



# LOCALLY PRIVATE GRAPH NEURAL NETWORKS

---

Sina Sajadmanesh

Daniel Gatica-Perez

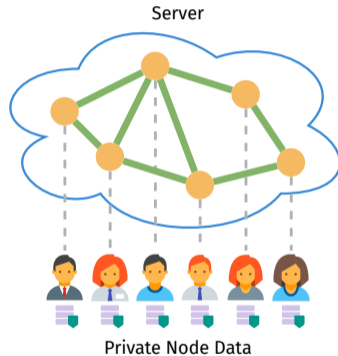
Idiap Research Institute  
EPFL

Graph Neural Networks User Group Meetup – July 29, 2021

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

### Assumptions:

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



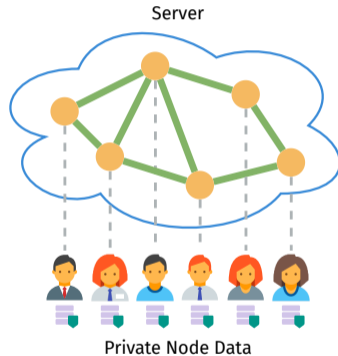
## Graph learning with **node data privacy**

### Assumptions:

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

### Problem:

- How to let the server train a GNN without giving up private node data?



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

### Assumptions:

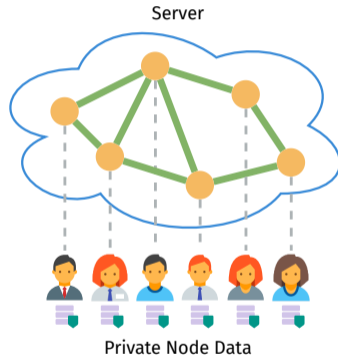
- Graph topology is public to the server
- Node data (features/labels) are private to nodes

### Problem:

- How to let the server train a GNN without giving up private node data?

### Solution:

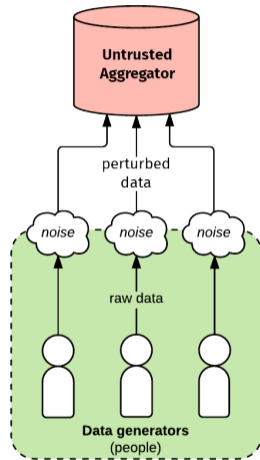
- Preserve the privacy of nodes using **Local Differential Privacy**



# LOCAL DIFFERENTIAL PRIVACY

## Procedure

- ▶ Data holders perturb their data using a **randomized mechanism**
- ▶ The aggregator **estimates** the target statistics by **aggregating** perturbed data
  - The noise cancels out through aggregation



# LOCAL DIFFERENTIAL PRIVACY

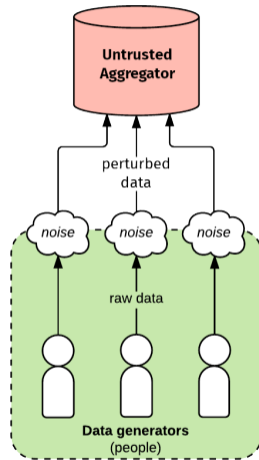
## Procedure

- ▶ Data holders perturb their data using a **randomized mechanism**
- ▶ The aggregator **estimates** the target statistics by **aggregating** perturbed data
  - The noise cancels out through aggregation

## 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:

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

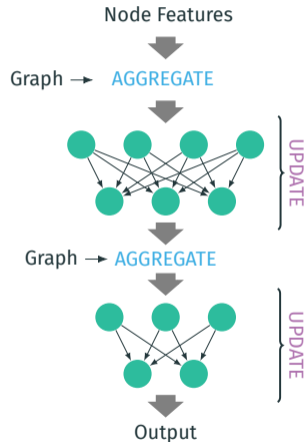


# WHY LOCAL DP?

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?

