

# LOCALLY PRIVATE GRAPH NEURAL NETWORKS

---

Sina Sajadmanesh  
Daniel Gatica-Perez

AI4Media Workshop on Explainability, Robustness and Privacy in AI  
June 2, 2021

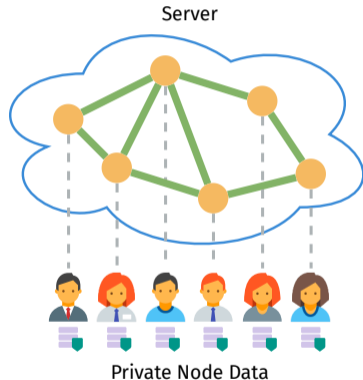
## Graph learning with **node data privacy**

### Setting:

- Graph topology is public to the server
- Node data (features and possibly labels) are private to nodes

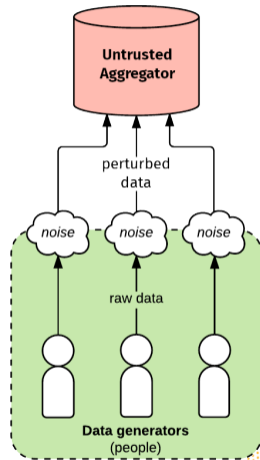
### Problem:

- How to learn a GNN without exposing private node data?



## Local Differential Privacy

- ▶ Every data holder perturbs their data using a **randomized mechanism**
- ▶ The aggregator collects and **aggregates** perturbed data to **estimate** the target statistics



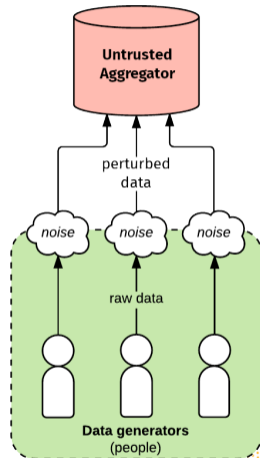
## Local Differential Privacy

- ▶ Every data holder perturbs their data using a **randomized mechanism**
- ▶ The aggregator collects and **aggregates** perturbed data to **estimate** the target statistics

## Definition

a randomized mechanism  $\mathcal{M}$  satisfies  $\epsilon$ -LDP if for all pairs of private data  $x_1$  and  $x_2$ , and for all outputs  $x'$  of  $\mathcal{M}$ , we have:

$$\Pr[\mathcal{M}(x_1) = x'] \leq e^\epsilon \Pr[\mathcal{M}(x_2) = x']$$

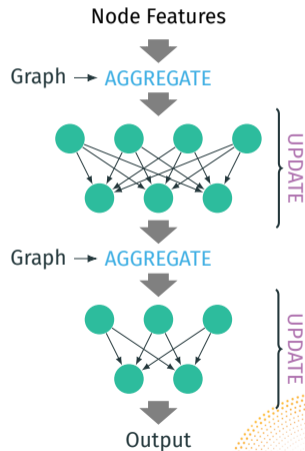


# WHY LOCAL DP?

GNNs are **message-passing** neural networks

**AGGREGATE**: nodes aggregate their neighbors' representation vector

**UPDATE**: a neural network generates new node representation from aggregated vectors



# WHY LOCAL DP?

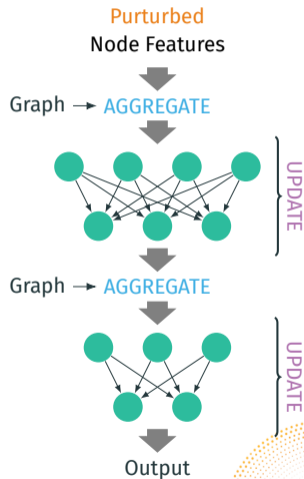
GNNs are **message-passing** neural networks

**AGGREGATE**: nodes aggregate their neighbors' representation vector

**UPDATE**: a neural network generates new node representation from aggregated vectors

Private neighborhood aggregation with LDP

- ▶ Node features are perturbed by **injecting noise**
- ▶ The neighborhood aggregation **cancels out** the noise



## High-dimensional features

- ▶ The total privacy budget of a node scales with the number of features
  - Keeping the total privacy budget small → **Too much noise!**

## High-dimensional features

- ▶ The total privacy budget of a node scales with the number of features
  - Keeping the total privacy budget small → **Too much noise!**

