# Revealing Network Structure, Confidentially (Improved Rates for Node-Private Graphon Estimation)

Ilias Zadik<sup>1</sup>, joint work with Christian Borgs<sup>2</sup>, Jennifer Chayes<sup>2</sup> and Adam Smith<sup>3</sup>

<sup>1</sup>Massachusetts Institute of Technology (MIT), <sup>2</sup>Microsoft Research (MSR) and <sup>3</sup>Boston University (BU)

59th Symposium of Foundations of Computer Science (FOCS), 2018

Borgs, Chayes, Smith, Zadik (MIT)

## Introduction

Large and complicated networks arise everywhere in society! For example,

- the Facebook graph,
- the disease transmission graph
- the collaboration graph
- and many others..

< 回 > < 回 > < 回 >

## Introduction

Large and complicated networks arise everywhere in society! For example,

- the Facebook graph,
- the disease transmission graph
- the collaboration graph
- and many others..

**Analysis of Networks:** Important across fields (sociology, medicine etc), rich in theory (random graphs, graph algorithms etc)

#### Introduction

Large and complicated networks arise everywhere in society! For example,

- the Facebook graph,
- the disease transmission graph
- the collaboration graph
- and many others..

**Analysis of Networks:** Important across fields (sociology, medicine etc), rich in theory (random graphs, graph algorithms etc)

**Privacy on Networks:** Huge concern (e.g. Cambridge Analytica Scandal) and also rich in theory. Many open questions for networks!!

# This work: Limits of Network Estimation under Privacy

New algorithms and impossibility results

for estimating complex network models, subject to rigorous **privacy constraints** (**node differentially privacy**.)

# This work: Limits of Network Estimation under Privacy

**New algorithms** and **impossibility results** for estimating complex network models, subject to rigorous **privacy constraints** (node differentially privacy.)

(1) Stochastic Block Model-Estimation of probability matrix:
 -new analysis of recent private algorithm (BCS'15)
 -matches in many regimes the optimal non-private estimation rate

# This work: Limits of Network Estimation under Privacy

New algorithms and impossibility results for estimating complex network models, subject to rigorous privacy constraints (node differentially privacy.)

(1) Stochastic Block Model-Estimation of probability matrix:
 -new analysis of recent private algorithm (BCS'15)
 -matches in many regimes the optimal non-private estimation rate

(2) Erdos-Renyi-Estimation of probability p:
 -Compute tightly the optimal estimation rate
 -Uses a novel extension lemma, potentially of broad use

ヘロト ヘヨト ヘヨト

# Node Differential Private Algorithms

Intuition: If n-vertex G, G' differ on one user's (node's) data then the outputs of the algorithm are close (in distribution).

- 4 回 ト 4 ヨ ト 4 ヨ ト

# Node Differential Private Algorithms

Intuition: If n-vertex G, G' differ on one user's (node's) data then the outputs of the algorithm are close (in distribution).

**Node-neighbors**: We call G, G' node-neighbors if they differ only on the neighborhood of one node.



# Node Differential Private Algorithms

Intuition: If n-vertex G, G' differ on one user's (node's) data then the outputs of the algorithm are close (in distribution).

**Node-neighbors**: We call G, G' node-neighbors if they differ only on the neighborhood of one node.



#### Definition

A randomized  ${\mathcal A}$  on n-vertex graphs is  $\epsilon\text{-node-DP}$  if for node-neighbors G, G'and S,

$$\exp\left(-\epsilon\right)\mathbb{P}\left(\mathcal{A}(\mathsf{G}')\in\mathsf{S}\right)\leq\mathbb{P}\left(\mathcal{A}(\mathsf{G})\in\mathsf{S}\right)\leq\exp\left(\epsilon\right)\mathbb{P}\left(\mathcal{A}(\mathsf{G}')\in\mathsf{S}\right).$$

**1-SBM (Erdos Renyi)** G(n, p): n nodes every edge appears independently with probability p.

イロト 不得下 イヨト イヨト

**1-SBM (Erdos Renyi)** G(n, p): n nodes every edge appears independently with probability p.



**k-SBM**, G(n, B), for sym.  $B \in [0, 1]^{k \times k}$ : n nodes, k **groups** (node choice u.a.r.), each edge between  $v_i, v_j$  with probability  $B_{group}(v_i), group(v_j)$ .



