# Differentially Private Analysis of Graphs

Adam Smith

Pennsylvania State University



#### Publishing information about graphs

Many types of data can be represented as graphs where

- nodes correspond to individuals
- edges capture relationships
  - "Friendships" in online social network
  - Financial transactions
  - Email communication
  - Health networks (of doctors and patients)
  - Romantic relationships



image source http://community.expressorsoftware.com/blogs/mtarallo/36-extracting-datafacebook-social-graph-expressor-tutorial.html

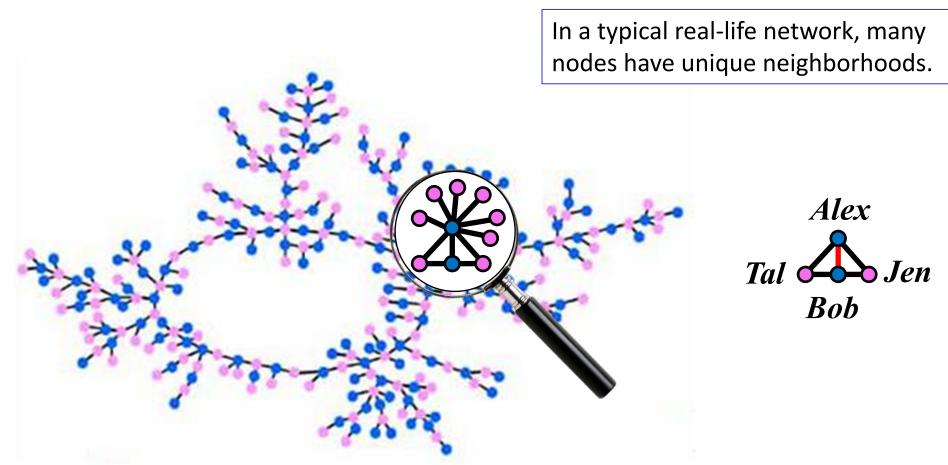


Such graphs contain potentially sensitive information.

image source http://www.queticointernetmarketing.com/new-amazing-facebook-photo-mapper/

### ``Anonymized'' graphs still pose privacy risk

- False dichotomy: personally identifying vs. non-personally identifying information.
- Links and any other information about individual can be used for de-anonymization.



Bearman, Moody, Stovel. Chains of affection: The structure of adolescent romantic and sexual networks, American J. Sociology, 2008

#### De-anonymization attacks

- Movie ratings [Narayanan, Shmatikov 08]
- Social networks

[Backstrom Dwork Kleinberg 07,

Narayanan Shmatikov 09, Narayanan Shi Rubinstein 12]

Computer networks

[Coull Wright Monrose Collins Reiter 07,

Ribeiro Chen Miklau Townsley 08]

#### Can reidentify individuals based on external sources.







# What information can be released

without violating privacy?

#### Two variants of differential privacy for graphs

Edge differential privacy



Two graphs are neighbors if they differ in one edge.

Node differential privacy



Two graphs are **neighbors** if one can be obtained from the other by deleting *a node and its adjacent edges*.

#### Work on weaker/incomparable privacy definitions

- Edge differential privacy
  - Number of triangles, MST cost [Nissim Raskhodnikova Smith 07]
  - ➤ Degree distribution [Hay Rastogi Miklau Suciu 09, Hay Li Miklau Jensen 09, Karwa Slavkovic 12, Kifer Lin 13]
  - Small subgraph counts [Karwa Raskhodnikova Smith Yaroslavtsev 11]
  - Cuts [Hardt Rothblum 10, Gupta Roth Ullman 12, Blocki Blum Datta Sheffet 12]
  - Kronecker graph model parameters [Mir Wright 12]
- Edge-private against Bayesian adversary (weaker privacy)
  - Small subgraph counts [Rastogi Hay Miklau Suciu 09]
- Node zero-knowledge private
  - Privacy only for bounded-degree graphs
  - > Average degree, distances to graph families [Gehrke Lui Pass 12]

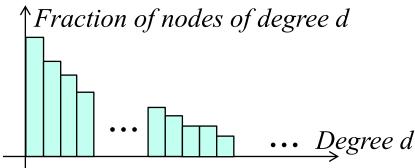
# What graph statistics can be computed accurately

## with node differential privacy?

Important: small error on sparse graphs.