## **Our solution:** Multi-bit mechanism for multidimensional perturbation

- ▶ **Multi-bit Encoder:** perturb a random subset of node features and compress the output
- ▶ **Multi-bit Rectifier:** uncompress and de-bias encoded features



## Small neighborhoods

- ▶ Lots of the nodes have **too few neighbors**
  - Noise won't cancel out if the neighborhood size is small

## Small neighborhoods

- ▶ Lots of the nodes have **too few neighbors**
  - Noise won't cancel out if the neighborhood size is small

## Our solution: KProp linear convolution

- ▶ Expands the neighborhood to the nodes that are up to **K-hops away**
- ▶ Applies  $K$  consecutive **AGGREGATE**
- ▶ Can be prepended to any GNN architecture as a **feature denoising mechanism**

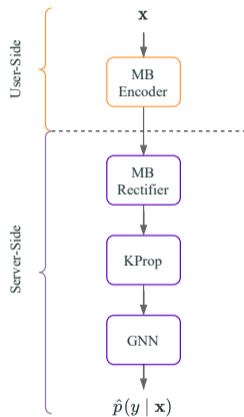
# LOCALLY PRIVATE GNN ARCHITECTURE

## User-Side:

1. Perturb node features using MB encoder
2. Send encoded features to server

## Server-Side:

3. De-bias encoded features with MB rectifier
4. De-noise rectifier's output using KProp
5. Train GNN on denoised features



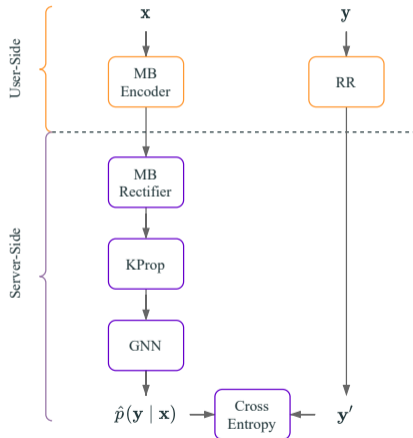
## Randomized Response for label differential privacy

- ▶ True label  $y$
- ▶ Perturbed label  $y'$
- ▶ Number of classes  $c$
- ▶ DP privacy budget  $\epsilon$

$$p(y' | y) = \begin{cases} \frac{e^\epsilon}{e^\epsilon + c - 1}, & \text{if } y' = y \\ \frac{1}{e^\epsilon + c - 1}, & \text{otherwise} \end{cases}$$

Trivial method: directly train GNN with noisy labels

- ▶ GNN severely overfits the noisy labels
- ▶ Poor generalization performance

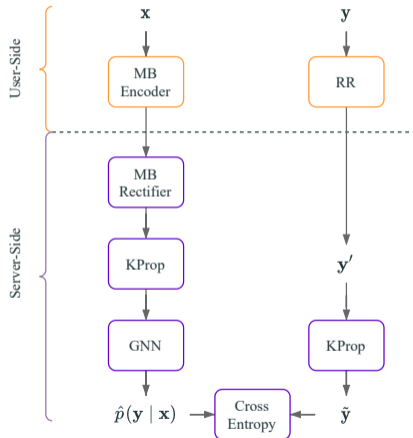


Trivial method: directly train GNN with noisy labels

- ▶ GNN severely overfits the noisy labels
- ▶ Poor generalization performance

**Key Idea:** use KProp to **denoise labels!**

- ▶ Apply KProp on one-hot encoded noisy labels
- ▶ Pick the label with highest value

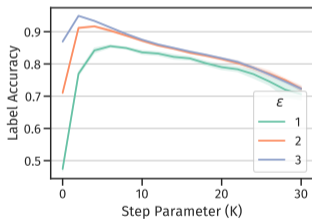


# DENOISING LABELS WITH KPROP

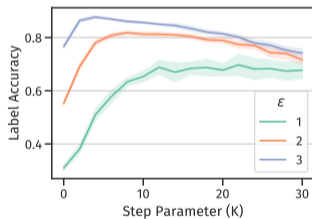
## Effect of KProp on label accuracy

- Accuracy between true label  $y$  and recovered label  $\tilde{y}$

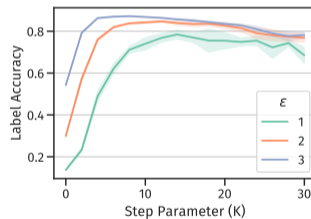
Facebook (4 classes)



Cora (7 classes)



LastFM (10 classes)

