

Differentially Private Analysis of Graphs and Social Networks

Sofya Raskhodnikova
BU Computer Science

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.expressor-software.com/blogs/mtarallo/36-extracting-data-facebook-social-graph-expressor-tutorial.html>

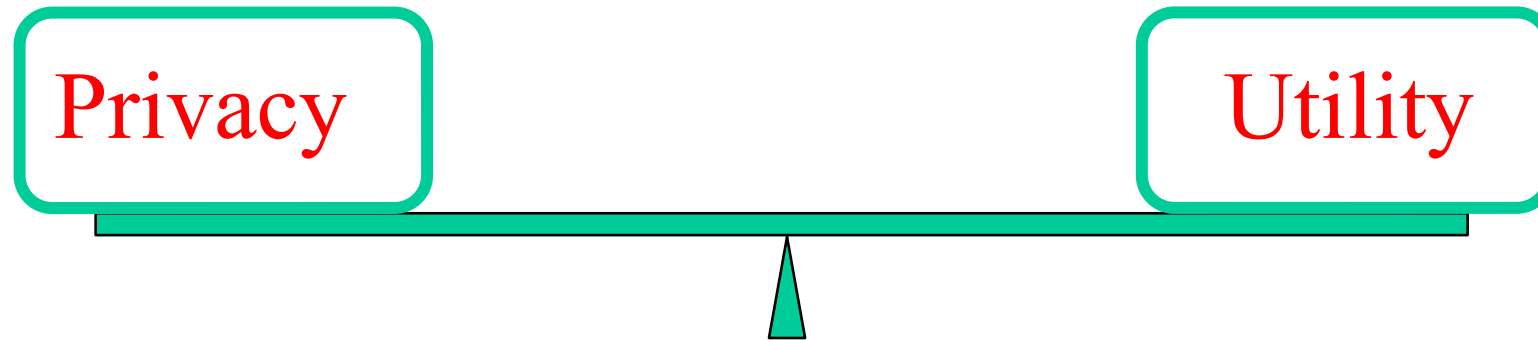


image source <http://www.queticointernetmarketing.com/new-amazing-facebook-photo-mapper/>

Such graphs contain potentially sensitive information.

Two conflicting goals

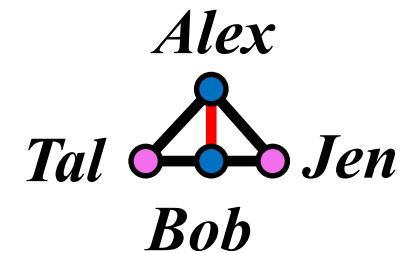
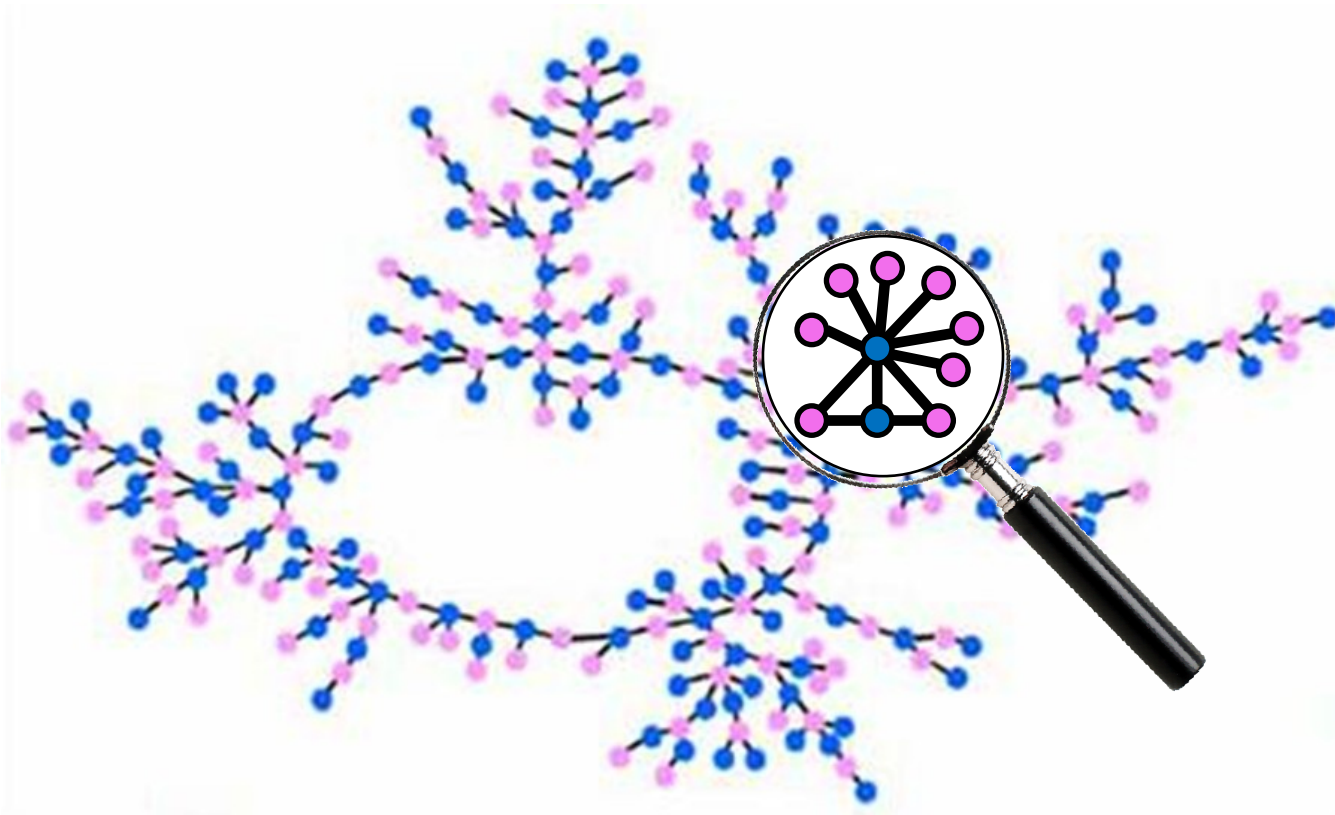
- **Privacy:** protecting information of individuals.
- **Utility:** drawing accurate conclusions about aggregate information.