#### Graph statistics

- [Blocki Blum Datta Sheffet 13, Kasiwiswanathan Nissim Raskhodnikova Smith 13, Chen Zhou 13]
  - ➤ Number of edges
  - Counts of small subgraphs (e.g., triangles, k-triangles, k-stars)
- Degree distribution [Kasiwiswanathan Nissim Raskhodnikova Smith 13, Raskhodnikova Smith 16]



#### Our algorithms for these statistics

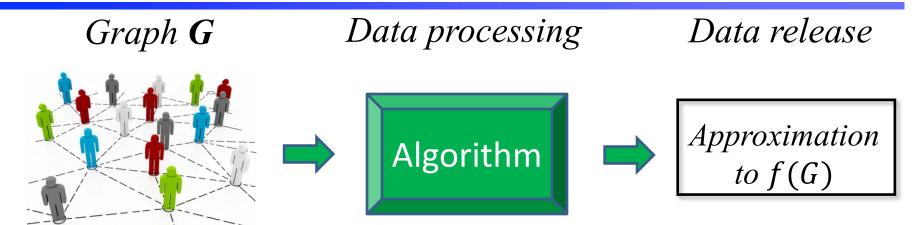
- node differentially private for all graphs
  - ➤ Accurate for a large subclass of graphs (including sparse graphs, scale-free graphs and Erdos-Renyi graphs)
  - $\rightarrow$  (1+ $o_n$ (1))-approximation

#### Techniques for node-private algorithms

- Previous work
  - Sensitivity analysis of simple projections [BBDS'13, KNRS'13]
  - > yielded generic reductions to privacy over bounded-degree graphs
  - Lipschitz extensions [BBDS'13, KNRS'13, CZ'13]
  - yielded much more accurate algorithms, but worked only for releasing 1-dimensional statistics (edge / small subgraph counts).
- This work: new techniques that improve accuracy
  - Lipschitz extensions for high-dimensional statistics
  - Generalized exponential mechanism for both methods

# Lipschitz, extensions

#### Basic question: how to compute a statistic f



How accurately can an  $\epsilon$ -differentially private algorithm compute f(G)?

### Basic technique: noise proportional to sensitivity

• Global sensitivity of a function f is

$$\partial f = \max_{\text{(node)neighbors } G,G'} |f(G) - f(G')|$$

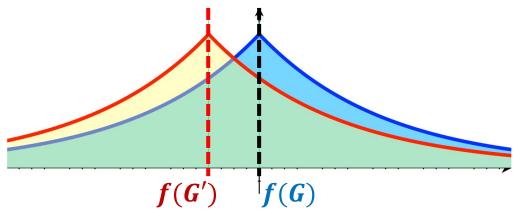
Also called **Lipschitz constant** of f.



• **GS Framework** [DMNS]: For every real-valued function f, there is an  $\epsilon$ -differentially private algorithm A such that

$$\mathbb{E}\left(|A(G) - f(G)|\right) = \frac{\partial f}{\epsilon}.$$

• Intuition: Adding noise  $\approx \frac{\partial f}{\epsilon}$  makes G, G' hard to distinguish



### Challenge for node privacy: high sensitivity

• Global sensitivity of a function f is

$$\partial f = \max_{\text{(node) neighbors } G,G'} |f(G) - f(G')|_1$$



- Examples:
- $\succ f_{\text{edge}}(G)$  is the number of edges in G.
- $\succ f_{\text{deg}}(G)$  is the degree list of G.

for graphs on n nodes:

$$\partial f_{\text{edge}} = n.$$

$$\partial f_{\text{deg}} = 2n.$$

*Problem:* high-degree nodes. <<

### Graphs of small degree

Let G = family of all graphs,

 $G_d$  = family of graphs of degree  $\leq d$ .