**k-SBM**, G(n, B), for sym.  $B \in [0, 1]^{k \times k}$ : n nodes, k **groups** (node choice u.a.r.), each edge between  $v_i, v_j$  with probability  $B_{group(v_i),group(v_i)}$ .

Constraint (!) : ( $\rho$ -sparse) k-SBM, G(n, B), where B  $\in [0, \rho]^{k \times k}$ .

(Vast literature - planted bisection, planted clique, graph limits etc)



・ ロ ト ・ 同 ト ・ 三 ト ・ 三 ト

**Task:** From one sample G from G(n, B) estimate B using an  $\epsilon$ -node-DP estimator A.

э

イロト イヨト イヨト イヨト

**Task:** From one sample G from G(n, B) estimate B using an  $\epsilon$ -node-DP estimator A.

Each A has (worst-case over B) error

$$\operatorname{err}(\mathcal{A}) = \max_{\mathsf{B} \in [0,\rho]^{k \times k}} \mathbb{E}_{\mathsf{G} \sim \mathsf{G}(\mathsf{n},\mathsf{B})} \left[ \frac{1}{\mathsf{k}^2} \| \mathcal{A}(\mathsf{G}) - \mathsf{B} \|_2^2 \right].$$

**Task:** From one sample G from G(n, B) estimate B using an  $\epsilon$ -node-DP estimator A.

Each A has (worst-case over B) error

$$err(\mathcal{A}) = \max_{B \in [0,\rho]^{k \times k}} \mathbb{E}_{G \sim G(n,B)} \left[ \frac{1}{k^2} \| \mathcal{A}(G) - B \|_2^2 \right].$$

The Estimation Rate

$$\mathsf{R}_{\mathsf{k}}\left(\epsilon\right) = \min_{\mathcal{A} \ \epsilon - \mathsf{node-DP}} \mathsf{err}(\mathcal{A}).$$

・ロト ・四ト・ モン・ モン

**Task:** From one sample G from G(n, B) estimate B using an  $\epsilon$ -node-DP estimator A.

Each  $\mathcal{A}$  has (worst-case over B) error

$$err(\mathcal{A}) = \max_{B \in [0,\rho]^{k \times k}} \mathbb{E}_{G \sim G(n,B)} \left[ \frac{1}{k^2} \| \mathcal{A}(G) - B \|_2^2 \right].$$

The Estimation Rate

$$\mathsf{R}_{\mathsf{k}}(\epsilon) = \min_{\mathcal{A} \ \epsilon - \mathsf{node-DP}} \mathsf{err}(\mathcal{A}).$$

(For agnostic learning see paper!)

< ロ > < 同 > < 回 > < 回 > < 回 > <

# k-SBM Upper Bound

#### Theorem (informal)

For any  $\epsilon > 0$ ,

$$\mathcal{R}_{\mathsf{k}}(\epsilon) = \mathsf{O}\left(\rho\left(\frac{\mathsf{k}^{2}}{\mathsf{n}^{2}} + \frac{\mathsf{log}\,\mathsf{k}}{\mathsf{n}}\right)\right) + \mathsf{O}\left(\rho^{2}\frac{(\mathsf{k}-1)^{2}\,\mathsf{log}\,\mathsf{n}}{\mathsf{n}\epsilon} + \frac{1}{\mathsf{n}^{2}\epsilon^{2}}\right).$$

э

イロト イヨト イヨト イヨト

# k-SBM Upper Bound

#### Theorem (informal)

For any  $\epsilon > 0$ ,

$$\mathcal{R}_{\mathsf{k}}(\epsilon) = \mathsf{O}\left(\rho\left(\frac{\mathsf{k}^2}{\mathsf{n}^2} + \frac{\mathsf{log}\,\mathsf{k}}{\mathsf{n}}\right)\right) + \mathsf{O}\left(\rho^2 \frac{(\mathsf{k}-1)^2 \,\mathsf{log}\,\mathsf{n}}{\mathsf{n}\epsilon} + \frac{1}{\mathsf{n}^2 \epsilon^2}\right).$$

- Intuition:  $\frac{k^2}{n^2}$  parametric rate for B,  $\frac{\log k}{n} = \frac{\log k^n}{n^2}$  combinatorial rate
- Via a new detailed analysis of an  $\epsilon$ -node-DP algorithm proposed in (BCS '15).