“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.

In a typical real-life network, many nodes have unique neighborhoods.



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 Monroe Collins Reiter 07,
Ribeiro Chen Miklau Townsley 08]

Can reidentify individuals based on external sources.



internet



social networks

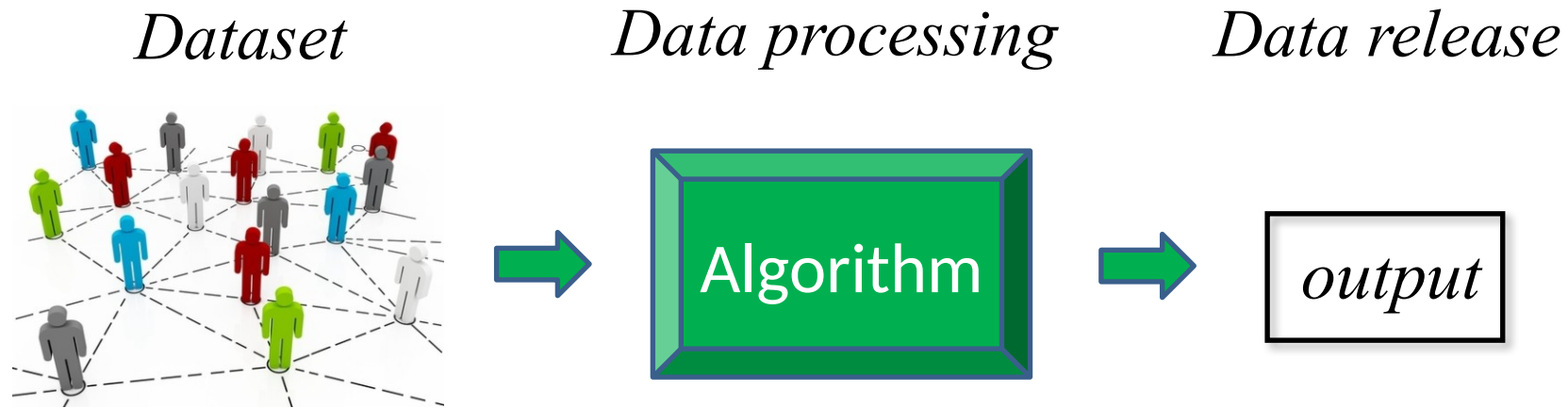


anonymized datasets

What information can be released

without violating privacy?