Notation.  $\partial f$  = global sensitivity of f over G.

 $\partial_d f$  = global sensitivity of f over  $G_d$ .

**Observation.**  $\partial_d f$  is low for many useful f.

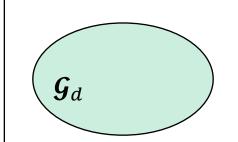
#### **Examples:**

- $\rightarrow$   $\partial_d f_{\text{edge}} = d$  (compare to  $\partial f_{\text{edge}} = n$ )
- $\rightarrow \partial_d f_{\text{deg}} = 2d$  (compare to  $\partial f_{\text{deg}} = 2n$ )

Goal: privacy for all graphs

Idea: "Extend" f from  $G_d$  to G for a carefully chosen  $d \in [n]$ .





#### Lipschitz extensions

A function f' is a Lipschitz extension of f from  $G_d$  to G if

 $\succ f'$  agrees with f on  $\boldsymbol{\mathcal{G}}_d$  and

$$> \partial f' = \partial_d f$$

• Release f' via GS framework [DMNS'06]

$$g$$

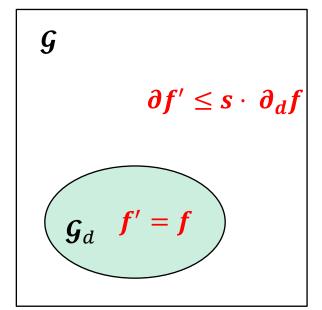
$$\partial f' = \partial_d f$$

$$g_d \quad f' = f$$

- All real-valued functions have Lipschitz extensions [McShane 34]
- Lipschitz extensions for subgraph counts that can be computed efficiently [Kasiviswanathan Nissim R Smith 13]

#### Lipschitz extensions: vector-valued functions

A function f' is a Lipschitz extension of f from  $\mathcal{G}_d$  to  $\mathcal{G}$  with stretch s if f' agrees with f on  $\mathcal{G}_d$  and  $f' \leq s \cdot \partial_d f$ 



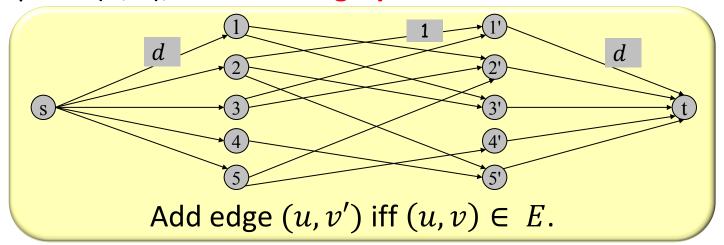
- We can still release f' via GS framework
- There exist functions  $f: \mathcal{G}_d \to \mathbb{R}^3$

 $(f: \mathbf{G}_d \to \mathbb{R}^2 \text{ if } \ell_2 \text{ is used as output metric instead of } \ell_1)$ that do not admit st pschitz extensions

 Lipschitz extensions of degree list and degree distribution (with small stretch) that can be computed efficiently

## Lipschitz extension of $f_{\text{edge}}$ : flow graph [KNRS'13]

For a graph G=(V, E), define **flow graph of G**:

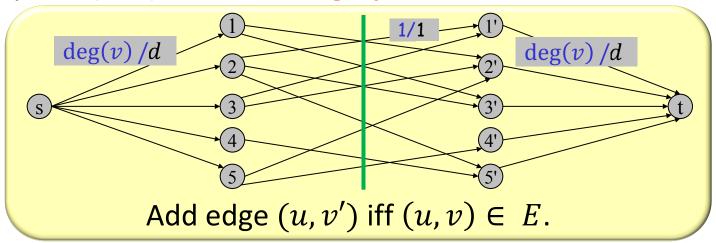


 $v_{\text{flow}}(G)$  is the value of the maximum flow in this graph.

