Computing Inconsistency Measures Under Differential Privacy

Shubhankar Mohapatra¹, Amir Gilad², Xi He¹, Benny Kimelfeld³

University of Waterloo¹, Hebrew University², Technion³

Challenges

are NP hard problems

the number of nodes

constraints **\Sigma**

Setup

I can't share

information before

you buy!

Private marketplace

Research Questions:

Are there errors in a private dataset?

How much effort is required to repair?

Background

Denial Constraints and Conflict Graphs

ID	Capital	Country	
	Ottawa	Canada	
2	Ottawa	Canada	
3	Ottawa	Canada	
4	Ottawa	Kanada	

 σ : Capital —> Country

"country of two tuples must be the same if their capital is the same"

Inconsistency Measures [1]

Given, dataset D and constraint set Σ , inconsistency measures are of the form $I(D, \Sigma) \to \mathbb{R}$:

- I.Drastic Measure $I_D(G)$ = existence of an edge \bigotimes
- 2. Minimal inconsistency measure $I_{MI}(G) = number of edges$
- 3. Problematic measure $I_P(G)$ = number of vertices with positive degree \bigcirc
- 4. Maximal consistency measure $I_{MC}(G) = number of maximal independent sets <math>(K)$
- 5. Optimal repair measure I_R (G) = minimum vertex cover size \bigcirc

Graph projection for I_P and I_{MI} [$\theta^* \ll \mathcal{O}(n)$] Optimize $|\Theta|$ using Tune to find a good θ^* knowledge from Σ Choose θ^* Generate Input dataset D, candidates conflict from Θ \rightarrow graph G privately Θ and

Is this data good for

me? 🤴

Data Buyer

Differentially Private Inconsistency Measures

I. Computational hardness: Minimum vertex cover (IR) and #maximal independent sets (IMC)

2. High sensitivity: Maximum change in output when G is replaced by G' is $\mathcal{O}(n)$, where n is

degree $heta^*$

DP vertex cover size for $I_R [2 \ll \mathcal{O}(n)]$

Each vertex in G_{θ^*}

has maximum

We analyse the 2-approximate vertex cover size algorithm and show that it has sensitivity of 2.

Project G to

 $G_{\theta^*}[3]$

Utility depends

on size of Θ and choice

of $heta^*$

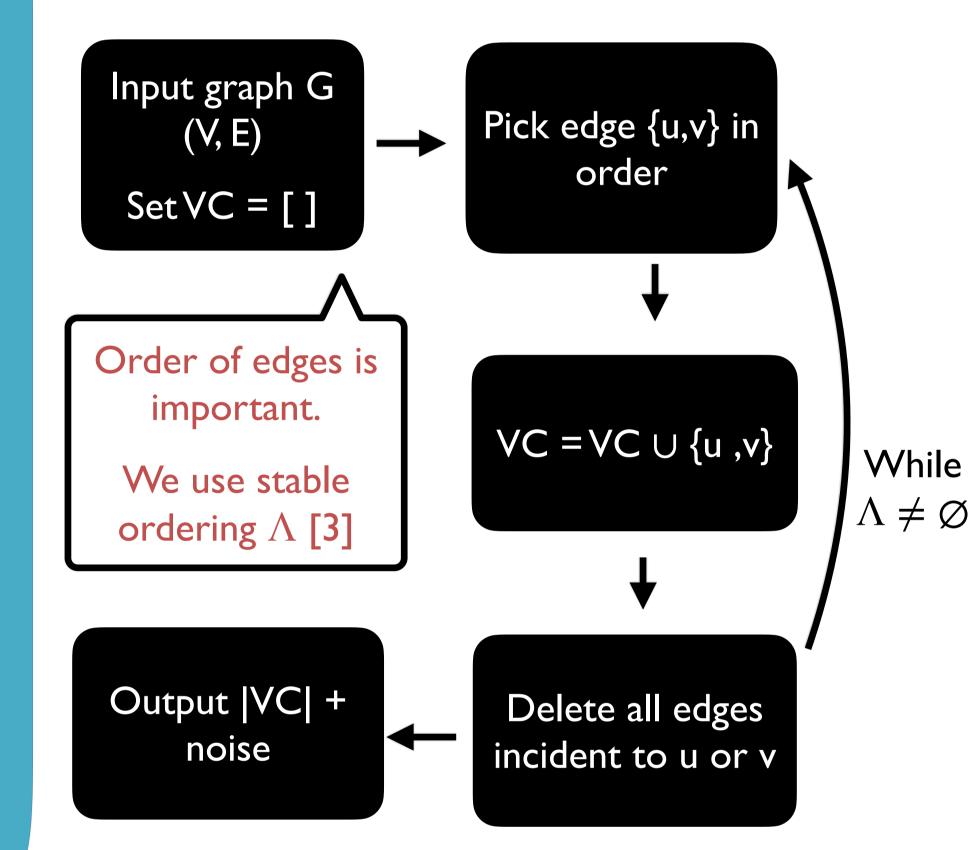
Compute

measure by

adding noise

Sensitivity of

measures $\propto \theta^*$



= our work

Differential Privacy [2]

A randomized algorithm $A: \mathcal{G} \to \mathcal{R}$ satisfies \mathcal{E} differential privacy (DP) if for any two adjacent graphs $G, G' \in \mathcal{G}$ that differ in a node and for any subset of outputs $o \subseteq \mathcal{R}$ it holds that :

