

HDSI | Harvard Data
Science Initiative



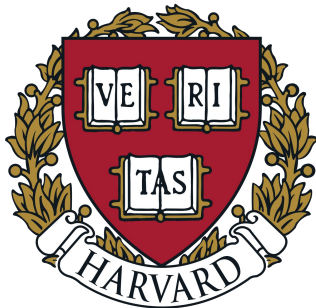
HARVARD
MEDICAL SCHOOL



GNNDelete: A General Strategy for Unlearning in Graph Neural Networks

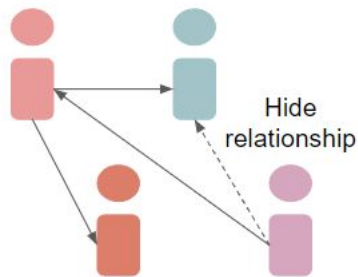
Jiali Cheng¹, George Dasoulas^{2,*}, Huan He^{2,*}, Chirag Agarwal³, Marinka Zitnik²

¹University of Massachusetts Lowell, ²Harvard University, ³Adobe

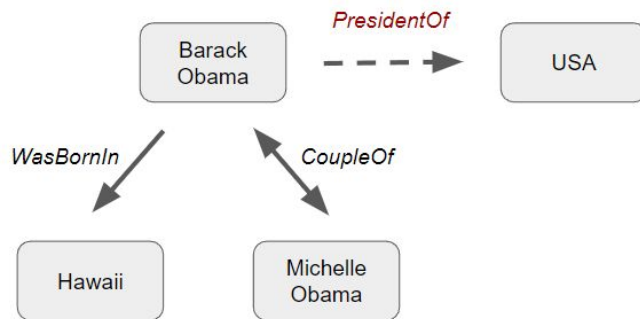


Motivation for Graph Unlearning

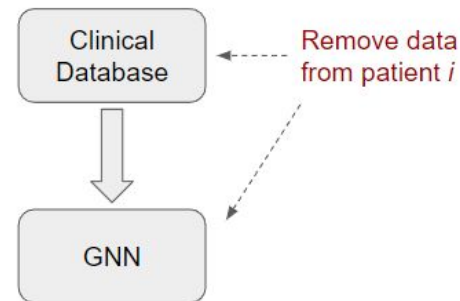
- Graph unlearning involves deleting graph elements such as nodes, node labels, and relationships from trained graph neural network (GNN) models
- Retraining models from scratch is not feasible. Needed are **efficient methods for model editing**



Graph elements become irrelevant or inaccurate



Underlying graphs evolving over time



Sensitive data and growing demands for privacy

Why is Graph Unlearning Challenging?

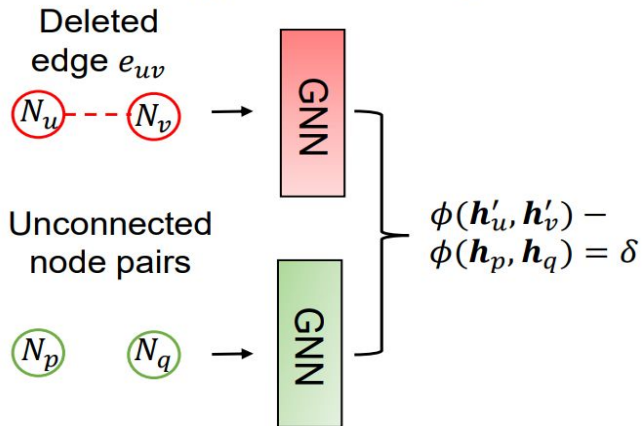
1. Graph elements exert strong influence on other elements with dependencies between nodes connected by edges
Challenge: Existing machine unlearning methods are unsuitable for data with underlying geometric and relational structure
2. Graph models make predictions by propagating messages across local neighborhoods
Challenge: Adversarial agents can infer the presence of graph elements from their local neighbors. Merely removing data from the graph is not sufficient
3. GNNs share model weights across many (often all) nodes or edges in the graph
Challenge: Naively perturbing model weights deteriorates model performance. Methods developed for other modalities are not suitable for graphs