**Lemma**.  $v_{\text{flow}}(G)/2$  is a Lipschitz extension of  $f_{\text{edge}}$ .

## Lipschitz extension of $f_{\text{edge}}$ : flow graph [KNRS'13]

For a graph G=(V, E), define **flow graph of G**:



 $v_{\text{flow}}(G)$  is the value of the maximum flow in this graph.

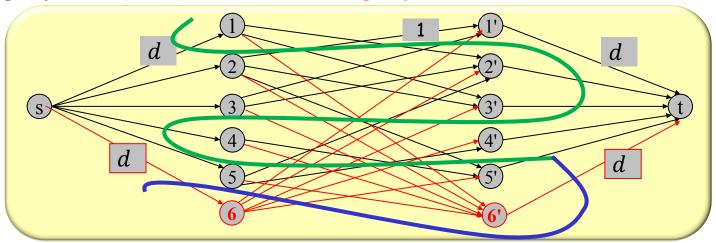
**Lemma**.  $v_{\text{flow}}(G)/2$  is a Lipschitz extension of  $f_{\text{edge}}$ .

Proof: (1) 
$$v_{\text{flow}}(G) = 2f_{\text{edge}}(G)$$
 for all  $G \in \mathcal{G}_d$ 

(2) 
$$\partial v_{\text{flow}} = 2 \cdot \partial_d f_{\text{edge}}$$

## Lipschitz extension of $f_{\text{edge}}$ : flow graph [KNRS'13]

For a graph G=(V, E), define **flow graph of G**:



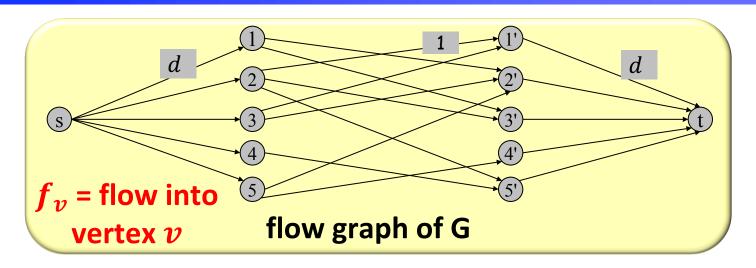
 $v_{\text{flow}}(G)$  is the value of the maximum flow in this graph.

**Lemma**.  $v_{\text{flow}}(G)/2$  is a Lipschitz extension of  $f_{\text{edge}}$ .

Proof: (1) 
$$v_{\text{flow}}(G) = 2f_{\text{edge}}(G)$$
 for all  $G \in \mathcal{G}_d$ 

(2) 
$$\partial v_{\text{flow}} = 2 \cdot \partial_d f_{\text{edge}} = 2d$$

#### Lipschitz extension of the degree list [RS '16]

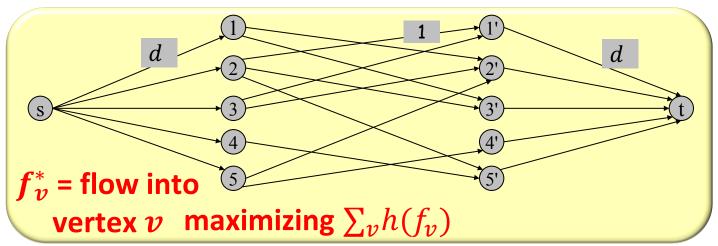


Can we use  $f_v$  as a proxy for degree of v?

*Issue:* max flow is not unique.

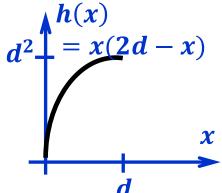
Want: unique flow that has low global sensitivity.