Differential privacy



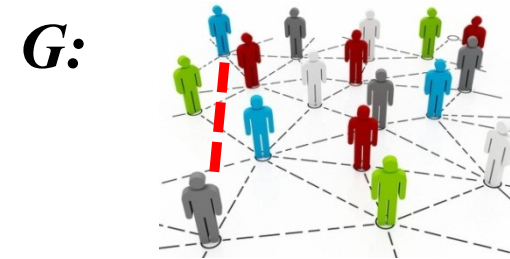
Differential privacy [Dwork McSherry Nissim Smith 06]

Intuition: Two datasets are **neighbors** if they differ in one individual's data.

An algorithm is **differentially private** if its output is roughly the same for all pairs of **neighbors**.

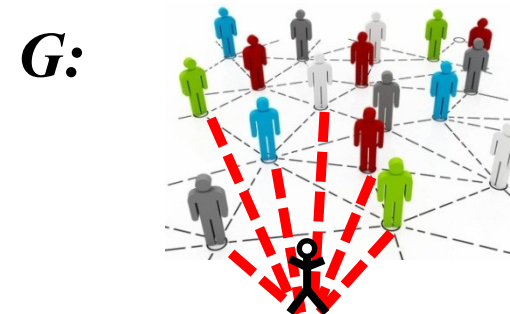
Two variants of differential privacy for graphs

- **Edge** differential privacy



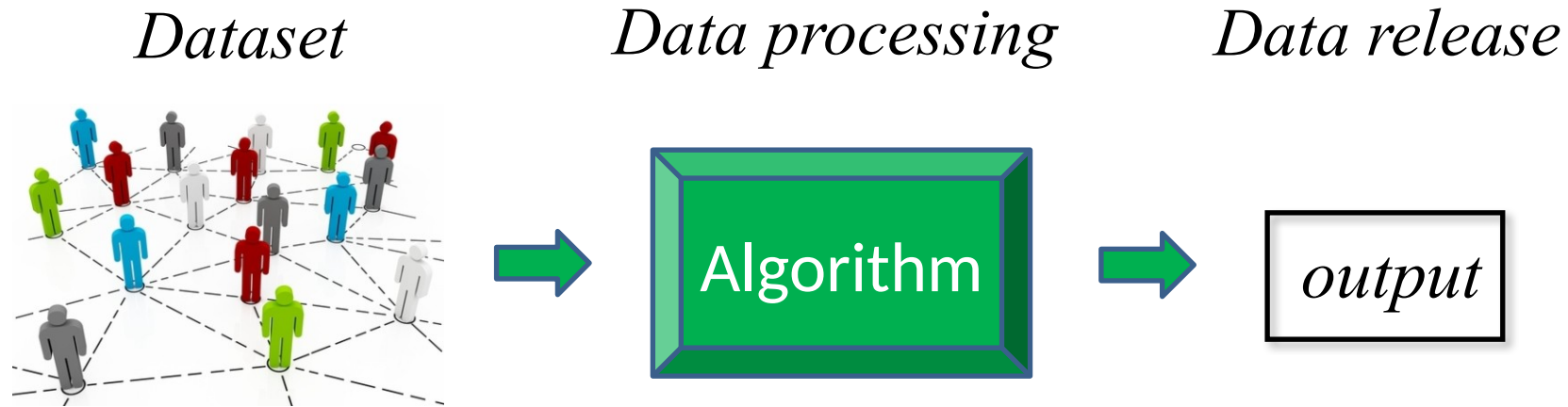
Two graphs are **neighbors** if they differ in **one edge**.

- **Node** differential privacy



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

Differential privacy (for graph data)



Differential privacy [Dwork McSherry Nissim Smith 06]

An algorithm A is **-differentially private** if for all pairs of **neighbors** and all possible sets of outputs S :

$$\Pr \left[\mathbf{A} \left(\mathbf{G} \right) \in \mathbf{S} \right] \leq e^{\epsilon} \Pr \left[\mathbf{A} \left(\mathbf{G}' \right) \in \mathbf{S} \right]$$

Important properties of differential privacy

- Robustness to side information

- Post-processing:

an algorithm accessing the data via an ϵ -differentially private subroutine is ϵ -differentially private.