Requirements for Successful Graph Deletion

Shown is a motivating example of deleting a single edge e_{uv}

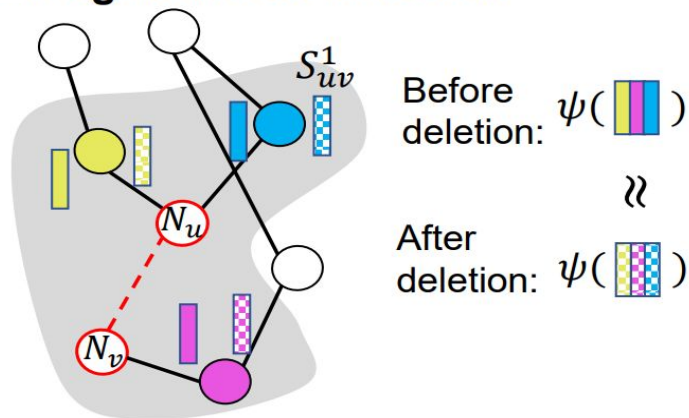
After deletion, GNNDelete treats e_{uv} as an unconnected node pair

Deleted Edge Consistency



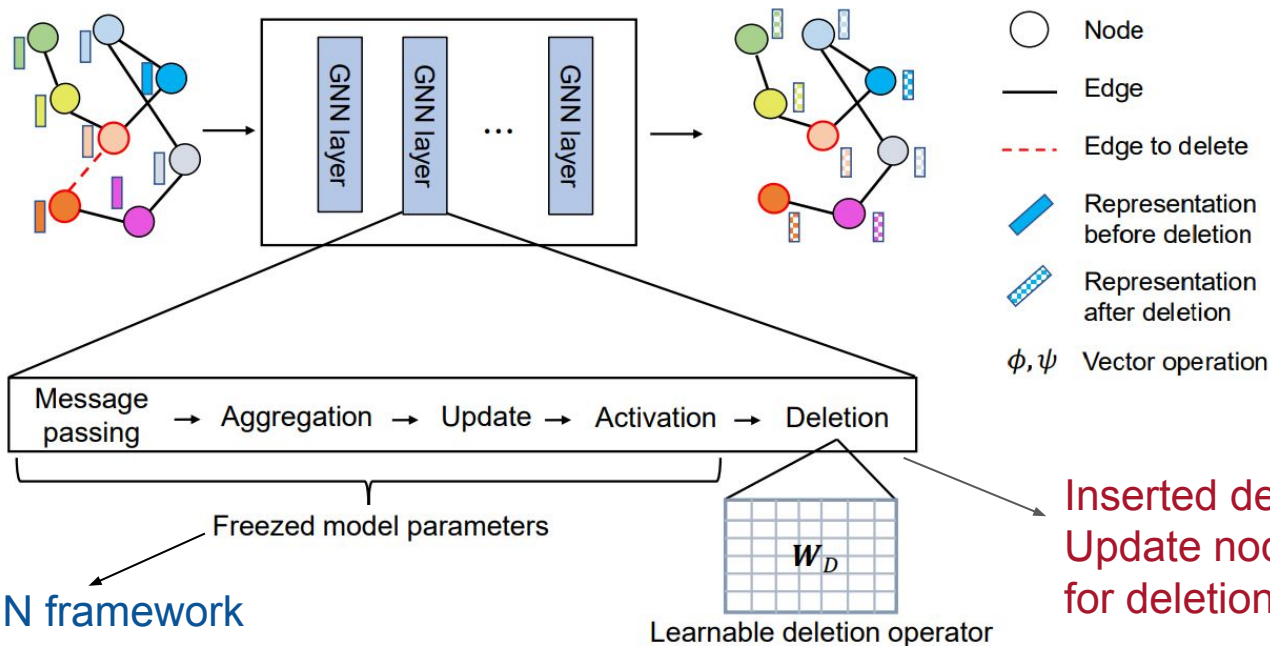
After deletion, GNNDelete keeps node embeddings close to original ones

