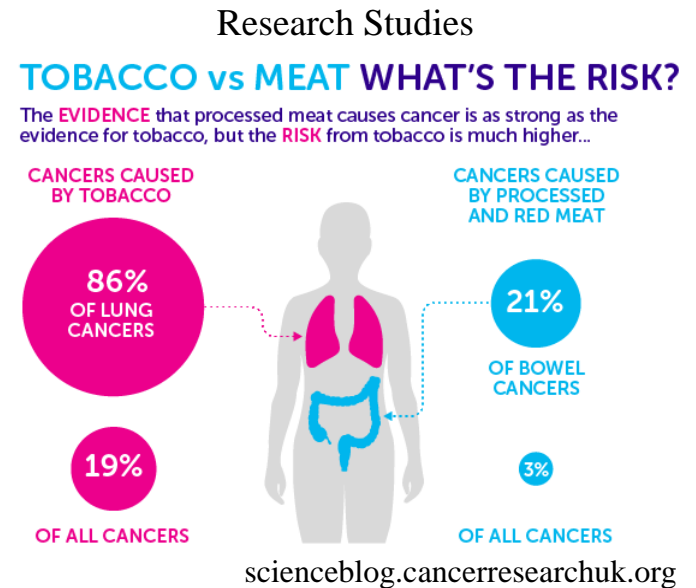
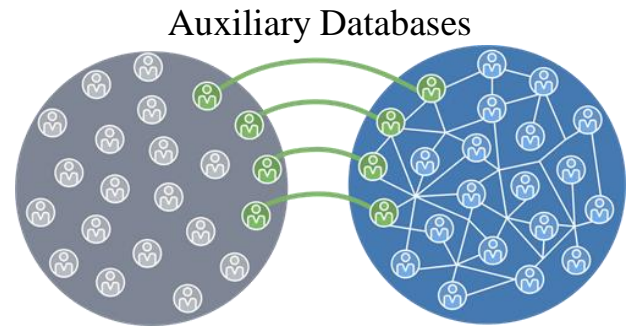
Optimal DP Mechanism

- Widely-used information loss function: $l_{ij} = l_{|j-i|}$
- A DP mechanism is called *optimal* if it **minimizes information loss** and preserves DP.
- Data managers solve the optimization problem for mechanism x_{ij}



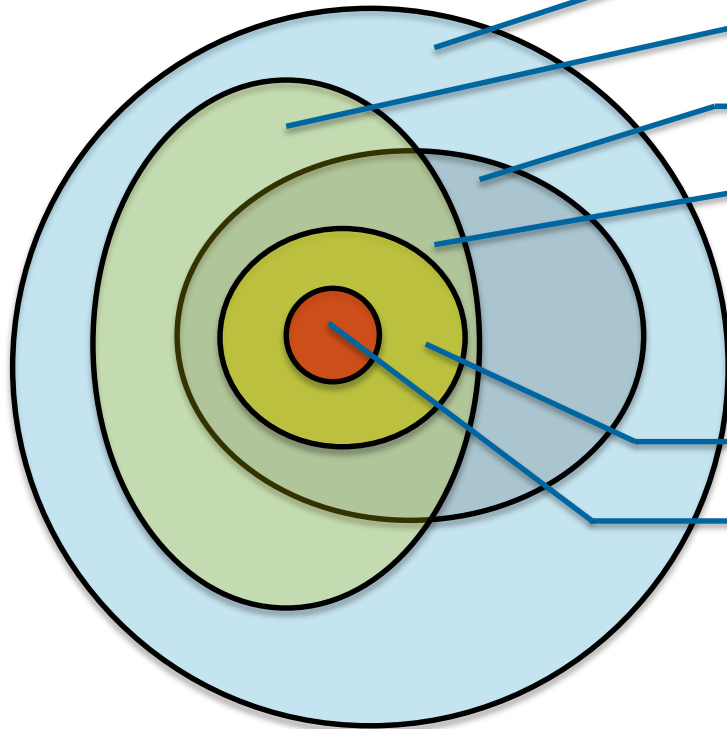
Presence of Side-Information

- ▶ Side-information exists everywhere...
 - auxiliary databases
 - research studies, common knowledge
 - mathematical theories
 - central limit theorem
 - transformations of random variables
- ▶ The presence of side-information is important and cannot be neglected.
- ▶ Side-information → Prior probability



State-of-the-Art and Open Questions

solution space = $(p_i, L, \Delta, l_{ij}) = l_{|j-i|}$



Optimal DP Mechanism (Bayesian)

Optimal DP Mechanism (Risk-Averse)

Optimal DP Mechanism (Bayesian, $\Delta = \Delta_{GS}$)

Optimal DP Mechanism (Risk-Averse, $\Delta = \Delta_{GS}$)

:Staircase Mechanism [Q. Geng '14]

Optimal in **Risk-Averse** model

Optimal for **unbounded domain L**

Universally Optimal DP Mechanism (**unknown**)

Universally Optimal DP Mechanism ($L \in \mathbb{Z}, \Delta=1$)

:Geometric Mechanism [M. Gupte '10] [A. Ghosh '12]

Universally optimal in both **Risk-Averse**
and **Bayesian** model

- ▶ A universally optimal mechanism is optimal for all priors p_i and all loss functions l_{ij} .