 $\Pr[A(G) \in o] \le e^{\varepsilon} \Pr[A(G') \in o]$

Experiments

<u>Setup</u>

- Inconsistency: Random typo to 1% rows
- Privacy: $\varepsilon = 1$
- Measure: True vs private by adding one typo at a time

<u>Datasets</u>

- Five real-world datasets with varying conflict graph densities
- We experiment on subset of 10k and repeat for 10 times and average

Dataset	#Tuples	#Attributes	#Constraints	Graph density
Adult	32k	15	3	9635
Flight	500k	20	13	1520
Hospital	I I 4k	15	7	793
Stock	I22k	7		I
Tax	IM	15	9	373

References:

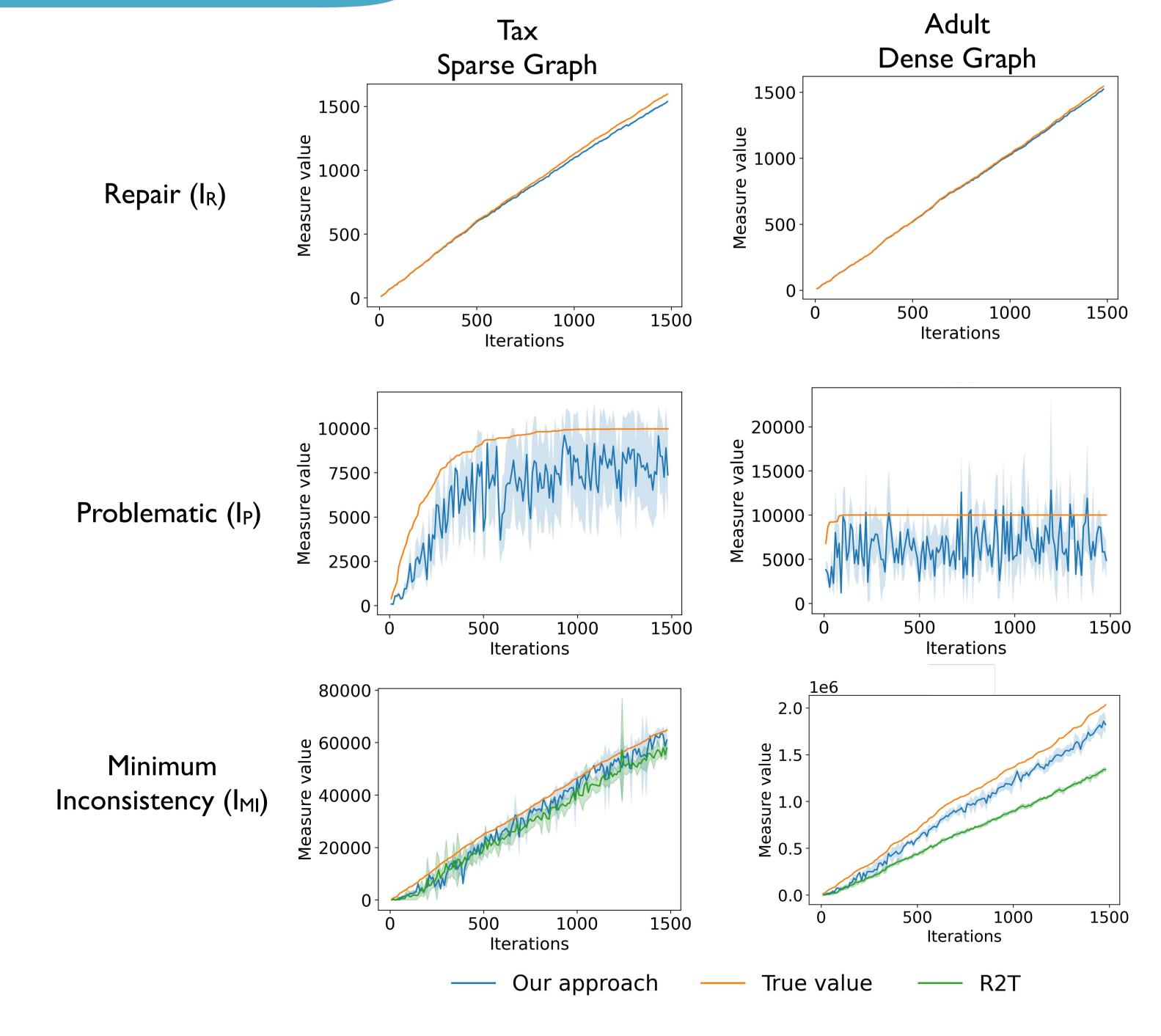
1. Livshits, Ester, et al. "Properties of inconsistency measures for databases." International Conference on Management of Data. 2021.

2. Dwork, Cynthia, et al. "Calibrating noise to sensitivity in private data analysis." TCC 2006

3. Day, Wei-Yen, et al. "Publishing graph degree distribution with node differential

4. Dong, Wei, et al. "R2T: Instance-optimal Truncation for Differentially Private Query Evaluation with Foreign Keys." ACM SIGMOD 2023





privacy." International Conference on Management of Data 2016.

DAVID R. CHERITON SCHOOL OF COMPUTER SCIENCE