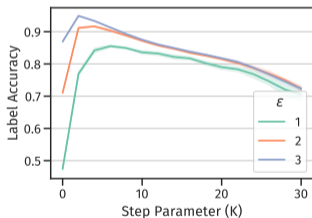


# DENOISING LABELS WITH KPROP

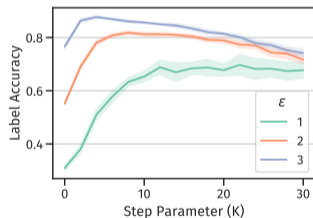
## Effect of KProp on label accuracy

- Accuracy between true label  $y$  and recovered label  $\tilde{y}$

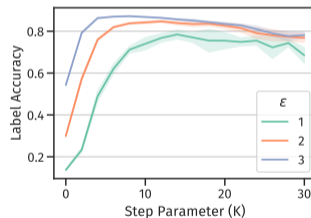
Facebook (4 classes)



Cora (7 classes)



LastFM (10 classes)



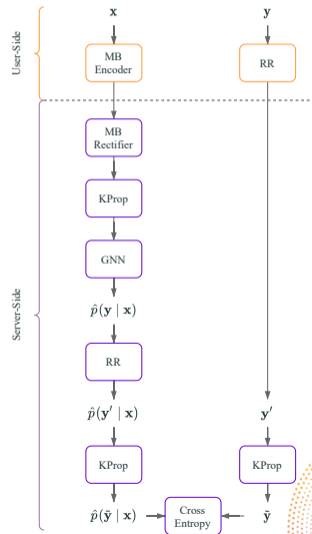
How to find best performing  $K$  **without clean validation data?**



# LABEL DENOISING WITH PROPAGATION

## Prevent absorbing noise in $\tilde{y}$

- ▶  $y$  is perturbed by RR and is given KProp to get  $\tilde{y}$
- ▶ Apply the **same process** on  $\hat{p}(y | x)$  to obtain  $\hat{p}(\tilde{y} | x)$
- ▶ Train  $\hat{p}(\tilde{y} | x)$  with  $\tilde{y}$



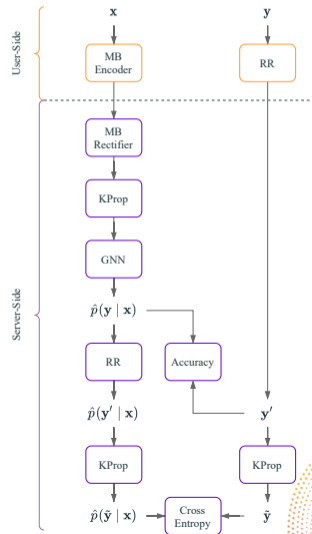
# LABEL DENOISING WITH PROPAGATION

## Prevent absorbing noise in $\tilde{y}$

- ▶  $y$  is perturbed by RR and is given KProp to get  $\tilde{y}$
- ▶ Apply the **same process** on  $\hat{p}(y | x)$  to obtain  $\hat{p}(\tilde{y} | x)$
- ▶ Train  $\hat{p}(\tilde{y} | x)$  with  $\tilde{y}$

## Prevent absorbing noise in $y'$

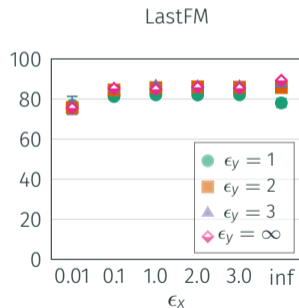
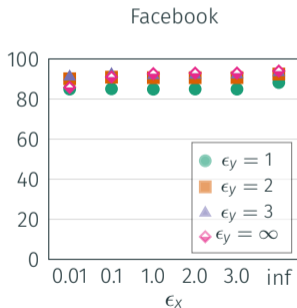
- ▶ RR gives an **upperbound** on label accuracy:  
$$Acc^* = p(y' = y) = \frac{e^c}{e^c + c - 1}$$
- ▶ Stop training when GNN's accuracy on  $y'$  goes beyond  $Acc^*$



# EXPERIMENTAL RESULTS

## LPGNN's performance under varying feature and label privacy budgets

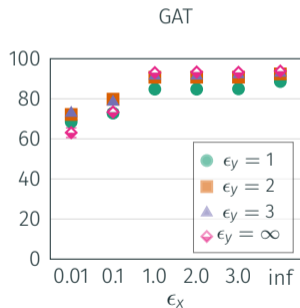
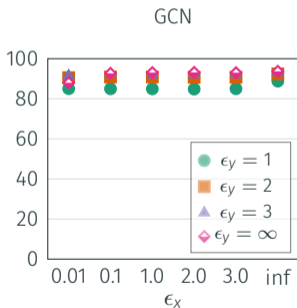
- Base GNN: GraphSAGE