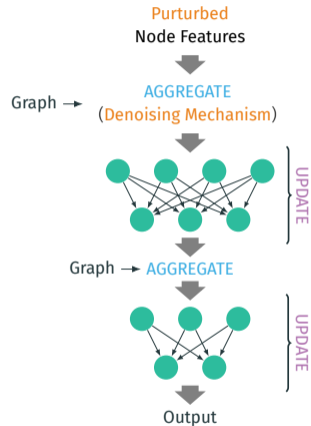
GNNs are **message-passing** neural networks

**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!**

## Small neighborhoods

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

# CHALLENGE: HIGH-DIMENSIONAL FEATURES

## Multi-bit mechanism for multidimensional perturbation

- ▶ **Multi-bit Encoder:** for feature selection, perturbation, and compression
  - Runs at **user-side**
  
- ▶ **Multi-bit Rectifier:** for decompression and de-biasing
  - Runs at **server-side**

# CHALLENGE: HIGH-DIMENSIONAL FEATURES

## Multi-bit mechanism for multidimensional perturbation

- ▶ **Multi-bit Encoder:** for feature selection, perturbation, and compression
  - Runs at **user-side**
  - Randomly perturb  **$m$  out of  $d$  features** with  $\epsilon/m$  privacy budget using **1-bit mechanism:**

$$x_{V,i}^* \sim \text{Bernoulli} \left( \frac{1}{e^{\epsilon/m} + 1} + \frac{x_{V,i} - \alpha}{\beta - \alpha} \cdot \frac{e^{\epsilon/m} - 1}{e^{\epsilon/m} + 1} \right)$$

- ▶ **Multi-bit Rectifier:** for decompression and de-biasing
  - Runs at **server-side**

# CHALLENGE: HIGH-DIMENSIONAL FEATURES

## Multi-bit mechanism for multidimensional perturbation

- ▶ **Multi-bit Encoder:** for feature selection, perturbation, and compression
  - Runs at **user-side**
  - Randomly perturb  **$m$  out of  $d$  features** with  $\epsilon/m$  privacy budget using **1-bit mechanism:**

$$x_{V,i}^* \sim \text{Bernoulli} \left( \frac{1}{e^{\epsilon/m} + 1} + \frac{x_{V,i} - \alpha}{\beta - \alpha} \cdot \frac{e^{\epsilon/m} - 1}{e^{\epsilon/m} + 1} \right)$$

- Map 1-bit output to either -1 or 1, return 0 for non-perturbed
- ▶ **Multi-bit Rectifier:** for decompression and de-biasing
  - Runs at **server-side**

# CHALLENGE: HIGH-DIMENSIONAL FEATURES

## Multi-bit mechanism for multidimensional perturbation

- ▶ **Multi-bit Encoder:** for feature selection, perturbation, and compression
  - Runs at **user-side**
  - Randomly perturb  **$m$  out of  $d$  features** with  $\epsilon/m$  privacy budget using **1-bit mechanism:**

$$x_{V,i}^* \sim \text{Bernoulli} \left( \frac{1}{e^{\epsilon/m} + 1} + \frac{x_{V,i} - \alpha}{\beta - \alpha} \cdot \frac{e^{\epsilon/m} - 1}{e^{\epsilon/m} + 1} \right)$$

- Map 1-bit output to either -1 or 1, return 0 for non-perturbed
- ▶ **Multi-bit Rectifier:** for decompression and de-biasing
  - Runs at **server-side**
  - Encoder's output is **biased**

# CHALLENGE: HIGH-DIMENSIONAL FEATURES

## Multi-bit mechanism for multidimensional perturbation

- ▶ **Multi-bit Encoder:** for feature selection, perturbation, and compression
  - Runs at **user-side**
  - Randomly perturb  **$m$  out of  $d$  features** with  $\epsilon/m$  privacy budget using **1-bit mechanism:**

$$x_{v,i}^* \sim \text{Bernoulli} \left( \frac{1}{e^{\epsilon/m} + 1} + \frac{x_{v,i} - \alpha}{\beta - \alpha} \cdot \frac{e^{\epsilon/m} - 1}{e^{\epsilon/m} + 1} \right)$$

- Map 1-bit output to either -1 or 1, return 0 for non-perturbed
- ▶ **Multi-bit Rectifier:** for decompression and de-biasing
  - Runs at **server-side**
  - Encoder's output is **biased**
  - **De-biases** encoded features by reversing the encoder's mapping:

$$x'_{v,i} = \frac{d(\beta - \alpha)}{2m} \cdot \frac{e^{\epsilon/m} + 1}{e^{\epsilon/m} - 1} \cdot x_{v,i}^* + \frac{\alpha + \beta}{2}$$

### Our solution: KProp denoising layer

- ▶ Expands the neighborhood to the nodes that are up to **K-hops away**
- ▶ Applies  $K$  consecutive linear **AGGREGATE** functions
  - No non-linearity in between
- ▶ Can be prepended to any GNN architecture