Neighborhood Influence



Layer-wise Deletion Operator in GNNDelete

Overview of GNNDELETE



Standard GNN framework
Frozen during training

$$W_D = \arg \min_{W_D} \lambda L_{DEC} + (1 - \lambda) L_{NI}$$

Graph Neural Network Model Unlearning

$$\mathcal{L}_{\text{DEC}}^l = \mathcal{L}(\{[\mathbf{h}'_u; \mathbf{h}'_v] | e_{uv} \in \mathcal{E}_d\}, \{[\mathbf{h}^l_u; \mathbf{h}^l_v] | u, v \in_R \mathcal{V}\})$$

Deleted edges

Similar to

Unconnected node pairs

$$\mathcal{L}_{\text{NI}}^l = \mathcal{L}(\|_w \{\mathbf{h}'_w | w \in \mathcal{S}_{uv}^l / e_{uv}\}, \|_w \{\mathbf{h}^l_w | w \in \mathcal{S}_{uv}^l\})$$

Embeddings after deletion

Similar to

Embeddings before deletion

$$\mathbf{W}_D^{l*} = \arg \min_{\mathbf{W}_D^l} \mathcal{L}^l = \arg \min_{\mathbf{W}_D^l} \lambda \mathcal{L}_{\text{DEC}}^l + (1 - \lambda) \mathcal{L}_{\text{NI}}^l \quad \text{Local update}$$

$$\text{DEL}^l = \begin{cases} \phi & \text{if } w \in \mathcal{S}_{uv}^l \\ \mathbb{1} & \text{otherwise} \end{cases}$$

Only train deletion operators.
Other parameters are frozen.

Results 1: Excellent Performance on Edge Deletion Benchmarks

GNNDelete outperforms all baselines by 30.7 (\mathcal{E}_t) and 25.1 (\mathcal{E}_d) on average, other graph unlearning methods by 24.1% (\mathcal{E}_t) and 28.5 (\mathcal{E}_d) on average

Predict test edges accurately

Distinguish between deleted and non-deleted edges

Reference

| Model | GCN | | GAT | | R-GCN | | R-GAT | |
|--------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| | \mathcal{E}_t | \mathcal{E}_d | \mathcal{E}_t | \mathcal{E}_d | \mathcal{E}_t | \mathcal{E}_d | \mathcal{E}_t | \mathcal{E}_d |
| RETRAIN | 0.964 \pm 0.003 | 0.506 \pm 0.013 | 0.956 \pm 0.002 | 0.525 \pm 0.012 | 0.800 \pm 0.005 | 0.580 \pm 0.006 | 0.891 \pm 0.005 | 0.783 \pm 0.009 |
| GRADASCENT | 0.555 \pm 0.066 | 0.594 \pm 0.063 | 0.501 \pm 0.020 | 0.592 \pm 0.017 | 0.490 \pm 0.001 | 0.502 \pm 0.002 | 0.490 \pm 0.001 | 0.492 \pm 0.003 |
| D2D | 0.500 \pm 0.000 | 0.500 \pm 0.000 | 0.500 \pm 0.000 | 0.500 \pm 0.000 | 0.500 \pm 0.000 | 0.500 \pm 0.000 | 0.500 \pm 0.000 | 0.500 \pm 0.000 |
| GRAPHERASER | 0.527 \pm 0.002 | 0.500 \pm 0.000 | 0.538 \pm 0.013 | 0.500 \pm 0.000 | 0.512 \pm 0.003 | 0.500 \pm 0.000 | 0.545 \pm 0.015 | 0.500 \pm 0.000 |
| GRAPHEEDITOR | 0.776 \pm 0.025 | 0.432 \pm 0.009 | - | - | N/A | N/A | N/A | N/A |
| CERTUNLEARN | 0.718 \pm 0.032 | 0.475 \pm 0.011 | - | - | N/A | N/A | N/A | N/A |
| GNNDELETE | 0.934 \pm 0.002 | 0.748 \pm 0.006 | 0.914 \pm 0.007 | 0.774 \pm 0.015 | 0.751 \pm 0.006 | 0.845 \pm 0.007 | 0.893 \pm 0.002 | 0.786 \pm 0.004 |

Other graph unlearning methods

GNNDelete achieves the highest AUROC on both settings using different architectures