Main Contributions

Propose **open questions** in DP mechanism design

For (Bayesian, $\Delta = \Delta_{GS}$), we propose a heuristic design

- optimal design for general priors is difficult
- we start with heuristic design, and it surprisingly leads to significant improvement in utility-privacy tradeoffs

Show via experiments, the importance of the optimal Bayesian mechanism design

- optimal Bayesian design is **non-trivial** when side-information substantially narrows down the outputs of the query

Experimental Context and Settings

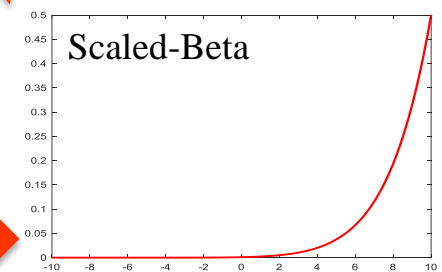
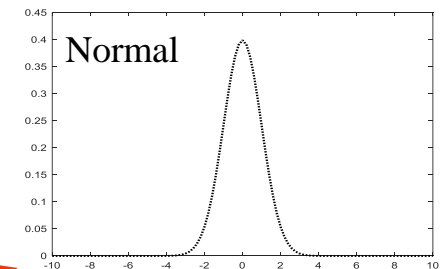
- ▶ Queries - **Mean** and **Max**
- ▶ Oblivious mechanism → database independent → synthetic data

– public information

- known: domain $L \in [-10,10]$
- each entity is independent and uniformly distributed