# EXPERIMENTAL RESULTS

## Comparison of base GNN architectures

► Dataset: **Facebook**



## Comparison of different LDP mechanisms

- ▶ Base GNN: GraphSAGE
- ▶  $\epsilon_y = \infty$

DATASET	MECHANISM	$\epsilon_x = 0.01$	$\epsilon_x = 0.1$	$\epsilon_x = 1$	$\epsilon_x = 2$
CORA	1B	45.8 $\pm$ 3.3	62.3 $\pm$ 1.5	59.9 $\pm$ 2.7	58.5 $\pm$ 2.9
	LP	43.2 $\pm$ 3.1	57.8 $\pm$ 2.3	61.9 $\pm$ 3.1	58.1 $\pm$ 2.1
	AG	59.7 $\pm$ 2.3	62.7 $\pm$ 2.8	67.5 $\pm$ 3.0	77.2 $\pm$ 1.9
	MB	<b>68.0 <math>\pm</math> 2.9</b>	<b>64.6 <math>\pm</math> 3.2</b>	<b>83.9 <math>\pm</math> 0.4</b>	<b>84.0 <math>\pm</math> 0.3</b>
FACEBOOK	1B	57.0 $\pm$ 3.4	76.3 $\pm$ 1.6	86.1 $\pm$ 0.6	84.0 $\pm$ 1.3
	LP	54.2 $\pm$ 2.9	72.5 $\pm$ 2.1	85.4 $\pm$ 0.4	84.8 $\pm$ 1.6
	AG	78.2 $\pm$ 1.4	85.6 $\pm$ 0.7	92.0 $\pm$ 0.1	92.4 $\pm$ 0.2
	MB	<b>85.8 <math>\pm</math> 0.4</b>	<b>91.0 <math>\pm</math> 0.4</b>	<b>92.7 <math>\pm</math> 0.1</b>	<b>92.9 <math>\pm</math> 0.1</b>

# EXPERIMENTAL RESULTS

## Comparison of different learning algorithms

- ▶ Base GNN: GraphSAGE
- ▶  $\epsilon_x = 1$

DATASET	$\epsilon_y$	CROSS ENTROPY	FORWARD CORRECTION	DROP
CORA	0.5	18.6 $\pm$ 1.3	18.6 $\pm$ 2.5	<b>42.9 <math>\pm</math> 1.5</b>
	1.0	25.5 $\pm$ 1.7	37.1 $\pm$ 2.5	<b>69.3 <math>\pm</math> 1.2</b>
	2.0	52.9 $\pm$ 2.1	75.1 $\pm$ 1.0	<b>78.4 <math>\pm</math> 0.7</b>
FACEBOOK	0.5	50.9 $\pm$ 4.2	68.9 $\pm$ 1.3	<b>75.1 <math>\pm</math> 0.6</b>
	1.0	55.2 $\pm$ 1.3	73.8 $\pm$ 1.1	<b>84.9 <math>\pm</math> 0.2</b>
	2.0	81.6 $\pm$ 1.2	88.9 $\pm$ 0.2	<b>90.7 <math>\pm</math> 0.1</b>
LASTFM	0.5	21.1 $\pm$ 4.6	44.9 $\pm$ 5.3	<b>70.0 <math>\pm</math> 3.0</b>
	1.0	28.4 $\pm$ 2.5	58.5 $\pm$ 3.6	<b>82.1 <math>\pm</math> 1.0</b>
	2.0	56.8 $\pm$ 2.8	79.2 $\pm$ 1.3	<b>85.7 <math>\pm</math> 0.7</b>

## Summary

- ▶ Proposed a privacy-preserving GNN based on local differential privacy
  - Multi-bit mechanism for high-dimensional feature perturbation
  - KProp for feature and label denoising
  - Drop algorithm for learning with noisy labels
- ▶ Demonstrated promising results in terms of accuracy-privacy trade-off

## Future Work

- ▶ Protect privacy of graph topology

# THANK YOU!

---

 [sajadmanesh](#)

 [sajadmanesh@idiap.ch](mailto:sajadmanesh@idiap.ch)