Results 2: Ablation Study

On the interplay of Deleted Edge Consistency (DEC) and Neighborhood Influence (NI)

$$\mathbf{W}_D^{l*} = \arg \min_{\mathbf{W}_D^l} \mathcal{L}^l = \arg \min_{\mathbf{W}_D^l} \lambda \mathcal{L}_{\text{DEC}}^l + (1 - \lambda) \mathcal{L}_{\text{NI}}^l$$

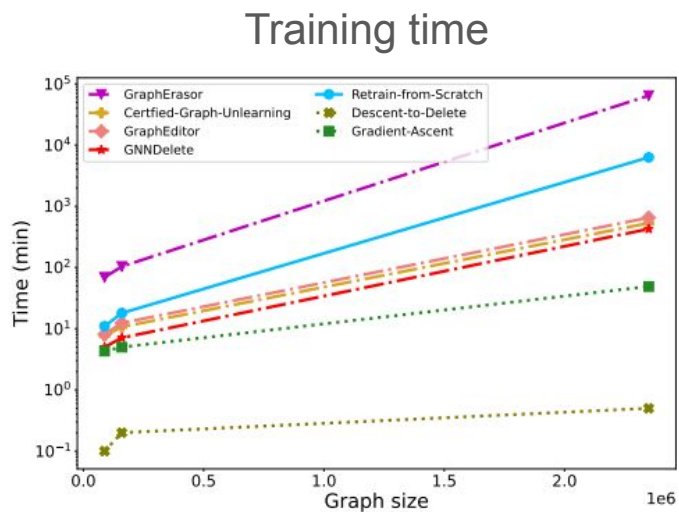
Considering performance on both \mathcal{E}_t and \mathcal{E}_d .

Deleted Edge Consistency
and **Neighborhood Influence**
are both necessary for
successful unlearning on
graphs

| λ | AUROC on \mathcal{E}_t | AUROC on \mathcal{E}_d | Avg. AUROC (Gap) |
|-----------|--------------------------|--------------------------|----------------------|
| 0.0 | 0.964 \pm 0.003 | 0.492 \pm 0.012 | 0.728 (0.473) |
| 0.2 | 0.961 \pm 0.003 | 0.593 \pm 0.011 | 0.777 (0.368) |
| 0.4 | 0.950 \pm 0.005 | 0.691 \pm 0.010 | 0.821 (0.259) |
| 0.5 | 0.934 \pm 0.002 | 0.748 \pm 0.006 | 0.841 (0.185) |
| 0.6 | 0.927 \pm 0.001 | 0.739 \pm 0.006 | 0.834 (0.188) |
| 0.8 | 0.893 \pm 0.003 | 0.759 \pm 0.008 | 0.823 (0.134) |
| 1.0 | 0.858 \pm 0.004 | 0.757 \pm 0.004 | 0.808 (0.101) |

Results 3: GNNDelete is Computationally Efficient

GNNDelete demonstrates efficiency in terms of both its training time and the number of trainable parameters it requires



Trainable parameters

| Model | OGB-Collab | OGB-BioKG |
|------------------|--------------|--------------|
| RETRAIN | 5,216 | 12,009,792 |
| GRADASCENT | 5,216 | 12,009,792 |
| D2D | 5,216 | 12,009,792 |
| GRAPHERASER | 52,160 | 120,097,920 |
| GRAPHEEDITOR | 5,216 | N/A |
| CERTUNLEARN | 5,216 | N/A |
| GNNDELETE | 5,120 | 5,120 |

Compared to GraphEraser: GNNDelete saves **~9x training time** and **~10x & ~23000x space**.

GNNDelete is a General Strategy for Graph Unlearning!

- GNNDelete is a novel deletion operator that is flexible and easy-to-use and **can be used with any graph neural network (GNN) model**
- We formulate two **key requirements that graph unlearning methods must satisfy, Deleted Edge Consistency** and **Neighborhood Influence** through which we can unlearn graph elements and retain strong predictive performance
- GNNDelete achieves **state-of-the-art performance across a wide range of deletion tasks including edge deletion, node deletion, and node feature unlearning**



openreview.net/pdf?id=X9yCkmT5Qrl



github.com/mims-harvard/GNNDelete