- Composition:

an algorithm that runs k subroutines that are ϵ -differentially private is $k\epsilon$ -differentially private.

Is differential privacy too strong?

- No weaker notion has been proposed that satisfies all three important properties.
- We can attain it for many useful statistics!

What graph statistics can be computed accurately with node differential privacy?

Important: small error on sparse graphs.

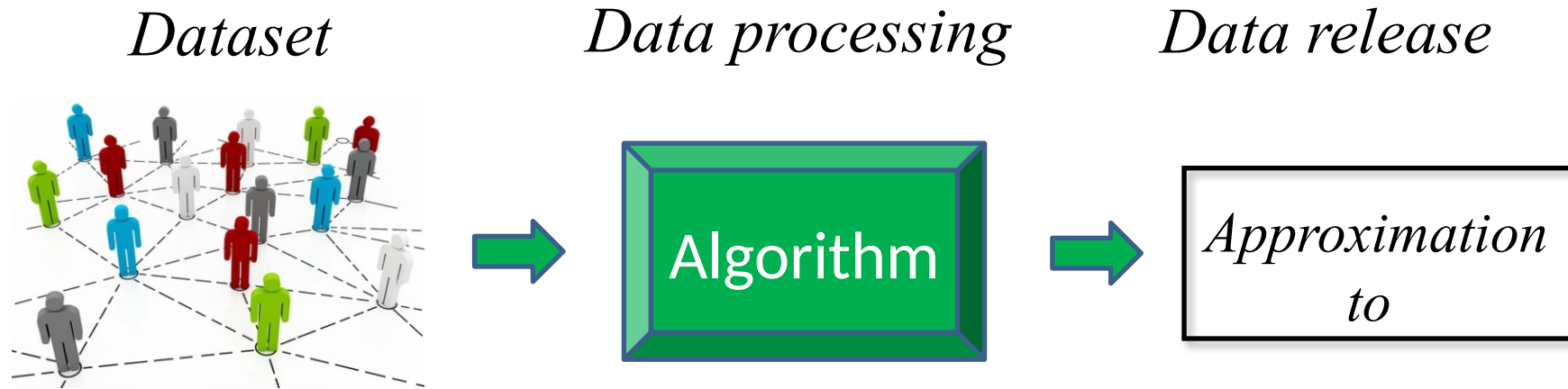
Tools used in differentially private graph algorithms

- Smooth sensitivity
 - A more nuanced notion of sensitivity than the one mentioned in Katrina's talk
- Sample and aggregate
- Maximum flow
- Linear and convex programming
- Random projections
- Iterative updates
- Postprocessing
- Noisy stochastic gradient descent

Differentially private graph analysis

A taste of techniques

Basic question: how to compute a statistic f

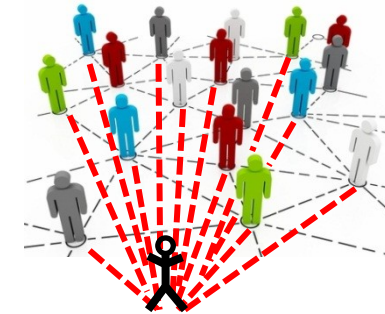


How accurately
can an ϵ -differentially private algorithm compute $f(G)$?

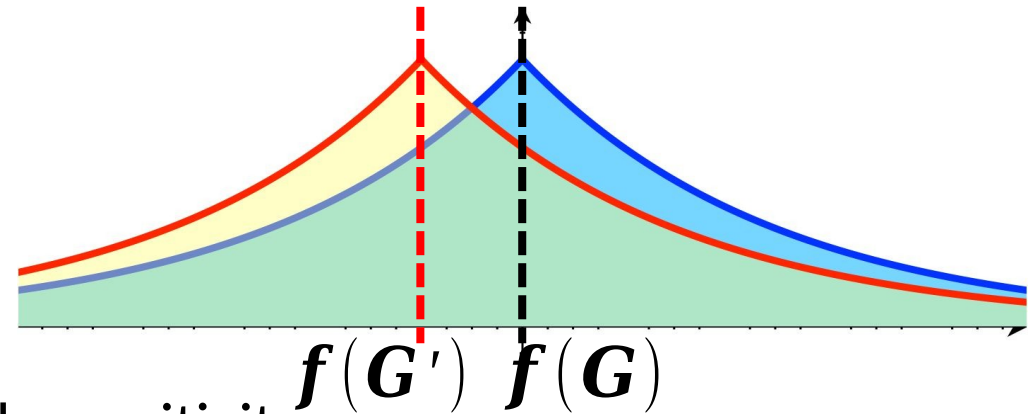
Challenge for node privacy: high sensitivity