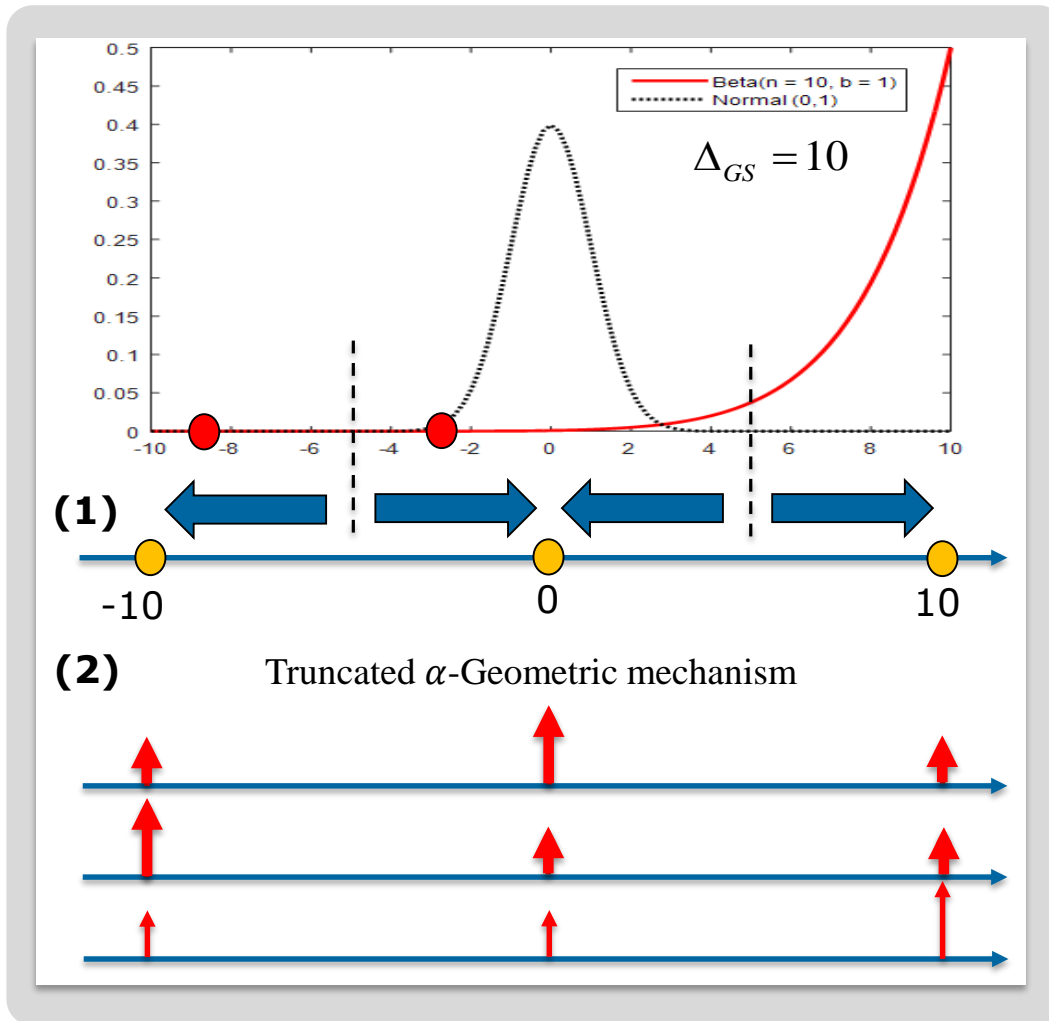
– mathematical theories

- central limit theorem
→ **mean** value is approximately *normal* distributed
- transformations of random variables
→ the **max** value is scaled-*beta* distributed over L



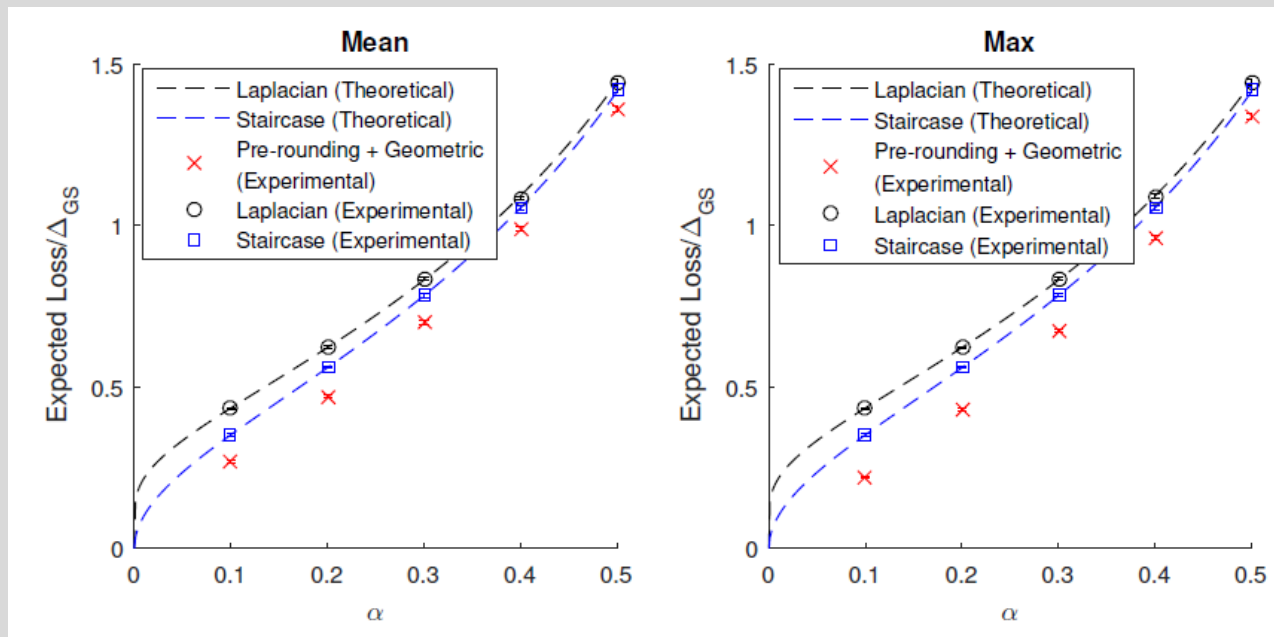
- ▶ Global sensitivity = 10
- ▶ Gaussian is truncated and normalized in L