## Lipschitz extension of $f_{\text{deg}}$ : convex programming



*Idea:* maximize  $\sum_{v} h(f_v)$  instead of  $\sum_{v} f_v$ , where h(x) is strictly concave.

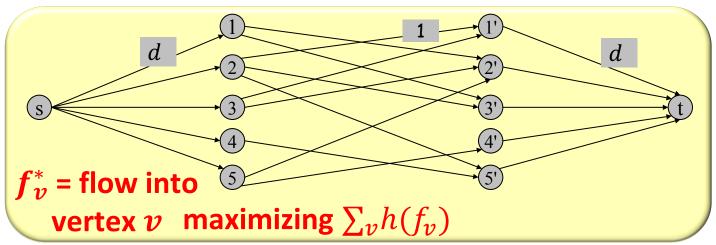
- Let f\* be the vector of s-out-flows.
  - $\triangleright f^*$  is unique, since h is strictly concave.
  - Poly-time computable [Lee Rao Srivastava 13].
- Lemma.  $f^*$  is a Lipschitz extension of degree-list with stretch 3/2.



## Lipschitz extension of $f_{\text{deg}}$ : combinatorial

- Subsequently simplified [W.-Y. Day, N. Li, M. Lyu, SIGNMOD 2016]
- **DLL Algorithm**: On input G = (V, E) and degree bound d
  - $\triangleright$  Order the edges lexicographically:  $E=e_1,e_2,\ldots,e_m$
  - $\succ G' = (V, \emptyset)$  //empty graph
  - $\triangleright$  For i=1 to m:
    - If (adding  $e_i$  to G' would not push maximum degree over d)
      - Add e to V
    - Else
      - Ignore  $e_i$
  - $\triangleright$  **Return** sort(degree-list(G'))
- **Lemma:** For any two graphs  $G_1$ ,  $G_2$  that differ in one node,  $\|DLL(G_1) DLL(G_2)\|_1 \le 2d + 1$

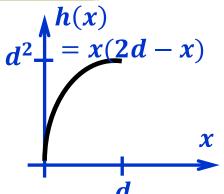
#### Releasing degree list: summary



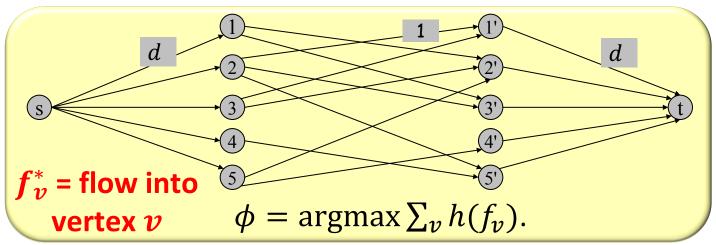
- 1. Construct flow graph of G.
- 2. Compute s-out-flows  $f^*$ .



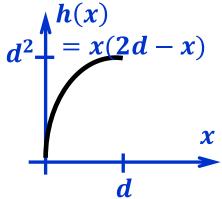




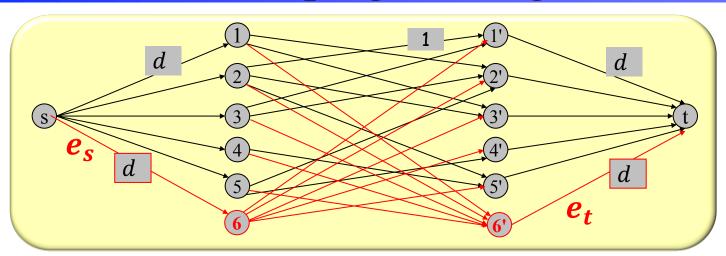
# Lipschitz, extension of degree list via convex programming



- If  $G \in \mathcal{G}_d$ , then  $f_v^* = \deg(v)$  for all v, since h is strictly increasing on [0,d].
- $\partial_d(degree\ list) = 2d$ : can add a node of degree  $\leq d$ .
- Lemma.  $\ell_1$  global sensitivity  $\partial f^* \leq 3d$ .