- **Sensitivity** of a function is

$$\partial f = \max_{(\text{node}) \text{ n e i g h b o r s } G, G'} |f(G) - f(G')|$$



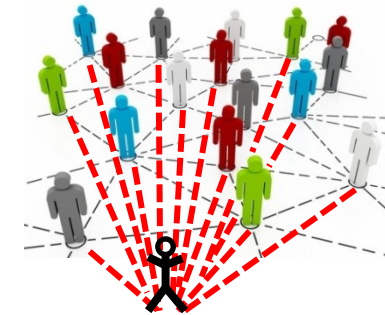
- **Katrina's talk:** if has low sensitivity, it can be computed accurately with differential privacy (by Laplace mechanism).
- **Intuition:** Adding noise makes , hard to distinguish



- **Challenge:** Even simple graph statistics have high sensitivity.

Challenge for node privacy: high sensitivity

- **Sensitivity** of a function is



- **Examples:**

- $f(G)$ is the number of edges in G .
- $f(G)$ is the number of triangles in G .

for graphs on n nodes:

$f(G) = \frac{n(n-1)}{2}$.

$f(G) = \frac{n(n-1)(n-2)}{6}$.



*Idea: project onto graphs
with low sensitivity.*

[Kasiviswanathan Nissim
Raskhodnikova Smith 13]

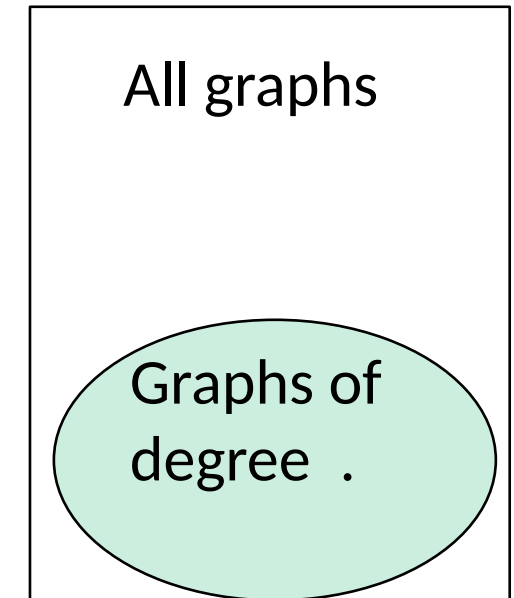
See also [Blocki Blum Datta Sheffet 13, Chen Zhou 13]

“Projections” on graphs of small degree

Consider graphs of degree d ,
where

- **Examples:**

- $e(G)$ is the number of edges in G .
- $t(G)$ is the number of triangles in G .



$e = \frac{1}{2} \sum d_i$

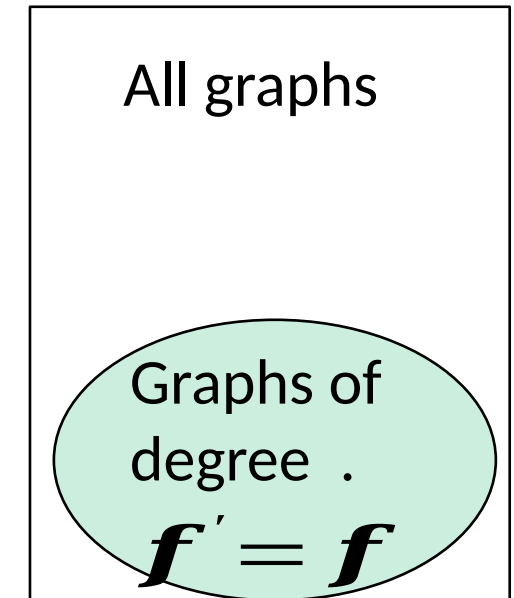
$t = \frac{1}{6} \sum d_i(d_i-1)(d_i-2)$