イロト イヨト イヨト ・

## k-SBM Upper Bound: Optimality in many regimes

Theorem (informal)  
For any 
$$\epsilon > 0$$
,  
 $\mathcal{R}_{k}(\epsilon) = \underbrace{O\left(\rho\left(\frac{k^{2}}{n^{2}} + \frac{\log k}{n}\right)\right)}_{optimal \text{ non-private rate }(KTV'17)} + O\left(\rho^{2}\frac{(k-1)^{2}\log n}{n\epsilon} + \frac{1}{n^{2}\epsilon^{2}}\right).$ 

Borgs, Chayes, Smith, Zadik (MIT)

▲ □ ▶ ▲ □ ▶ ▲ □ ▶

# k-SBM Upper Bound: Optimality in many regimes

 $\begin{aligned} & \text{Theorem (informal)} \\ & \text{For any } \epsilon > 0, \\ & \mathcal{R}_{k}(\epsilon) = \underbrace{O\left(\rho\left(\frac{k^{2}}{n^{2}} + \frac{\log k}{n}\right)\right)}_{optimal \text{ non-private rate } (KTV'17)} + O\left(\rho^{2}\frac{(k-1)^{2}\log n}{n\epsilon} + \frac{1}{n^{2}\epsilon^{2}}\right). \end{aligned}$ 

Comments:

- (GLZ'14), (MS'17), (KTV'17): Optimal *e*-independent part.
- Many regimes (e.g.  $\epsilon$ , k constant and  $\frac{1}{n} < \rho < \frac{1}{\log n}$ ): -(BCS'15) algorithm, optimal over **all** algorithms! -**No additional error** with privacy!

(日)

# A lower bound for $k\geq 2$

Suppose each node  $i \in [n]$  chooses the group in a close to uniform way. (Say each group has probability in  $[\frac{1}{4k}, \frac{4}{k}]$ .)

Proposition (informal)

For  $k \ge 2$  and any  $\epsilon > 0$ ,

$$\mathcal{R}^*_{\mathsf{k}}(\epsilon) = \Omega\left(rac{1}{\mathsf{n}^2\epsilon^2}
ight)$$
 ,

where  $\mathcal{R}_{k}^{*}$  stands for the rate for the new variant of the SBM.

▲ □ ▶ ▲ □ ▶ ▲ □ ▶

# A lower bound for $k\geq 2$

Suppose each node  $i\in[n]$  chooses the group in a close to uniform way. (Say each group has probability in  $[\frac{1}{4k},\frac{4}{k}]$ .)

Proposition (informal)

For  $k \ge 2$  and any  $\epsilon > 0$ ,

$$\mathcal{R}^*_{\mathsf{k}}(\epsilon) = \Omega\left(rac{1}{\mathsf{n}^2\epsilon^2}
ight)$$
 ,

where  $\mathcal{R}_{k}^{*}$  stands for the rate for the new variant of the SBM.

Proof: reduction to privately estimating  $q \in [0, 1]$  out of n samples from Bern(q).

イロト イヨト イヨト ・

Observe simply a G(n, p): estimate **privately** p!

< □ > < 同 > < 回 > < 回 > < 回 >

Observe simply a G(n, p): estimate **privately** p!

$$\mathcal{R}_1(\epsilon) = ?$$

Observe simply a G(n, p): estimate privately p!

$$\Omega\left(\frac{1}{\mathsf{n}^2} + \frac{1}{\mathsf{n}^4\epsilon^2}\right) = \mathcal{R}_1(\epsilon) = O\left(\frac{1}{\mathsf{n}^2} + \frac{1}{\mathsf{n}^2\epsilon^2}\right).$$

**Upper bound** by main result, Laplace noise to edge density, median of degrees etc.

Lower bounds, by vanilla methods such as packing arguments.

Observe simply a G(n, p): estimate **privately** p!

$$\Omega\left(\frac{1}{\mathsf{n}^2} + \frac{1}{\mathsf{n}^4\epsilon^2}\right) = \mathcal{R}_1(\epsilon) = O\left(\frac{1}{\mathsf{n}^2} + \frac{1}{\mathsf{n}^2\epsilon^2}\right).$$

**Upper bound** by main result, Laplace noise to edge density, median of degrees etc.

Lower bounds, by vanilla methods such as packing arguments.

What is the true  $\epsilon$ -dependent rate?!

The case 
$$k = 1$$
:  $\frac{1}{n^4 \epsilon^2} \le \epsilon - \text{dep.} \le \frac{1}{n^2 \epsilon^2}$ 



Many novel techniques including a general extension lemma (next slide)!