#### Lipschitz extension of degree list via convex programming



Lemma.  $\ell_1$  global sensitivity  $\partial f^* \leq 3d$ .

**Proof sketch:** Consider  $g = \phi_{new} - \phi_{old}$ .

g is a union of simple s-t-paths and cycles of several types:

- s-t-paths using  $e_t$ .
- 3. Cycles using  $e_t$ .
- 4. Remaining paths and cycles.

1. 
$$s$$
- $t$ -paths and cycles using  $e_s$ . Contribute  $\leq 2d$  to  $|f_{new}^* - f_{old}^*|_1$   
2.  $s$ - $t$ -paths using  $e_t$ .  $\leq d$   
3. Cycles using  $e_t$ .

Do not exist. Use strict concavity of h

# Generalized Exponential Mechanism

Choosing the cutoff degree d

#### Evaluating Cutoff Degrees

Given: candidate real-valued Lipschitz extensions  $f_d$  and their sensitivities  $\partial f_d$  (Specifically, we will try d= powers of 2 in [n])

• If we approximate f by releasing  $f_d$  in GS framework, the expected error is roughly  $f(G) - f_d(G) + \frac{\partial f_d}{\epsilon}$ 

Want:  $f_d$  with (approximately) smallest

score 
$$q_d(G) := -f_d(G) + \frac{\partial f_d}{\epsilon}$$

# Differentially private algorithms for choosing the item with the smallest score

- Exponential Mechanism [McSherry Talwar 07]
  - ➤ Initial motivation: auction design.
  - Subsequent applications:
    - ➤ Learning discrete classifiers [KLNRS'08]
    - ➤ Synthetic data generation [BLR'08,...,HLM'10]
    - ➤ Convex Optimization [CM'08,CMS'10]
    - Frequent Pattern Mining [BLST'10]
    - ➤ Genome-wide association studies [FUS'11]
    - ➤ High-dimensional sparse regression [KST'12]
    - >...

# Differentially private algorithms for choosing the item with the smallest score

- Exponential Mechanism [McSherry Talwar 07]
  - ➤ Additive error in the score is proportional to **maximum** (over items) **sensitivity of a score** function

- Our Generalized Exponential Mechanism
  - ➤ Additive error in the score is proportional to the **sensitivity of the optimal score** function

#### Exponential Mechanism [McSherry Talwar 07]

Given: database x from universe U, parameter  $\epsilon > 0$ , score functions  $q_i \colon U \to \mathbb{R}$  with  $\delta_i = \partial q_i$  for  $i \in [k]$ 

Want: index  $i^* = \arg\min_i q_i(x)$ 

#### Algorithm EM

- 1. Set  $\delta = \max_i \delta_i$
- 2. Output  $\hat{i}$ , set to i with probability  $\propto \exp(\epsilon \cdot q_i(x)/\delta)$ , normalized so that probabilities sum to 1, for  $i \in [k]$ .

Guarantees: (1) EM is 
$$2\epsilon$$
-differentially private (2)  $\forall \beta \in (0,1)$ , w.p.  $\geq \beta$ ,  $q_{\hat{\imath}}(x) \leq q_{i^*}(x) + \delta \cdot \frac{2 \ln(k/\beta)}{\epsilon}$ 

Goal: Guarantee with  $\delta_{i^*}$  instead of  $\delta_{i^*}$ 

#### Generalized Exponential Mechanism [RS16]

Given: database x from universe U,  $\epsilon > 0$ ,  $\beta \in (0,1)$ , score functions  $q_i \colon U \to \mathbb{R}$  with  $\delta_i = \partial q_i$  for  $i \in [k]$ 