Lipschitz extensions

A function is a **Lipschitz extension** of f if

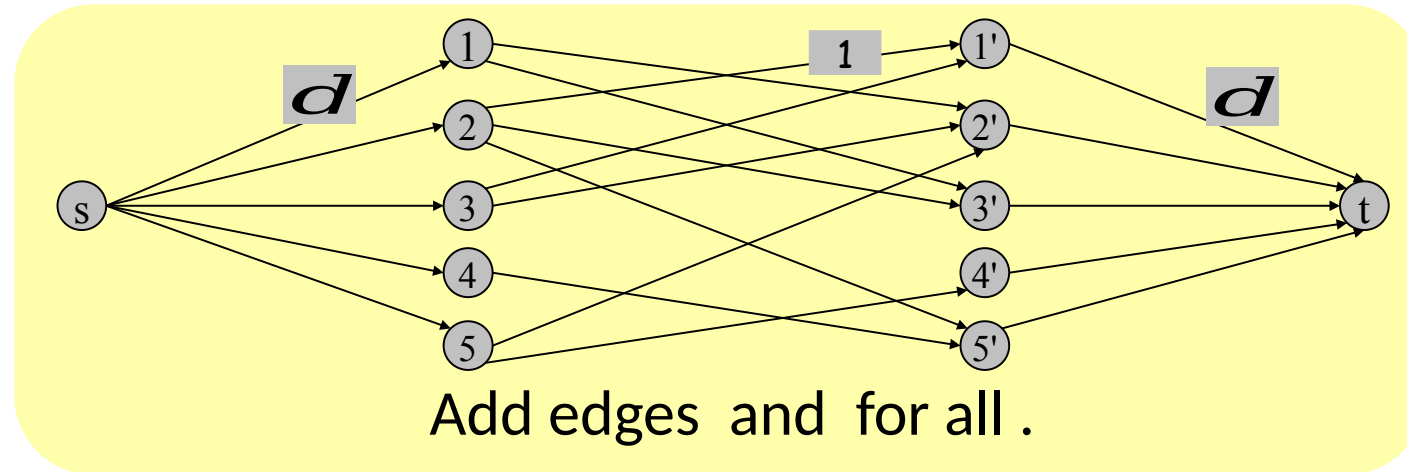
- agrees with f on graphs of degree d
- Sensitivity of f is low

- Release f using Laplace mechanism
- There exist Lipschitz extensions for all real-valued functions [McShane 34]
- We design Lipschitz extensions that can be computed efficiently using linear and convex programming for
 - subgraph counts [Kasiviswanathan Nissim Raskhodnikova Smith 13]
 - degree distribution [Raskhodnikova Smith 16, Day Li Lyu 16]
 - number of connected components [Kalemaj Raskhodnikova Smith Tsourakakis]