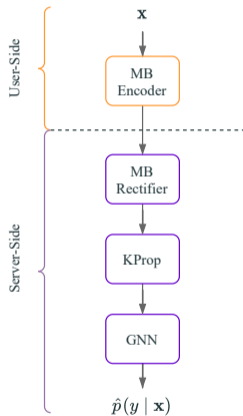
# LOCALLY PRIVATE GNN ARCHITECTURE

## 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



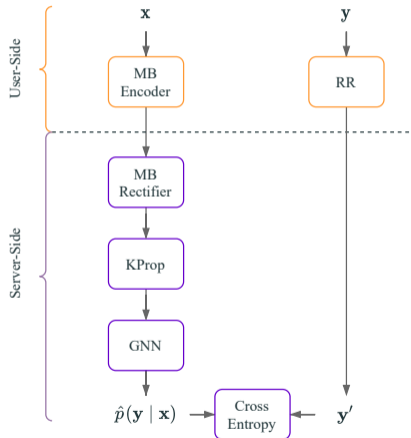
## Randomized Response for label differential privacy

- ▶ 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

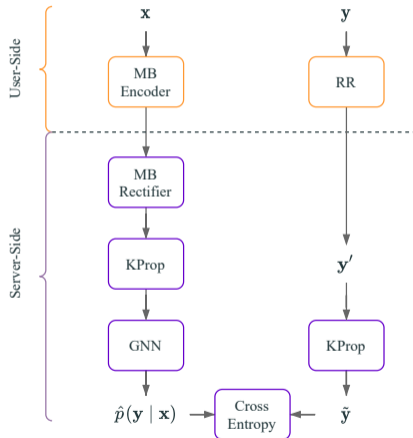


Trivial method: directly train GNN with noisy labels

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

**Idea:** use KProp for **label denoising**

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

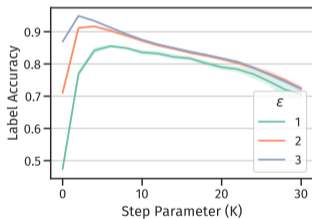


# DENOISING LABELS WITH KPROP

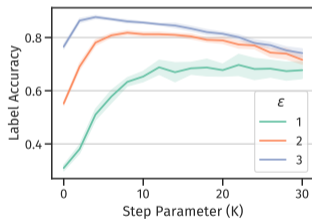
## Effect of KProp on label accuracy

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

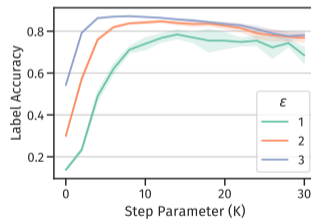
Facebook (4 classes)



Cora (7 classes)



LastFM (10 classes)

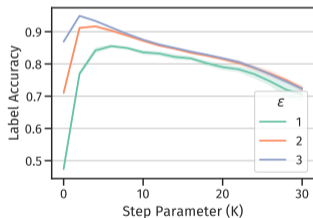


# DENOISING LABELS WITH KPROP

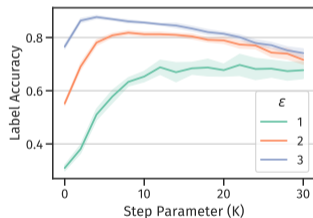
## Effect of KProp on label accuracy

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

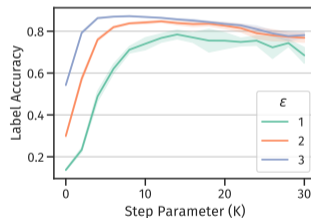
Facebook (4 classes)



Cora (7 classes)



LastFM (10 classes)

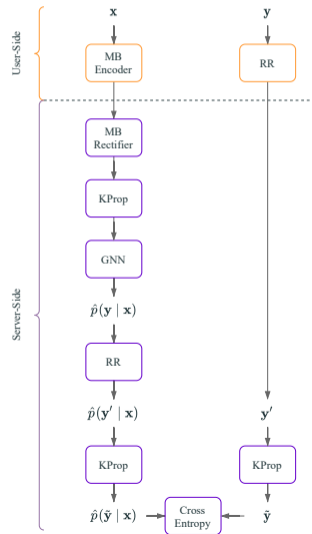


How to find best performing  $K$  **without clean validation data?**

# LABEL DENOISING WITH PROPAGATION

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

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



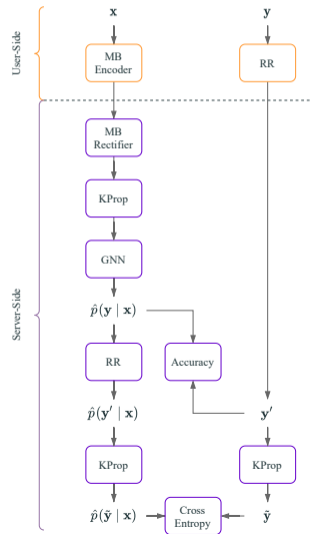
# LABEL DENOISING WITH PROPAGATION

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

- ▶  $y$  is perturbed by RR and is given to KProp to yield  $\tilde{y}$
- ▶ 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 $y'$