Borgs, Chayes, Smith, Zadik (MIT)

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

The case 
$$k = 1$$
:  $\frac{1}{n^4 \epsilon^2} \le \epsilon - \text{dep.} \le \frac{1}{n^2 \epsilon^2}$ 

#### Theorem

For 
$$\epsilon > \frac{\log n}{n}$$
,

$$\mathcal{R}_1(\epsilon) = O(\frac{1}{n^2} + \frac{\log n}{n^3\epsilon^2}).$$

Many novel techniques including a general extension lemma (next slide)!

#### Proposition $(n^3 \text{ is tight!})$

Furthermore, if G is sampled u.a.r. from graphs with a fixed number of edges (conditional Erdos Renyi) for  $\epsilon$  constant,

$$\mathcal{R}_1'(\epsilon) = \Omega(\frac{1}{\mathsf{n}^3\epsilon^2}).$$

Borgs, Chayes, Smith, Zadik (MIT)

イロト イヨト イヨト イヨト

## The extension lemma: beyond networks

Technical challenge with *designing* differential private algorithms:

- Privacy constraint should hold for any pair of datasets
- Accuracy guarantee suffice to hold for typical datasets of our input distribution.

<日<br />
<</p>

## The extension lemma: beyond networks

Technical challenge with *designing* differential private algorithms:

- Privacy constraint should hold for any pair of datasets
- Accuracy guarantee suffice to hold for typical datasets of our input distribution.

*Key contribution:* Suffices to be **private** only for **typical** datasets of our input distribution!

# The extension lemma: beyond networks

Technical challenge with *designing* differential private algorithms:

- Privacy constraint should hold for any pair of datasets
- Accuracy guarantee suffice to hold for typical datasets of our input distribution.

*Key contribution:* Suffices to be **private** only for **typical** datasets of our input distribution!

Proposition ("Extending Private Algorithms at  $\epsilon$ -cost")

Let  $\hat{\mathcal{A}} \in DP$  on a subset of the input space  $\mathcal{H} \subseteq \mathcal{M}$ . Then there exists  $\mathcal{A}$  defined on  $\mathcal{M}$  which is 1)  $2\epsilon$ -DP on  $\mathcal{M}$  and 2) for every  $D \in \mathcal{H}$ ,  $\mathcal{A}(D) \stackrel{d}{=} \hat{\mathcal{A}}(D)$ .

Generalizes "extensions" from (KNRS'13), (BBDS'13), (CZ'13), (BCS'15), (RS'15).

Borgs, Chayes, Smith, Zadik (MIT)

(1) We focus on optimal private estimation of **Stochastic Block Model** and **Erdos Renyi** models.

< □ > < 同 > < 回 > < 回 > < 回 >

- (1) We focus on optimal private estimation of **Stochastic Block Model** and **Erdos Renyi** models.
- (2) Stochastic Block Model: new analysis of existing algorithm (BCS'15) matches optimal non-private rate in many regimes. Graphons (k-SBM for k → +∞) and agnostic learning in the paper!

<日<br />
<</p>

- (1) We focus on optimal private estimation of **Stochastic Block Model** and **Erdos Renyi** models.
- (2) Stochastic Block Model: new analysis of existing algorithm (BCS'15) matches optimal non-private rate in many regimes. Graphons (k-SBM for k → +∞) and agnostic learning in the paper!
- (3) Erdos-Renyi: "almost" tight optimal rate.

<日<br />
<</p>

- (1) We focus on optimal private estimation of **Stochastic Block Model** and **Erdos Renyi** models.
- (2) Stochastic Block Model: new analysis of existing algorithm (BCS'15) matches optimal non-private rate in many regimes. Graphons (k-SBM for k → +∞) and agnostic learning in the paper!
- (3) Erdos-Renyi: "almost" tight optimal rate.
- (4) Proved an extension lemma potentially of broad use.

- (1) We focus on optimal private estimation of **Stochastic Block Model** and **Erdos Renyi** models.
- (2) Stochastic Block Model: new analysis of existing algorithm (BCS'15) matches optimal non-private rate in many regimes. Graphons (k-SBM for k → +∞) and agnostic learning in the paper!
- (3) Erdos-Renyi: "almost" tight optimal rate.
- (4) Proved an extension lemma potentially of broad use.

# Thank you!!

Borgs, Chayes, Smith, Zadik (MIT)

- ロ ト ・ 同 ト ・ 三 ト ・ 三 ト - -