Our Heuristic DP Mechanism



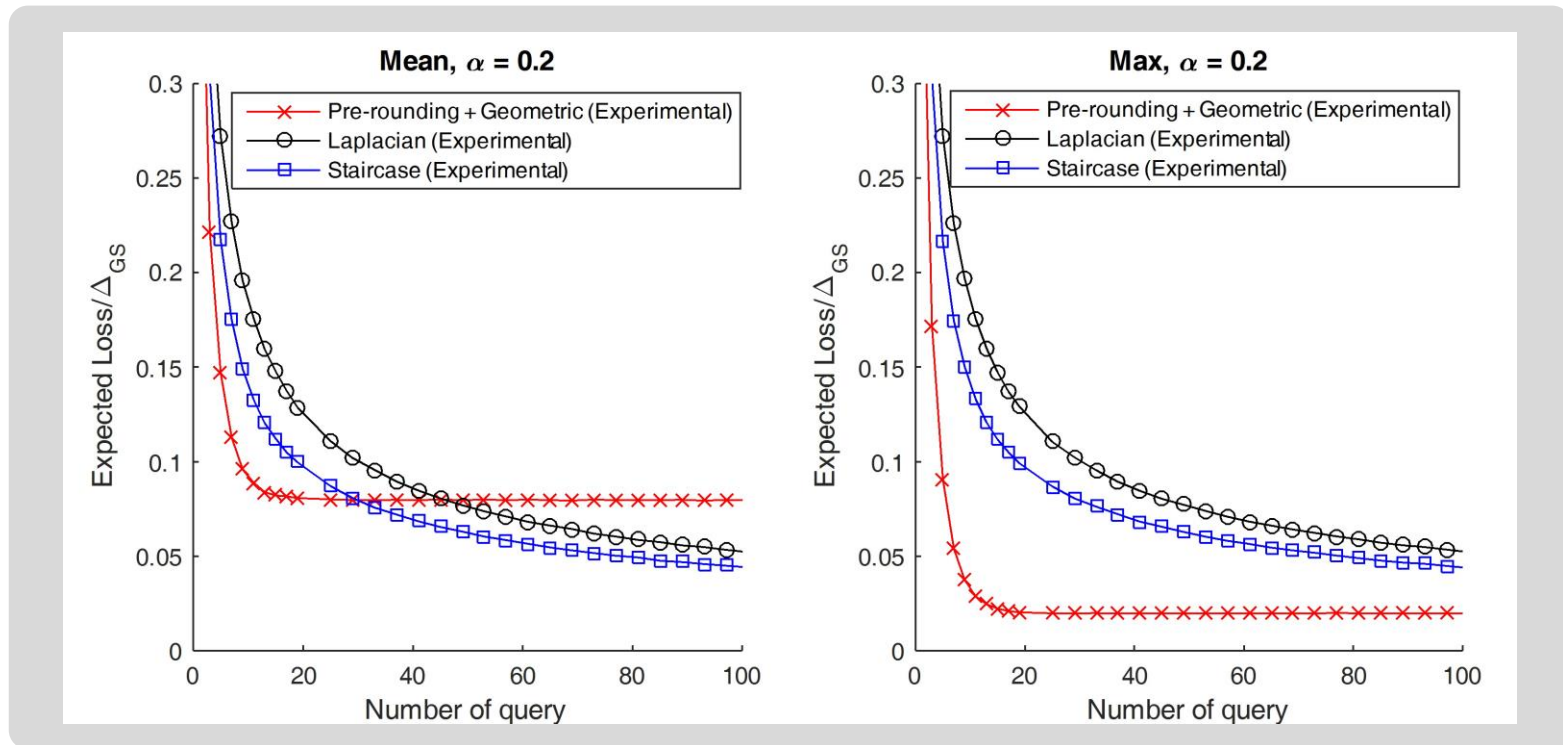
- ▶ (1): Pre-rounding
 - only outputs $\{-10, 0, 10\}$
- ▶ (2): Add truncated α -Geometric Noise ($\alpha = e^{-\epsilon}$)
 - $P_r[X > 10]$ goes to $P_r[X = 10]$
 - $P_r[X < -10]$ goes to $P_r[X = -10]$
- ▶ The heuristic mechanism satisfies ϵ -DP ($\alpha = e^{-\epsilon}$)
- ▶ Mechanism designed for **low-variance** priors

Utility-Privacy Tradeoff Performance



- ▶ Significant improvement in low & intermediate privacy regime (*the red 'x'*).
- ▶ In the high privacy regime tend towards convergence
 - DP mechanism adds extremely large noise to maintain privacy
 - noise dominates the performance

Our Mechanism is Collusion-Proof !



- ▶ Users collude in perturbed results (based on MLE)
- ▶ The heuristic design is collusion-proof (*the red curve*)

Design Insights

When query outputs are substantially narrowed down by side-information, discretizing the domain and adding truncated Geometric noise is a good idea

A robust, simple, and efficient Bayesian design is possible!

A collusion-proof Bayesian design is also feasible

Future Directions

Optimal Bayesian design mechanism

- so that we know how good our design is
- new heuristic methods and design insights
- studies of implementation complexity

Applications of the optimal Bayesian design

- applying Bayesian design to practical problems with side-information
- many practical issues will be involved

Optimal Bayesian design in approximate DP

- more efficient, but less robust

Thank you!

Email: chienlun@usc.edu