Want: index  $i^* = \arg\min_i q_i(x)$ 

#### Algorithm **GEM**

- 1.  $t \leftarrow 2 \ln(k/\beta)/\epsilon$
- 2.  $q'_i(x) \leftarrow q_i(x) + t\delta_i$  for all  $i \in [k]$
- 3.  $s_i(x) \leftarrow \max_j \frac{q_i'(x) q_j'(x)}{\delta_i + \delta_j}$  for all  $i \in [k]$
- 4. Output  $\hat{i}$ , set to i with probability  $\propto \exp(\epsilon \cdot s_i(x))$ , normalized so that probabilities sum to 1, for  $i \in [k]$ .

Note: each  $q_i'$  has sensitivity  $\delta_i$  each  $s_i$  has sensitivity 1

#### Generalized Exponential Mechanism [RS16]

Given: database x from universe U,  $\epsilon > 0$ ,  $\beta \in (0,1)$ , score functions  $q_i \colon U \to \mathbb{R}$  with  $\delta_i = \partial q_i$  for  $i \in [k]$ 

Want: index  $i^* = \arg\min_i q_i(x)$ 

Guarantees: (1) GEM is 
$$2\epsilon$$
-differentially private 
$$(2) \ \forall \beta \in (0,1), \ \text{w.p.} \ \geq \beta,$$
 
$$q_{\hat{\imath}}(x) \leq q_{i^*}(x) + \delta_{i^*} \cdot \frac{4 \ln(k/\beta)}{\epsilon}$$

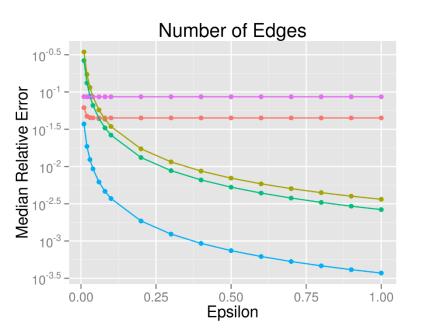
#### Summary of techniques we saw

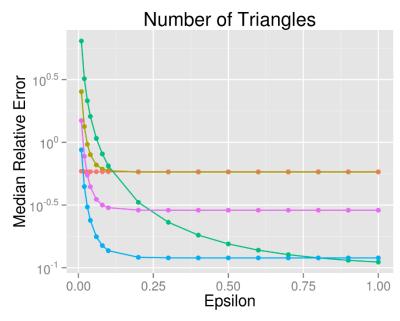
- 1. Lipschitz extensions [BBDS13, KNRS13]
- ➤ Releasing number of edges via max flow [KNRS13]
- > Releasing degree list: via convex programming [RS16]
- > The most accurate known method
- 2. Generalized exponential mechanism [SR16]
- ➤ For choosing among objects with score functions of different sensitivities
- > For choosing the cutoff degree

#### **Conclusions**

- It is possible to design node differentially private algorithms with good utility for a large class of graphs
  - $\triangleright$  One can choose a "good" value of d privately
- Directions for future work
  - ➤ Understanding which functions have efficiently computable Lipschitz extensions with small stretch
  - ➤ Node-private algorithms for releasing other graph statistics
  - Node-private synthetic graphs
- Open Question: Is there a node-differentially private algorithm for releasing the cost of all graph cuts with worst-case error  $o(n \cdot \text{max-degree}(G))$ ?

#### Experiments for the flow and LP method [Lu]





	Graph	# nodes	# edges	Max degree	Time, secs # edges	Time, secs # Δs
-	CA-GrQc	5,242	28,992	81	0.02	7
-	CA-HepTh	9,877	51,996	65	0.68	0.5
-	CA-AstroPh	18,772	396,220	504	0.34	10,222
-	com-dblp-ungraph	317,080	2,099,732	343	2	2128
-	com-youtube-ungraph	1,134,890	5,975,248	28,754	9	94