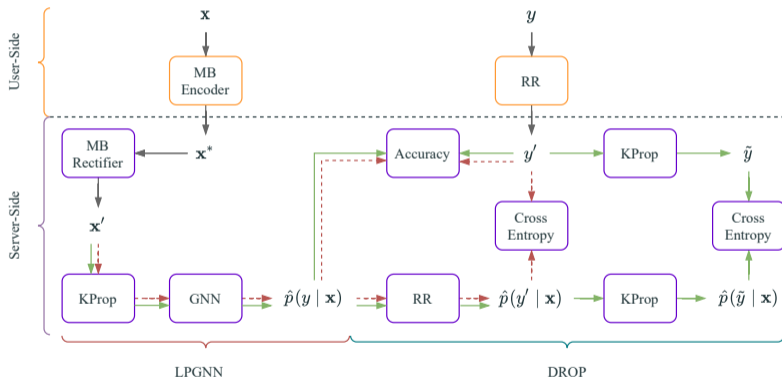
- ▶ RR gives an **upperbound** on label accuracy:  
$$Acc^* = p(y' = y) = \frac{e^c}{e^c + c - 1}$$
- ▶ Stop training when GNN's accuracy on  $y'$  goes beyond  $Acc^*$





# DROP ALGORITHM

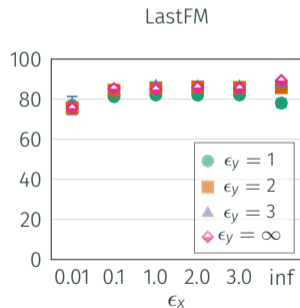
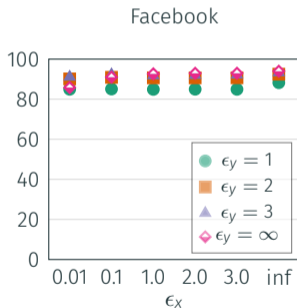
## Label Denoising with Propagation (Drop)



# EXPERIMENTAL RESULTS

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

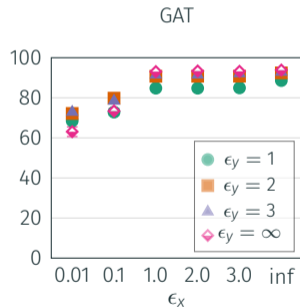
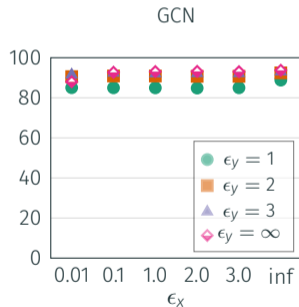
- ▶ Base GNN: **GraphSAGE**



# EXPERIMENTAL RESULTS

## Comparison of base GNN architectures

- Dataset: **Facebook**



## Comparison against ad-hoc features

- ▶ Base GNN: **GraphSAGE**
- ▶  $\epsilon_y = 1$

FEATURE	CORA	PUBMED	FACEBOOK	LASTFM
ONES	22.6 $\pm$ 5.0	38.9 $\pm$ 0.4	29.0 $\pm$ 1.4	19.6 $\pm$ 1.8
OHD	44.4 $\pm$ 3.5	52.5 $\pm$ 5.7	77.2 $\pm$ 0.3	66.4 $\pm$ 1.6
RND	26.4 $\pm$ 3.0	56.0 $\pm$ 1.3	35.2 $\pm$ 5.6	32.3 $\pm$ 6.3
MBM ( $\epsilon_x = 1$ )	<b>69.3 <math>\pm</math> 1.2</b>	<b>74.9 <math>\pm</math> 0.3</b>	<b>84.9 <math>\pm</math> 0.2</b>	<b>82.1 <math>\pm</math> 1.0</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
- ▶ GNN models demonstrate reasonable accuracy-privacy trade-off
  - Feature privacy almost comes for free in simpler models
  - Label privacy with low to moderate privacy budget gives acceptable accuracy

# THANK YOU!

---

 [sajadmanesh](#)

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