Lipschitz extension of : flow graph

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

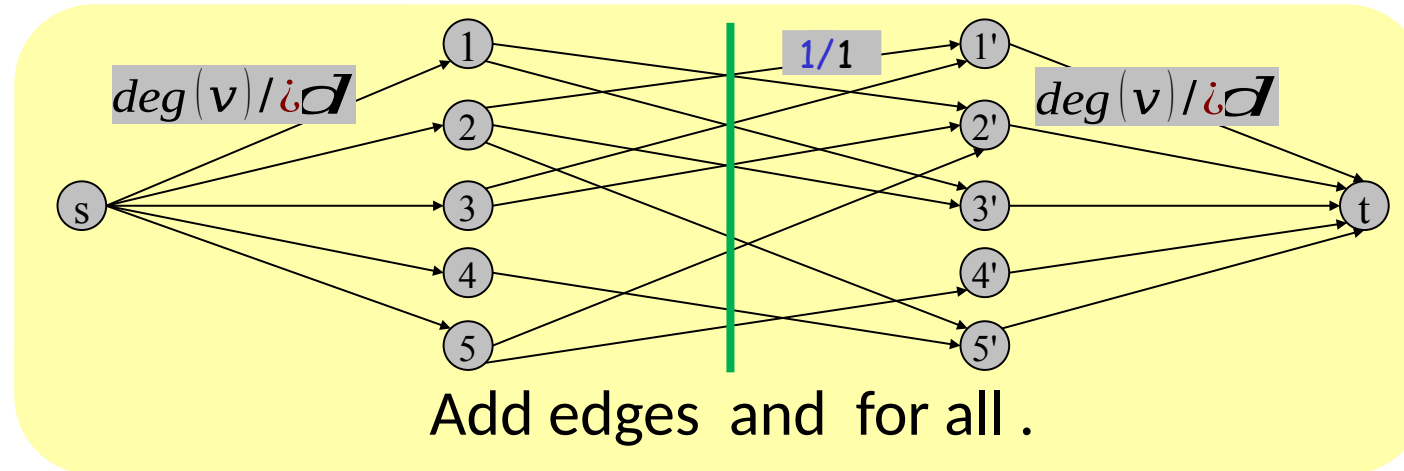


(G) is the value of the maximum flow in this graph.

Lemma. $(G)/2$ is a Lipschitz extension of .

Lipschitz extension of : flow graph

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



(G) is the value of the maximum flow in this graph.

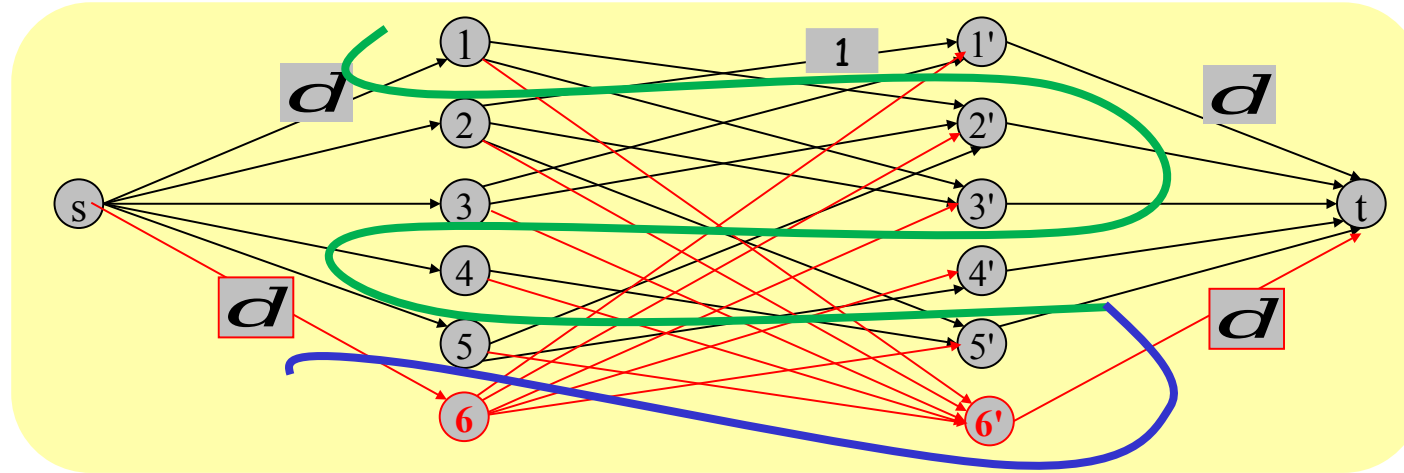
Lemma. $(G)/2$ is a Lipschitz extension of .

Proof: (1) $(G) =$ for all G of degree at most

(2) is small

Lipschitz extension of : flow graph

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



(G) is the value of the maximum flow in this graph.

Lemma. $(G)/2$ is a Lipschitz extension of .

Proof: (1) $(G) =$ for all G of degree at most

(2) $= 2$

Lipschitz extensions via linear programs

For a graph $G=(V, E)$, define **LP** with variables for all triangles :

$$\begin{aligned} & \text{Maximize} \sum_{T = \triangle \text{ of } G} x_T \\ & \text{for all triangles} \\ & \text{for all nodes} \end{aligned}$$

(G) is the value of **LP**.

Lemma. **(G)** is a Lipschitz extension of .

- Computable efficiently.
- Can be generalized to other counting queries.

Summary

- Accurate subgraph counts for realistic graphs can be computed by node-private algorithms
 - Use Lipschitz extensions and linear programming
 - One can choose a “good” value of cut-off degree privately
 - It is one example of many graph statistics that node-private algorithms do well on.